UNIOMA: UNIFIED OPTIMAL-TRANSPORT MULTI-MODAL STRUCTURAL ALIGNMENT FOR ROBOT PERCEPTION

Anonymous authorsPaper under double-blind review

ABSTRACT

Achieving generalizable and well-aligned multimodal representation remains a core challenge in artificial intelligence. While recent approaches have attempted to align modalities by modeling conditional or higher-order statistical dependencies, they often fail to capture the structural coherence across modalities. In this work, we propose a novel multimodal alignment method that augments existing contrastive losses with a geometry-aware Gromov-Wasserstein (GW) distance-based regularization. To this end, we encode intra-modality geometry with modalityspecific similarity matrices and extend the GW distance to minimize their discrepancies from a dynamically learned barycenter, thereby enforcing structural alignment across modalities beyond what is captured by InfoNCE-like mutual information objectives. We apply this optimal-transport-based alignment strategy to robot perception tasks involving underexplored modalities such as force and tactile signals, where modality data often exhibit varying sample densities. Experimental results show that our method yields superior inter-modal coherence and significantly improves downstream robot perception tasks such as robot and environment state prediction. Moreover, our GW-based augmentation term is versatile and can be seamlessly integrated into most InfoNCE-like objectives.

1 Introduction

The integration of information from diverse sources or modalities has received increasing attention across a wide range of AI applications, including image/video/text generation (Rombach et al., 2022; Mirza & Osindero, 2014), healthcare (Acosta et al., 2022), autonomous systems (Feng et al., 2021), and scientific discovery (Steyaert et al., 2023). Recent advances in contrastive self-supervised learning (CSSL) (He et al., 2020; Chen et al., 2020; Grill et al., 2020; Chen & He, 2021), particularly those leveraging InfoNCE losses (Oord et al., 2018), have shown strong performance in aligning heterogeneous modalities into a shared representation space (Radford et al., 2021). Such alignment has enabled zero-shot cross-modal retrieval, transfer, generation, and completion (Radford et al., 2021; Girdhar et al., 2023; Chen et al., 2023; Zhu et al., 2023; Luo et al., 2022). By maximizing agreement between paired modalities of the same instance while minimizing similarity between distinct instances, CSSL encourages the learning of invariant and semantically meaningful features.

While effective, InfoNCE-style objectives operate as binary classification losses that only discriminate positives from negatives (Wang & Isola, 2020), without explicitly modeling the continuous pairwise distance geometry within each modality. In multimodal alignment, this limitation produces what we call a *structural alignment gap* (Liang et al., 2022): embeddings may appear statistically aligned across modalities yet fail to preserve their intrinsic structural topologies. Our key insight is that multimodal alignment should not be limited to maximizing *population-level* statistical dependence between distributions of modality representations. It must also preserve *instance-level* geometric relationships within each modality. In other words, if x_i is close to x_j , then their counterparts y_i and y_j should also remain close. Classic InfoNCE objectives, which are essentially a lower bound of Shannon's mutual information (Kraskov et al., 2004; Poole et al., 2019), rely on binary discrimination between positive and negative pairs. While effective at capturing population-level dependence, this approach is theoretically incapable of preserving intra-modal geometry, often leading to representations that are statistically aligned but structurally inconsistent.

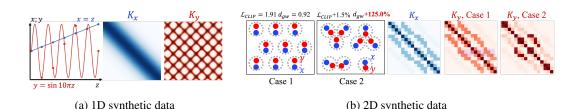


Figure 1: Structural alignment gap. (a) 1D synthetic data. Although x and y have high mutual information and thus a low InfoNCE loss \mathcal{L}_c , their intra-modal similarity matrices differ: (b) 2D synthetic data with instances (dashed gray circles). Blue and red denote two modalities. Pointwise correspondences are close in both cases (thus InfoNCE-like loss $\mathcal{L}_{\text{CLIP}}$ changes only +1.5%), but the GW distance jumps by +125%. Also, K_y in Case 2 shows block structure absent from K_x .

We illustrate the gap with two synthetic examples in Fig. 1. (a) Let the latent variable be $z \sim \text{Uniform}[0,1]$, from which we generate two modalities x=z and $y=\sin(10\pi z)$. Although x and y are highly dependent, their intra-modal geometries differ markedly. In x, distances are simply $|x_i-x_j|$, whereas the high-frequency oscillation in y disrupts local neighborhoods, so nearby x can map to distant y, leading to dissimilar intra-modal similarity matrices. (b) modality x forms a regular grid, while modality y is either a globally shifted/noisy copy (case 1, left figure) or an unevenly shifted version that clusters points into triplets (case 2, right figure). Both cases preserve pointwise correspondences, leading to a lower InfoNCE loss. However, case 2 distorts the global structure, which is reflected in block patterns in K_y (kernel similarity matrix) that are absent in K_x .

This *structural alignment gap* is particularly critical in robotics, where multimodal sensor streams are neither i.i.d. nor structureless: trajectories form subclusters (Sermanet et al., 2017), contact events induce discontinuities (Stewart & Trinkle; Guo et al., 2023), and proprioceptive signals follow physical constraints (Lee et al., 2020; Welch & Bishop, 1995). Failing to account for these structures limits the effectiveness of learned representations for downstream robotic tasks.

To address the identified *structural alignment gap*, we introduce **UniOMA**—a **Uni**fied **O**ptimal-transport **M**ulti-modal structural **A**lignment framework that scales naturally to three or more modalities. UniOMA augments contrastive learning with a structure-aware regularization based on Gromov–Wasserstein (GW) distances and barycenters (Peyré et al., 2016; Gong et al., 2022). In our formulation, observations from each modality are represented as a metric space through intra-modal similarity matrices. A dynamic GW barycenter is then computed as the structural consensus across modalities, and each modality is softly aligned to this consensus by minimizing weighted GW distances. The modality weights are optimized end-to-end alongside encoder parameters, enabling adaptive contributions of different modalities to the structural consensus. This barycentric formulation avoids pairwise couplings across modalities, reducing the complexity from $O(M^2)$ to O(M), where M is the number of modalities, and thus scales naturally to three or more modalities.

In summary, our main contributions are:

- C1 We propose UniOMA, a structure-aware multimodal alignment framework based on Gromov–Wasserstein distance and barycenters, which naturally scales to 3+ modalities.
- C2 We identify and formalize the structural alignment gap, demonstrating why InfoNCE-style objectives fail to preserve intra-modal geometry, supported by synthetic analysis.

We evaluate UniOMA on diverse robotic benchmarks across vision, audio, tactile, force, and proprioception modalities, including robot state prediction, environment state prediction, and cross-modal retrieval. Comprehensive experiments show that UniOMA improves downstream performance and preserves intra-modal structural consistency across diverse modalities.

2 BACKGROUND AND RELATED WORK

In this section, we first introduce the background of contrastive learning-based multimodal alignment and review its extensions to settings with three or more modalities, highlighting their inherent

connections and limitations. We then briefly review existing approaches to multimodal representation learning in robotics, with a focus on multimodal fusion.

2.1 ALIGNMENT VIA INFONCE AND EXTENSIONS TO MORE THAN TWO MODALITIES

Unlike multimodal fusion (Lu et al., 2019; Li et al., 2019), which typically requires all modalities to be present at inference, alignment into a shared embedding space remains functional even if some modalities are missing, enabling zero-shot retrieval, generation, and modality completion (Jia et al., 2021). A representative example is CLIP (Radford et al., 2021), which trains modality-specific encoders $f_{\theta}^{(1)}$, $f_{\theta}^{(2)}$ using an InfoNCE-style objective to identify the correct cross-modal pair among N candidates:

$$\ell_{\text{CLIP}}^{(1\rightarrow 2)}(\theta) = -\frac{1}{N} \sum_{i=1}^{N} \log \frac{\exp\left(\text{sim}(\mathbf{z}_{i}^{(1)}, \mathbf{z}_{i}^{(2)})/\tau\right)}{\sum_{j=1}^{N} \exp\left(\text{sim}(\mathbf{z}_{i}^{(1)}, \mathbf{z}_{j}^{(2)})/\tau\right)},\tag{1}$$

where $sim(\cdot, \cdot)$ is the similarity between the embeddings $\mathbf{z}^{(m)} = f_{\theta}^{(m)}(\mathbf{x}^{(m)}), m = 1, 2$ and τ denotes a temperature parameter. The final CLIP objective symmetrizes Eq. (1) by taking the average:

$$\mathcal{L}_{\text{CLIP}}^{(1,2)}(\theta) = \frac{1}{2} (\ell_{\text{CLIP}}^{(1 \to 2)}(\theta) + \ell_{\text{CLIP}}^{(2 \to 1)}(\theta)), \tag{2}$$

where $\mathcal{L}_{\text{CLIP}}^{(2\to 1)}$ is the reverse direction $2\to 1$. In general, this InfoNCE-based objective captures the statistical correlation, providing lower-bound of the mutual information (MI; Kraskov et al. (2004); Poole et al. (2019)) between the anchor modality 1 $\mathcal{X}^{(1)}$ and modality 2 $\mathcal{X}^{(2)}$

$$I(\mathcal{X}^{(1)}; \mathcal{X}^{(2)}) \ge \log N - 2\mathcal{L}_{\text{CLIP}}^{(1,2)}(\theta). \tag{3}$$

Despite their success, InfoNCE-like objectives reduce continuous similarity structure among samples into a binary signal (positive vs. negative), leading to the learned embedding space containing modality-wise co-located yet structurally isolated instances, neglecting intra-modal geometry.

