
Group Equivariance Meets Mechanistic Interpretability: Equivariant Sparse Autoencoders

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Adapting sparse autoencoders (SAEs) to domains beyond language, such as sci-
2 entific data with group symmetries, introduces challenges that can hinder their
3 effectiveness. We show that incorporating such group symmetries into the SAEs
4 yields features more useful in downstream tasks. More specifically, we train au-
5 toencoders on synthetic images and find that a single matrix can explain how their
6 activations transform as the images are rotated. Building on this, we develop *adap-*
7 *tively equivariant SAEs* that can adapt to the base model’s level of equivariance.
8 These adaptive SAEs discover features that lead to superior probing performance
9 compared to regular SAEs, demonstrating the value of incorporating symmetries
10 in mechanistic interpretability tools.

11 1 Introduction

12 Sparse autoencoders (SAEs) are increasingly used in domains beyond language, particularly with
13 scientific data such as proteins [1, 2, 3, 4] and cell images [5, 6]. Mechanistic interpretability of
14 scientific models can help us detect missing labels in our datasets [2], steer model outputs such as
15 generated proteins [1], and discover novel predictors for quantities of interest [3]. Scientific data is
16 frequently characterized by underlying symmetries: transformations such as rotations or translations
17 that alter particular attributes in particular ways. Although accounting for those symmetries can lead
18 to more data-efficient models [7], applications of SAEs in such settings overlook those symmetries.
19 In this paper, we present early results suggesting that building SAEs that automatically adapt to
20 symmetries in the data can greatly improve their performance in downstream tasks.

21 A set of symmetries such as rotations can be modeled as a **group** G . Groups **act on** sets such as
22 images; e.g. with $g \in G$, gx can denote the rotated version of an image x . Transformations of the
23 same element x form an **orbit** $\{gx : g \in G\}$. We can then split the features of x into those that
24 are **invariant** with respect to G and those that are **equivariant**. For example, the types of atoms
25 in a molecule would be features invariant under 3D rotations, but the force vectors acting on each
26 atom would be equivariant features, rotating along with the molecule. More generally, we define
27 invariant features as those shared across an orbit, while equivariant features depend on the particular
28 transformation applied to x . We provide a more detailed background on groups and symmetries in
29 Appendix B.

30 This reveals two pitfalls for training SAEs on data with group symmetries. *First*, the optimally sparse
31 solution learns one latent per transformation for each equivariant feature, requiring $O(|G|)$ latents
32 per semantic feature, which is undesirable for larger groups. *Second*, since we do not know a priori
33 the degree of equivariance in base model activations, we should *adapt* to what the model has learned
34 rather than prematurely enforcing exact symmetries. We will demonstrate that designing SAEs
35 while avoiding these two pitfalls can lead to improved performance in downstream applications.
36 More specifically:

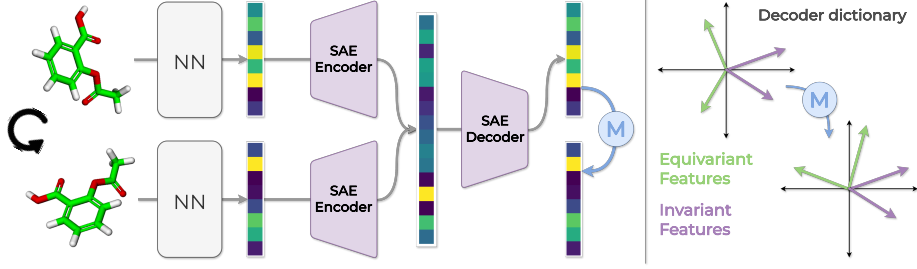


Figure 1: **Overview of our approach.** *Left:* We train an invariant SAE that maps activations of the transformed inputs to the same latents, and optimize the matrix \mathbf{M} to estimate how the activations transform. *Right:* Transforming the decoder dictionary $\mathbf{D} \mapsto \mathbf{MD}$ allows us to observe which features discovered by the SAE are **invariant** or **equivariant** with respect to input transformations.

