

REPRESENTATION LEARNING FOR EQUIVARIANT INFERENCE WITH GUARANTEES

Anonymous authors

Paper under double-blind review

ABSTRACT

In many real-world applications of regression, conditional probability estimation, and uncertainty quantification, exploiting symmetries rooted in physics or geometry can dramatically improve generalization and sample efficiency. While geometric deep learning has made significant empirical advances by incorporating group-theoretic structure, less attention has been given to statistical learning guarantees. In this paper, we introduce an equivariant representation learning framework that simultaneously addresses regression, conditional probability estimation, and uncertainty quantification while providing first-of-its-kind non-asymptotic statistical learning guarantees. Grounded in operator and group representation theory, our framework approximates the spectral decomposition of the conditional expectation operator, building representations that are both equivariant and disentangled along independent symmetry quotient groups. Empirical evaluations on synthetic datasets and real-world robotics applications confirm the potential of our approach, matching or outperforming existing equivariant baselines in regression while additionally providing well-calibrated parametric uncertainty estimates.¹

1 INTRODUCTION

A central problem in machine learning is modeling conditional probabilities—understanding how the distribution of a target variable y changes with an observed variable x . This underpins robust reasoning under uncertainty in critical applications such as medicine, finance, robotics, and physics (Izbicki, 2025; Smith, 2013; Wasserman, 2007). However, estimating conditional distributions remains challenging in high-dimensional settings without strong inductive biases (Scott, 1991; Nagler and Czado, 2016; Izbicki and B. Lee, 2017).

Symmetry priors, in the form of principled assumptions about invariance or equivariance in the underlying data-generating process, offer a compelling way to reduce sample complexity and improve generalization (Kashinath et al., 2021; Wang, 2021; Bronstein et al., 2021; Elesedy, 2023). These priors naturally arise in inference tasks in chemistry and particle physics (Dresselhaus et al., 2007), set-&-graph structured data (Bronstein et al., 2021), computer graphics (Mitra et al., 2013; Kovar et al., 2002), and dynamical systems with group-invariant/equivariant laws of motion, which are ubiquitous in fields like physics (Dresselhaus et al., 2007), fluid dynamics (Olver, 1993), and robotics (Ordoñez-Apaez et al.; Zhu et al., 2022).

Over the past few years, Geometric Deep Learning (GDL) has produced a rich ecosystem of architectures that encode symmetries, achieving strong empirical performance across various supervised (Bronstein et al., 2021; Weiler et al., 2023; van der Pol et al., 2020) and unsupervised tasks (Dangovski et al.; Keurti et al., 2023; Wang et al., 2024a). However, the field remains focused on application specific designs and architectural innovation, with limited understanding of how symmetry priors can be leveraged to *learn representations with provable generalization guarantees*.

In this work, we take a different route: rather than proposing new architectures or solving specific inference tasks, we ask *how to systematically learn symmetry-aware representations that best capture conditional structure in the data*. Specifically, how should equivariant networks be trained so that their learned features reveal conditional distributions, and how does the quality of these representations affect performance in downstream tasks such as regression and uncertainty quantification?

¹All experiments and examples are provided in the open-access repository [symm_rep_learn](#)

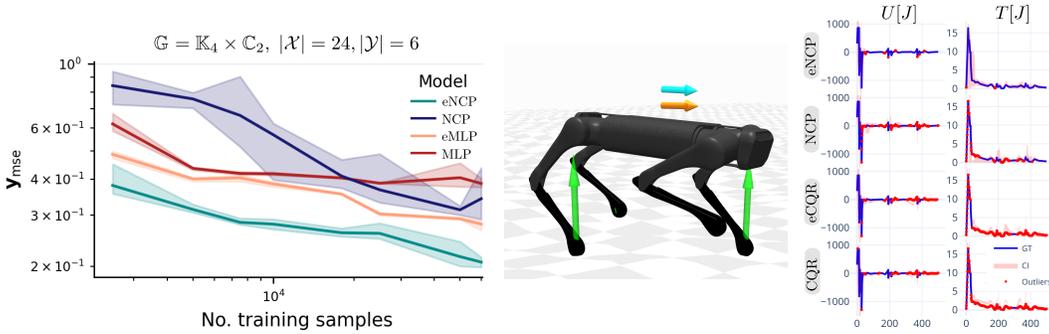


Figure 1: **Left:** Test set sample efficiency for \mathbb{G} -equivariant regression (MSE vs. training samples) when predicting the \mathbb{G} -equivariant linear and angular momentum of a quadruped robot’s center of mass (CoM) from noisy joint positions and velocities. **Right:** Uncertainty quantification via \mathbb{G} -equivariant prediction of 90% confidence intervals (CI, light-red area) for the robot’s instantaneous work U_t and kinetic energy T_t during locomotion over rough terrain for our method (eNCP) and competitors. The figure shows a trajectory with a strong initial disturbance, where blue markers denote samples within the predicted CI and red markers denote those outside. Note that only eNCP is able to predict well-calibrated CI intervals that cover both the disturbance and recovery phases.

To answer these questions, we bridge two fields rarely studied together: spectral contrastive learning (Zou et al., 2013), a self-supervised approach that learns deep representations of data via operator-theoretic modeling of conditional expectations (Tsai et al., 2021; Kostic et al., 2024a), and GDL (Bronstein et al., 2021), which enforces symmetry priors as architectural constraints in Neural Networks (NNs). Our approach shows how symmetry constraints shape the representation space and enhance generalization, opening new avenues for cross-fertilization between these fields. Concretely, we demonstrate that our method outperforms GDL techniques on regression tasks while providing reliable uncertainty quantification on a challenging robot locomotion task (see Fig. 1).

Contributions (1) Methodological framework: We introduce **Equivariant Neural Conditional Probability (eNCP)**, the first framework to combine equivariant neural networks with operator-theoretic estimation of conditional distributions. (2) Task-agnostic representation learning: We show that any \mathbb{G} -equivariant architecture can be used to learn *disentangled, symmetry-respecting representations* that generalize across diverse downstream inference tasks. (3) Learning guarantees: By linking the representation quality directly to sample complexity, we provide the *first non-asymptotic statistical learning guarantees* for equivariant conditional models, including regression and uncertainty quantification. (4) Empirical results: On both synthetic and real-world robotics tasks, eNCP consistently outperforms baselines, including contrastive methods **Neural Conditional Probability (NCP)** (Kostic et al., 2024a) and current equivariant models. In particular, eNCP achieves state-of-the-art performance in the challenging task of contact force inference in quadruped locomotion.

Paper structure Sec. 2 reviews modeling conditional probabilities with linear operators and NCP. Sec. 3 formally presents the symmetry priors we consider. Sec. 4 introduces our eNCP learning framework. Sec. 5 outlines our theoretical learning guarantees. Sec. 6 showcases experiments on synthetic and real-world data. Furthermore, because the paper involves complex notation from probability, operator theory, and group theory, the appendices include a glossary of notation (Sec. A) as well as detailed expositions on representation theory (Sec. I), symmetric function spaces (Sec. J), and equivariant linear operators (Sec. K). Finally, Sec. C offers an in-depth discussion of related work, contrasting our framework with the literature across these rich fields.

2 BACKGROUND

We briefly review the operator-theoretic framework for modeling conditional probabilities, which underpins both NCP and our proposed eNCP method. We denote a random variable by \mathbf{x} , its realizations by $x \in \mathcal{X}$, its probability distribution by $\mathbb{P}(\mathbf{x})$, and measure by $P_{\mathbf{x}}$. Expectations are written as $\mathbb{E}_{\mathbf{x}}[f(\mathbf{x})] = \int_{\mathcal{X}} f(x) P_{\mathbf{x}}(dx)$. The same notation applies to other random variables such as y .

Operator-theoretic modeling of conditional probabilities Kostic et al. (2024a) proposed to model conditional probabilities by approximating the *conditional expectation operator* (Baker,

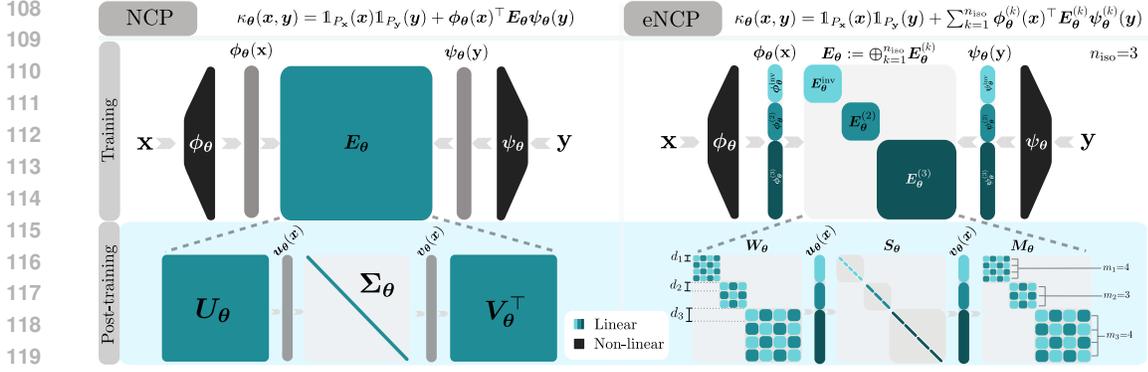


Figure 2: **Left:** NCP’s bilinear NN architecture. **Right:** eNCP’s \mathbb{G} -equivariant bilinear NN architecture, featuring ϕ_θ and ψ_θ as \mathbb{G} -equivariant NNs and E_θ as a \mathbb{G} -equivariant block-diagonal matrix. Each block is equivariant to a quotient group $\mathbb{G}^{(k)} \subseteq \mathbb{G}$ and is constrained to have singular spaces of dimension at least d_k —the minimal dimension for a representation of the action of $\mathbb{G}^{(k)}$.

1973; Song et al., 2009; Ryu et al., 2024), $E_{y|x}: \mathcal{L}_y^2 \rightarrow \mathcal{L}_x^2$, a linear integral operator acting on the Hilbert spaces $\mathcal{L}_x^2 := \mathcal{L}_{P_x}^2(\mathcal{X}, \mathbb{R})$ and $\mathcal{L}_y^2 := \mathcal{L}_{P_y}^2(\mathcal{Y}, \mathbb{R})$ of square-integrable functions of the random variables \mathbf{x} and \mathbf{y} , respectively. The action of this operator on any function $h \in \mathcal{L}_y^2$ returns the function’s conditional expectation:

$$[E_{y|x}h](\mathbf{x}) = \mathbb{E}[h(\mathbf{y})|\mathbf{x}=\mathbf{x}] := \int_{\mathcal{Y}} h(\mathbf{y})P_{y|x}(d\mathbf{y}|\mathbf{x}) = \int_{\mathcal{Y}} h(\mathbf{y})\frac{P_{\mathbf{xy}}(d\mathbf{y}, d\mathbf{x})}{P_x(d\mathbf{x})} = \int_{\mathcal{Y}} h(\mathbf{y})\kappa(\mathbf{x}, \mathbf{y})P_y(d\mathbf{y}), \quad (1)$$

where $P_{y|x}$ is the conditional probability measure, and $\kappa(\mathbf{x}, \mathbf{y}) := \frac{P_{\mathbf{xy}}(d\mathbf{x}, d\mathbf{y})}{P_x(d\mathbf{x})P_y(d\mathbf{y})}$ is the **Pointwise Mutual Dependency (PMD)** (Sugiyama et al., 2012) kernel defining $E_{y|x}$, obtained as the Radon-Nikodym derivative of the joint measure to the product of marginal measures (see Fig. 3 and Sec. H).

The conditional expectation operator is significant because it provides an infinite-dimensional linear model—in a nonlinear representation space—for computing conditional probabilities and expectations. To see this, note that for any $\mathbf{x} \in \mathcal{X}$ and any measurable set $\mathbb{B} \subset \mathcal{Y}$ we have that:

$$\mathbb{P}(\mathbf{y} \in \mathbb{B}|\mathbf{x}=\mathbf{x}) := \int_{\mathcal{Y}} \mathbb{1}_{\mathbb{B}}(\mathbf{y})P_{y|x}(d\mathbf{y}|\mathbf{x}) = [E_{y|x}\mathbb{1}_{\mathbb{B}}](\mathbf{x}), \quad \text{and} \quad \mathbb{E}[\mathbf{y}|\mathbf{x}=\mathbf{x}] := [E_{y|x}\mathbf{y}](\mathbf{x}). \quad (2)$$

Therefore, to estimate conditional probabilities and expectations, NCP seeks the best finite-dimensional approximation of $E_{y|x}$. As we explain next, this gives rise to a representation learning problem (Oord et al., 2018), in which the optimal representations of \mathbf{x} and \mathbf{y} are given by the top left and right singular functions of $E_{y|x}$ (Kostic et al., 2024b; Ryu et al., 2024).

Spectral representation learning The problem of approximating the conditional expectation operator $E_{y|x}$ as a rank- r operator E_θ with matrix representation $\mathbf{E}_\theta \in \mathbb{R}^{r \times r}$ is defined as

$$\arg \min_{\theta} \|\mathbb{E}_{y|x} - E_\theta\|_{\text{HS}}^2 = \mathbb{E}_x \mathbb{E}_y (\kappa(\mathbf{x}, \mathbf{y}) - \kappa_\theta(\mathbf{x}, \mathbf{y}))^2, \quad \text{s.t. } \mathbb{E}_x \mathbb{E}_y \kappa_\theta(\mathbf{x}, \mathbf{y}) = 1 \text{ and } \text{rank}(\mathbf{E}_\theta) \leq r. \quad (3)$$

The optimal solution, denoted E_* , is the r -truncated **Singular Value Decomposition (SVD)** of $E_{y|x}$ (Eckart and Young, 1936; Weidmann, 2012), namely

$$[E_*f](\mathbf{x}) = \sum_{i=0}^r \sigma_i \langle f, v_i \rangle_{P_y} u_i(\mathbf{x}), \quad \text{with } \sigma_i u_i(\mathbf{x}) = [E_{y|x}v_i](\mathbf{x}), \quad \forall i \in [r], \quad (4)$$

where (σ_i, u_i, v_i) denotes the i^{th} singular value and left/right singular functions of $E_{y|x}$, with $(\sigma_0=1, u_0=\mathbb{1}_{P_x}, v_0=\mathbb{1}_{P_y})$ being the constant functions supported on P_x and P_y (Baker, 1973).

Consequently, NCP parameterizes E_θ by a bilinear model $\kappa_\theta(\mathbf{x}, \mathbf{y}) = 1 + \phi_\theta(\mathbf{x})^\top \mathbf{E}_\theta \psi_\theta(\mathbf{y})$, composed of two encoder NNs $\phi_\theta: \mathcal{X} \rightarrow \mathbb{R}^r$ and $\psi_\theta: \mathcal{Y} \rightarrow \mathbb{R}^r$ that aim to approximate the *span* of the top r (non-constant) left and right singular functions of $E_{y|x}$. See Fig. 2-left.

As κ is generally unavailable analytically, (3) is solved via the regularized **contrastive low-rank (cLoRa)** loss:

$$\begin{aligned} \mathcal{L}_\gamma(\theta) &= -2\mathbb{E}_{\mathbf{xy}} \kappa_\theta(\mathbf{x}, \mathbf{y}) + \mathbb{E}_x \mathbb{E}_y \kappa_\theta(\mathbf{x}, \mathbf{y})^2 + 2\gamma (\|\mathbb{E}_x \phi_\theta(\mathbf{x})\|_F^2 + \|\mathbb{E}_y \psi_\theta(\mathbf{y})\|_F^2) \\ &\quad + \|\text{Cov}(\phi_\theta) - \mathbf{I}_r\|_F^2 + \|\text{Cov}(\psi_\theta) - \mathbf{I}_r\|_F^2, \end{aligned} \quad (5)$$

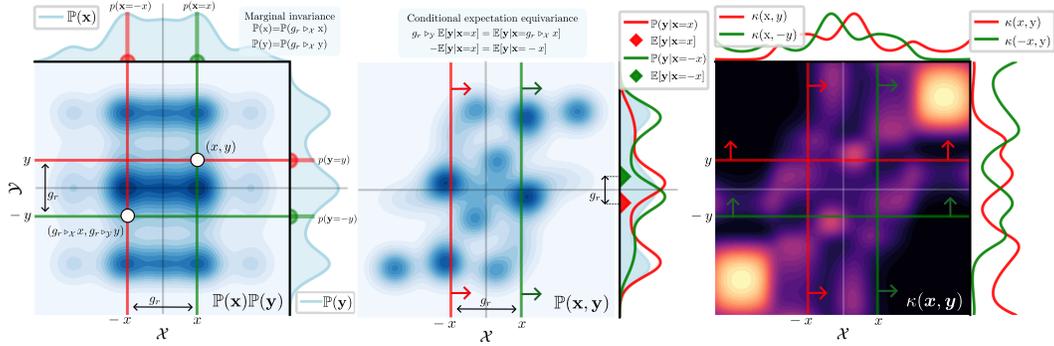


Figure 3: Example of symmetric random variables $(x, y) \sim \mathcal{X} \times \mathcal{Y} \subset \mathbb{R} \times \mathbb{R}$, whose marginals $\mathbb{P}(x)$ and $\mathbb{P}(y)$; joint $\mathbb{P}(x, y)$; and conditional $\mathbb{P}(y|x)$ distributions are invariant to reflections of the data: $g_r \triangleright_x x = -x$ and $g_r \triangleright_y y = -y$, where g_r denotes the reflection element of the reflection symmetry group $\mathbb{C}_2 := \{e, g_r | g_r^2 = e\}$. Consequently, the PMD $\kappa(x, y)$ is \mathbb{C}_2 -invariant.

where the first two regularization terms center the learned representations, ensuring that $\mathbb{E}_x \mathbb{E}_y \kappa_\theta(x, y) \approx 1$ (Kostic et al., 2024a), while the last two enforce approximate orthonormality of the learned bases in $\mathcal{F}_x^\theta := \text{span}(\phi_\theta) \subset \mathcal{L}_x^2$ and $\mathcal{F}_y^\theta := \text{span}(\psi_\theta) \subset \mathcal{L}_y^2$ (Izibicki and B. Lee, 2017). A key property of NCP is that the learned representations enables reliable regression and conditional probability estimation—and thus uncertainty quantification—via (2) (see Tab. 3).

3 PROBLEM FORMULATION

This paper tackles the problem of estimating the conditional expectation $\mathbb{E}[y|x = \cdot]$, and, more generally, conditional distribution $\mathbb{P}(y|x)$, for random variables $x \in \mathcal{X}$ and $y \in \mathcal{Y}$, under the assumption that $\mathbb{P}(y|x)$ and $\mathbb{P}(x)$ are \mathbb{G} -invariant under symmetry transformations of the data (see Fig. 3), i.e:

$$\mathbb{P}(y|x) = \mathbb{P}(g \triangleright_y y | g \triangleright_x x), \quad \mathbb{P}(x) = \mathbb{P}(g \triangleright_x x), \quad \forall g \in \mathbb{G}, \quad (6)$$

where \mathbb{G} denotes a finite symmetry group (Thm. I.1) acting on the data spaces \mathcal{X} and \mathcal{Y} via the group actions, $\triangleright_x: \mathbb{G} \times \mathcal{X} \rightarrow \mathcal{X}$, and $\triangleright_y: \mathbb{G} \times \mathcal{Y} \rightarrow \mathcal{Y}$, with $g \triangleright_x x \in \mathcal{X}$ and $g \triangleright_y y \in \mathcal{Y}$ denoting linear, invertible transformations of x and y defined by $g \in \mathbb{G}$ (see Fig. 3 and Def. I.2).

These priors imply the \mathbb{G} -invariance of the joint distribution $\mathbb{P}(x, y)$ and of y 's marginal distribution $\mathbb{P}(y)$, as well as the \mathbb{G} -equivariance of conditional expectations (see Fig. 3-middle and Thm. D.1):

$$g \triangleright_y \mathbb{E}[y|x=x] = \mathbb{E}[y|x=g \triangleright_x x] \quad \forall g \in \mathbb{G}, x \in \mathcal{X}. \quad (7)$$

Note that (7) implies the \mathbb{G} -equivariance of the regression function $x \mapsto \mathbb{E}[y|x=x]$. Therefore, the symmetry priors (6) are satisfied whenever we approximate an equivariant/invariant function—that is, in virtually all applications of GDL (Bronstein et al., 2021).

The above symmetry priors represent a strong inductive bias for the conditional expectation operator (2), as they lead the PMD kernel defining the operator to be \mathbb{G} -invariant (see Fig. 3-right):

$$\kappa(x, y) = \kappa(g \triangleright_x x, g \triangleright_y y) \quad \forall g \in \mathbb{G}, x \in \mathcal{X}, y \in \mathcal{Y}. \quad (8)$$

In Sec. 4, we extend the NCP framework (Kostic et al., 2024a) by leveraging (8) to incorporate symmetry priors. This enables efficient estimation of \mathbb{G} -invariant conditional probabilities (6) and \mathbb{G} -equivariant regression (7) using GDL architectures, via (2), with strong learning guarantees.

4 ENCP METHOD FOR EQUIVARIANT REPRESENTATION LEARNING

In this section, we show how to incorporate the symmetry priors (6) into NCP's representation learning framework. First, we analyze the symmetry constraints on the infinite-dimensional conditional expectation operator and prove that, for symmetric random variables x and y , the optimal solution

²Throughout, with some abuse of notation we denote by $\mathbb{P}(x)$ and $\mathbb{P}(y|x)$ both the probability and conditional probability, respectively, as well as the corresponding densities, when they exist.

of (3) yields \mathbb{G} -equivariant representations ϕ_θ and ψ_θ and approximates the operator with a \mathbb{G} -equivariant matrix \mathbf{E}_θ . Then, we explain how to embed these structural constraints into the bilinear neural network architecture of NCP using any type of equivariant NNs.

Symmetric function spaces The assumption of \mathbb{G} -invariance of the marginal probabilities (Sec. 3) implies that the function spaces \mathcal{L}_x^2 and \mathcal{L}_y^2 are symmetric Hilbert spaces of \mathbb{G} -equivariant functions, as these inherit unitary group actions $\triangleright_{\mathcal{L}_x^2}: \mathbb{G} \times \mathcal{L}_x^2 \rightarrow \mathcal{L}_x^2$ and $\triangleright_{\mathcal{L}_y^2}: \mathbb{G} \times \mathcal{L}_y^2 \rightarrow \mathcal{L}_y^2$ defined via the push-forward of symmetry transformations of the data spaces (see details in Sec. J and in Fig. 16):

$$g \triangleright_{\mathcal{L}_x^2} f(\cdot) := f(g^{-1} \triangleright_x \cdot) \in \mathcal{L}_x^2, \quad g \triangleright_{\mathcal{L}_y^2} h(\cdot) := h(g^{-1} \triangleright_y \cdot) \in \mathcal{L}_y^2, \quad \forall g \in \mathbb{G}. \quad (9)$$

A fundamental property of \mathbb{G} -symmetric Hilbert spaces is their orthogonal decomposition into $n_{\text{iso}} \leq |\mathbb{G}|$ subspaces referred to as *isotypic subspaces*: $\mathcal{L}_x^2 = \bigoplus_{k \in [1, n_{\text{iso}}]} \mathcal{L}_x^{2(k)}$, and $\mathcal{L}_y^2 = \bigoplus_{k \in [1, n_{\text{iso}}]} \mathcal{L}_y^{2(k)}$ (see Thm. I.16). Where each $\mathcal{L}_x^{2(k)}$ and $\mathcal{L}_y^{2(k)}$ denote the spaces of $\mathbb{G}^{(k)}$ -equivariant functions of \mathbf{x} and \mathbf{y} , with $\mathbb{G}^{(k)} := \mathbb{G}/\mathbb{N}_k$ being a **quotient subgroup**, generated by a normal subgroup \mathbb{N}_k defined by the kernel of the group's k^{th} irreducible representation (see Subsec. I.2). This standard result from harmonic analysis (Mackey, 1980) enables us to express any \mathbb{G} -equivariant function as a sum of its projections onto the isotypic subspaces:

$$f(\cdot) = f^{\text{inv}}(\cdot) + \sum_{k=2}^{n_{\text{iso}}} f^{(k)}(\cdot), \quad h(\cdot) = h^{\text{inv}}(\cdot) + \sum_{k=2}^{n_{\text{iso}}} h^{(k)}(\cdot), \quad \text{s.t. } f^{(k)} \in \mathcal{L}_x^{2(k)}, h^{(k)} \in \mathcal{L}_y^{2(k)}, \forall k \in [n_{\text{iso}}], \quad (10)$$

where $f^{(k)}$ and $h^{(k)}$ denote the $\mathbb{G}^{(k)}$ -equivariant components of f and h . Moreover, by convention, we associate the first subspace ($k = 1$) with the space of \mathbb{G} -invariant functions, i.e., $\mathbb{G}^{(1)} = \mathbb{G}^{\text{inv}} = \{e\}$ (see Thm. J.4 in the Appendix).

Equivariant conditional expectation operator The \mathbb{G} -invariance of the PMD kernel (8), implies that $\mathbb{E}_{y|x}$ is a \mathbb{G} -equivariant linear operator (see Thm. K.1). This means that $\mathbb{E}_{y|x}$ commutes with the group action on the function spaces, and consequently, can be decomposed (disentangled) into a direct sum of operators acting on the corresponding isotypic subspaces (see details in Sec. K):

$$g \triangleright_{\mathcal{L}_x^2} [\mathbb{E}_{y|x} h](\cdot) = \mathbb{E}_{y|x} [g \triangleright_{\mathcal{L}_y^2} h](\cdot) \iff [\mathbb{E}_{y|x} h](\cdot) = \sum_{k=1}^{n_{\text{iso}}} [\mathbb{E}_{y|x}^{(k)} h^{(k)}](\cdot) \quad \forall h \in \mathcal{L}_y^2, g \in \mathbb{G}, \quad (11)$$

where each $\mathbb{E}_{y|x}^{(k)}: \mathcal{L}_y^{2(k)} \rightarrow \mathcal{L}_x^{2(k)}$ models the conditional expectation for $\mathbb{G}^{(k)}$ -equivariant functions.

Equivariant disentangled representation learning The \mathbb{G} -equivariant structure of $\mathbb{E}_{y|x}$ and its disentanglement (11) into isotypic components suggests that computing the conditional expectation of a \mathbb{G} -equivariant function is equivalent to summing the conditional expectations of its $\mathbb{G}^{(k)}$ -equivariant components for all $k \in [n_{\text{iso}}]$. Therefore, the loss function of problem (3), where $\mathbb{E}_{y|x}$ is approximated in finite dimensional spaces \mathcal{F}_x^θ and \mathcal{F}_y^θ , decouples into n_{iso} independent (disentangled) components:

$$\begin{aligned} \arg \min_{\theta} \|\mathbb{E}_{y|x} - \mathbf{E}_\theta\|_{\text{HS}}^2 &= \sum_{k=1}^{n_{\text{iso}}} \|\mathbb{E}_{y|x}^{(k)} - \mathbf{E}_\theta^{(k)}\|_{\text{HS}}^2 = \mathbb{E}_x \mathbb{E}_y \sum_{k=1}^{n_{\text{iso}}} (\kappa^{(k)}(\mathbf{x}, \mathbf{y}) - \kappa_\theta^{(k)}(\mathbf{x}, \mathbf{y}))^2, \\ \text{s.t. } \mathbb{E}_x \mathbb{E}_y \kappa_\theta(\mathbf{x}, \mathbf{y}) &= 1, \quad \text{and } \kappa_\theta(g \triangleright_x \mathbf{x}, g \triangleright_y \mathbf{y}) = \kappa_\theta(\mathbf{x}, \mathbf{y}), \quad \forall g \in \mathbb{G}, (\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{Y}. \end{aligned} \quad (12)$$

Importantly, to satisfy the \mathbb{G} -invariance constraint on κ_θ , the truncated operator $\mathbf{E}_\theta: \mathcal{F}_y^\theta \rightarrow \mathcal{F}_x^\theta$ must act on symmetric finite-dimensional Hilbert spaces that are stable under \mathbb{G} , i.e., $g \triangleright_{\mathcal{L}_x^2} f \in \mathcal{F}_x^\theta$ and $g \triangleright_{\mathcal{L}_y^2} h \in \mathcal{F}_y^\theta$ for all $f \in \mathcal{F}_x^\theta$ and $h \in \mathcal{F}_y^\theta$, and have a \mathbb{G} -equivariant matrix representation \mathbf{E}_θ (see Thm. K.2). Thus, as in the infinite-dimensional case, these finite-dimensional spaces decompose into n_{iso} isotypic subspaces $\mathcal{F}_x^\theta = \bigoplus_{k=1}^{n_{\text{iso}}} \mathcal{F}_x^{\theta(k)}$ and $\mathcal{F}_y^\theta = \bigoplus_{k=1}^{n_{\text{iso}}} \mathcal{F}_y^{\theta(k)}$. Accordingly, the truncated operator matrix decomposes block-diagonally into n_{iso} blocks, i.e., $\mathbf{E}_\theta = \bigoplus_{k=1}^{n_{\text{iso}}} \mathbf{E}_\theta^{(k)}$, where each k -th block needs to be a $\mathbb{G}^{(k)}$ -equivariant matrix approximating the restriction of the conditional expectation operator on its corresponding isotypic subspace (Salova et al., 2019; Ordoñez-Apraéz et al., 2024), see the upper block of Fig. 2-right.

Moreover, analogous to (4), the optimal truncation of $\mathbb{E}_{y|x}^{(k)}$ is its truncated SVD (Weidmann, 2012). However, the $\mathbb{G}^{(k)}$ -equivariance of each disentangled component imposes specific structure on the SVD, since any symmetry-transformed singular function remains within the same singular space (Ordoñez-Apraéz et al., 2024). Consequently, each singular space must have at least dimension

d_k —the smallest real-vector-space dimension in which $\mathbb{G}^{(k)}$ can be faithfully represented, as depicted in Fig. 2-right and formalized in Thm. K.3.

To solve \mathbb{G} -equivariant regression and estimate \mathbb{G} -invariant conditional probabilities via (12), we propose the disentangled bilinear NN setup in Fig. 2-right. This approach supports any GDL \mathbb{G} -equivariant backbone (CNN, GNN, Transformer), adapting to diverse data modalities and tasks. Our framework can handle *continuous* compact groups via group discretization, and *non-compact* groups by selecting appropriate backbones, as it is done with \mathbb{G} -steerable CNNs for image processing with rotation/reflection and translation equivariance (Cesa et al., 2022).

Our equivariant and disentangled bilinear NN architecture directly models the isotypic decomposition of $\mathbb{E}_{y|x}$. We use two \mathbb{G} -equivariant encoder NNs, $\phi_\theta : \mathcal{X} \rightarrow \mathbb{R}^r$ and $\psi_\theta : \mathcal{Y} \rightarrow \mathbb{R}^r$, parameterized to expose the isotypic subspaces, i.e., $\phi_\theta(\cdot) = [\phi_\theta^{\text{inv}}(\cdot)^\top, \dots, \phi_\theta^{(n_{\text{iso}})}(\cdot)^\top]^\top$ and $\psi_\theta(\cdot) = [\psi_\theta^{\text{inv}}(\cdot)^\top, \dots, \psi_\theta^{(n_{\text{iso}})}(\cdot)^\top]^\top$. Here, each component $\phi_\theta^{(k)} : \mathcal{X} \rightarrow \mathbb{R}^{r_k}$ and $\psi_\theta^{(k)} : \mathcal{Y} \rightarrow \mathbb{R}^{r_k}$ spawns the approximated isotypic subspace of $\mathbb{G}^{(k)}$ -equivariant functions, i.e., $\mathcal{F}_x^{\theta^{(k)}} := \text{span}(\phi_\theta^{(k)}) \subset \mathcal{L}_x^{2(k)}$ and $\mathcal{F}_y^{\theta^{(k)}} := \text{span}(\psi_\theta^{(k)}) \subset \mathcal{L}_y^{2(k)}$. Further, the truncated operator’s matrix is parameterized in block-diagonal form $\mathbf{E}_\theta = \bigoplus_{k=1}^{n_{\text{iso}}} \mathbf{E}_\theta^{(k)}$, with each block an $r_k \times r_k$ $\mathbb{G}^{(k)}$ -equivariant matrix³. The corresponding approximated PMD kernel is given by:

$$\kappa_\theta(\mathbf{x}, \mathbf{y}) = \mathbb{1}_{P_x}(\mathbf{x})\mathbb{1}_{P_y}(\mathbf{y}) + \sum_{k=1}^{n_{\text{iso}}} \kappa_\theta^{(k)}(\mathbf{x}, \mathbf{y}), \quad \kappa_\theta^{(k)}(\mathbf{x}, \mathbf{y}) := \phi_\theta^{(k)}(\mathbf{x})^\top \mathbf{E}_\theta^{(k)} \psi_\theta^{(k)}(\mathbf{y}), \quad (13)$$

where $\mathbb{1}_{P_x}(\mathbf{x})\mathbb{1}_{P_y}(\mathbf{y})$ arises since the first singular functions of $\mathbb{E}_{y|x}$ are constant, see (4).

This parameterization inherently preserves the symmetry constraints of each operator’s singular functions, which we leverage in both theory and practice (see Sec. E and Subsec. K.2 for details).

Disentangled training loss Having introduced the equivariance constraints on the truncated operator matrix, we decompose the contrastive loss (5) to reflect the separability of the optimization arising from the operator’s isotypic decomposition (11):

$$\mathcal{L}_\gamma(\theta) := \sum_{k=1}^{n_{\text{iso}}} (-2\mathbb{E}_{\mathbf{x}\mathbf{y}} \kappa_\theta^{(k)}(\mathbf{x}, \mathbf{y}) + \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{y}} \kappa_\theta^{(k)}(\mathbf{x}, \mathbf{y})^2 + \gamma \Omega^{(k)}(\theta)) + 2\gamma (\|\mathbb{E}_{\mathbf{x}} \phi_\theta^{\text{inv}}(\mathbf{x})\|_F^2 + \|\mathbb{E}_{\mathbf{y}} \psi_\theta^{\text{inv}}(\mathbf{y})\|_F^2). \quad (14)$$

This decomposes learning \mathbb{G} -equivariant representations of \mathbf{x} and \mathbf{y} into learning n_{iso} less constrained $\mathbb{G}^{(k)}$ -equivariant representations for distinct quotient groups of \mathbb{G} . Such representations are known in the literature as *disentangled* representations (Higgins et al., 2018) (see Thm. I.18).

Moreover, we improve the estimates of the regularization terms in (5) by leveraging our symmetry priors to: (i) tighten the centering regularization (14) given that functions in $\mathcal{F}_x^{(k)}$ and $\mathcal{F}_y^{(k)}$ are centered by construction for $k \neq \text{inv}$ (see Thm. L.4)—and (ii) exploit the orthogonality between isotypic subspaces (10) to independently regularize orthonormality for each isotypic subspace (see example in Fig. 9), leading to better covariance estimates (Shah and Chandrasekaran, 2012):

$$\Omega^{(k)}(\theta) := \sum_{k=1}^{n_{\text{iso}}} \|\text{Cov}(\phi^{(k)}) - \mathbf{I}_{r_k}\|_F^2 + \|\text{Cov}(\psi^{(k)}) - \mathbf{I}_{r_k}\|_F^2. \quad (15)$$

Given a batch $\{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^N$ and their corresponding embeddings $\{(\phi_\theta(\mathbf{x}_n), \psi_\theta(\mathbf{y}_n))\}_{n=1}^N$, the empirical unregularized loss is estimated via U-statistics, yielding an unbiased estimate with an effective sample size of N^2 (Wang et al., 2022b; Tsai et al., 2020).

$$\widehat{\mathcal{L}}_0(\theta) = \sum_{k \in [n_{\text{iso}}]} \left[\frac{1}{N} \sum_{n \in [N]} \kappa_\theta^{(k)}(\mathbf{x}_n, \mathbf{y}_n) + \frac{1}{N(N-1)} \sum_{a \in [N]} \sum_{b \in [N] \setminus \{a\}} \kappa_\theta^{(k)}(\mathbf{x}_a, \mathbf{y}_b)^2 \right]. \quad (16)$$

Similarly, we use U-statistics to obtain unbiased estimates for orthonormal regularization in (15), achieving an effective sample size of $d_k N^2$ per isotypic subspace (see Subsec. F.2). Consequently, standard NN optimization methods can be employed to learn equivariant representations via the approximate model of $\mathbb{E}_{y|x}$, enabling downstream inference tasks described in the next section.

5 INFERENCE AND LEARNING GUARANTEES

Once training is complete, the learned \mathbb{G} -invariant PMD from (13) can be used, via (2), for \mathbb{G} -equivariant regression and \mathbb{G} -invariant conditional probability estimation. These estimates are obtained using a NN architecture composed of ϕ_θ , \mathbf{E}_θ , and a final linear layer that delivers the basis expansion coefficients of the target variable in the \mathbf{y} representation space $\mathcal{F}_y^\theta = \text{span}(\psi_\theta)$. For a summary of the estimates and their learning guarantees refer to Tab. 1.

³We chose square matrices for notational convenience. Dimensions for \mathbf{y} and \mathbf{x} spaces need not match

324	Task	$f(\mathbf{x}) := \mathbb{E}_{\mathbf{y}}[\mathbf{y} \mathbf{x}=\mathbf{x}] \approx \hat{f}_{\theta}(\mathbf{x})$	$\mathbb{P}[\mathbf{y} \in \mathbb{B} \mathbf{x} \in \mathbb{A}] \approx \hat{\mathbb{P}}_{\theta}[\mathbf{y} \in \mathbb{B} \mathbf{x} \in \mathbb{A}]$
325			
326	Estimate	$\hat{\mathbb{E}}_{\mathbf{y}}[\mathbf{y}] + \phi_{\theta}(\mathbf{x})^{\top} \mathbf{E}_{\theta} \hat{\mathbb{E}}_{\mathbf{y}}[\psi_{\theta}(\mathbf{y}) \otimes \mathbf{y}]$	$\hat{\mathbb{E}}_{\mathbf{y}}[\mathbb{1}_{\mathbb{B}}] + \frac{\hat{\mathbb{E}}_{\mathbf{x}}[\mathbb{1}_{\mathbb{A}}(\mathbf{x}) \otimes \phi_{\theta}(\mathbf{x})]^{\top} \mathbf{E}_{\theta} \hat{\mathbb{E}}_{\mathbf{y}}[\mathbb{1}_{\mathbb{B}}(\mathbf{y}) \otimes \psi_{\theta}(\mathbf{y})]}{\hat{\mathbb{E}}_{\mathbf{x}}[\mathbb{1}_{\mathbb{A}}(\mathbf{x})]}$
327			
328	Guarantees	$\ f - \hat{f}_{\theta}\ _{\mathcal{L}_{\mathbf{x}}^2} \lesssim \sqrt{\text{Var}[\ \mathbf{y}\]} \left(\mathcal{E}_{\theta}^r + \frac{\ln(n_{\text{iso}}/\delta)}{(d_{\text{iso}}N)^{\frac{\alpha}{1+2\alpha}}} \right)$	$ \mathbb{P} - \hat{\mathbb{P}}_{\theta} \lesssim \sqrt{\frac{\mathbb{P}[\mathbf{y} \in \mathbb{B}]}{\mathbb{P}[\mathbf{x} \in \mathbb{G}_{\triangleright \mathbf{x}} \mathbb{A}]}} \left(\mathcal{E}_{\theta}^r + \frac{\ln(n_{\text{iso}}/\delta)}{(d_{\text{iso}}N)^{\frac{\alpha}{1+2\alpha}}} \right)$

330 Table 1: Statistical guarantees for eNCP. The error bounds are shaped by (i) the structure of the
 331 symmetry group \mathbb{G} —the number of isotypic subspaces n_{iso} and their minimum singular space di-
 332 mensions $d_{\text{iso}} = \sum_{k=1}^{n_{\text{iso}}} d_k$ (see Fig. 2), which enlarge the effective sample size—, (ii) the quality
 333 of the learned representations $\mathcal{E}_{\theta}^r = \|\mathbb{E}_{\mathbf{y}|\mathbf{x}} - \mathbf{E}_{\theta}\|_{\text{op}} \leq \sqrt{\mathcal{L}_{\gamma}(\theta) - \mathcal{L}_{\gamma}(\star)}$, and (iii) the operator’s
 334 singular-value decay rate $\alpha > 0$. Note that $\mathbb{G}_{\triangleright \mathcal{X}} \mathbb{A} := \cup_{g \in \mathbb{G}} g \triangleright_{\mathcal{X}} \mathbb{A}$ denotes the group orbit of \mathbb{A} .

335 Both estimates are derived from the general problem of vector-valued regression of a target function
 336 $\mathbf{z}: \mathcal{X} \rightarrow \mathcal{Z}$ defined by the conditional expectation of an observable $\mathbf{h} \in \mathcal{L}_{\mathcal{Y}}^2(\mathcal{Y}, \mathcal{Z})$, that is, $\mathbf{z}(\mathbf{x}) :=$
 337 $\mathbb{E}_{\mathbf{y}}[\mathbf{h}(\mathbf{y})|\mathbf{x} = \mathbf{x}] = [\mathbb{E}_{\mathbf{y}|\mathbf{x}}\mathbf{h}](\mathbf{x})$, where \mathcal{Z} is a symmetry-endowed vector space. Using the learned
 338 model, we estimate \mathbf{z} by

$$339 \quad [\mathbb{E}_{\mathbf{y}|\mathbf{x}}\mathbf{h}](\mathbf{x}) \approx \hat{\mathbf{z}}_{\theta}(\mathbf{x}) := \hat{\mathbb{E}}_{\mathbf{y}}[\mathbf{h}(\mathbf{y})] + \phi_{\theta}(\mathbf{x})^{\top} \mathbf{E}_{\theta} \hat{\mathbb{E}}_{\mathbf{y}}[\psi_{\theta}(\mathbf{y}) \otimes \mathbf{h}(\mathbf{y})], \quad (17)$$

341 where $\hat{\mathbb{E}}_{\mathbf{y}}[\psi_{\theta}(\mathbf{y}) \otimes \mathbf{h}(\mathbf{y})]$ represents the basis expansion coefficients of \mathbf{h} in the learned basis of
 342 $\mathcal{F}_{\mathbf{y}}^{\theta} \subset \mathcal{L}_{\mathbf{y}}^2$. Here, $\hat{\mathbb{E}}_{\mathbf{x}}: \mathcal{L}_{\mathbf{x}}^2 \rightarrow \mathbb{R}$ and $\hat{\mathbb{E}}_{\mathbf{y}}: \mathcal{L}_{\mathbf{y}}^2 \rightarrow \mathbb{R}$ are the \mathbb{G} -invariant empirical expectations defined
 343 by: $\hat{\mathbb{E}}_{\mathbf{x}}[f(\mathbf{x})] = \frac{1}{|\mathbb{G}|N} \sum_{g \in \mathbb{G}} \sum_{n=1}^N f(g \triangleright_{\mathcal{X}} \mathbf{x}_n)$ and $\hat{\mathbb{E}}_{\mathbf{y}}[h(\mathbf{y})] = \frac{1}{|\mathbb{G}|N} \sum_{g \in \mathbb{G}} \sum_{n=1}^N h(g \triangleright_{\mathcal{Y}} \mathbf{y}_n)$.

344 Hence, our method learns representations of \mathbf{x} and \mathbf{y} that transform *nonlinear regression of ob-*
 345 *servables* into *linear regression* in the learned representation space. For example, assuming \mathbf{y} has
 346 bounded variance and setting $\mathbf{h}(\mathbf{y}) = \mathbf{y}$, we recover the standard (\mathbb{G} -equivariant) regression so-
 347 lution (see Tab. 1-left). Equally important, by letting $\mathbf{h} = \mathbb{1}_{\mathbb{B}}$ —the indicator of a measurable set
 348 $\mathbb{B} \subseteq \mathcal{Y}$ —the model estimates conditional probabilities (see Tab. 1-right), thereby supporting both
 349 regression and uncertainty quantification (e.g., conditional quantiles, covariance; see Sec. 6 and
 350 (Kostic et al., 2024b)). The following learning bounds cover this general setting.

352 **Theorem 5.1.** *Let the symmetry priors in (6) hold, $\mathbb{E}_{\mathbf{y}|\mathbf{x}}$ be a $(1/\alpha)$ -Schatten-class \mathbb{G} -equivariant*
 353 *conditional expectation operator, as in (11), and \mathbf{E}_{θ} be a truncated approximation, defined as in (13)*
 354 *and trained via (14), with rank r and parameters $\theta \in \Theta$. Then, given an appropriate truncation*
 355 *dimension $r \asymp (N/d_{\text{iso}}^{\alpha})^{1/(1+2\alpha)}$, the following results hold for any \mathbb{G} -equivariant/invariant $\mathbf{h} \in$*
 356 *$\mathcal{L}_{\mathcal{Y}}^2(\mathcal{Y}, \mathcal{Z})$, measurable set $\mathbb{A} \subset \mathcal{X}$, and $\mathbb{G}' \leq \mathbb{G}$; with probability at least $1 - \delta$ w.r.t. an iid draw of*
 357 *$\mathbb{D}_N = \{(\mathbf{x}_n, \mathbf{y}_n) \sim P_{\mathbf{x}\mathbf{y}}\}_{n=1}^N$:*

$$358 \quad \|\mathbf{z} - \hat{\mathbf{z}}_{\theta}\|_{\mathcal{L}_{P_{\mathbf{x}}(\mathcal{X}, \mathcal{Z})}^2} \lesssim \sqrt{\text{Var}[\|\mathbf{h}(\mathbf{y})\|_{\mathcal{Z}}]} \left(\mathcal{E}_{\theta}^r + \frac{\ln(n_{\text{iso}}/\delta)}{(d_{\text{iso}}N)^{\frac{\alpha}{1+2\alpha}}} \right) \quad \text{and} \quad (18a)$$

$$360 \quad \|\mathbf{z}(\mathbb{A}) - \hat{\mathbf{z}}_{\theta}(\mathbb{A})\|_{\mathcal{Z}} \lesssim \frac{\sqrt{1 + (|\mathbb{G}'| - 1)\gamma_{\mathbb{G}'}(\mathbb{A})} \sqrt{\text{Var}[\|\mathbf{h}(\mathbf{y})\|_{\mathcal{Z}}]}}{\sqrt{|\mathbb{G}'|\mathbb{P}(\mathbf{x} \in \mathbb{A})}} \left(\mathcal{E}_{\theta}^r + \frac{\ln(n_{\text{iso}}/\delta)}{(d_{\text{iso}}N)^{\frac{\alpha}{1+2\alpha}}} \right). \quad (18b)$$

362 Where $\mathcal{E}_{\theta}^r = \|\mathbb{E}_{\mathbf{y}|\mathbf{x}} - \mathbf{E}_{\theta}\|_{\text{op}}$ denotes representation learning error (12), $r = \sum_{k=1}^{n_{\text{iso}}} d_k m_k$ defines the
 363 representation space dimension, with (d_k, m_k) denoting dimension and multiplicity of the group’s
 364 irreducible representation of type k , and $d_{\text{iso}} := \sum_{k=1}^{n_{\text{iso}}} d_k \geq n_{\text{iso}}$ (see Fig. 2).

367 *Proof.* \mathbb{G} -invariance of $P_{\mathbf{x}}$ and $P_{\mathbf{y}}$ allows us to control both bias (Thm. M.2) and variance (M.3) of
 368 $\hat{\mathbf{z}}_{\theta}$. A simple balancing of r yields the final bound on the error. \square

370 We conclude by highlighting key theoretical and practical implications of Thm. 5.1.

371 As point-wise guarantees are essential in several applications we present (18b), which offers learning
 372 bound when considering conditioning on measurable sets $\mathbb{A} \subseteq \mathcal{X}$, leading to the estimate $\mathbf{z}(\mathbb{A}) :=$
 373 $\mathbb{E}_{\mathbf{y}}[\mathbf{h}(\mathbf{y})|\mathbf{x} \in \mathbb{A}] \approx \hat{\mathbb{E}}_{\mathbf{x}}[\hat{\mathbf{z}}_{\theta}(\mathbf{x})]/\hat{\mathbb{E}}_{\mathbf{x}}[\mathbb{1}_{\mathbb{A}}(\mathbf{x})]$. In this context, symmetry exploitation becomes crucial in
 374 alleviating the bottlenecks of rare event estimation. To capture this effect, we introduce the *symmetry*
 375 *index* of the set \mathbb{A} w.r.t $P_{\mathbf{x}}$, which quantifies the degree of symmetry in \mathbb{A} and shapes (18b):

$$376 \quad \gamma_{\mathbb{G}}(\mathbb{A}) = \frac{1}{|\mathbb{G}| - 1} \sum_{g \in \mathbb{G} \setminus \{e\}} \frac{\mathbb{P}(\mathbf{x} \in \mathbb{A} \cap g \triangleright_{\mathcal{X}} \mathbb{A})}{\mathbb{P}(\mathbf{x} \in \mathbb{A})}, \quad \text{where } \gamma_{\mathbb{G}}(\mathbb{A}) \in [0, 1] \quad \forall \mathbb{A} \subseteq \mathcal{X}. \quad (19)$$

Observe that $\gamma_{\mathbb{G}}(\mathbb{A}) = 1$ if \mathbb{A} is \mathbb{G} -invariant (e.g., the vertical and horizontal reflection planes in Fig. 3), while $\gamma_{\mathbb{G}}(\mathbb{A}) = 0$ if \mathbb{A} is a \mathbb{G} -asymmetric set, that is, if $g \triangleright_{\mathcal{X}} \mathbb{A} \cap \mathbb{A} = \emptyset$ for all $g \in \mathbb{G}$ (e.g., any set disjoint from the reflection planes in Fig. 3). In particular, the effective rarity of $\mathbf{x} \in \mathbb{A}$ is captured by $\gamma_{\mathbb{G}'}(\mathbb{A})$, yielding a maximal gain of $|\mathbb{G}|\mathbb{P}[\mathbf{x} \in \mathbb{A}] \gg \mathbb{P}[\mathbf{x} \in \mathbb{A}]$ when \mathbb{A} is asymmetric.

Equivariant disentangled representations boost the effective sample size to $N \ll n_{\text{iso}}N \leq d_{\text{iso}}N \leq |\mathbb{G}|N$, providing not only the expected n_{iso} gain from disentanglement but also a remarkable $d_{\text{iso}} = \sum_{k=1}^{n_{\text{iso}}} d_k$ boost (see Fig. 2-right)—achieved by exploiting structural constraints of the singular spaces associated with the group’s irreducible representations present in the chosen representation space. Furthermore, note that when no symmetry prior exist, that is, when $\mathbb{G} = \{e\}$ and $|\mathbb{G}| = n_{\text{iso}} = d_{\text{iso}} = 1$, our framework recovers the baseline results of (Kostic et al., 2024b).

Lastly, note that the parameter α quantifies the problem regularity via the rate of decay of the operator’s singular values, $\sum_{i \in \mathbb{N}} \sigma_i^{1/\alpha} < \infty$, with special cases including finite-rank operators ($\alpha = \infty$), compact operators ($\alpha = 0$), trace class operators ($\alpha = 1$), and Hilbert-Schmidt operators ($\alpha = 1/2$, equivalent to $\kappa \in \mathcal{L}_{\mathcal{P}_{\mathbf{x}} \times \mathcal{P}_{\mathbf{y}}}^2(\mathcal{X} \times \mathcal{Y})$; see Sec. M). Hence, our results cover learning rates ranging from arbitrarily slow (as $\alpha \rightarrow 0$) to fast rates $[d_{\text{iso}}N]^{-1/2}$ as $\alpha \rightarrow \infty$.

6 EXPERIMENTS

We present three experiments evaluating our method in (i) approximating the conditional expectation operator and the use of the learned operator for (ii) \mathbb{G} -equivariant regression and (iii) symmetry-aware uncertainty quantification. For additional experiments, evidence, and details refer to Sec. G.