Real-world systems, particularly in robotics, often involve three or more modalities. Aligning these multimodal sources within a shared embedding space enables richer cross-modal interactions. Existing approaches typically extend CLIP to three modalities by summing all pairwise contrastive losses (Tian et al., 2020; Girdhar et al., 2023; Akbari et al., 2021; Chen et al., 2023; Alayrac et al., 2020; Chen et al., 2021; Liu et al., 2024; Huang et al., 2023; Mai et al., 2022; Moon et al., 2022; Shvetsova et al., 2022; Xue et al., 2022; Guzhov et al., 2022):

$$\mathcal{L}_{\text{CMC}}^{(1,2,3)}(\theta) = \mathcal{L}_{\text{CLIP}}^{(1,2)}(\theta) + \mathcal{L}_{\text{CLIP}}^{(1,3)}(\theta) + \mathcal{L}_{\text{CLIP}}^{(2,3)}(\theta). \tag{4}$$

Such pairwise extensions neglect higher-order dependencies among modalities. To address this issue, Symile (Saporta et al., 2024) formulates triple-wise contrastive objectives as:

$$\mathcal{L}_{\text{Symile}}^{(1,2,3)}(\theta) = \frac{1}{3} [\ell^{(1\to2,3)}(\theta) + \ell^{(2\to1,3)}(\theta) + \ell^{(3\to1,2)}(\theta)]. \tag{5}$$

Here, $\ell^{(1\to 2,3)}$ is the InfoNCE-like loss for one positive triple and N-1 negative triples given by

$$\ell^{(1\to 2,3)}(\theta) = -\frac{1}{N} \sum_{i=1}^{N} \log \frac{\exp(\langle \mathbf{z}_{i}^{(1)}, \mathbf{z}_{i}^{(2)}, \mathbf{z}_{i}^{(3)} \rangle / \tau)}{\sum_{j=1}^{N} \exp(\langle \mathbf{z}_{i}^{(1)}, \mathbf{z}_{j}^{(2)}, \mathbf{z}_{j}^{(3)} \rangle / \tau)}, \tag{6}$$

where each term $\ell^{(1\to 2,3)}$ compares one positive triple against N-1 negatives, $\langle \cdot, \cdot, \cdot \rangle$ is the coordinate-wise sum of the element-wise product. More recently, GRAM (Cicchetti et al., 2024) replaces the dot product similarity with the Gramian volume spanned by embeddings from multiple modalities, providing a higher-order, groupwise compatibility score (rather than pairwise similarity)

$$\mathcal{L}_{\text{GRAM}}^{(1,\dots,M)}(\theta) = \frac{1}{2} (\ell_{\text{D2A}}^{(1\to 2,\dots,M)}(\theta) + \ell_{\text{A2D}}^{(1\to 2,\dots,M)}(\theta)) + \lambda \ell_{\text{DAM}}(\theta), \tag{7}$$

$$\ell_{\text{D2A}}^{(1 \to 2, \dots, M)}(\theta) = -\frac{1}{N} \sum_{i=1}^{N} \log \frac{\exp(-\text{Vol}(\mathbf{z}_{i}^{(1)}, \mathbf{z}_{i}^{(2)}, \dots, \mathbf{z}_{i}^{(M)})/\tau)}{\sum_{j=1}^{N} \exp(-\text{Vol}(\mathbf{z}_{j}^{(1)}, \mathbf{z}_{i}^{(2)}, \dots, \mathbf{z}_{i}^{(M)})/\tau)},$$
(8)

$$\ell_{\text{A2D}}^{(1\to 2,\dots,M)}(\theta) = -\frac{1}{N} \sum_{i=1}^{N} \log \frac{\exp(-\text{Vol}(\mathbf{z}_{i}^{(1)}, \mathbf{z}_{i}^{(2)}, \dots, \mathbf{z}_{i}^{(M)})/\tau)}{\sum_{j=1}^{N} \exp(-\text{Vol}(\mathbf{z}_{i}^{(1)}, \mathbf{z}_{j}^{(2)}, \dots, \mathbf{z}_{j}^{(M)})/\tau)}.$$
 (9)

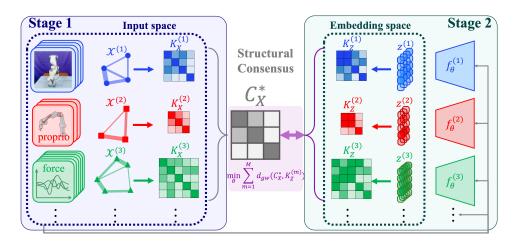


Figure 2: **UniOMA in two stages.** Stage 1 (left): for each modality $\mathcal{X}^{(m)}$ we form an input-space similarity matrix $\mathbf{K}_{\mathbf{x}}^{(m)}$ and estimate a GW barycenter $\mathbf{C}_{\mathbf{x}}^*$ as the structural consensus. Stage 2 (right): encoders produce embeddings $\mathbf{z}^{(m)}$ inducing $\mathbf{K}_{\mathbf{z}}^{(m)}$, which are aligned to the consensus by minimizing $\sum_{m} \lambda_m \, d_{gw}(\mathbf{C}_{\mathbf{x}}^*, \mathbf{K}_{\mathbf{z}}^{(m)})$ (with a standard contrastive loss; omitted). Aligning each modality to a single consensus avoids pairwise $O(M^2)$ couplings and scales to $M \geq 3$.

where \mathcal{L}_{D2A} , \mathcal{L}_{A2D} are the GRAM contrastive loss (data-to-anchor for D2A and anchor-to-data for A2D) with modality 1 as the anchor. \mathcal{L}_{DAM} is the data-caption matching loss to match the modality labels (Cicchetti et al., 2024). Vol (\cdot, \dots, \cdot) is the volume of the M-dimensional parallelotope formed by the embedding vectors $\mathbf{z}^{(m)}$.

These methods mark progress toward multi-modal (M>2) alignment but still remain limited to instance-level dependencies, overlooking intra-modal structure. Zhu & Luo (2024) address this by adding an optimal transport (OT; Villani et al. (2008)) regularizer to enforce cross-modal consistency. Yet, their approach still treats modalities as holistic distributions, ignoring relational structures within each modality, and applies OT directly on embeddings rather than raw data geometry, limiting interpretability and flexibility.

2.2 MULTIMODAL REPRESENTATION LEARNING IN ROBOTICS

Robotics is inherently multimodal: vision, force—torque, tactile sensing, and proprioception provide complementary views of the robot—environment system. While multimodal representation learning has been extensively studied in vision—language settings, its exploration in robotics remains limited. Existing work, including the recent Vision—Language—Action (VLA) model, has primarily focused on modality fusion or transfer (Lee et al., 2019a;b; Shridhar et al., 2020; Brohan et al., 2022; Driess et al., 2023; Kim et al., 2024; Octo Model Team et al., 2024; Intelligence et al., 2025).

By comparison, alignment of robotic perception modalities into a shared space remains underexplored. Recent efforts (Wojcik et al., 2024; Dutta et al., 2024) demonstrate cross-modal retrieval and perception, while Zambelli et al. (2021); Sermanet et al. (2017) demonstrate how cross-modal or cross-temporal alignment can yield transferable representations. These developments underscore that robot perception data is highly structured (trajectories, contact events, physical constraints), motivating alignment methods that preserve intra-modal geometry across modalities rather than relying solely on fusion.

3 Method

Our proposed UniOMA aligns three or more heterogeneous modalities by preserving both statistical correspondence and structural coherence across modalities. Leveraging the optimal transport geometry, UniOMA augments contrastive-based binary instance-wise correlations (positive or neg-

ative pair) with structural properties by minimizing the Gromov-Wasserstein (GW) distance across the modalities. In this section, we first explain the multimodal alignment problem, and then we define the intra-modal structure information and the cross-modal structure consensus, followed by the definition of the UniOMA objective and the alignment algorithm.

3.1 PROBLEM STATEMENT

Let $\mathcal{X}^{(1)}, \mathcal{X}^{(2)}, \dots, \mathcal{X}^{(M)}$ denote M modalities. The goal of multimodal alignment is to learn modality-specific encoders $f^{(m)}: \mathcal{X}^{(m)} \to \mathbb{R}^d$, $m=1,\dots,M$ that project inputs $\mathbf{x}^{(m)} \in \mathcal{X}^{(m)}$ into a shared latent space $\mathbf{z}^{(m)} = f^{(m)}(\mathbf{x}^{(m)}) \in \mathbb{R}^d$. The key objective is that embeddings of the same underlying instance across modalities map to nearby latent vectors, i.e., $\mathbf{z}_i^{(1)} \approx \mathbf{z}_i^{(2)} \approx \dots \approx \mathbf{z}_i^{(M)}$. In robotics, the modalities $\mathcal{X}^{(m)}$ may include vision (third-person or wrist-mounted), audio commands, force-torque signals, proprioception (joints, inertial measurement unit (IMU), end-effector pose), tactile sensing, and environment states (e.g., object pose). Aligning them into a shared latent space enables cross-modal reasoning and zero-shot transfer: for example, vision of an end-effector trajectory should yield embeddings consistent with the same trajectory from proprioception or touch. Such unified representations enable downstream tasks such as robot/environment state prediction, action prediction or generation, and modality completion when certain sensor streams are missing.

3.2 Gromov-Wasserstein Distance

The Gromov–Wasserstein (GW) distance (Peyré et al., 2016; Gong et al., 2022) is a natural extension of Optimal Transport (OT) (Villani et al., 2008) to settings where distributions lie in different metric spaces. The classic OT problem seeks the minimum cost of transporting one probability measure into another within the same metric space. Given two measures μ and ν and a cost function $c: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$, the Kantorovich formulation of the OT problem is

$$d_w(\mu, \nu) := \inf_{\pi \in \Pi(\mu, \nu)} \int_{\mathcal{X} \times \mathcal{X}} c(\mathbf{x}, \mathbf{y}) d\pi(\mathbf{x}, \mathbf{y}), \tag{10}$$

where π is a transport plan with marginals μ and ν . When both measures are supported on the same space \mathcal{X} , $c(\cdot, \cdot)$ is a distance metric (e.g., ℓ_2), and Eq. 10 defines the Wasserstein distance.

However, in multimodal learning the two distributions often live in different spaces (e.g., images vs. tactile signals). In such cases, defining a cross-modal cost $c(\mathbf{x}, \mathbf{y})$ is generally impossible. The GW distance addresses this by replacing the direct cross-modal cost with a relational cost that compares intra-modal similarities.

Definition 1 (Gromov-Wasserstein Distance). Let $\mathcal{X}_{d_{\mathbf{x}},\mu}$ and $\mathcal{Y}_{d_{\mathbf{y}},\nu}$ be two metric–measure spaces (mm-spaces), with distance metrics $d_{\mathbf{x}}$, $d_{\mathbf{y}}$ and probability measures μ,ν . The GW distance between them is defined as:

$$d_{gw}(\mu, \nu) := \inf_{\pi \in \Pi(\mu, \nu)} \int_{\mathcal{X}^2 \times \mathcal{Y}^2} c\left(d_{\mathbf{x}}(\mathbf{x}, \mathbf{x}'), d_{\mathbf{y}}(\mathbf{y}, \mathbf{y}')\right) d\pi(\mathbf{x}, \mathbf{y}) d\pi(\mathbf{x}', \mathbf{y}'),$$

where $c(d_{\mathbf{x}}(\mathbf{x}, \mathbf{x}'), d_{\mathbf{y}}(\mathbf{y}, \mathbf{y}'))$ is relational distance measuring the discrepancy between the sample pairs $(\mathbf{x}, \mathbf{x}')$ and $(\mathbf{y}, \mathbf{y}')$.

Intuitively, minimizing GW distance aligns two distributions by matching their relational geometry (pairwise structures), rather than raw coordinates. This is crucial in robotics, where modalities such as vision and force-torque are in incomparable metric spaces but have meaningful internal geometries. For discrete samples, consider the two mm-spaces $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^I$ and $\mathcal{Y} = \{\mathbf{y}_j\}_{j=1}^J$ with uniform sample distributions $\hat{\mathbf{p}}_{\mathbf{x}} = \frac{1}{I}\mathbf{1}_I$ and $\hat{\mathbf{p}}_{\mathbf{y}} = \frac{1}{J}\mathbf{1}_J$, we calculate the empirical GW distance (Gong et al., 2022) in the following definition.