- We show that a single matrix \mathbf{M} acting on the activations explains more than 98% of the variance in the transformed activations on MLP and CNN-based autoencoders trained on a synthetic image dataset transformed under the group of 90° rotations.
- We build on this observation to design **adaptively equivariant SAEs** consisting of an *invariant* autoencoder (avoiding the pitfall of exploding number of latents) which is made *equivariant* with the addition of \mathbf{M} (avoiding the pitfall of unnecessarily-exact equivariance).
- We demonstrate **our adaptively equivariant SAEs learn more useful features that outperform regular SAEs in a set of binary probing tasks over our synthetic dataset** despite lagging behind in the reconstruction/sparsity frontier.

2 SAEs with adaptive equivariance

We consider groups G where all transformations can be obtained as powers of a generator $g \in G$, i.e. $G = \{g^p\}_{p=1}^{|G|}$. Our design consists of a group-invariant TopK SAE [8, 9] with a two-layer encoder E and linear decoder D , and a matrix \mathbf{M} that learns to fit how the base model’s activations transform as its inputs are transformed with the action of group G . Thus, with canonical (one representative per orbit) inputs $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^n$ and model activations $\psi(\mathbf{x}) \in \mathbb{R}^d$, for $p = 1, \dots, |G|$,

$$D(E(\psi(g^p \mathbf{x}))) = \psi(\mathbf{x}), \quad \text{and} \quad \mathbf{M}^p D(E(\psi(g^p \mathbf{x}))) = \psi(g^p \mathbf{x}). \quad (1)$$

First, the SAE reconstructs all activations $\psi(g^p \mathbf{x})$ of the transformed inputs as the canonical activation $\psi(\mathbf{x})$. Then this reconstruction is transformed with \mathbf{M} to obtain $\mathbf{M}^p \psi(\mathbf{x}) = \psi(g^p \mathbf{x})$.

Invariant SAE. We make our SAE invariant to group transformations of the base model’s inputs by training it with the following *invariance loss*:

$$\mathcal{L}_{\text{inv}} := \mathbb{E}_{\mathbf{x} \in \mathcal{X}, p=1, \dots, |G|} \|\psi(\mathbf{x}) - D(E(\psi(g^p \mathbf{x})))\|_2^2 \quad (2)$$

After observing that a linear encoder may fail to learn to be invariant, we build our encoder with two linear layers, but keep the decoder to one layer. Therefore, while the encoder is expressive enough to learn to be invariant, the canonical activations are still reconstructed as sparse linear combinations of the dictionary vectors.

Activation transformation matrix \mathbf{M} . To map the canonical reconstructions back to their original forms, we need to learn how the base model’s activations transform as its inputs are transformed. While closed-form solutions might be possible for certain cases, they are not practical for arbitrary neural networks. Instead, we hypothesize that a linear transformation should be able to explain to the transformation in the model’s activations, since many group actions we care about, including rotations, can be represented as linear transformations. Thus, we optimize $\mathbf{M} \in \mathbb{R}^{d \times d}$ to minimize

$$\mathcal{L}_M := \mathbb{E}_{\mathbf{x} \in \mathcal{X}, p=1, \dots, |G|} \|\psi(g^p \mathbf{x}) - \mathbf{M}^p \psi(\mathbf{x})\|_2^2. \quad (3)$$

We initialize \mathbf{M} as the identity matrix so that it fits perfectly right away if the model has learned invariant representations, and we optimize it using the Adam optimizer [10].

3 Evaluation

Dataset and models. We create a synthetic image dataset where each image contains four shapes. There are 8 possible shapes, each can be in one of four positions (see Figure 2). Applying 90° rotations to the images yields either two or four possible orientations for each shape. We train MLP and CNN-based autoencoders as base models where the task is to compress and reconstruct the image (see Appendix C for training details).



Figure 2: Example images.

SAE setup and baselines. We train all SAEs in our experiments over the 256-dimensional middle-layer activations of the base models, and compare our equivariant SAE with two regular TopK SAEs (linear encoder and decoder). The equivariant SAE and the first regular SAE both have an expansion factor of 8, resulting in 2,048 latents. The second regular SAE has $|G|$ times the number of latents, corresponding to learning separate latents for each orientation of equivariant features for a total of 8,192. The regular SAEs are trained by augmenting the dataset with 90° rotations.