Conditional expectation operator learning This experiment *directly* quantifies the **Mean Squared Error (MSE)** of approximating $\mathbb{E}_{\mathbf{y}|\mathbf{x}}$, i.e., $\kappa_{\text{mse}} := \mathbb{E}_{\mathbf{x}}\mathbb{E}_{\mathbf{y}}(\kappa(\mathbf{x}, \mathbf{y}) - \kappa_{\theta}(\mathbf{x}, \mathbf{y}))^2$. To achieve this, we extend the **Conditional Gaussian Mixture Model (cGMM)** of Gilardi et al. (2002) to parametrically construct symmetric vector-valued random variables $\mathbf{x} \in \mathcal{X}$ and $\mathbf{y} \in \mathcal{Y}$ that satisfy the symmetry priors in (6) for arbitrary finite symmetry groups (see example in Fig. 3). **CGMM** possess an analytical **PMD** ratio κ , enabling direct estimation of κ_{mse} , usually impossible for real-world datasets.

The results in Fig. 4 compare our **eNCP** against **NCP** (Kostic et al., 2024b), baseline **Density Ratio Fitting (DRF)** (Tsai et al., 2020), and our **Invariant Density Ratio Fitting (iDRF)** adaptation. Note that baselines approximate κ as a single **NN**, $\kappa_{\theta}^{\text{drf}}: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^+$, via contrastive loss (5)—without enforcing separable form (13)—and hence *cannot* be used for regression or conditional probability estimation. We evaluate across **cGMMs** with diverse symmetry groups and (\mathbf{x}, \mathbf{y}) dimensionality. In all cases, **eNCP** exhibits superior performance and sample efficiency versus **iDRF** and symmetry-agnostic models (**NCP**, **DRF**), which struggle to approximate ratio \mathbb{G} -invariance (see Fig. 5). Results highlight the capacity of **eNCP** and **NCP** to approximate the conditional expectation operator and underscore the importance of exploiting symmetry for improved approximation.

\mathbb{G} -Equivariant regression To test our model’s potential for performing \mathbb{G} -equivariant regression, we address the robot’s **Center of Mass (CoM)** momenta regression task of (Ordoñez-Apraéz et al.). The goal is to predict a quadruped robot’s **CoM** linear $\mathbf{l} \in \mathbb{R}^3$ and angular momenta $\mathbf{k} \in \mathbb{R}^3$ given the noisy observations of the robot’s generalized positions $\mathbf{q} \in \mathbb{R}^{12}$ and velocity coordinates $\dot{\mathbf{q}} \in \mathbb{R}^{12}$, i.e., $[\mathbf{l}^T, \mathbf{k}^T]^T = h_{\text{CoM}}(\mathbf{q} + \epsilon_{\mathbf{q}}, \dot{\mathbf{q}} + \epsilon_{\dot{\mathbf{q}}})$ (see details in Subsec. G.2 and Fig. 6). We compare **eNCP** against **NCP** and two baselines—a standard **Multi-Layer Perceptron (MLP)** and an

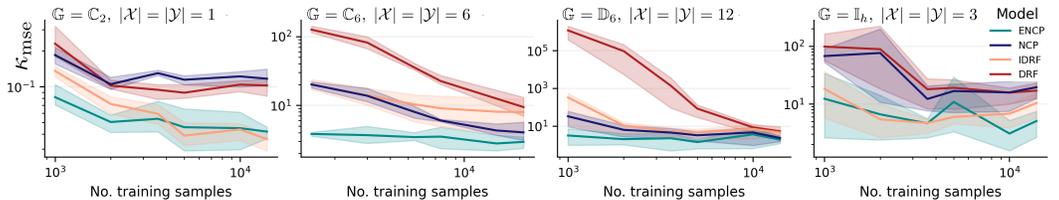


Figure 4: Sample efficiency plots comparing the test set **PMD MSE** $\kappa_{\text{mse}} := \mathbb{E}_{\mathbf{x}}\mathbb{E}_{\mathbf{y}}(\kappa(\mathbf{x}, \mathbf{y}) - \kappa_{\theta}(\mathbf{x}, \mathbf{y}))^2$ versus the number of training samples, in log scales. Each plot corresponds to a symmetric **cGMM** with distinct symmetry groups and (\mathbf{x}, \mathbf{y}) dimensionality. The tested groups are the cyclic groups \mathbb{C}_2 and \mathbb{C}_6 , the Dihedral group \mathbb{D}_6 (order 12), and the Icosahedral group \mathbb{I}_h (order 60).

432 **Equivariant MLP (eMLP)**—all with equivalent architectural footprint. Where **NCP** and **eNCP** are
 433 trained using (5) and (14), respectively, while **MLP** and **eMLP** are trained using standard **MSE**.
 434

435 The results in Fig. 1 demonstrate that our **eNCP** model outperforms all other baselines in both per-
 436 formance and sample complexity. Consistent with (Kostic et al., 2024b), the **NCP** model shows
 437 poorer sample complexity than **MLP** and **eMLP** due to its indirect approach to regression, via ap-
 438 proximation of $E_{y|x}$. However, by incorporating symmetry priors **eNCP** to mitigate this limitation.

439 **Symmetry aware uncertainty quantification** Finally, we demonstrate the practical impact of our
 440 approach on a core robotics problem: providing robust uncertainty quantification for unavailable yet
 441 crucial state observables for robot control and state estimation (Bledt et al., 2018; Maravgakis et al.,
 442 2023). Specifically, we use proprioceptive sensor readings to provide 90% confidence intervals for
 443 the robot’s **Ground Reaction Forces (GRF)** $\tau_{\text{grf}} \in \mathbb{R}^{12}$, the instantaneous work exerted or subtracted
 444 to the robot $U(\mathbf{q}, \dot{\mathbf{q}}, \boldsymbol{\tau}) \in \mathbb{R}$, and the kinetic energy $T(\mathbf{q}, \dot{\mathbf{q}}) \in \mathbb{R}$, while the robot traverses rough
 445 terrain (see Subsubsec. G.4.1). Reliable probabilistic estimates of these quantities are of crucial
 446 relevance for optimal control (Bledt et al., 2018), contact detection (Maravgakis et al., 2023), state
 447 estimation (Nisticò et al., 2025), and system identification (Gautier, 1997).

448 This task tests our model’s ability to learn conditional
 449 distributions from high-dimensional data, considering
 450 that for the **eNCP** and **NCP** models, quantile estima-
 451 tion is done by regressing the **Conditional Cumulative**
 452 **Distribution Function (CCDF)** for each dimension of
 453 $\mathbf{y} = [y_1, \dots]$ and then applying a linear search to ex-
 454 tract quantiles (see Figs. 8 and 13). This is achieved by
 455 discretizing the range of each y_i into N_b bins and es-
 456 timating $\mathbb{P}(y_i \in \mathbb{A}_{i,n} | \mathbf{x} = \cdot) := [E_{y|x} \mathbb{1}_{\mathbb{A}_{i,n}}](\cdot)$ for all
 457 $n \in [N_b]$ (see Sec. 5), where $\mathbb{A}_{i,n}$ consists of the first n
 458 bins. In practice, this means regressing $|\mathcal{Y}| \times N_b$ con-
 459 ditional probabilities corresponding to sets of varying
 460 sizes in a single forward pass. By contrast, the baseline
 461 **Conditional Quantile Regression (CQR)** (Feldman et al., 2023) and its equivariant adaptation **Equi-**
 462 **variant CQR (eCQR)** directly regress quantiles for a fixed coverage level (i.e., the probability that an
 463 event lies within the predicted confidence interval) and need retraining for different coverage values.

463 The results in Tab. 2, Fig. 1 (for U and T) and in Fig. 14 (for τ_{grf}) show **eNCP** as the only model
 464 capable of providing robust uncertainty quantification, as it is the only model with an empirical
 465 coverage on the test set close to the desired value, rendering other models unreliable for practical
 466 applications. This underscores **eNCP**’s potential for conditional probability estimation.

467 7 CONCLUSIONS

468 We introduce a novel framework for equivariant representation learning that enables estimation of
 469 equivariant regression and conditional probabilities with statistical learning guarantees. Our ap-
 470 proach builds on a recent contrastive representation learning method that approximates the spectral
 471 decomposition of the conditional expectation operator. By incorporating symmetry priors, we im-
 472 pose additional structural constraints that further decompose the conditional expectation operator
 473 and enhance the effective sample size. We demonstrate the benefits of our approach through both
 474 theoretical learning bounds and empirical experiments. Notably, we provide the first theoretical
 475 learning guarantees for equivariant regression using neural network features, thereby bridging spec-
 476 tral representation learning and geometric deep learning.
 477

478 **Limitations and future work** A potential limitation of our method is its reliance on fully specified
 479 symmetry priors. In practice, symmetries may only be partially known or subject to misspecification.
 480 As such, promising directions for future work include extending our framework to accommodate
 481 partial or uncertain symmetry information, exhibit robustness to symmetry violations, or statistically
 482 test for the presence of symmetry in data.
 483
 484
 485

	r-Coverage \uparrow	Coverage \uparrow
eNCP	99.5 \pm 0.1%	95.0 \pm 0.4%
NCP	99.5 \pm 0.0%	56.9 \pm 0.3%
eCQR	84.2 \pm 0.7%	6.7 \pm 1.2%
CQR	80.5 \pm 3.7%	8.5 \pm 0.9%

Table 2: Relaxed coverage, see (28), and Coverage, see (27), for the test-set confidence intervals in quadruped loco- motion uncertainty estimation of $\mathbf{y} = [\tau_{\text{grf}}^\top, U, T]^\top$. Target coverage is: 90%.

486 **Reproducibility statement** Detailed experimental setup is provided in [Sec. G](#), and all experiments
487 and plots are reproducible using the code available at the (anonymous) repository: [symm_rep_learn](#)
488

489 **LLM use** Large language models were used to polish writing and grammar. Furthermore, agen-
490 tic LLMs were used during coding to help via autocompletion, automatic generation of unit tests,
491 debugging, and code refactoring.
492

493 REFERENCES

494 Charles R Baker. Joint measures and cross-covariance operators. *Trans. Am. Math. Soc.*, 186:
495 273–273, 1973.
496

497 Han Bao, Yoshihiro Nagano, and Kento Nozawa. On the surrogate gap between contrastive and
498 supervised losses. In *International conference on machine learning*, pages 1585–1606. PMLR,
499 2022.

500 Arash Behboodi, Gabriele Cesa, and Taco S Cohen. A pac-bayesian generalization bound for equiv-
501 ariant networks. *Advances in Neural Information Processing Systems*, 35:5654–5668, 2022.
502

503 Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new
504 perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828,
505 2013.

506 Bernard Bercu, Bernard Delyon, and Emmanuel Rio. *Concentration Inequalities for Sums and*
507 *Martingales*. SpringerBriefs in Mathematics. Springer, 2015.
508

509 Alberto Bietti, Luca Venturi, and Joan Bruna. On the sample complexity of learning under geometric
510 stability. *Advances in neural information processing systems*, 34:18673–18684, 2021.

511 Gerardo Bledt, Patrick M Wensing, Sam Ingersoll, and Sangbae Kim. Contact model fusion for
512 event-based locomotion in unstructured terrains. In *2018 IEEE International Conference on*
513 *Robotics and Automation (ICRA)*, pages 4399–4406. IEEE, 2018.

514 Johannes Brandstetter, Rianne van den Berg, Max Welling, and Jayesh K Gupta. Clifford neural
515 layers for pde modeling. *arXiv preprint arXiv:2209.04934*, 2022.
516

517 Michael M Bronstein, Joan Bruna, Taco Cohen, and Petar Veličković. Geometric deep learning:
518 Grids, groups, graphs, geodesics, and gauges. *arXiv preprint arXiv:2104.13478*, 2021.

519 Élie Cartan. *La théorie des groupes finis et continus et l’analyse situs*. Number 42 in *Mémorial des*
520 *sciences mathématiques*. Gauthier-Villars, 1952.
521

522 Gabriele Cesa, Leon Lang, and Maurice Weiler. A program to build $e(n)$ -equivariant steerable cnns.
523 In *International conference on learning representations*, 2022.
524

525 Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for
526 contrastive learning of visual representations. In *International conference on machine learning*,
527 pages 1597–1607. PmLR, 2020.

528 Elijah Cole, Xuan Yang, Kimberly Wilber, Oisín Mac Aodha, and Serge Belongie. When does
529 contrastive visual representation learning work? In *Proceedings of the IEEE/CVF Conference on*
530 *Computer Vision and Pattern Recognition (CVPR)*, pages 14755–14764, June 2022.

531 Rumén Dangovski, Li Jing, Charlotte Loh, Seungwook Han, Akash Srivastava, Brian Cheung, Pulkit
532 Agrawal, and Marin Soljačić. Equivariant self-supervised learning: Encouraging equivariance in
533 representations. In *International Conference on Learning Representations*.

534 Rumén Dangovski, Li Jing, Charlotte Loh, Seungwook Han, Akash Srivastava, Brian Cheung, Pulkit
535 Agrawal, and Marin Soljačić. Equivariant contrastive learning. In *International Conference on*
536 *Learning Representations*, 2022.
537

538 Konstantin Donhauser, Mingqi Wu, and Fanny Yang. How rotational invariance of common kernels
539 prevents generalization in high dimensions. In *International Conference on Machine Learning*,
pages 2804–2814. PMLR, 2021.

- 540 Mildred S Dresselhaus, Gene Dresselhaus, and Ado Jorio. *Group theory: application to the physics*
541 *of condensed matter*. Springer Science & Business Media, 2007.
- 542
- 543 Carl Eckart and Gale Young. The approximation of one matrix by another of lower rank. *Psychome-*
544 *trika*, 1(3):211–218, 1936.
- 545 B Elesedy. *Symmetry and generalisation in machine learning*. PhD thesis, University of Oxford,
546 2023.
- 547
- 548 Bryn Elesedy. Provably strict generalisation benefit for invariance in kernel methods. In *Ad-*
549 *vances in Neural Information Processing Systems*, volume 34, pages 17273–17283. Curran
550 Associates, Inc., 2021. URL [https://proceedings.neurips.cc/paper_files/
551 paper/2021/file/8fe04df45a22b63156ebabbb064fcd5e-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/8fe04df45a22b63156ebabbb064fcd5e-Paper.pdf).
- 552 Bryn Elesedy and Sheheryar Zaidi. Provably strict generalisation benefit for equivariant mod-
553 els. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Confer-*
554 *ence on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages
555 2959–2969. PMLR, 18–24 Jul 2021. URL [https://proceedings.mlr.press/v139/
556 elesedy21a.html](https://proceedings.mlr.press/v139/elesedy21a.html).
- 557
- 558 Shai Feldman, Stephen Bates, and Yaniv Romano. Calibrated multiple-output quantile regression
559 with representation learning. *Journal of Machine Learning Research*, 24(24):1–48, 2023.
- 560 Kenji Fukumizu, Francis R Bach, and Michael I Jordan. Dimensionality reduction for supervised
561 learning with reproducing kernel hilbert spaces. *Journal of Machine Learning Research*, 5(Jan):
562 73–99, 2004.
- 563
- 564 Maxime Gautier. Dynamic identification of robots with power model. In *Proceedings of interna-*
565 *tional conference on robotics and automation*, volume 3, pages 1922–1927. IEEE, 1997.
- 566 Nicolas Gilardi, Samy Bengio, and Mikhail Kanevski. Conditional gaussian mixture models for
567 environmental risk mapping. In *Proceedings of the 12th IEEE workshop on neural networks for*
568 *signal processing*, pages 777–786. IEEE, 2002.
- 569 Martin Golubitsky, Ian Stewart, and David G Schaeffer. *Singularities and Groups in Bifurcation*
570 *Theory: Volume II*, volume 69. Springer Science & Business Media, 2012.
- 571
- 572 Sharut Gupta, Joshua Robinson, Derek Lim, Soledad Villar, and Stefanie Jegelka. Structuring repre-
573 sentation geometry with rotationally equivariant contrastive learning. In *The Twelfth International*
574 *Conference on Learning Representations*, 2023.
- 575
- 576 Jeff Z. HaoChen, Colin Wei, Adrien Gaidon, and Tengyu Ma. Provable guarantees for
577 self-supervised deep learning with spectral contrastive loss. In *Advances in Neural In-*
578 *formation Processing Systems*, volume 34, pages 5000–5011. Curran Associates, Inc.,
579 2021. URL [https://proceedings.neurips.cc/paper_files/paper/2021/
580 file/27debb435021eb68b3965290b5e24c49-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/27debb435021eb68b3965290b5e24c49-Paper.pdf).
- 581 Jeff Z HaoChen, Colin Wei, Ananya Kumar, and Tengyu Ma. Beyond separability: Analyzing the
582 linear transferability of contrastive representations to related subpopulations. *Advances in neural*
583 *information processing systems*, 35:26889–26902, 2022.
- 584
- 585 Olivier Henaff. Data-efficient image recognition with contrastive predictive coding. In *International*
586 *conference on machine learning*, pages 4182–4192. PMLR, 2020.
- 587 Irina Higgins, David Amos, David Pfau, Sebastien Racaniere, Loic Matthey, Danilo Rezende,
588 and Alexander Lerchner. Towards a definition of disentangled representations. *arXiv preprint*
589 *arXiv:1812.02230*, 2018.
- 590 Rafael Izbicki. *Machine Learning Beyond Point Predictions: Uncertainty Quantification*. imprint,
591 1st edition, 2025. ISBN 978-65-01-20272-3.
- 592
- 593 Rafael Izbicki and Ann B. Lee. Converting high-dimensional regression to high-dimensional condi-
tional density estimation. 2017.

- 594 Daniel D Johnson, Ayoub El Hanchi, and Chris J Maddison. Contrastive learning can find an optimal
595 basis for approximately view-invariant functions. *arXiv preprint arXiv:2210.01883*, 2022.
596
- 597 Karthik Kashinath, M Mustafa, Adrian Albert, JL Wu, C Jiang, Soheil Esmailzadeh, Kamyar Aziz-
598 zadenesheli, R Wang, Ashesh Chattopadhyay, A Singh, et al. Physics-informed machine learning:
599 case studies for weather and climate modelling. *Philosophical Transactions of the Royal Society*
600 *A*, 379(2194):20200093, 2021.
- 601 Hamza Keurti, Hsiao-Ru Pan, Michel Besserve, Benjamin F Grewe, and Bernhard Schölkopf. Ho-
602 momorphism autoencoder–learning group structured representations from observed transitions.
603 In *International Conference on Machine Learning*, pages 16190–16215. PMLR, 2023.
- 604 Anthony W. Knap. *Representation Theory of Semisimple Groups, An Overview Based on Examples*
605 *(PMS-36)*. Princeton University Press, Princeton, 1986.
606
- 607 Vladimir R Kostic, Pietro Novelli, Riccardo Grazi, Karim Lounici, and Massimiliano Pontil. Learn-
608 ing invariant representations of time-homogeneous stochastic dynamical systems. In *The Twelfth*
609 *International Conference on Learning Representations*, 2024a.
- 610 Vladimir R Kostic, Gregoire Pacreau, Giacomo Turri, Pietro Novelli, Karim Lounici, and Massim-
611 iliano Pontil. Neural conditional probability for uncertainty quantification. In *The Thirty-eighth*
612 *Annual Conference on Neural Information Processing Systems*, 2024b.
613
- 614 Nikola Kovachki, Zongyi Li, Burigede Liu, Kamyar Aizzadenesheli, Kaushik Bhattacharya, An-
615 drew Stuart, and Anima Anandkumar. Neural operator: Learning maps between function spaces
616 with applications to pdes. *Journal of Machine Learning Research*, 24(89):1–97, 2023.
- 617 Lucas Kovar, Michael Gleicher, and Frédéric Pighin. Motion graphs. *ACM Trans. Graph.*, 21(3):
618 473–482, July 2002. URL <https://doi.org/10.1145/566654.566605>.
- 619 Phuc H Le-Khac, Graham Healy, and Alan F Smeaton. Contrastive representation learning: A
620 framework and review. *Ieee Access*, 8:193907–193934, 2020.
621
- 622 Lilang Lin, Jiahang Zhang, and Jiaying Liu. Mutual information driven equivariant contrastive
623 learning for 3d action representation learning. *IEEE Transactions on Image Processing*, 2024.
624
- 625 Zicheng Liu, Steven J Gortler, and Michael F Cohen. Hierarchical spacetime control. In *Proceedings*
626 *of the 21st annual conference on Computer graphics and interactive techniques*, pages 35–42,
627 1994.
- 628 George W. Mackey. Harmonic analysis as the exploitation of symmetry—a historical survey. *Bulletin*
629 *(New Series) of the American Mathematical Society*, 3(1.P1):543 – 698, 1980.
- 630 Michael Maravgakis, Despina-Ekaterini Argiropoulos, Stylianos Piperakis, and Panos Trahanias.
631 Probabilistic contact state estimation for legged robots using inertial information. In *2023 IEEE*
632 *International Conference on Robotics and Automation (ICRA)*, pages 12163–12169. IEEE, 2023.
633
- 634 Giovanni Luca Marchetti, Gustaf Tegnér, Anastasiia Varava, and Danica Kragic. Equivariant repre-
635 sentation learning via class-pose decomposition. In *International Conference on Artificial Intel-*
636 *ligence and Statistics*, pages 4745–4756. PMLR, 2023.
- 637 Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Learning with invariances in random
638 features and kernel models. In *Conference on Learning Theory*, pages 3351–3418. PMLR, 2021.
639
- 640 Niloy J. Mitra, Mark Pauly, Michael Wand, and Duygu Ceylan. Symmetry in 3d geometry:
641 Extraction and applications. *Computer Graphics Forum*, 32(6):1–23, February 2013. URL
642 <http://dx.doi.org/10.1111/cgf.12010>.
- 643 Youssef Mroueh, Stephen Voinea, and Tomaso A Poggio. Learning with group invariant features: A
644 kernel perspective. *Advances in neural information processing systems*, 28, 2015.
645
- 646 Thomas Nagler and Claudia Czado. Evading the curse of dimensionality in nonparametric density
647 estimation with simplified vine copulas. *Journal of Multivariate Analysis*, 151:69–89, October
2016. URL <http://dx.doi.org/10.1016/j.jmva.2016.07.003>.

- 648 Ylenia Nisticò, João Carlos Virgolino Soares, Lorenzo Amatucci, Geoff Fink, and Claudio Sem-
649 ini. Muse: A real-time multi-sensor state estimator for quadruped robots. *IEEE Robotics and*
650 *Automation Letters*, 2025.
- 651 Peter J Olver. *Applications of Lie groups to differential equations*, volume 107. Springer Science &
652 Business Media, 1993.
- 654 Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predic-
655 tive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- 657 Daniel Ordoñez-Apraez, Vladimir Kostic, Giulio Turrisi, Pietro Novelli, Carlos Mastalli, Claudio
658 Semini, and Massimiliano Pontil. Dynamics harmonic analysis of robotic systems: Application in
659 data-driven koopman modelling. In *6th Annual Learning for Dynamics & Control Conference*,
660 pages 1318–1329. PMLR, 2024.
- 661 Daniel Ordoñez-Apraez, Giulio Turrisi, Vladimir Kostic, Mario Martin, Antonio Agudo, Francesc
662 Moreno-Noguer, Massimiliano Pontil, Claudio Semini, and Carlos Mastalli. Morphological
663 symmetries in robotics. *The International Journal of Robotics Research*. doi: 10.1177/
664 02783649241282422. URL <https://doi.org/10.1177/02783649241282422>.
- 666 Sherjil Ozair, Corey Lynch, Yoshua Bengio, Aaron van den Oord, Sergey Levine, and Pierre
667 Sermanet. Wasserstein dependency measure for representation learning. In H. Wal-
668 lach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors,
669 *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.,
670 2019. URL [https://proceedings.neurips.cc/paper_files/paper/2019/
671 file/f9209b7866c9f69823201c1732cc8645-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/f9209b7866c9f69823201c1732cc8645-Paper.pdf).
- 672 Dipan Pal, Ashwin Kannan, Gautam Arakalgud, and Marios Savvides. Max-margin invariant fea-
673 tures from transformed unlabelled data. In *Advances in Neural Information Processing Systems*,
674 volume 30. Curran Associates, Inc., 2017.
- 675 Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier
676 Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-learn:
677 Machine learning in python. *the Journal of machine Learning research*, 12:2825–2830, 2011.
- 679 J Jon Ryu, Xiangxiang Xu, HS Erol, Yuheng Bu, Lizhong Zheng, and Gregory W Wornell. Operator
680 svd with neural networks via nested low-rank approximation. *arXiv preprint arXiv:2402.03655*,
681 2024.
- 682 Anastasiya Salova, Jeffrey Emenheiser, Adam Rupe, James P Crutchfield, and Raissa M D’Souza.
683 Koopman operator and its approximations for systems with symmetries. *Chaos: An Interdisci-
684 plinary Journal of Nonlinear Science*, 29(9), 2019.
- 686 DAVID W. Scott. Feasibility of multivariate density estimates. *Biometrika*, 78(1):197–205, 1991.
687 URL <http://dx.doi.org/10.1093/biomet/78.1.197>.
- 688 Parikshit Shah and Venkat Chandrasekaran. Group symmetry and covariance regularization. In *2012*
689 *46th Annual Conference on Information Sciences and Systems (CISS)*, pages 1–6. IEEE, 2012.
- 691 Ralph C. Smith. *Uncertainty Quantification: Theory, Implementation, and Applications*. Society for
692 Industrial and Applied Mathematics, January 2013. URL [http://dx.doi.org/10.1137/
693 1.9781611973228](http://dx.doi.org/10.1137/1.9781611973228).
- 694 Le Song, Jonathan Huang, Alex Smola, and Kenji Fukumizu. Hilbert space embeddings of condi-
695 tional distributions with applications to dynamical systems. In *Proceedings of the 26th Annual*
696 *International Conference on Machine Learning*, pages 961–968, 2009.
- 698 Masashi Sugiyama, Taiji Suzuki, and Takafumi Kanamori. *Density ratio estimation in machine*
699 *learning*. Cambridge University Press, 2012.
- 700 Behrooz Tahmasebi and Stefanie Jegelka. The exact sample complexity gain from invariances for
701 kernel regression. *Advances in Neural Information Processing Systems*, 36, 2023.

- 702 Christopher Tosh, Akshay Krishnamurthy, and Daniel Hsu. Contrastive learning, multi-view redun-
703 dancy, and linear models. In *Proceedings of the 32nd International Conference on Algorithmic*
704 *Learning Theory*, volume 132 of *Proceedings of Machine Learning Research*, pages 1179–1206.
705 PMLR, 16–19 Mar 2021.
- 706 Yao-Hung Hubert Tsai, Han Zhao, Makoto Yamada, Louis-Philippe Morency, and Russ R Salakhut-
707 dinov. Neural methods for point-wise dependency estimation. In *Advances in Neural Information*
708 *Processing Systems*, volume 33, pages 62–72. Curran Associates, Inc., 2020.
- 709 Yao-Hung Hubert Tsai, Martin Q Ma, Muqiao Yang, Han Zhao, Louis-Philippe Morency, and Rus-
710 lan Salakhutdinov. Self-supervised representation learning with relative predictive coding. In
711 *International Conference on Learning Representations*, 2021.
- 712 Giulio Turrisi, Valerio Modugno, Lorenzo Amatucci, Dimitrios Kanoulas, and Claudio Semini. On
713 the benefits of gpu sample-based stochastic predictive controllers for legged locomotion. In *2024*
714 *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 13757–
715 13764, 2024. doi: 10.1109/IROS58592.2024.10801698.
- 716 Elise van der Pol, Daniel Worrall, Herke van Hoof, Frans Oliehoek, and Max Welling. Mdp homo-
717 morphic networks: Group symmetries in reinforcement learning. In *Advances in Neural Informa-*
718 *tion Processing Systems*, volume 33, pages 4199–4210. Curran Associates, Inc., 2020.
- 719 Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices.
720 *arXiv:1011.3027*, 2011.
- 721 Hiroki Waida, Yuichiro Wada, Léo Andéol, Takumi Nakagawa, Yuhui Zhang, and Takafumi
722 Kanamori. Towards understanding the mechanism of contrastive learning via similarity struc-
723 ture: A theoretical analysis. In *Joint European Conference on Machine Learning and Knowledge*
724 *Discovery in Databases*, pages 709–727. Springer, 2023.
- 725 R Wang. Incorporating symmetry into deep dynamics models for improved generalization. In
726 *International Conference on Learning Representations (ICLR)*, 2021.
- 727 Rui Wang, Robin Walters, and Rose Yu. Incorporating symmetry into deep dynamics models for
728 improved generalization. *arXiv preprint arXiv:2002.03061*, 2020.
- 729 Rui Wang, Robin Walters, and Rose Yu. Approximately equivariant networks for imperfectly sym-
730 metric dynamics. In *International Conference on Machine Learning*, pages 23078–23091. PMLR,
731 2022a.
- 732 Xin Wang, Hong Chen, Si’ao Tang, Zihao Wu, and Wenwu Zhu. Disentangled representation learn-
733 ing. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024a.
- 734 Yifei Wang, Kaiwen Hu, Sharut Gupta, Ziyu Ye, Yisen Wang, and Stefanie Jegelka. Understanding
735 the role of equivariance in self-supervised learning. In *The Thirty-eighth Annual Conference on*
736 *Neural Information Processing Systems*, 2024b.
- 737 Ziyu Wang, Yucen Luo, Yueru Li, Jun Zhu, and Bernhard Schölkopf. Spectral representation learn-
738 ing for conditional moment models. *arXiv preprint arXiv:2210.16525*, 2022b.
- 739 Larry Wasserman. *All of Nonparametric Statistics*. Springer Texts in Statistics. Springer, New York,
740 NY, 1 edition, May 2007.
- 741 Joachim Weidmann. *Linear operators in Hilbert spaces*, volume 68. Springer Science & Business
742 Media, 2012.
- 743 Maurice Weiler, Patrick Forré, Erik Verlinde, and Max Welling. *Equivari-*
744 *ant and Coordinate Independent Convolutional Networks*. World Scien-
745 tific, 2023. URL [https://maurice-weiler.gitlab.io/cnn_book/](https://maurice-weiler.gitlab.io/cnn_book/EquivariantAndCoordinateIndependentCNNs.pdf)
746 [EquivariantAndCoordinateIndependentCNNs.pdf](https://maurice-weiler.gitlab.io/cnn_book/EquivariantAndCoordinateIndependentCNNs.pdf).
- 747 Jianke Yang, Nima Dehmamy, Robin Walters, and Rose Yu. Latent space symmetry discovery. In
748 *International Conference on Machine Learning*, 2023.

756 Thomas Yerxa, Jenelle Feather, Eero Simoncelli, and SueYeon Chung. Contrastive-equivariant self-
757 supervised learning improves alignment with primate visual area it. *Advances in neural informa-*
758 *tion processing systems*, 37:96045–96070, 2024.

759 Xupeng Zhu, Dian Wang, Ondrej Biza, Guanang Su, Robin Walters, and Robert Platt. Sample
760 Efficient Grasp Learning Using Equivariant Models. In *Proceedings of Robotics: Science and*
761 *Systems*, New York City, NY, USA, June 2022.

762 James Y Zou, Daniel J Hsu, David C Parkes, and Ryan P Adams. Contrastive learning using spectral
763 methods. In *Advances in Neural Information Processing Systems*, volume 26. Curran Associates,
764 Inc., 2013.

765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

Appendix

Table of Contents

A Symbols and notation	17
B Acronyms	18
C Related work	19
D Symmetry constraints on conditional expectations	20
E \mathbb{G}-Equivariant bilinear NN architecture	21
F Symmetry aware orthonormalization of disentangled representations	21
G Experimental setup	23
H Conditional probability modeling via the conditional expectation operator	32
I Background on group and representation theory	33
J Representation theory of symmetric function spaces	39
K \mathbb{G}-equivariant linear integral operators	42
L Relevant \mathbb{G}-equivariant operators in probability theory	47
M Statistical Learning Theory	49

The supplementary material is organized as follows.

- [Sec. A](#) summarizes the notations used, while [Sec. B](#) provides a glossary.
- [Sec. C](#) offers a detailed discussion of related work on contrastive learning, equivariant representations, and statistical learning theory with symmetry priors.
- [Sec. D](#)-[Sec. G](#) detail our methodological contributions. In particular, [Sec. E](#) and [Sec. F](#) explain how our method leverages equivariant neural networks, and [Sec. G](#) complements [Sec. 6](#) with additional experimental details and studies.
- We provide self-contained sections with unified notation covering essential preliminaries: group theory ([Sec. I](#)), representation theory in function spaces ([Sec. J](#)), equivariant linear operators ([Sec. K](#)), and the symmetries of covariance operators central to our work ([Sec. L](#)).
- Finally, [Sec. M](#) presents our theoretical contributions, summarized in [Theorem Thm. 5.1](#). We first establish approximation error bounds using operator theory and equivariant representations, then combine group theory with concentration inequalities to derive estimation error bounds that fully expose the benefits of symmetry priors.

864	A SYMBOLS AND NOTATION	
865		
866		Numbers and Arrays
867	x	A scalar, or scalar function $x(\cdot)$
868	\mathbf{x}	A vector, or vector-valued function $\mathbf{x}(\cdot)$
869	$\mathbf{x}_1 \oplus \mathbf{x}_2$	Direct sum (stacking) of vectors, such that $\mathbf{x}_1 \oplus \mathbf{x}_2 := \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}$
870	\mathbf{K}	A matrix
871	$\mathbf{A} \oplus \mathbf{B}$	Direct sum of matrices, such that $\mathbf{A} \oplus \mathbf{B} := \begin{bmatrix} \mathbf{A} & \mathbf{O} \\ \mathbf{O} & \mathbf{B} \end{bmatrix}$
872	\mathbf{K}	A linear operator
873	\mathbf{I}	Identity matrix
874	$\delta_{i,j}$	The Kronecker function, equal to 1 when $i = j$, and 0 when $i \neq j$
875		Sets, Vector Spaces, and Function Spaces
876	$\mathcal{X}, \mathcal{Z}, \mathcal{H}, \mathcal{F}$	A vector or Hilbert space
877	$\bar{\mathcal{X}}, \bar{\mathcal{Z}}, \bar{\mathcal{V}}$	An irreducible \mathbb{G} -stable space (Thm. I.6)
878	\mathbb{R}, \mathbb{C}	The set of real and complex numbers
879	$\mathcal{X} \oplus \mathcal{Y}$	Direct sum of vector spaces \mathcal{X} and \mathcal{Y} , such that if $\mathbf{x} \in \mathcal{X}$ and $\mathbf{y} \in \mathcal{Y}$, then $\mathbf{x} \oplus \mathbf{y} \in \mathcal{X} \oplus \mathcal{Y}$
880		
881	$\mathcal{L}_x^2 := \mathcal{L}_{P_x}^2(\mathcal{X}, \mathbb{R})$	The space of square-integrable functions on \mathcal{X} . That is $\{f \mid \int_{\mathcal{X}} f(\mathbf{x}) ^2 P_x(d\mathbf{x}) < \infty, f: \mathcal{X} \rightarrow \mathbb{R}\}$
882		
883	$\langle f, f' \rangle_{P_x}$	Inner product $\langle f, f' \rangle_{P_x} := \int_{\mathcal{X}} f(\mathbf{x}) f'(\mathbf{x}) P_x(d\mathbf{x})$
884		
885		Group and Representation theory
886	\mathbb{G}	A symmetry group
887	g, g_1, g_a	A symmetry group element
888	$g \triangleright \mathbf{x}$	The (left) group action of g on \mathbf{x} defined by $g \triangleright \mathbf{x} := \rho_{\mathcal{X}}(g)\mathbf{x}$
889	$\rho_{\mathcal{X}}$	A representation of the group \mathbb{G} on the vector space \mathcal{X} , defined for a chosen basis of \mathcal{X}
890	$\bar{\rho}_k$	An irreducible representation (Thm. I.8) of the group \mathbb{G}
891	$d_k := \bar{\rho}_k $	Dimensionality of the irreducible representation $\bar{\rho}_k$ (see Fig. 2)
892	m_k	Multiplicity of the irreducible representation $\bar{\rho}_k$ in a given larger representation (see Fig. 2)
893		
894	n_{iso}	Number of distinct irreps present in a given larger representation.
895	$\rho_{\mathcal{X}}(g)$	Representation of the group element g on the vector space \mathcal{X}
896	$\rho_{\mathcal{X}} \oplus \rho_{\mathcal{Y}}$	Direct sum of group representations, such that $\rho_{\mathcal{X}}(g) \oplus \rho_{\mathcal{Y}}(g) := \begin{bmatrix} \rho_{\mathcal{X}}(g) & \\ & \rho_{\mathcal{Y}}(g) \end{bmatrix}$
897		
898	$\mathbb{G}\mathbf{x}$	The group orbit of \mathbf{x} , defined as $\mathbb{G}\mathbf{x} := \{g \triangleright \mathbf{x} \mid g \in \mathbb{G}\}$
899	$\gamma_{\mathbb{G}'}(A)$	The symmetry index of a set $A \subseteq \mathcal{X}$ w.r.t. probability distribution on \mathcal{X} and group elements $\mathbb{G}' \subseteq \mathbb{G}$
900		
901	$\mathbb{G}_a \times \mathbb{G}_b$	Direct product of groups \mathbb{G}_a and \mathbb{G}_b
902	$\mathbf{U}(\mathcal{X})$	Unitary group on the vector space \mathcal{X}
903	$\mathbf{GL}(\mathcal{X})$	General Linear group on the vector space \mathcal{X} , a.k.a the space of invertible matrices in $\mathbb{R}^{ \mathcal{X} \times \mathcal{X} }$
904		
905	\mathbb{C}_n	Cyclic group of order n
906	\mathbb{K}_4	Klein four-group
907		Probability Theory
908	$\mathbf{x} \sim \mathbb{P}(\mathbf{x})$	Random vector $\mathbf{x} \in \mathcal{X}$ has distribution $\mathbb{P}(\mathbf{x})$
909	$P_{\mathbf{x}}$	A probability measure on the space \mathcal{X}
910	$\mathbb{E}_{\mathbf{x}}[f(\mathbf{x})]$	Expectation of $f(\mathbf{x})$ with respect to $P_{\mathbf{x}}$
911	$\text{Cov}(f(\mathbf{x}))$	Variance of $f(\mathbf{x})$ with respect to $P_{\mathbf{x}}$, define as $\mathbb{E}_{\mathbf{x}}(f(\mathbf{x}) - \mathbb{E}_{\mathbf{x}}f(\mathbf{x}))^2$
912	$\text{Cov}(f(\mathbf{x}), h(\mathbf{y}))$	Covariance of $f(\mathbf{x})$ and $h(\mathbf{y})$ with respect to the joint distribution $P_{\mathbf{xy}}$, defined as $\mathbb{E}_{\mathbf{xy}}(f(\mathbf{x}) - \mathbb{E}_{\mathbf{x}}f(\mathbf{x}))(h(\mathbf{y}) - \mathbb{E}_{\mathbf{y}}h(\mathbf{y}))$
913		
914	$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$	Gaussian distribution over \mathbf{x} with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$
915		
916		
917		

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

B ACRONYMS

CCDF Conditional Cumulative Distribution Function.

cGMM Conditional Gaussian Mixture Model (Gilardi et al., 2002): A parametric model for benchmark conditional density estimation datasets. Generates random variables \mathbf{x} and \mathbf{y} of arbitrary dimensions with varying mutual information. Enables analytical computation of the **PMD** density ratio (see Sec. 2), unavailable in real-world datasets, allowing direct quantification of approximation error for the conditional expectation operator $E_{\mathbf{y}|\mathbf{x}}$ and its **PMD** density ratio .

cLoRa Contrastive Low-Rank loss from (Kostic et al., 2024b; Ryu et al., 2024) for operator and representation learning. Used in density-ratio fitting (Sugiyama et al., 2012), representation learning (Wang et al., 2022b; HaoChen et al., 2021), and mutual information estimation (Tsai et al., 2020) .

CoM Center of Mass.

CQR Conditional Quantile Regression (Feldman et al., 2023): A multivariate neural network approach for regressing upper and lower quantiles at a specified miscoverage level α using pinball loss. Confidence intervals are typically calibrated post-training via conformal prediction, but calibration is omitted here for fair model comparison .

DoF degree of freedom.

DRF Density Ratio Fitting (Tsai et al., 2020): A density ratio **NN** architecture that parameterizes the approximated **PMD** $\kappa_{\theta} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_+$ as a single **NN**. Consequently, this model cannot be used for downstream conditional probability estimation and regression—it is limited to estimating the mutual information between \mathbf{x} and \mathbf{y} (Tsai et al., 2020) .

eCQR Version of eCQR where the upper and lower parametric quantile functions are parameterized by \mathbb{G} -equivariant NNs .

eMLP Equivariant MLP.

eNCP Equivariant Neural Conditional Probability: Our proposed model integrating the symmetry priors Eq. (6) into the NCP deep representation learning algorithm .

GDL Geometric Deep Learning: A field of machine learning that incorporates geometric priors into deep learning models (Bronstein et al., 2021).

GRF Ground Reaction Forces.

iDRF Invariant Density Ratio Fitting: This is a \mathbb{G} -invariant adaptation of the DRF model (Tsai et al., 2020) that parameterizes the approximated **PMD** κ_{θ} as a \mathbb{G} -invariant **NN** .

MLP Multi-Layer Perceptron.

MSE Mean Squared Error.

NCP Neural Conditional Probability: A deep representation learning framework (Kostic et al., 2024b) for conditional probability estimation and regression with statistical guarantees via operator theory (Baker, 1973). This framework is symmetry-agnostic .

NN Neural Network.

PMD Pointwise Mutual Dependency Tsai et al. (2020): A pointwise dependency measure between random variables \mathbf{x} and \mathbf{y} , defined as $\kappa(\mathbf{x}, \mathbf{y}) = \frac{dP_{\mathbf{x}\mathbf{y}}(\mathbf{x}, \mathbf{y})}{d(P_{\mathbf{x}}(\mathbf{x}) \times P_{\mathbf{y}}(\mathbf{y}))} = \exp(\text{MI}(\mathbf{x}, \mathbf{y}))$.

SVD Singular Value Decomposition.

C RELATED WORK

C.1 CONTRASTIVE REPRESENTATION LEARNING

Contrastive representation learning obtains high-dimensional representations from unlabeled data by contrasting positive and negative sample pairs via a noise contrastive loss (similar to Eq. (5)) (Le-Khac et al., 2020; Waida et al., 2023; Bao et al., 2022). Most works in this field aim to learn representations in a self-supervised fashion that transfer well to downstream classification tasks Johnson et al. (2022); Cole et al. (2022); Tosh et al. (2021); Oord et al. (2018); Tsai et al. (2021); HaoChen et al. (2021). In contrast, our approach targets representations that effectively transfer to (equivariant) regression and uncertainty quantification, as in Kostic et al. (2024a). Given a dataset $\mathbb{D} = \{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^N$ from a target (stochastic) function $\mathbf{y} = \mathbf{f}(\mathbf{x})$, we treat positive pairs as drawn from the joint distribution $(\mathbf{x}, \mathbf{y}) \sim \mathbb{P}(\mathbf{x}, \mathbf{y})$ and negative pairs as drawn from the product of the marginals $(\mathbf{x}, \mathbf{y}) \sim \mathbb{P}(\mathbf{x})\mathbb{P}(\mathbf{y})$. In this setting, our contrastive loss aims to learn representations that approximate the PMD ratio $\kappa(\mathbf{x}, \mathbf{y}) = \frac{\mathbb{P}(\mathbf{x}, \mathbf{y})}{\mathbb{P}(\mathbf{x})\mathbb{P}(\mathbf{y})}$, Kostic et al. (2024a) or equivalently, the pointwise mutual information $\ln(\kappa(\mathbf{x}, \mathbf{y}))$ (Oord et al., 2018; Lin et al., 2024; Henaff, 2020; Ozair et al., 2019). Crucially, our work is the first study this problem when there is prior knowledge of the invariance of κ under the action of a compact symmetry group, which occurs in most applications of GDL.

Linear transferability The goal of contrastive representation learning is to acquire representations that transfer to diverse downstream inference tasks Bengio et al. (2013); Waida et al. (2023). While empirical studies demonstrate that contrastive learning can outperform supervised methods Cole et al. (2022); Oord et al. (2018); Henaff (2020), theoretical works aim to establish *linear separability/transferability* guarantees (HaoChen et al., 2022)⁴. That is, showing that linear functionals of the (frozen) learned representations suffice for regression/classification inference.

In the context of **classification**, Waida et al. (2023); Bao et al. (2022); Johnson et al. (2022) show that contrastive learning losses serve as surrogates for standard supervised classification losses (e.g., the cross-entropy). Where the gap between the surrogate and supervised loss diminishes with the number of negative samples Bao et al. (2022) (N^2 for the loss in Eq. (16)). To provide these transferability guarantees, these work assume $\mathcal{X} = \mathcal{Y}$, so that the PMD ratio κ becomes a positive definite kernel. Consequently, kernel method guarantees can be transferred to the classification task, even when the representations are parameterized by NNs (Johnson et al., 2022; Bao et al., 2022; HaoChen et al., 2022).

Considerably fewer works have studied contrastive representation learning in the context of downstream **regression** tasks (Yerxa et al., 2024; Kostic et al., 2024a). Crucially, Kostic et al. (2024a) show that a contrastive learning loss serves as surrogate to the MSE regression loss (A summary of this method appears in Sec. 2 and in Tab. 3). While, to the best of our knowledge, (Yerxa et al., 2024) is the only work empirically studying contrastive learning for regression in the presence of symmetries.

Task	$\mathbf{f}(\mathbf{x}) := \mathbb{E}_{\mathbf{y}}[\mathbf{y} \mathbf{x}=\mathbf{x}] \approx \hat{\mathbf{f}}_{\boldsymbol{\theta}}(\mathbf{x})$	$\mathbb{P}[\mathbf{y} \in \mathbb{B} \mathbf{x} \in \mathbb{A}] \approx \hat{\mathbb{P}}_{\boldsymbol{\theta}}[\mathbf{y} \in \mathbb{B} \mathbf{x} \in \mathbb{A}]$
Estimate	$\hat{\mathbb{E}}_{\mathbf{y}}[\mathbf{y}] + \boldsymbol{\phi}_{\boldsymbol{\theta}}(\mathbf{x})^{\top} \mathbf{E}_{\boldsymbol{\theta}} \hat{\mathbb{E}}_{\mathbf{y}}[\boldsymbol{\psi}_{\boldsymbol{\theta}}(\mathbf{y}) \otimes \mathbf{y}]$	$\hat{\mathbb{E}}_{\mathbf{y}}[\mathbb{1}_{\mathbb{B}}] + \frac{\hat{\mathbb{E}}_{\mathbf{x}}[\mathbb{1}_{\mathbb{A}}(\mathbf{x}) \otimes \boldsymbol{\phi}_{\boldsymbol{\theta}}(\mathbf{x})]^{\top} \mathbf{E}_{\boldsymbol{\theta}} \hat{\mathbb{E}}_{\mathbf{y}}[\mathbb{1}_{\mathbb{B}}(\mathbf{y}) \otimes \boldsymbol{\psi}_{\boldsymbol{\theta}}(\mathbf{y})]}{\hat{\mathbb{E}}_{\mathbf{x}}[\mathbb{1}_{\mathbb{A}}(\mathbf{x})]}$
Guarantees	$\ \mathbf{f} - \hat{\mathbf{f}}_{\boldsymbol{\theta}}\ _{\mathcal{L}_{\mathbf{x}}^2} \lesssim \sqrt{\text{var}[\ \mathbf{y}\]} \left(\mathcal{E}_{\boldsymbol{\theta}}^r + \frac{\ln(1/\delta)}{N \frac{\alpha}{1+2\alpha}} \right)$	$ \mathbb{P} - \hat{\mathbb{P}}_{\boldsymbol{\theta}} \lesssim \sqrt{\frac{\mathbb{P}[\mathbf{y} \in \mathbb{B}]}{\mathbb{P}[\mathbf{x} \in \mathbb{A}]}} \left(\mathcal{E}_{\boldsymbol{\theta}}^r + \frac{\ln(1/\delta)}{N \frac{\alpha}{1+2\alpha}} \right)$

Table 3: Statistical learning guarantees of NCP (Kostic et al., 2024a) for regression and conditional probability estimation. The bounds are shaped by the quality of the learned representations $\mathcal{E}_{\boldsymbol{\theta}}^r = \|\mathbb{E}_{\mathbf{y}|\mathbf{x}} - \mathbf{E}_{\boldsymbol{\theta}}\|_{\text{op}} \leq \sqrt{\mathcal{L}_{\gamma}(\boldsymbol{\theta}) - \mathcal{L}_{\gamma}(\star)}$ (see (5)), the sample size N , and the decay rate of $\mathbb{E}_{\mathbf{y}|\mathbf{x}}$ singular-values $\alpha > 0$, which quantifies the difficulty of the problem.

C.2 EQUIVARIANT REPRESENTATION LEARNING

Equivariant contrastive representation learning Dangovski et al. (2022); Wang et al. (2024b) aims to learn representations that are equivariant—instead of invariant—to data transformations. For example, Marchetti et al. (2023); Gupta et al. (2023); Lin et al. (2024) provide empirical evidence that

⁴Also referred to as linear evaluation protocol (Chen et al., 2020)

representations of 3D scenes, images, and human body poses that are equivariant to translations, rotations, or reflections yield improved performance in *classification* tasks. Additionally, Yerxa et al. (2024) show that rotation- and reflection-aware image representations enhance the *regression* of neural responses in the macaque inferior temporal cortex, while also providing theoretical justification that such equivariant representations mirror the known structure of animal visual perception. By introducing these transformations via data-augmentation of the training set, these methods inherently enforce symmetries in the data distributions, which are the fundamental priors assumed in Sec. 3.

Disentangled representations In equivariant representation learning, disentangled representations have been extensively studied (Wang et al., 2024a). Initially, (Bengio et al., 2013) defined disentanglement as decomposing representations into components that capture distinct, independently varying factors. Later, using group theory, Higgins et al. (2018) formalized that a representation is disentangled if its space decomposes into orthogonal subspaces reflecting a symmetry group decomposition, with each subspace influenced exclusively by one quotient group (see Thm. I.18). As discussed in Sec. I, this aligns with the isotypic decomposition of a Hilbert space (Mackey, 1980): $\mathcal{H} = \bigoplus_{k=1}^{\perp} \mathcal{H}^{(k)}$ —known in dynamical systems (Golubitsky et al., 2012)—when the symmetry group decomposes as $\mathbb{G} = \prod_{k=1}^{n_{\text{iso}}} \mathbb{G}^{(k)}$. Orthogonality between subspaces follows from Schur’s orthogonality relations via Cartan’s and Peter-Weyl’s theorems (Cartan, 1952). This symmetric structure is the cause of the architectural constraints imposed in the eNCP architecture Fig. 2.

Several empirical works have explored disentanglement in representation learning. For instance, Keurti et al. (2023) proposed an autoencoder-based method to learn disentangled equivariant representations by using loss regularization to enforce latent space equivariance and sparsity for separating latent group actions. Unlike our approach, their method does not assume prior knowledge of the symmetry group and relies entirely on loss regularization rather than architectural constraints. Similarly, works such as (see e.g. Yang et al., 2023; Dangovski et al.) have investigated various symmetry priors in latent space by examining the emergence of disentangled structures and enforcing algebraic constraints. Notably, in fields like molecular dynamics, physics, computer graphics, and robotics, symmetry priors are intrinsic to the task or system (Ordoñez-Apaez et al.; Lin et al., 2024; Marchetti et al., 2023), making them natural assumptions. In a similar spirit to our work, Marchetti et al. (2023) leverage the known \mathbb{SO}_3 symmetries of the 3D world to learn \mathbb{SO}_3 -disentangled equivariant representations using contrastive learning, thereby demonstrating the empirical advantages of symmetry-aware, disentangled representations for object classification.

C.3 SYMMETRY-AWARE STATISTICAL LEARNING THEORY

Existing literature on symmetry-aware learning focuses on group-invariant regression via kernel methods (Mroueh et al., 2015; Tahmasebi and Jegelka, 2023; Elesedy, 2021; Elesedy and Zaidi, 2021; Elesedy, 2023; Pal et al., 2017; Mei et al., 2021; Bietti et al., 2021; Donhauser et al., 2021). Most of these methods cannot be directly transferred to modern GDL architectures.