Theorem 1 (Empirical GW Distance). Let the kernel matrices $\mathbf{K_x} \in \mathbb{R}^{I \times I}$ and $\mathbf{K_y} \in \mathbb{R}^{J \times J}$ be the similarity matrices conducted by the samples \mathbf{x}, \mathbf{y} from two mm-spaces \mathcal{X}, \mathcal{Y} , the empirical GW distance between the samples is:

$$\hat{d}_{gw}(\mathbf{K}_{\mathbf{x}}, \mathbf{K}_{\mathbf{y}}) := \max_{\mathbf{T} \in \Pi(\hat{\mathbf{p}}_{\mathbf{x}}, \hat{\mathbf{p}}_{\mathbf{y}})} \operatorname{tr}(\mathbf{K}_{\mathbf{x}}^{\top} \mathbf{T}^{\top} \mathbf{K}_{\mathbf{y}} \mathbf{T}),$$

where ${f T}$ is the doubly-stochastic matrix to model the transport between the two sets of samples.

See Appx. A.1 for the proof. In practice, we estimate \mathbf{T}^* via iterative OT solvers (Alg. 2), and compute $\hat{d}_{gw}(\mathbf{K_x}, \mathbf{K_y}) = \operatorname{tr}(\mathbf{K_x}^{\top} \mathbf{T^*}^{\top} \mathbf{K_y} \mathbf{T^*})$. This formulation enables cross-modal alignment directly from intra-modal similarity structures, without the need of an explicit cross-modal cost function or extra neural potential models (Korotin et al., 2022b;a; Arjovsky et al., 2017).

3.3 STRUCTURAL CONSENSUS

To preserve intra-modal structure during alignment, we treat each modality $\mathcal{X}^{(m)}$ as a metric space and represent its geometry via a kernel matrix $\mathbf{K}_{\mathbf{x}}^{(m)} \in \mathbb{R}^{N_m \times N_m}$, where $(K_{\mathbf{x}}^{(m)})_{ij} = \sin(\mathbf{x}_i^{(m)}, \mathbf{x}_j^{(m)})$ encodes the pairwise similarity between samples. Such kernel matrices provide a unified representation of relational structure across heterogeneous modalities, independent of raw dimensionality. The construction of $\mathbf{K}_{\mathbf{x}}^{(m)}$ depends on the modality: for visual signals (e.g., RGB or depth), we embed inputs with a pretrained encoder and compute similarities using an RBF kernel; for sequential or time-series modalities common in robotics (e.g., force-torque), we adopt a time-series clustering kernel (TCK; Mikalsen et al. (2018)) to better capture temporal structure. Additional details are provided in Appx. A.4.

The central idea is to identify a structural consensus: a latent geometry that captures the common similarity patterns across all modalities. Formally, we define it as a Gromov–Wasserstein (GW) barycenter (Gong et al., 2022) of the intra-modal structures.

Definition 2 (Structural Consensus of Multimodal Data). *Given intra-modal kernel matrices* $\{\mathbf{K}_{\mathbf{x}}^{(m)}\}_{m=1}^{M}$, the structural consensus is defined as the barycenter:

$$\mathbf{C}_{\mathbf{x}}^* = \arg\min_{\mathbf{C}_{\mathbf{x}} \in \mathcal{M}} \sum_{m=1}^{M} \lambda_m \cdot d_{gw}(\mathbf{C}_{\mathbf{x}}, \mathbf{K}_{\mathbf{x}}^{(m)}), \tag{11}$$

where \mathcal{M} denotes the space of symmetric positive definite (SPD) matrices, d_{gw} is the GW distance (Def. 1), and λ_m are learnable modality weights.

Practically, $\mathbf{C}_{\mathbf{x}}^*$ is estimated via an iterative optimization scheme (Alg. 3 in Appx. A.3). During training, we align each modality by minimizing the GW discrepancy between its embedding-induced kernel $\mathbf{K}_{\mathbf{z}}^{(m)}$ and the consensus $\mathbf{C}_{\mathbf{x}}^*$, as described in the next section.

3.4 UNIOMA OBJECTIVE AND ALGORITHM

Given the batch-wise structural consensus $\mathbf{C}^*_{\mathbf{x}}$ in Sec. 3.3, UniOMA augments a standard contrastive term with a structure-aware regularizer

$$\mathcal{L}_{\text{UniOMA}}(\theta) = \mathcal{L}_{c}(\theta) + \alpha \sum_{m=1}^{M} \lambda_{m} \cdot d_{gw}(\mathbf{C}_{\mathbf{x}}^{*}, \mathbf{K}_{\mathbf{z}}^{(m)}), \tag{12}$$

where $\mathbf{K}_{\mathbf{z}}^{(m)}$ is the embedding-space similarity matrix of $\mathbf{z}^{(m)} = f_{\theta}^{(m)}(\mathbf{x}^{(m)})$. The scalar α balances contrastive discrimination and structural coherence, and the learnable weights $\{\lambda_m\}$ quantify each modality's contribution to the consensus. Implementation details for estimating \mathbf{C}_x^* and evaluating $d_{qw}(\cdot,\cdot)$ are in Appx. A.3–A.1 (see also Fig. 2).

Why this design? (1) Scalable to $M \geq 3$. Aligning every modality to one consensus avoids $O(M^2)$ pairwise couplings. (2) Flexible to heterogeneous and asynchronous modalities. GW distance compares intra-modal similarity matrices, not raw coordinates, thus is naturally robust to modalities with different dimensionalities. Also, GW barycenter naturally handles unequal sample counts across modalities, which is particularly advantageous in robot perception. We empirically validate (3) in Sec. 4.5.

```
Algorithm 1 UniOMA Training(\{\mathcal{X}^{(m)}\}_{m=1}^{M}, \gamma, \alpha)
```

Input: Multimodal dataset $\{\mathcal{X}^{(m)}\}_{m=1}^{M}$, learning rate γ , structural weight α , entropy weight α' Initialize encoders $\{f^{(m)}(\cdot)\}_{m=1}^{M}$, modality weights $\{\lambda_m\}_{m=1}^{M}$, while not converged **do**

```
 \begin{cases} \text{// Stage 1: structural consensus estimation} \\ \text{Sample a batch } \{\mathbf{x}_i^{(m)}\}_{i=1}^{N_m} \text{ for each modality } \{\mathcal{X}^{(m)}\}_{m=1}^{M} \\ \text{for } m \leftarrow 1 \text{ to } M \text{ do} \\ & \bot \text{ Compute the structural information } \mathbf{K}_{\mathbf{x}}^{(m)} \in \mathbb{R}^{N_m \times N_m} \text{ for the batch } \{\mathbf{x}_i^{(m)}\}_{i=1}^{N_m} \\ \text{Estimate the structural consensus } \mathbf{C}_{\mathbf{x}}^* \text{ via Alg. 3} \end{cases}
```

return $\{f_{\theta}^{(m)}\}_{m=1}^{M}, \{\lambda_{m}\}_{m=1}^{M}$

The training procedure is summarized in Alg. 1. Each iteration proceeds in two stages:

Stage 1 (Consensus Estimation): Compute kernel matrices $K_x^{(m)}$ from a mini-batch using modality-specific similarity measures (e.g., RBF kernel for images, TCK for time series), then estimate the batch-wise consensus C_x^* via an iterative GW barycenter algorithm (Appx. A.3).

Stage 2 (Alignment Update): Encode the same mini-batch into $\mathbf{z}^{(m)}$, form kernel matrices $\mathbf{K}_{\mathbf{z}}^{(m)}$, and compute their GW distances to the consensus. The UniOMA loss is then minimized by stochastic gradient descent, jointly updating encoder parameters θ and modality weights λ_m .

4 EXPERIMENTS

We evaluate UniOMA on four multimodal robot perception settings: (i) VFD (Vision–Force–Depth) from the Vision&Touch dataset (Lee et al., 2019b; Liang et al., 2021); (ii) VFP (Vision–Force–Proprioception) from the same source; (iii) MuJoCo Push (Lee et al., 2020; Todorov et al., 2012) (Vision–Force–End-effector pose); and (iv) VAT (Vision–Audio–Tactile) derived from ObjectFolder 2.0 (Gao et al., 2022; Wojcik et al., 2024). Downstream tasks includes regression, classification, and cross-modal retrieval. To avoid architecture confounds, we use the same backbones and training schedule across methods (fusion baselines necessarily differ in fusion heads); details are in Appx. A.4.

4.1 TASKS AND DATASETS

VFD (Vision–Force–Depth). We evaluate two tasks: (1). Next-step end-effector orientation prediction (regression, 4D): Inputs are third-person RGB ($[b\times3\times128\times128]$), force–torque histories ($[b\times32\times6]$), and depth ($[b\times1\times128\times128]$). (2). Modality-consistency discrimination (classification, real vs. fake): given two triplets, identify the coherent instance. inputs are third-person RGB ($[b\times3\times128\times128]$), force–torque histories ($[b\times32\times6]$), and depth ($[b\times1\times128\times128]$). The positives are synchronized triplets from the same timestep/trajectory; negatives are constructed via shuffles of the index. We report Top-1 accuracy in Table 1.

Table 1: Comparative results on downstream tasks (regression, classification, and cross-modal retrieval). Performance is measured by MSE ($\times 10^{-3} \downarrow$), Top-1 Acc. (% \uparrow), and MAP (\uparrow). Arrows denote retrieval direction. Gray rows are baselines augmented with our GW regularizer.

	Regression \downarrow (×10 ⁻³)		Classification \uparrow (%)		VAT MAP Score ↑		
Method	V&F&D	MuJoCo	V&F&D	V&F&P	Vis→Aud	Vis→Tact	Tact→Aud
Pairwise	1.27±0.14	0.44±0.07	89.59±0.05	94.51±0.02	0.25±0.07	0.41±0.11	0.10±0.01
Pairwise+GW	1.22±0.12	0.38 ±0.09	92.44 ±0.02	94.68 ±0.03	0.36 ±0.05	0.60±0.03	0.12 ±0.02
Symile	2.81±0.10	0.28±0.04	90.02±0.04	93.94±0.06	0.10±0.02	0.21 ±0.05	0.08±0.01
Symile+GW	2.15±0.08	0.23 ±0.02	92.81 ±0.02	93.87±0.03	0.13 ±0.03	0.15±0.03	0.14 ±0.03
GRAM	3.37±0.09	0.52±0.07	92.47±0.04	93.65±0.05	0.13±0.02	0.34±0.05	0.15±0.01
GRAM+GW	2.31 ±0.05	0.30 ±0.06	93.30±0.01	93.91 ±0.04	0.79±0.10	0.58 ±0.04	0.16 ±0.01
CoMM	1.51±0.05	0.26 ± 0.04	92.39±0.01	94.13±0.03	_	_	_

VFP (Vision–Force–Proprioception). We evaluate next-step contact prediction (classification, binary). Inputs are RGB, force–torque histories, and end-effector pose ($[b \times 7]$). We classify the end-effector is in contact to the object or not.