Probing. We define 180 binary probing tasks to evaluate the downstream usefulness of the features discovered by SAEs (see Appendix C). The tasks are split into four subsets, based on if a shape is in an image (**S**), and in a specific position (**SP**), in a specific orientation (**SO**), and in both a position and an orientation (**SPO**). Note that only the **S** tasks are invariant to rotations. It is desirable that a small number of latents encode the concepts used to separate the images in the binary probing task. Thus we first identify a small set of latents that maximally differ between the two classes by filtering the SAE latents with the highest absolute difference between the two classes [11]. Then for each task, we train three different probes over the truncated latents as well as the activations, and report results from the best performing probe averaged across all tasks. The probing methods are kNN, logistic regression, and XGBoost, with the XGBoost performing the best overall.

3.1 Results

RESULT 1: M can be learned effectively. Over the middle layer activations of both of our base autoencoders, we optimize M for 150 epochs and observe that it can be learned with an average $R^2 > 0.98$ between the ground truth and predicted activations across all transformations. As a naive baseline, setting $M = I$ results in average R^2 values 0.05 and 0.49 for CNN and MLP autoencoders, respectively. These results support our hypothesis that the activation-space transformation can be explained by a linear transformation.

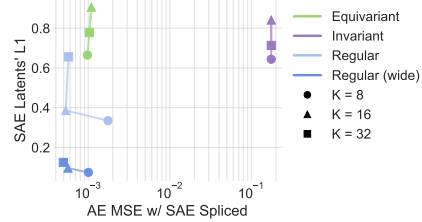


Figure 3: **Sparsity/reconstruction performance of SAEs** for varying TopK values. The x axis corresponds to the base autoencoders’ reconstruction performance when their intermediate activations are passed through the SAEs, and the y axis denotes the L1 norm of the SAE latent activations.

RESULT 2: Equivariant SAEs lag behind in the reconstruction/sparsity frontier. Figure 3 displays sparsity (L1) and the reconstruction (loss when SAE is spliced into the base AE) performance of the SAEs in our experiments. While for small K values the equivariant SAE matches the wide SAE’s reconstruction performance, regular SAEs generally have sparser latents. The improvement in reconstruction from the invariant to the equivariant SAE further shows the effect of learning M .

RESULT 3: Despite worse overall reconstructions, equivariant SAEs learn more useful directions. Figure 4 shows the classification performance on the 180 tasks of the XGBoost probes with the CNN autoencoder for $K = 16$ and truncation lengths of 8 and 32. (see Appendix E for results with different K and the MLP autoencoder, and Appendix D for ablations). The main outcomes are as follows:

- When probing over the SAE latents, the invariant and equivariant SAEs outperform regular SAEs on the invariant **S** group of tasks. Their performance unsurprisingly drops for equivariant tasks since the latent activations are learned to be invariant, although in some cases they might still outperform regular SAEs such as when the truncation length is 8.
- Across all tasks and setups, the invariant and equivariant SAEs perform the best when probes are trained over the truncated reconstructions, even matching the performance of probes trained

over the base model activations (taken as the upper bounds). This is our primary result, as it demonstrates that despite the lower overall reconstruction performance (Figure 3), **the equivariant SAE learns more informative directions in activation space that can lead to increased performance on group-structured probing tasks.**

4 Related work

Our work is one of the first to bridge ideas from the sparse dictionary learning literature and the equivariant representation learning literature with a particular application towards mechanistic interpretability. Learning approximate equivariance via objectives similar to Equation 2 has been proposed in [12, 13], and our general approach of learning an invariant encoder and a separate mapping from the canonical outputs to their original forms follows that of [14]. Our main difference is that we adapt these approaches to the space of neural network activations, where the symmetries are induced by the input transformations and are not well-defined.

Group-equivariant sparse dictionary learning methods have also been proposed [15, 16] although such exact symmetries cannot be enforced over neural network activations as we do not know to what degree they are equivariant. Finally, [17] proposes group crosscoders to analyze how the features learned by a neural network change as its inputs are transformed, constructing each dictionary vector with G blocks each the size of the inputs. The size of the dictionary thus scales linearly with $|G|$ unlike our approach where the number of parameters is constant with respect to $|G|$.