In contrast, in deep learning and GDL, while many works offer a group-theoretical analysis and empirical evidence of the benefits of incorporating symmetry priors (Kashinath et al., 2021; Wang et al., 2020; 2022a; Brandstetter et al., 2022), none, to our knowledge, provide statistical learning guarantees for equivariant **regression**. The only exception is (Behboodi et al., 2022), which derives generalization bounds for a MLP architecture in the context of n -class **classification** task using a margin loss. In contrast, our work provides statistical learning guarantees for equivariant **regression** and symmetry-aware **uncertainty quantification**, both as corollaries of Thm. 5.1.

D SYMMETRY CONSTRAINTS ON CONDITIONAL EXPECTATIONS

Under the assumed symmetry priors in (6) the conditional expectation of \mathbf{y} is a \mathbb{G} -equivariant function/map. This property is depicted in Fig. 3-center and proved in the following proposition.

Proposition D.1 (\mathbb{G} -equivariant conditional expectations). *Let $\mathbf{x} \in \mathcal{X}$ and $\mathbf{y} \in \mathcal{Y}$ be two vector valued random variables satisfying the symmetry priors of Eq. (6). Then, the conditional expectation*

of \mathbf{y} given \mathbf{x} is \mathbb{G} -equivariant, since, for every $g \in \mathbb{G}$, $\mathbf{x} \in \mathcal{X}$,

$$\begin{aligned} \mathbb{E}[\mathbf{y}|\mathbf{x} = g \triangleright_{\mathcal{X}} \mathbf{x}] &= g \triangleright_{\mathcal{Y}} \mathbb{E}[\mathbf{y}|\mathbf{x} = \mathbf{x}] \\ &= \int_{\mathcal{Y}} g \triangleright_{\mathcal{Y}} \mathbf{y} P_{\mathbf{y}|\mathbf{x}}(d\mathbf{y}|\mathbf{x}) \\ &= \int_{\mathcal{Y}} \mathbf{y} P_{\mathbf{y}|\mathbf{x}}(g^{-1} \triangleright_{\mathcal{Y}} d\mathbf{y}|\mathbf{x}) \\ &= \int_{\mathcal{Y}} \mathbf{y} P_{\mathbf{y}|\mathbf{x}}(d\mathbf{y}|g \triangleright_{\mathcal{X}} \mathbf{x}) \quad (\text{by Eq. (6)}) \\ &= \mathbb{E}[\mathbf{y}|\mathbf{x} = g \triangleright_{\mathcal{X}} \mathbf{x}]. \end{aligned}$$

E \mathbb{G} -EQUIVARIANT BILINEAR NN ARCHITECTURE

This section outlines how to construct a \mathbb{G} -equivariant disentangled representation for the random variables \mathbf{x} and \mathbf{y} using **any** type of \mathbb{G} -equivariant NN architecture backbone, such as MLP, CNNs, Transformers, and others.

Let $\mathbf{f}_{\theta} : \mathcal{X} \mapsto \mathbb{R}^r$ and $\mathbf{h}_{\theta} : \mathcal{Y} \mapsto \mathbb{R}^r$ be two \mathbb{G} -equivariant NNs, whose outputs will be interpreted as the basis functions of the truncated symmetric function spaces $\mathcal{F}_{\mathbf{x}} \subset \mathcal{L}_{\mathbf{x}}^2$ and $\mathcal{F}_{\mathbf{y}} \subset \mathcal{L}_{\mathbf{y}}^2$. Assume, the group representations on $\mathcal{F}_{\mathbf{x}}$ and $\mathcal{F}_{\mathbf{y}}$ are constructed from multiplicities of the group’s regular representation, $\rho_{\mathcal{F}_{\mathbf{x}}} = \bigoplus_{n=1}^{r/|\mathbb{G}|} \rho_{\text{reg}}$ and $\rho_{\mathcal{F}_{\mathbf{y}}} = \bigoplus_{n=1}^{r/|\mathbb{G}|} \rho_{\text{reg}}$ —as done usually in practice (Weiler et al., 2023). Since for (most) finite groups, the decomposition of ρ_{reg} into *irreps* is known or can be computed, we have access to the analytical change of basis $\mathbf{Q}_{\mathbf{x}} : \mathcal{F}_{\mathbf{x}} \mapsto \mathcal{F}_{\mathbf{x}}$ and $\mathbf{Q}_{\mathbf{y}} : \mathcal{F}_{\mathbf{y}} \mapsto \mathcal{F}_{\mathbf{y}}$ to transition to the isotypic basis. Consequently, we can directly parameterize the representations of the random variables in disentangled form as:

$$\phi_{\theta}(\cdot) = \mathbf{Q}_{\mathbf{x}}^{\top}(\mathbf{f}_{\theta}(\cdot) - \mathbb{E}_{\mathbf{x}}[\mathbf{f}_{\theta}(\mathbf{x})]), \quad \psi_{\theta}(\cdot) = \mathbf{Q}_{\mathbf{y}}^{\top}(\mathbf{h}_{\theta}(\cdot) - \mathbb{E}_{\mathbf{y}}[\mathbf{h}_{\theta}(\mathbf{y})]). \quad (20)$$

Given that during training these representations are not orthogonal, the truncated operator is parameterized as the trainable \mathbb{G} -equivariant matrix $\mathbf{E}_{\theta} = \bigoplus_k^{n_{\text{iso}}} \mathbf{E}_{\theta}^{(k)} = \bigoplus_k^{n_{\text{iso}}} \sum_{b \in \mathbb{B}} \Theta_b^{(k)} \otimes \Psi(b)$, where \mathbb{B} is the endomorphism’s basis of the irreducible representation $\bar{\rho}_k$, $\Psi : \mathbb{B} \rightarrow \mathbb{R}^{d_k \times d_k}$ is a mapping from the endomorphism’s basis to the endomorphism space, and $\Theta_b^{(k)} \in \mathbb{R}^{m_k^y \times m_k^x}$ represent the block’s trainable parameters for each $b \in \mathbb{B}$, see details in Thm. I.13.

Note that after training, the SVD of the learned operator can be computed by exploiting the constraints imposed by the operator’s \mathbb{G} -equivariance (see Thm. K.4 and Fig. 2).

F SYMMETRY AWARE ORTHONORMALIZATION OF DISENTANGLED REPRESENTATIONS

This section covers how to compute unbiased empirical estimates of the orthonormalization and centering regularization terms in Eq. (14) in the presence of symmetries.

Let $\mathbb{E}_{\mathbf{y}|\mathbf{x}} : \mathcal{L}_{\mathbf{y}}^2 \mapsto \mathcal{L}_{\mathbf{x}}^2$ be the conditional expectation operator and $\mathbb{E}_{\theta} : \mathcal{F}_{\mathbf{y}} \mapsto \mathcal{F}_{\mathbf{x}}$ be its r -rank approximation on the spaces $\mathcal{F}_{\mathbf{x}} = \text{span}(\{\phi_i\}_{i=1}^r)$ and $\mathcal{F}_{\mathbf{y}} = \text{span}(\{\psi_i\}_{i=1}^r)$. Denote by $\kappa(\mathbf{x}, \mathbf{y}) := \frac{P_{\mathbf{x}\mathbf{y}}(\mathbf{x}, \mathbf{y})}{P_{\mathbf{x}}(\mathbf{x})P_{\mathbf{y}}(\mathbf{y})}$ and $\kappa_{\theta}(\mathbf{x}, \mathbf{y}) := \sum_{i,j=1}^r [\mathbf{E}_{\theta}]_{i,j} \phi_i(\mathbf{x}) \psi_j(\mathbf{y}) = \phi(\mathbf{x})^{\top} \mathbf{E}_{\theta} \psi(\mathbf{y})$ the kernel functions of the full and restricted operator, respectively. Then we have that:

$$\|\mathbb{E}_{\mathbf{y}|\mathbf{x}} - \mathbb{E}_{\theta}\|_{\text{HS}}^2 \leq -2\langle \mathbb{E}_{\mathbf{y}|\mathbf{x}}, \mathbb{E}_{\theta} \rangle_{\text{HS}} + \|\mathbb{E}_{\theta}\|_{\text{HS}}^2, \quad (21a)$$

$$\begin{aligned} &\leq -2 \int_{\mathcal{X} \times \mathcal{Y}} \kappa(\mathbf{x}, \mathbf{y}) \kappa_{\theta}(\mathbf{x}, \mathbf{y}) P_{\mathbf{x}}(d\mathbf{x}) P_{\mathbf{y}}(d\mathbf{y}) + \int_{\mathcal{X} \times \mathcal{Y}} \kappa_{\theta}(\mathbf{x}, \mathbf{y})^2 P_{\mathbf{x}}(d\mathbf{x}) P_{\mathbf{y}}(d\mathbf{y}) \\ &\leq -2 \int_{\mathcal{X} \times \mathcal{Y}} \kappa_{\theta}(\mathbf{x}, \mathbf{y}) P_{\mathbf{x}\mathbf{y}}(d\mathbf{x}, d\mathbf{y}) + \int_{\mathcal{X} \times \mathcal{Y}} \kappa_{\theta}(\mathbf{x}, \mathbf{y})^2 P_{\mathbf{x}}(d\mathbf{x}) P_{\mathbf{y}}(d\mathbf{y}) \\ &\leq -2\mathbb{E}_{\mathbf{x}\mathbf{y}} \kappa_{\theta}(\mathbf{x}, \mathbf{y}) + \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{y}} \kappa_{\theta}(\mathbf{x}, \mathbf{y})^2. \end{aligned} \quad (21b)$$

For the purpose of our representation learning problem, we consider the scenario in which the chosen basis sets include the constant function, and all other basis functions are centered by construction.

That is, $\mathbb{I}_{\mathcal{F}_x} = \{\mathbb{1}_{P_x}\} \cup \{\phi_i \mid \langle \phi_i, \mathbb{1}_{P_x} \rangle_{\mathbf{x}} = 0\}_{i=1}^r$ and $\mathbb{I}_{\mathcal{F}_y} = \{\mathbb{1}_{P_y}\} \cup \{\psi_i \mid \langle \psi_i, \mathbb{1}_{P_y} \rangle_{\mathbf{y}} = 0\}_{i=1}^r$. This results in the $(r + 1)$ -dimensional matrices:

$$\mathbf{V}_x := \begin{bmatrix} 1 & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_x \end{bmatrix}, \quad \mathbf{V}_y := \begin{bmatrix} 1 & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_y \end{bmatrix}, \quad (22)$$

where $\mathbf{C}_x = \text{Cov}(\phi(\mathbf{x}), \phi(\mathbf{x})) \in \mathbb{R}^{r \times r}$, $\mathbf{C}_y = \text{Cov}(\psi(\mathbf{y}), \psi(\mathbf{y})) \in \mathbb{R}^{r \times r}$ denote the matrix forms of the truncated covariance operators $\mathbf{C}_x : \mathcal{F}_x \mapsto \mathcal{F}_x$ and $\mathbf{C}_y : \mathcal{F}_y \mapsto \mathcal{F}_y$ (see [Thm. L.5](#)), respectively. Then the orthonormality regularization of [Eq. \(5\)](#) becomes:

$$\|\mathbf{V}_x - \mathbf{I}\|_{\mathbb{F}}^2 = \|\mathbf{C}_x - \mathbf{I}_r\|_{\mathbb{F}}^2 + 2\|\mathbb{E}_{P_x} \phi(\mathbf{x})\|^2 \quad \|\mathbf{V}_y - \mathbf{I}\|_{\mathbb{F}}^2 = \|\mathbf{C}_y - \mathbf{I}_r\|_{\mathbb{F}}^2 + 2\|\mathbb{E}_{P_y} \psi(\mathbf{y})\|^2. \quad (23)$$

Since $\|\mathbf{C}_x\|_{\mathbb{F}}^2 = \text{tr}(\mathbf{C}_x^2)$ involves products of covariance matrices, we compute its empirical value using unbiased estimators. For generality, we present the unbiased estimator for the cross-covariance.

Unbiased estimation of Frobenious norm of cross-covariance operators Since $\|\mathbf{C}_{xy}\|_{\mathbb{F}}^2 = \text{tr}(\mathbf{C}_{xy}^2)$ involves products of covariance matrices, we obtain unbiased estimates from finite samples by computing the metric using two independent sampling sets from P_{xy} , by:

$$\begin{aligned} \|\mathbf{C}_{xy}\|_{\mathbb{F}}^2 &= \text{tr}(\mathbf{C}_{xy}^2) = \sum_{i=1}^r [\mathbf{C}_{xy}^2]_{i,i} = \sum_{i=1}^r \sum_{j=1}^r [\mathbf{C}_{xy}]_{i,j} [\mathbf{C}_{xy}]_{j,i} \\ &= \sum_{i=1}^r \sum_{j=1}^r \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim P_{xy}} [\phi_{c,i}(\mathbf{x}) \psi_{c,j}(\mathbf{y})] \mathbb{E}_{(\mathbf{x}', \mathbf{y}') \sim P_{xy}} [\phi_{c,j}(\mathbf{x}') \psi_{c,i}(\mathbf{y}')] \\ &= \mathbb{E}_{(\mathbf{x}, \mathbf{y}, \mathbf{x}', \mathbf{y}') \sim P_{xy}} \left[\sum_{i=1}^r \phi_{c,i}(\mathbf{x}) \psi_{c,i}(\mathbf{y}') \sum_{j=1}^r \phi_{c,j}(\mathbf{x}') \psi_{c,j}(\mathbf{y}) \right] \\ &= \mathbb{E}_{(\mathbf{x}, \mathbf{y}, \mathbf{x}', \mathbf{y}') \sim P_{xy}} [(\phi_c(\mathbf{x})^\top \psi_c(\mathbf{y}')) (\phi_c(\mathbf{x}')^\top \psi_c(\mathbf{y}))] \\ &\approx \frac{1}{N^2} \sum_{n=1}^N \sum_{m=1}^N (\phi_c(\mathbf{x}_n)^\top \psi_c(\mathbf{y}'_m)) (\phi_c(\mathbf{x}'_m)^\top \psi_c(\mathbf{y}_n)), \end{aligned} \quad (24)$$

where $\phi_c(\mathbf{x}) = \phi(\mathbf{x}) - \mathbb{E}_{P_x} \phi(\mathbf{x})$ denotes the centered basis functions, and $((\mathbf{x}, \mathbf{y}), (\mathbf{x}', \mathbf{y}')) \sim P_{xy}$ indicates two independent sampling sets from P_{xy} used for the unbiased estimation of $\|\mathbf{C}_x\|_{\mathbb{F}}^2$. The final equation then provides the unbiased empirical estimator computed on a dataset $\mathbb{D} = \{(\mathbf{x}_n, \mathbf{y}_n) \sim P_{xy}\}_{n=1}^N$ and any random permutation of it, denoted as $\mathbb{D}' = \{(\mathbf{x}'_n, \mathbf{y}'_n) \sim P_{xy}\}_{n=1}^N$.

F.1 UNBIASED ESTIMATION OF ORTHONORMAL REGULARIZATION

The regularization term for optimizing the loss [\(5\)](#) involves encouraging the basis sets to be orthonormal. The metric quantifying the orthogonality of the basis sets is defined by:

$$\begin{aligned} \|\mathbf{V}_x - \mathbf{I}\|_{\mathbb{F}}^2 &= \|\mathbf{C}_x - \mathbf{I}_r\|_{\mathbb{F}}^2 + 2\|\mathbb{E}_{P_x} \phi(\mathbf{x})\|^2 = \text{tr}(\mathbf{C}_x^2) - 2\text{tr}(\mathbf{C}_x) + r + 2\|\mathbb{E}_{P_x} \phi(\mathbf{x})\|^2, \\ \|\mathbf{V}_y - \mathbf{I}\|_{\mathbb{F}}^2 &= \|\mathbf{C}_y - \mathbf{I}_r\|_{\mathbb{F}}^2 + 2\|\mathbb{E}_{P_y} \psi(\mathbf{y})\|^2 = \text{tr}(\mathbf{C}_y^2) - 2\text{tr}(\mathbf{C}_y) + r + 2\|\mathbb{E}_{P_y} \psi(\mathbf{y})\|^2. \end{aligned} \quad (25)$$

Hence given a dataset of samples $\mathbb{D} = \{(\mathbf{x}_n, \mathbf{y}_n) \sim P_{xy}\}_{n=1}^N$, and any random permutation of the dataset order $\mathbb{D}' = \{(\mathbf{x}'_n, \mathbf{y}'_n) \sim P_{xy}\}_{n=1}^N$ we can derive unbiased empirical estimates of [\(25\)](#) as:

$$\begin{aligned} \|\mathbf{V}_x - \mathbf{I}\|_{\mathbb{F}}^2 &\approx \widehat{\mathbb{E}}_{(\mathbf{x}, \mathbf{x}') \sim P_x} [(\phi_c(\mathbf{x})^\top \phi_c(\mathbf{x}'))^2] - 2\widehat{\mathbb{E}}_{P_x} [\phi_c(\mathbf{x})^\top \phi_c(\mathbf{x})] + r + 2\|\widehat{\mathbb{E}}_{P_x} \phi(\mathbf{x})\|^2 \\ &\approx \frac{1}{N^2} \sum_{n=1}^N \sum_{m=1}^N (\phi_c(\mathbf{x}_n)^\top \phi_c(\mathbf{x}'_m))^2 - 2\frac{1}{N} \sum_{n=1}^N \phi_c(\mathbf{x}_n)^\top \phi_c(\mathbf{x}_n) + r + 2\left\| \frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \right\|^2, \\ \|\mathbf{V}_y - \mathbf{I}\|_{\mathbb{F}}^2 &\approx \widehat{\mathbb{E}}_{(\mathbf{y}, \mathbf{y}') \sim P_y} [(\psi_c(\mathbf{y})^\top \psi_c(\mathbf{y}'))^2] - 2\widehat{\mathbb{E}}_{P_y} [\psi_c(\mathbf{y})^\top \psi_c(\mathbf{y})] + r + 2\|\widehat{\mathbb{E}}_{P_y} \psi(\mathbf{y})\|^2 \\ &\approx \frac{1}{N^2} \sum_{n=1}^N \sum_{m=1}^N (\psi_c(\mathbf{y}_n)^\top \psi_c(\mathbf{y}'_m))^2 - 2\frac{1}{N} \sum_{n=1}^N \psi_c(\mathbf{y}_n)^\top \psi_c(\mathbf{y}_n) + r + 2\left\| \frac{1}{N} \sum_{n=1}^N \psi(\mathbf{y}_n) \right\|^2. \end{aligned} \quad (26)$$

F.2 ORTHONORMAL REGULARIZATION OF SYMMETRIC HILBERT SPACES

Since the covariance operators $C_{\mathbf{x}} : \mathcal{L}_{\mathbf{x}}^2 \mapsto \mathcal{L}_{\mathbf{x}}^2$ and $C_{\mathbf{y}} : \mathcal{L}_{\mathbf{y}}^2 \mapsto \mathcal{L}_{\mathbf{y}}^2$ are \mathbb{G} -equivariant (see [Thm. L.6](#)), their matrix representations in the isotypic basis are constrained to be block-diagonal, with each block being composed of irrep endomorphisms ([Thm. I.13](#)). Hence (25) becomes:

$$\begin{aligned}
\|\mathbf{V}_{\mathbf{x}} - \mathbf{I}\|_{\mathbb{F}}^2 &= \|C_{\mathbf{x}} - \mathbf{I}_r\|_{\mathbb{F}}^2 + 2\|\mathbb{E}_{P_{\mathbf{x}}}\phi(\mathbf{x})\|^2 \\
&= \|\bigoplus_{k=1}^{n_{\text{iso}}} C_{\mathbf{x}}^{(k)} - \mathbf{I}_r\|_{\mathbb{F}}^2 + 2\|\mathbb{E}_{P_{\mathbf{x}}}\phi^{\text{inv}}(\mathbf{x})\|^2, \\
&= \sum_{k=1}^{n_{\text{iso}}} \|C_{\mathbf{x}}^{(k)} - \mathbf{I}_r^{(k)}\|_{\mathbb{F}}^2 + 2\|\mathbb{E}_{P_{\mathbf{x}}}\phi^{\text{inv}}(\mathbf{x})\|^2 \\
&= \sum_{k=1}^{n_{\text{iso}}} \left(\|C_{\mathbf{x}}^{(k)}\|_{\mathbb{F}}^2 - 2\text{tr}(C_{\mathbf{x}}^{(k)}) + r_k \right) + 2\|\mathbb{E}_{P_{\mathbf{x}}}\phi(\mathbf{x})\|^2 \\
&= \sum_{k=1}^{n_{\text{iso}}} \left(\left\| \sum_{b \in \mathbb{B}} \Theta_b^{(k)} \otimes \Psi_k(b) \right\|_{\mathbb{F}}^2 - 2\text{tr} \left(\sum_{b \in \mathbb{B}} \Theta_b^{(k)} \otimes \Psi_k(b) \right) + r_k \right) + 2\|\mathbb{E}_{P_{\mathbf{x}}}\phi(\mathbf{x})\|^2, \quad \text{by (42)} \\
&= 2\|\mathbb{E}_{P_{\mathbf{x}}}\phi(\mathbf{x})\|^2 + \sum_{k=1}^{n_{\text{iso}}} \left\| \sum_{b \in \mathbb{B}} \Theta_b^{(k)} \otimes \Psi_k(b) \right\|_{\mathbb{F}}^2 - 2 \sum_{b \in \mathbb{B}} \text{tr}(\Theta_b^{(k)}) \text{tr}(\Psi_k(b)) + r_k \\
&= 2\|\mathbb{E}_{P_{\mathbf{x}}}\phi(\mathbf{x})\|^2 + \sum_{k=1}^{n_{\text{iso}}} \left\| \sum_{b \in \mathbb{B}} \Theta_b^{(k)} \otimes \Psi_k(b) \right\|_{\mathbb{F}}^2 - 2d_k \sum_{b \in \mathbb{B}} \text{tr}(\Theta_b^{(k)}) + r_k \quad \text{by (39)} \\
&= 2\|\mathbb{E}_{P_{\mathbf{x}}}\phi(\mathbf{x})\|^2 + \sum_{k=1}^{n_{\text{iso}}} \sum_{b \in \mathbb{B}} \|\Theta_b^{(k)}\|_{\mathbb{F}}^2 \|\Psi_k(b)\|_{\mathbb{F}}^2 - 2d_k \sum_{b \in \mathbb{B}} \text{tr}(\Theta_b^{(k)}) + r_k \quad \text{by (39)} \\
&= 2\|\mathbb{E}_{P_{\mathbf{x}}}\phi(\mathbf{x})\|^2 + \sum_{k=1}^{n_{\text{iso}}} d_k \sum_{b \in \mathbb{B}} \|\Theta_b^{(k)}\|_{\mathbb{F}}^2 - 2d_k \sum_{b \in \mathbb{B}} \text{tr}(\Theta_b^{(k)}) + r_k, \quad \text{s.t } \Psi_k(b)\Psi_k(b)^\top = \mathbf{I}_{d_k} \\
&= 2\|\mathbb{E}_{P_{\mathbf{x}}}\phi(\mathbf{x})\|^2 + d_k \sum_{k=1}^{n_{\text{iso}}} \sum_{b \in \mathbb{B}} \left(\|\Theta_b^{(k)}\|_{\mathbb{F}}^2 - 2\text{tr}(\Theta_b^{(k)}) \right) + r_k,
\end{aligned}$$

Where the set of matrices $\{\Theta_b^{(k)} \in \mathbb{R}^{m_k \times m_k}\}_{b \in \mathbb{B}}$ define the free degrees of freedom of the covariance operator (see [Thm. I.13](#)), and $\{\Psi_k(b)\}_{b \in \mathbb{B}}$ denote the basis of endomorphisms of the irreducible representation $\bar{\rho}_k$ (see (39)). Crucially, $\|\Theta_b^{(k)}\|_{\mathbb{F}}^2$ features an unbiased U-statistic estimator as in (24).

G EXPERIMENTAL SETUP

In this section we provide details on the experimental setup. We first describe general design choices and hyperparameters and then provide details for each experiment.

Sample efficiency experiments For both the conditional expectation operator approximation and the \mathbb{G} -equivariant regression experiments, we evaluate model performance by measuring sample efficiency/complexity. To do so, we partition the dataset $\mathbb{D} = \{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^N$ into training, validation, and testing splits in proportions of 70%, 15%, and 15%, respectively. With fixed validation and testing sets, we iteratively train the models on increasing portions of the training set and report the test performance for each size.

For each training set size, we select the model checkpoint with the best validation loss to compute the test performance. Thus, these experiments quantify the generalization error (or true risk) and its evolution as a function of the training set size.

NNS architectures and hyperparameters To compare our equivariant representation learning framework with other contrastive and supervised methods, all (inference) models share a similar fixed architectural footprint. For the baseline models, the only hyperparameter tuned is the learning rate, whereas for the **NCP** and **eNCP** models we additionally tune the regularization weight γ in [Eqs. \(5\)](#) and [\(14\)](#). Further details for each experiment are provided in the corresponding sections below.

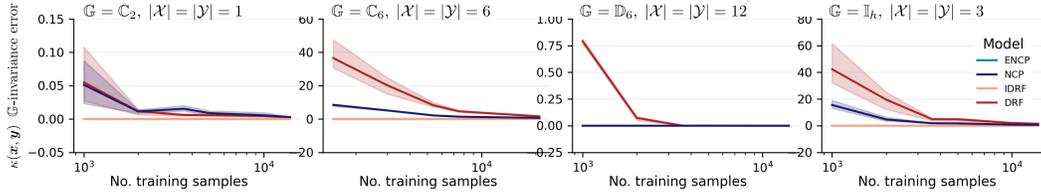


Figure 5: Sample efficiency plots comparing test set regression mean-square-error of the density ratio $\kappa(\mathbf{x}, \mathbf{y}) = P_{\mathbf{x}\mathbf{y}}(\mathbf{x}, \mathbf{y})/P_{\mathbf{x}}(\mathbf{x})P_{\mathbf{y}}(\mathbf{y})$ (log-scale) vs. the number of samples in the training set (log-scale). Each plot represents a different symmetric GMM with varying symmetry groups \mathbb{G} and dimensionalities of the random variables $|\mathcal{X}|$ and $|\mathcal{Y}|$. The groups tested are the cyclic groups \mathbb{C}_2 and \mathbb{C}_6 , the Dihedral group \mathbb{D}_6 of order 12 and the Icosahedral group \mathbb{I}_h of order 60.

Code reproducibility All experiments, plots and examples are provided in the open-access repository and python package [symm_rep_learn](#).

G.1 CONDITIONAL EXPECTATION OPERATOR APPROXIMATION

In this experiment, we extend the conditional Gaussian Mixture Model (GMM) proposed by Gilardi et al. (2002) to parametrically construct symmetric random variables taking values in arbitrary data spaces \mathcal{X} and \mathcal{Y} and with arbitrary finite symmetry groups \mathbb{G} . The GMM is defined by

$$\mathbf{z} := \mathbf{x} \oplus \mathbf{y} \sim \sum_{g \in \mathbb{G}} \sum_{c=1}^{n_g} \mathcal{N}(\boldsymbol{\rho}_z(g) \mu_{z,c}, \boldsymbol{\rho}_z(g) \Sigma_{z,c} \boldsymbol{\rho}_z(g)^\top),$$

where $\boldsymbol{\rho}_z(g) := \boldsymbol{\rho}_x(g) \oplus \boldsymbol{\rho}_y(g)$ are arbitrary group representations of \mathbb{G} and n_g is the number of unique Gaussians with randomly sampled means $\mu_z := \mu_x \oplus \mu_y$ and block-diagonal covariances $\Sigma_z := \Sigma_x \oplus \Sigma_y$. Since every Gaussian appears in group orbits, this symmetric GMM has \mathbb{G} -invariant marginal distributions and an analytical expression for the conditional expectation operator kernel $\kappa(\mathbf{x}, \mathbf{y}) = P_{\mathbf{x}\mathbf{y}}(\mathbf{x}, \mathbf{y})/P_{\mathbf{x}}(\mathbf{x})P_{\mathbf{y}}(\mathbf{y})$ (see 2D example in Fig. 3). Consequently, we can directly estimate the approximation of the conditional expectation operator (Eq. (5)) as the mean squared error between the true and learned density ratios, i.e., $\kappa_{\text{mse}} := \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{y}} \|\kappa(\mathbf{x}, \mathbf{y}) - \kappa_\theta(\mathbf{x}, \mathbf{y})\|^2$.

To the best of our knowledge, this is the first synthetic experiment that directly estimates the truncation error of the conditional expectation operator in an inference task-agnostic setting, serving as a benchmark for future work.

Fig. 4 compares sample efficiency using κ_{mse} , while Fig. 5 shows the error in the \mathbb{G} -invariant of the learned κ ratio versus sample size, highlighting that symmetry-aware methods encode this property as an architectural constraint, ensuring a strictly \mathbb{G} -invariant learned ratio.

G.2 \mathbb{G} -EQUIVARIANT REGRESSION OF ROBOT’S CoM MOMENTA

In this experiment, we evaluate the quality of the learned representations using the contrastive loss Eqs. (5) and (14) alongside supervised learning baselines trained with the standard MSE loss. The task is a \mathbb{G} -equivariant benchmark in robotics presented in (Ordoñez-Apraez et al.), with the goal of predicting a quadruped robot’s CoM linear $\mathbf{l} \in \mathbb{R}^3$ and angular momenta $\mathbf{k} \in \mathbb{R}^3$ from noisy observations of the robot’s generalized positions $\mathbf{q} \in \mathbb{R}^{12}$ and velocity coordinates $\dot{\mathbf{q}} \in \mathbb{R}^{12}$. Consequently, the random variables are defined as $\mathbf{x} = \mathbf{q} + \epsilon_q \oplus \dot{\mathbf{q}} + \epsilon_{\dot{q}}$ and $\mathbf{y} = \mathbf{l} \oplus \mathbf{k}$, where $\epsilon_q \in \mathbb{R}^{12}$ and $\epsilon_{\dot{q}} \in \mathbb{R}^{12}$ are independent Gaussian noise terms that model sensor noise. The function computing the CoM momenta from these proprioceptive observations is highly non-linear and \mathbb{G} -equivariant whenever \mathbb{G} is a morphological symmetry group of the robot (see Fig. 6 and (Ordoñez-Apraez et al.) for details).

The robot considered is the quadruped robot Solo (Fig. 6-right), which possesses a symmetry group of order 8: $\mathbb{G} = \mathbb{K}_4 \times \mathbb{C}_2$, as depicted in this animation showing 8 symmetric robot configurations along with their corresponding linear and angular momenta vectors.

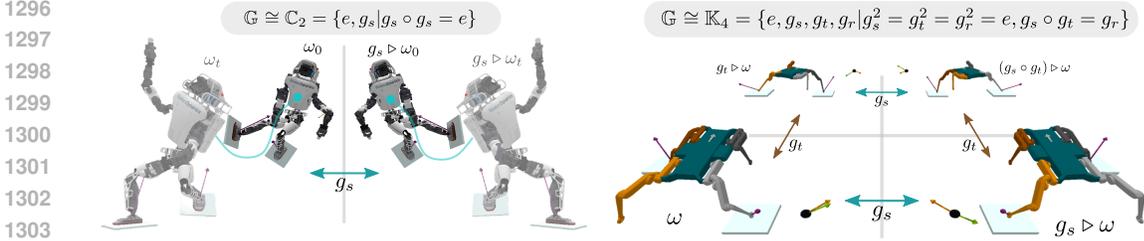


Figure 6: Example of morphological finite symmetry in robotics. **Left:** A humanoid robot with the reflectional symmetry group $\mathbb{G} \equiv \mathbb{C}_2$. **Right:** The quadruped robot Solo with the symmetry group $\mathbb{G} = \mathbb{K}_4 \times \mathbb{C}_2$ (only \mathbb{K}_4 is shown for clarity). The robot’s center of mass linear $l \in \mathbb{R}^3$ and angular $k \in \mathbb{R}^3$ momentum are depicted as orange and green vectors, respectively, for each symmetric configuration. Images adapted from [Ordoñez-Apaez et al.](#) with author approval.

NN architectures We configure all models under consideration (eNCP, NCP, eMLP, and MLP) to have an inference-time NN architecture with a similar footprint. In particular, the encoder network for \mathbf{x} in NCP and eNCP is designed similarly to the NN used in MLP/eMLP. The idea is to test how a model with the same capacity performs on the downstream task of classification when trained using either the representation learning loss or a supervised learning loss. The backbone of all architectures is a standard multilayer perceptron consisting of three hidden layers, each with 512 units, followed by a final hidden layer containing 128 units. This final layer encodes the feature vector r for the NCP and eNCP models. Crucially, since \mathbb{G} -equivariance enforces weight sharing in the NN architecture, the encoder NN for eNCP and eMLP comprises $\times 8$ fewer parameters than their symmetry-agnostic counterparts.

G.3 UNCERTAINTY QUANTIFICATION VIA CONDITIONAL QUANTILE REGRESSION

The goal of these experiments benchmark is to learn the family of conditional distributions $\mathbb{P}(\mathbf{y} | \mathbf{x} = \cdot)$ for a bivariate random variable $\mathbf{y} = [y_0, y_1] \in \mathbb{R}^2$ given a scalar covariate $\mathbf{x} \in \mathbb{R}$. Once $\mathbb{P}(\mathbf{y} | \mathbf{x})$ is recovered, the practitioner can estimate conditional $(1 - \alpha)$ -confidence regions by regressing the lower and upper conditional quantiles $q_{\alpha/2}(\mathbf{x})$, $q_{1-\alpha/2}(\mathbf{x})$ for any desired miscoverage level $\alpha \in (0, 1)$. In particular, a 95% confidence region corresponds to $\alpha = 0.05$, so the two quantiles of interest are $q_{0.025}(\mathbf{x})$ and $q_{0.975}(\mathbf{x})$. See Fig. 7 for a visual representation of the problem.

Conditional quantile regression models We compare the NCP and proposed eNCP models to a standard baseline for parametric NN conditional quantile regression, namely CQR [Feldman et al. \(2023\)](#), which uses two separate NNs to predict lower and upper quantiles, trained with pinball loss. Both models use MLP backbones with similar parameter counts, ensuring improvements are solely due to loss functions.

Furthermore, CQR can only be trained for specific confidence intervals, requiring retraining for different quantiles. In contrast, the NCP and eNCP models, trained using the deep representation learning approach of [Secs. 2 and 4](#), regress the CCDF of each dimension of \mathbf{y} given \mathbf{x} . Thus, they can estimate conditional quantiles for any confidence interval via the quantile estimation algorithm from the CCDF described in [Kostic et al. \(2024a\)](#) without retraining. See details in [Fig. 8](#).

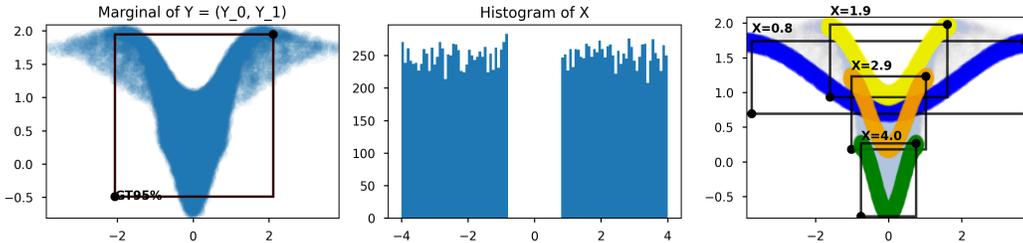


Figure 7: Synthetic uncertainty quantification experiment adapted from [Feldman et al. \(2023\)](#). Task: predict 95% confidence intervals (black boxes) of $\mathbf{y} \in \mathbb{R}^2$ given $\mathbf{x} \in \mathbb{R}$. **Left:** Marginal $\mathbb{P}(\mathbf{y})$. **Middle:** Marginal $\mathbb{P}(\mathbf{x})$. **Right:** Conditional distributions $\mathbb{P}(\mathbf{y} | \mathbf{x} = \cdot)$ for different values.

Evaluation metrics: coverage and set size Let $\mathbb{C}_{1-\alpha}(\mathbf{x}) \subseteq \mathbb{R}^d$ denote a *prediction set* of nominal level $(1 - \alpha)$ produced by a conditional quantile regression model for the response $\mathbf{y} \in \mathbb{R}^d$ given the covariate $\mathbf{x} \in \mathbb{R}^p$. In all experiments we assess two complementary metrics.

- **Coverage.** The conditional *coverage* of $\mathbb{C}_{1-\alpha}$ is the probability that the true response is captured by the predicted region,

$$c_{1-\alpha}(\mathbf{x}) := \mathbb{P}(\mathbf{y} \in \mathbb{C}_{1-\alpha}(\mathbf{x}) \mid \mathbf{x}), \quad \text{with the target } c_{1-\alpha}(\mathbf{x}) \approx 1 - \alpha \quad \forall \mathbf{x}. \quad (27)$$

In practice we report the *marginal coverage* $\widehat{\mathbb{E}}_{\mathbf{x}}[c_{1-\alpha}(\mathbf{x})]$, estimated on a large held-out sample; values above (resp. below) $1 - \alpha$ indicate over- (resp. under-) coverage.

- **Relaxed Coverage (r-Coverage).** The conditional *relaxed coverage* of $\mathbb{C}_{1-\alpha}$ is defined as the probability that each scalar component of the response lies within its corresponding predicted confidence interval. Formally, if $\mathbf{y} = [y_1, \dots, y_d]$ and $\mathbb{C}_{1-\alpha}(\mathbf{x})$ has corresponding marginal intervals $\mathbb{C}_{1-\alpha}^{(i)}(\mathbf{x})$ for $i \in \{1, \dots, d\}$, then

$$rc_{1-\alpha}(\mathbf{x}) := \prod_{i=1}^d \mathbb{P}(y_i \in \mathbb{C}_{1-\alpha}^{(i)}(\mathbf{x}) \mid \mathbf{x}), \quad (28)$$

with the target $rc_{1-\alpha}(\mathbf{x}) \approx 1 - \alpha$ for all \mathbf{x} . As with coverage, we report the *marginal relaxed coverage* $\widehat{\mathbb{E}}_{\mathbf{x}}[rc_{1-\alpha}(\mathbf{x})]$.

- **Set size.** To quantify how informative the region is, we measure its *size* (volume) under the Lebesgue measure λ^d :

$$\text{Size}_{1-\alpha}(\mathbf{x}) := \text{vol}(\mathbb{C}_{1-\alpha}(\mathbf{x})). \quad (29)$$

Smaller sets correspond to sharper uncertainty estimates, provided the required coverage is met. For multidimensional responses the volume is expressed in the natural units of \mathbb{R}^d ; for $d = 1$ it reduces to the interval length. As with coverage, we report the marginal expectation $\widehat{\mathbb{E}}_{\mathbf{x}}[\text{Size}_{1-\alpha}(\mathbf{x})]$ so that models can be compared fairly across the entire input distribution.

Data generation The data is generated from a conditional distribution $\mathbb{P}(\mathbf{y} \mid \mathbf{x} = \cdot)$ for a bivariate random variable $\mathbf{y} = [y_0, y_1] \in \mathbb{R}^2$ given a scalar covariate $\mathbf{x} \in \mathbb{R}$. Adapted from [Feldman et al. \(2023\)](#), the covariate is sampled uniformly: $\mathbf{x} \sim \text{Unif}([0.8, 4.0] \cup [-4.0, -0.8])$, and the response

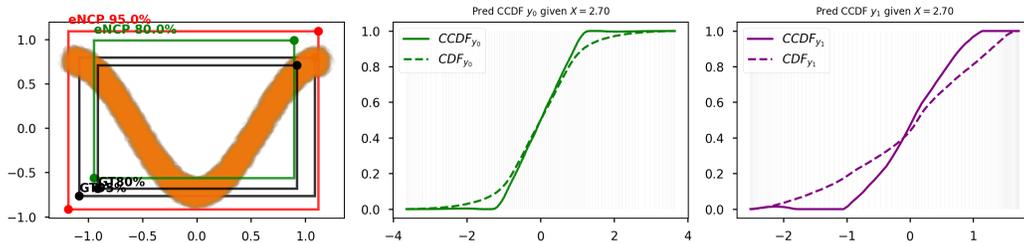


Figure 8: Prediction of the 80% and 95% confidence intervals for the random variable \mathbf{y} in experiment [Subsec. G.3](#) using the proposed eNCP model. The model estimates the CCDF by discretizing each dimension of $\mathbf{y} = [y_1, y_2]$ into 100 bins and computing the conditional probabilities $\mathbb{P}(y_i \in \mathbb{A}_n \mid \mathbf{x} = \cdot) := [\mathbb{E}_{\mathbf{y} \mid \mathbf{x}} \mathbb{1}_{\mathbb{A}_n}](\cdot)$ for all $n \in [100]$ based on the learned conditional expectation operator $\kappa_{\theta}(\mathbf{x}, \mathbf{y})$ (see [Sec. 5](#)). Here, \mathbb{A}_n comprises the bins from the 0-th to the n -th. This yields the estimated CCDF for y_1 (center) and y_2 (right) at $\mathbf{x} = 2.7$. The CCDFs can then be used to estimate upper and lower quantiles for any confidence interval ([Kostic et al., 2024a](#)). In practice, the eNCP model regresses 2×100 variables in a single forward pass. Thus, the final layer of the conditional quantile regression model is a linear layer of size $r \times (2 \times 100)$, where r is the number of features in the \mathbf{y} representation (see [Sec. 2](#)). Note that the gray vertical lines in the CCDF plots indicate the discretization bins, estimated using a ‘quantile.transformer’, from `sklearn` ([Pedregosa et al., 2011](#)), fitted on the marginal distribution of each dimension of \mathbf{y} . This discretization procedure ensures high-resolution bins in high-density regions of \mathbf{y} and coarser bins in low-density regions, while being able to cover the entire empirical support of \mathbf{y} from the training set.

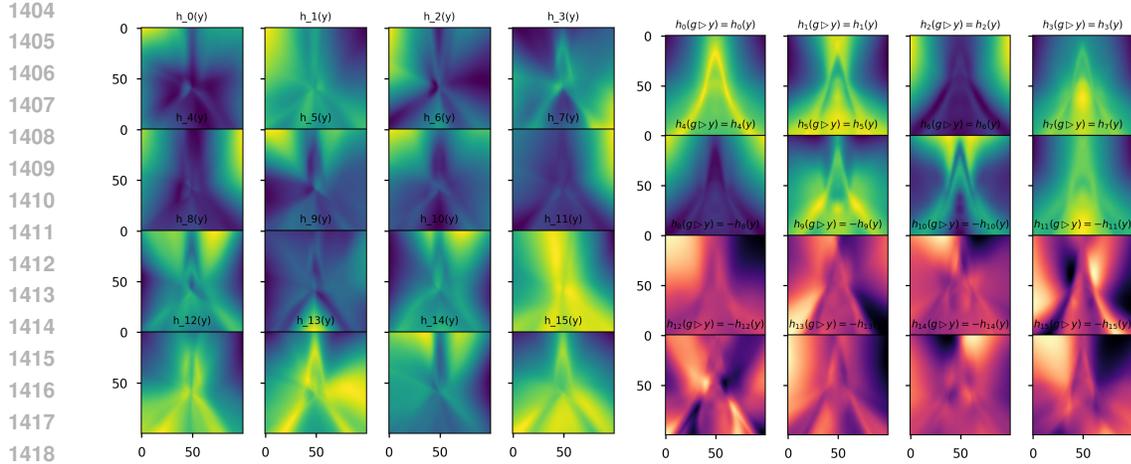


Figure 9: **Left:** Learned basis functions from the **NCP** model for $\mathbf{y} = [y_0, y_1]$. **Right:** Learned basis functions from the **eNCP** model for \mathbf{y} . The marginal distribution of \mathbf{y} exhibits reflection symmetry $g_r \triangleright \mathbf{y} = [-y_0, y_1]$ under $\mathbb{G} = \mathbb{C}_2$. Incorporating this prior, the **eNCP** model decomposes its latent space as $\mathcal{F}_{\mathbf{y}} = \mathcal{F}_{\mathbf{y}}^{\text{inv}} \oplus \mathcal{F}_{\mathbf{y}}^{(2)}$, with the first subspace capturing \mathbb{C}_2 -invariant functions and the second capturing those that change sign under reflection. The orthogonality of these subspaces allows independent optimization of the basis functions.

variable \mathbf{y} is produced by a non-linear transformation of auxiliary latent variables (see Fig. 7):

$$\begin{aligned} y_0 &= \frac{z}{\beta |\mathbf{x}|} + r \cos \phi, & z &\sim \text{Unif}(-\pi, \pi), \\ y_1 &= \frac{1}{2}(-\cos z + 1) + r \sin \phi + \sin |\mathbf{x}|, & \phi &\sim \text{Unif}(0, 2\pi), \\ & & r &\sim \text{Unif}(-0.1, 0.1). \end{aligned}$$

Here, $\beta > 0$ is a scaling constant. Both marginal distributions (see Fig. 7) and the conditional probability distribution are invariant under the reflection symmetry group $\mathbb{G} = \mathbb{C}_2 = \{e, g_r \mid g_r \circ g_r = e\}$, acting on \mathbf{y} as $g_r \triangleright \mathbf{y} = [-y_0, y_1]$ and on \mathbf{x} as $g_r \triangleright_{\mathcal{X}} \mathbf{x} = -\mathbf{x}$.

Results The experiment results are depicted in Fig. 10. Where the **NCP** and **eNCP** models outperform the baseline **CQR** model in terms of both coverage and set size. Furthermore, Fig. 9 illustrates the basis functions learned by the **NCP** and **eNCP** models for the random variable $\mathbf{y} = [y_0, y_1]$. In contrast to the standard **NCP** model, the **eNCP** model incorporates symmetry priors, enabling a clean separation of its latent representation into two orthogonal subspaces: one corresponding to \mathbb{C}_2 -invariant functions and the other to functions that change sign under reflection.

G.4 \mathbb{G} -EQUIVARIANT SYNTHETIC REGRESSION WITH DIFFERENT NOISE TYPES.

This synthetic experiment demonstrates how the **NCP** and **eNCP** frameworks can train a NNs for regression tasks under diverse noise conditions, particularly in applications where the conditional expectation (i.e., regression) may be uninformative due to skewed and asymmetric conditional distributions (see Fig. 11). Crucially, in such settings, the learned representations from **NCP** and **eNCP** can be reused without retraining to predict conditional quantiles at arbitrary coverage levels (see Eq. (27)), which is not possible with standard supervised **MSE** regression training.

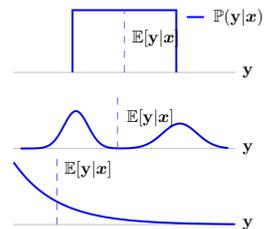


Figure 11: Conditional distributions with uninformative expected values.

Data generation. The dataset is composed of a 1-dimensional samples $\mathbf{x} \in \mathbb{R}$ and $\mathbf{y} \in \mathbb{R}$ drawn from a piecewise conditional distribution with three distinct regimes (see Fig. 12):

- Skewed distribution ($|x| \leq 1$): This regime features a skewed conditional distribution with an exponential tail, with heteroscedastic noise.
- Symmetric distribution ($1 < |x| \leq 2$): This regime has a symmetric conditional distribution, with heteroscedastic noise.

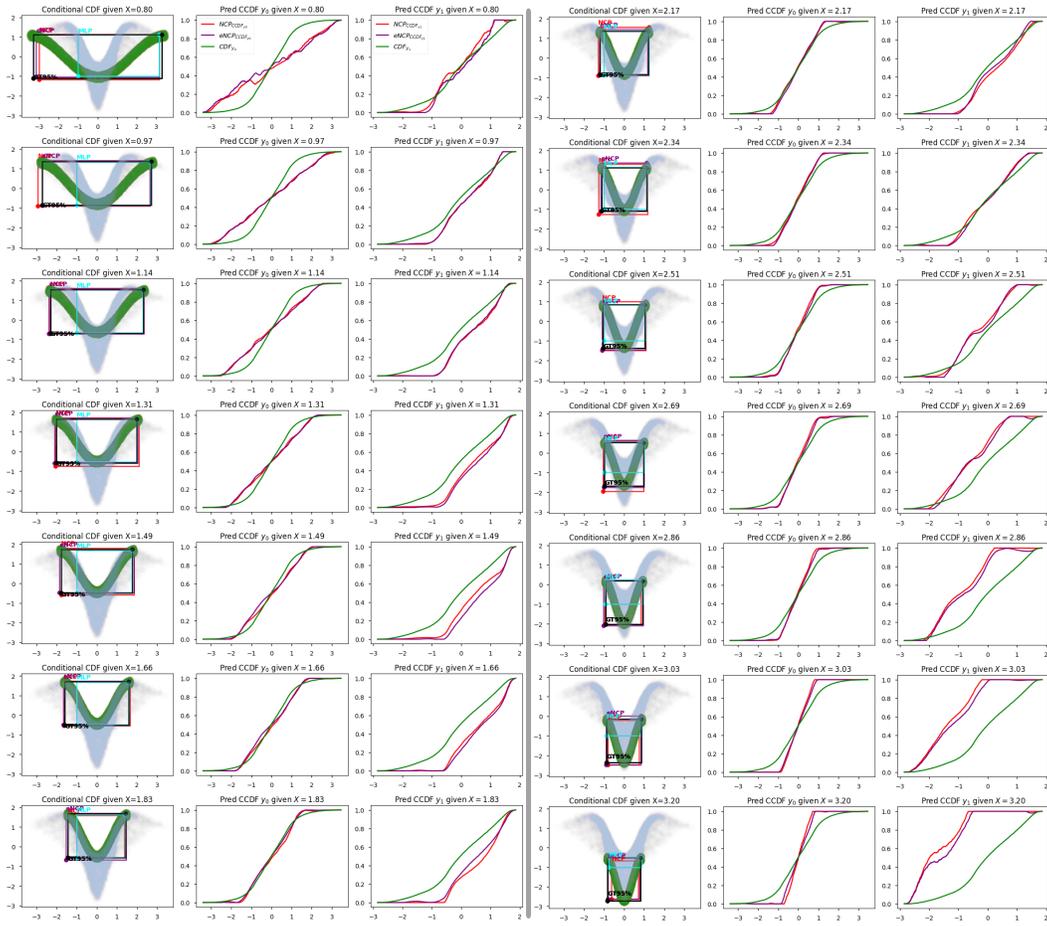


Figure 10: Results of a synthetic uncertainty quantification experiment comparing CQR, NCP, and eNCP models for predicting 95% confidence intervals of $\mathbf{y} \in \mathbb{R}^2$ given $\mathbf{x} \in \mathbb{R}$. Left and fourth columns show conditional distributions $\mathbb{P}(\mathbf{y}|\mathbf{x} = \cdot)$ for different conditioning values. Second-third and fifth-sixth columns display CCDF predictions by eNCP and NCP models. While CQR directly regresses quantiles and requires retraining for different confidence levels, NCP and eNCP estimate the full CCDF, enabling adaptation to any confidence interval without retraining.

- Bimodal distribution ($|x| > 2$): This regime exhibits a bimodal conditional distribution, with heteroscedastic noise.