MuJoCo Push. A planar pushing task with a Franka Emika Panda arm interacting with a puck. Inputs are low-resolution gray-scale image ($[b \times 1 \times 32 \times 32]$), current force—torque ($[b \times 6]$), and end-effector pose ($[b \times 7]$). The task is to predict the next-step object's 2-D position on the table.

VAT (Vision–Audio–Tactile). We evaluate **cross-modal retrieval** using mean average precision (MAP). Queries and retrievals are built across {Vis, Aud, Tact}; we report direction-specific MAP (e.g., Vis→Tact). The dataset provides per-object visual, sound, and tactile observations.

4.2 IMPLEMENTATION DETAILS

Encoders, optimizer, temperature, and schedules are shared across methods (fusion heads differ in CoMM). We compute input-space kernels $\{\mathbf{K}_{\mathbf{x}}^{(m)}\}$ (RBF for images with tuned γ ; TCK for time-series/force; RBF for other signals) and estimate the batch-wise consensus $\mathbf{C}_{\mathbf{x}}^*$ using iterative barycenter updates (Appx. A.3). We then align embedding-space kernels $\{\mathbf{K}_{\mathbf{z}}^{(m)}\}$ to $\mathbf{C}_{\mathbf{x}}^*$ via the UniOMA loss. Hyperparameters, TCK settings, and convergence diagnostics are detailed in Appx. A.4–A.4.

4.3 Comparisons on Downstream Tasks

We compare against: (i) **Pairwise** (CMC) (Tian et al., 2020) using summed pairwise InfoNCE; (ii) **Symile** (Saporta et al., 2024) using triple-wise InfoNCE variants; (iii) **GRAM** (Cicchetti et al., 2024) using Gramian volume similarity for $M \geq 3$; and (iv) **CoMM** (Dufumier et al., 2024) as a strong fusion-based baseline. For (i)–(iii) we also report "+GW" variants by adding our GW regularizer to show the marginal value of structural alignment. We match optimizer, batch size, temperature, and training epochs across comparable methods; see Appx. A.4.

Table 1 summarizes results across the tasks in Sec. 4.1. Overall, UniOMA with its GW-augmented variants consistently outperform purely contrastive objectives. In particular, adding our GW regularizer (+GW) yields stable gains across all objectives, confirming that structure-aware alignment provides benefits orthogonal to instance discrimination. In the two cells where a baseline is slightly higher (Symile on VFP classification and Vis—Tact), the GW term trades a bit of contrastive correlation for structural coherence. All hyperparameters were kept fixed across methods.

4.4 Analysis of Learned modality weights

UniOMA learns modality weights $\{\lambda_m\}$ that quantify each modality's contribution to the consensus (Appx. A.3). Fig 3 summarizes trends: vision dominates VAT retrieval (high discriminative content);

proprioception dominates VFP contact prediction (contact reasoning); depth is critical for VFD orientation regression, while force contributes marginally. In Sec. 4.5, we downsample one modality (e.g., b=32) while keeping others at b=64. UniOMA maintains both downstream performance and consistent weight patterns; weights shift modestly towards more informative modalities, supporting interpretability under asynchrony.

4.5 ABLATION STUDY: UNEQUAL MODALITY SAMPLING

To evaluate UniOMA's robustness to realistic asynchrony in robot perception, we perform an ablation on the VFD classification task. Specifically, we downsample one modality per batch (vision, force, or depth) by a factor of two, thereby inducing unequal sample counts and breaking strict one-to-one pairing across modalities. We compare UniOMA against its contrastive-only variant (without GW regularization).

Results. Fig. 3(f) shows that UniOMA (Pairwise+GW) outperforms the contrastive-only baseline (Pairwise) across all cases. This confirms that aligning each modality to the GW barycenter consensus, rather than enforcing pairwise matches, enables the model to effectively leverage heterogeneous streams even under sampling-rate mismatch.

Interpretability. Beyond accuracy, UniOMA provides insights into modality importance through its learned weights. Figure 3(e) visualizes the weight distributions under each downsampling setting, showing how the framework adaptively shifts reliance toward intact modalities while still retaining useful signal from the under-sampled one. For comparison, Figure 3(a-d) aggregates the learned weights across the four benchmark datasets (VFP, VFD, MuJoCo, VAT), illustrating task-dependent modality dominance. These results highlight UniOMA's ability to not only maintain structural alignment under unequal sampling but also to yield interpretable modality relevance.

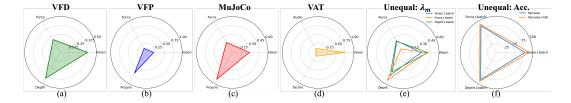


Figure 3: (a–d) Final learned modality weights $\{\lambda_m\}$ for each benchmark (VFD, VFP, MuJoCo Push, VAT). Each radar chart shows per-modality weights that sum to 1, highlighting dataset-specific salience (e.g., depth in VFD, proprioception in VFP) and the interpretability of UniOMA's structural-consensus weighting. (e) ablation on VFD. One modality is downsampled by $\times \frac{1}{2}$ per batch. The plot shows UniOMA's adaptive redistribution of $\{\lambda_m\}$ toward intact modalities while retaining signal from the undersampled one. (f) Accuracy under the same ablation (Top-1, %). Pairwise vs. Pairwise+GW (UniOMA). The outer polygon indicates consistent gains from the GW regularizer across all downsampled cases.

5 DISCUSSION AND CONCLUSION

We revisit multimodal alignment through the lens of *structural* consistency: even when pointwise correspondences are statistically strong, the intra-modal geometries can disagree across modalities. UniOMA closes this gap by combining standard contrastive learning with a GW-barycenter regularizer that can align 3+ modalities to a shared structural consensus. Across VFP, VFD, MuJoCo Push, and VAT, UniOMA improves regression, classification, and cross-modal retrieval while learning interpretable, dataset-specific modality weights. Limitations include the cost of GW subproblems and sensitivity to kernel choices, where our mini-batch barycenter and network-simplex method mitigate these but do not eliminate the costs. We see promising directions in (i) Large-scale real-robot perception alignment with heterogeneous sampling rates, and (ii) extending the consensus to asymmetric similarity matrices (e.g., directed kernels), enabling reasoning about causal dependencies.

6 REPRODUCIBILITY STATEMENT

We will release an anonymized code repository shortly after submission to reproduce all tables and figures end-to-end. For the used datasets, Appx. B provides a complete description of preprocessing and splits for VFP, VFD, MuJoCo Push, and VAT. For theory, Appx. A.1–A.3 contain clear assumptions, derivations, and algorithmic details used in UniOMA.

REFERENCES

- Julián N. Acosta, Guido J. Falcone, Pranav Rajpurkar, and Eric J. Topol. Multimodal biomedical ai. *Nature Medicine*, 28(9):1773–1784, September 2022. ISSN 1546-170X. doi: 10.1038/s41591-022-01981-2.
- Hassan Akbari, Liangzhe Yuan, Rui Qian, Wei-Hong Chuang, Shih-Fu Chang, Yin Cui, and Boqing Gong. Vatt: Transformers for multimodal self-supervised learning from raw video, audio and text. *Advances in neural information processing systems*, 34:24206–24221, 2021.
- Jean-Baptiste Alayrac, Adria Recasens, Rosalia Schneider, Relja Arandjelović, Jason Ramapuram, Jeffrey De Fauw, Lucas Smaira, Sander Dieleman, and Andrew Zisserman. Self-supervised multimodal versatile networks. *Advances in neural information processing systems*, 33:25–37, 2020.
- Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International conference on machine learning*, pp. 214–223. PMLR, 2017.
- Nicolas Bonneel, Michiel van de Panne, Sylvain Paris, and Wolfgang Heidrich. Displacement interpolation using lagrangian mass transport. *ACM Transactions on Graphics*, 30(6):1–12, December 2011. ISSN 1557-7368. doi: 10.1145/2070781.2024192.
- Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Tomas Jackson, Sally Jesmonth, Nikhil J Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Kuang-Huei Lee, Sergey Levine, Yao Lu, Utsav Malla, Deeksha Manjunath, Igor Mordatch, Ofir Nachum, Carolina Parada, Jodilyn Peralta, Emily Perez, Karl Pertsch, Jornell Quiambao, Kanishka Rao, Michael Ryoo, Grecia Salazar, Pannag Sanketi, Kevin Sayed, Jaspiar Singh, Sumedh Sontakke, Austin Stone, Clayton Tan, Huong Tran, Vincent Vanhoucke, Steve Vega, Quan Vuong, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. Rt-1: Robotics transformer for real-world control at scale, 2022.
- Brian Chen, Andrew Rouditchenko, Kevin Duarte, Hilde Kuehne, Samuel Thomas, Angie Boggust, Rameswar Panda, Brian Kingsbury, Rogerio Feris, David Harwath, et al. Multimodal clustering networks for self-supervised learning from unlabeled videos. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 8012–8021, 2021.
- Sihan Chen, Handong Li, Qunbo Wang, Zijia Zhao, Mingzhen Sun, Xinxin Zhu, and Jing Liu. Vast: A vision-audio-subtitle-text omni-modality foundation model and dataset. *Advances in Neural Information Processing Systems*, 36:72842–72866, 2023.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pp. 1597–1607. PmLR, 2020.
- Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 15750–15758, 2021.
- Giordano Cicchetti, Eleonora Grassucci, Luigi Sigillo, and Danilo Comminiello. Gramian multimodal representation learning and alignment. *arXiv preprint arXiv:2412.11959*, 2024.
- Danny Driess, Fei Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, Wenlong Huang, Yevgen Chebotar, Pierre Sermanet, Daniel Duckworth, Sergey Levine, Vincent Vanhoucke, Karol Hausman, Marc Toussaint, Klaus Greff, Andy Zeng, Igor Mordatch, and Pete Florence. Palm-e: an embodied multimodal language model. In *Proceedings of the 40th International Conference on Machine Learning*, ICML'23. JMLR.org, 2023.