5 Conclusion

We have presented early results showing that adding domain-specific properties such as group equivariance to sparse autoencoders can significantly increase their utility in domains beyond language. Our first result is that a single matrix can explain more than 98% of the variance in how the activations of a neural network transform as its inputs are transformed with the action of a discrete cyclic group. We then used this result to design equivariant SAEs that discover features that lead to better probing performance than regular SAEs, indicating that they are more useful in downstream applications.

Limitations and future work. Although our results are promising, they are so far limited to a synthetic task, relatively small models, and a small discrete group, and thus many important questions remain to establish the practical usefulness of adaptively equivariant SAEs: Can \mathbf{M} be learned as effectively in larger models such as frontier foundation models? Can the optimization of \mathbf{M} be improved by better utilizing the group structure, e.g. by constraining the optimization to $\mathbf{M}^{|G|} = \mathbf{I}$? Does the two-layer encoder qualitatively alter what kind of features are discovered, and how can the features best be labeled incorporating the knowledge of how they transform with \mathbf{M} ? Finally, can the reconstruction performance of the equivariant SAE be improved to match that of regular SAEs?

While the concept of interpretability is domain-agnostic, progress in mechanistic interpretability has largely been driven by its applications in language, which has led to certain concepts such as group equivariance being represented far less prominently than they are in the broader ML literature. Our results highlight the potential benefit of bridging that gap and tailoring mechanistic interpretability tools for domains beyond language, despite being early results for a work in progress.

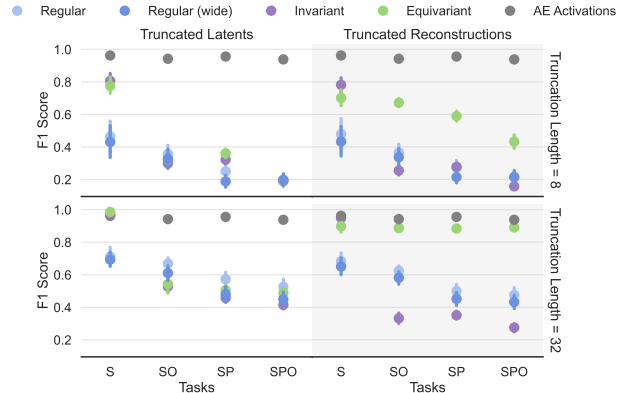


Figure 4: **Latent and reconstruction probing performance**, with $K = 16$ and the CNN autoencoder. Performance of probes over the base autoencoder activations are duplicated for easier comparisons.