Formally the piecewise conditional distribution is defined as follows:

$$\mathbb{P}(\mathbf{y} | \mathbf{x}) = \begin{cases} f_c(x) + \varepsilon_{\text{exp}}, & |x| \leq 1, \\ f_c(x) + \sigma_h(|x|)Z, & 1 < |x| \leq 2, \\ S a(|x|) + \sigma_p(|x|)Z', & |x| > 2, \end{cases} \quad f_c(x) = \frac{1}{2} \cos\left(\frac{2\pi}{3}x\right) + \frac{1}{3} \cos\left(\frac{8\pi}{3}x\right) + \frac{1}{4},$$

Where f_c is a deterministic function $\varepsilon_{\text{exp}} \sim \text{Exp}(\text{scale} = s_{\text{exp}}(|x|))$, $Z, Z' \sim \mathcal{N}(0, 1)$, $S \in \{-1, +1\}$ is equiprobable. Furthermore s_{exp} , σ_p and σ_h introduce heteroscedasticity (variance is conditioned on values of \mathbf{x}). The function is design such that the conditional distribution is \mathbb{G} -invariant under the reflection group $\mathbb{G} = \mathbb{C}_2$ acting as $g_r \triangleright \mathbf{x} = -\mathbf{x}$ and $g_r \triangleright \mathbf{y} = \mathbf{y}$. Consequently, the target \mathbb{G} -invariant function defined by the conditional expectation is given as (see Fig. 12):

$$z(x) = \mathbb{E}[\mathbf{y} | \mathbf{x} = x] = \begin{cases} f_c(x) + \mathbb{E}[\varepsilon_{\text{exp}}], & |x| \leq 1, \\ f_c(x), & 1 < |x| \leq 2, \\ 0, & |x| > 2, \end{cases}$$

Training details The dataset is generated with $N = 20,000$ sample pairs. The baseline MLP, NCP, and proposed eNCP share the same architecture backbone, a standard 3 layer deep MLP with

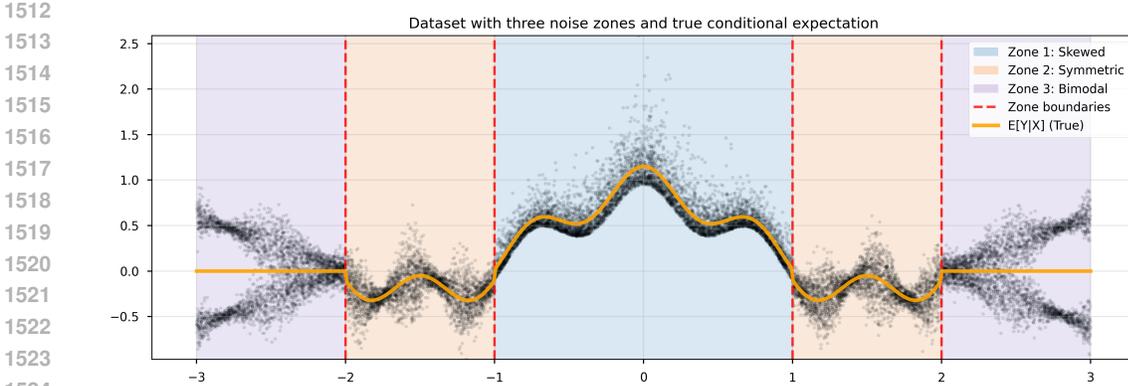


Figure 12: Synthetic regression experiment featuring three distinct zones with different families of conditional distributions for which the expected value is uninformative. Zone 1 ($|x| \leq 1$) exhibits a exponential skewed distribution, Zone 2 ($1 < |x| \leq 2$) has symmetric, and Zone 3 ($|x| > 2$) features a bimodal distributon. All zones are affected by heteroscedastic noise.

64 hidden neurons in each layer and ELU activations. The representation space dimension for **NCP** and **eNCP** is set to 32 (and consequently the last hidden layer of **MLP** is also set to 32 dimensions). Both contrastive models use orthogonality regularization 10^{-2} , a momentum 0.995 for computing the orthonormality statistics using exponential moving average estimates. All models are trained with early stopping on the validation objective, choosing the checkpoint with best validation loss to compute the test performance.

After training, for both contrastive models, we fit two linear decoders from the learned representations to predict the conditional expectation \hat{z}_θ , as described in (17), and to recover the **CCDF** with 200 support points for quantile extraction, as described in Fig. 8.

Results The results, shown in Fig. 12, illustrate how all models perform similarly in terms of predicting the conditional expectation $\mathbb{E}[y | x]$. However, the key advantage of the **NCP** and **eNCP** relies in their capacity to model the conditional probability distribution, which enables uncertainty quantification, shown here via prediction of conditional upper and lower quantiles with coverage of 90% and 70%. Note that recovering these quantiles is not possible when training the **NN** backbone using standard **MSE** supervised training.

While both **NCP** and **eNCP** perform comparatively well and predict well calibrated confidence intervals, the **eNCP** model achieves a perfect empirical coverage tracking for coverage levels from 10% to 90%, as shown in Fig. 13 (middle). Furthermore, **eNCP** consistently predicts confidence intervals of smaller size (bottom-middle), indicating accurate non-conservative estimates.

G.4.1 UNCERTAINTY QUANTIFICATION IN QUADRUPED LEGGED LOCOMOTION

We test how well conditional-quantile models can recover the conditional 95% confidence regions of three physically meaningful observables produced by a simulated AlienGo quadruped walking over rough terrain (see Fig. 1) under varying friction coefficients. The dataset was collected using the **Quadruped-PyMPC** simulation framework and model predictive controller from Turrisi et al. (2024).

The observables for which state-dependent uncertainty estimates are desired are $\mathbf{y}_t = [U_t, T_t, \boldsymbol{\tau}_t^{\text{grf}}]^\top$, with each component defined as follows:

- **G-invariant Kinetic Energy.** $T(\mathbf{q}, \dot{\mathbf{q}}) = \frac{1}{2} \dot{\mathbf{q}}^\top M(\mathbf{q}) \dot{\mathbf{q}} \in \mathbb{R}$, where $M(\mathbf{q})$ is the configuration-dependent inertia matrix. Noise is introduced through sensor measurement errors on the robot’s **degree of freedom (DoF)** position $\mathbf{q} \in \mathbb{R}^{12}$ and velocity $\dot{\mathbf{q}} \in \mathbb{R}^{12}$.
- **G-invariant Instantaneous Mechanical Work.** $U(\mathbf{q}, \dot{\mathbf{q}}, \boldsymbol{\tau}) \in \mathbb{R}$, representing the instantaneous mechanical work exerted or absorbed by the robot. This quantity depends on the actuator torques (typically measured with noisy, biased sensors) as well as the external forces (e.g. gravity, contact forces) that are not reliably measurable due to unobserved terrain parameters.

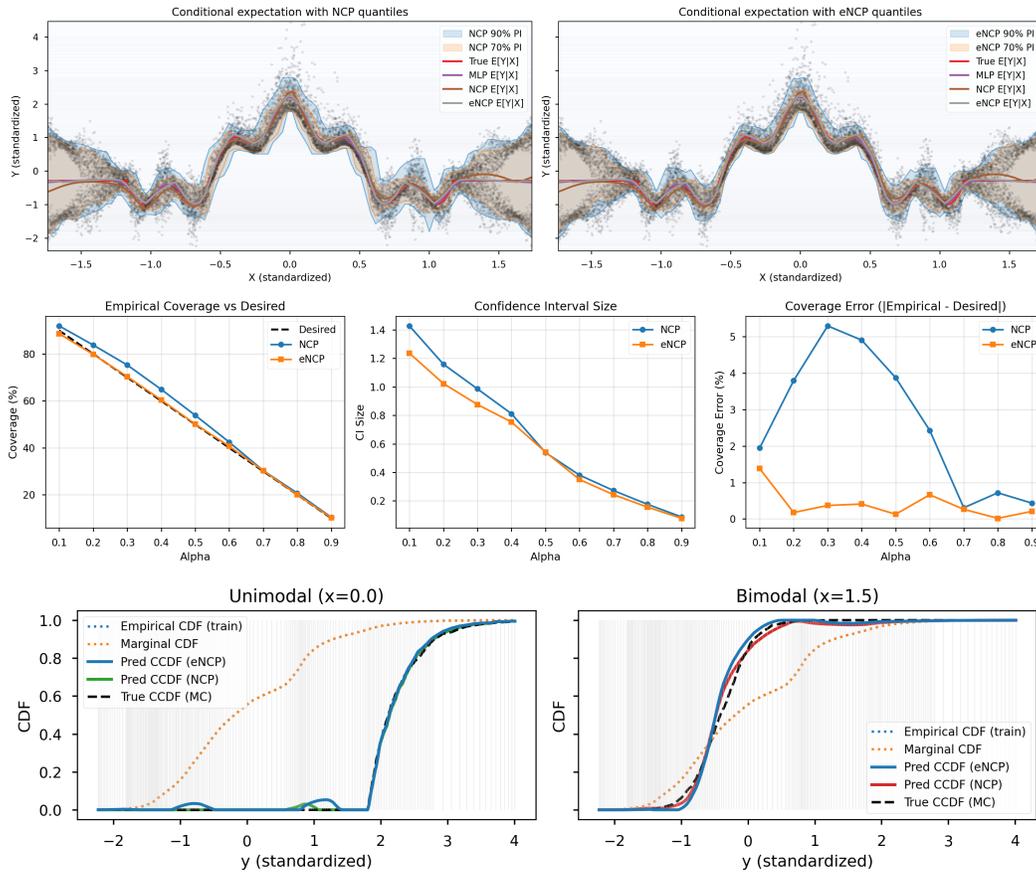


Figure 13: **Top:** Comparison of conditional expectation predictions $\mathbb{E}[y | x]$ and uncertainty quantification via 90% and 70% confidence intervals for baseline MLP, NCP, and proposed eNCP models. **Middle:** Empirical coverage (27) tracking and confidence intervals set size (29) for coverage levels from 10% to 90%. **Bottom:** Example CCDFs predicted by NCP and eNCP on the unimodal skewed distribution regime and the bimodal regime (see details in Fig. 8).

- **G-equivariant Ground-Reaction Forces** $\tau_{\text{grf}} \in \mathbb{R}^{12}$, a fundamental quantity in quadruped control, whose reliable estimation and uncertainty quantification are critical for downstream tasks in robotics (Nisticò et al., 2025; Liu et al., 1994).

The observables of interest are predicted using a suit of onboard proprioceptive sensory signals available at time t :

$$\mathbf{x}_t = \left[\mathbf{q}_t, \dot{\mathbf{q}}_t, \mathbf{a}_t, \mathbf{v}_t, \mathbf{v}_{t,\text{err}}, \boldsymbol{\omega}_t, \boldsymbol{\omega}_{t,\text{err}}, \mathbf{g}_t, \dot{\mathbf{p}}_{t,\text{feet}}, \boldsymbol{\tau}_t^{\text{cmd}} \right]^\top.$$

Specifically, $\mathbf{q}_t \in \mathbb{R}^{n_q}$ and $\dot{\mathbf{q}}_t \in \mathbb{R}^{n_q}$ are the joint positions and velocities, respectively; $\mathbf{a}_t \in \mathbb{R}^3$ is the linear acceleration of the robot’s base frame measured by the IMU; $\mathbf{v}_t \in \mathbb{R}^3$ is the base linear velocity, while $\mathbf{v}_{t,\text{err}} \in \mathbb{R}^3$ the command error base linear velocity; $\boldsymbol{\omega}_t \in \mathbb{R}^3$ and $\boldsymbol{\omega}_{t,\text{err}} \in \mathbb{R}^3$ are the base angular velocity and its command error; $\mathbf{g}_t \in \mathbb{R}^3$ is the gravity vector expressed in the base frame; $\dot{\mathbf{p}}_{t,\text{feet}} \in \mathbb{R}^{12}$ stacks the linear velocities of the four feet (three components each); and $\boldsymbol{\tau}_t^{\text{cmd}} \in \mathbb{R}^{n_q}$ contains the commanded joint torques.

Hence we design the experiments to compare models of similar footprint in number of parameters, while the loss used for training differs between the NCP and eNCP models w.r.t to the CQR and eCQR models.

NN architectures We configure all models (eNCP, NCP, eCQR, and CQR) with similar inference-time NN architectures. The backbone consists of three hidden layers with 512 units each, followed

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

	Validation			Test		
	r-Coverage \uparrow	Coverage \uparrow	Set Size \downarrow	r-Coverage \uparrow	Coverage \uparrow	Set Size \downarrow
eNCP	99.3 \pm 0.0%	94.1 \pm 0.4%	2.4 \pm 0.4 $\times 10^{10}$	99.5 \pm 0.1%	95.0 \pm 0.4%	4.3 \pm 3.6 $\times 10^9$
NCP	96.4 \pm 0.0%	56.9 \pm 0.1%	3.9 \pm 4.5 $\times 10^{10}$	99.5 \pm 0.0%	56.9 \pm 0.3%	2.6 \pm 1.4 $\times 10^{10}$
eCQR	70.7 \pm 0.6%	7.3 \pm 1.7%	3.7 \pm 2.6 $\times 10^8$	84.2 \pm 0.7%	6.7 \pm 1.2%	1.7 \pm 1.7 $\times 10^7$
CQR	67.6 \pm 1.8%	7.6 \pm 0.4%	2.5 \pm 2.4 $\times 10^9$	80.5 \pm 3.7%	8.5 \pm 0.9%	1.4 \pm 0.1 $\times 10^8$

Table 4: Validation and test metrics for 95% confidence intervals on quadruped robot observables traversing rough terrain (see Subsubsec. G.4.1). Metrics: (i) relaxed coverage (r-Coverage) (28), (ii) coverage (27), and (iii) set size (29). Best results in blue. While eCQR and CQR produce smaller confidence intervals, they fail to achieve expected 95% coverage on validation and test sets. The eNCP model achieves best overall coverage for reliable uncertainty quantification. Importantly, eNCP and NCP models can provide confidence intervals at any coverage level **without retraining**, whereas CQR and eCQR require retraining for each new level.

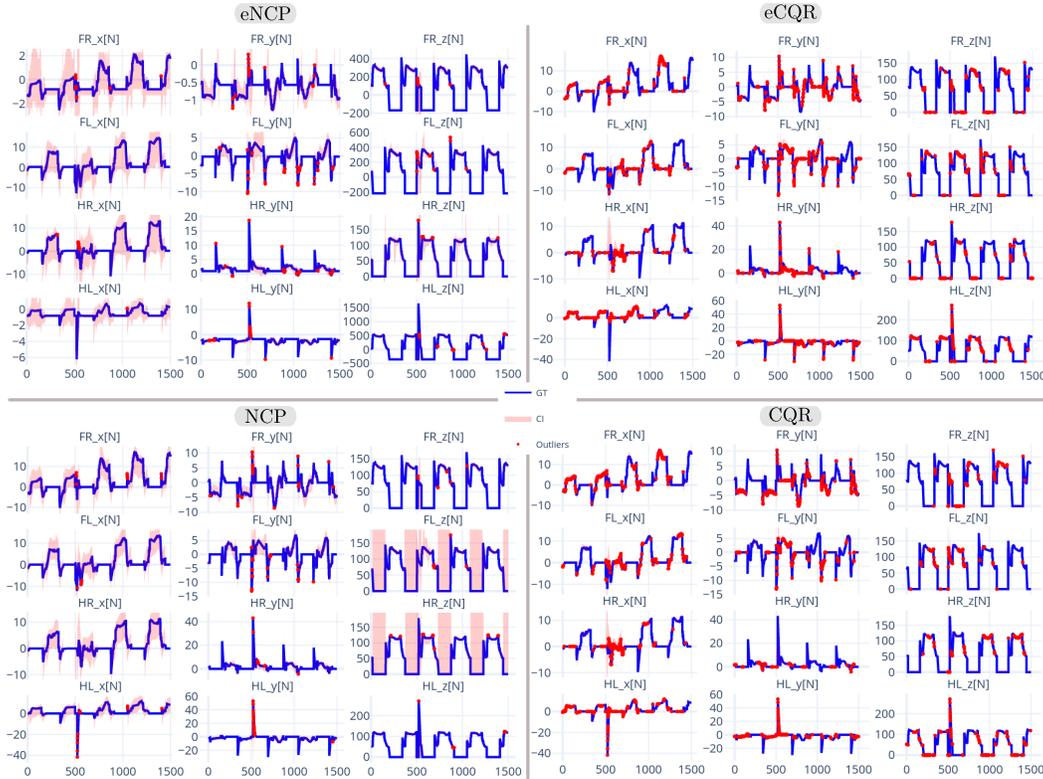


Figure 14: Prediction of 90% confidence intervals for ground-reaction forces $\tau_{\text{grf}} \in \mathbb{R}^{12}$ of a quadruped robot on rough terrain with varying friction. We compare eNCP, NCP, eCQR, and CQR models. CIs are computed for each leg—front-right (FR), front-left (FL), hind-right (HR), and hind-left—along x , y , and z axes. Forces outside the CI are highlighted in red, those within in blue. Terrain variations cause significant variability in x and y components due to surface orientation and friction differences, while z components are mainly influenced by local height changes affecting contact timing and producing short-duration high-impact forces.

by a final 128-unit layer that encodes the feature vector r for NCP and eNCP models. Due to \mathbb{G} -equivariance weight sharing, eNCP and eCQR have $\times 2$ fewer trainable parameters than their symmetry-agnostic counterparts.

Results. Given sensory input \mathbf{x} , the model predicts a set $\mathbb{C}_{0.95}(\mathbf{x}) \subseteq \mathbb{R}^{14}$ satisfying $\mathbb{P}(\mathbf{y} \in \mathbb{C}_{0.95}(\mathbf{x}) \mid \mathbf{x}) \approx 0.95$, while minimizing its volume $\widehat{\mathbb{E}}_{\mathbf{x}}[\text{vol}(\mathbb{C}_{0.95}(\mathbf{x}))]$. High coverage implies that the true \mathbb{G} -invariant kinetic energy, instantaneous mechanical work, and \mathbb{G} -equivariant 12-dimensional ground-reaction forces lie within the predicted confidence set. Relaxed coverage (r-

Coverage) quantifies the reliability of estimates on a per-dimension basis. Tab. 4 summarizes validation and test results for all models, and Fig. 14 illustrates a trajectory of GRF with respective 90% confidence intervals. Both CQR and eCQR produce smaller confidence intervals but fail to achieve desired coverage on the test set, implying unreliable confidence intervals requiring further calibration through retraining or conformal calibration (Feldman et al., 2023). In contrast, the eNCP model achieves desired coverage on the test set while producing larger confidence intervals, hence yielding reliable uncertainty estimates.

H CONDITIONAL PROBABILITY MODELING VIA THE CONDITIONAL EXPECTATION OPERATOR

This section introduces the modelling of conditional probabilities for two random variables via the **conditional expectation operator**. Our goal is to understand conditional expectation from an operator-theoretic perspective. We begin by describing the marginal, joint, and conditional probabilities of the random variables within a measure-theoretic framework. This discussion extends the exposition of Kostic et al. (2024b).

Given two random variables (\mathbf{x}, \mathbf{y}) taking values in the measure spaces $(\mathcal{X}, \Sigma_{\mathcal{X}}, P_{\mathbf{x}})$ and $(\mathcal{Y}, \Sigma_{\mathcal{Y}}, P_{\mathbf{y}})$, we have that the marginal probability of any set $\mathbb{A} \in \Sigma_{\mathcal{X}}$ and $\mathbb{B} \in \Sigma_{\mathcal{Y}}$ are given by

$$\mathbb{P}(\mathbf{x} \in \mathbb{A}) = \int_{\mathcal{X}} \mathbb{1}_{\mathbb{A}}(\mathbf{x}) P_{\mathbf{x}}(d\mathbf{x}) = \int_{\mathbb{A}} P_{\mathbf{x}}(d\mathbf{x}) \quad \text{and} \quad \mathbb{P}(\mathbf{y} \in \mathbb{B}) = \int_{\mathcal{Y}} \mathbb{1}_{\mathbb{B}}(\mathbf{y}) P_{\mathbf{y}}(d\mathbf{y}) = \int_{\mathbb{B}} P_{\mathbf{y}}(d\mathbf{y}), \quad (30)$$

where $\mathbb{1}_{\mathbb{A}} \in \mathcal{L}_{\mathbf{x}}^2$ and $\mathbb{1}_{\mathbb{B}} \in \mathcal{L}_{\mathbf{y}}^2$ denote the characteristic functions of sets \mathbb{A} and \mathbb{B} , respectively.

Furthermore, under the reasonable assumption that the joint probability measure is absolutely continuous w.r.t to the product of the marginals $P_{\mathbf{xy}} \ll P_{\mathbf{x}} \times P_{\mathbf{y}}$, we have that there exist a Radon-Nikodym derivative $\kappa : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_+$ such that $P_{\mathbf{xy}}(d\mathbf{x}, d\mathbf{y}) = \kappa(\mathbf{x}, \mathbf{y}) P_{\mathbf{x}}(d\mathbf{x}) P_{\mathbf{y}}(d\mathbf{y})$. Note that κ is a kernel function that pointwise deforms the product of the marginals to produce the joint distribution Sugiyama et al. (2012) (see Fig. 3). This kernel function enable us to express the joint probability by:

$$\mathbb{P}(\mathbf{x} \in \mathbb{A}, \mathbf{y} \in \mathbb{B}) = \int_{\mathcal{X} \times \mathcal{Y}} \mathbb{1}_{\mathbb{A}}(\mathbf{x}) \mathbb{1}_{\mathbb{B}}(\mathbf{y}) \underbrace{\kappa(\mathbf{x}, \mathbf{y}) P_{\mathbf{y}}(d\mathbf{y}) P_{\mathbf{x}}(d\mathbf{x})}_{P_{\mathbf{xy}}(d\mathbf{x}, d\mathbf{y})} = \int_{\mathbb{A} \times \mathbb{B}} \kappa(\mathbf{x}, \mathbf{y}) P_{\mathbf{x}}(d\mathbf{x}) P_{\mathbf{y}}(d\mathbf{y}). \quad (31)$$

Furthermore, given that $\mathbb{P}(\mathbf{y} \in \mathbb{B} | \mathbf{x} \in \mathbb{A}) = \mathbb{P}(\mathbf{x} \in \mathbb{A}, \mathbf{y} \in \mathbb{B}) / \mathbb{P}(\mathbf{x} \in \mathbb{A})$, the conditional probability of any set $\mathbb{B} \in \Sigma_{\mathcal{Y}}$ given a value of the random variable $\mathbf{x} = \mathbf{x}$ is given by:

$$\mathbb{P}(\mathbf{y} \in \mathbb{B} | \mathbf{x} = \mathbf{x}) = \int_{\mathcal{Y}} \mathbb{1}_{\mathbb{B}}(\mathbf{y}) P_{\mathbf{y} | \mathbf{x}}(d\mathbf{y} | \mathbf{x}) = \int_{\mathcal{Y}} \mathbb{1}_{\mathbb{B}}(\mathbf{y}) \kappa(\mathbf{x}, \mathbf{y}) P_{\mathbf{y}}(d\mathbf{y}) = \int_{\mathbb{B}} \kappa(\mathbf{x}, \mathbf{y}) P_{\mathbf{y}}(d\mathbf{y}), \quad (32)$$

where $P_{\mathbf{y} | \mathbf{x}} : \Sigma_{\mathcal{Y}} \times \mathcal{X} \mapsto [0, 1]$ denotes the **conditional probability measure**. This is a well-defined probability measure considering that:

$$\mathbb{P}(\mathbf{x} \in \mathbb{A}) := \mathbb{P}(\mathbf{x} \in \mathbb{A}, \mathbf{y} \in \mathcal{Y}) = \int_{\mathbb{A}} \underbrace{\left(\int_{\mathcal{Y}} \kappa(\mathbf{x}, \mathbf{y}) P_{\mathbf{y}}(d\mathbf{y}) \right)}_{\mathbb{E} P_{\mathbf{y} | \mathbf{x}}(d\mathbf{y} | \mathbf{x} = \mathbf{x}) = 1 \quad \forall \mathbf{x} \in \mathcal{X}} P_{\mathbf{x}}(d\mathbf{x}) = \int_{\mathbb{A}} P_{\mathbf{x}}(d\mathbf{x}).$$

The operator perspective Every measurable function $h \in \mathcal{L}_{\mathbf{y}}^2$ can be approximated by simple functions—that is, as a combination of characteristic functions on measurable sets: $h(\cdot) \approx \sum_{i \in \mathbb{N}} \beta_i \mathbb{1}_{\mathbb{A}_i}(\cdot)$. Thus, Eq. (32) is a special case of the more general problem of approximating the conditional expectation of any function $h \in \mathcal{L}_{\mathbf{y}}^2$ given \mathbf{x} . This conditional expectation is captured by the action of a linear integral operator:

Definition H.1 (Conditional expectation operator). *Let (\mathbf{x}, \mathbf{y}) be two random variables defined on the measure spaces $(\mathcal{X}, \Sigma_{\mathcal{X}}, P_{\mathbf{x}})$ and $(\mathcal{Y}, \Sigma_{\mathcal{Y}}, P_{\mathbf{y}})$, respectively, and let $\mathcal{L}_{\mathbf{x}}^2$ and $\mathcal{L}_{\mathbf{y}}^2$ denote the corresponding spaces of square-integrable functions. The conditional expectation operator $\mathbb{E}_{\mathbf{y} | \mathbf{x}} : \mathcal{L}_{\mathbf{y}}^2 \rightarrow \mathcal{L}_{\mathbf{x}}^2$ is the linear integral operator—defined via the PMD Radon-Nikodym derivative $\kappa(\mathbf{x}, \mathbf{y}) = P_{\mathbf{xy}}(d\mathbf{x}, d\mathbf{y}) / P_{\mathbf{x}}(d\mathbf{x}) P_{\mathbf{y}}(d\mathbf{y})$ —which acts on any function $h \in \mathcal{L}_{\mathbf{y}}^2$ by computing its conditional expectation:*

$$[\mathbb{E}_{\mathbf{y} | \mathbf{x}} h](\mathbf{x}) = \mathbb{E}[h(\mathbf{y}) | \mathbf{x} = \mathbf{x}] := \int_{\mathcal{Y}} h(\mathbf{y}) P_{\mathbf{y} | \mathbf{x}}(d\mathbf{y} | \mathbf{x}) = \int_{\mathcal{Y}} h(\mathbf{y}) \frac{P_{\mathbf{xy}}(d\mathbf{y}, \mathbf{x})}{P_{\mathbf{x}}(d\mathbf{x})} = \int_{\mathcal{Y}} h(\mathbf{y}) \kappa(\mathbf{x}, \mathbf{y}) P_{\mathbf{y}}(d\mathbf{y}).$$

From a learning perspective, approximating the conditional expectation operator sufficiently well for a relevant set of functions in \mathcal{L}_y^2 implies that we can approximate the conditional probability measure of any set $\mathbb{A} \in \Sigma_y$. This enables both regression *and* uncertainty quantification applications with a single model (see Eq. (2)).

I BACKGROUND ON GROUP AND REPRESENTATION THEORY

Group actions and representations This section provides a concise overview of the fundamental concepts in group and representation theory, which are used to define the symmetries of the random variables we consider in this work. For a comprehensive background on these topics in finite-dimensional vector spaces, see Weiler et al. (2023); for the infinite-dimensional case, consult Knapp (1986). These concepts will be referenced as needed in the main text. To begin, we define a group as an abstract mathematical object.

Definition I.1 (Group). *A group is a set \mathbb{G} , endowed with a binary composition operator defined as:*

$$\begin{aligned} (\circ) : \mathbb{G} \times \mathbb{G} &\longrightarrow \mathbb{G} \\ (g_1, g_2) &\longrightarrow g_1 \circ g_2, \end{aligned} \quad (33a)$$

such that the following axioms hold:

$$\text{Associativity: } (g_1 \circ g_2) \circ g_3 = g_1 \circ (g_2 \circ g_3), \quad \forall g_1, g_2, g_3 \in \mathbb{G}, \quad (33b)$$

$$\text{Identity: } \exists e \in \mathbb{G} \text{ such that } e \circ g = g = g \circ e, \quad \forall g \in \mathbb{G}, \quad (33c)$$

$$\text{Inverses: } \forall g \in \mathbb{G}, \exists g^{-1} \in \mathbb{G} \text{ such that } g \circ g^{-1} = e = g^{-1} \circ g. \quad (33d)$$

We are primarily interested in symmetry groups, i.e., groups of transformations acting on a set \mathcal{X} . Each transformation is a bijection that leaves a fundamental property of the element of the set invariant. For example, if \mathcal{X} represents states of a dynamical system, the invariant property is the state energy (see Fig. 6); if \mathcal{X} is a data space, the preserved quantity is typically the probability density/distribution (see Fig. 3).

Definition I.2 (Group action on a set (Weiler et al., 2023)). *Let \mathcal{X} be a set endowed with symmetry group \mathbb{G} . The (left) group action of the group \mathbb{G} on the set \mathcal{X} is a map:*

$$\begin{aligned} (\triangleright) : \mathbb{G} \times \mathcal{X} &\longrightarrow \mathcal{X} \\ (g, \mathbf{x}) &\longrightarrow g \triangleright \mathbf{x} \end{aligned} \quad (34a)$$

that is compatible with the group composition and identity element $e \in \mathbb{G}$, in the sense that:

$$\text{Identity: } e \triangleright \mathbf{x} = \mathbf{x}, \quad \forall \mathbf{x} \in \mathcal{X} \quad (34b)$$

$$\text{Associativity: } (g_1 \circ g_2) \triangleright \mathbf{x} = g_1 \triangleright (g_2 \triangleright \mathbf{x}), \quad \forall g_1, g_2 \in \mathbb{G}, \forall \mathbf{x} \in \mathcal{X}. \quad (34c)$$

For our practical purposes, we focus on symmetry transformations acting on sets with a vector space structure. In most cases, the group action on a vector space is linear, which allows us to express symmetry transformations as (orthogonal) matrix-vector operations, once a basis for the space is chosen.

Definition I.3 (Linear group representation). *Let \mathcal{X} be a vector space endowed with symmetry group \mathbb{G} . A linear representation of \mathbb{G} on \mathcal{X} is a map, denoted by $\rho_{\mathcal{X}}$, between symmetry transformation and invertible linear maps on \mathcal{X} (i.e., elements of the general linear group $\mathbb{GL}(\mathcal{X})$):*

$$\begin{aligned} \rho_{\mathcal{X}} : \mathbb{G} &\longrightarrow \mathbb{GL}(\mathcal{X}) \\ g &\longrightarrow \rho_{\mathcal{X}}(g), \end{aligned} \quad (35a)$$

such that the following properties hold:

$$\text{composition : } \rho_{\mathcal{X}}(g_1 \circ g_2) = \rho_{\mathcal{X}}(g_1) \rho_{\mathcal{X}}(g_2), \quad \forall g_1, g_2 \in \mathbb{G}, \quad (35b)$$

$$\text{inversion : } \rho_{\mathcal{X}}(g^{-1}) = \rho_{\mathcal{X}}(g)^{-1}, \quad \forall g \in \mathbb{G}. \quad (35c)$$

$$\text{identity : } \rho_{\mathcal{X}}(g \circ g^{-1}) = \rho_{\mathcal{X}}(e) = \mathbf{I}, \quad (35d)$$

Whenever the vector space is of finite dimension $n < \infty$, linear maps admit a matrix form $\rho_{\mathcal{X}}(g) \in \mathbb{R}^{n \times n}$, once a basis set $\mathbb{I}_{\mathcal{X}}$ for the vector space \mathcal{X} is chosen. In this case, Eqs. (35b) to (35d) show

1782 how the composition and inversion of symmetry transformations translate to matrix multiplication
 1783 and inversion, respectively. Moreover, $\rho_{\mathcal{X}}$ allows to express a (linear) group action (Def. I.2) as a
 1784 matrix-vector multiplication:

$$1785 \quad (\triangleright) : \mathbb{G} \times \mathcal{X} \longrightarrow \mathcal{X} \\
 1786 \quad \quad \quad (g, \mathbf{x}) \longrightarrow g \triangleright \mathbf{x} := \rho_{\mathcal{X}}(g)\mathbf{x}. \quad (35e) \\
 1787$$

1788 We will often study linear maps that preserve a vector space’s symmetry structure by commuting
 1789 with the group action. These maps are known as **endomorphisms** and include all change of basis
 1790 and affine transformations that do not break the symmetry.

1791 **Definition I.4** (Endomorphism). *Let \mathcal{X} be a vector space endowed with symmetry group \mathbb{G} , with
 1792 the group action $\triangleright_{\mathcal{X}}: \mathbb{G} \times \mathcal{X} \mapsto \mathcal{X}$. A linear map $\mathbf{A} : \mathcal{X} \mapsto \mathcal{X}$ is said to be an endomorphism if it
 1793 commutes with the group action, such that:*

$$1794 \quad \rho_{\mathcal{X}}(g)\mathbf{A} = \mathbf{A}\rho_{\mathcal{X}}(g), \quad \forall g \in \mathbb{G} \quad \iff \quad \begin{array}{ccc} \mathcal{X} & \xrightarrow{\triangleright_{\mathcal{X}}} & \mathcal{X} \\ \downarrow \mathbf{A} & & \downarrow \mathbf{A} \\ \mathcal{X} & \xrightarrow{\triangleright_{\mathcal{X}}} & \mathcal{X} \end{array}$$

1795 We will denote the space of all endomorphisms of \mathcal{X} as $End_{\mathbb{G}}(\mathcal{X})$, such that any $\mathbf{A} \in End_{\mathbb{G}}(\mathcal{X})$
 1796 satisfies the above commutation property.

1802 So far, we have studied symmetric vector spaces \mathcal{X} endowed with a group action $\triangleright_{\mathcal{X}}$, which can be
 1803 represented in matrix form via a group representation $\rho_{\mathcal{X}}$ once a basis for \mathcal{X} is chosen. However,
 1804 while the choice of basis alters the group representation $\rho_{\mathcal{X}}$, the underlying group action $\triangleright_{\mathcal{X}}$ remains
 1805 invariant. This observation leads us to the concept of equivalent group representations.

1806 **Definition I.5** (Equivalent group representations). *Let \mathcal{X} be a vector space endowed with symmetry
 1807 group \mathbb{G} , and let $\rho'_{\mathcal{X}}$ and $\rho_{\mathcal{X}}$ be two group representations of \mathbb{G} on \mathcal{X} . They are said to be equivalent,
 1808 denoted by $\rho'_{\mathcal{X}} \sim \rho_{\mathcal{X}}$, if there exists an invertible change of basis $\mathbf{Q} \in End_{\mathbb{G}}(\mathcal{X})$ such that*

$$1809 \quad \rho'_{\mathcal{X}}(g) = \mathbf{Q}\rho_{\mathcal{X}}(g)\mathbf{Q}^{-1}, \quad \forall g \in \mathbb{G}. \quad (36) \\
 1810$$

1811 *Equivalent representations arise when the same group action $(\triangleright) : \mathbb{G} \times \mathcal{X} \rightarrow \mathcal{X}$ is expressed in dif-*
 1812 *ferent coordinate frames or bases. For instance, let $\mathbb{A}_{\mathcal{X}}$ and $\mathbb{B}_{\mathcal{X}}$ be two bases for $\mathcal{X} = span(\mathbb{A}_{\mathcal{X}}) =$
 1813 $span(\mathbb{B}_{\mathcal{X}})$, and let $\mathbf{Q}_{\mathbb{A}}^{\mathbb{B}} : \mathcal{X} \rightarrow \mathcal{X}$ denote the change of basis from $\mathbb{A}_{\mathcal{X}}$ to $\mathbb{B}_{\mathcal{X}}$, so that $\mathbf{x}^{\mathbb{B}} = \mathbf{Q}_{\mathbb{A}}^{\mathbb{B}}\mathbf{x}^{\mathbb{A}}$
 1814 for all $\mathbf{x}^{\mathbb{A}} \in \mathcal{X}$. Then the group action admits equivalent representations, $\rho_{\mathcal{X}}^{\mathbb{A}} \sim \rho_{\mathcal{X}}^{\mathbb{B}}$, since*

$$1815 \quad g \triangleright \mathbf{x}^{\mathbb{B}} := \mathbf{Q}_{\mathbb{A}}^{\mathbb{B}}(g \triangleright \mathbf{x}^{\mathbb{A}}), \quad \forall g \in \mathbb{G}, \\
 1816 \quad \rho_{\mathcal{X}}^{\mathbb{B}}(g)\mathbf{x}^{\mathbb{B}} = \mathbf{Q}_{\mathbb{A}}^{\mathbb{B}}(\rho_{\mathcal{X}}^{\mathbb{A}}(g)\mathbf{x}^{\mathbb{A}}) = \left(\mathbf{Q}_{\mathbb{A}}^{\mathbb{B}}\rho_{\mathcal{X}}^{\mathbb{A}}(g)\mathbf{Q}_{\mathbb{A}}^{\mathbb{B}-1}\right)\mathbf{x}^{\mathbb{B}}, \quad (37) \\
 1817 \quad \rho_{\mathcal{X}}^{\mathbb{B}}(g) = \mathbf{Q}_{\mathbb{A}}^{\mathbb{B}}\rho_{\mathcal{X}}^{\mathbb{A}}(g)\mathbf{Q}_{\mathbb{A}}^{\mathbb{B}-1}. \\
 1818 \\
 1819$$

1820 To reveal the modular structure of symmetric vector spaces, we often change bases to decompose
 1821 them into subspaces stable under the action of the group \mathbb{G} , termed \mathbb{G} -stable subspaces. This de-
 1822 composition mirrors how a symmetry group can be broken down into products and direct products
 1823 of smaller groups and is essential for analyzing and simplifying group representations. We introduce
 1824 the following definition.

1825 **Definition I.6** (\mathbb{G} -stable and irreducible subspaces). *Let \mathcal{X} be a vector space endowed with a group
 1826 action (\triangleright) of the symmetry group \mathbb{G} . A subspace $\mathcal{X}' \subseteq \mathcal{X}$ is said to be \mathbb{G} -stable if the action of any
 1827 group element on any vector in the subspace remains within the subspace, that is,*

$$1828 \quad g \triangleright \mathbf{x} \in \mathcal{X}', \quad \forall \mathbf{x} \in \mathcal{X}' \subseteq \mathcal{X}, \forall g \in \mathbb{G}. \\
 1829$$

1830 *If the only \mathbb{G} -stable subspaces of \mathcal{X} are $\{\mathbf{0}\}$ and \mathcal{X} itself, then \mathcal{X} is a irreducible \mathbb{G} -stable space.
 1831 We will denote irreducible \mathbb{G} -stable spaces with an over bar, e.g., $\bar{\mathcal{V}}$.*

1832 **Building blocks of symmetric vector spaces** A recurrent theme in this work and in geometric
 1833 deep learning is to study and process symmetric vector spaces by decomposing them in terms of ir-
 1834 reducible \mathbb{G} -stable spaces. This process directly corresponds to the decomposition of the associated
 1835 group representation $\rho_{\mathcal{X}}$ into smaller representations acting on these \mathbb{G} -stable subspaces.

Definition I.7 (Decomposable representation). Let \mathcal{X} be a vector space with a group action (\triangleright) defined by the representation $\rho_{\mathcal{X}}$ in a chosen basis $\mathbb{A}_{\mathcal{X}}$. The representation is decomposable if it is equivalent to a direct sum of two lower-dimensional representations, $\rho_{\mathcal{X}} \sim \rho_{\mathcal{X}_1} \oplus \rho_{\mathcal{X}_2}$, where \mathcal{X}_1 and \mathcal{X}_2 are \mathbb{G} -stable subspaces of \mathcal{X} . Equivalently, there exists a change of basis $\mathbf{Q}_{\mathbb{A}}^{\mathbb{B}} : \mathcal{X} \rightarrow \mathcal{X}$ such that

$$\rho_{\mathcal{X}}^{\mathbb{B}} = \begin{bmatrix} \rho_{\mathcal{X}_1} & 0 \\ 0 & \rho_{\mathcal{X}_2} \end{bmatrix} = \mathbf{Q}_{\mathbb{A}}^{\mathbb{B}} \rho_{\mathcal{X}} \mathbf{Q}_{\mathbb{A}}^{\mathbb{B}-1}, \text{ and } g \triangleright \mathbf{x}^{\mathbb{B}} := \rho_{\mathcal{X}}^{\mathbb{B}}(g) \mathbf{x}^{\mathbb{B}} = \begin{bmatrix} \rho_{\mathcal{X}_1}(g) \mathbf{x}_1^{\mathbb{B}} \\ \rho_{\mathcal{X}_2}(g) \mathbf{x}_2^{\mathbb{B}} \end{bmatrix}, \text{ where } \mathbf{Q}_{\mathbb{A}}^{\mathbb{B}} \mathbf{x} = \begin{bmatrix} \mathbf{x}_1^{\mathbb{B}} \in \mathcal{X}_1 \\ \mathbf{x}_2^{\mathbb{B}} \in \mathcal{X}_2 \end{bmatrix}$$

Hence, the representation's decomposition $\rho_{\mathcal{X}} \sim \rho_{\mathcal{X}_1} \oplus \rho_{\mathcal{X}_2}$ corresponds to **decomposing the vector space** into \mathbb{G} -stable subspaces, $\mathcal{X} = \mathcal{X}_1 \oplus \mathcal{X}_2$.

When iteratively applying the decomposition process, we eventually reach representations that cannot be further decomposed. These are known as irreducible representations, or *irreps*, and they serve as the fundamental building blocks for all representations of a compact symmetry group \mathbb{G} . From a vector space perspective, the irreducible \mathbb{G} -stable subspaces (Thm. I.6) associated with these irreps are the elementary subspaces that comprise any symmetric vector space, analogous to how one-dimensional subspaces are the fundamental components of standard vector spaces.

Definition I.8 (Irreducible representation). Let \mathcal{X} be a vector space endowed with a group action (\triangleright) of a symmetry group \mathbb{G} . A representation $\rho_{\mathcal{X}}$ of \mathbb{G} on \mathcal{X} is said to be irreducible if it cannot be decomposed into smaller representations acting on proper \mathbb{G} -stable subspaces (Thm. I.6). That is, the only \mathbb{G} -stable subspaces $\mathcal{X}' \subseteq \mathcal{X}$ are $\mathcal{X}' = \{\mathbf{0}\}$ and $\mathcal{X}' = \mathcal{X}$ itself.

To differentiate irreps from decomposable representations we will denote the formers and their associated irreducible \mathbb{G} -stable spaces with an over bar: $\bar{\rho}_{\mathcal{V}} : \mathbb{G} \rightarrow \mathbb{GL}(\bar{\mathcal{V}})$.

Crucially, any compact symmetry group \mathbb{G} has a unique set of countably many irreps, denoted by $\{\bar{\rho}_k\}_{k \in [1, n_{\text{iso}}]}$, where k denotes the irrep type and $n_{\text{iso}} \leq |\mathbb{G}|$ denotes the number of unique irreps of \mathbb{G} . A fundamental property of these irreps is that any two non-equivalent irreps act on vector spaces that are mutually **orthogonal**. This implies that whenever we decompose symmetric vector spaces into their irreducible subspaces we are inherently decomposing the space into orthogonal \mathbb{G} -stable subspaces, which will greatly simplify numerical and theoretical analyses. Formally, these orthogonality relations are a consequence of Schur's lemma, which we state below in its original form for the case of complex irreps and discuss its adaptation to real irreps.

Lemma I.9 (Schur's Lemma for unitary (complex) representations (Knapp, 1986, Prop 1.5)). Consider two complex Hilbert spaces, \mathcal{H} and \mathcal{H}' , endowed with the (complex) irreducible unitary representations $\bar{\phi} : \mathbb{G} \mapsto \mathbb{U}(\mathcal{H})$ and $\bar{\phi}' : \mathbb{G} \mapsto \mathbb{U}(\mathcal{H}')$, respectively. Let $\mathbf{T} : \mathcal{H} \mapsto \mathcal{H}'$ be a linear map commuting with the group actions, such that $\mathbf{T} \in \text{Homo}_{\mathbb{C}}(\mathcal{H}, \mathcal{H}')$. Then, if the irreducible representations are not equivalent, i.e., $\bar{\phi} \not\sim \bar{\phi}'$, then \mathbf{T} is the trivial (or zero) map. Conversely, if $\bar{\phi} \sim \bar{\phi}'$, then \mathbf{T} is a constant multiple of an isomorphism. Denoting \mathbf{I} as the identity operator, this can be expressed as:

$$\bar{\phi} \not\sim \bar{\phi}' \iff \mathbf{0}_{\mathcal{H}'} = \mathbf{T} \mathbf{h} \mid \forall \mathbf{h} \in \mathcal{H} \quad (38a)$$

$$\bar{\phi} \sim \bar{\phi}' \iff \mathbf{T} = \alpha \mathbf{U}, \alpha \in \mathbb{C}, \mathbf{U} \cdot \mathbf{U}^H = \mathbf{I} \quad (38b)$$

$$\bar{\phi} = \bar{\phi}' \iff \mathbf{T} = \alpha \mathbf{I}, \alpha \in \mathbb{C} \quad (38c)$$

The most common interpretation of Schur's lemma is that whenever the irreps are equivalent, $\bar{\phi} \sim \bar{\phi}'$, their associated spaces are isomorphic, $\mathcal{H} \sim \mathcal{H}'$ and \mathbf{T} is an element of the endomorphism space $\text{End}_{\mathbb{C}}^{\mathbb{C}}(\bar{\mathcal{H}})$, with $\bar{\mathcal{H}} \sim \mathcal{H} \sim \mathcal{H}'$ (see Thm. I.4). Consequently, (38b) implies that the endomorphism space is one-dimensional, i.e., $\dim(\text{End}_{\mathbb{C}}^{\mathbb{C}}(\bar{\mathcal{H}})) = 1$, with $\alpha \in \mathbb{C}$ denoting the only degree of freedom, and (38c) denotes the scenario in which the basis sets for the two spaces are identical.

However, this result holds only for *complex* irreducible representations and requires adaptation for the case of real irreducible representations. The main difference stems from the fact that if $\bar{\rho} : \mathbb{G} \rightarrow \mathbb{GL}(\bar{\mathcal{V}})$ is a real irrep, then the space of (real) endomorphisms, $\text{End}_{\mathbb{R}}(\bar{\mathcal{V}})$, is no longer one-dimensional, but rather it can be 1, 2, or 4 dimensional, depending on whether $\text{End}_{\mathbb{R}}(\bar{\mathcal{V}})$ is isomorphic to the real algebra ($\mathbb{R} = \text{span}\{1\}$), complex algebra ($\mathbb{C} = \text{span}\{1, i \mid i^2 = -1\}$), or quaternionic algebra ($\mathbb{H} = \text{span}\{1, i, j, k \mid i^2 = j^2 = k^2 = -1\}$), respectively. Denoting by $\Psi : \mathbb{K} \rightarrow \text{End}_{\mathbb{R}}(\bar{\mathcal{V}})$ the isomorphism of basis elements of $\mathbb{K} \in \{\mathbb{R}, \mathbb{C}, \mathbb{H}\}$ with the basis elements

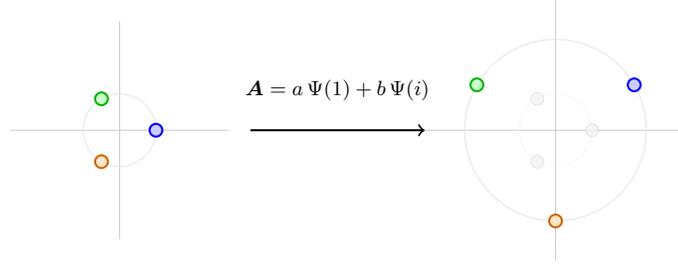


Figure 15: Example of an endomorphism acting on a \mathbb{C}_3 -stable irreducible 2D space. The irreducible representation is of complex type, with endomorphism space $\text{End}_{\mathbb{C}_3}(\mathbb{R}^2) \sim \mathbb{C} = \text{span}\{1, i\}$, comprising all transformations that uniformly scale and rotate/reflect the plane.

of $\text{End}_{\mathbb{G}}(\bar{\mathcal{V}})$, we can summarize the basis sets of the three cases as follows (see Cesa et al. (2022, Appendix C) for details):

$$\begin{aligned}
 \mathbb{R} &\sim \text{End}_{\mathbb{G}}(\mathcal{V}) = \text{span}\{\Psi(1) = \mathbf{I}_{\dim \bar{\phi}}\} \\
 \mathbb{C} &\sim \text{End}_{\mathbb{G}}(\mathcal{V}) = \text{span}\{\Psi(1) = \begin{bmatrix} \mathbf{I}_n & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_n \end{bmatrix}, \Psi(i) = \begin{bmatrix} \mathbf{0} & -\mathbf{I}_n \\ \mathbf{I}_n & \mathbf{0} \end{bmatrix}\} \\
 \mathbb{H} &\sim \text{End}_{\mathbb{G}}(\mathcal{V}) = \text{span} \left\{ \begin{array}{l} \Psi(1) = \begin{bmatrix} \mathbf{I}_n & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_n & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I}_n & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{I}_n \end{bmatrix}, \Psi(i) = \begin{bmatrix} \mathbf{0} & \mathbf{0} & -\mathbf{I}_n & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & -\mathbf{I}_n \\ \mathbf{I}_n & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_n & \mathbf{0} & \mathbf{0} \end{bmatrix} \\ \Psi(j) = \begin{bmatrix} \mathbf{0} & -\mathbf{I}_n & \mathbf{0} & \mathbf{0} \\ \mathbf{I}_n & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{I}_n \\ \mathbf{0} & \mathbf{0} & -\mathbf{I}_n & \mathbf{0} \end{bmatrix}, \Psi(k) = \begin{bmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} & -\mathbf{I}_n \\ \mathbf{0} & \mathbf{0} & \mathbf{I}_n & \mathbf{0} \\ \mathbf{0} & -\mathbf{I}_n & \mathbf{0} & \mathbf{0} \\ \mathbf{I}_n & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix} \end{array} \right\} \quad (39)
 \end{aligned}$$

While this result might appear complex, its interpretation is straightforward: given a \mathbb{G} -stable irreducible space $\bar{\mathcal{V}}$, the space of linear maps from the space to itself that preserve the symmetry structure consists of *linear transformations that scale all dimensions of $\bar{\mathcal{V}}$ uniformly, and possibly rotate or reflect the space*. Algebraically, this implies that any element of the algebra has a unique singular space, with a single singular value determined by the element's coefficients in the basis of (39). We summarize this result in the following proposition.

Proposition L.10 (A real endomorphism has a single singular space). *Let \mathbb{G} be a compact symmetry group, $(\bar{\rho}, \bar{\mathcal{V}})$ be an irreducible representation and its associated \mathbb{G} -stable space, and let $\text{End}_{\mathbb{G}}(\bar{\mathcal{V}})$ denote the space endomorphism algebra. Then every $\mathbf{A} \in \text{End}_{\mathbb{G}}(\bar{\mathcal{V}})$ admits an SVD*

$$\mathbf{A} = \mathbf{U} \gamma \mathbf{I}_d \mathbf{V}^\top,$$

where $\gamma \in \mathbb{R}_{\geq 0}$ is the single singular value, repeated with multiplicity $d = |\bar{\rho}| = |\mathcal{V}_k|$. The right singular basis \mathbf{V} can, without loss of generality, be taken as the canonical orthonormal basis of $\bar{\mathcal{V}}$, while the left singular basis is then $\mathbf{U} = \gamma^{-1} \mathbf{A} \mathbf{V}$, which is an orthogonal rotation/reflection of $\bar{\mathcal{V}}$.