- Benoit Dufumier, Javiera Castillo-Navarro, Devis Tuia, and Jean-Philippe Thiran. What to align in multimodal contrastive learning? *arXiv preprint arXiv:2409.07402*, 2024.
- Anirvan Dutta, Etienne Burdet, and Mohsen Kaboli. Visuo-tactile based predictive cross modal perception for object exploration in robotics, 2024.
- Di Feng, Christian Haase-Schütz, Lars Rosenbaum, Heinz Hertlein, Claudius Gläser, Fabian Timm, Werner Wiesbeck, and Klaus Dietmayer. Deep multi-modal object detection and semantic segmentation for autonomous driving: Datasets, methods, and challenges. *IEEE Transactions on Intelligent Transportation Systems*, 22(3):1341–1360, 2021. doi: 10.1109/TITS.2020.2972974.
- Rémi Flamary, Nicolas Courty, Alexandre Gramfort, Mokhtar Z. Alaya, Aurélie Boisbunon, Stanislas Chambon, Laetitia Chapel, Adrien Corenflos, Kilian Fatras, Nemo Fournier, Léo Gautheron, Nathalie T.H. Gayraud, Hicham Janati, Alain Rakotomamonjy, Ievgen Redko, Antoine Rolet, Antony Schutz, Vivien Seguy, Danica J. Sutherland, Romain Tavenard, Alexander Tong, and Titouan Vayer. Pot: Python optimal transport. *Journal of Machine Learning Research*, 22(78):1–8, 2021. URL http://jmlr.org/papers/v22/20-451.html.
- Ruohan Gao, Zilin Si, Yen-Yu Chang, Samuel Clarke, Jeannette Bohg, Li Fei-Fei, Wenzhen Yuan, and Jiajun Wu. Objectfolder 2.0: A multisensory object dataset for sim2real transfer. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10598–10608, 2022.
- Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. Imagebind: One embedding space to bind them all. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 15180–15190, 2023.
- Fengjiao Gong, Yuzhou Nie, and Hongteng Xu. Gromov-wasserstein multi-modal alignment and clustering. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*, CIKM '22, pp. 603–613. ACM, October 2022. doi: 10.1145/3511808. 3557339.
- Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H. Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Koray Kavukcuoglu, Rémi Munos, and Michal Valko. Bootstrap your own latent a new approach to self-supervised learning. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS '20, Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.
- Michelle Guo, Yifeng Jiang, Andrew Everett Spielberg, Jiajun Wu, and Karen Liu. Benchmarking rigid body contact models. In Nikolai Matni, Manfred Morari, and George J. Pappas (eds.), Proceedings of The 5th Annual Learning for Dynamics and Control Conference, volume 211 of Proceedings of Machine Learning Research, pp. 1480–1492. PMLR, 15–16 Jun 2023. URL https://proceedings.mlr.press/v211/guo23b.html.
- Andrey Guzhov, Federico Raue, Jörn Hees, and Andreas Dengel. Audioclip: Extending clip to image, text and audio. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 976–980. IEEE, 2022.
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 9729–9738, 2020.
- Jingjia Huang, Yinan Li, Jiashi Feng, Xinglong Wu, Xiaoshuai Sun, and Rongrong Ji. Clover: Towards a unified video-language alignment and fusion model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14856–14866, 2023.
- Physical Intelligence, Kevin Black, Noah Brown, James Darpinian, Karan Dhabalia, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Manuel Y. Galliker, Dibya Ghosh, Lachy Groom, Karol Hausman, Brian Ichter, Szymon Jakubczak, Tim Jones, Liyiming Ke, Devin LeBlanc, Sergey Levine, Adrian Li-Bell, Mohith Mothukuri, Suraj Nair, Karl Pertsch, Allen Z. Ren, Lucy Xiaoyang Shi, Laura Smith, Jost Tobias Springenberg, Kyle Stachowicz, James Tanner, Quan Vuong, Homer Walke, Anna Walling, Haohuan Wang, Lili Yu, and Ury Zhilinsky. 0.5: a vision-language-action model with open-world generalization, 2025.

Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 4904–4916. PMLR, 18–24 Jul 2021. URL https://proceedings.mlr.press/v139/jia21b.html.

- Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, Quan Vuong, Thomas Kollar, Benjamin Burchfiel, Russ Tedrake, Dorsa Sadigh, Sergey Levine, Percy Liang, and Chelsea Finn. Openvla: An open-source vision-language-action model, 2024.
- Alexander Korotin, Daniil Selikhanovych, and Evgeny Burnaev. Kernel neural optimal transport. *arXiv preprint arXiv:2205.15269*, 2022a.
- Alexander Korotin, Daniil Selikhanovych, and Evgeny Burnaev. Neural optimal transport. *arXiv* preprint arXiv:2201.12220, 2022b.
- Alexander Kraskov, Harald Stögbauer, and Peter Grassberger. Estimating mutual information. *Physical Review E*, 69(6):066138, June 2004. ISSN 1550-2376. doi: 10.1103/physreve.69.066138.
- Jet-Tsyn Lee, Danushka Bollegala, and Shan Luo. "touching to see" and "seeing to feel": Robotic cross-modal sensory data generation for visual-tactile perception. In *2019 International Conference on Robotics and Automation (ICRA)*, pp. 4276–4282. IEEE, May 2019a. doi: 10.1109/icra.2019.8793763.
- Michelle A Lee, Yuke Zhu, Krishnan Srinivasan, Parth Shah, Silvio Savarese, Li Fei-Fei, Animesh Garg, and Jeannette Bohg. Making sense of vision and touch: Self-supervised learning of multimodal representations for contact-rich tasks. In *2019 International conference on robotics and automation (ICRA)*, pp. 8943–8950. IEEE, 2019b.
- Michelle A Lee, Brent Yi, Roberto Martín-Martín, Silvio Savarese, and Jeannette Bohg. Multimodal sensor fusion with differentiable filters. In 2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pp. 10444–10451. IEEE, IEEE, October 2020. doi: 10.1109/iros45743.2020.9341579.
- Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language, 2019.
- Paul Pu Liang, Yiwei Lyu, Xiang Fan, Zetian Wu, Yun Cheng, Jason Wu, Leslie Chen, Peter Wu, Michelle A. Lee, Yuke Zhu, Ruslan Salakhutdinov, and Louis-Philippe Morency. Multibench: Multiscale benchmarks for multimodal representation learning, 2021.
- Weixin Liang, Yuhui Zhang, Yongchan Kwon, Serena Yeung, and James Zou. Mind the gap: understanding the modality gap in multi-modal contrastive representation learning. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA, 2022. Curran Associates Inc. ISBN 9781713871088.
- Jing Liu, Sihan Chen, Xingjian He, Longteng Guo, Xinxin Zhu, Weining Wang, and Jinhui Tang. Valor: Vision-audio-language omni-perception pretraining model and dataset. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. *ViLBERT: pretraining task-agnostic visiolinguistic representations for vision-and-language tasks*. Curran Associates Inc., Red Hook, NY, USA, 2019.
- Huaishao Luo, Lei Ji, Ming Zhong, Yang Chen, Wen Lei, Nan Duan, and Tianrui Li. Clip4clip: An empirical study of clip for end to end video clip retrieval and captioning. *Neurocomputing*, 508: 293–304, 2022.
- Sijie Mai, Ying Zeng, Shuangjia Zheng, and Haifeng Hu. Hybrid contrastive learning of tri-modal representation for multimodal sentiment analysis. *IEEE Transactions on Affective Computing*, 14 (3):2276–2289, 2022.

Karl Øyvind Mikalsen, Filippo Maria Bianchi, Cristina Soguero-Ruiz, and Robert Jenssen. Time series cluster kernel for learning similarities between multivariate time series with missing data. *Pattern Recognition*, 76:569–581, April 2018. ISSN 0031-3203. doi: 10.1016/j.patcog.2017.11. 030.

Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets, 2014.

Seungwhan Moon, Andrea Madotto, Zhaojiang Lin, Alireza Dirafzoon, Aparajita Saraf, Amy Bearman, and Babak Damavandi. Imu2clip: Multimodal contrastive learning for imu motion sensors from egocentric videos and text. *arXiv preprint arXiv:2210.14395*, 2022.

Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Charles Xu, Jianlan Luo, Tobias Kreiman, You Liang Tan, Lawrence Yunliang Chen, Pannag Sanketi, Quan Vuong, Ted Xiao, Dorsa Sadigh, Chelsea Finn, and Sergey Levine. Octo: An open-source generalist robot policy. In *Proceedings of Robotics: Science and Systems*, Delft, Netherlands, 2024.

Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding, 2018.

Gabriel Peyré, Marco Cuturi, and Justin Solomon. Gromov-wasserstein averaging of kernel and distance matrices. In *International conference on machine learning*, pp. 2664–2672. PMLR, 2016.

Ben Poole, Sherjil Ozair, Aaron Van Den Oord, Alex Alemi, and George Tucker. On variational bounds of mutual information. In *International conference on machine learning*, pp. 5171–5180, 2010

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PmLR, 2021.

Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 10684–10695, June 2022.

Adriel Saporta, Aahlad Manas Puli, Mark Goldstein, and Rajesh Ranganath. Contrasting with symile: Simple model-agnostic representation learning for unlimited modalities. *Advances in Neural Information Processing Systems*, 37:56919–56957, 2024.

Pierre Sermanet, Corey Lynch, Yevgen Chebotar, Jasmine Hsu, Eric Jang, Stefan Schaal, and Sergey Levine. Time-contrastive networks: Self-supervised learning from video. 2018 IEEE International Conference on Robotics and Automation (ICRA), pp. 1134–1141, 2017. URL https://api.semanticscholar.org/CorpusID:3997350.

Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke Zettlemoyer, and Dieter Fox. Alfred: A benchmark for interpreting grounded instructions for everyday tasks. In 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). IEEE, June 2020. doi: 10.1109/cvpr42600.2020.01075.

Nina Shvetsova, Brian Chen, Andrew Rouditchenko, Samuel Thomas, Brian Kingsbury, Rogerio S Feris, David Harwath, James Glass, and Hilde Kuehne. Everything at once-multi-modal fusion transformer for video retrieval. In *Proceedings of the ieee/cvf conference on computer vision and pattern recognition*, pp. 20020–20029, 2022.

D. Stewart and J.C. Trinkle. An implicit time-stepping scheme for rigid body dynamics with coulomb friction. In *Proceedings 2000 ICRA*. *Millennium Conference*. *IEEE International Conference on Robotics and Automation*. *Symposia Proceedings (Cat. No.00CH37065)*, volume 1 of *ROBOT-00*, pp. 162–169. IEEE. doi: 10.1109/robot.2000.844054.