References

- [1] Nithin Parsan, David J. Yang, and John Jingxuan Yang. Towards Interpretable Structure Prediction With Sparse Autoencoders. In *ICLR 2025 Workshop on Machine Learning Multiscale Processes*, March 2025.
- [2] Elana Simon and James Zou. InterPLM: Discovering Interpretable Features in Protein Language Models via Sparse Autoencoders. page 2024.11.14.623630, January 2025.
- [3] Etowah Adams, Liam Bai, Minji Lee, Yiyang Yu, and Mohammed AlQuraishi. From Mechanistic Interpretability to Mechanistic Biology: Training, Evaluating, and Interpreting Sparse Autoencoders on Protein Language Models. page 2025.02.06.636901, February 2025.
- [4] Edith Natalia Villegas Garcia and Alessio Ansuini. Interpreting and Steering Protein Language Models through Sparse Autoencoders. (arXiv:2502.09135), February 2025.
- [5] Muhammed Furkan Dasdelen, Hyesu Lim, Michele Buck, Katharina S. Götze, Carsten Marr, and Steffen Schneider. CytoSAE: Interpretable Cell Embeddings for Hematology. (arXiv:2507.12464), July 2025.
- [6] Konstantin Donhauser, Kristina Ulicna, Gemma Elyse Moran, Aditya Ravuri, Kian Kenyon-Dean, Cian Eastwood, and Jason Hartford. Towards scientific discovery with dictionary learning: Extracting biological concepts from microscopy foundation models. (arXiv:2412.16247), February 2025.
- [7] Johann Brehmer, Sönke Behrends, Pim de Haan, and Taco Cohen. Does equivariance matter at scale? (arXiv:2410.23179), October 2024.
- [8] Leo Gao, Tom Dupré la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya Sutskever, Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders. (arXiv:2406.04093), June 2024.
- [9] Alireza Makhzani and Brendan Frey. K-Sparse Autoencoders. (arXiv:1312.5663), March 2014.
- [10] Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. (arXiv:1412.6980), January 2017.
- [11] Subhash Kantamneni, Joshua Engels, Senthoooran Rajamanoharan, Max Tegmark, and Neel Nanda. Are Sparse Autoencoders Useful? A Case Study in Sparse Probing. (arXiv:2502.16681), February 2025.
- [12] Ahmed A. Elhag, T. Konstantin Rusch, Francesco Di Giovanni, and Michael Bronstein. Relaxed Equivariance via Multitask Learning. (arXiv:2410.17878), October 2024.
- [13] Yinzhu Jin, Aman Shrivastava, and P. T. Fletcher. Learning Group Actions on Latent Representations. *Advances in Neural Information Processing Systems*, 37:127273–127295, December 2024.
- [14] Robin Winter, Marco Bertolini, Tuan Le, Frank Noe, and Djork-Arné Clevert. Unsupervised Learning of Group Invariant and Equivariant Representations. In *Advances in Neural Information Processing Systems*, May 2022.
- [15] Christian Shewmake, Nina Miolane, and Bruno Olshausen. Group Equivariant Sparse Coding. In Frank Nielsen and Frédéric Barbaresco, editors, *Geometric Science of Information*, volume 14071, pages 91–101. Springer Nature Switzerland, Cham, 2023.
- [16] Christian Shewmake, Domas Buracas, Hansen Lillemark, Jinho Shin, Erik Bekkers, Nina Miolane, and Bruno Olshausen. Visual Scene Representation with Hierarchical Equivariant Sparse Coding. December 2023.
- [17] Liv Gorton. Group Crosscoders for Mechanistic Analysis of Symmetry. (arXiv:2410.24184), November 2024.

- 221 [18] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
222 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas
223 Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy,
224 Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An Imperative Style,
225 High-Performance Deep Learning Library. (arXiv:1912.01703), December 2019.
- 226 [19] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion,
227 Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vander-
228 plas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Édouard
229 Duchesnay. Scikit-learn: Machine Learning in Python. *J. Mach. Learn. Res.*, 12(null):2825–
230 2830, November 2011.
- 231 [20] Tianqi Chen and Carlos Guestrin. XGBoost: A Scalable Tree Boosting System. In *Proceedings*
232 *of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Min-*
233 *ing*, KDD ’16, pages 785–794, New York, NY, USA, August 2016. Association for Computing
234 Machinery.
- 235 [21] G. Bradski. The OpenCV Library. *Dr. Dobb’s Journal of Software Tools*, 2000.
- 236 [22] Michael M. Bronstein, Joan Bruna, Taco Cohen, and Petar Veličković. Geometric Deep Learn-
237 ing: Grids, Groups, Graphs, Geodesics, and Gauges. (arXiv:2104.13478), May 2021.
- 238 [23] Johannes Gasteiger, Janek Groß, and Stephan Günnemann. Directional Message Passing for
239 Molecular Graphs. (arXiv:2003.03123), April 2022.
- 240 [24] Omri Puny, Matan Atzmon, Heli Ben-Hamu, Ishan Misra, Aditya Grover, Edward J. Smith,
241 and Yaron Lipman. Frame Averaging for Invariant and Equivariant Network Design.
242 (arXiv:2110.03336), March 2022.
- 243 [25] Alexandre Duval, Simon V. Mathis, Chaitanya K. Joshi, Victor Schmidt, Santiago Miret,
244 Fragkiskos D. Malliaros, Taco Cohen, Pietro Liò, Yoshua Bengio, and Michael Bronstein. A
245 Hitchhiker’s Guide to Geometric GNNs for 3D Atomic Systems. (arXiv:2312.07511), March
246 2024.