(\mathbb{R}) **Real case.** Any $\mathbf{A} \in \text{End}_{\mathbb{G}}(\bar{\mathcal{V}})$ is of the form

$$\mathbf{A} = a \Psi(1) = a \mathbf{I}_d, \quad a \in \mathbb{R}.$$

Hence

$$\mathbf{A}^\top \mathbf{A} = a^2 \mathbf{I}_d, \quad \sigma(\mathbf{A}) = \{|a|\}^{\times d}.$$

(\mathbb{C}) **Complex case.** Every element can be written as

$$\mathbf{A} = a \Psi(1) + b \Psi(i) = \begin{bmatrix} a\mathbf{I}_n & -b\mathbf{I}_n \\ b\mathbf{I}_n & a\mathbf{I}_n \end{bmatrix}, \quad a, b \in \mathbb{R}.$$

Using $\Psi(i)^\top = -\Psi(i)$ and $\Psi(i)^\top \Psi(i) = \mathbf{I}$,

$$\mathbf{A}^\top \mathbf{A} = (a^2 + b^2) \mathbf{I}_d, \quad \sigma(\mathbf{A}) = \{\sqrt{a^2 + b^2}\}^{\times d}.$$

(\mathbb{H}) **Quaternionic case.** Each element admits the expansion

$$\mathbf{A} = a \Psi(1) + b \Psi(i) + c \Psi(j) + d \Psi(k), \quad a, b, c, d \in \mathbb{R},$$

where $\Psi(i), \Psi(j), \Psi(k)$ are the quaternionic structure matrices from (39), satisfying $\Psi(\alpha)^\top = -\Psi(\alpha)$, $\Psi(\alpha)^2 = -\mathbf{I}$, and $\Psi(\alpha)^\top \Psi(\alpha) = \mathbf{I}$, with the usual anti-commutation rules. Consequently,

$$\mathbf{A}^\top \mathbf{A} = (a^2 + b^2 + c^2 + d^2) \mathbf{I}_d, \quad \sigma(\mathbf{A}) = \{\sqrt{a^2 + b^2 + c^2 + d^2}\}^{\times d}.$$

As an intuitive low-dimensional example, we can consider the case of a 2D rotational irrep of the cyclic group \mathbb{C}_3 . The dimension of the irreducible \mathbb{G} -stable subspace is $|\text{barvs}V| = 2$, and the irreducible representation is of complex type, $\bar{\rho} : \mathbb{C}_3 \rightarrow \mathbb{GL}(\bar{V})$, with $\text{End}_{\mathbb{G}}(\bar{V}) \sim \mathbb{C} = \text{span}\{1, i\}$ denoting the space of all rotations/reflections and uniform scaling of the plane, see Fig. 15.

I.1 MAPS BETWEEN SYMMETRIC VECTOR SPACES

We will frequently study and use linear and non-linear maps between symmetric vector spaces. Our focus is on maps that preserve entirely or partially the group structure of the vector spaces. These types of maps can be classified as \mathbb{G} -equivariant, \mathbb{G} -invariant maps:

Definition I.11 (\mathbb{G} -equivariant and \mathbb{G} -invariant maps). *Let \mathcal{X} and \mathcal{Y} be two vector spaces endowed with the same symmetry group \mathbb{G} , with the respective group actions $\triangleright_{\mathcal{X}}$ and $\triangleright_{\mathcal{Y}}$. A map $f : \mathcal{X} \mapsto \mathcal{Y}$ is said to be \mathbb{G} -equivariant if it commutes with the group action, such that:*

$$g \triangleright_{\mathcal{Y}} \mathbf{y} = g \triangleright_{\mathcal{Y}} f(\mathbf{x}) = f(g \triangleright_{\mathcal{X}} \mathbf{x}), \quad \forall \mathbf{x} \in \mathcal{X}, g \in \mathbb{G}. \quad \iff \quad \begin{array}{ccc} \mathcal{X} & \xrightarrow{\triangleright_{\mathcal{X}}} & \mathcal{X} \\ \downarrow f & & \downarrow f \\ \mathcal{Y} & \xrightarrow{\triangleright_{\mathcal{Y}}} & \mathcal{Y} \end{array} \quad (40a)$$

$$\rho_{\mathcal{Y}}(g)f(\mathbf{x}) = f(\rho_{\mathcal{X}}(g)\mathbf{x})$$

A specific case of \mathbb{G} -equivariant maps are the \mathbb{G} -invariant ones, which are maps that commute with the group action and have trivial output group actions $\triangleright_{\mathcal{Y}}$ such that $\rho_{\mathcal{Y}}(g) = \mathbf{I}$ for all $g \in \mathbb{G}$. That is:

$$\mathbf{y} = g \triangleright_{\mathcal{Y}} f(\mathbf{x}) = f(g \triangleright_{\mathcal{X}} \mathbf{x}), \quad \forall \mathbf{x} \in \mathcal{X}, g \in \mathbb{G}. \quad \iff \quad \begin{array}{ccc} \mathcal{X} & \xrightarrow{\triangleright_{\mathcal{X}}} & \mathcal{X} \\ \searrow f & & \downarrow f \\ & & \mathcal{Y} \\ & & \curvearrowright \triangleright_{\mathcal{Y}} \end{array} \quad (40b)$$

$$\mathbf{y} = \rho_{\mathcal{Y}}(g)f(\mathbf{x}) = f(\rho_{\mathcal{X}}(g)\mathbf{x})$$

Structure of \mathbb{G} -equivariant linear maps

Definition I.12 (Homomorphism and Isomorphism). *Let \mathcal{X} and \mathcal{Y} be two vector spaces endowed with the same symmetry group \mathbb{G} , with the respective group actions $\triangleright_{\mathcal{X}}: \mathbb{G} \times \mathcal{X} \mapsto \mathcal{X}$ and $\triangleright_{\mathcal{Y}}: \mathbb{G} \times \mathcal{Y} \mapsto \mathcal{Y}$. The spaces are said to be \mathbb{G} -homomorphic if there exists a linear map $\mathbf{A} : \mathcal{X} \mapsto \mathcal{Y}$ that commutes with the group action, such that $g \triangleright_{\mathcal{Y}} (\mathbf{A}\mathbf{x}) = \mathbf{A}(g \triangleright_{\mathcal{X}} \mathbf{x})$ for all $\mathbf{x} \in \mathcal{X}$. They are said to be \mathbb{G} -isomorphic if the linear map is invertible. Graphically, \mathcal{X} and \mathcal{Y} are \mathbb{G} -homomorphic or \mathbb{G} -isomorphic if the following diagrams commute:*

$$\underbrace{\begin{array}{ccc} \mathcal{X} & \xrightarrow{\triangleright_{\mathcal{X}}} & \mathcal{X} \\ \downarrow \mathbf{A} & & \downarrow \mathbf{A} \\ \mathcal{Y} & \xrightarrow{\triangleright_{\mathcal{Y}}} & \mathcal{Y} \end{array}}_{\text{Homomorphism}} \quad \mathbf{A} \in \text{Homo}_{\mathbb{G}}(\mathcal{X}, \mathcal{Y}) \quad \text{or} \quad \underbrace{\begin{array}{ccc} \mathcal{X} & \xrightarrow{\triangleright_{\mathcal{X}}} & \mathcal{X} \\ \mathbf{A}^{-1} \updownarrow \mathbf{A} & & \mathbf{A}^{-1} \updownarrow \mathbf{A} \\ \mathcal{Y} & \xrightarrow{\triangleright_{\mathcal{Y}}} & \mathcal{Y} \end{array}}_{\text{Isomorphism}} \quad \mathbf{A} \in \text{Iso}_{\mathbb{G}}(\mathcal{X}, \mathcal{Y}). \quad (41)$$

Here, $\text{Homo}_{\mathbb{G}}(\mathcal{X}, \mathcal{Y})$ denotes the space of \mathbb{G} -equivariant linear maps between \mathcal{X} and \mathcal{Y} , and $\text{Iso}_{\mathbb{G}}(\mathcal{X}, \mathcal{Y})$ denotes the space of \mathbb{G} -equivariant invertible linear maps between \mathcal{X} and \mathcal{Y} .

Proposition I.13 (Structure of \mathbb{G} -homomorphisms / interwiners / \mathbb{G} -equivariant linear maps). *Let \mathbb{G} be a compact group and $\mathbf{A} \in \text{Homo}_{\mathbb{G}}(\mathcal{X}, \mathcal{Y})$ be a \mathbb{G} -equivariant linear map between two (real) \mathbb{G} -symmetric vector spaces \mathcal{X} and \mathcal{Y} , with isotypic decompositions:*

$$\mathcal{X} = \bigoplus_{k=1}^{n_{\text{iso}}} \mathcal{X}^{(k)} = \bigoplus_{k=1}^{n_{\text{iso}}} \bigoplus_{i=1}^{m_k^x} \mathcal{X}_i^{(k)} \quad \text{and} \quad \mathcal{Y} = \bigoplus_{k=1}^{n_{\text{iso}}} \mathcal{Y}^{(k)} = \bigoplus_{k=1}^{n_{\text{iso}}} \bigoplus_{j=1}^{m_k^y} \mathcal{Y}_j^{(k)},$$

where n_{iso} denotes the number of isotypic subspaces, and m_k^x and m_k^y denote the multiplicities of the irreducible representation $\bar{\rho}_k : \mathbb{G} \rightarrow \mathbb{GL}(\bar{\mathcal{V}}_k)$ in \mathcal{X} and \mathcal{Y} , respectively. Each $\mathcal{X}_i^{(k)}$ and $\mathcal{Y}_j^{(k)}$ is isometrically isomorphic to $\bar{\mathcal{V}}_k$ (see [Thm. I.16](#)). Hence, in the isotypic bases, the map \mathbf{A} decomposes block-diagonally into n_{iso} blocks corresponding to homomorphisms between isotypic subspaces of the same type, that is:

$$\mathbf{A} = \bigoplus_{k=1}^{n_{\text{iso}}} \mathbf{A}^{(k)} \quad \text{where} \quad \mathbf{A}^{(k)} \in \text{Homo}_{\mathbb{G}}(\mathcal{X}^{(k)}, \mathcal{Y}^{(k)}).$$

Furthermore, the map $\mathbf{A}^{(k)}$ decomposes into $m_k^x \times m_k^y$ blocks of endomorphisms of the irreducible subspace $\bar{\mathcal{V}}_k$. That is:

$$\mathbf{A}^{(k)} = \begin{bmatrix} \mathbf{A}_{1,1}^{(k)} & \cdots & \mathbf{A}_{1,m_k^x}^{(k)} \\ \vdots & \ddots & \vdots \\ \mathbf{A}_{m_k^y,1}^{(k)} & \cdots & \mathbf{A}_{m_k^y,m_k^x}^{(k)} \end{bmatrix} \quad \text{where} \quad \mathbf{A}_{i,j}^{(k)} \in \text{End}_{\mathbb{G}}(\bar{\mathcal{V}}_k), \forall i \in [1, m_k^y], j \in [1, m_k^x].$$

Consequently, depending of the type of irreducible representation $\mathbb{K} \in \{\mathbb{R}, \mathbb{C}, \mathbb{H}\}$, each sub-block is constrained to be in the span of the corresponding basis elements in [Eq. \(39\)](#). Consequently, if we denote by \mathbb{B} the basis set of \mathbb{K} , we have that the map $\mathbf{A}^{(k)}$ can be expressed in tensor product form as:

$$\mathbf{A}^{(k)} = \sum_{b \in \mathbb{B}} \Theta_b^{(k)} \otimes \Psi_k(b), \quad \text{where} \quad \Theta_b^{(k)} \in \mathbb{R}^{m_k^y \times m_k^x}, \Psi_k : \mathbb{B} \rightarrow \text{End}_{\mathbb{G}}(\bar{\mathcal{V}}_k) \quad (42)$$

With $[\Theta_b^{(k)}]_{i,j} = \langle \mathbf{A}_{i,j}^{(k)}, \Psi_k(b) \rangle$ denoting the basis expansion coefficient of the i -th, j -th endomorphism sub-block with the basis element $\Psi_k(b)$.

I.2 ISOTYPIC DECOMPOSITION AND DISENTANGLED REPRESENTATIONS

We now have all the necessary tools to decompose symmetric vector spaces into their fundamental building blocks: irreducible \mathbb{G} -stable subspaces ([Thm. I.6](#)). This decomposition forms the theoretical foundation for disentangled representations, a concept introduced by [Higgins et al. \(2018\)](#) in the context of group theory and generalized here within the framework of representation theory.

By Maschke's theorem ([Knapp, 1986](#)), we have that irreducible representations are the fundamental building blocks of the representations of a compact symmetry group \mathbb{G} , given that any group representation $\rho_{\mathcal{X}} : \mathbb{G} \rightarrow \mathbb{GL}(\mathcal{X})$ can be decomposed into a direct sum of irreducible representations, $\rho_{\mathcal{X}} \sim \bigoplus_{i=1}^{n_{\text{iso}}} \rho_{\mathcal{X}_i}$, where each $\rho_{\mathcal{X}_i}$ is isomorphic to one of the group's $n_{\text{iso}} \leq |\mathbb{G}|$ irreducible representations ([Thms. I.7](#) and [I.8](#)).

This decomposition will play a crucial role in facilitating numerical and theoretical analysis of operations on symmetric vector spaces. Therefore, we will frequently choose a convenient basis of the symmetric vector space which readily exposes this decomposition, termed isotypic basis. For the sake of generality we consider below the more general case of (finite and infinite dimensional) separable Hilbert spaces, which will enable us to extend these results to function spaces.

Definition I.14 (Isotypic Basis). *Let $\rho_{\mathcal{H}} : \mathbb{G} \rightarrow \mathbb{U}(\mathcal{H})$ be a unitary group representation of a compact group \mathbb{G} on a separable Hilbert Space \mathcal{H} . The representation is said to be defined in an isotypic basis if it is defined by a direct sum of irreducible representations grouped by their type, that is, if:*

$$\rho_{\mathcal{H}} = \bigoplus_{k=1}^{n_{\text{iso}}} \bigoplus_{p=1}^{m_k} \bar{\rho}_k, \quad (43)$$

where $\{\bar{\rho}_k : \mathbb{G} \rightarrow \mathbb{U}(\bar{\mathcal{H}}_k)\}_{k=1}^{n_{\text{iso}}}$ are the $n_{\text{iso}} \leq |\mathbb{G}|$ irreducible representations of \mathbb{G} , and $m_k \leq \infty$ is the multiplicity (i.e., number of copies) of the irrep type k in the representation $\rho_{\mathcal{H}}$.

Remark I.15. Note that multiple isotypic bases exist for a given representation $\rho_{\mathcal{H}}$, as both the irreducible ordering and each irrep's multiplicity ordering can be arbitrarily permuted.

The utility of an isotypic basis stems from Schur's orthogonality relations ([Thm. I.9](#)), which ensure that any symmetric vector space decomposes into at most n_{iso} orthogonal subspaces.

Theorem I.16 (Isotypic decomposition of symmetric Hilbert spaces ([Knapp, 1986](#))). *Let \mathbb{G} be a compact group and \mathcal{H} a separable Hilbert space with a unitary group representation $\rho_{\mathcal{H}} :$*

$\mathbb{G} \rightarrow \mathbb{U}(\mathcal{H})$. Then we can identify $n_{\text{iso}} \leq |\mathbb{G}|$ irreducible representations $\bar{\rho}_k : \mathbb{G} \rightarrow \mathbb{U}(\bar{\mathcal{H}}_k)$ that allow us to decompose \mathcal{H} into a sum of orthogonal subspaces, denoted isotypic subspaces: $\mathcal{H} = \bigoplus_{1 \leq k \leq n_{\text{iso}}}^{\perp} \mathcal{H}^{(k)}$ where each $\mathcal{H}^{(k)} = \bigoplus_{p=1}^{m_k} \mathcal{H}_p^{(k)}$ is the sum of at most $m_k \leq \infty$ countably many subspaces isometrically isomorphic to $\bar{\mathcal{H}}_k$.

Remark I.17. Note that the isotypic decomposition process can be simply interpreted as a change of basis to any isotypic basis (Thm. I.14) of the space.

Disentangled representations The concept of isotypic decomposition is intricately linked to the idea of disentangled representations, introduced by Higgins et al. (2018) in the representation learning literature, restated below for completeness.

Definition I.18 (Disentangled representation (Higgins et al. (2018))). A vector representation is called a disentangled representation with respect to a particular decomposition of a symmetry group into subgroups, if it decomposes into independent subspaces, where each subspace is affected by the action of a single subgroup, and the actions of all other subgroups leave the subspace unaffected.

Note that the independent subspaces of Thm. I.18 refer to the orthogonal isotypic subspaces $\{\mathcal{H}^{(k)}\}_{k=1}^{n_{\text{iso}}}$, each of which is acted upon by a unique quotient group⁵ defined by the kernel of the irrep acting on that isotypic subspace:

$$\mathbb{G}^{(k)} = \mathbb{G}/\mathbb{N}_k, \quad \text{where} \quad \mathbb{N}_k := \ker(\bar{\rho}_k) = \{g \in \mathbb{G} \mid \bar{\rho}_k(g) = \mathbf{I}_{d_k}\}. \quad (44)$$

Where each $\mathbb{G}^{(k)}$ is a well defined group of cosets generated by the normal subgroup \mathbb{N}_k . In practice, each $\mathbb{G}^{(k)}$ is isomorphic to the effective (matrix) group encoded by each irreducible representation $\bar{\rho}_k : \mathbb{G} \mapsto \mathbb{U}(\bar{\mathcal{H}}_k)$.

J REPRESENTATION THEORY OF SYMMETRIC FUNCTION SPACES

In this section, we study symmetry group actions on infinite-dimensional function spaces and specify the conditions needed to approximate these spaces in finite dimensions. Specifically, given a set \mathcal{X} with a compact symmetry group \mathbb{G} acting via (\triangleright) (Def. I.2), the space of scalar-valued functions on \mathcal{X} , $\mathcal{F} = \{f \mid f : \mathcal{X} \mapsto \mathbb{R}\}$, becomes a symmetric function space. The action of a symmetry transformation on a function is defined as:

Definition J.1 (Group action on a function space). Let \mathcal{X} be a set endowed with the symmetry group \mathbb{G} , and let \mathcal{F} be the space of scalar-valued functions on \mathcal{X} . The (left) action of \mathbb{G} on a function $f \in \mathcal{F}$ is defined as the composition of f with the inverse of the group element g^{-1} :

$$\begin{aligned} (\triangleright_{\mathcal{F}}) : \mathbb{G} \times \mathcal{F} &\longrightarrow \mathcal{F} \\ (g, f) &\longrightarrow [g \triangleright_{\mathcal{F}} f](\mathbf{x}) := [f \circ g^{-1}](\mathbf{x}) = f(g^{-1} \triangleright \mathbf{x}), \quad \forall \mathbf{x} \in \mathcal{X}. \end{aligned} \quad (45a)$$

In other words, the point-wise evaluation of f on a g^{-1} -transformed set \mathcal{X} is equivalent to the evaluation of the transformed function $g \triangleright_{\mathcal{F}} f \in \mathcal{F}$ on the original set \mathcal{X} (see simple examples in Fig. 16). Any function space that is stable under the group action Eq. (45a) is referred to as a symmetric function space. Note that this action is compatible with the group composition and identity element $e \in \mathbb{G}$, such that the following properties hold:

$$\text{Identity: } e \triangleright_{\mathcal{F}} f(\cdot) = f(\cdot), \quad (45b)$$

$$\text{Associativity: } [(g_2 \circ g_1) \triangleright_{\mathcal{F}} f](\cdot) = [g_2 \triangleright_{\mathcal{F}} [g_1 \triangleright_{\mathcal{F}} f]](\cdot), \quad \forall g_1, g_2 \in \mathbb{G}. \quad (45c)$$

Remark J.2. From an algebraic perspective, the inversion g^{-1} (contragredient representation) emerges to ensure that the associativity property of the group action (Eq. (45c)) holds:

$$\begin{aligned} [(g_2 \circ g_1) \triangleright_{\mathcal{F}} f](\mathbf{x}) &= [g_2 \triangleright_{\mathcal{F}} [g_1 \triangleright_{\mathcal{F}} f]](\mathbf{x}), \quad \forall \mathbf{x} \in \mathcal{X} \\ f((g_2 \circ g_1)^{-1} \triangleright \mathbf{x}) &= [g_1 \triangleright_{\mathcal{F}} f](g_2^{-1} \triangleright \mathbf{x}) = f(g_1^{-1} \triangleright (g_2^{-1} \triangleright \mathbf{x})) \\ f((g_2 \circ g_1)^{-1} \triangleright \mathbf{x}) &= f((g_1 \circ g_2)^{-1} \triangleright \mathbf{x}). \end{aligned}$$

In the context of this work, we will study the scenario where the function space \mathcal{F} is a separable Hilbert space and the group action of \mathbb{G} on \mathcal{F} is unitary, i.e., it preserves the inner product of the function space. This setup is crucial to enable us to approximate \mathcal{F} and the group action on \mathcal{F} in finite dimensions.

⁵The original definition of disentangled representations refers to subgroups, but in general the quotient group $\mathbb{G}^{(k)}$ need not be a **subgroup** of \mathbb{G} .

J.1 UNITARY GROUP REPRESENTATION ON FUNCTION SPACES

Assume our symmetric set \mathcal{X} is endowed with a measure space structure $(\mathcal{X}, \Sigma_{\mathcal{X}}, P_{\mathbf{x}})$, where $P_{\mathbf{x}} : \Sigma_{\mathcal{X}} \mapsto \mathbb{R}$ is the space measure. Then, consider a function space with a separable Hilbert space structure $\mathcal{F} := \mathcal{L}_{P_{\mathbf{x}}}^2 \mathcal{X}, \mathbb{R}$, and inner product $\langle f_1, f_2 \rangle_{P_{\mathbf{x}}} = \int_{\mathcal{X}} f_1(\mathbf{x}) f_2(\mathbf{x}) P_{\mathbf{x}}(d\mathbf{x})$ for all $f_1, f_2 \in \mathcal{F}$. Then, the action $\triangleright_{\mathcal{F}}$ of the group \mathbb{G} on the function space \mathcal{F} is termed unitary if it preserves the inner product of the function space:

$$\begin{aligned} \langle f_1, f_2 \rangle_{P_{\mathbf{x}}} &= \langle g \triangleright_{\mathcal{F}} f_1, g \triangleright_{\mathcal{F}} f_2 \rangle_{P_{\mathbf{x}}} \quad \forall f_1, f_2 \in \mathcal{F}, g \in \mathbb{G} \\ \int_{\mathcal{X}} f_1(\mathbf{x}) f_2(\mathbf{x}) P_{\mathbf{x}}(d\mathbf{x}) &= \int_{\mathcal{X}} (g \triangleright_{\mathcal{F}} f_1)(\mathbf{x}) (g \triangleright_{\mathcal{F}} f_2)(\mathbf{x}) P_{\mathbf{x}}(d\mathbf{x}) \\ &= \int_{\mathcal{X}} f_1(g^{-1} \triangleright \mathbf{x}) f_2(g^{-1} \triangleright \mathbf{x}) P_{\mathbf{x}}(d\mathbf{x}) \\ &= \int_{g \triangleright \mathcal{X} = \mathcal{X}} f_1(\mathbf{x}) f_2(\mathbf{x}) P_{\mathbf{x}}(g \triangleright d\mathbf{x}). \end{aligned} \quad (46)$$

That is, the group action is unitary if $P_{\mathbf{x}}$ is a \mathbb{G} -invariant measure $P_{\mathbf{x}}(g \triangleright d\mathbf{x}) = P_{\mathbf{x}}(d\mathbf{x})$, $\forall g \in \mathbb{G}, d\mathbf{x} \subseteq \mathcal{X}$. Note that an \mathbb{G} -invariant measure (and inner product) exists whenever \mathbb{G} is finite, because for any measure $\eta : \Sigma_{\Omega} \mapsto \mathbb{R}$, we can use the group-average trick to obtain one, given by $P_{\mathbf{x}}(\mathbb{X}) = \Sigma_{g \in \mathbb{G}} \eta(g \triangleright \mathbb{X})$.⁶

The importance of the Hilbert space structure is that it enables the definition of a unitary group representation. Unitary representations have a well-studied modular structure that allows their decomposition (Thm. I.16) into \mathbb{G} -stable subspaces (Thm. I.6), which is crucial for approximating symmetric function spaces using a finite set of basis elements. Let $\mathbb{I}_{\mathcal{F}} = \{\phi_i \mid \phi_i \in \mathcal{L}_{P_{\mathbf{x}}}^2\}_{i \in \mathbb{N}}$ be an orthogonal basis for the function space $\mathcal{F} = \text{span}(\mathbb{I}_{\mathcal{F}})$, so that any function $f \in \mathcal{F}$ can be represented by its basis expansion coefficients $\alpha = [\langle \phi_i \rangle_{P_{\mathbf{x}}} f]_{i \in \mathbb{N}}$, since $f_{\alpha}(\mathbf{x}) = \sum_{i \in \mathbb{N}} \langle \phi_i, f \rangle_{P_{\mathbf{x}}} \phi_i(\mathbf{x})$. In this basis, the group action of \mathbb{G} on \mathcal{F} defines a unitary group representation mapping group elements to unitary linear integral operators on \mathcal{F} , which can be expressed in matrix form.

⁶Such a \mathbb{G} -invariant measure exists for any (finite or continuous) compact group. See discussion.

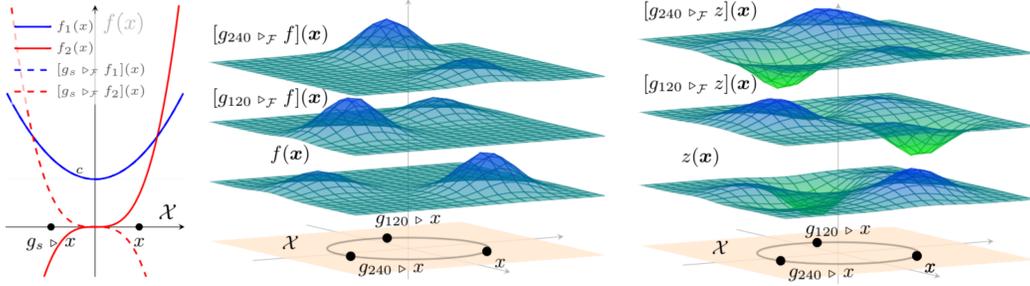


Figure 16: **Left:** Diagram of the group action $\triangleright_{\mathcal{F}}$ on functions $f_1(x) = x^2 + c$ and $f_2(x) = -x^3$ defined on the domain $\mathcal{X} := \mathbb{R}$ endowed with the reflectional symmetry group $\mathbb{G} := \mathbb{C}_2 = \{e, g_s\}$, with the reflection action acting on the domain by $g_s \triangleright x = -x$ and on the function space $\mathcal{F} := \{f \mid f : \mathcal{X} \mapsto \mathbb{R}\}$ by $[g \triangleright_{\mathcal{F}} f](\mathbf{x}) = f(g \triangleright_{\mathcal{X}} \mathbf{x}) = f(-\mathbf{x})$. Hence we have that f_1 is a \mathbb{G} -invariant function, $g_s \triangleright_{\mathcal{F}} f_1(x) = f_1(x)$ and f_2 a \mathbb{G} -equivariant function $g_s \triangleright_{\mathcal{F}} f_2(x) = -x^3$. **Center:** Diagram representing the action $\triangleright_{\mathcal{F}}$ on the (arbitrarily chosen) function $f(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \mathbf{c}_1, 2) + \mathcal{N}(\mathbf{x}; \mathbf{c}_2, 1)$ defined over the symmetric domain $\mathcal{X} = \mathbb{R}^2$ with the cyclic symmetry group $\mathbb{G} = \mathbb{C}_3 = \{e, g_{120}, g_{240}\}$ and group action $g \triangleright \mathbf{x} = \rho_{\mathcal{X}}(g)\mathbf{x} = \mathbf{R}_g \mathbf{x}$, where \mathbf{R}_g is a rotation matrix in 2D. Here, $g_{120} \triangleright_{\mathcal{F}} f$ is equivalent to evaluating f on a domain rotated by -120° . The same holds for $g_{240} \triangleright_{\mathcal{F}} f$. Note that the z -offsets are added for visualization purposes. **Right:** Diagram representing the action $\triangleright_{\mathcal{F}}$ on the function $z \in \widehat{\mathcal{F}}$, defined to be a member of the finite-dimensional symmetric function space $\widehat{\mathcal{F}} := \text{span}(\mathbb{I}_{\widehat{\mathcal{F}}})$, constructed from a basis set composed of the group orbit of the (arbitrarily chosen) function $f \in \mathcal{F}$, that is $\mathbb{I}_{\widehat{\mathcal{F}}} := \mathbb{G}f = \{f, g_{120} \triangleright_{\mathcal{F}} f, g_{240} \triangleright_{\mathcal{F}} f\}$. This function space is \mathbb{G} -stable by construction, since $\mathbb{G}\mathbb{I}_{\widehat{\mathcal{F}}} = \mathbb{I}_{\widehat{\mathcal{F}}}$. Note that the z -offsets are added for visualization purposes.

Definition J.3 (Unitary group representation on a function space). Let $\mathcal{F} = \mathcal{L}_{P_x}^2 \mathcal{X}, \mathbb{R}$ be a separable Hilbert space of scalar-valued functions on a set \mathcal{X} endowed with the symmetry group \mathbb{G} . Let $\mathbb{I}_{\mathcal{F}}$ be an orthogonal basis set spanning \mathcal{F} . Then, the group action of \mathbb{G} on \mathcal{F} (Thm. J.1) defines a unitary group representation mapping group elements to unitary linear integral operators on \mathcal{F} :

$$\rho_{\mathcal{F}} : \begin{array}{ccc} \mathbb{G} & \longrightarrow & \mathbb{U}(\mathcal{F}) \\ g & \longrightarrow & \rho_{\mathcal{F}}(g), \end{array} \quad \text{s.t.} \quad \rho_{\mathcal{F}}(g)^* = \rho_{\mathcal{F}}(g^{-1}). \quad (47)$$

Each unitary operator $\rho_{\mathcal{F}}(g) : \mathcal{F} \mapsto \mathcal{F}$ admits an infinite-dimensional matrix representation with entries $[\rho_{\mathcal{F}}(g)]_{i,j} := \langle \hat{f}_i, g \triangleright_{\mathcal{F}} \hat{f}_j \rangle_{P_x}$, which characterize how the group action transforms the chosen basis functions. Consequently, once the group representation for a chosen basis set is defined, the group action on a function $f_{\alpha} \in \mathcal{F}$ can be expressed as an (infinite-dimensional) matrix transformation of its basis expansion coefficients α , given by:

$$[g \triangleright_{\mathcal{F}} f_{\alpha}](\cdot) := \sum_{i \in \mathbb{N}} \langle \hat{f}_i, g \triangleright_{\mathcal{F}} f_{\alpha} \rangle_{P_x} \hat{f}_i(\cdot) = \sum_{i \in \mathbb{N}} \left(\sum_{j \in \mathbb{N}} \langle \hat{f}_i, g \triangleright_{\mathcal{F}} \hat{f}_j \rangle_{P_x} \underbrace{\langle \hat{f}_j, f \rangle_{P_x}}_{\alpha_j} \right) \hat{f}_i(\cdot). \quad (48)$$

Example J.4 (Isotypic decomposition of symmetric function space). Let $(\mathcal{X}, \Sigma_{\mathcal{X}}, P_x)$ be a symmetric 2D measure space with domain $\mathcal{X} \sim \mathbb{R}^2$ and cyclic symmetry group $\mathbb{G} := \mathbb{C}_3 = \{e, g_{120}, g_{240}\}$, acting on the 2D plane by 120° rotations (Fig. 17). Define the finite-dimensional function space $\mathcal{F}_x \subset \mathcal{L}_x^2$ with basis $\mathbb{I}_{\mathcal{F}_x} = \{\phi, g_{120} \triangleright \phi, g_{240} \triangleright \phi\}$, where $\phi \in \mathcal{F}_x$ is an arbitrary measurable function (Fig. 17-left). In this basis, the group action $\triangleright_{\mathcal{F}_x}$ for any function $z_{\alpha} \in \mathcal{F}_x$ is given by the regular representation $\rho_{\mathcal{F}_x} = \rho_{\text{reg}}$ acting on the coefficient vector $\alpha \in \mathbb{R}^3$ (Fig. 6-right).

$$[g \triangleright_{\mathcal{F}_x} z_{\alpha}](\cdot) = \sum_{i=1}^3 \langle \phi_i, g \triangleright_{\mathcal{F}_x} z_{\alpha} \rangle_{P_x} \phi_i(\cdot) \equiv (\rho_{\text{reg}}(g)\alpha)^{\top} \begin{bmatrix} \phi(\cdot) \\ g_{120} \triangleright \phi(\cdot) \\ g_{240} \triangleright \phi(\cdot) \end{bmatrix}, \quad \rho_{\text{reg}}(g) = \begin{cases} \mathbf{I}_3, & \text{if } g = e \\ \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}, & \text{if } g = g_{120} \\ \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}, & \text{if } g = g_{240} \end{cases} \quad (49)$$

The group \mathbb{C}_3 possesses two types of (real-valued) irreducible representations, $n_{\text{iso}} = 2$: the trivial irreducible representation $\bar{\rho}_{\text{inv}}$ and a 2D rotation representation $\bar{\rho}_{2\pi/3}$, defined by:

$$\bar{\rho}_{\text{inv}}(g) = \mathbf{I}_1, \forall g \in \mathbb{C}_3, \quad \text{and} \quad \bar{\rho}_{2\pi/3}(g) = \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix}, \quad \text{s.t. } \theta = \begin{cases} 0^\circ, & \text{if } g = e \\ 120^\circ & \text{if } g = g_{120} \\ 240^\circ & \text{if } g = g_{240} \end{cases} \quad (50)$$

Applying the appropriate change of basis, we decompose the regular representation into a direct sum of the group's irreducible representations: $\rho_{\text{reg}} = \mathbf{Q}(\bar{\rho}_{\text{inv}} \oplus \bar{\rho}_{2\pi/3})\mathbf{Q}^{-1}$, where \mathbf{Q} transitions from the regular basis to the isotypic basis of \mathcal{F}_x . Since \mathbb{C}_3 is abelian, \mathbf{Q} corresponds to the linear map defining the Fourier transform.

By Thm. I.16, this results in the orthogonal decomposition of the finite-dimensional function space into two orthogonal subspaces; $\mathcal{F}_x = \mathcal{F}_x^{\text{inv}} \oplus^{\perp} \mathcal{F}_x^{(2)}$, where $\mathcal{F}_x^{\text{inv}}$ denotes the 1-dimensional subspace of \mathbb{G} -invariant functions ($\mathbb{G}^{(0)} = \{e\}$), and $\mathcal{F}_x^{(2)}$ is the 2-dimensional subspace with group actions defined by the 2D irreducible representation $\bar{\rho}_{2\pi/3}$ ($\mathbb{G}^{(1)} = \mathbb{C}_3$). We can construct the basis set in the isotypic basis given:

$$\mathbb{I}_{\mathcal{F}_x}^{\text{iso}} = \mathbf{Q} \begin{bmatrix} \phi(\cdot) \\ g_{120} \triangleright \phi(\cdot) \\ g_{240} \triangleright \phi(\cdot) \end{bmatrix} = \begin{bmatrix} u^{\text{inv}}(\cdot) \\ u_1^{(2)}(\cdot) \\ u_2^{(2)}(\cdot) \end{bmatrix} \quad \text{s.t. } \mathbf{Q} = \begin{bmatrix} 1/\sqrt{3} & 1/\sqrt{3} & 1/\sqrt{3} \\ 2/\sqrt{6} & -1/\sqrt{6} & -1/\sqrt{6} \\ 0 & 1/\sqrt{2} & -1/\sqrt{2} \end{bmatrix} \quad (51)$$

The new basis functions in the isotypic basis are depicted in Fig. 17-right, and elucidate that the symmetry constraints on this 3-dimensional function space, result in $m = 2$ unique functions, each associated with a unique irreducible representation.

Assuming P_x is a \mathbb{G} -invariant probability measure, we compute the expected value of each basis function. In the regular basis, functions related by a symmetry transformation share the same expected value, i.e., $\mathbb{E}_x \phi = \mathbb{E}_x g \triangleright \phi$ for all $g \in \mathbb{C}_3$. In the isotypic basis, functions lacking a \mathbb{G} -invariant component (i.e., $u_1^{(2)}, u_2^{(2)}$) are centered: $\mathbb{E}_x u_1^{(2)} = \mathbb{E}_x u_2^{(2)} = 0$. In our example this constraint becomes clear from the nature of the change of basis \mathbf{Q} . Eq. (51).

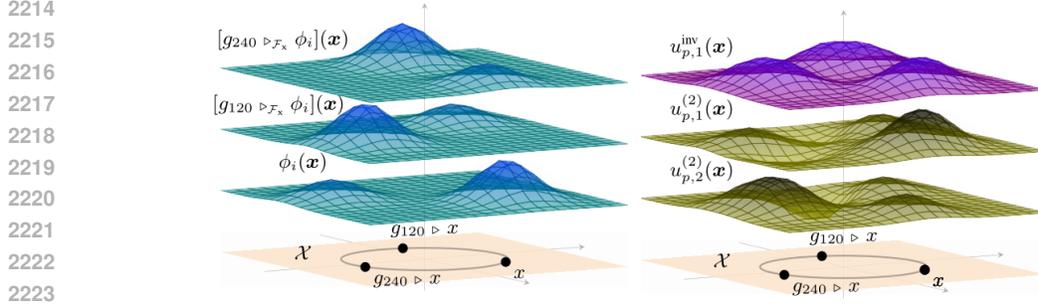


Figure 17: Visualization of the basis functions in the finite-dimensional symmetric function space \mathcal{F}_x from [Thm. J.4](#). **Left:** Depiction of the basis functions in the regular basis $\mathbb{I}_{\mathcal{F}_x} = \{\phi, g_{120} \triangleright \phi, g_{240} \triangleright \phi\}$, generated by the action of the cyclic group \mathbb{C}_3 on an arbitrary function $\phi \in \mathcal{F}_x$. **Right:** Depiction of the basis functions in the isotypic basis $\mathbb{I}_{\mathcal{F}_x}^{\text{iso}} = \{u^{\text{inv}}, u_1^{(2)}, u_2^{(2)}\}$, obtained via the change of basis matrix Q . The first basis function u^{inv} corresponds to the \mathbb{G} -invariant subspace $\mathcal{F}_x^{\text{inv}}$ and is visually invariant under the action of \mathbb{C}_3 on \mathcal{X} . The other two basis functions $u_1^{(2)}, u_2^{(2)}$ are constrained to span a \mathbb{G} -stable subspace of \mathcal{L}_x^2 , denoted by $\mathcal{F}_x^{(2)}$ that transform according to the irreducible representation $\bar{\rho}_{2\pi/3}$. Meaning for any function $f \in \mathcal{F}_x$, the group action $g \triangleright_{\mathcal{F}_x} f$ can be computed by a linear transformation of its basis expansion coefficients.

K \mathbb{G} -EQUIVARIANT LINEAR INTEGRAL OPERATORS

This section gives an overview of \mathbb{G} -equivariant linear integral operators between symmetric function spaces. We define these operators, discuss their properties, and specify conditions under which they commute with group actions. In [Subsec. K.1](#) we examine their infinite-dimensional matrix form and the resulting algebraic constraints from \mathbb{G} -equivariance. In [Subsec. K.2](#) we then show how to exploit these constraints in a finite-rank approximation.

Let \mathbb{G} be a compact group acting on two measure spaces $(\mathcal{X}, \Sigma_{\mathcal{X}}, P_{\mathbf{x}})$ and $(\mathcal{Y}, \Sigma_{\mathcal{Y}}, P_{\mathbf{y}})$ via the group actions $\triangleright_{\mathcal{X}}$ and $\triangleright_{\mathcal{Y}}$ (see [Def. I.2](#)). Assume that the measures $P_{\mathbf{x}}$ and $P_{\mathbf{y}}$ are \mathbb{G} -invariant, i.e., $P_{\mathbf{x}}(g \triangleright_{\mathcal{X}} \mathbb{B}) = P_{\mathbf{x}}(\mathbb{B})$ and $P_{\mathbf{y}}(g \triangleright_{\mathcal{Y}} \mathbb{A}) = P_{\mathbf{y}}(\mathbb{A})$ for all $g \in \mathbb{G}$, $\mathbb{B} \in \Sigma_{\mathcal{X}}$, and $\mathbb{A} \in \Sigma_{\mathcal{Y}}$ (see [Thm. I.11](#)).

Let $\mathcal{L}_{\mathbf{x}}^2 = \{f : \mathcal{X} \mapsto \mathbb{R} \mid \|f\|_{P_{\mathbf{x}}} < +\infty\}$ and $\mathcal{L}_{\mathbf{y}}^2 = \{h : \mathcal{Y} \mapsto \mathbb{R} \mid \|h\|_{P_{\mathbf{y}}} < +\infty\}$ be the Hilbert spaces of square-integrable functions with respect to $P_{\mathbf{x}}$ and $P_{\mathbf{y}}$, respectively. Since \mathcal{X} and \mathcal{Y} have a \mathbb{G} -action, the spaces $\mathcal{L}_{\mathbf{x}}^2$ and $\mathcal{L}_{\mathbf{y}}^2$ inherit group actions defined by $[g \triangleright_{\mathcal{L}_{\mathbf{x}}^2} f](\mathbf{x}) = f(g^{-1} \triangleright_{\mathcal{X}} \mathbf{x})$, $[g \triangleright_{\mathcal{L}_{\mathbf{y}}^2} h](\mathbf{y}) = h(g^{-1} \triangleright_{\mathcal{Y}} \mathbf{y})$, for all $f \in \mathcal{L}_{\mathbf{x}}^2$ and $h \in \mathcal{L}_{\mathbf{y}}^2$ (see [Thm. J.1](#)).

We consider linear integral operators $\mathbb{T} : \mathcal{L}_{\mathbf{x}}^2 \mapsto \mathcal{L}_{\mathbf{y}}^2$ defined by

$$h(\mathbf{y}) = [\mathbb{T}f](\mathbf{y}) = \int_{\mathcal{X}} \kappa(\mathbf{x}, \mathbf{y}) f(\mathbf{x}) P_{\mathbf{x}}(d\mathbf{x}), \quad (52)$$

where $k : \mathcal{X} \times \mathcal{Y} \mapsto \mathbb{R}$ is the kernel function of \mathbb{T} . In this work we focus on those operators whose kernels are \mathbb{G} -invariant such operators are called \mathbb{G} -equivariant.

Definition K.1 (\mathbb{G} -equivariant linear integral operators). *Let $(\mathcal{X}, \Sigma_{\mathcal{X}}, P_{\mathbf{x}})$ and $(\mathcal{Y}, \Sigma_{\mathcal{Y}}, P_{\mathbf{y}})$ be two measure spaces endowed with group actions $\triangleright_{\mathcal{X}}$ and $\triangleright_{\mathcal{Y}}$ and \mathbb{G} -invariant measures $P_{\mathbf{x}}$ and $P_{\mathbf{y}}$ for a given compact symmetry group \mathbb{G} . Let $\mathbb{T} : \mathcal{L}_{\mathbf{x}}^2 \mapsto \mathcal{L}_{\mathbf{y}}^2$ be a linear integral operator between the spaces of square-integrable functions defined on the two measure spaces. The operator \mathbb{T} is said to*

be \mathbb{G} -equivariant if it commutes with the group action, that is $\forall f \in \mathcal{L}_{\mathbf{x}}^2, g \in \mathbb{G}$ and $\mathbf{y} \in \mathcal{Y}$:

$$[\mathbb{T}[g \triangleright_{\mathcal{L}_{\mathbf{x}}^2} f]](\mathbf{y}) = [g \triangleright_{\mathcal{L}_{\mathbf{y}}^2} [\mathbb{T}f]](\mathbf{y}) \quad (53a)$$

$$\int_{\mathcal{X}} \kappa(\mathbf{x}, \mathbf{y}) f(g^{-1} \triangleright_{\mathcal{X}} \mathbf{x}) P_{\mathbf{x}}(d\mathbf{x}) = g \triangleright_{\mathcal{L}_{\mathbf{y}}^2} \left(\int_{\mathcal{X}} \kappa(\mathbf{x}, \mathbf{y}) f(\mathbf{x}) P_{\mathbf{x}}(d\mathbf{x}) \right)$$

$$\int_{\mathcal{X}} \kappa(g \triangleright_{\mathcal{X}} \mathbf{x}, \mathbf{y}) f(\mathbf{x}) P_{\mathbf{x}}(g \triangleright_{\mathcal{X}} d\mathbf{x}) = \int_{\mathcal{X}} \kappa(\mathbf{x}, g^{-1} \triangleright_{\mathcal{Y}} \mathbf{y}) f(\mathbf{x}) P_{\mathbf{x}}(d\mathbf{x}) \quad \text{s.t. } g \triangleright_{\mathcal{X}} \mathcal{X} := \mathcal{X}$$

$$\int_{\mathcal{X}} \kappa(g \triangleright_{\mathcal{X}} \mathbf{x}, \mathbf{y}) f(\mathbf{x}) P_{\mathbf{x}}(d\mathbf{x}) = \int_{\mathcal{X}} \kappa(\mathbf{x}, g^{-1} \triangleright_{\mathcal{Y}} \mathbf{y}) f(\mathbf{x}) P_{\mathbf{x}}(d\mathbf{x}) \quad \text{s.t. } P_{\mathbf{x}}(g \triangleright_{\mathcal{X}} d\mathbf{x}) = P_{\mathbf{x}}(d\mathbf{x})$$

$$k(g \triangleright_{\mathcal{X}} \mathbf{x}, \mathbf{y}) = \kappa(\mathbf{x}, g^{-1} \triangleright_{\mathcal{Y}} \mathbf{y}) \iff k(g \triangleright_{\mathcal{X}} \mathbf{x}, g \triangleright_{\mathcal{Y}} \mathbf{y}) = \kappa(\mathbf{x}, \mathbf{y}). \quad (53b)$$

Notice that the \mathbb{G} -equivariance of the operator \mathbb{T} is linked to the \mathbb{G} -invariance of its kernel function, which is required to satisfy Eq. (53b).

Multiple approaches exist to parameterize and approximate linear integral operators with finite resources (Kovachki et al., 2023, sec. 4). Here, we assume that both the input and output function spaces are separable Hilbert spaces, so that the operator can be represented as an infinite-dimensional matrix once appropriate basis sets are chosen. Its finite-dimensional (truncated or finite-rank) approximation is then obtained by selecting a finite number of basis functions in each space.

K.1 INFINITE-DIMENSIONAL MATRIX FORM OF THE OPERATOR

Since $\mathcal{L}_{\mathbf{x}}^2$ and $\mathcal{L}_{\mathbf{y}}^2$ are Hilbert spaces with inner products $\langle \cdot, \cdot \rangle_{P_{\mathbf{x}}}$ and $\langle \cdot, \cdot \rangle_{P_{\mathbf{y}}}$ respectively, we can choose orthogonal bases for both spaces: $\mathbb{I}_{\mathcal{L}_{\mathbf{x}}^2} = \{\phi_i \mid \phi_i \in \mathcal{L}_{\mathbf{x}}^2\}_{i \in \mathbb{N}}$ and $\mathbb{I}_{\mathcal{L}_{\mathbf{y}}^2} = \{\psi_j \mid \psi_j \in \mathcal{L}_{\mathbf{y}}^2\}_{j \in \mathbb{N}}$. This choice allows any function $f \in \mathcal{L}_{\mathbf{x}}^2$ and $h \in \mathcal{L}_{\mathbf{y}}^2$ to be represented by their infinite-dimensional coefficient vectors $\boldsymbol{\alpha} = [\langle \phi_i, f \rangle_{P_{\mathbf{x}}}]_{i \in \mathbb{N}}$ and $\boldsymbol{\beta} = [\langle \psi_j, h \rangle_{P_{\mathbf{y}}}]_{j \in \mathbb{N}}$, so that:

$$f(\mathbf{x}) := f_{\boldsymbol{\alpha}}(\mathbf{x}) = \sum_{i=1}^{\infty} \langle \phi_i, f \rangle_{P_{\mathbf{x}}} \phi_i(\mathbf{x}) \equiv \boldsymbol{\alpha}^T \boldsymbol{\phi}(\mathbf{x}) \quad h(\mathbf{y}) := h_{\boldsymbol{\beta}}(\mathbf{y}) = \sum_{j=1}^{\infty} \langle \psi_j, h \rangle_{P_{\mathbf{y}}} \psi_j(\mathbf{y}) \equiv \boldsymbol{\beta}^T \boldsymbol{\psi}(\mathbf{y}) \quad (54)$$

Here, $\boldsymbol{\alpha}^T \boldsymbol{\phi}(\mathbf{x})$ and $\boldsymbol{\beta}^T \boldsymbol{\psi}(\mathbf{y})$ represent the function as the dot product of its expansion coefficients with the basis evaluations $\boldsymbol{\phi}(\mathbf{x}) = [\phi_i(\mathbf{x})]_{i \in \mathbb{N}}$ and $\boldsymbol{\psi}(\mathbf{y}) = [\psi_j(\mathbf{y})]_{j \in \mathbb{N}}$. This notation is useful when we later select a finite number of basis functions to form a finite-dimensional approximation of \mathbb{T} .

With the chosen bases, the action of a linear integral operator $\mathbb{T} : \mathcal{L}_{\mathbf{y}}^2 \rightarrow \mathcal{L}_{\mathbf{x}}^2$ on any $f \in \mathcal{L}_{\mathbf{x}}^2$ is determined by its action on the basis functions:

$$[\mathbb{T}f_{\boldsymbol{\alpha}}](\mathbf{y}) = \int_{\mathcal{X}} \kappa(\mathbf{x}, \mathbf{y}) \left(\sum_{i \in \mathbb{N}} \alpha_i \phi_i(\mathbf{x}) \right) P_{\mathbf{x}}(d\mathbf{x}) = \sum_{i \in \mathbb{N}} \alpha_i \int_{\mathcal{X}} \kappa(\mathbf{x}, \mathbf{y}) \phi_i(\mathbf{x}) P_{\mathbf{x}}(d\mathbf{x}) = \sum_{i \in \mathbb{N}} \alpha_i [\mathbb{T}\phi_i](\mathbf{y}) \quad (55)$$

Since $[\mathbb{T}\phi_i] \in \mathcal{L}_{\mathbf{y}}^2$, each $[\mathbb{T}\phi_i](\mathbf{y})$ can be expanded using the output basis as $[\mathbb{T}\phi_i](\mathbf{y}) = \sum_{j \in \mathbb{N}} \langle \psi_j, [\mathbb{T}\phi_i] \rangle_{P_{\mathbf{y}}} \psi_j(\mathbf{y})$. Thus, the operator \mathbb{T} can be represented by the infinite-dimensional matrix \mathbf{T} with entries $\mathbf{T}_{ij} = \langle \psi_i, [\mathbb{T}\phi_j] \rangle_{P_{\mathbf{y}}}$. Therefore, the action of \mathbb{T} on any $f_{\boldsymbol{\alpha}} \in \mathcal{L}_{\mathbf{x}}^2$ is given by the matrix-vector product $\boldsymbol{\beta} = \mathbf{T}\boldsymbol{\alpha}$, i.e.,

$$\begin{aligned} [\mathbb{T}f_{\boldsymbol{\alpha}}](\mathbf{y}) &= \sum_{j \in \mathbb{N}} \alpha_j [\mathbb{T}\phi_j](\mathbf{y}) = \sum_{j \in \mathbb{N}} \alpha_j \sum_{i \in \mathbb{N}} \langle \psi_i, [\mathbb{T}\phi_j] \rangle_{P_{\mathbf{y}}} \psi_i(\mathbf{y}) \\ &= \sum_{i \in \mathbb{N}} \sum_{j \in \mathbb{N}} \mathbf{T}_{ij} \alpha_j \psi_i(\mathbf{y}) \equiv (\mathbf{T}\boldsymbol{\alpha})^T \boldsymbol{\psi}(\mathbf{y}) \end{aligned} \quad (56)$$

Eq. (56) shows that knowing the action of \mathbb{T} on the bases $\mathbb{I}_{\mathcal{L}_{\mathbf{x}}^2}$ and $\mathbb{I}_{\mathcal{L}_{\mathbf{y}}^2}$ determines its action on any function in $\mathcal{L}_{\mathbf{x}}^2$. In the sections that follow, we describe how symmetry constrains this action by requiring the bases to be \mathbb{G} -stable and by imposing \mathbb{G} -equivariance on \mathbf{T} , thereby introducing exploitable algebraic constraints for improved finite-rank approximations.

K.1.1 \mathbb{G} -EQUIVARIANT MATRIX FORM OF THE OPERATOR

Whenever the function spaces carry a symmetry group \mathbb{G} , the group action on their bases $\mathbb{I}_{\mathcal{L}_{\mathbf{x}}^2}$ and $\mathbb{I}_{\mathcal{L}_{\mathbf{y}}^2}$ is defined by the unitary representations $\boldsymbol{\rho}_{\mathcal{L}_{\mathbf{x}}^2} : \mathbb{G} \rightarrow \mathbb{U}(\mathcal{L}_{\mathbf{x}}^2)$ and $\boldsymbol{\rho}_{\mathcal{L}_{\mathbf{y}}^2} : \mathbb{G} \rightarrow \mathbb{U}(\mathcal{L}_{\mathbf{y}}^2)$ (see Thm. J.3).