- Sandra Steyaert, Marija Pizurica, Divya Nagaraj, Priya Khandelwal, Tina Hernandez-Boussard, Andrew J. Gentles, and Olivier Gevaert. Multimodal data fusion for cancer biomarker discovery with deep learning. *Nature Machine Intelligence*, 5(4):351–362, April 2023. ISSN 2522-5839. doi: 10.1038/s42256-023-00633-5.
- Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. In *European conference on computer vision*, pp. 776–794. Springer, 2020.
- Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In 2012 IEEE/RSJ International Conference on Intelligent Robots and Systems. IEEE, October 2012. doi: 10.1109/iros.2012.6386109.
- Cédric Villani et al. Optimal transport: old and new, volume 338. Springer, 2008.
- Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *Proceedings of the 37th International Conference on Machine Learning*, ICML'20. JMLR.org, 2020.
- Greg Welch and Gary Bishop. An introduction to the kalman filter. Technical report, USA, 1995.
- Jagoda Wojcik, Jiaqi Jiang, Jiacheng Wu, and Shan Luo. A case study on visual-audio-tactile cross-modal retrieval. In 2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pp. 12472–12478. IEEE, 2024.
- Hongwei Xue, Yuchong Sun, Bei Liu, Jianlong Fu, Ruihua Song, Houqiang Li, and Jiebo Luo. Clipvip: Adapting pre-trained image-text model to video-language representation alignment. *arXiv* preprint arXiv:2209.06430, 2022.
- Martina Zambelli, Yusuf Aytar, Francesco Visin, Yuxiang Zhou, and Raia Hadsell. Learning rich touch representations through cross-modal self-supervision. In Jens Kober, Fabio Ramos, and Claire Tomlin (eds.), *Proceedings of the 2020 Conference on Robot Learning*, volume 155 of *Proceedings of Machine Learning Research*, pp. 1415–1425. PMLR, 16–18 Nov 2021. URL https://proceedings.mlr.press/v155/zambelli21a.html.
- Bin Zhu, Bin Lin, Munan Ning, Yang Yan, Jiaxi Cui, HongFa Wang, Yatian Pang, Wenhao Jiang, Junwu Zhang, Zongwei Li, et al. Languagebind: Extending video-language pretraining to n-modality by language-based semantic alignment. *arXiv preprint arXiv:2310.01852*, 2023.
- Sidan Zhu and Dixin Luo. *Enhancing Multi-modal Contrastive Learning via Optimal Transport-Based Consistent Modality Alignment*, pp. 157–171. Springer Nature Singapore, November 2024. ISBN 9789819787951. doi: 10.1007/978-981-97-8795-1_11.

A THEORETICAL DETAILS AND IMPLEMENTATIONS

A.1 EMPIRICAL GW DISTANCE

 Theorem 1 (Empirical GW Distance). Let the kernel matrices $\mathbf{K_x} \in \mathbb{R}^{I \times I}$ and $\mathbf{K_y} \in \mathbb{R}^{J \times J}$ be the similarity matrices conducted by the samples \mathbf{x}, \mathbf{y} from two mm-spaces \mathcal{X}, \mathcal{Y} , the empirical GW distance between the samples is:

$$\hat{d}_{gw}(\mathbf{K}_{\mathbf{x}}, \mathbf{K}_{\mathbf{y}}) := \max_{\mathbf{T} \in \Pi(\hat{\mathbf{p}}_{\mathbf{x}}, \hat{\mathbf{p}}_{\mathbf{y}})} \operatorname{tr}(\mathbf{K}_{\mathbf{x}}^{\top} \mathbf{T}^{\top} \mathbf{K}_{\mathbf{y}} \mathbf{T}),$$

where T is the doubly-stochastic matrix to model the transport between the two sets of samples.

Proof. Let $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^I$ and $\mathcal{Y} = \{\mathbf{y}_j\}_{j=1}^J$ be the two finite mm-spaces with uniform empirical marginals $\hat{\mathbf{p}}_{\mathbf{x}} = \frac{1}{I}\mathbf{1}_I$ and $\hat{\mathbf{p}}_{\mathbf{y}} = \frac{1}{J}\mathbf{1}_J$. Denote their intra-modal similarity matrices by $\mathbf{K}_{\mathbf{x}} \in \mathbb{R}^{I \times I}$ and $\mathbf{K}_{\mathbf{y}} \in \mathbb{R}^{J \times J}$, where $(K_{\mathbf{x}})_{ii'} = \sin(\mathbf{x}_i, \mathbf{x}_{i'})$ and $(K_{\mathbf{y}})_{jj'} = \sin(\mathbf{y}_j, \mathbf{y}_{j'})$. A cross-domain soft matching is a coupling

$$\mathbf{T} \in \Pi(\hat{\mathbf{p}}_{\mathbf{x}}, \hat{\mathbf{p}}_{\mathbf{y}}) := \left\{ \mathbf{T} \ge 0 \mid \mathbf{T} \mathbf{1}_{J} = \hat{\mathbf{p}}_{\mathbf{x}}, \ \mathbf{T}^{\top} \mathbf{1}_{I} = \hat{\mathbf{p}}_{\mathbf{y}} \right\}.$$

The empirical GW distance can be written as the minimum expected squared discrepancy of withindomain relations:

$$\hat{d}_{gw}^{2}(\mathbf{K}_{\mathbf{x}}, \mathbf{K}_{\mathbf{y}}) = \min_{\mathbf{T} \in \Pi(\hat{\mathbf{p}}_{\mathbf{x}}, \hat{\mathbf{p}}_{\mathbf{y}})} \sum_{i,i'} \sum_{j,j'} \left((K_{\mathbf{x}})_{ii'} - (K_{\mathbf{y}})_{jj'} \right)^{2} \mathbf{T}_{ij} \mathbf{T}_{i'j'}. \tag{13}$$

Expand the square in Eq. 13 and group terms:

$$\sum_{i,i',j,j'} \left((K_{\mathbf{x}})_{ii'} - (K_{\mathbf{y}})_{jj'} \right)^2 \mathbf{T}_{ij} \mathbf{T}_{i'j'} = A + B - 2 \sum_{i,i',j,j'} (K_{\mathbf{x}})_{ii'} (K_{\mathbf{y}})_{jj'} \mathbf{T}_{ij} \mathbf{T}_{i'j'},$$

where A, B are constants

$$A = \sum_{i,i'} (K_{\mathbf{x}})_{ii'}^2 \, \hat{\mathbf{p}}_{\mathbf{x}}(i) \, \hat{\mathbf{p}}_{\mathbf{x}}(i'), \qquad B = \sum_{j,j'} (K_{\mathbf{y}})_{jj'}^2 \, \hat{\mathbf{p}}_{\mathbf{y}}(j) \, \hat{\mathbf{p}}_{\mathbf{y}}(j').$$

Therefore, minimizing Eq. 13 is equivalent to maximizing the quadratic term

$$\max_{\mathbf{T} \in \Pi(\hat{\mathbf{p}}_{\mathbf{x}}, \hat{\mathbf{p}}_{\mathbf{y}})} \ \sum_{i,i',j,j'} (K_{\mathbf{x}})_{ii'} \, (K_{\mathbf{y}})_{jj'} \, \mathbf{T}_{ij} \, \mathbf{T}_{i'j'}.$$

In matrix notation, this becomes the quadratic type objective as is in Thrm. 1

$$\hat{d}_{gw}(\mathbf{K}_{\mathbf{x}}, \mathbf{K}_{\mathbf{y}}) = \max_{\mathbf{T} \in \Pi(\hat{\mathbf{p}}_{\mathbf{x}}, \hat{\mathbf{p}}_{\mathbf{y}})} \operatorname{tr}(\mathbf{K}_{\mathbf{x}}^{\top} \mathbf{T}^{\top} \mathbf{K}_{\mathbf{y}} \mathbf{T}).$$
(14)

Consequently, given an optimal plan T^* estimated by Alg. 2,

$$\hat{d}_{aw}(\mathbf{K}_{\mathbf{x}}, \mathbf{K}_{\mathbf{v}}) = \operatorname{tr}(\mathbf{K}_{\mathbf{x}}^{\top} \mathbf{T}^{*\top} \mathbf{K}_{\mathbf{v}} \mathbf{T}^{*}). \tag{15}$$

A.2 OPTIMAL TRANSPORT PLAN ESTIMATION

```
{f Algorithm~2} OTEstimation(\hat{{f K}},{f K})
Input: Kernel matrices \hat{\mathbf{K}} \in \mathbb{R}^{\hat{N} \times \hat{N}}, \mathbf{K} \in \mathbb{R}^{N \times N}
Output: Optimal transport matrix T*
```

Initialize $\mathbf{p} \leftarrow \frac{1}{N} \mathbf{1}_N, \quad \hat{\mathbf{p}} \leftarrow \frac{1}{\hat{N}} \mathbf{1}_{\hat{N}}, \quad \mathbf{T} \leftarrow \hat{\mathbf{p}} \mathbf{p}^{\top}$

while not converged do

```
// Apply Network simplex algorithm:
\hat{\mathbf{T}} \leftarrow \arg \max_{\mathbf{T} \in \Pi(\hat{\mathbf{p}}, \mathbf{p})} \operatorname{tr}(\hat{\mathbf{K}}^{\top} \mathbf{T}^{\top} \mathbf{K} \mathbf{T})
// Line search method to find the minimum:
a \leftarrow -2\operatorname{tr}(\hat{\mathbf{K}}^{\top}\hat{\mathbf{T}}^{\top}\mathbf{K}\mathbf{T})
b \leftarrow \operatorname{tr}((\hat{\mathbf{K}} \odot \hat{\mathbf{K}})\hat{\mathbf{p}}\mathbf{p}^{\top} + \hat{\mathbf{p}}\mathbf{p}^{\top}(\mathbf{K} \odot \mathbf{K})^{\top})
c \leftarrow -2\left(\operatorname{tr}(\hat{\mathbf{K}}^{\top}\mathbf{T}^{\top}\mathbf{K}\hat{\mathbf{T}}) + \operatorname{tr}(\hat{\mathbf{K}}^{\top}\hat{\mathbf{T}}^{\top}\mathbf{K}\mathbf{T})\right)
if a > 0 then
\tau \leftarrow \begin{cases} 1, & \text{if } a+b+c < 0, \\ 0, & \text{otherwise.} \end{cases}
```

return T

 $\mathbf{T} \leftarrow (1 - \tau)\mathbf{T} + \tau \hat{\mathbf{T}}$

Algorithm 2 computes an empirical OT plan T by solving the quadratic program

$$\max_{\mathbf{T} \in \Pi(\hat{\mathbf{p}}, \mathbf{p})} f(\mathbf{T}) := \operatorname{tr}(\hat{\mathbf{K}}^{\top} \mathbf{T}^{\top} \mathbf{K} \mathbf{T}),$$

where $\hat{\mathbf{K}}, \mathbf{K} \in \mathbb{R}^{N \times N}$ are intra-domain similarity (or distance) matrices and $\Pi(\hat{\mathbf{p}}, \mathbf{p}) = \{\mathbf{T} \geq$ $0 \mid \mathbf{T}\mathbf{1} = \hat{\mathbf{p}}, \ \mathbf{T}^{\mathsf{T}}\mathbf{1} = \mathbf{p}$ is the transportation polytope (doubly-stochastic when $\hat{\mathbf{p}} = \frac{1}{\hat{N}}\mathbf{1}_{\hat{N}}, \mathbf{p} =$ $\frac{1}{N}\mathbf{1}_N$). Here \odot is the Hadamard product, so $(\hat{\mathbf{K}}\odot\hat{\mathbf{K}})$ and $(\mathbf{K}\odot\mathbf{K})$ are elementwise squares of the corresponding kernels, which makes b compact. We initialize with the independent coupling $\mathbf{T} = \hat{\mathbf{p}}\mathbf{p}^{\top}$ and iterate a Conditional Gradient (Frank–Wolfe; FW) update.