247 A Reproducibility

248 We make our SAE and probing implementations public at [https://anonymous.4open.science/](https://anonymous.4open.science/r/equivariant-sae)
249 [r/equivariant-sae](https://anonymous.4open.science/r/equivariant-sae). We implement our models from scratch using PyTorch [18], and use the
250 scikit-learn [19] and XGBoost [20] packages for the various probing methods. We utilize the
251 OpenCV package [21] to create our dataset.

252 B Background

253 B.1 Symmetries & group equivariance

254 Data in scientific problems often involve various **symmetries**, transformations that preserve certain
255 properties of the data while some properties can transform along with the symmetries. For exam-
256 ple, rotating a molecule, moving it in space, or permuting its identical atoms does not change the
257 molecule’s identity, but rotating it can change the orientation of certain vector quantities such as the
258 forces acting on each atom.

259 The frame of reference we associate with such data, e.g. the particular 3D coordinates we assign
260 to atoms in a molecule, is not an inherent property of the molecule or the world but an artifact of
261 our observational bias. Thus to model physical phenomena more faithfully, we would prefer to be
262 independent of particular reference frames, and developing such tools has become an active research
263 area [22].

264 Symmetries share certain properties such as being composable (subsequent rotations can also be
265 modeled as a single rotation) and invertible (any rotation can be inverted). Moreover, the identity
266 transform is a trivial symmetry for any object, and the order of composition of three transforma-
267 tions does not change the end result. These notions are unified by the definition of a **group** that
268 characterizes symmetry transformations:

269 **Definition 1** (Group). A **group** $(G, *)$ is a set G along with an operation $*$: $G \times G \rightarrow G$ such that
270 the following axioms are satisfied:

- 271 • (Associativity) For all $g, h, j \in G$, it holds that $(g * h) * j = g * (h * j)$.
- 272 • (Identity element) There exists $e \in G$ such that $e * g = g$ for all $g \in G$.
- 273 • (Inverses) Any $g \in G$ has an inverse $g^{-1} \in G$ such that $g * g^{-1} = e$.

274 A group is further called **abelian** if its group operation is commutative, i.e. $g * h = h * g$ for all
275 $g, h \in G$.

276 Groups are often denoted by their set, e.g. as G alone, omitting the operation. Groups can be
277 *discrete*, such as permutation groups S_n of n objects and the cyclic groups C_n corresponding to
278 rotations of an n -gon, or *continuous*, such as the group of rotations in n -dimensional space, defined
279 as $SO(n) := \{\mathbf{R} \in \mathbb{R}^{n \times n} : \mathbf{R}\mathbf{R}^T = \mathbf{I}, \det \mathbf{R} = 1\}$.

280 Groups transform sets of objects such as images or molecules via their **actions**:

281 **Definition 2** (Group action). A (left) **action** of the group $(G, *)$ on set X is a map $\alpha : G \times X \rightarrow X$,
282 denoted

$$(g, x) \mapsto \alpha(g, x) = g \cdot x,$$

283 that satisfies these axioms for all $x \in X$:

- 284 • $\alpha(e, x) = x$ with e the identity element in G .
- 285 • $\alpha(g, \alpha(h, x)) = \alpha(g * h, x)$ for all $g, h \in G$.

286 A right group action can similarly be defined, and the set X is said to be a **G -set**.

287 **Definition 3** (Orbit). Let X be a G -set. The **orbit** of $x \in X$ is the set of all points in X reachable
288 by transforming x with G , denoted

$$Gx := \{\alpha(g, x) : g \in G\}.$$



Figure 5: **Set of available shapes in our dataset.** None of the shapes is rotation-invariant, with the horizontal rectangle and diagonal line having two orientations and the other six shapes having four orientations.

Functions mapping between G -sets can then be **invariant** or **equivariant** depending on how they behave in each orbit:

Definition 4 (Invariant and equivariant functions). *For G -sets X and Y with associated actions α_X, α_Y , a function $f : X \rightarrow Y$ is G -invariant if for all $g \in G, x \in X$,*

$$f(\alpha_X(g, x)) = f(x),$$

and G -equivariant if

$$f(\alpha_X(g, x)) = \alpha_Y(g, f(x)).$$

Using neural networks, approximate invariance or equivariance can be achieved by augmenting the data with symmetric inputs, or explicitly via additional loss terms [12]. Straightforward ways of achieving exact invariance include limiting the model’s inputs to invariant properties of the data, such as internal bond angles in a molecule that are rotation-invariant scalars [23], or averaging the outputs over each orbit [24]. Achieving exact equivariance requires more careful consideration of how the input features are processed in each layer of the neural network, but there exists a wide-ranging literature of equivariant models for various groups and data types [25].