As in Eq. (56), these representations can be expressed in (infinite-dimensional) matrix form so that the group action is given by a matrix-vector product:

$$\begin{aligned} [g \triangleright_{\mathcal{L}_x^2} f_\alpha](\cdot) &\equiv (\rho_{\mathcal{L}_x^2}(g)\alpha)^T \phi(\cdot), \quad \forall f_\alpha \in \mathcal{L}_x^2, g \in \mathbb{G} \\ [g \triangleright_{\mathcal{L}_y^2} h_\beta](\cdot) &\equiv (\rho_{\mathcal{L}_y^2}(g)\beta)^T \psi(\cdot), \quad \forall h_\beta \in \mathcal{L}_y^2, g \in \mathbb{G} \end{aligned} \quad (57)$$

Since the operator \mathbb{T} is \mathbb{G} -equivariant by construction (Eq. (53a)), the matrix form \mathbf{T} of the operator must also be \mathbb{G} -equivariant with respect to the group representations $\rho_{\mathcal{L}_x^2}$ and $\rho_{\mathcal{L}_y^2}$:

$$\begin{aligned} [\mathbb{T}[g \triangleright_{\mathcal{L}_x^2} f_\alpha]](\mathbf{y}) &= [g \triangleright_{\mathcal{L}_y^2} [\mathbb{T}f_\alpha]](\mathbf{y}) \quad \forall f_\alpha \in \mathcal{L}_x^2, g \in \mathbb{G}, \mathbf{y} \in \mathcal{Y} \\ (\mathbf{T}\rho_{\mathcal{L}_x^2}(g)\alpha)^T \psi(\mathbf{y}) &= (\rho_{\mathcal{L}_y^2}(g)\mathbf{T}\alpha)^T \psi(\mathbf{y}) \quad \text{s.t. Eqs. (56) and (57)} \\ \mathbf{T}\rho_{\mathcal{L}_x^2}(g) &= \rho_{\mathcal{L}_y^2}(g)\mathbf{T} \end{aligned} \quad (58)$$

With bases $\mathbb{I}_{\mathcal{L}_x^2}$ and $\mathbb{I}_{\mathcal{L}_y^2}$ for \mathcal{L}_x^2 and \mathcal{L}_y^2 , the kernel (Thm. K.1) can be written as $\kappa(\mathbf{x}, \mathbf{y}) = \sum_{i,j \in \mathbb{N}} \mathbf{T}_{i,j} \phi_j(\mathbf{x}) \psi_i(\mathbf{y})$. Hence, the \mathbb{G} -invariance condition (Eq. (53b)) on the kernel directly implies that the matrix \mathbf{T} is \mathbb{G} -equivariant, as stated in the following proposition:

Proposition K.2 (\mathbb{G} -invariant kernel implies \mathbb{G} -equivariant matrix form). *Let $\mathbb{T} : \mathcal{L}_x^2 \mapsto \mathcal{L}_y^2$ be a \mathbb{G} -equivariant operator between symmetric function spaces endowed with the group actions $\triangleright_{\mathcal{L}_x^2}$ and $\triangleright_{\mathcal{L}_y^2}$ of a compact symmetry group \mathbb{G} . Let $\rho_{\mathcal{L}_x^2}$ and $\rho_{\mathcal{L}_y^2}$ be the group representation of the on the input/output function spaces on the chosen basis sets $\mathbb{I}_{\mathcal{L}_x^2}$ and $\mathbb{I}_{\mathcal{L}_y^2}$. Then the \mathbb{G} -invariance of the operator's kernel function (Eq. (53b)) implies that the matrix form of the operator, in the chosen basis sets, is \mathbb{G} -equivariant w.r.t the group representations $\rho_{\mathcal{L}_x^2}$ and $\rho_{\mathcal{L}_y^2}$ (Eq. (58)).*

Proof. The proof follows by choosing appropriate \mathbb{G} -stable basis sets $\{\phi_i\} \subset \mathcal{L}_x^2$ and $\{\psi_j\} \subset \mathcal{L}_y^2$, so that for all $g \in \mathbb{G}$ we have $g \triangleright_{\mathcal{L}_x^2} \phi_i = \phi_{g \triangleright i}$ and $g \triangleright_{\mathcal{L}_y^2} \psi_j = \psi_{g \triangleright j}$ with $g \triangleright i, g \triangleright j \in \mathbb{N}$. This basis sets the \mathbb{G} -invariance of the kernel translates into algebraic constraints on the matrix form \mathbf{T} .

$$\begin{aligned} k(\mathbf{x}, \mathbf{y}) &= \kappa(g^{-1} \triangleright_x \mathbf{x}, g^{-1} \triangleright_y \mathbf{y}) \quad \forall g \in \mathbb{G}, \mathbf{x} \in \mathcal{X}, \mathbf{y} \in \mathcal{Y} \\ \sum_{i \in \mathbb{N}} \sum_{j \in \mathbb{N}} \mathbf{T}_{i,j} \phi_i(\mathbf{x}) \psi_j(\mathbf{y}) &= \sum_{i \in \mathbb{N}} \sum_{j \in \mathbb{N}} \mathbf{T}_{i,j} [g \triangleright_{\mathcal{L}_x^2} \phi_i](\mathbf{x}) [g \triangleright_{\mathcal{L}_y^2} \psi_j](\mathbf{y}) = \sum_{i \in \mathbb{N}} \sum_{j \in \mathbb{N}} \mathbf{T}_{i,j} \phi_{g \triangleright i}(\mathbf{x}) \psi_{g \triangleright j}(\mathbf{y}) \end{aligned} \quad (59)$$

That is, the kernel is \mathbb{G} -equivariant if the operator's matrix satisfies $\mathbf{T}_{i,j} = \mathbf{T}_{g \triangleright i, g \triangleright j}$ for all $g \in \mathbb{G}$, $i, j \in \mathbb{N}$. This condition exactly characterizes the \mathbb{G} -equivariance of the matrix form.

$$\begin{aligned} \mathbf{T}_{i,j} = \langle \psi_i, \mathbb{T} \phi_j \rangle_{P_y} &= \langle \psi_{g \triangleright i}, \mathbb{T} \phi_{g \triangleright j} \rangle_{P_y} = \mathbf{T}_{g \triangleright i, g \triangleright j} \quad \forall g \in \mathbb{G}, i, j \in \mathbb{N} \\ &= \langle g \triangleright_{\mathcal{L}_y^2} \psi_i, \mathbb{T} [g \triangleright_{\mathcal{L}_x^2} \phi_j] \rangle_{P_y} \\ &= \langle g \triangleright_{\mathcal{L}_y^2} \psi_i, g \triangleright_{\mathcal{L}_y^2} [\mathbb{T} \phi_j] \rangle_{P_y} \quad \text{s.t. Eq. (53a)} \\ &= \langle \psi_i, \mathbb{T} \phi_j \rangle_{P_y} = \mathbf{T}_{i,j} \quad \text{s.t. Eq. (46)} \end{aligned} \quad (60)$$

□

K.1.2 BLOCK-DIAGONAL STRUCTURE OF THE OPERATOR MATRIX FORM

According to Thm. I.16, a Hilbert space with a compact symmetry group \mathbb{G} decomposes into n_{iso} orthogonal subspaces—one for each irreducible representation type—yielding an orthogonal decomposition of the operator's input and output spaces:

$$\mathcal{L}_x^2 := \bigoplus_{1 \leq k \leq n_{\text{iso}}} \mathcal{L}_x^{2(k)}, \quad \text{and} \quad \mathcal{L}_y^2 := \bigoplus_{1 \leq k \leq n_{\text{iso}}} \mathcal{L}_y^{2(k)}, \quad (61)$$

where $\mathcal{L}_x^{2(k)}$ and $\mathcal{L}_y^{2(k)}$ denote the k -th isotypic subspaces of \mathcal{L}_x^2 and \mathcal{L}_y^2 , respectively. Such that any function in these spaces can be decomposed into a sum of its projections onto the isotypic subspaces:

$$f(\mathbf{x}) = \sum_{k=1}^{n_{\text{iso}}} f^{(k)}(\mathbf{x}), \quad h(\mathbf{y}) = \sum_{k=1}^{n_{\text{iso}}} h^{(k)}(\mathbf{y}) \quad \text{with} \quad f^{(k)} \in \mathcal{L}_x^{2(k)}, h^{(k)} \in \mathcal{L}_y^{2(k)}. \quad (62)$$

The orthogonal decomposition of the function spaces implies there exist unitary operators $\mathbf{A} : \mathcal{L}_x^2 \rightarrow \mathcal{L}_x^2$ and $\mathbf{B} : \mathcal{L}_y^2 \rightarrow \mathcal{L}_y^2$ (with matrix forms \mathbf{A} and \mathbf{B}), that describe a change of basis from the canonical basis to an *isotypic basis*, $\mathbb{I}_{\mathcal{L}_x^2}^{\text{iso}} = \bigcup_{k=1}^{n_{\text{iso}}} \mathbb{I}_{\mathcal{L}_x^2}^{2(k)} = \mathbf{A} \mathbb{I}_{\mathcal{L}_x^2}$ and $\mathbb{I}_{\mathcal{L}_y^2}^{\text{iso}} = \bigcup_{k=1}^{n_{\text{iso}}} \mathbb{I}_{\mathcal{L}_y^2}^{2(k)} = \mathbf{B} \mathbb{I}_{\mathcal{L}_y^2}$, where

the group’s representations decompose into a direct sum of representations per isotypic subspace (see [Thm. I.5](#)):

$$\rho_{\mathcal{L}_x^2}^{\text{iso}}(\cdot) := \mathbf{A} \rho_{\mathcal{L}_x^2}(\cdot) \mathbf{A}^* = \bigoplus_{k=1}^{n_{\text{iso}}} \rho_{\mathcal{L}_x^{2(k)}}^{\text{iso}}(\cdot) \quad \text{and} \quad \rho_{\mathcal{L}_y^2}^{\text{iso}}(\cdot) := \mathbf{B} \rho_{\mathcal{L}_y^2}(\cdot) \mathbf{B}^* = \bigoplus_{k=1}^{n_{\text{iso}}} \rho_{\mathcal{L}_y^{2(k)}}^{\text{iso}}(\cdot). \quad (63)$$

Then, denoting the matrix form of \mathbb{T} in the isotypic basis by $\mathbf{T}^{\text{iso}} = \mathbf{B}^* \mathbf{T} \mathbf{A}$, the \mathbb{G} -equivariance of \mathbb{T} results in the matrix form of the operator in the isotypic basis being block-diagonal, with each block being \mathbb{G} -equivariant with respect to the group representations on the isotypic subspaces:

$$\begin{aligned} \mathbf{T}^{\text{iso}} &= \rho_{\mathcal{L}_y^2}^{\text{iso}}(g) \mathbf{T}^{\text{iso}} \rho_{\mathcal{L}_x^2}^{\text{iso}}(g^{-1}) \\ &= \bigoplus_{k=1}^{n_{\text{iso}}} \rho_{\mathcal{L}_y^{2(k)}}^{\text{iso}}(g) \mathbf{T}^{\text{iso}} \bigoplus_{k=1}^{n_{\text{iso}}} \rho_{\mathcal{L}_x^{2(k)}}^{\text{iso}}(g^{-1}) \quad \text{s.t. Eqs. (58) and (63)} \\ \mathbf{T}^{(k)} &= \rho_{\mathcal{L}_y^{2(k)}}(g) \mathbf{T}^{(k)} \rho_{\mathcal{L}_x^{2(k)}}(g^{-1}), \quad \forall k = 1, \dots, n_{\text{iso}} \quad \mathbf{T}^{\text{iso}} = \bigoplus_{k=1}^{n_{\text{iso}}} \mathbf{T}^{(k)} = \begin{bmatrix} \mathbf{T}^{(1)} & & \\ & \ddots & \\ & & \mathbf{T}^{(n_{\text{iso}})} \end{bmatrix}. \end{aligned} \quad (64)$$

Each $\mathbf{T}^{(k)}$ represents the matrix form of the operator $\mathbb{T}^{(k)} : \mathcal{L}_x^{2(k)} \mapsto \mathcal{L}_y^{2(k)}$ in the isotypic basis, acting on the isotypic subspaces of type k in the input and output spaces. This shows that \mathbb{G} -equivariant operators preserve the structure of isotypic subspaces without mixing functions from different types.

This property is crucial for the finite-rank approximation of the operator \mathbb{T} , as it reduces the problem to approximating lower-rank operators $\mathbb{T}^{(k)} : \mathcal{L}_x^{2(k)} \mapsto \mathcal{L}_y^{2(k)}$, for $k \in [1, n_{\text{iso}}]$. Moreover, the block diagonal structure of \mathbf{T}^{iso} allows us to rewrite [Eq. \(56\)](#) in the isotypic basis in terms of the action of each $\mathbb{T}^{(k)}$ on the projection $f^{(k)}$ of the function onto the k^{th} isotypic subspace, see [\(62\)](#), such that:

$$[\mathbb{T}f_\alpha](\mathbf{y}) = \sum_{k=1}^{n_{\text{iso}}} [\mathbb{T}^{(k)} f^{(k)}](\mathbf{y}) \equiv \sum_{k=1}^{n_{\text{iso}}} (\mathbf{T}^{(k)} \boldsymbol{\alpha}^{(k)})^\top \boldsymbol{\psi}^{(k)}(\mathbf{y}). \quad \boldsymbol{\psi}^{(k)}(\cdot) = [\psi_j^{(k)}(\cdot)]_{j \in \mathbb{N}}, \forall \psi_j^{(k)} \in \mathbb{I}_{\mathcal{L}_y^{2(k)}}. \quad (65)$$

In the isotypic basis $\mathbb{I}_{\mathcal{L}_x^2}^{\text{iso}} = \bigcup_{k=1}^{n_{\text{iso}}} \mathbb{I}_{\mathcal{L}_x^{2(k)}}$, the expansion coefficient vector $\boldsymbol{\alpha} = \bigoplus_{k=1}^{n_{\text{iso}}} \boldsymbol{\alpha}^{(k)}$ is formed from the projections of f onto each isotypic subspace: $\boldsymbol{\alpha}^{(k)} = [\langle \phi_i^{(k)}, f \rangle_{P_x}]_{i \in \mathbb{N}}$. The block-diagonal structure of \mathbf{T}^{iso} is only one of the algebraic constraints imposed on the matrix form of \mathbb{T} by the \mathbb{G} -equivariance condition. The next section describes the further structural constraints on each block.

K.1.3 STRUCTURE OF OPERATORS BETWEEN ISOTYPIC SUBSPACES

In this section, we shift the focus from the input and output function spaces, \mathcal{L}_x^2 and \mathcal{L}_y^2 ; and the operator $\mathbb{T} : \mathcal{L}_x^2 \mapsto \mathcal{L}_y^2$, to their individual isotypic subspaces, $\mathcal{L}_x^{2(k)}$ and $\mathcal{L}_y^{2(k)}$ for $k \in [1, n_{\text{iso}}]$, and the operators $\mathbb{T}^{(k)} : \mathcal{L}_x^{2(k)} \mapsto \mathcal{L}_y^{2(k)}$ ([Eq. \(61\)](#)).

Recall from [Thm. I.16](#), that each isotypic subspace possesses unitary group representations that decompose into direct sums of (infinitely many) multiplicities of the irreducible representation of type k ; that is:

$$\rho_{\mathcal{L}_x^{2(k)}}(g) \sim \bigoplus_{p=1}^{\infty} \bar{\rho}_k(g) \quad \text{and} \quad \rho_{\mathcal{L}_y^{2(k)}}(g) \sim \bigoplus_{p=1}^{\infty} \bar{\rho}_k(g). \quad (66)$$

This implies that each isotypic subspace further decomposes into (infinitely many) finite-dimensional \mathbb{G} -stable subspaces: $\mathcal{L}_x^{2(k)} := \bigoplus_{p=1}^{\infty} \mathcal{L}_x^{2k,p}$ and $\mathcal{L}_y^{2(k)} := \bigoplus_{p=1}^{\infty} \mathcal{L}_y^{2k,p}$. Each subspace $\mathcal{L}_x^{2k,p}$ (and similarly $\mathcal{L}_y^{2k,p}$) has finite dimension $d_k \leq \infty$ and its elements transform according to the irreducible representation $\bar{\rho}_k$ of the group \mathbb{G} .

The modular structure of the isotypic subspaces results in a similar structure of $\mathbb{T}^{(k)}$ as it decomposes into \mathbb{G} -equivariant maps between the multiplicities of finite-dimensional, irreducible, \mathbb{G} -stable subspaces $\mathbb{T}_{i,j}^{(k)} : \mathcal{L}_x^{2k,i} \mapsto \mathcal{L}_y^{2k,j}$ for $i, j \in \mathbb{N}$, such that, in such basis the infinite-dimensional matrix form of $\mathbb{T}^{(k)}$ can be expressed as:

$$\mathbf{T}^{(k)} = \begin{bmatrix} \mathbf{T}_{1,1}^{(k)} & \mathbf{T}_{1,2}^{(k)} & \dots \\ \mathbf{T}_{2,1}^{(k)} & \ddots & \dots \\ \vdots & \vdots & \ddots \end{bmatrix} \quad \text{s.t.} \quad \mathbf{T}_{i,j}^{(k)} \in \text{End}_{\mathbb{G}}(\bar{\mathcal{V}}_k), \forall i, j \in \mathbb{N}, \quad (67)$$

Where each $\mathbf{T}_{i,j}^{(k)}$ is constrained to be an (real) endomorphism of the vector space $\bar{\mathcal{V}}_k$ associated to the irreducible representation $\bar{\rho}_k : \mathbb{G} \rightarrow \mathbb{U}(\bar{\mathcal{V}}_k)$.

In practical terms, this implies that each operator $\mathbb{T}^{(k)}$ is constructed from a combination of infinitely many elements of the endomorphism space $\text{End}_{\mathbb{G}}(\bar{\mathcal{V}}_k)$ (Thm. I.4). Crucially, for compact symmetry groups and real irreducible representations, the space of endomorphisms is 1, 2 or 4 dimensional, and possess an analytical basis set, defined by the map $\Psi_k : \mathbb{K} \rightarrow \text{End}_{\mathbb{G}}(\bar{\mathcal{V}}_k)$, where $\mathbb{K} \in \{\mathbb{R}, \mathbb{C}, \mathbb{H}\}$ denotes the basis algebra associated to the irreducible representation $\bar{\rho}_k$ (see details in Thm. I.13).

Such a constraint implies that each $\mathbb{T}_{i,j}^{(k)}$ can be expressed as a linear combination of endomorphism basis elements: $\mathbb{T}_{i,j}^{(k)} = \sum_{b \in \mathbb{K}} \alpha_{b,i,j} \Psi_k(b)$. Since every element of the endomorphism basis $\{\Psi_k(b)\}_{b \in \mathbb{K}}$ is a finite-dimensional orthogonal matrix, this structure can be interpreted as a constraint on the dimensionality of the singular spaces of the operator $\mathbb{T}^{(k)}$ to be of dimension larger than $d_k = |\bar{\rho}_k|$, as summarized in the following proposition:

Proposition K.3 (Minimum dimensionality of singular space of \mathbb{G} -equivariant operators between isotypic subspaces). *Let $\mathbb{T}^{(k)} : \mathcal{L}_{\mathbf{x}}^{2(k)} \mapsto \mathcal{L}_{\mathbf{y}}^{2(k)}$ be a \mathbb{G} -equivariant operator between isotypic subspaces $\mathcal{L}_{\mathbf{x}}^{2(k)}$ and $\mathcal{L}_{\mathbf{y}}^{2(k)}$ of type k . Then, the minimum dimension of a singular space of the operator is d_k .*

Proof. The proof follows naturally from Thm. I.10, which characterizes the singular spaces of each irrep endomorphism. \square

K.2 FINITE-RANK APPROXIMATION OF \mathbb{G} -EQUIVARIANT OPERATORS

In practical applications, infinite-dimensional operators are approximated by finite-dimensional ones to enable computation. For any linear integral operator $\mathbb{T} : \mathcal{L}_{\mathbf{x}}^2 \mapsto \mathcal{L}_{\mathbf{y}}^2$, the optimal rank- r approximation in the Hilbert-Schmidt norm is obtained by truncating its SVD to the top r singular values and associated left/right singular functions. Let $\{\sigma_i\}_{i=1}^{\infty}$ be the singular values of \mathbb{T} in decreasing order and let $\{u_i\}_{i=1}^{\infty} \subset \mathcal{L}_{\mathbf{x}}^2$, $\{v_i\}_{i=1}^{\infty} \subset \mathcal{L}_{\mathbf{y}}^2$ be the corresponding singular functions satisfying $\langle v_i, \mathbb{T}u_i \rangle_{P_{\mathbf{y}}} = \sigma_i$ for each $i \in \mathbb{N}$ and $\langle v_i, \mathbb{T}u_j \rangle_{P_{\mathbf{y}}} = 0$ when $i \neq j$. The best rank- r approximation of \mathbb{T} is then given by (Eckart and Young, 1936):

$$\mathbb{T}_r f = \sum_{i=1}^r \sigma_i \langle u_i, f \rangle_{P_{\mathbf{x}}} v_i, \quad \forall f \in \mathcal{L}_{\mathbf{x}}^2, \quad \iff \quad \kappa(\mathbf{x}, \mathbf{y}) \approx \sum_{i=1}^r \sigma_i u_i(\mathbf{x}) v_i(\mathbf{y}). \quad (68)$$

Since the left and right singular functions form orthonormal bases for $\mathcal{L}_{\mathbf{y}}^2$ and $\mathcal{L}_{\mathbf{x}}^2$, a rank- r approximation reduces these infinite-dimensional spaces to the r -dimensional subspaces $\mathcal{F}_{\mathbf{x}} = \text{span}(\{u_i\}_{i=1}^r)$ and $\mathcal{F}_{\mathbf{y}} = \text{span}(\{v_i\}_{i=1}^r)$.

When $\mathcal{L}_{\mathbf{x}}^2$ and $\mathcal{L}_{\mathbf{y}}^2$ are symmetric function spaces with group actions $\triangleright_{\mathcal{L}_{\mathbf{x}}^2}$ and $\triangleright_{\mathcal{L}_{\mathbf{y}}^2}$ of a compact group \mathbb{G} , and \mathbb{T} is \mathbb{G} -equivariant, the finite-rank approximation $\mathbb{T}_r : \mathcal{F}_{\mathbf{x}} \rightarrow \mathcal{F}_{\mathbf{y}}$ must satisfy that for all $f \in \mathcal{F}_{\mathbf{x}}$, $h \in \mathcal{F}_{\mathbf{y}}$, and $g \in \mathbb{G}$, both $g \triangleright_{\mathcal{L}_{\mathbf{x}}^2} f \in \mathcal{F}_{\mathbf{x}}$ and $g \triangleright_{\mathcal{L}_{\mathbf{y}}^2} h \in \mathcal{F}_{\mathbf{y}}$. This ensures that $g \triangleright_{\mathcal{L}_{\mathbf{y}}^2} [\mathbb{T}_r f] = \mathbb{T}_r [g \triangleright_{\mathcal{L}_{\mathbf{x}}^2} f]$ (see Sec. J).

Moreover, since $\mathcal{L}_{\mathbf{x}}^2$ and $\mathcal{L}_{\mathbf{y}}^2$ decompose orthogonally into isotypic subspaces, $\mathcal{L}_{\mathbf{x}}^2 = \bigoplus_{1 \leq k \leq n_{\text{iso}}} \mathcal{L}_{\mathbf{x}}^{2(k)}$ and $\mathcal{L}_{\mathbf{y}}^2 = \bigoplus_{1 \leq k \leq n_{\text{iso}}} \mathcal{L}_{\mathbf{y}}^{2(k)}$, the operator \mathbb{T} is completely determined by the n_{iso} operators $\mathbb{T}^{(k)} : \mathcal{L}_{\mathbf{x}}^{2(k)} \rightarrow \mathcal{L}_{\mathbf{y}}^{2(k)}$ (see Subsubsec. K.1.2). Thus, the \mathbb{G} -equivariance of \mathbb{T}_r depends on that of each finite-rank operator $\mathbb{T}_{r_k}^{(k)} : \mathcal{F}_{\mathbf{x}}^{(k)} \rightarrow \mathcal{F}_{\mathbf{y}}^{(k)}$, which requires the approximated subspaces $\mathcal{F}_{\mathbf{x}}^{(k)}$ and $\mathcal{F}_{\mathbf{y}}^{(k)}$ to be \mathbb{G} -stable. For simplicity, we assume $|\mathcal{F}_{\mathbf{x}}^{(k)}| = |\mathcal{F}_{\mathbf{y}}^{(k)}| = r_k$, although this equality need not hold in general.

Constraints on the dimensionality of truncation Each approximation of an isotypic subspace $\mathcal{L}_{\mathbf{x}}^{2(k)}$ (and similarly $\mathcal{L}_{\mathbf{y}}^{2(k)}$) is \mathbb{G} -stable if the group representation is defined using a truncated multiplicity $m_k < \infty$ for the k^{th} irreducible representation, i.e. $\rho_{\mathcal{F}_{\mathbf{x}}^{(k)}} \sim \bigoplus_{p=1}^{m_k} \bar{\rho}_k$ and $\rho_{\mathcal{L}_{\mathbf{y}}^{2(k)}} \sim \bigoplus_{p=1}^{m_k} \bar{\rho}_k$. Consequently, the dimension of the approximated subspaces is multiple of the irreducible representation's dimension: $r_k = d_k m_k$ (see Subsubsec. K.1.3).

Spectral decomposition Given the block-diagonal structure of the operator \mathbb{T} in the isotypic basis (Eq. (64)), the truncated SVD of \mathbb{T} reduces to performing the truncated SVD of each per-isotypic operator $\mathbb{T}^{(k)}$. Let $\mathcal{F}_{\mathbf{x}} \subset \mathcal{L}_{\mathbf{x}}^2$ and $\mathcal{F}_{\mathbf{y}} \subset \mathcal{L}_{\mathbf{y}}^2$ be the \mathbb{G} -stable finite-dimensional approximations of the input/output spaces of \mathbb{T} , endowed with group representations $\rho_{\mathcal{F}_{\mathbf{x}}} = \bigoplus_{k=1}^{n_{\text{iso}}} \rho_{\mathcal{F}_{\mathbf{x}}^{(k)}} = \bigoplus_{k=1}^{n_{\text{iso}}} \bigoplus_{p=1}^{m_k} \bar{\rho}_k$ and $\rho_{\mathcal{F}_{\mathbf{y}}} = \bigoplus_{k=1}^{n_{\text{iso}}} \rho_{\mathcal{F}_{\mathbf{y}}^{(k)}} = \bigoplus_{k=1}^{n_{\text{iso}}} \bigoplus_{p=1}^{m_k} \bar{\rho}_k$. Here, $m_k \in \mathbb{N}$ denotes the multiplicity of the irreducible

representation of type k , and $d_k := |\bar{\rho}_k|$ is its dimension. Then, the structural constraints on the SVD of the restriction of \mathbb{T} to these spaces are summarized in the following theorem:

Theorem K.4 (Isotypic-spectral basis). *Let \mathbb{T} be a \mathbb{G} -equivariant operator and let $\mathbb{T}_\star: \mathcal{F}_y \rightarrow \mathcal{F}_x$ be its \mathbb{G} -equivariant restriction in finite dimensions. Then, the singular value decomposition of the restricted operator matrix representation \mathbb{T}_\star reduces to:*

$$\mathbb{T}_\star = \bigoplus_{k=1}^{n_{\text{iso}}} \mathbb{T}_\star^{(k)} = \bigoplus_{k=1}^{n_{\text{iso}}} \mathbb{W}_\star^{(k)} \mathbb{S}_\star^{(k)} \mathbb{M}_\star^{(k)\top} = \bigoplus_{k=1}^{n_{\text{iso}}} \mathbb{U}_\star^{(k)} (\mathbb{\Sigma}_\star^{(k)} \otimes \mathbb{I}_{d_k}) \mathbb{V}_\star^{(k)\top}$$

Where \mathbb{I}_{d_k} denotes the identity matrix in d_k -dimensions and $\mathbb{\Sigma}_\star^{(k)}$ denotes the infinite-dimensional diagonal matrix singular values per irreducible \mathbb{G} -stable subspace.

Thm. K.4 shows that symmetries force each isotypic subspace’s singular space to have dimension at least d_k , which is the minimum required for a faithful representation of $\mathbb{G}^{(k)}$ (see **Thm. I.8**). Because in practice our goal is to approximate the top r singular spaces of \mathbb{T} , this result precisely characterizes the constraints imposed by \mathbb{G} -equivariance on the optimal rank- r truncation’s spectral basis and corresponding kernel function in **Eq. (13)**, as summarized in the following corollary:

Corollary K.5 (Symmetry constraints on the spectral basis). *Let \mathbb{T} be a \mathbb{G} -equivariant operator and let $\mathbb{T}_\star: \mathcal{F}_y \rightarrow \mathcal{F}_x$ be its \mathbb{G} -equivariant restriction in r -dimensions. Then, the spectral basis of \mathbb{T}_\star is given by:*

$$\kappa_\star(\mathbf{x}, \mathbf{y}) = \sum_{k=1}^{n_{\text{iso}}} \kappa_\star^{(k)}(\mathbf{x}, \mathbf{y}) = \sum_{k=1}^{n_{\text{iso}}} \sum_{s=1}^{r_k} \sigma_s^{(k)} \sum_{i=1}^{d_k} u_{s,i}^{(k)}(\mathbf{x}) v_{s,i}^{(k)}(\mathbf{y}), \quad (69)$$

where $\{u_{s,i}^{(k)}\}_{i \in [d_k]}$ and $\{v_{s,i}^{(k)}\}_{i \in [d_k]}$ are the left and right singular basis sets of the s^{th} singular space of $\mathbb{T}^{(k)}$. Note that the truncated dimension is restricted by the dimensionality and multiplicities of the individual irreducible representations $r = \sum_{k=1}^{d_{\text{iso}}} r_k = \sum_{k=1}^{d_{\text{iso}}} d_k m_k$.

L RELEVANT \mathbb{G} -EQUIVARIANT OPERATORS IN PROBABILITY THEORY

In this section we study the properties of expectations and covariances of functions of symmetric random variables in the presence our assumed symmetry priors **Eq. (6)**. In a nutshell, we characterize how expectations of observables of symmetric random variables are invariant to the group action, and that the covariance and cross-covariance matrices in these spaces are \mathbb{G} -equivariant and hence inherit rich structural constraints that can aid in empirical estimation.

Let (\mathbf{x}, \mathbf{y}) be two vector-valued random variables over the probability spaces $(\mathcal{X}, \Sigma_{\mathcal{X}}, P_{\mathbf{x}})$ and $(\mathcal{Y}, \Sigma_{\mathcal{Y}}, P_{\mathbf{y}})$, with $\mathcal{L}_{\mathbf{x}}^2$ and $\mathcal{L}_{\mathbf{y}}^2$ being the corresponding square-integrable function spaces and $\mathbb{1}_{P_{\mathbf{x}}} \in \mathcal{L}_{\mathbf{x}}^2, \mathbb{1}_{P_{\mathbf{y}}} \in \mathcal{L}_{\mathbf{y}}^2$ the characteristic functions of sets with nonzero probability.

When $\mathcal{L}_{\mathbf{x}}^2$ and $\mathcal{L}_{\mathbf{y}}^2$ are symmetric function spaces (see **Sec. J**), denote their orthogonal isotypic decompositions by $\mathcal{L}_{\mathbf{x}}^2 := \bigoplus_{k=1}^{n_{\text{iso}}} \mathcal{L}_{\mathbf{x}}^{2(k)}$ and $\mathcal{L}_{\mathbf{y}}^2 := \bigoplus_{k=1}^{n_{\text{iso}}} \mathcal{L}_{\mathbf{y}}^{2(k)}$ (cf. **Thm. I.16**). Any function $f \in \mathcal{L}_{\mathbf{x}}^2$ or $h \in \mathcal{L}_{\mathbf{y}}^2$ decomposes as $f = \sum_{k=1}^{n_{\text{iso}}} f^{(k)}$ and $h = \sum_{k=1}^{n_{\text{iso}}} h^{(k)}$ (see **Eq. (62)**). By convention, the first isotypic subspace corresponds to the trivial group action. Thus, we write $\mathcal{L}_{\mathbf{x}}^{2\text{inv}} := \mathcal{L}_{\mathbf{x}}^{2(1)} \subset \mathcal{L}_{\mathbf{x}}^2$ and denote the \mathbb{G} -invariant component of f by $f^{\text{inv}} := f^{(1)}$ (and similarly for $\mathcal{L}_{\mathbf{y}}^2$).

L.1 THE EXPECTATION OPERATOR

The expected value of a function $f \in \mathcal{F} := \mathcal{L}_{\mathbf{x}}^2$ can be interpreted as the result of applying a linear integral operator that projects each $f \in \mathcal{F}$ to a constant function evaluating to the function’s expected value $\mathbb{E}_{P_{\mathbf{x}}} f$.

Definition L.1 (Expectation operator). *Let $\mathcal{F} \subseteq \mathcal{L}_{\mathbf{x}}^2$ be a function space. The expectation operator $\mathbb{E}_{\mathbf{x}}: \mathcal{F} \mapsto \mathcal{F}$ is a linear integral operator defined by a constant kernel function $k_{\mathbb{E}}(\mathbf{x}, \mathbf{x}') = \mathbb{1}_{P_{\mathbf{x}}}(\mathbf{x}) \mathbb{1}_{P_{\mathbf{x}}}(\mathbf{x}')$ for all $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$, such that this operator maps any function f to a constant function that evaluates to the function’s expected value $\mathbb{1}_{P_{\mathbf{x}}}(\cdot) \mathbb{E}_{P_{\mathbf{x}}} f$, that is:*

$$[\mathbb{E}_{\mathbf{x}} f](\mathbf{x}') = \int_{\mathcal{X}} k_{\mathbb{E}}(\mathbf{x}, \mathbf{x}') f(\mathbf{x}) \mu(d\mathbf{x}) = \mathbb{1}_{P_{\mathbf{x}}}(\mathbf{x}') \int_{\mathcal{X}} f(\mathbf{x}) \mu(d\mathbf{x}) \equiv \mathbb{1}_{P_{\mathbf{x}}}(\mathbf{x}') \mathbb{E}_{P_{\mathbf{x}}} f. \quad (70)$$

Whenever \mathcal{F} is a symmetric function space, the operator $\mathbb{E}_{\mathbf{x}}$ commutes with the group action and is \mathbb{G} -invariant (**Thm. I.11**):

Proposition L.2 (\mathbb{G} -invariant expectation operator). *Let \mathcal{F} be a symmetric function space with the action $\triangleright_{\mathcal{F}}$ of a compact symmetry group \mathbb{G} . Then, the expectation operator commutes with the group action and is a \mathbb{G} -invariant operator $\mathbb{E}_{\mathbf{x}} : \mathcal{F} \mapsto \mathcal{F}^{\text{inv}} \subseteq \mathcal{F}$:*

$$\mathbb{E}_{\mathbf{x}}[g \triangleright_{\mathcal{F}} f] = g \triangleright_{\mathcal{F}} [\mathbb{E}_{\mathbf{x}} f] \quad \text{and} \quad \mathbb{E}_{\mathbf{x}} f = \mathbb{E}_{\mathbf{x}}[g \triangleright_{\mathcal{F}} f] \in \mathcal{F}^{\text{inv}}, \quad \forall f \in \mathcal{F}, g \in \mathbb{G}. \quad (71)$$

Proof. The operator $\mathbb{E}_{\mathbf{x}}$ commutes with the group action as its kernel function $k_{\mathbb{E}}$ is constant and therefore \mathbb{G} -invariant (Thm. K.1). Furthermore since the image of the expectation operator are constant functions, these functions belong to the subspace of \mathbb{G} -invariant functions, \mathcal{F}^{inv} . \square

As an operator that commutes with the group action, the expectation operator decomposes into $\mathbb{E}_{\mathbf{x}} := \bigoplus_{k=1}^{n_{\text{iso}}} \mathbb{E}_{\mathbf{x}}^{(k)}$, where $\mathbb{E}_{\mathbf{x}}^{(k)} : \mathcal{F}^{(k)} \mapsto \mathcal{F}^{(k)}$ denotes the restriction of $\mathbb{E}_{\mathbf{x}}$ to the isotypic subspace $\mathcal{F}^{(k)}$ (Subsubsec. K.1.2). However, since the image of the operator lies in the subspace of \mathbb{G} -invariant functions, $\mathfrak{S}(\mathbb{E}_{\mathbf{x}}) \subset \mathcal{F}^{\text{inv}}$, it follows that $\mathbb{E}_{\mathbf{x}}^{(k)} = \mathbf{0}$ for every $k \neq \text{inv}$. Consequently, we obtain the following:

Corollary L.3 (Expectation of a function depends only on its \mathbb{G} -invariant component). *For any function $f \in \mathcal{F}$, the expectation depends only on its \mathbb{G} -invariant component:*

$$[\mathbb{E}_{\mathbf{x}} f](\cdot) = \sum_{k=1}^{n_{\text{iso}}} [\mathbb{E}_{\mathbf{x}}^{(k)} f^{(k)}](\cdot) = [\mathbb{E}_{\mathbf{x}}^{\text{inv}} f^{\text{inv}}](\cdot) := \mathbb{1}_{\mu}(\cdot) \mathbb{E}_{\mu} f^{\text{inv}}. \quad (72)$$

Corollary L.4 (Functions without a \mathbb{G} -invariant component are centered). *Any function $f = \sum_{k=1}^{n_{\text{iso}}} f^{(k)} \in \mathcal{L}_{\mathbf{x}}^2$ without a \mathbb{G} -invariant component, i.e., $f^{\text{inv}} = 0$, is centered:*

$$[\mathbb{E}_{\mathbf{x}} f](\cdot) = \sum_{k=2}^{n_{\text{iso}}} [\mathbb{E}_{\mathbf{x}}^{(k)} f^{(k)}](\cdot) = \mathbb{1}_{\mu}(\cdot) 0, \quad \iff \quad \mathbb{E}_{\mu} f = 0, \quad \forall f \in \mathcal{L}_{\mathbf{x}}^{2\text{inv}\perp}. \quad (73)$$

To better comprehend these concepts we refer the reader to Thm. J.4.

L.2 THE CROSS-COVARIANCE OPERATOR

Given two vector-valued random variables ($\mathbf{x} = [x_1, \dots, x_n], \mathbf{y} = [y_1, \dots, y_m]$) defined on the measure spaces $(\mathcal{X}, \Sigma_{\mathcal{X}}, P_{\mathbf{x}})$ and $(\mathcal{Y}, \Sigma_{\mathcal{Y}}, P_{\mathbf{y}})$, a key statistic assessing the linear relationship between scalar components is the covariance:

$$\text{Cov}(x_i, y_j) = \mathbb{E}_{P_{\mathbf{x}\mathbf{y}}}[x_i - \mathbb{E}_{\mathbf{x}}[x_i]](y_j - \mathbb{E}_{\mathbf{y}}[y_j])] = \mathbb{E}_{P_{\mathbf{x}\mathbf{y}}}[x_i y_j] - \mathbb{E}_{\mathbf{x}}[x_i] \mathbb{E}_{\mathbf{y}}[y_j].$$

For vector-valued random variables, the cross-covariance matrix $\text{Cov}(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{n \times m}$ is defined entrywise by $\text{Cov}(\mathbf{x}, \mathbf{y})_{i,j} := \text{Cov}(x_i, y_j)$. The cross-covariance operator is the extension of this concept to the Hilbert spaces of functions $\mathcal{L}_{\mathbf{x}}^2$ and $\mathcal{L}_{\mathbf{y}}^2$.

Definition L.5 (Cross-covariance operator (Fukumizu et al., 2004)). *Let $\mathcal{F}_{\mathbf{x}} \subseteq \mathcal{L}_{\mathbf{x}}^2$ and $\mathcal{L}_{\mathbf{y}}^2 \subseteq \mathcal{L}_{\mathbf{y}}^2$ be two Hilbert spaces of functions defined on the random variables \mathbf{x} and \mathbf{y} , which take values in the measure spaces $(\mathcal{X}, \Sigma_{\mathcal{X}}, P_{\mathbf{x}})$ and $(\mathcal{Y}, \Sigma_{\mathcal{Y}}, P_{\mathbf{y}})$, respectively. The cross-covariance operator $\mathbf{C}_{\mathbf{x}\mathbf{y}} : \mathcal{L}_{\mathbf{y}}^2 \mapsto \mathcal{L}_{\mathbf{x}}^2$ is a linear integral operator defined by*

$$\langle f, \mathbf{C}_{\mathbf{x}\mathbf{y}} h \rangle_{P_{\mathbf{x}}} := \text{Cov}(f, h) = \mathbb{E}_{P_{\mathbf{x}\mathbf{y}}}[f(\mathbf{x})h(\mathbf{y})] - \mathbb{E}_{\mathbf{x}}[f(\mathbf{x})] \mathbb{E}_{\mathbf{y}}[h(\mathbf{y})], \quad \forall f \in \mathcal{L}_{\mathbf{x}}^2, h \in \mathcal{L}_{\mathbf{y}}^2. \quad (74)$$

Choosing separable basis sets for the two spaces, $\mathbb{I}_{\mathcal{L}_{\mathbf{x}}^2} = \{\phi_i\}_{i \in \mathbb{N}}$ and $\mathbb{I}_{\mathcal{L}_{\mathbf{y}}^2} = \{\psi_i\}_{i \in \mathbb{N}}$, the matrix representation of the cross-covariance operator has entries $[\mathbf{C}_{\mathbf{x}, \mathbf{y}}]_{i,j} := \langle \phi_i, \mathbf{C}_{\mathbf{x}\mathbf{y}} \psi_j \rangle_{P_{\mathbf{x}}} = \text{Cov}(\phi_i, \psi_j)$, where the covariance is computed with respect to the joint measure $P_{\mathbf{x}\mathbf{y}}$ and the marginals $P_{\mathbf{x}}$ and $P_{\mathbf{y}}$. Given a dataset of N samples from the joint distribution $(\mathbf{x}, \mathbf{y}) \sim P_{\mathbf{x}\mathbf{y}}$, the empirical estimate of the matrix form of the cross-covariance operator is

$$\widehat{\mathbf{C}}_{\mathbf{x}\mathbf{y}} = \frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \psi(\mathbf{y}_n)^{\top} - \widehat{\mathbb{E}}_{\mathbf{x}}[\phi(\mathbf{x}_n)] \widehat{\mathbb{E}}_{\mathbf{y}}[\psi(\mathbf{y}_n)]^{\top}, \quad \phi(\cdot) = [\phi(\cdot)]_{i \in \mathbb{N}}, \psi(\cdot) = [\psi(\cdot)]_{i \in \mathbb{N}}. \quad (75)$$

Note that the adjoint of the operator is defined by $\mathbf{C}_{\mathbf{y}\mathbf{x}}^* = \mathbf{C}_{\mathbf{y}\mathbf{x}} : \mathcal{L}_{\mathbf{x}}^2 \mapsto \mathcal{L}_{\mathbf{y}}^2$. In the case $\mathcal{L}_{\mathbf{x}}^2 = \mathcal{L}_{\mathbf{y}}^2$, the cross-covariance operator reduces to the covariance operator, and has an analog definition to Thm. L.5.

Covariance and cross-covariance operators of symmetric Hilbert spaces of functions Whenever \mathcal{L}_x^2 and \mathcal{L}_y^2 are symmetric function spaces, and the joint probability measure is \mathbb{G} -invariant, the cross-covariance operator C_{xy} commutes with the group action and is \mathbb{G} -equivariant (Subsec. I.1):

Proposition L.6 (\mathbb{G} -equivariant cross-covariance operator). *Let $\mathcal{L}_x^2 \subseteq \mathcal{L}_x^2$ and $\mathcal{L}_y^2 \subseteq \mathcal{L}_y^2$ be symmetric Hilbert spaces of functions endowed with the group actions $\triangleright_{\mathcal{L}_x^2}$ and $\triangleright_{\mathcal{L}_y^2}$ of a compact symmetry group \mathbb{G} . Then, whenever the joint probability measure is \mathbb{G} -invariant, i.e., $P_{xy}(\mathbb{B}, \mathbb{A}) = P_{xy}(g \triangleright_{\mathcal{X}} \mathbb{B}, g \triangleright_{\mathcal{Y}} \mathbb{A})$ for all $g \in \mathbb{G}, \mathbb{B} \in \Sigma_{\mathcal{X}}, \mathbb{A} \in \Sigma_{\mathcal{Y}}$, the cross-covariance operator $C_{xy} : \mathcal{L}_y^2 \mapsto \mathcal{L}_x^2$ (Thm. L.5) commutes with the group actions and is a \mathbb{G} -equivariant operator (Thm. K.1):*

$$g \triangleright_{\mathcal{L}_x^2} [C_{xy} h] = C_{xy} [g \triangleright_{\mathcal{L}_y^2} h], \quad \forall h \in \mathcal{L}_y^2, g \in \mathbb{G}. \quad (76)$$

Proof. To prove that the operator is \mathbb{G} -equivariant we must show its kernel function is \mathbb{G} -invariant (see Thm. K.1). The proof follows naturally in any regular basis of the input and output functions spaces $\mathbb{I}_{\mathcal{L}_x^2} = \{\phi_i\}_{i \in \mathbb{N}}$ and $\mathbb{I}_{\mathcal{L}_y^2} = \{\psi_j\}_{j \in \mathbb{N}}$, in which the group action on basis functions acts by permutations of basis functions, such that, $g \triangleright_{\mathcal{L}_x^2} \phi_i \equiv \phi_{g \triangleright i} \in \mathbb{I}_{\mathcal{L}_x^2}$ and $g \triangleright_{\mathcal{L}_y^2} \psi_j \equiv \psi_{g \triangleright j} \in \mathbb{I}_{\mathcal{L}_y^2}$, where $g \triangleright i, g \triangleright j \in \mathbb{N}$. Then we must show that that:

$$\begin{aligned} k(\mathbf{x}, \mathbf{y}) &= \kappa(g^{-1} \triangleright_{\mathcal{X}} \mathbf{x}, g^{-1} \triangleright_{\mathcal{Y}} \mathbf{y}) && \forall g \in \mathbb{G}, \mathbf{x} \in \mathcal{X}, \mathbf{y} \in \mathcal{Y} \\ \sum_{i \in \mathbb{N}} \sum_{j \in \mathbb{N}} [C_{xy}]_{i,j} \phi_i(\mathbf{x}) \psi_j(\mathbf{y}) &= \sum_{i \in \mathbb{N}} \sum_{j \in \mathbb{N}} [C_{xy}]_{i,j} [g \triangleright_{\mathcal{L}_x^2} \phi_i](\mathbf{x}) [g \triangleright_{\mathcal{L}_y^2} \psi_j](\mathbf{y}) && \text{s.t. Thms. J.1 and L.5} \\ \sum_{i \in \mathbb{N}} \sum_{j \in \mathbb{N}} \text{Cov}(\phi_i, \psi_j) \phi_i(\mathbf{x}) \psi_j(\mathbf{y}) &= \sum_{i \in \mathbb{N}} \sum_{j \in \mathbb{N}} \text{Cov}(\phi_i, \psi_j) \phi_{g \triangleright i}(\mathbf{x}) \psi_{g \triangleright j}(\mathbf{y}). \end{aligned} \quad (77)$$

Hence, the cross-covariance operator's kernel function is \mathbb{G} -invariant only if the covariance is \mathbb{G} -invariant:

$$\begin{aligned} \text{Cov}(\phi_i, \psi_j) &= \text{Cov}(g \triangleright_{\mathcal{L}_x^2} \phi_i, g \triangleright_{\mathcal{L}_y^2} \psi_j) && \forall g \in \mathbb{G}, i, j \in \mathbb{N} \\ \mathbb{E}_{P_{xy}}[\phi_i(\mathbf{x}) \psi_j(\mathbf{y})] &= \mathbb{E}_{P_{xy}}[\phi_i(g^{-1} \triangleright_{\mathcal{X}} \mathbf{x}) \psi_j(g^{-1} \triangleright_{\mathcal{Y}} \mathbf{y})] && \mathbb{E}_{\mu} f = \mathbb{E}_{\mu} g \triangleright f \\ \int_{\mathcal{X} \times \mathcal{Y}} \phi_i(\mathbf{x}) \psi_j(\mathbf{y}) P_{xy}(d\mathbf{x}, d\mathbf{y}) &= \int_{\mathcal{X} \times \mathcal{Y}} \phi_i(g^{-1} \triangleright_{\mathcal{X}} \mathbf{x}) \psi_j(g^{-1} \triangleright_{\mathcal{Y}} \mathbf{y}) P_{xy}(d\mathbf{x}, d\mathbf{y}) \\ &= \int_{\mathcal{X} \times \mathcal{Y}} \phi_i(\mathbf{x}) \psi_j(\mathbf{y}) P_{xy}(g \triangleright d\mathbf{x}, g \triangleright d\mathbf{y}) \\ &= \text{Cov}(\phi_i, \psi_j). \end{aligned} \quad (78)$$

□

An equivalent result follows for covariance operators of symmetric Hilbert spaces.

M STATISTICAL LEARNING THEORY

This section provides the development and proofs of the statistical learning guarantees in Thm. 5.1 for regression and conditional probability estimation using our proposed model.

Recall that regression and conditional probabilities can be expressed in terms of the conditional expectation operator $\mathbb{E}_{y|x} : \mathcal{L}_y^2 \rightarrow \mathcal{L}_x^2$ (see Eqs. (1) and (2)). Given that the operator is compact (Kostic et al., 2024b), it admits a singular value decomposition. Hence, the kernel function defining the operator Eq. (1) can be expanded in terms of the operator spectral basis:

$$\kappa(\mathbf{x}, \mathbf{y}) := \frac{dP_{xy}(\mathbf{x}, \mathbf{y})}{d(P_x(\mathbf{x}) \times P_y(\mathbf{y}))} = \sum_{i=0}^{\infty} \sigma_i u_i(\mathbf{x}) v_i(\mathbf{y}). \quad (79)$$

Where $(\sigma_i)_{i \in \mathbb{N}}$ denotes the operator's singular values, and $(u_i)_{i \in \mathbb{N}}$ and $(v_i)_{i \in \mathbb{N}}$ denote the left and right singular functions, which form complete orthonormal basis sets for \mathcal{L}_x^2 and \mathcal{L}_y^2 , respectively. Given that the operator's first singular value is $\sigma_0 = 1$, associated with the constant functions $u_0 = \mathbb{1}_{\mathcal{X}}, v_0 = \mathbb{1}_{\mathcal{Y}}$, the conditional expectation operator can be defined as:

$$\mathbb{E}_{y|x} = \sum_{i=1}^{\infty} \sigma_i u_i \langle v_i, \cdot \rangle_{P_y} = \mathbb{1}_{\mathcal{X}} \langle \mathbb{1}_{\mathcal{Y}}, \cdot \rangle_{P_y} + \underbrace{\sum_{i=1}^{\infty} \sigma_i u_i \langle v_i, \cdot \rangle_{P_y}}_{D_{y|x}}. \quad (80)$$

Where $D_{y|x}$ denotes the *deflated* operator, excluding the first eigen triplet (σ_0, u_0, v_0) . Leveraging the SVD of $E_{y|x}$, we approximate the operator’s action for any $h \in \mathcal{L}_y^2$ using a rank- r ($1 < r < \infty$) operator given by:

$$\begin{aligned} \mathbb{E}[h(\mathbf{y})|\mathbf{x}=\mathbf{x}] &= [E_{y|x}h](\mathbf{x}) \approx \mathbb{E}[h(\mathbf{y})] + \sum_{i=1}^r \sigma_i u_i^\theta(\mathbf{x}) \mathbb{E}[v_i^\theta(\mathbf{y})h(\mathbf{y})], \\ \text{s.t. } \mathbb{E}[u_i^\theta(\mathbf{x})] &= \mathbb{E}[v_i^\theta(\mathbf{y})] = 0, \forall i \geq 1. \end{aligned} \quad (81)$$

Where $(u_i^\theta)_{i=1}^r$ and $(v_i^\theta)_{i=1}^r$ denote parametrizations of the top- r left and right singular functions. Given that the operator’s kernel Eq. (79) preserves the probability mass, that is $\int_{\mathcal{X} \times \mathcal{Y}} \kappa(\mathbf{x}, \mathbf{y}) dP_{\mathbf{x}}(\mathbf{x}) dP_{\mathbf{y}}(\mathbf{y}) = 1$, every non-constant singular function is constrained to be centered, as described in the r.h.s of Eq. (81).