Network simplex algorithm. At each iteration, we linearize f and solve

$$\hat{\mathbf{T}} \in \arg\max_{\mathbf{T} \in \Pi(\hat{\mathbf{p}}, \mathbf{p})} \langle \mathbf{T}, \nabla f(\mathbf{T}) \rangle.$$

For $f(\mathbf{T}) = \operatorname{tr}(\hat{\mathbf{K}}^{\top} \mathbf{T}^{\top} \mathbf{K} \mathbf{T})$, we use the gradient form

$$\nabla f(\mathbf{T}) = \mathbf{K} \, \mathbf{T} \, \hat{\mathbf{K}} + \mathbf{K}^{\top} \, \mathbf{T} \, \hat{\mathbf{K}}^{\top},$$

which reduces to $2 \mathbf{K} \mathbf{T} \hat{\mathbf{K}}$ when $\mathbf{K}, \hat{\mathbf{K}}$ are symmetric. The oracle is a linear transportation problem. We implement it using a network simplex (Flamary et al., 2021; Bonneel et al., 2011).

Line search. Define the search segment $\mathbf{T}(\tau) = (1 - \tau)\mathbf{T} + \tau \hat{\mathbf{T}}, \tau \in [0, 1]$. Substituting $\mathbf{T}(\tau)$ into f yields a univariate quadratic $f(\tau) = a \tau^2 + b \tau + c$ whose coefficients admit closed forms. The code computes (a, b, c) and picks the maximizer on [0, 1]: $\tau^* = \min(1, \max(0, -(b+c)/(2a)))$ if a>0, otherwise $\tau^{\star}\in\{0,1\}$ by comparing endpoints. We then set $\mathbf{T}=\mathbf{T}(\tau^{\star})$.

A.3 GW BARYCENTER ESTIMATION

Algorithm 3 GW Barycenter Estimation (mini-batch)

Input: Intra-modal similarity matrices $\{\mathbf{K}_{\mathbf{x}}^{(m)}\}_{m=1}^{M}$ (batch size N_m per modality with $\min\{N_m\}=N$), modality weights $\{\lambda_m\}_{m=1}^{M}$ with $\lambda_m\geq 0$, $\sum_m\lambda_m=1$, uniform marginal $\hat{\mathbf{p}}=\frac{1}{N}\mathbf{1}_N$, $\mathbf{p}^{(m)}=\frac{1}{N_m}\mathbf{1}_{N_m}$, max iters T_{\max}

Output: Batch-wise structural consensus (GW barycenter) $\mathbf{C}^*_{\mathbf{x}} \in \mathbb{R}^{N \times N}$

Initialize $\mathbf{C}_{\mathbf{x}}$ as the weighted average of $\mathbf{K}_{\mathbf{x}}^{(m)}$

return $C_{\mathbf{x}}^* \leftarrow C_{\mathbf{x}}$

Consider the barycenter objective (Def. 2):

$$\mathbf{C}_{\mathbf{x}}^* = \arg\min_{\mathbf{C}_{\mathbf{x}} \in \mathcal{M}} \sum_{m=1}^{M} \lambda_m \cdot d_{gw}(\mathbf{C}_{\mathbf{x}}, \mathbf{K}_{\mathbf{x}}^{(m)}), \qquad \lambda_m \ge 0, \ \sum_{m=1}^{M} \lambda_m = 1.$$

According to the discrete empirical GW distance form (Thrm. 1), each term differs from a constant by a (negative) maximized trace. Fix couplings $\{\mathbf{T}^{(m)}\}_{m=1}^{M}$ with $\mathbf{T}^{(m)} \in \Pi(\hat{\mathbf{p}}, \mathbf{p}^{(m)})$ for the current consensus C_x , and define

$$\mathbf{A}^{(m)} := \mathbf{T}^{(m)} \mathbf{K}_{\mathbf{x}}^{(m)} \mathbf{T}^{(m)}^{\top} \in \mathbb{R}^{N \times N}.$$

as C_x -independent constants, the objective reduces to

$$\mathcal{J}(\mathbf{C}_{\mathbf{x}}) = -2 \sum_{m=1}^{M} \lambda_m \operatorname{tr}(\mathbf{C}_{\mathbf{x}}^{\top} \mathbf{A}^{(m)}).$$

Following the standard GW-barycenter normalization (as in Eq. (8) of Gong et al. (2022)), we take the derivative with respect to C and set it to zero

$$\frac{\partial \mathcal{J}(\mathbf{C}_{\mathbf{x}})}{\partial \mathbf{C}_{\mathbf{x}}} = \mathbf{0} \quad \Rightarrow \quad \mathbf{C}_{\mathbf{x}} = \left(\sum_{m=1}^{M} \lambda_m \, \mathbf{A}^{(m)}\right) \, \oslash \, \left(\hat{\mathbf{p}} \, \hat{\mathbf{p}}^{\top}\right),$$

i.e.

$$\mathbf{C}_{\mathbf{x}} \leftarrow \widetilde{\mathbf{C}} \oslash (\hat{\mathbf{p}} \, \hat{\mathbf{p}}^{\top}), \quad \widetilde{\mathbf{C}} = \sum_{m=1}^{M} \lambda_m \, \mathbf{T}^{(m)} \, \mathbf{K}_{\mathbf{x}}^{(m)} \, \mathbf{T}^{(m)}^{\top}.$$
 (16)

Here ⊘ denotes the element-wise division.

A.4 IMPLEMENTATION DETAILS

Implementation: Time-Series Cluster Kernel We use the Time-series Cluster Kernel (TCK; Mikalsen et al. (2018)) to build intra-modal similarity matrices for time-series modalities (e.g., force/torque). TCK fits an ensemble of diagonal covariance Gaussian mixture models (GMMs) with informative priors and computes a posterior membership vector per sample

$$\Pi_i(q) = (\pi_1^{(i)}(q), \dots, \pi_{G_q}^{(i)}(q))^\top, \qquad \sum_{q=1}^{G_q} \pi_g^{(i)}(q) = 1,$$

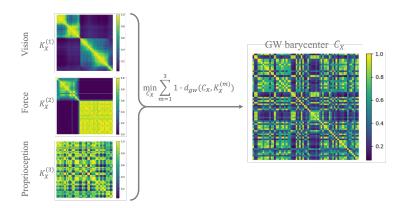


Figure 4: GW barycenter of input-space kernels on the VFP dataset. Left: intra-modal similarity matrices $K_{\mathbf{x}}^{(1)}$ (Vision), $K_{\mathbf{x}}^{(2)}$ (Force), and $K_{\mathbf{x}}^{(3)}$ (Proprioception), each min-max normalized for display. Right: the batch-wise structural consensus $\mathbf{C}_{\mathbf{x}}^*$ obtained by solving

$$\min_{\mathbf{C}_{\mathbf{x}}} \sum_{m=1}^{3} \lambda_{m} \, d_{gw} \big(\mathbf{C}_{\mathbf{x}}, K_{\mathbf{x}}^{(m)} \big) \text{ (with } \lambda_{m} = 1 \text{ here). The barycenter preserves recurrent block/trajec-}$$

tory patterns shared across modalities while smoothing modality-specific artifacts, and is later used to regularize the embedding-space geometry in Stage 2.

where each component $\pi_g^{(i)}(q)$ is the posterior responsibility of mixture g for sequence i under the g-th GMM, i.e.

$$\pi_g^{(i)}(q) \ = \ p\Big(z{=}g \ \Big| \ \mathbf{x}_i^{(q)}; \ \hat{\theta}_q\Big) \,,$$

where z is the latent mixture index, $\hat{\theta}_q$ is the MAP-EM estimate of the q-th model parameters, and $\mathbf{x}_i^{(q)}$ is the subsequence of i restricted to the time window and variable subset chosen by that ensemble member. The final kernel is the sum of posterior inner products over the ensemble:

$$(K_{\text{TCK}})_{ij} \; \leftarrow \; \sum_{q \in \mathcal{Q}} \mathbf{\Pi}_i(q)^\top \mathbf{\Pi}_j(q),$$

which is positive semidefinite as a sum of linear kernels. In practice, for time–critical training, we precompute the full TCK kernel for the entire force dataset (about 10^5 sequences) to get a single symmetric matrix $\mathbf{K}_{\text{TCK}} \in \mathbb{R}^{N \times N}$. During mini-batch training, the intra-modal similarity submatrix for an index set $\mathcal{I} \subseteq \{1,\dots,N\}$ is obtained by simple indexing

$$\mathbf{K}_{\text{batch}} = K_{\text{TCK}}[\mathcal{I}, \mathcal{I}],$$

thus avoiding repeated TCK fits inside the inner learning loop. We follow the original TCK protocol to induce ensemble diversity (random time windows and variable subsets, random initializations, and varying mixture counts), and we cache per-member posteriors to enable fast posterior lookups at test time. See $\S4.1-4.4$ of Mikalsen et al. (2018) for modeling details. In practice, we set the maximal number of mixtures C and the number of randomizations Q as the only user-set hyperparameters. We set C=30 and Q=15 for force/torque signals in the VFP and VFD settings.

Implementation: Pre-trained features We use pre-trained feature extractors for some modalities to produce modality-specific features whose pairwise similarities form the input-space kernels $\{\mathbf{K}_{\mathbf{x}}^{(m)}\}_{m=1}^{M}$ used by our structural consensus $\mathbf{C}_{\mathbf{x}}^{*}$. For the time-series modality (e.g., force/torque), we directly use the TCK method to obtain the input-space kernels.

Pre-trained feature extractors (frozen).

- Vision / Depth / Tactile: Vision Transformer (ViT-B/16; ?) via timm (?), taking the final [CLS] embedding. Single-channel inputs (e.g., depth) are replicated to 3 channels before preprocessing.
- Force: Time-Series Cluster Kernel (TCK; Mikalsen et al. 2018) directly forms $\mathbf{K}_{\mathbf{x}}^{(force)}$ (Sec. A.4).