C Experimental details

C.1 Dataset and probing

Figure 5 displays the base shapes in our dataset. When rotated in increments of 90° , the rectangle and the diagonal shapes have two orientations, and the other six shapes have four. Each image then contains a randomly sampled shape in each of its four quadrants. Precise definitions of our binary probing tasks are then as follows, with a shape’s position denoting which of the four quadrants it is in in an image, and its orientation denoting which of the four or two orientations it is in:

- **S**(s): Does the image contain shape s in any position or orientation?
- **SO**(s, o): Does the image contain shape s in orientation o and any position?
- **SP**(s, p): Does the image contain shape s in position p and any orientation?
- **SPO**(s, p, o): Does the image contain shape s in position p and orientation o ?

This results in a total of 8 **S** (one for each shape), 28 **SO** (2 shapes \times 2 orientations + 6 shapes \times 4 orientations), 32 **SP** (8 shapes \times 4 orientations), and 112 **SPO** (2 shapes \times 2 orientations \times 4 positions + 6 shapes \times 4 orientations \times 4 positions) tasks, for a total of 180 tasks. Note the tasks are not balanced and contain more negative examples than positive examples, and hence we report F1 scores rather than raw accuracies.

The probe we ultimately report the results from, XGBoost [20], consists of 100 estimators with a maximum depth of 6, and is trained with the learning rate 0.3 and L2 regularization.

C.2 Base models

We train our base autoencoders for 100 epochs over 10,000 randomly generated samples from our dataset and augmenting with random 90° rotations with a batch size of 64 using Adam [10] with learning rate $1e-3$. Their architectures are detailed in Table 1.

Table 1: **Architectures of the MLP and CNN autoencoders.** The first section of each models corresponds to the encoder and the second section to the decoder. We train our SAEs over the pre-activation encoder outputs.

MLP	CNN
Input: 4096 (64×64)	Input: $1 \times 64 \times 64$
Linear(4096, 256)	Conv2d(1, 16, 3×3 , stride=2, pad=1)
ReLU	ReLU
Linear(256, 256)	Conv2d(16, 32, 3×3 , stride=2, pad=1)
ReLU	ReLU
	Conv2d(32, 256, 16×16)
	ReLU
Linear(256, 256)	ConvTranspose2d(256, 32, 16×16)
ReLU	ReLU
Linear(256, 4096)	ConvTranspose2d(32, 16, 3×3 , stride=2, pad=1, out_pad=1)
	ReLU
	ConvTranspose2d(16, 1, 3×3 , stride=2, pad=1, out_pad=1)

323 C.3 SAEs

324 The regular SAEs used in our comparisons are typical TopK SAEs with linear encoders and decoders.
325 The equivariant SAEs also have linear decoders and use the TopK activation, but their encoders
326 consists of two linear layers with a ReLU activation in between and hidden dimension of 512. We
327 train our sparse autoencoders for 500 epochs over 10,000 samples from our dataset with batch size
328 64 using Adam [10] with a learning rate of $1e-3$.

329 D Ablations

330 Figure 6 displays the performance of XGBoost probes trained over the truncated latent activations
331 and the reconstructions of the equivariant SAE, and after separately using a linear instead of a two-
332 layer encoder, and replacing \mathcal{L}_{inv} with the typical reconstruction loss.

333 We observe that when probing over the truncated latent activations, the equivariant SAE latents result
334 in the best probes for the invariant **S** tasks, and ablating the invariance loss increases the performance
335 in the equivariant tasks. This is expected and in a similar line with the results in Section 3.

336 When probing over the truncated reconstructions however, the equivariant SAE results in the best-
337 performing probes for all tasks except in the $K = 32$, truncation length = 32 setup, where ablating
338 the invariance loss leads to slightly better probes.