In the context of symmetries, we note that $D_{y|x}$ admits a block-diagonal structure w.r.t. to isotypic basis of associated \mathcal{L}^2 spaces. Indeed we have the following from Thm. K.4.

$$Q_{\mathbf{x}}^* D_{y|x} Q_{\mathbf{y}} = \bigoplus_{k=1}^{n_{\text{iso}}} Q_{\mathbf{x}}^{(k)*} D_{y|x}^{(k)} Q_{\mathbf{y}}^{(k)} = \bigoplus_{k=1}^{n_{\text{iso}}} [(\mathbf{U}^{(k)} \mathbf{S}^{(k)} \mathbf{V}^{(k)*}) \otimes \mathbf{I}_{d_k}]. \quad (82)$$

Where the unitary operators $Q_{\mathbf{x}}: \mathcal{L}_{\mathbf{x}}^2 \rightarrow \mathcal{L}_{\mathbf{x}}^2$ and $Q_{\mathbf{y}}: \mathcal{L}_{\mathbf{y}}^2 \rightarrow \mathcal{L}_{\mathbf{y}}^2$ change the basis to the isotypic decompositions $\mathbb{I}_{\mathcal{L}_{\mathbf{x}}^2} = \{\phi_{i,j}^{(k)}\}_{k \in [n_{\text{iso}}], i \in [m_k], j \in [d_k]}$ and $\mathbb{I}_{\mathcal{L}_{\mathbf{y}}^2} = \{\psi_{i,j}^{(k)}\}_{k \in [n_{\text{iso}}], i \in [m_k], j \in [d_k]}$, with i indexing each irreducible \mathbb{G} -stable subspace and j indexing the dimensions within that subspace (see Subsec. K.2).

Further, by Thm. K.4, the SVD of $D_{y|x}$ forces each isotypic subspace to have dimension at least $d_k = \bar{\rho}_k$ for every $k \in [n_{\text{iso}}]$.

$$Q_{\mathbf{x}}^{(k)*} D_{y|x}^{(k)} Q_{\mathbf{y}}^{(k)} = [\mathbf{U}^{(k)} \otimes \mathbf{I}_{d_k}] [\mathbf{S}^{(k)} \otimes \mathbf{I}_{d_k}] [\mathbf{V}^{(k)} \otimes \mathbf{I}_{d_k}]^*, k \in [n_{\text{iso}}], \quad (83)$$

where $Q_{\mathbf{x}}^{(k)*} Q_{\mathbf{x}}^{(k)}$ and $Q_{\mathbf{y}}^{(k)*} Q_{\mathbf{y}}^{(k)}$ are orthogonal projectors on k -th isotypic subspace, and

$$Q_{\mathbf{x}}^* D_{y|x} Q_{\mathbf{y}} = [\mathbf{I}_{n_{\text{iso}}} \otimes \mathbf{U}^{(k)} \otimes \mathbf{I}_{d_k}] [\mathbf{I}_{n_{\text{iso}}} \otimes \mathbf{S}^{(k)} \otimes \mathbf{I}_{d_k}] [\mathbf{I}_{n_{\text{iso}}} \otimes \mathbf{V}^{(k)} \otimes \mathbf{I}_{d_k}]^*. \quad (84)$$

Further, observe that the singular values of $D_{y|x}$ are elements of positive diagonal operators $\mathbf{S}^{(k)}$, denoted as $(\mathbf{S}^{(k)})_i = \sigma_i^{(k)}$, while the left and right singular functions are $u_i^{(k)} \otimes e_j^{d_k}$ and $v_i^{(k)} \otimes e_j^{d_k}$, respectively, for $i \in \mathbb{N}$, $j \in [d_k]$ and $k \in [n_{\text{iso}}]$, where e_j^d is j -th vector of standard basis of \mathbb{R}^d .

Given the constraints on the spectral basis of \mathbb{G} -equivariant operators (see Thm. K.5), our representation learning procedure approach results in feature maps:

$$\begin{aligned} \mathbf{u}_\theta(\cdot) &= \sum_{k \in [n_{\text{iso}}], i \in [m], j \in [d_k]} [e_k^{n_{\text{iso}}} \otimes e_i^m \otimes e_j^{d_k}] u_{i,j}^{\theta^{(k)}}(\cdot): \mathcal{X} \rightarrow \mathbb{R}^{r_m} \\ \mathbf{v}_\theta(\cdot) &= \sum_{k \in [n_{\text{iso}}], i \in [m], j \in [d_k]} [e_k^{n_{\text{iso}}} \otimes e_i^m \otimes e_j^{d_k}] v_{i,j}^{\theta^{(k)}}(\cdot): \mathcal{X} \rightarrow \mathbb{R}^{r_m}, \end{aligned} \quad (85)$$

which can further be separated into n_{iso} orthogonal blocks $\mathbf{u}_\theta^{(k)} = \sum_{i \in [m], j \in [d_k]} \phi_{i,j}^{\theta^{(k)}}$ and $\mathbf{v}_\theta^{(k)} = \sum_{i \in [m], j \in [d_k]} \psi_{i,j}^{\theta^{(k)}}$ as

$$\mathbf{u}_\theta^{(k)} = \sum_{i \in [m], j \in [d_k]} [e_i^m \otimes e_j^{d_k}] u_{i,j}^{\theta^{(k)}}(\cdot) \quad \text{and} \quad \mathbf{v}_\theta^{(k)} = \sum_{i \in [m], j \in [d_k]} [e_i^m \otimes e_j^{d_k}] v_{i,j}^{\theta^{(k)}}(\cdot). \quad (86)$$

In addition, the singular value matrices have a tensor form $\mathbf{S}_\theta = \text{diag}(\mathbf{S}_\theta^{(1)}, \dots, \mathbf{S}_\theta^{(n_{\text{iso}})})$, where $\mathbf{S}_\theta^{(k)} = \text{diag}(\sigma_1^{\theta^{(k)}}, \dots, \sigma_m^{\theta^{(k)}}) \otimes \mathbf{I}_{d_k}$ and $\sigma_i^{\theta^{(k)}} \in [0, 1]$, $i \in [m]$, $k \in [n_{\text{iso}}]$. Thus, we obtain the operator $D_\theta = E_\theta - \mathbb{1}_{P_{\mathbf{x}}} \otimes \mathbb{1}_{P_{\mathbf{y}}}$ in block form, $D_\theta = \bigoplus_{k \in [n_{\text{iso}}]} D_\theta^{(k)}$, where each $D_\theta^{(k)}$ acts on the k -th isotypic subspace as

$$[D_\theta^{(k)} f](\mathbf{x}) := \mathbf{u}_\theta^{(k)}(\mathbf{x})^\top \mathbf{S}_\theta^{(k)} \mathbb{E}_{\mathbf{y}}[\mathbf{v}_\theta^{(k)}(\mathbf{y}) f^{(k)}(\mathbf{y})], \quad f \in \mathcal{L}_{\mathbf{y}}^2, \quad (87)$$

and hence

$$[D_\theta f](\mathbf{x}) := \mathbf{u}_\theta(\mathbf{x})^\top \mathbf{S}_\theta \mathbb{E}_{\mathbf{y}}[\mathbf{v}_\theta(\mathbf{y}) f(\mathbf{y})], \quad f \in \mathcal{L}_{\mathbf{y}}^2. \quad (88)$$

Finally, we extend the definition of D_θ to vector-valued observables $\mathbf{h}: \mathcal{Y} \rightarrow \mathcal{Z}$ via basis expansions.

$$[D_\theta \mathbf{h}](\mathbf{x}) := \sum_{\ell} \mathbf{u}_\theta(\mathbf{x})^\top \mathbf{S}_\theta \mathbb{E}_y[\mathbf{v}_\theta(\mathbf{y}) \langle \mathbf{h}(\mathbf{y}), \mathbf{z}_\ell \rangle_{\mathcal{Z}} \mathbf{z}_\ell], \quad \mathbf{h} \in \mathcal{L}_y^2(\mathcal{Y}, \mathcal{Z}) \quad (89)$$

where $(\mathbf{z}_i)_{i \in [n_{\mathcal{Z}}]}$ is the orthonormal basis of \mathcal{Z} .

By doing so, we ensure that D_θ and, consequently, E_θ are \mathbb{G} -equivariant operators for both the scalar map $\mathcal{L}_y^2 \rightarrow \mathcal{L}_x^2$ and the vector-valued map $\mathcal{L}_y^2(\mathcal{Y}, \mathcal{Z}) \rightarrow \mathcal{L}_x^2(\mathcal{X}, \mathcal{Z})$. Moreover, a direct consequence of (89) is as follows.

Proposition M.1. *Let with \mathcal{Z} being a real Euclidean space endowed with symmetry group \mathbb{G} , and let $E_\theta: \mathcal{L}_{P_y}^2(\mathcal{Y}, \mathcal{Z}) \mapsto \mathcal{L}_{P_x}^2(\mathcal{X}, \mathcal{Z})$ be given by $E_\theta \mathbf{f} = \mathbb{E}_y[\mathbf{f}(\mathbf{y})] + D_\theta \mathbf{f}$. Then for every \mathbb{G} -equivariant $\mathbf{f} \in \mathcal{L}_{P_y}^2(\mathcal{Y}, \mathcal{Z})$ and every $\mathbf{x} \in \mathcal{X}$*

$$[E_\theta \mathbf{f}](g \triangleright_x \mathbf{x}) = \mathbb{E}_y[\mathbf{f}(\mathbf{y})] + [D_\theta \mathbf{f}](g \triangleright_x \mathbf{x}) = \mathbb{E}_y[\mathbf{f}(\mathbf{y})] + g \triangleright_z [D_\theta \mathbf{f}](\mathbf{x}) = g \triangleright_z [E_\theta \mathbf{f}](\mathbf{x}). \quad (90)$$

Proof. Since D_θ is \mathbb{G} -equivariant, for every $g \in \mathbb{G}$ we have that

$$[D_\theta \mathbf{h}](g^{-1} \triangleright_x \mathbf{x}) = [D_\theta[\mathbf{h}(g^{-1} \triangleright_y \cdot)]](\mathbf{x}) = \sum_i \mathbf{u}_\theta(\mathbf{x})^\top \mathbf{S}_\theta \mathbb{E}_y[\mathbf{v}_\theta(\mathbf{y}) \langle \mathbf{h}(g^{-1} \triangleright_y \mathbf{y}), \mathbf{z}_i \rangle_{\mathcal{Z}} \mathbf{z}_i],$$

which, using g instead of g^{-1} and the assumption that f is \mathbb{G} -equivariant, implies

$$\begin{aligned} [D_\theta \mathbf{h}](g \triangleright_x \mathbf{x}) &= \sum_i (\mathbf{u}_\theta(\mathbf{x})^\top \mathbf{S}_\theta \mathbb{E}_y[\mathbf{v}_\theta(\mathbf{y}) \langle g \triangleright_z \mathbf{h}(\mathbf{y}), \mathbf{z}_i \rangle_{\mathcal{Z}} \mathbf{z}_i]) \\ &= \sum_i (\mathbf{u}_\theta(\mathbf{x})^\top \mathbf{S}_\theta \mathbb{E}_y[\mathbf{v}_\theta(\mathbf{y}) \langle \mathbf{h}(\mathbf{y}), g^{-1} \triangleright_z \mathbf{z}_i \rangle_{\mathcal{Z}} \mathbf{z}_i]). \end{aligned}$$

Thus, changing the basis to $(g^{-1} \triangleright_z \mathbf{z}_i)_{i \in [n_{\mathcal{Z}}]}$ we obtain the result when $\mathbb{E}_y[\mathbf{h}(\mathbf{y})] = 0$. But since $\mathbb{1}_{\mathcal{X}}(g \triangleright_x \mathbf{x}) = 1$ for every $\mathbf{x} \in \mathcal{X}$ and $g \in \mathbb{G}$, the same holds for E_θ . \square

Recall that for the effective latent dimension m the true latent dimension is constrained by the dimensionality of the singular spaces, i.e., $r_m = \sum_{k \in [n_{\text{iso}}]} r_k = \sum_{k \in [n_{\text{iso}}]} m[k] d_k$. Further, given a measurable set $\mathbb{A} \subseteq \mathcal{X}$ and collection of group elements $\mathbb{G}' \subseteq \mathbb{G}$, let us define the following symmetry index of a set \mathbb{A} w.r.t. probability distribution of random variable \mathbf{x}

$$\gamma_{\mathbb{G}'}(\mathbb{A}) = \frac{1}{|\mathbb{G}'|(|\mathbb{G}'| - 1)} \sum_{\substack{g_1, g_2 \in \mathbb{G}' \\ g_1 \neq g_2}} \frac{\mathbb{P}[\mathbf{x} \in g_1 \triangleright \mathbb{A} \cap g_2 \triangleright \mathbb{A}]}{\mathbb{P}[\mathbf{x} \in \mathbb{A}]}, \quad (91)$$

which in the case when \mathbb{G}' is a subgroup of \mathbb{G} simplifies as

$$\gamma_{\mathbb{G}'}(\mathbb{A}) = \frac{1}{|\mathbb{G}'| - 1} \sum_{\substack{g \in \mathbb{G}' \\ g \neq e}} \frac{\mathbb{P}[\mathbf{x} \in \mathbb{A} \cap g \triangleright \mathbb{A}]}{\mathbb{P}[\mathbf{x} \in \mathbb{A}]}. \quad (92)$$

Observe that always $\gamma_{\mathbb{G}'}(\mathbb{A}) \in [0, 1]$, where extremes correspond to the cases $\gamma_{\mathbb{G}'}(\mathbb{A}) = 1$ when set \mathbb{A} is \mathbb{G}' invariant, and $\gamma_{\mathbb{G}'}(\mathbb{A}) = 0$ when \mathbb{A} equals its coset w.r.t. \mathbb{G}' , that is $g \triangleright \mathbb{A} \cap \mathbb{A} = \emptyset$ for all $g \in \mathbb{G}'$, meaning that the set is fully asymmetric w.r.t transformations $g \in \mathbb{G}'$.

We first generalize the approximation error bound in Lemma 1 from Kostic et al. (2024b) to the case of vector valued functions in the presence of symmetries.

Theorem M.2 (Approximation error). *Given a group of symmetries \mathbb{G} , let \mathcal{X} , \mathcal{Y} and \mathcal{Z} be Hilbert spaces endowed with symmetry group \mathbb{G} , and let P_x , P_y and P_{xy} be \mathbb{G} -invariant probability distributions on \mathcal{X} , \mathcal{Y} and $\mathcal{X} \times \mathcal{Y}$. Then, for every $\mathbf{h} \in \mathcal{L}_y^2(\mathcal{Y}, \mathcal{Z})$ it holds that*

$$\|\mathbb{E}_y[\mathbf{h}(\mathbf{y}) | \mathbf{x} = \cdot] - E_\theta \mathbf{h}\|_{\mathcal{L}_{P_x}^2(\mathcal{X}, \mathcal{Z})} \leq (\sigma_{r_m+1}^* + \|[\mathbb{D}_{y|x}]_{r_m} - D_\theta\|) \|\mathbf{h}\|_{\mathcal{L}_y^2(\mathcal{Y}, \mathcal{Z})}. \quad (93)$$

Moreover, denoting

$$E_\theta[f(\mathbf{y}) | \mathbf{x} \in \mathbb{A}] = \mathbb{E}_y[f] + \frac{\mathbb{E}_x[\mathbb{1}_{\mathbb{A}}(\mathbf{x}) [D_\theta f](\mathbf{x})]}{\mathbb{P}[\mathbf{x} \in \mathbb{A}]}, \quad (94)$$

if \mathbf{h} is either \mathbb{G}' -invariant or \mathbb{G}' -equivariant for some $\mathbb{G}' \subseteq \mathbb{G}$, then for every measurable set \mathbb{A}

$$\|\mathbb{E}[\mathbf{h}(\mathbf{y}) | \mathbf{x} \in \mathbb{A}] - E_\theta[\mathbf{h}(\mathbf{y}) | \mathbf{x} \in \mathbb{A}]\|_{\mathcal{Z}} \leq (\sigma_{r_m+1}^* + \|[\mathbb{D}_{y|x}]_{r_m} - D_\theta\|) \frac{\|\mathbf{h}\|_{\mathcal{L}_y^2(\mathcal{Y}, \mathcal{Z})}}{\sqrt{\mathbb{P}[\mathbf{x} \in \mathbb{A}]}} \sqrt{\frac{1 + (|\mathbb{G}'| - 1) \gamma_{\mathbb{G}'}(\mathbb{A})}{|\mathbb{G}'|}}. \quad (95)$$

2754 *Proof.* Start by observing that

$$2755 \quad \|\mathbb{E}[\mathbf{h}(\mathbf{y}) | \mathbf{x} = \cdot] - \mathbf{E}_\theta \mathbf{h}\|_{\mathcal{L}_{P_{\mathbf{x}}}^2(\mathcal{X}, \mathcal{Z})} \leq \|D_{y|x} - D_\theta\|_{\mathcal{L}_{\mathcal{Y}}^2(\mathcal{Y}, \mathcal{Z}) \rightarrow \mathcal{L}_{P_{\mathbf{x}}}^2(\mathcal{X}, \mathcal{Z})} \|\mathbf{h}\|_{\mathcal{L}_{\mathcal{Y}}^2(\mathcal{Y}, \mathcal{Z})}$$

$$2757 \quad = \|D_{y|x} - D_\theta\|_{\mathcal{L}_{P_{\mathbf{y}}}^2(\mathcal{Y}) \rightarrow \mathcal{L}_{P_{\mathbf{x}}}^2(\mathcal{X})} \|\mathbf{h}\|_{\mathcal{L}_{\mathcal{Y}}^2(\mathcal{Y}, \mathcal{Z})},$$

2759 where the equality holds since we extended operators $D_{y|x}$ and D_θ to vector valued setting as integral
2760 operators with the same scalar kernel. Hence, (93) readily follows.

2761 To prove (95), start with noting

$$2763 \quad \mathbb{E}[\mathbf{h}(\mathbf{y}) | \mathbf{x} \in \mathbb{A}] - \mathbf{E}_\theta[\mathbf{h}(\mathbf{y}) | \mathbf{x} \in \mathbb{A}] = \frac{\mathbb{E}_{\mathbf{x}}[\mathbb{1}_{\mathbb{A}}(\mathbf{x})[(D_{y|x} - D_\theta)\mathbf{h}](\mathbf{x})]}{\mathbb{P}[\mathbf{x} \in \mathbb{A}]}.$$

2765 Then, if \mathbf{h} is \mathbb{G} -equivariant, then, using that invariance of the probability distribution $P_{\mathbf{x}}$, \mathbb{G} -
2766 equivariance of $D_{y|x}$ and, due to Proposition M.1, of D_θ , we have that for every $g \in \mathbb{G}' \subseteq \mathbb{G}$

$$2767 \quad \mathbb{E}_{\mathbf{x}}[\mathbb{1}_{\mathbb{A}}(\mathbf{x})[(D_{y|x} - D_\theta)\mathbf{h}](\mathbf{x})] = \mathbb{E}_{\mathbf{x}}[\mathbb{1}_{\mathbb{A}}(g \triangleright_{\mathcal{X}} \mathbf{x})[(D_{y|x} - D_\theta)\mathbf{h}](g \triangleright_{\mathcal{X}} \mathbf{x})]$$

$$2769 \quad = \mathbb{E}_{\mathbf{x}}[\mathbb{1}_{g^{-1} \triangleright_{\mathcal{X}} \mathbb{A}}(\mathbf{x}) g \triangleright_{\mathcal{Z}} [(D_{y|x} - D_\theta)\mathbf{h}](\mathbf{x})]$$

$$2770 \quad = \mathbb{E}_{\mathbf{x}}[\mathbb{1}_{g^{-1} \triangleright_{\mathcal{X}} \mathbb{A}}(\mathbf{x}) \bar{\rho}_{\mathcal{Z}}(g) [(D_{y|x} - D_\theta)\mathbf{h}](\mathbf{x})].$$

2771 Hence, averaging over \mathbb{G}' we obtain

$$2773 \quad \mathbb{E}[\mathbf{h}(\mathbf{y}) | \mathbf{x} \in \mathbb{A}] - \mathbf{E}_\theta[\mathbf{h}(\mathbf{y}) | \mathbf{x} \in \mathbb{A}] = \mathbb{E}_{\mathbf{x}}[\mathbf{H}(\mathbf{x})\bar{\mathbf{z}}(\mathbf{x})],$$

2774 where

$$2775 \quad \mathbf{H}(\mathbf{x}) = \frac{1}{|\mathbb{G}'| \mathbb{P}[\mathbf{x} \in \mathbb{A}]} \sum_{g \in \mathbb{G}'} \mathbb{1}_{g^{-1} \triangleright_{\mathcal{X}} \mathbb{A}}(\mathbf{x}) \bar{\rho}_{\mathcal{Z}}(g) \quad \text{and} \quad \bar{\mathbf{z}}(\mathbf{x}) = [(D_{y|x} - D_\theta)\mathbf{h}](\mathbf{x}).$$

2778 Since due to Cauchy-Schwartz inequality we have

$$2780 \quad \|\mathbb{E}_{\mathbf{x}}[\mathbf{H}(\mathbf{x})\bar{\mathbf{z}}(\mathbf{x})]\|_{\mathcal{Z}}^2 \leq [\mathbb{E}_{\mathbf{x}} \|\mathbf{H}(\mathbf{x})\|_{\mathcal{Z} \rightarrow \mathcal{Z}}^2] [\mathbb{E}_{\mathbf{x}} \|\bar{\mathbf{z}}(\mathbf{x})\|_{\mathcal{Z}}^2] = \|\bar{\mathbf{z}}\|_{\mathcal{L}_{P_{\mathbf{x}}}^2(\mathcal{X}, \mathcal{Z})}^2 [\mathbb{E}_{\mathbf{x}} \|\mathbf{H}(\mathbf{x})\|_{\mathcal{Z} \rightarrow \mathcal{Z}}^2]$$

2782 and $\|\bar{\mathbf{z}}\|_{\mathcal{L}_{P_{\mathbf{x}}}^2(\mathcal{X}, \mathcal{Z})} \leq \|D_{y|x} - D_\theta\|_{\mathcal{L}_{\mathcal{Y}}^2(\mathcal{Y}, \mathcal{Z}) \rightarrow \mathcal{L}_{P_{\mathbf{x}}}^2(\mathcal{X}, \mathcal{Z})} \|\mathbf{h}\|_{\mathcal{L}_{\mathcal{Y}}^2(\mathcal{Y}, \mathcal{Z})}$, it remains to bound
2783 $\mathbb{E}_{\mathbf{x}} \|\mathbf{H}(\mathbf{x})\|_{\mathcal{Z} \rightarrow \mathcal{Z}}^2$. But, the group actions in the vector spaces are unitary, so using the \mathbb{G} -invariance
2784 of the distribution of \mathbf{x} we obtain

$$2785 \quad \mathbb{E}_{\mathbf{x}} \|\mathbf{H}(\mathbf{x})\|_{\mathcal{Z} \rightarrow \mathcal{Z}}^2 \leq \mathbb{E}_{\mathbf{x}} \left[\frac{1}{|\mathbb{G}'| \mathbb{P}[\mathbf{x} \in \mathbb{A}]} \sum_{g \in \mathbb{G}'} \mathbb{1}_{g^{-1} \triangleright_{\mathcal{X}} \mathbb{A}}(\mathbf{x}) \right]^2$$

$$2786 \quad = \frac{1}{|\mathbb{G}'|^2 \mathbb{P}[\mathbf{x} \in \mathbb{A}]^2} \sum_{g, g' \in \mathbb{G}'} \mathbb{E}_{\mathbf{x}}[\mathbb{1}_{g^{-1} \triangleright_{\mathcal{X}} \mathbb{A}}(\mathbf{x}) \mathbb{1}_{g'^{-1} \triangleright_{\mathcal{X}} \mathbb{A}}(\mathbf{x})]$$

$$2788 \quad = \frac{1}{|\mathbb{G}'|^2 \mathbb{P}[\mathbf{x} \in \mathbb{A}]^2} \sum_{g, g' \in \mathbb{G}'} \mathbb{E}_{\mathbf{x}}[\mathbb{1}_{g \triangleright_{\mathcal{X}} \mathbb{A} \cap g' \triangleright_{\mathcal{X}} \mathbb{A}}(\mathbf{x})]$$

$$2790 \quad = \frac{1}{|\mathbb{G}'|^2 \mathbb{P}[\mathbf{x} \in \mathbb{A}]} \sum_{g, g' \in \mathbb{G}'} \frac{\mathbb{P}[\mathbf{x} \in g \triangleright_{\mathcal{X}} \mathbb{A} \cap g' \triangleright_{\mathcal{X}} \mathbb{A}]}{\mathbb{P}[\mathbf{x} \in \mathbb{A}]}$$

$$2792 \quad = \frac{1}{\mathbb{P}[\mathbf{x} \in \mathbb{A}]} \frac{1 + (|\mathbb{G}'| - 1) \gamma_{\mathbb{G}'}(\mathbb{A})}{|\mathbb{G}'|},$$

2799 which completes the proof of (95) for \mathbb{G}' -equivariant functions. Finally, if f is \mathbb{G}' -invariant, the
2800 proof follows the same lines by replacing group actions ($\triangleright_{\mathcal{Z}}$) by their respective group representation
2801 $\rho_{\mathcal{Z}}$ (see Thm. I.3) with identity. \square

2802 Next we analyze the errors when, instead of applying learned operators \mathbf{E}_θ , we apply their empirical
2803 counterparts in inference tasks. To that end, we define now estimators of $\mathbb{E}[\mathbf{h}(\mathbf{x})]$ and $\mathbb{E}[\bar{\mathbf{z}}(\mathbf{y})]$
2804 exploiting the \mathbb{G} -invariance of the distributions of \mathbf{x} and \mathbf{y} . First, define the empirical \mathbb{G} -invariant
2805 distributions

$$2806 \quad \hat{\mathbb{P}}_{\mathbf{x}} := \frac{1}{|\mathbb{G}|N} \sum_{i=1}^N \sum_{g \in \mathbb{G}} \delta_{g \triangleright_{\mathcal{X}} \mathbf{x}_i}(\cdot), \quad \hat{\mathbb{P}}_{\mathbf{y}} := \frac{1}{|\mathbb{G}|N} \sum_{i=1}^N \sum_{g \in \mathbb{G}} \delta_{g \triangleright_{\mathcal{Y}} \mathbf{y}_i}(\cdot).$$

Hence we can define the equivariant empirical mean of any function $f \in \mathcal{L}_{\mathbf{x}}^2$, $h \in \mathcal{L}_{\mathbf{y}}^2$ as

$$\widehat{\mathbb{E}}_{\mathbf{x}}[f] = \frac{1}{|\mathbb{G}|N} \sum_{i=1}^N \sum_{g \in \mathbb{G}} f(g \triangleright_{\mathbf{x}} \mathbf{x}_i), \quad \widehat{\mathbb{E}}_{\mathbf{y}}[h] = \frac{1}{|\mathbb{G}|N} \sum_{i=1}^N \sum_{g \in \mathbb{G}} h(g \triangleright_{\mathbf{y}} \mathbf{y}_i). \quad (96)$$

This extends naturally to operator on a function space $\mathcal{L}_{\mathbf{y}}^2(\mathcal{Y}, \mathcal{Z})$ where \mathcal{Z} is endowed with an inner product $\langle \cdot, \cdot \rangle = \langle \cdot, \cdot \rangle_{\mathcal{Z}}$. If the distribution of \mathbf{y} is \mathbb{G}' -invariant, then for any $\mathbf{h} \in \mathcal{L}_{\mathbf{y}}^2(\mathcal{Y}, \mathcal{Z})$, we use the estimator $\widehat{\mathbb{E}}_{\mathbf{y}}[\mathbf{h}(\mathbf{y})]$ in (96) as an estimator of $\mathbb{E}[\mathbf{h}(\mathbf{y})]$:

$$\widehat{\mathbb{E}}_{\mathbf{y}}[\mathbf{h}] = \frac{1}{|\mathbb{G}|N} \sum_{i=1}^N \sum_{g \in \mathbb{G}} \mathbf{h}(g \triangleright_{\mathbf{y}} \mathbf{y}_i). \quad (97)$$

In this notation, we define our empirical estimators

$$[\mathbb{E}_{\theta} \mathbf{h}](\mathbf{x}) \approx \widehat{[\mathbb{E}_{\theta} \mathbf{h}]}(\mathbf{x}) = \widehat{\mathbb{E}}_{\mathbf{y}}[\mathbf{h}(\mathbf{y})] + \sum_{k \in [n_{\text{iso}}]} \sum_{i \in [m]} \sum_{j \in [d_k]} \sigma_i^{\theta^{(k)}} u_{i,j}^{\theta^{(k)}}(\mathbf{x}) \widehat{\mathbb{E}}_{\mathbf{y}}[v_{i,j}^{\theta^{(k)}} \mathbf{h}]$$

and

$$\mathbb{E}_{\theta}[\mathbf{h}(\mathbf{y}) \mid \mathbf{x} \in \mathbb{A}] \approx \widehat{\mathbb{E}}_{\theta}[\mathbf{h}(\mathbf{y}) \mid \mathbf{x} \in \mathbb{A}] = \widehat{\mathbb{E}}_{\mathbf{y}}[\mathbf{h}] + \sum_{k \in [n_{\text{iso}}]} \sum_{i \in [m]} \sum_{j \in [d_k]} \sigma_i^{\theta^{(k)}} \frac{\widehat{\mathbb{E}}_{\mathbf{x}}[u_{i,j}^{\theta^{(k)}} \mathbb{1}_{\mathbb{A}}]}{\widehat{\mathbb{E}}_{\mathbf{x}}[\mathbb{1}_{\mathbb{A}}]} \widehat{\mathbb{E}}_{\mathbf{y}}[v_{i,j}^{\theta^{(k)}} \mathbf{h}].$$

and, by choosing $\mathbf{h} = \mathbb{1}_{\mathbb{B}}$,

$$P[\mathbf{y} \in \mathbb{B} \mid \mathbf{x} \in \mathbb{A}] \approx \widehat{P}_{\theta}[\mathbf{y} \in \mathbb{B} \mid \mathbf{x} \in \mathbb{A}] = \widehat{\mathbb{E}}_{\mathbf{y}}[\mathbb{1}_{\mathbb{B}}] + \sum_{k \in [n_{\text{iso}}]} \sum_{i \in [m]} \sum_{j \in [d_k]} \sigma_i^{\theta^{(k)}} \frac{\widehat{\mathbb{E}}_{\mathbf{x}}[u_{i,j}^{\theta^{(k)}} \mathbb{1}_{\mathbb{A}}]}{\widehat{\mathbb{E}}_{\mathbf{x}}[\mathbb{1}_{\mathbb{A}}]} \widehat{\mathbb{E}}_{\mathbf{y}}[v_{i,j}^{\theta^{(k)}} \mathbb{1}_{\mathbb{B}}].$$

Direct consequence of the above construction which ensures that $\widehat{P}_{\mathbf{x}}$ and $\widehat{P}_{\mathbf{y}}$ are \mathbb{G} -invariant is the following result.

Proposition M.3. *Let $P_{\mathbf{x}}$ and $P_{\mathbf{y}}$ are \mathbb{G} -invariant, and D_{θ} from (88) is \mathbb{G} -equivariant model, and let $z \in \mathcal{L}_{\mathbf{x}}^2(\mathcal{X}, \mathbb{R})$ and $\mathbf{h} \in \mathcal{L}_{\mathbf{y}}^2(\mathcal{Y}, \mathcal{Z})$ be arbitrary. If for every $k \in [n_{\text{iso}}]$*

$$\left\{ \|\mathbb{D}_{\mathbf{y}|\mathbf{x}}^{(k)} - D_{\theta}^{(k)}\|, \left\| \mathbb{E}_{\mathbf{x}}[\mathbf{u}_{\theta^{(k)}}^{(1)}(\mathbf{x}) \mathbf{u}_{\theta^{(k)}}^{(1)}(\mathbf{x})^{\top}] - \mathbf{I}_m \right\|, \left\| \mathbb{E}_{\mathbf{y}}[\mathbf{v}_{\theta^{(k)}}^{(1)}(\mathbf{y}) \mathbf{v}_{\theta^{(k)}}^{(1)}(\mathbf{y})^{\top}] - \mathbf{I}_m \right\| \right\} \leq \mathcal{E}_{\theta}^{(k)}$$

holds with $\mathbf{u}_{\theta^{(k)}}^{(1)} = [u_{1,1}^{\theta^{(k)}} | \dots | u_{m,1}^{\theta^{(k)}}]^{\top} \in \mathbb{R}^m$ and $\mathbf{v}_{\theta^{(k)}}^{(1)} = [v_{1,1}^{\theta^{(k)}} | \dots | v_{m,1}^{\theta^{(k)}}]^{\top} \in \mathbb{R}^m$, and if

$$\begin{aligned} & \left\| \frac{\widehat{\mathbb{E}}_{\mathbf{x}}[\mathbf{u}_{\theta^{(k)}}^{(1)} z_1^{(k)}] - \mathbb{E}_{\mathbf{x}}[\mathbf{u}_{\theta^{(k)}}^{(1)}(\mathbf{x}) z_1^{(k)}(\mathbf{x})]}{\|z_1^{(k)}\|_{\mathcal{L}_{\mathbf{x}}^2}} \right\| \leq A(\mathbf{u}_{\theta}, z), \\ & \left\| \frac{\widehat{\mathbb{E}}_{\mathbf{y}}[\mathbf{v}_{\theta^{(k)}}^{(1)} \otimes \mathbf{h}_1^{(k)}] - \mathbb{E}_{\mathbf{y}}[\mathbf{v}_{\theta^{(k)}}^{(1)}(\mathbf{y}) \otimes \mathbf{h}_1^{(k)}(\mathbf{y})]}{\|\mathbf{h}_1^{(k)}\|_{\mathcal{L}_{\mathbf{y}}^2}} \right\| \leq A(\mathbf{v}_{\theta}, \mathbf{h}), \end{aligned} \quad (98)$$

where $z = \sum_{k \in [n_{\text{iso}}]} \sum_{j \in [d_k]} z_j^{(k)}$ and $\mathbf{h} = \sum_{k \in [n_{\text{iso}}]} \sum_{j \in [d_k]} \mathbf{h}_j^{(k)}$ are isospectral decompositions, then

$$\left\| \mathbb{E}_{\theta} \mathbf{h} - \widehat{\mathbb{E}}_{\theta} \mathbf{h} \right\|_{\mathcal{L}_{P_{\mathbf{x}}}^2(\mathcal{X}, \mathcal{Z})}^2 \leq \left\| \mathbb{E}_{\mathbf{y}}[\mathbf{h}(\mathbf{y})] - \widehat{\mathbb{E}}_{\mathbf{y}}[\mathbf{h}] \right\|_{\mathcal{Z}}^2 + \left[1 + \max_{k \in [n_{\text{iso}}]} \mathcal{E}_{\theta}^{(k)} \right]^3 \|\mathbf{h}\|_{\mathcal{L}_{\mathbf{y}}^2(\mathcal{Y}, \mathcal{Z})}^2 [A(\mathbf{v}_{\theta}, \mathbf{h})]^2. \quad (99)$$

Moreover, the empirical estimation error is upper bounded by

$$\begin{aligned} & \left\| \mathbb{E}_{\mathbf{x}}[z(\mathbf{x})[\mathbb{D}_{\theta} \mathbf{h}](\mathbf{x})] - \widehat{\mathbb{E}}_{\mathbf{x}}[z[\widehat{\mathbb{D}}_{\theta} \mathbf{h}]] \right\|_{\mathcal{Z}}^2 \leq \\ & (1 + \mathcal{E}_{\theta})^3 \left[A(\mathbf{u}_{\theta}, z) + A(\mathbf{v}_{\theta}, \mathbf{h}) + A(\mathbf{u}_{\theta}, z)A(\mathbf{v}_{\theta}, \mathbf{h}) \right]^2 \|z\|_{\mathcal{L}_{P_{\mathbf{x}}}^2(\mathcal{X})}^2 \|\mathbf{h}\|_{\mathcal{L}_{\mathbf{y}}^2(\mathcal{Y}, \mathcal{Z})}^2. \end{aligned} \quad (100)$$

2862 *Proof.* First, observe that due to \mathbb{G} -invariance of distribution $P_{\mathbf{x}\mathbf{y}}$ and \mathbb{G} -equivariance of E_θ and D_θ
 2863 we have that

$$2864 \mathbf{E}_\theta \mathbf{h} = \mathbb{E}_\mathbf{y}[\mathbf{h}^{(1)}(\mathbf{y})] + \sum_{k \in [n_{\text{iso}}]} D_\theta^{(k)} \mathbf{h}^{(k)}, \quad (101)$$

2865 and

$$2866 \mathbb{E}_\mathbf{x}[z(\mathbf{x})[\mathbf{E}_\theta \mathbf{h}](\mathbf{x})] = \mathbb{E}_\mathbf{x}[z^{(1)}(\mathbf{x})\mathbb{E}_\mathbf{y}[\mathbf{h}^{(1)}(\mathbf{y})] + \sum_{k \in [n_{\text{iso}}]} \mathbb{E}_\mathbf{x}[z^{(k)}(\mathbf{x})[D_\theta^{(k)} \mathbf{h}^{(k)}](\mathbf{x})]. \quad (102)$$

2870 In the same way, since the empirical distributions $\hat{P}_\mathbf{x}$ and $\hat{P}_\mathbf{y}$ are \mathbb{G} -invariant, we have that

$$2872 \hat{\mathbf{E}}_\theta \mathbf{h} = \hat{\mathbb{E}}_\mathbf{y}[\mathbf{h}^{(1)}] + \sum_{k \in [n_{\text{iso}}]} \hat{D}_\theta^{(k)} \mathbf{h}^{(k)}, \quad (103)$$

2873 and

$$2874 \hat{\mathbb{E}}_\mathbf{x}[z[\hat{\mathbf{E}}_\theta \mathbf{h}]] = \hat{\mathbb{E}}_\mathbf{x}[z^{(1)}]\hat{\mathbb{E}}_\mathbf{y}[\mathbf{h}^{(1)}] + \sum_{k \in [n_{\text{iso}}]} \hat{\mathbb{E}}_\mathbf{x}[z^{(k)}[\hat{D}_\theta^{(k)} \mathbf{h}^{(k)}]], \quad (104)$$

2875 where

$$2876 [\hat{D}_\theta^{(k)} \mathbf{h}^{(k)}](\mathbf{x}) = \mathbf{u}_\theta^{(k)}(\mathbf{x})^\top \mathbf{S}_\theta^{(k)} \hat{\mathbb{E}}_\mathbf{y}[\mathbf{v}_\theta^{(k)} \otimes \mathbf{h}^{(k)}]. \quad (105)$$

2877 Therefore, combining (101) and (103), we obtain that

$$2878 \begin{aligned} 2879 \mathbb{E}_\theta \mathbf{h}(\mathbf{x}) - \hat{\mathbb{E}}_\theta \mathbf{h}(\mathbf{x}) &= \left(\mathbb{E}_\mathbf{y}[\mathbf{h}^{(1)}(\mathbf{y})] - \hat{\mathbb{E}}_\mathbf{y}[\mathbf{h}^{(1)}] \right) \mathbb{1}_\mathcal{X}(\mathbf{x}) \\ 2880 &+ \sum_{k \in [n_{\text{iso}}]} \left([D_\theta^{(k)} \mathbf{h}^{(k)}](\mathbf{x}) - [\hat{D}_\theta^{(k)} \mathbf{h}^{(k)}](\mathbf{x}) \right), \end{aligned}$$

2881 which after taking the norm in $\mathcal{L}_{\hat{P}_\mathbf{x}}^2(\mathcal{X}, \mathcal{Z})$, due to orthonormality of isotypic subspaces gives

$$2882 \begin{aligned} 2883 \|\mathbb{E}_\theta \mathbf{h} - \hat{\mathbb{E}}_\theta \mathbf{h}\|_{\mathcal{L}_{\hat{P}_\mathbf{x}}^2(\mathcal{X}, \mathcal{Z})}^2 &= \left\| \mathbb{E}_\mathbf{y}[\mathbf{h}^{(1)}(\mathbf{y})] - \hat{\mathbb{E}}_\mathbf{y}[\mathbf{h}^{(1)}] \right\|_{\mathcal{Z}}^2 \\ 2884 &+ \sum_{k \in [n_{\text{iso}}]} \left\| [D_\theta^{(k)} \mathbf{h}^{(k)}] - [\hat{D}_\theta^{(k)} \mathbf{h}^{(k)}] \right\|_{\mathcal{L}_{\hat{P}_\mathbf{x}}^2(\mathcal{X}, \mathcal{Z})}^2. \end{aligned}$$

2885 Now, observe that, since

$$2886 [D_\theta^{(k)} \mathbf{h}^{(k)}](\mathbf{x}) - [\hat{D}_\theta^{(k)} \mathbf{h}^{(k)}](\mathbf{x}) = \mathbf{u}_\theta^{(k)}(\mathbf{x})^\top \mathbf{S}_\theta^{(k)} \left(\mathbb{E}_\mathbf{y}[\mathbf{v}_\theta^{(k)} \otimes \mathbf{h}^{(k)}] - \hat{\mathbb{E}}_\mathbf{y}[\mathbf{v}_\theta^{(k)} \otimes \mathbf{h}^{(k)}] \right)$$

2887 applying the norm we have that $\left\| [D_\theta^{(k)} \mathbf{h}^{(k)}] - [\hat{D}_\theta^{(k)} \mathbf{h}^{(k)}] \right\|_{\mathcal{L}_{\hat{P}_\mathbf{x}}^2(\mathcal{X}, \mathcal{Z})}^2$ equals

$$2888 \left(\mathbb{E}_\mathbf{y}[\mathbf{v}_\theta^{(k)} \otimes \mathbf{h}^{(k)}] - \hat{\mathbb{E}}_\mathbf{y}[\mathbf{v}_\theta^{(k)} \otimes \mathbf{h}^{(k)}] \right)^\top \mathbf{S}_\theta^{(k)} \left(\mathbb{E}_\mathbf{x}[\mathbf{u}_\theta^{(k)}(\mathbf{x})\mathbf{u}_\theta^{(k)}(\mathbf{x})^\top] \right) \mathbf{S}_\theta^{(k)} \left(\mathbb{E}_\mathbf{y}[\mathbf{v}_\theta^{(k)} \otimes \mathbf{h}^{(k)}] - \hat{\mathbb{E}}_\mathbf{y}[\mathbf{v}_\theta^{(k)} \otimes \mathbf{h}^{(k)}] \right)$$

2889 which using constraints within each isotypic block and

$$2900 \mathbb{E}_\mathbf{x}[\mathbf{u}_\theta^{(k)}(\mathbf{x})\mathbf{u}_\theta^{(k)}(\mathbf{x})^\top] \preceq \left\| \mathbb{E}_\mathbf{x}[\mathbf{u}_\theta^{(k)}(\mathbf{x})\mathbf{u}_\theta^{(k)}(\mathbf{x})^\top] \right\| \mathbf{I}_m \leq (1 + \mathcal{E}_\theta^{(k)}) \mathbf{I}_m,$$

2901 implies, due to (98), that

$$2902 \begin{aligned} 2903 \left\| D_\theta^{(k)} \mathbf{h}^{(k)} - \hat{D}_\theta^{(k)} \mathbf{h}^{(k)} \right\|_{\mathcal{L}_{\hat{P}_\mathbf{x}}^2(\mathcal{X}, \mathcal{Z})}^2 &\leq d_k (1 + \mathcal{E}_\theta^{(k)}) (\sigma_1^{\theta^{(k)}})^2 \\ 2904 &\cdot \left\| \mathbb{E}_\mathbf{y}[\mathbf{v}_\theta^{(k)}(1)(\mathbf{y}) \otimes \mathbf{h}_1^{(k)}(\mathbf{y})] - \hat{\mathbb{E}}_\mathbf{y}[\mathbf{v}_\theta^{(k)}(1) \otimes \mathbf{h}_1^{(k)}] \right\|_{\mathbb{R}^m \times \mathcal{Z}}^2 \\ 2905 &\leq d_k (1 + \mathcal{E}_\theta^{(k)}) (\sigma_1^{\theta^{(k)}})^2 [A(\mathbf{v}_\theta, \mathbf{h})]^2 \left\| \mathbf{h}_1^{(k)} \right\|_{\mathcal{L}_{\hat{P}_\mathbf{x}}^2(\mathcal{Y}, \mathcal{Z})}^2. \end{aligned}$$

2906 Therefore, bounding $\sigma_1^{\theta^{(k)}} \leq \sigma_1^{(k)} + |\sigma_1^{(k)} - \sigma_1^{\theta^{(k)}}| \leq 1 + \left\| D_{\mathbf{y}|\mathbf{x}}^{(k)} - D_\theta^{(k)} \right\|$ and summing over isotypic
 2907 components, since $\|\mathbf{h}\|_{\mathcal{L}_{\hat{P}_\mathbf{x}}^2(\mathcal{Y}, \mathcal{Z})}^2 = \sum_{k \in [n_{\text{iso}}], j \in [d_k]} \|\mathbf{h}_j^{(k)}\|_{\mathcal{L}_{\hat{P}_\mathbf{x}}^2(\mathcal{Y}, \mathcal{Z})}^2 = \sum_{k \in [n_{\text{iso}}]} d_k \|\mathbf{h}_1^{(k)}\|_{\mathcal{L}_{\hat{P}_\mathbf{x}}^2(\mathcal{Y}, \mathcal{Z})}^2$,
 2908 we complete the proof of (99).
 2909
 2910
 2911
 2912
 2913
 2914
 2915

To show (100), we combine (102) and (104), and obtain that $\mathbb{E}_{\mathbf{x}}[z(\mathbf{x})[\mathbf{D}_{\theta}\mathbf{h}](\mathbf{x})] - \widehat{\mathbb{E}}_{\mathbf{x}}[z[\widehat{\mathbf{D}}_{\theta}\mathbf{h}]]$ can be written as

$$\sum_{k \in [n_{\text{iso}}]} d_k \left[\mathbb{E}_{\mathbf{x}}[\mathbf{u}_{\theta}^{(k)}(\mathbf{x})z_1^{(k)}(\mathbf{x})]^{\top} \mathbf{S}_{\theta}^{(k)} \mathbb{E}_{\mathbf{y}}[\mathbf{v}_{\theta}^{(k)}(\mathbf{y}) \otimes \mathbf{h}_1^{(k)}(\mathbf{y})] - \widehat{\mathbb{E}}_{\mathbf{x}}[\mathbf{u}_{\theta}^{(k)}z_1^{(k)}]^{\top} \mathbf{S}_{\theta}^{(k)} \widehat{\mathbb{E}}_{\mathbf{y}}[\mathbf{v}_{\theta}^{(k)}(\mathbf{y}) \otimes \mathbf{h}_1^{(k)}] \right].$$

Adding and subtracting mixed terms we then obtain for each isotypic component, $\frac{1}{d_k} \mathbb{E}_{\mathbf{x}}[z(\mathbf{x})[\mathbf{D}_{\theta}^{(k)}\mathbf{h}](\mathbf{x})] - \widehat{\mathbb{E}}_{\mathbf{x}}[z[\widehat{\mathbf{D}}_{\theta}^{(k)}\mathbf{h}]]$ can be expressed as

$$\begin{aligned} & \mathbb{E}_{\mathbf{x}}[\mathbf{u}_{\theta}^{(k)}(\mathbf{x})z_1^{(k)}(\mathbf{x})]^{\top} \mathbf{S}^{(k)} \left(\mathbb{E}_{\mathbf{y}}[\mathbf{v}_{\theta}^{(k)}(\mathbf{y}) \otimes \mathbf{h}_1^{(k)}(\mathbf{y})] - \widehat{\mathbb{E}}_{\mathbf{y}}[\mathbf{v}_{\theta}^{(k)}(\mathbf{y}) \otimes \mathbf{h}_1^{(k)}] \right) \\ & + \left(\mathbb{E}_{\mathbf{x}}[\mathbf{u}_{\theta}^{(k)}(\mathbf{x})z_1^{(k)}(\mathbf{x})] - \widehat{\mathbb{E}}_{\mathbf{x}}[\mathbf{u}_{\theta}^{(k)}z_1^{(k)}] \right)^{\top} \mathbf{S}^{(k)} \mathbb{E}_{\mathbf{y}}[\mathbf{v}_{\theta}^{(k)}(\mathbf{y}) \otimes \mathbf{h}_1^{(k)}(\mathbf{y})] \\ & + \left(\mathbb{E}_{\mathbf{x}}[\mathbf{u}_{\theta}^{(k)}(\mathbf{x})z_1^{(k)}(\mathbf{x})] - \widehat{\mathbb{E}}_{\mathbf{x}}[\mathbf{u}_{\theta}^{(k)}z_1^{(k)}] \right)^{\top} \mathbf{S}^{(k)} \left(\mathbb{E}_{\mathbf{y}}[\mathbf{v}_{\theta}^{(k)}(\mathbf{y}) \otimes \mathbf{h}_1^{(k)}(\mathbf{y})] - \widehat{\mathbb{E}}_{\mathbf{y}}[\mathbf{v}_{\theta}^{(k)}(\mathbf{y}) \otimes \mathbf{h}_1^{(k)}] \right), \end{aligned}$$

and consequently bounded using (98) as

$$\begin{aligned} \left\| \mathbb{E}_{\mathbf{x}}[z(\mathbf{x})[\mathbf{D}_{\theta}^{(k)}\mathbf{h}](\mathbf{x})] - \widehat{\mathbb{E}}_{\mathbf{x}}[z[\widehat{\mathbf{D}}_{\theta}^{(k)}\mathbf{h}]] \right\|_{\mathcal{Z}} & \leq d_k \sigma_1^{\theta^{(k)}} \left[A(\mathbf{u}_{\theta}, z) + A(\mathbf{v}_{\theta}, \mathbf{h}) \right. \\ & \left. + A(\mathbf{u}_{\theta}, z)A(\mathbf{v}_{\theta}, \mathbf{h}) \right] \left\| z_1^{(k)} \right\|_{\mathcal{L}_{P_{\mathbf{x}}}^2(\mathcal{X})} \left\| \mathbf{h}_1^{(k)} \right\|_{\mathcal{L}_{\mathcal{Y}}^2(\mathcal{Y}, \mathcal{Z})}. \end{aligned}$$

Summing across isotypic components and bounding $\sigma_1^{\theta^{(k)}}$ as before, we complete the proof. \square

First note that coupling (99) with (93) ensures that we can prove regression bound via concentration result ensuring (98). To obtain similar result for set-wise regression, we set $z = \mathbb{1}_{\mathbb{A}}$ and use (100) to obtain the following.