- Audio (VAT): A frozen Audio Spectrogram Transformer (AST-B, AudioSet-pretrained; ?) on log-mel spectrograms; we take the [CLS] embedding and build $\mathbf{K}_{\mathbf{x}}^{(\text{aud})}$ with a simple similarity (cosine or RBF).
- Other modalities: RBF kernel on frozen features.

Implementation: Modality Encoders To avoid architectural confounds, all methods share identical backbones and training schedules. In UniOMA (Stage 2), each modality encoder $\mathcal{E}_{\theta}^{(m)}$ produces a feature $\mathbf{h}^{(m)} \in \mathbb{R}^{d_h}$, which is passed through a modality-specific MLP projector $g_{\theta}^{(m)}$ to a shared embedding size $d{=}256$:

$$\mathbf{z}^{(m)} = g_{\theta}^{(m)} (\mathcal{E}_{\theta}^{(m)}(\mathbf{x}^{(m)})) \in \mathbb{R}^d, \quad f_{\theta}^{(m)} = g_{\theta}^{(m)} \cdot \mathcal{E}_{\theta}^{(m)}$$

and the embedding-space kernel within a mini-batch is

$$\left(\mathbf{K}_{\mathbf{z}}^{(m)}\right)_{ij} = \exp\left(-\gamma \|\mathbf{z}_{i}^{(m)} - \mathbf{z}_{j}^{(m)}\|_{2}^{2}\right), \qquad \gamma = \frac{20}{d},$$

unless stated otherwise. (Stage 1 input-space kernels $\{\mathbf{K}_{\mathbf{x}}^{(m)}\}$ are computed independently using frozen extractors; see Sec. A.4.)

Backbones.

- Vision / Depth / Tactile. A 2D CNN (ResNet-18). Single-channel inputs (e.g., depth, some tactile) are replicated to 3 channels.
- Force (time series). A 1D temporal ConvNet (temporal Conv–GELU–pool stack).
- **Proprioception.** A 3-layer MLP with ReLU activations.
- Audio (VAT). 1D CNN with three convolutional blocks (channels $1 \rightarrow 64 \rightarrow 128 \rightarrow 256$, kernel size 5, stride 2) each followed by ReLU, then AdaptiveAvgPool1d(1), flatten(), and a final Linear($256 \rightarrow d_b$).

We fix the projector output to d=256, use the same temperature τ for the contrastive term, and share optimizer, batch size, and schedule across methods. UniOMA augments the contrastive loss with a GW-barycenter regularizer (weight α) and learnable modality weights $\{\lambda_m\}$ (softmax-parameterized to enforce $\lambda_m \geq 0$ and $\sum_m \lambda_m = 1$). Encoders and projectors are trained end-to-end with the UniOMA objective; the structure-aware term is computed on $\{\mathbf{K}_{\mathbf{z}}^{(m)}\}$, while Stage 1 kernels $\{\mathbf{K}_{\mathbf{x}}^{(m)}\}$ remain fixed within each epoch.

Implementation: Hyper-parameters Shared training. Unless otherwise noted, all methods use the same backbone-projector settings. We optimize with AdamW (learning rate 3×10^{-4} , weight decay 10^{-4} , β_1 =0.9, β_2 =0.999), batch size 64, and temperature τ =0.1. Each modality head outputs a d=256-dimensional embedding via a lightweight MLP projector (shared width across modalities). We train for 200 epochs with early stopping on the validation metric when applicable, and report mean \pm std over 10 independent seeds.

Stage-1 input-space kernels. Pre-trained feature extractors for vision/depth/tactile (ViT-B/16 via timm) are frozen to compute $\{\mathbf{K}_{\mathbf{x}}^{(m)}\}$. For force/torque we use TCK with max mixtures $C{=}30$ and randomizations $Q{=}15$ following §A.4. For VAT audio, we use AST-B as in Sec. A.4 to form features and then an RBF kernel. To avoid repeated online estimation during Stage 2, we compute force's full dataset kernel once and cache it; mini-batch kernels $\mathbf{K}_{\text{batch}}^{(\text{force})}$ are obtained by submatrix indexing.

Stage-2 embedding-space kernels. All modalities use the same Gaussian kernel

$$\left(\mathbf{K}_{\mathbf{z}}^{(m)}\right)_{ij} = \exp\left(-\gamma \|\mathbf{z}_{i}^{(m)} - \mathbf{z}_{j}^{(m)}\|_{2}^{2}\right),\,$$

with a shared, modality-invariant scale $\gamma = 20/d$, d = 256.

UniOMA-specific. The GW regularization weight $\alpha = 1000$. Modality weights $\{\lambda_m\}$ are learnable with a softmax parameterization $(\lambda_m \ge 0, \sum_m \lambda_m = 1)$ and initialized uniformly. For the coupling

oracle in OTEstimation we use a Frank-Wolfe linearization; the linear subproblem is solved with a network-simplex transportation solver. The line search on the FW segment uses the closed-form quadratic coefficients (a,b,c) derived in Appx. A.2. GW barycenter iterations are run with a maximum of $T_{\rm max}=5$ per inner-loop (in §C.1 we analyze the solidity of this choice).

B DATASETS AND PREPROCESSING

We detail the exact splitting, windowing, and per-modality preprocessing used in our experiments. Unless specified, all randomization uses a fixed seed (seed=42), and splits are performed at the file/trajectory level to avoid leakage.

VFD / VFP (Vision–Force–Depth / Vision–Force–Proprioception). We use test_ratio = 0.2 at the file level with seed = 42 (train vs. test); validation set shares the test set. Each episode has a length 32. For time step t, we form a fixed history window of length L for force (default L=32) and read targets at t+1. RGB images are center–cropped to 128×128 , normalized by ImageNet statistics (mean [0.485, 0.456, 0.406], std [0.229, 0.224, 0.225]). Depth is stored as (128, 128, 1), normalized by mean 0.5/std 0.5, and used as single–channel tensors. Force–torque histories are truncated to the last L steps. The resulting tensor has shape $[b\times L\times 6]$. Proprioception is parsed from the first 7 pose components (end effector position/orientation) in the loader and returned as $[b\times 7]$ at the current step.

Tasks. For **VFD**, we follow the main text: (1) next-step end-effector orientation regression (4D), using $(RGB_t, F/T_{t-L+1:t}, Depth_t)$ as inputs; (2) modality-consistency discrimination with negatives produced by cross-time/trajectory shuffles at 50/50 balance. For **VFP**, we perform next-step contact prediction (binary) using $(RGB_t, F/T_{t-L+1:t}, Proprio_t)$; class balance is enforced by uniform sampling across trajectories.

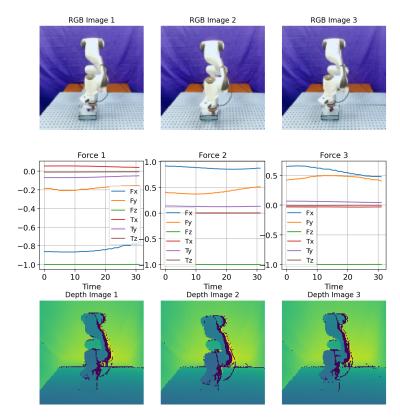


Figure 5: **VFD overview.** Synchronized windows of RGB images, force–torque signals (last L=32 steps), and depth camera images. Images are center–cropped to 128×128 and normalized. Depth images are normalized with mean/std 0.5.

MuJoCo Push. A planar pushing task with a Franka Panda arm interacting with a puck. *Image modality:* we use sequences of grayscale frames. Each sample contains a length S=32 subsequence of 32×32 frames, forming tensors of shape (B,S,1,32,32). Force-torque modality uses the current signal to form tensors of shape (B,6), and end-effector pose modality forms (B,7).

VAT (Vision–Audio–Tactile). We assemble object–level triplets from per–class folders. We use predefined train/val/test directory structures over a fixed object list. Labels for the retrieval tasks are integer–encoded. Visual and tactile images are resized to 246×246 and normalized by ImageNet statistics. Audio is loaded at its native sampling rate; at test time, the raw waveform is truncated to TARGET_LENGTH = 132,300 samples. The final shape of the tensors is (B,132,000)

Task. Cross-modal retrieval with relevance at the object identity level; we report direction-specific MAP on the test set.

C ADDITIONAL EXPERIMENTS

C.1 HYPER-PARAMETER ANALYSIS

RBF kernel scale γ . We use an RBF kernel in the embedding space:

$$\left(\mathbf{K}_{\mathbf{z}}^{(m)}\right)_{ij} = \exp\left(-\gamma_m \|\mathbf{z}_i^{(m)} - \mathbf{z}_j^{(m)}\|_2^2\right).$$

Because distance scales differ by modality, we set γ_m per modality based on empirical pairwise distances at convergence: $\gamma_{\rm vision/depth/tactile} = 5$, $\gamma_{\rm proprio} = 20$, and $\gamma = 10$ for other learnable streams unless stated. Performance is stable within a $\times 0.5 \sim \times 2$ range; very small γ over-smooths similarities, while very large γ over-peaks them.

Number of GW barycenter iterations T_{\max} . Let $\mathbf{C}^{(t)}$ be the consensus at inner-loop iteration t in Alg. 3. We monitor the relative Frobenius change $\Delta_t = \|\mathbf{C}^{(t)} - \mathbf{C}^{(t-1)}\|_F / \|\mathbf{C}^{(t-1)}\|_F$ and the trace objective $\sum_m \lambda_m \operatorname{tr} (\mathbf{C}^{(t)\top} \mathbf{T}^{(m)\top} \mathbf{K}_{\mathbf{x}}^{(m)} \mathbf{T}^{(m)})$. Both stabilize rapidly; after t=5 further changes are negligible $(\Delta_t < 10^{-3})$. We therefore fix $T_{\max} = 5$ for all reported results.

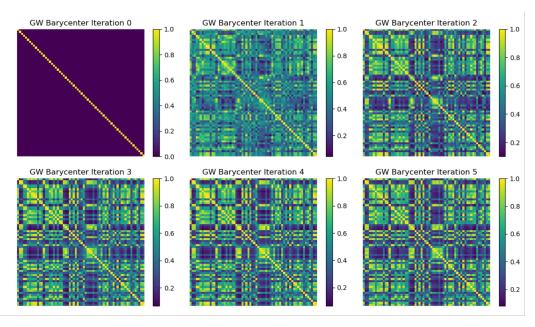


Figure 6: GW barycenter inner-loop: structural consensus across iterations t in Alg. 3. By t=5, both geometry and objective are effectively stable, thus we choose $T_{\rm max}=5$.

D LLM USAGE STATEMENT

This work does not incorporate large language models (LLMs) as a key, novel, or unconventional component of the method, experiments, or analysis. Any LLM assistance was limited to the writing refinement (grammar, clarity, and copy-editing). All technical formulation, algorithms, proofs, hyperparameters, implementations, and results were created and validated by the authors.