Proposition M.4. *Under the assumptions of Proposition M.3, let $A(\mathbf{u}_{\theta}, \mathbb{1}_{\mathbb{A}})A(\mathbf{v}_{\theta}, \mathbf{h}) \leq A(\mathbf{u}_{\theta}, \mathbb{1}_{\mathbb{A}}) + A(\mathbf{v}_{\theta}, \mathbf{h})$. If*

$$\left| \mathbb{E}_{\mathbf{x}}[\mathbb{1}_{\mathbb{A}}(\mathbf{x})] - \widehat{\mathbb{E}}_{\mathbf{x}}[\mathbb{1}_{\mathbb{A}}] \right| / \mathbb{E}_{\mathbf{x}}[\mathbb{1}_{\mathbb{A}}(\mathbf{x})] \leq \eta_{\mathbb{A}} \quad (106)$$

and $\eta_{\mathbb{A}} < 1/2$, then

$$\begin{aligned} \left\| \mathbb{E}_{\theta}[\mathbf{h}(\mathbf{y}) | \mathbf{x} \in \mathbb{A}] - \widehat{\mathbb{E}}_{\theta}[\mathbf{h}(\mathbf{y}) | \mathbf{x} \in \mathbb{A}] \right\|_{\mathcal{Z}} & \leq \left\| \mathbb{E}_{\mathbf{y}}[\mathbf{h}] - \widehat{\mathbb{E}}_{\mathbf{y}}[\mathbf{h}] \right\|_{\mathcal{Z}} + \frac{2 \|\mathbf{h}\|_{\mathcal{L}_{P_{\mathbf{x}}}^2(\mathcal{Y}, \mathcal{Z})}}}{\sqrt{P[\mathbf{x} \in \mathbb{A}]}} \\ & \times \left[2(1 + \mathcal{E}_{\theta}) \left(A(\mathbf{u}_{\theta}, \mathbb{1}_{\mathbb{A}}) + A(\mathbf{v}_{\theta}, \mathbf{h}) \right) + \eta_{\mathbb{A}} \right], \end{aligned} \quad (107)$$

and for $\mathbf{h} = \mathbb{1}_{\mathbb{B}}$

$$\begin{aligned} |P[\mathbf{y} \in \mathbb{B} | \mathbf{x} \in \mathbb{A}] - \widehat{P}_{\theta}[\mathbf{y} \in \mathbb{B} | \mathbf{x} \in \mathbb{A}]| & \leq \left\| \mathbb{E}_{\mathbf{y}}[\mathbf{h}] - \widehat{\mathbb{E}}_{\mathbf{y}}[\mathbf{h}] \right\|_{\mathcal{Z}} \\ & + \frac{2}{\widehat{\mathbb{E}}_{\mathbf{x}}[\mathbf{h}]} \sqrt{\frac{P[\mathbf{y} \in \mathbb{B}]}{P[\mathbf{x} \in \mathbb{A}]}} \left[2(1 + \mathcal{E}_{\theta}) [A(\mathbf{u}_{\theta}, \mathbb{1}_{\mathbb{A}}) + A(\mathbf{v}_{\theta}, \mathbb{1}_{\mathbb{B}})] + \eta_{\mathbb{A}} \right]. \end{aligned} \quad (108)$$

Proof. Leveraging the representations in (102) and (104) with $z = \mathbb{1}_{\mathbb{A}}$, we get

$$\begin{aligned} \mathbb{E}_{\theta}[\mathbf{h}(\mathbf{y}) | \mathbf{x} \in \mathbb{A}] - \widehat{\mathbb{E}}_{\theta}[\mathbf{h}(\mathbf{y}) | \mathbf{x} \in \mathbb{A}] & = \mathbb{E}_{\mathbf{y}}[\mathbf{h}] - \widehat{\mathbb{E}}_{\mathbf{y}}[\mathbf{h}] + \frac{\mathbb{E}_{\mathbf{x}}[\mathbb{1}_{\mathbb{A}}(\mathbf{x})[\mathbf{D}_{\theta}\mathbf{h}](\mathbf{x})]}{\mathbb{E}_{\mathbf{x}}[\mathbb{1}_{\mathbb{A}}]} - \frac{\widehat{\mathbb{E}}_{\mathbf{x}}[\mathbb{1}_{\mathbb{A}}[\widehat{\mathbf{D}}_{\theta}\mathbf{h}]]}{\widehat{\mathbb{E}}_{\mathbf{x}}[\mathbb{1}_{\mathbb{A}}]} = \\ \mathbb{E}_{\mathbf{y}}[\mathbf{h}] - \widehat{\mathbb{E}}_{\mathbf{y}}[\mathbf{h}] + \mathbb{E}_{\mathbf{x}}[\mathbb{1}_{\mathbb{A}}(\mathbf{x})[\mathbf{D}_{\mathbf{y}|\mathbf{x}}\mathbf{h}](\mathbf{x})] & \left(\frac{1}{\mathbb{E}_{\mathbf{x}}[\mathbb{1}_{\mathbb{A}}(\mathbf{x})]} - \frac{1}{\widehat{\mathbb{E}}_{\mathbf{x}}[\mathbb{1}_{\mathbb{A}}]} \right) + \frac{\mathbb{E}_{\mathbf{x}}[\mathbb{1}_{\mathbb{A}}(\mathbf{x})[\mathbf{D}_{\theta}\mathbf{h}](\mathbf{x})] - \widehat{\mathbb{E}}_{\mathbf{x}}[\mathbb{1}_{\mathbb{A}}[\widehat{\mathbf{D}}_{\theta}\mathbf{h}]]}{\widehat{\mathbb{E}}_{\mathbf{x}}[\mathbb{1}_{\mathbb{A}}]}. \end{aligned}$$

2970 By triangular inequality applied to the norm in \mathcal{Z} , we get

$$\begin{aligned}
2971 & \left\| \mathbb{E}_\theta[\mathbf{h}(\mathbf{y}) | \mathbf{x} \in \mathbb{A}] - \widehat{\mathbb{E}}_\theta[\mathbf{h}(\mathbf{y}) | \mathbf{x} \in \mathbb{A}] \right\|_{\mathcal{Z}} \\
2972 & \leq \left\| \mathbb{E}_\mathbf{y}[\mathbf{h}] - \widehat{\mathbb{E}}_\mathbf{y}[\mathbf{h}] \right\|_{\mathcal{Z}} + \left\| \mathbb{E}[\mathbb{1}_\mathbb{A}(\mathbf{x})[\mathbf{D}_\theta \mathbf{f}(\mathbf{x})]] \right\|_{\mathcal{Z}} \left| \frac{1}{\mathbb{P}[\mathbf{x} \in \mathbb{A}]} - \frac{1}{\widehat{\mathbb{E}}_\mathbf{x}[\mathbb{1}_\mathbb{A}]} \right| + \frac{\left\| \mathbb{E}_\mathbf{x}[\mathbb{1}_\mathbb{A}(\mathbf{x})[\mathbf{D}_\theta \mathbf{h}(\mathbf{x})] - \widehat{\mathbb{E}}_\mathbf{x}[\mathbb{1}_\mathbb{A}[\widehat{\mathbf{D}}_\theta \mathbf{h}]] \right\|_{\mathcal{Z}}}{\widehat{\mathbb{E}}_\mathbf{x}[\mathbb{1}_\mathbb{A}]} \\
2973 & \leq \left\| \mathbb{E}_\mathbf{y}[\mathbf{h}] - \widehat{\mathbb{E}}_\mathbf{y}[\mathbf{h}] \right\|_{\mathcal{Z}} + \left\| \mathbb{E}[\mathbb{1}_\mathbb{A}(\mathbf{x})[\mathbf{D}_\theta \mathbf{h}(\mathbf{x})]] \right\|_{\mathcal{Z}} \frac{2\eta_\mathbb{A}}{\mathbb{P}[\mathbf{x} \in \mathbb{A}]} + \frac{\left\| \mathbb{E}_\mathbf{x}[\mathbb{1}_\mathbb{A}(\mathbf{x})[\mathbf{D}_\theta \mathbf{h}(\mathbf{x})] - \widehat{\mathbb{E}}_\mathbf{x}[\mathbb{1}_\mathbb{A}[\widehat{\mathbf{D}}_\theta \mathbf{h}]] \right\|_{\mathcal{Z}}}{\widehat{\mathbb{E}}_\mathbf{x}[\mathbb{1}_\mathbb{A}]},
\end{aligned}$$

2978 where we have used Condition (106) in the last line to get that

$$\left| \frac{1}{\mathbb{P}[\mathbf{x} \in \mathbb{A}]} - \frac{1}{\widehat{\mathbb{E}}_\mathbf{x}[\mathbb{1}_\mathbb{A}]} \right| \leq \frac{\eta_\mathbb{A}}{(1 - \eta_\mathbb{A})\mathbb{P}[\mathbf{x} \in \mathbb{A}]} \leq \frac{2\eta_\mathbb{A}}{\mathbb{P}[\mathbf{x} \in \mathbb{A}]}.$$

2982 From Proposition M.3 and Condition (106) we get that

$$\frac{1}{\widehat{\mathbb{E}}_\mathbf{x}[\mathbb{1}_\mathbb{A}]} \left\| \mathbb{E}_\mathbf{x}[\mathbb{1}_\mathbb{A}(\mathbf{x})[\mathbf{D}_\theta \mathbf{h}(\mathbf{x})] - \widehat{\mathbb{E}}_\mathbf{x}[\mathbb{1}_\mathbb{A}(\mathbf{x})[\widehat{\mathbf{D}}_\theta \mathbf{h}]] \right\|_{\mathcal{Z}} \leq \frac{2(1+\mathcal{E}_\theta) \left[A(\mathbf{u}_\theta, \mathbb{1}_\mathbb{A}) A(\mathbf{v}_\theta, \mathbf{h}) \right]}{\mathbb{P}[\mathbf{x} \in \mathbb{A}]} \|\mathbb{1}_\mathbb{A}\|_{\mathcal{L}_{P_\mathbf{x}}^2} \|\mathbf{h}\|_{\mathcal{L}_{P_\mathbf{y}}^2(\mathcal{Y}, \mathcal{Z})}$$

2987 Cauchy's Schwarz's inequality again and $\|\mathbf{D}_\theta\| \leq 1$ give

$$\left\| \mathbb{E}[\mathbb{1}_\mathbb{A}(\mathbf{x})[\mathbf{D}_\theta \mathbf{h}(\mathbf{x})]] \right\|_{\mathcal{Z}} \leq \|\mathbb{1}_\mathbb{A}\|_{\mathcal{L}_{P_\mathbf{x}}^2} \|\mathbf{D}_\theta\| \|\mathbf{h}\|_{\mathcal{L}_{P_\mathbf{y}}^2} \leq \|\mathbb{1}_\mathbb{A}\|_{\mathcal{L}_{P_\mathbf{x}}^2} \|\mathbf{h}\|_{\mathcal{L}_{P_\mathbf{y}}^2} = \sqrt{\mathbb{P}[\mathbf{x} \in \mathbb{A}]} \|\mathbf{h}\|_{\mathcal{L}_{P_\mathbf{y}}^2(\mathcal{Y}, \mathcal{Z})}.$$

2990 Combining the last four displays give the first result. The second result follows immediately for $\mathbf{h} = \mathbb{1}_\mathbb{B}$. \square

2993 Consequence of this result is that we can bound the error in probability as we can derive concentra-
2994 tion inequalities on the terms in (98) and (106). Then an union bound gives the estimation result for
2995 regression conditional on sets.

2996 Next, we recall that $\mathbb{E}_{\mathbf{y}|\mathbf{x}}$ being $(1/\alpha)$ -Schatten class operator, implies:

2997 **Assumption M.5.** *Let there exist some constant $c > 0$ such that for $\alpha > 0$, any $i \geq 1$ and any*
2998 *$k \in [n_{\text{iso}}]$, we have $\sigma_i^{(k)} \leq c i^{-\alpha}$.*

3000 Further, for any $\mathbf{h} \in \mathcal{L}_\mathbf{y}^2(\mathcal{Y}, \mathcal{Z})$, we define $\bar{\mathbf{h}}(\mathbf{y}) = \mathbf{h}(\mathbf{y}) - \mathbb{E}[\mathbf{h}(\mathbf{y})]$ and

$$\gamma_{\mathbb{G}'}(\mathbf{h}) := \frac{1}{|\mathbb{G}'| - 1} \sum_{\substack{g \in \mathbb{G}' \\ g \neq e}} \mathbb{E}[\langle \bar{\mathbf{h}}(\mathbf{y}), \bar{\mathbf{h}}(g \triangleright \mathbf{y}) \rangle]. \quad (109)$$

3005 In the following, we consider observables h satisfying the following condition (that is clearly satis-
3006 fied for an indicator of a set of positive measure)

3007 **Assumption M.6.** *Let there exists an absolute constant $C_0 \geq 1$ such that $(|\mathbb{G}'| - 1)\gamma_{\mathbb{G}'}(h) \leq$
3008 $C_0 \mathbb{E}[\|\mathbf{h}(\mathbf{y})\|_{\mathcal{Z}}^2]$.*

3010 Define

$$\eta_\mathbb{A} = \eta_\mathbb{A}(\delta) := \left(\frac{1 - \mathbb{P}[\mathbf{x} \in \mathbb{G} \triangleright \mathbb{A}]}{\mathbb{P}[\mathbf{x} \in \mathbb{G} \triangleright \mathbb{A}]} \right) \frac{\log 2\delta^{-1}}{N} + \sqrt{2 \frac{\log 2\delta^{-1}}{N}} \sqrt{\frac{1 - \mathbb{P}[\mathbf{x} \in \mathbb{G} \triangleright \mathbb{A}]}{\mathbb{P}[\mathbf{x} \in \mathbb{G} \triangleright \mathbb{A}]}}.$$

3014 **Theorem M.7.** *Let Assumptions M.6 and M.5 be satisfied. Let $P_\mathbf{x}$ and $P_\mathbf{y}$ are \mathbb{G} -invariant, and \mathbf{D}_θ
3015 from (88) is \mathbb{G} -equivariant model, and let $\mathbf{h} \in \mathcal{L}_\mathbf{y}^2(\mathcal{Y}, \mathcal{Z})$ and $\mathbf{f} \in \mathcal{L}_\mathbf{x}^2(\mathcal{X}, \mathcal{Z})$ (with values in \mathcal{Z}) be
3016 subGaussian random variables. Assume in addition that the event \mathbb{A} is anti-symmetric for \mathbb{G} and
3017 that $m_k = m$ for all $k \in [n_{\text{iso}}]$. Assume that $N \geq |\mathbb{G}|$. Then for any $\delta \in (0, 1)$, it holds w.p.a.1 $1 - \delta$*

$$\left\| \mathbb{E}[\mathbf{h}(\mathbf{y}) | \mathbf{x} \in \mathbb{A}] - \widehat{\mathbb{E}}_\theta[\mathbf{h}(\mathbf{y}) | \mathbf{x} \in \mathbb{A}] \right\|_{\mathcal{Z}} \lesssim_{C_0} \frac{\|\mathbf{h}\|_{\mathcal{L}_\mathbf{y}^2(\mathcal{Y}, \mathcal{Z})}}{\sqrt{\mathbb{P}[\mathbf{x} \in \mathbb{G} \triangleright \mathbb{A}]}} \left(\mathcal{E}_\theta + \frac{\log(2n_{\text{iso}}\delta^{-1})}{(d_{\text{iso}}N)^{\frac{1}{1+2\alpha}}} \right),$$

3021 and

$$\left| \mathbb{P}(\mathbf{y} \in \mathbb{B} | \mathbf{x} \in \mathbb{A}) - \widehat{\mathbb{P}}_\theta(\mathbf{y} \in \mathbb{B} | \mathbf{x} \in \mathbb{A}) \right| \lesssim_{C_0} \sqrt{\frac{\mathbb{P}[\mathbf{y} \in \mathbb{B}]}{\mathbb{P}[\mathbf{x} \in \mathbb{G} \triangleright \mathbb{A}]}} \left(\mathcal{E}_\theta + \frac{\log(2n_{\text{iso}}\delta^{-1})}{(d_{\text{iso}}N)^{\frac{1}{1+2\alpha}}} + \sqrt{|\mathbb{G}|\eta_\mathbb{A}} \right).$$

Proof. This result follows immediately from Propositions M.3 and M.4 combined with Lemmas M.9 and Lemma M.10. Set

$$A(\mathbf{u}_\theta, \mathbf{f}) := C \sqrt{\frac{1}{|\mathbb{G}'|N}} \sqrt{C_0 \vee \frac{|\mathbb{G}'|}{N}} \log(2n_{\text{iso}}\delta^{-1}),$$

$$A(\mathbf{v}_\theta, \mathbf{h}) := C \sqrt{\frac{\max_{k \in [n_{\text{iso}}]} \{m_k\}}{|\mathbb{G}'|N}} \sqrt{C_0 \vee \frac{|\mathbb{G}'|}{N}} \log(2n_{\text{iso}}\delta^{-1}),$$

for some large enough absolute constant $C > 0$.

Then an union bound based on Lemmas M.9 and M.10 guarantees that (107) is satisfied w.p.a.l. $1 - \delta$ (up to a rescaling of the constant C):

$$\left\| \mathbb{E}_\theta[\mathbf{h}(\mathbf{y}) | \mathbf{x} \in \mathbb{A}] - \widehat{\mathbb{E}}_\theta[\mathbf{h}(\mathbf{y}) | \mathbf{x} \in \mathbb{A}] \right\|_{\mathcal{Z}} \leq C \frac{\|\mathbf{h}\|_{\mathcal{L}_{\mathbf{x}}^2(\mathcal{X}, \mathcal{Z})}}{\sqrt{P[\mathbf{x} \in \mathbb{A}]}} \left[2(1 + \mathcal{E}_\theta) \left(\sqrt{\frac{\max_{k \in [n_{\text{iso}}]} \{m_k\}}{|\mathbb{G}'|N}} \sqrt{C_0 \vee \frac{|\mathbb{G}'|}{N}} \log(2n_{\text{iso}}\delta^{-1}) \right) \right].$$

Next we use our bound on the representation bias in (95)

$$\|\mathbb{E}[\mathbf{y}(\mathbf{y}) | \mathbf{x} \in \mathbb{A}] - \mathbb{E}_\theta[\mathbf{y}(\mathbf{y}) | \mathbf{x} \in \mathbb{A}]\|_{\mathcal{Z}} \leq (\sigma_{r_m+1}^* + \mathcal{E}_\theta) \frac{\|\mathbf{h}\|_{\mathcal{L}_{\mathbf{y}}^2(\mathcal{Y}, \mathcal{Z})}}{\sqrt{P[\mathbf{x} \in \mathbb{A}]}} \sqrt{\frac{1 + (|\mathbb{G}'| - 1)\gamma_{\mathbb{G}}(\mathbb{A})}{|\mathbb{G}'|}}. \quad (110)$$

Recall that $\mathcal{E}_\theta = \max_{k \in [n_{\text{iso}}]} \{\mathcal{E}_\theta^{(k)}\}$. Under Assumption M.5, we have $\|\llbracket D_{\mathbf{y}|\mathbf{x}} \rrbracket_{r_m} - D_\theta\| \leq \frac{1}{(d_{\text{iso}}m)^\alpha}$. In addition, $(|\mathbb{G}'| - 1)\gamma_{\mathbb{G}}(\mathbb{A}) \leq C_0$ under Assumption M.6.

Combining the last two display gives w.p.a.l $1 - \delta$

$$\left\| \mathbb{E}[\mathbf{h}(\mathbf{y}) | \mathbf{x} \in \mathbb{A}] - \widehat{\mathbb{E}}_\theta[\mathbf{h}(\mathbf{y}) | \mathbf{x} \in \mathbb{A}] \right\|_{\mathcal{Z}} \lesssim_{C_0} \frac{\|\mathbf{h}\|_{\mathcal{L}_{\mathbf{y}}^2(\mathcal{Y}, \mathcal{Z})}}{\sqrt{P[\mathbf{x} \in \mathbb{G}_{\triangleright, \mathcal{X}} \mathbb{A}]}} \left(\mathcal{E}_\theta + \frac{1}{(d_{\text{iso}}m)^\alpha} + \sqrt{\frac{m}{N}} \log(2n_{\text{iso}}\delta^{-1}) \right).$$

Balancing the previous display w.r.t. dimension m , we get that $m \asymp (d_{\text{iso}}^{-2\alpha} N)^{\frac{1}{1+2\alpha}}$ and the first result follows.

The bound for the conditional probability follows by picking $\mathbf{y} = \mathbf{1}_{\mathbb{B}}$. \square

M.1 QUADRATIC ERROR REGRESSION BOUND

Our goal is to estimate the conditional expectation function

$$z(\mathbf{x}) = \mathbb{E}[\mathbf{h}(\mathbf{y}) | \mathbf{x} = \mathbf{x}] = \mathbb{E}[\mathbf{h}(\mathbf{y})] + [D_{\mathbf{y}|\mathbf{x}} \mathbf{h}](\mathbf{x}).$$

Our estimator is

$$\widehat{\mathbf{z}}_\theta(\cdot) = \widehat{\mathbb{E}}_\theta[\mathbf{y}] + [\widehat{D}_\theta \mathbf{h}](\cdot).$$

Theorem M.8. *Assume that Y is a sub-Gaussian random vector. Let Assumption M.5 be satisfied. Assume in addition that $\mathcal{E}_\theta \leq 1$, $m_k = m$ for all $k \in [n_{\text{iso}}]$. Then for any $\delta \in (0, 1)$ such that $N \geq (c_u \vee c_v)^2 m \log(e\delta^{-1}n_{\text{iso}}) \vee |\mathbb{G}|$, it holds w.p.a.l. $1 - \delta$*

$$\|\mathbf{z} - \widehat{\mathbf{z}}_\theta\|_{\mathcal{L}_{\mathbf{x}}^2(\mathcal{X}, \mathcal{Z})}^2 \lesssim \text{Tr}(\text{Cov}(Y)) \left(\mathcal{E}_\theta^2 + (d_{\text{iso}} |\mathbb{G}| N)^{\frac{-2\alpha}{1+2\alpha}} \log^2(\delta^{-1}n_{\text{iso}}) \right). \quad (111)$$

Discussion When the training of the NN is successful, we expect the statistical rate to dominate the optimization error $\max_{k \in [n_{\text{iso}}]} \{\mathcal{E}_\theta^{(k)}\}$ for large enough sample size N . For distribution containing symmetry invariants with large isotopic components (m is large), we observe that exploiting this information in the construction of the NCP operator yields a substantial improvement in the statistical error rate as we go from a rate $N^{-\frac{\alpha}{1+2\alpha}}$ for standard NCP to $(Nm)^{-\frac{\alpha}{1+2\alpha}}$ for eNCP.

Proof. Combining (99) with Lemma M.9 gives w.p.a.l. $1 - \delta$

$$\begin{aligned} \left\| \mathbb{E}_\theta \mathbf{y} - \widehat{\mathbb{E}}_\theta \mathbf{y} \right\|_{\mathcal{L}_{\mathbf{x}}^2(\mathcal{X})}^2 &\lesssim (1 + \mathcal{E}_\theta)^3 \text{Tr}(\text{Cov}(\mathbf{y})) \frac{m}{|\mathbb{G}|N} \log^2(2n_{\text{iso}}\delta^{-1}) \\ &\lesssim \text{Tr}(\text{Cov}(\mathbf{y})) \frac{m}{|\mathbb{G}|N} \log^2(n_{\text{iso}}\delta^{-1}), \end{aligned}$$

3078 provide that $\mathcal{E}_\theta \leq 1$. We derived in (93) an upper bound on the bias term

$$3079 \quad \|\mathbb{E}_{\mathbf{y}|\mathbf{x}}[\mathbf{y} | \mathbf{x} = \cdot] - \mathbb{E}_\theta \mathbf{y}\|_{\mathcal{L}_{P_{\mathbf{x}}}^2(\mathcal{X}, \mathcal{Z})}^2 \leq \text{Tr}(\text{Cov}(\mathbf{y})) \left(\frac{1}{(d_{\text{iso}} m)^{2\alpha}} + \mathcal{E}_\theta^2 \right). \quad (112)$$

3082 Balancing the two bounds in the last two displays w.r.t. $m \asymp (|\mathbb{G}| d_{\text{iso}} N)^{\frac{1}{1+2\alpha}}$, we get the result. \square

3084 M.2 AUXILIARY RESULTS.

3086 Consider the function space $\mathcal{L}_{\mathbf{y}}^2(\mathcal{Y}, \mathcal{Z})$ where \mathcal{Z} is endowed with an inner product $\langle \cdot, \cdot \rangle = \langle \cdot, \cdot \rangle_{\mathcal{Z}}$. If
3087 the distribution of \mathbf{y} is \mathbb{G}' -invariant, then for any $h \in \mathcal{L}_{\mathbf{y}}^2(\mathcal{Y}, \mathcal{Z})$, we use the estimator $\widehat{\mathbb{E}}_{\mathbf{y}}[h]$ in (96)
3088 as an estimator of $\mathbb{E}[h(\mathbf{y})]$.

3089 **Lemma M.9.** *Assume that the distribution $P_{\mathbf{y}}$ of \mathbf{y} is \mathbb{G} -invariant and let $\mathbb{G}' \leq \mathbb{G}$. Let there exists
3090 a function $\mathbf{h} \in \mathcal{L}_{\mathbf{y}}^2(\mathcal{Y}, \mathcal{Z})$ such that $\mathbf{h}(\mathbf{y})$ is subGaussian. Then there exists an absolute constant
3091 $C > 0$ such that for any $\delta \in (0, 1)$, it holds w.p.a.l. $1 - \delta$*

$$3092 \quad \left\| \widehat{\mathbb{E}}_{\mathbf{y}}[\mathbf{h}] - \mathbb{E}[\mathbf{h}(\mathbf{y})] \right\|_{\mathcal{Z}} \leq C \sqrt{\frac{\log^2 2\delta^{-1}}{|\mathbb{G}'| N}} \sqrt{\mathbb{E}[\|\overline{\mathbf{h}}(\mathbf{y})\|_{\mathcal{Z}}^2] + (|\mathbb{G}'| - 1)\gamma_{\mathbb{G}'}(\mathbf{h}) + \frac{|\mathbb{G}'| \mathbb{E}[\|\overline{\mathbf{h}}(\mathbf{y})\|_{\mathcal{Z}}^2]}{N}}.$$

3096 Assume in addition that there exists an absolute constant $C_0 \geq 1$ such that $(|\mathbb{G}'| - 1)\gamma_{\mathbb{G}'}(\overline{\mathbf{h}}) \leq$
3097 $C_0 \mathbb{E}[\|\mathbf{h}(\mathbf{y})\|_{\mathcal{Z}}^2]$. Then for any $\delta \in (0, 1)$, it holds w.p.a.l. $1 - \delta$

$$3098 \quad \left\| \widehat{\mathbb{E}}_{\mathbf{y}}[\mathbf{h}] - \mathbb{E}[\mathbf{h}(\mathbf{y})] \right\|_{\mathcal{Z}} \leq C \sqrt{\frac{\mathbb{E}[\|\overline{\mathbf{h}}(\mathbf{y})\|_{\mathcal{Z}}^2]}{|\mathbb{G}'| N}} \sqrt{(1 + C_0) + \frac{|\mathbb{G}'|}{N} \log 2\delta^{-1}}.$$

3101 Note that similar bounds hold valid for the \mathbb{G} -invariant distribution $P_{\mathbf{x}}$ and any function $\mathbf{f} \in$
3102 $\mathcal{L}_{P_{\mathbf{x}}}^2(\mathcal{X}, \mathcal{Z})$ such that $\mathbf{f}(\mathbf{x})$ is subGaussian.

3104 *Proof.* We note that

$$3105 \quad \widehat{\mathbb{E}}_{\mathbf{y}}[\mathbf{h}] - \mathbb{E}[\mathbf{h}(\mathbf{y})] = \frac{1}{N} \sum_{i=1}^N Z_i \quad \text{with} \quad Z_i = \frac{1}{|\mathbb{G}'|} \sum_{g \in \mathbb{G}'} \mathbf{h}(g \triangleright_{\mathbf{y}} y_i) - \mathbb{E}_{y_i}[\mathbf{h}(g \triangleright_{\mathbf{y}} y_i)], \quad \forall i \in [N].$$

3109 Define

$$3110 \quad Z := \frac{1}{|\mathbb{G}'|} \sum_{g \in \mathbb{G}'} \mathbf{h}(g \triangleright_{\mathbf{y}} \mathbf{y}) - \mathbb{E}_{\mathbf{y}}[\mathbf{h}(g \triangleright_{\mathbf{y}} \mathbf{y})], \quad (113)$$

3112 and, for brevity, set $\|z\| = \|z\|_{\mathcal{Z}} = \sqrt{\langle z, z \rangle_{\mathcal{Z}}}$ for any $z \in \mathcal{Z}$. We apply Proposition M.12, to get
3113 w.p.a.l. $1 - \delta$

$$3114 \quad \left\| \widehat{\mathbb{E}}_{\mathbf{y}}[\mathbf{h}] - \mathbb{E}_{\mathbf{y}}[\mathbf{h}(\mathbf{y})] \right\| \leq \frac{4\sqrt{2}}{\sqrt{N}} \sqrt{\text{Var}_{\mathbf{y}}(\|Z\|) + \frac{\|Z\|_{\psi_2}^2}{N} \log \frac{2}{\delta}}.$$

3117 Using the triangular inequality successively on $\|\cdot\|$ and $\|\cdot\|_{\psi_2}$ and the \mathbb{G}' -invariance of $P_{\mathbf{y}}$,
3118 $\|\overline{\mathbf{h}}(g \triangleright_{\mathbf{y}} \mathbf{y})\|_{\psi_2} = \|\overline{\mathbf{h}}(\mathbf{y})\|_{\psi_2}$ for any $g \in \mathbb{G}'$, we get that

$$3119 \quad \|\|Z\|\|_{\psi_2} \lesssim \|\|\overline{\mathbf{h}}(\mathbf{y})\|\|_{\psi_2}.$$

3120 We note next that $\|\overline{\mathbf{h}}(\mathbf{y})\|$ is subGaussian. Consequently the well-known property of equivalence
3121 of moments for subGaussian distributions gives $\|Z\|_{\psi_2} \lesssim \|\|\overline{\mathbf{h}}(\mathbf{y})\|\|_{\psi_2} \lesssim \mathbb{E}[\|\overline{\mathbf{h}}(\mathbf{y})\|^2]$. We derive
3122 now a control on $\text{Var}_{\mathbf{y}}(\|Z\|) \leq \mathbb{E}[\|Z\|^2]$. Using the \mathbb{G}' -invariance of $P_{\mathbf{y}}$, we get

$$3123 \quad \text{Var}(\|Z\|) \leq \frac{\mathbb{E}[\|\overline{\mathbf{h}}(\mathbf{y})\|^2]}{|\mathbb{G}'|} + \frac{1}{|\mathbb{G}'|} \sum_{\substack{g \in \mathbb{G}' \\ g \neq e}} \mathbb{E}[\langle \mathbf{h}(\mathbf{y}) - \mathbb{E}[\mathbf{h}(\mathbf{y})], \mathbf{h}(g \triangleright_{\mathbf{y}} \mathbf{y}) - \mathbb{E}[\mathbf{h}(\mathbf{y})] \rangle]$$

$$3124 \quad = \frac{\mathbb{E}[\|\overline{\mathbf{h}}(\mathbf{y})\|^2]}{|\mathbb{G}'|} + \frac{(|\mathbb{G}'| - 1)\gamma_{\mathbb{G}'}(\mathbf{h})}{|\mathbb{G}'|} \leq (1 + C_0) \frac{\mathbb{E}[\|\overline{\mathbf{h}}(\mathbf{y})\|^2]}{|\mathbb{G}'|}. \quad (114)$$

Hence we get the result. \square

We focus now on a concentration bound for indicator functions $z = \mathbb{1}_A$ for any event $A \in \Sigma_{\mathcal{X}}$. We define

$$\begin{aligned} Z_A &:= \widehat{\mathbb{E}}_{\mathbf{x}}[\mathbb{1}_A] - \mathbb{P}[\mathbf{x} \in A] = \frac{1}{|\mathbb{G}'|} \sum_{g \in \mathbb{G}'} (\mathbb{1}_{g^{-1} \triangleright_{\mathcal{X}} A}(\mathbf{x}) - \mathbb{E}[\mathbb{1}_{g^{-1} \triangleright_{\mathcal{X}} A}(\mathbf{x})]) \\ &= \frac{1}{|\mathbb{G}'|} \sum_{g \in \mathbb{G}'} (\mathbb{1}_{g^{-1} \triangleright_{\mathcal{X}} A}(\mathbf{x}) - \mathbb{P}[\mathbf{x} \in A]) = \left(\frac{1}{|\mathbb{G}'|} \sum_{g \in \mathbb{G}'} \mathbb{1}_{g^{-1} \triangleright_{\mathcal{X}} A}(\mathbf{x}) \right) - \mathbb{P}[\mathbf{x} \in A]. \end{aligned} \quad (115)$$

Note that we always have $|Z_A| \leq 1$ but this bound can be quite conservative as we could get a much sharper bound for some events A . We denote by $\gamma_{\mathbb{G}', \infty}(A)$ the smallest deterministic upper-bound on $\frac{1}{|\mathbb{G}'|} \sum_{g \in \mathbb{G}'} \mathbb{1}_{g^{-1} \triangleright_{\mathcal{X}} A}(\mathbf{x})$ (For instance when A is an antisymmetric event, then we have $\gamma_{\mathbb{G}', \infty}(A) = 1/|\mathbb{G}'|$). Then we have

$$-\mathbb{P}[\mathbf{x} \in A] \leq Z_A \leq \gamma_{\mathbb{G}', \infty}(A) - \mathbb{P}[\mathbf{x} \in A]. \quad (116)$$

Define also

$$\Upsilon_{\mathbb{G}', \mathcal{X}}(A) := \mathbb{P}(\mathbf{x} \in A)(1 - \mathbb{P}(\mathbf{x} \in A)) + (|\mathbb{G}'| - 1) (\gamma_{\mathbb{G}'}(A) - \mathbb{P}[\mathbf{x} \in A])\mathbb{P}[\mathbf{x} \in A]. \quad (117)$$

Lemma M.10. *Let the distribution of \mathbf{x} be \mathbb{G}' -invariant. Then for any $A \in \Sigma_{\mathcal{X}}$ and any $\delta \in (0, 1)$, it holds w.p.a.l. $1 - \delta$*

$$|\widehat{\mathbb{E}}_{\mathbf{x}}[\mathbb{1}_A] - \mathbb{E}[\mathbb{1}_A(\mathbf{x})]| \leq |\gamma_{\mathbb{G}', \infty}(A) - \mathbb{P}[\mathbf{x} \in A]| \frac{\log 2\delta^{-1}}{N} + \sqrt{\frac{\Upsilon_{\mathbb{G}', \mathcal{X}}(A)}{|\mathbb{G}'|}} \sqrt{2 \frac{\log 2\delta^{-1}}{N}}.$$

Assume in addition that $g \triangleright A \cap A = \emptyset$ for all $g \in \mathbb{G}' \setminus \{e\}$. Then it holds w.p.a.l. $1 - \delta$

$$\frac{|\widehat{\mathbb{E}}_{\mathbf{x}}[\mathbb{1}_A] - \mathbb{E}[\mathbb{1}_A(\mathbf{x})]|}{\mathbb{E}[\mathbb{1}_A(\mathbf{x})]} \leq \left(\frac{1 - \mathbb{P}[\mathbf{x} \in \mathbb{G}' \triangleright A]}{\mathbb{P}[\mathbf{x} \in \mathbb{G}' \triangleright A]} \right) \frac{\log 2\delta^{-1}}{N} + \sqrt{2 \frac{\log 2\delta^{-1}}{N}} \sqrt{\frac{1 - \mathbb{P}[\mathbf{x} \in \mathbb{G}' \triangleright A]}{\mathbb{P}[\mathbf{x} \in \mathbb{G}' \triangleright A]}}.$$

If the distribution of \mathbf{y} is \mathbb{G}' -invariant, then an identical result is immediately available for \mathbf{y} by the same proof argument.

Remark M.11. Using the standard empirical mean estimator that does not take advantage of \mathbb{G} -invariance, we obtain a concentration bound with a slower rate. For example, for an antisymmetric event A , we would achieve, w.p.a.l. $1 - \delta$, the following result:

$$\frac{|\widehat{\mathbb{E}}_{\mathbf{x}}[\mathbb{1}_A] - \mathbb{E}[\mathbb{1}_A(\mathbf{x})]|}{\mathbb{E}[\mathbb{1}_A(\mathbf{x})]} \leq \left(\frac{1 - \mathbb{P}[\mathbf{x} \in A]}{\mathbb{P}[\mathbf{x} \in A]} \right) \frac{\log 2\delta^{-1}}{N} + \sqrt{2 \frac{\log 2\delta^{-1}}{N}} \sqrt{\frac{1 - \mathbb{P}[\mathbf{x} \in A]}{\mathbb{P}[\mathbf{x} \in A]}}.$$

Specifically, leveraging \mathbb{G}' -invariance allows us to replace $\mathbb{P}[\mathbf{x} \in A]$ with $\mathbb{P}[\mathbf{x} \in \mathbb{G}' \triangleright_{\mathcal{X}} A]$, which represents the probability of the entire orbit of A under the action of \mathbb{G}' . This becomes particularly interesting when $\mathbb{P}[\mathbf{x} \in A] \ll \mathbb{P}[\mathbf{x} \in \mathbb{G}' \triangleright_{\mathcal{X}} A]$, especially in the case of rare events where $\mathbb{P}[\mathbf{x} \in A] \approx 0$.

Proof. Since $P_{\mathbf{x}}$ is \mathbb{G}' -invariant, we have $\mathbb{E}[\mathbb{1}_{g^{-1} \triangleright_{\mathcal{X}} A}(\mathbf{x})] = \mathbb{P}[\mathbf{x} \in A]$ and $\text{Var}(\mathbb{1}_{g^{-1} \triangleright_{\mathcal{X}} A}(\mathbf{x})) = \text{Var}(\mathbb{1}_A(\mathbf{x})) = \mathbb{P}[\mathbf{x} \in A](1 - \mathbb{P}[\mathbf{x} \in A])$, for any $g \in \mathbb{G}'$. Hence

$$\widehat{\mathbb{E}}_{\mathbf{x}}[\mathbb{1}_A] - \mathbb{E}[\mathbb{1}_A(\mathbf{x})] = \frac{1}{N} \sum_{i=1}^N Z_i \quad \text{with} \quad Z_i = \frac{1}{|\mathbb{G}'|} \sum_{g \in \mathbb{G}'} \mathbb{1}_{g^{-1} \triangleright_{\mathcal{X}} A}(\mathbf{x}_i) - \mathbb{E}[\mathbb{1}_{g^{-1} \triangleright_{\mathcal{X}} A}(\mathbf{x}_i)], \quad \forall i \in [N].$$

The Z_i 's are i.i.d. copies of $Z = Z_A$. In view of (116), we can apply Hoeffding's inequality [Bercu et al. \(2015, Theorem 2.16\)](#). We get for any $\delta \in (0, 1)$ w.p.a.l $1 - \delta$

$$|\widehat{\mathbb{E}}_{\mathbf{x}}[\mathbb{1}_A] - \mathbb{E}[\mathbb{1}_A(\mathbf{x})]| \leq \gamma_{\mathbb{G}', \infty}(A) \sqrt{\frac{\log 2\delta^{-1}}{2N}}. \quad (118)$$

We propose to prove another bound based on application of Bernstein's inequality. We first prove an improved bound on $\text{Var}(Z)$ as compared to the standard empirical mean estimator which does not exploit \mathbb{G} -invariance. Indeed we have

$$\begin{aligned} \text{Var}(Z) &= \frac{1}{|\mathbb{G}'|^2} \left(\sum_{g \in \mathbb{G}'} \text{Var}(\mathbb{1}_{g^{-1} \triangleright \mathbb{A}}(\mathbf{x})) + \sum_{g \neq g'} \text{Cov} \left(\mathbb{1}_{g^{-1} \triangleright \mathbb{A}}(\mathbf{x}), \mathbb{1}_{(g')^{-1} \triangleright \mathbb{A}}(\mathbf{x}) \right) \right) \\ &= \frac{\mathbb{P}(\mathbf{x} \in \mathbb{A})(1 - \mathbb{P}(\mathbf{x} \in \mathbb{A}))}{|\mathbb{G}'|} + \frac{1}{|\mathbb{G}'|^2} \sum_{g \neq g'} \text{Cov} \left(\mathbb{1}_{g^{-1} \triangleright \mathbb{A}}(\mathbf{x}), \mathbb{1}_{(g')^{-1} \triangleright \mathbb{A}}(\mathbf{x}) \right). \end{aligned}$$

Next, using again that \mathbb{P}_X is \mathbb{G} -invariant, we get for any $g, g' \in \mathbb{G}'$

$$\begin{aligned} \text{Cov} \left(\mathbb{1}_{g^{-1} \triangleright \mathbb{A}}(\mathbf{x}), \mathbb{1}_{(g')^{-1} \triangleright \mathbb{A}}(\mathbf{x}) \right) &= \\ &= \mathbb{P}[\mathbf{x} \in g^{-1} \triangleright \mathbb{A} \cap (g')^{-1} \triangleright \mathbb{A}] - \mathbb{P}[\mathbf{x} \in g^{-1} \triangleright \mathbb{A}] \mathbb{P}[\mathbf{x} \in (g')^{-1} \triangleright \mathbb{A}] \\ &= \mathbb{P}[\mathbf{x} \in g^{-1} \triangleright \mathbb{A} \cap (g')^{-1} \triangleright \mathbb{A}] - \mathbb{P}[\mathbf{x} \in \mathbb{A}]^2. \end{aligned} \quad (119)$$

Using again the invariance assumption, we note that

$$\sum_{g \neq g'} \text{Cov}(\mathbb{1}_{g^{-1} \triangleright \mathbb{A}}(\mathbf{x}), \mathbb{1}_{(g')^{-1} \triangleright \mathbb{A}}(\mathbf{x})) = |\mathbb{G}'| \left(\sum_{g \in \mathbb{G}', g \neq e} \mathbb{P}[\mathbf{x} \in \mathbb{A} \cap g \triangleright \mathbb{A}] - |\mathbb{G}'|(|\mathbb{G}'| - 1) \mathbb{P}[\mathbf{x} \in \mathbb{A}]^2 \right)$$

Consequently by definition of $\gamma_{\mathbb{G}'}(A)$ in (91) and (92), we get

$$\sum_{g \in \mathbb{G}', g \neq e} \mathbb{P}[\mathbf{x} \in \mathbb{A} \cap g \triangleright \mathbb{A}] = (|\mathbb{G}'| - 1) \gamma_{\mathbb{G}'}(A) \mathbb{P}(\mathbf{x} \in \mathbb{A}).$$

Combining the last four displays, we get

$$\text{Var}(Z) = \frac{\mathbb{P}(\mathbf{x} \in \mathbb{A})(1 - \mathbb{P}(\mathbf{x} \in \mathbb{A})) + (|\mathbb{G}'| - 1) (\gamma_{\mathbb{G}'}(A) - \mathbb{P}[\mathbf{x} \in \mathbb{A}]) \mathbb{P}[\mathbf{x} \in \mathbb{A}]}{|\mathbb{G}'|} = \frac{\Upsilon_{\mathbb{G}', \mathbf{x}}(A)}{|\mathbb{G}'|}. \quad (121)$$

We note that for any $p \geq 3$

$$\sum_{i=1}^N \mathbb{E}[(\max(0, Z_i))^p] \leq \frac{p!}{2} \max(0, \gamma_{\mathbb{G}', \infty}(\mathbb{A}) - \mathbb{P}[\mathbf{x} \in \mathbb{A}])^{p-2} N \text{Var}(Z).$$

Then [Bercu et al. \(2015, Theorem 2.1\)](#) gives w.p.a.l. $1 - \delta$

$$\widehat{\mathbb{E}}_{\mathbf{x}}[\mathbb{1}_{\mathbb{A}}] - \mathbb{E}[\mathbb{1}_{\mathbb{A}}(\mathbf{x})] \leq \max(0, \gamma_{\mathbb{G}', \infty}(\mathbb{A}) - \mathbb{P}[\mathbf{x} \in \mathbb{A}]) \frac{\log \delta^{-1}}{N} + \sqrt{\text{Var}(Z)} \sqrt{2 \frac{\log \delta^{-1}}{N}}.$$

Applying the same reasoning to variables $-Z_1, \dots, -Z_N$ and an union bound gives gives w.p.a.l. $1 - 2\delta$

$$|\widehat{\mathbb{E}}_{\mathbf{x}}[\mathbb{1}_{\mathbb{A}}] - \mathbb{E}[\mathbb{1}_{\mathbb{A}}(\mathbf{x})]| \leq |\gamma_{\mathbb{G}', \infty}(\mathbb{A}) - \mathbb{P}[\mathbf{x} \in \mathbb{A}]| \frac{\log \delta^{-1}}{N} + \sqrt{\text{Var}(Z)} \sqrt{2 \frac{\log \delta^{-1}}{N}}. \quad (122)$$

Next, we note that when $g \triangleright \mathbb{A} \cap \mathbb{A} = \emptyset$ for all $g \in \mathbb{G}' \setminus \{e\}$, then $\gamma_{\mathbb{G}'}(\mathbb{A}) = 0$ and $\mathbb{P}[\mathbf{x} \in \mathbb{A}] = \mathbb{P}[\mathbf{x} \in \mathbb{G}' \triangleright \mathbb{A}] / |\mathbb{G}'|$. Consequently we get

$$\Upsilon_{\mathbb{G}', \mathbf{x}}(\mathbb{A}) = \frac{\mathbb{P}[\mathbf{x} \in \mathbb{G}' \triangleright \mathbb{A}](1 - \mathbb{P}[\mathbf{x} \in \mathbb{G}' \triangleright \mathbb{A}])}{|\mathbb{G}'|} \quad \text{and} \quad \frac{\gamma_{\mathbb{G}', \infty}(\mathbb{A}) - \mathbb{P}[\mathbf{x} \in \mathbb{A}]}{\mathbb{P}[\mathbf{x} \in \mathbb{A}]} = \frac{1}{\mathbb{P}[\mathbf{x} \in \mathbb{G}' \triangleright \mathbb{A}]} - 1.$$

Hence under the additional assumptions, dividing by $\mathbb{E}[\mathbb{1}_{\mathbb{A}}(\mathbf{x})] = \mathbb{P}[\mathbf{x} \in \mathbb{A}]$ gives w.p.a.l. $1 - 2\delta$

$$\frac{|\widehat{\mathbb{E}}_{\mathbf{x}}[\mathbb{1}_{\mathbb{A}}] - \mathbb{E}[\mathbb{1}_{\mathbb{A}}(\mathbf{x})]|}{\mathbb{E}[\mathbb{1}_{\mathbb{A}}(\mathbf{x})]} \leq \left(\frac{1}{\mathbb{P}[\mathbf{x} \in \mathbb{G}' \triangleright \mathbb{A}]} - 1 \right) \frac{\log \delta^{-1}}{N} + \sqrt{2 \frac{\log \delta^{-1}}{N}} \sqrt{\frac{1 - \mathbb{P}[\mathbf{x} \in \mathbb{G}' \triangleright \mathbb{A}]}{\mathbb{P}[\mathbf{x} \in \mathbb{G}' \triangleright \mathbb{A}]}}.$$

Replacing δ by $\delta/2$ gives w.p.a.l. $1 - \delta$

$$\frac{|\widehat{\mathbb{E}}_{\mathbf{x}}[\mathbb{1}_{\mathbb{A}}] - \mathbb{E}[\mathbb{1}_{\mathbb{A}}(\mathbf{x})]|}{\mathbb{E}[\mathbb{1}_{\mathbb{A}}(\mathbf{x})]} \leq \left(\frac{1}{\mathbb{P}[\mathbf{x} \in \mathbb{G}' \triangleright \mathbb{A}]} - 1 \right) \frac{\log 2\delta^{-1}}{N} + \sqrt{2 \frac{\log 2\delta^{-1}}{N}} \sqrt{\frac{1 - \mathbb{P}[\mathbf{x} \in \mathbb{G}' \triangleright \mathbb{A}]}{\mathbb{P}[\mathbf{x} \in \mathbb{G}' \triangleright \mathbb{A}]}}. \quad (123)$$

□

Proposition M.12. Let $A_i, i \in [N]$ be i.i.d copies of a random variable A in a separable Hilbert space with norm $\|\cdot\|$. If there exist constants $L > 0$ and $\sigma > 0$ such that for every $m \geq 2$, $\mathbb{E} \|A\|^m \leq \frac{1}{2} m! L^{m-2} \sigma^2$, then with probability at least $1 - \delta$

$$\left\| \frac{1}{N} \sum_{i \in [N]} A_i - \mathbb{E} A \right\| \leq \frac{4\sqrt{2}}{\sqrt{N}} \sqrt{\sigma^2 + \frac{L^2}{N} \log \frac{2}{\delta}}. \quad (124)$$

Lemma M.13 ((Sub-Gaussian random variable) Lemma 5.5. in Vershynin (2011)). Let Z be a random variable. Then, the following assertions are equivalent with parameters $K_i > 0$ differing from each other by at most an absolute constant factor.

1. Tails: $\mathbb{P}\{|Z| > t\} \leq \exp(1 - t^2/K_1^2)$ for all $t \geq 0$;
2. Moments: $(\mathbb{E}|Z|^p)^{1/p} \leq K_2 \sqrt{p}$ for all $p \geq 1$;
3. Super-exponential moment: $\mathbb{E} \exp(Z^2/K_3^2) \leq 2$.

A random variable Z satisfying any of the above assertions is called a sub-Gaussian random variable. We will denote by K_3 the sub-Gaussian norm.

Consequently, a sub-Gaussian random variable satisfies the following equivalence of moments property. There exists an absolute constant $c > 0$ such that for any $m \geq 2$,

$$(\mathbb{E}|Z|^m)^{1/m} \leq cK_3 \sqrt{m} (\mathbb{E}|Z|^2)^{1/2}.$$

Lemma M.14. Assume that Y is sub-Gaussian with sub-Gaussian norm K . We set $\sigma_\theta^2(Y) := \text{Var}(\|Y - \mathbb{E}[y]\|)$. Then there exists an absolute constant $C > 0$ such that for any $\delta \in (0, 1)$, it holds w.p.a.l. $1 - \delta$

$$\left\| \widehat{\mathbb{E}}_y[\mathbf{y}] - \mathbb{E}[\mathbf{y}] \right\| \leq \frac{C}{\sqrt{N}} \sqrt{\sigma^2(\mathbf{y}) + \frac{K^2}{N} \log(2\delta^{-1})}.$$

Proof. Set $Z := \|\mathbf{y} - \mathbb{E}\mathbf{y}\|$ and we recall that $\sigma^2(\mathbf{y}) := \text{Var}(\|\mathbf{y} - \mathbb{E}[\mathbf{y}]\|)$. We check that the moment condition,

$$\mathbb{E} Z^m \leq \frac{1}{2} m! L^{m-2} \sigma^2(\mathbf{y})^2, \quad \forall m \geq 2,$$

for some constant $L > 0$ to be specified.

The condition is obviously satisfied for $m = 2$. Next for any $m \geq 3$, the Cauchy-Schwarz inequality and the equivalence of moment property give

$$\mathbb{E} Z^m \leq \left(\mathbb{E} Z^{2(m-2)} \right)^{1/2} (\mathbb{E} Z^4)^{1/2} \leq 4K_3^2 \sigma_\theta^2(Y)^2 \left(\mathbb{E} Z^{2(m-2)} \right)^{1/2}.$$

Next, by homogeneity, rescaling Z to Z/K_1 we can assume that $K_1 = 1$ in Lemma M.13. We recall that if Z is in addition non-negative random variable, then for every integer $p \geq 1$, we have

$$\mathbb{E} Z^p = \int_0^\infty \mathbb{P}\{Z \geq t\} p t^{p-1} dt \leq \int_0^\infty e^{1-t^2} p t^{p-1} dt = \left(\frac{ep}{2}\right) \Gamma\left(\frac{p}{2}\right).$$

With $p = 2(m-2)$, we get that $\mathbb{E} Z^p \leq e(m-2) \Gamma(m-2) = e(m-2)! = em!/2$. Using again Lemma M.13, we can take $L = cK$ for some large enough absolute constant $c > 0$. Then Proposition M.12 gives the result. \square