339 Overall, these results are evidence for the hypothesis that the better performance of the reconstruc-
340 tion probes we have observed in Section 3 requires both the invariance loss and the two-layer en-
341 coder.

342 E Further probing results

343 Figure 7 displays further probing results with TopK values 8 and 32. Results generally agree
344 with those presented in Section 3. When probing over the truncated reconstructions, the invari-
345 ant/equivariant SAEs result in the most accurate probes over all tasks while the performance of the
346 latent activation probes drop with the equivariant tasks as expected.

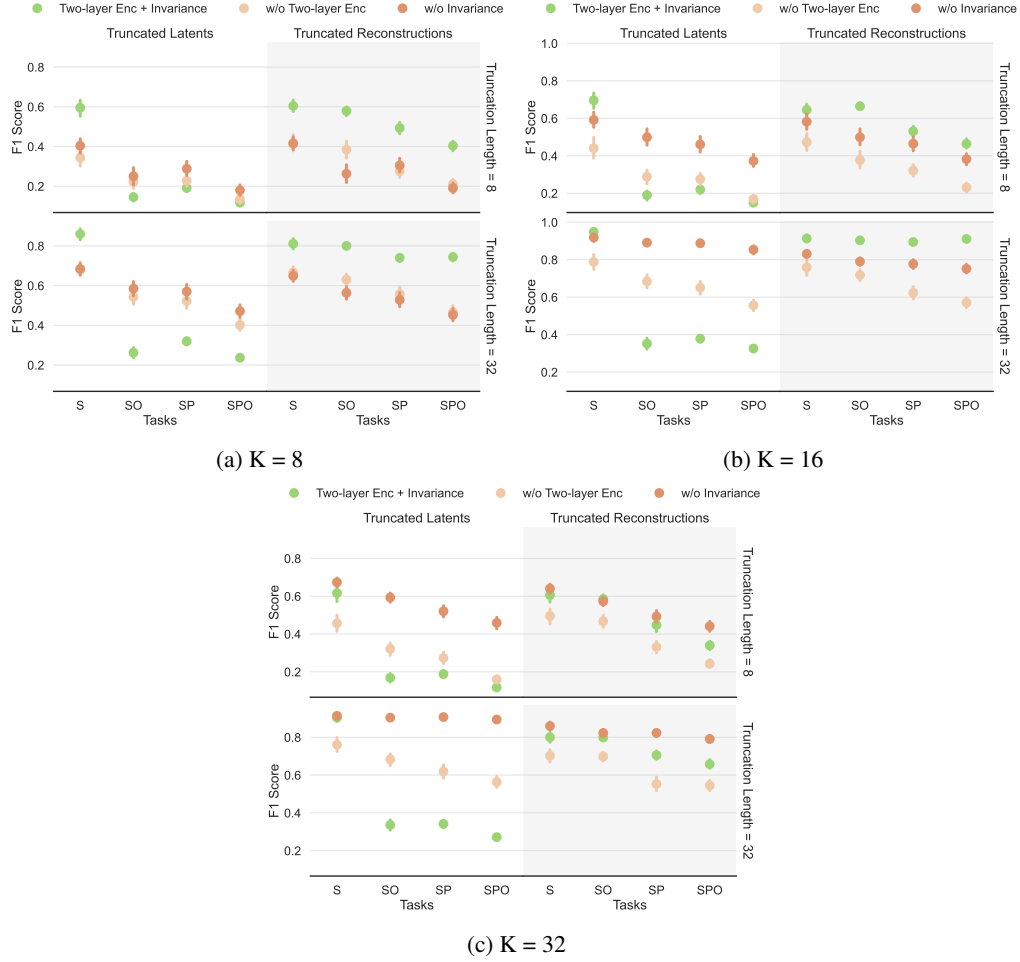


Figure 6: **Probing results after separately ablating the two-layer encoder and the invariance loss from our equivariant SAE.** Ablating either components reduces the performance of latent activation probing for the invariant **S** task, and likewise the performance of the reconstruction probes over all tasks except in the $K = 32$, truncation length = 32 setup.

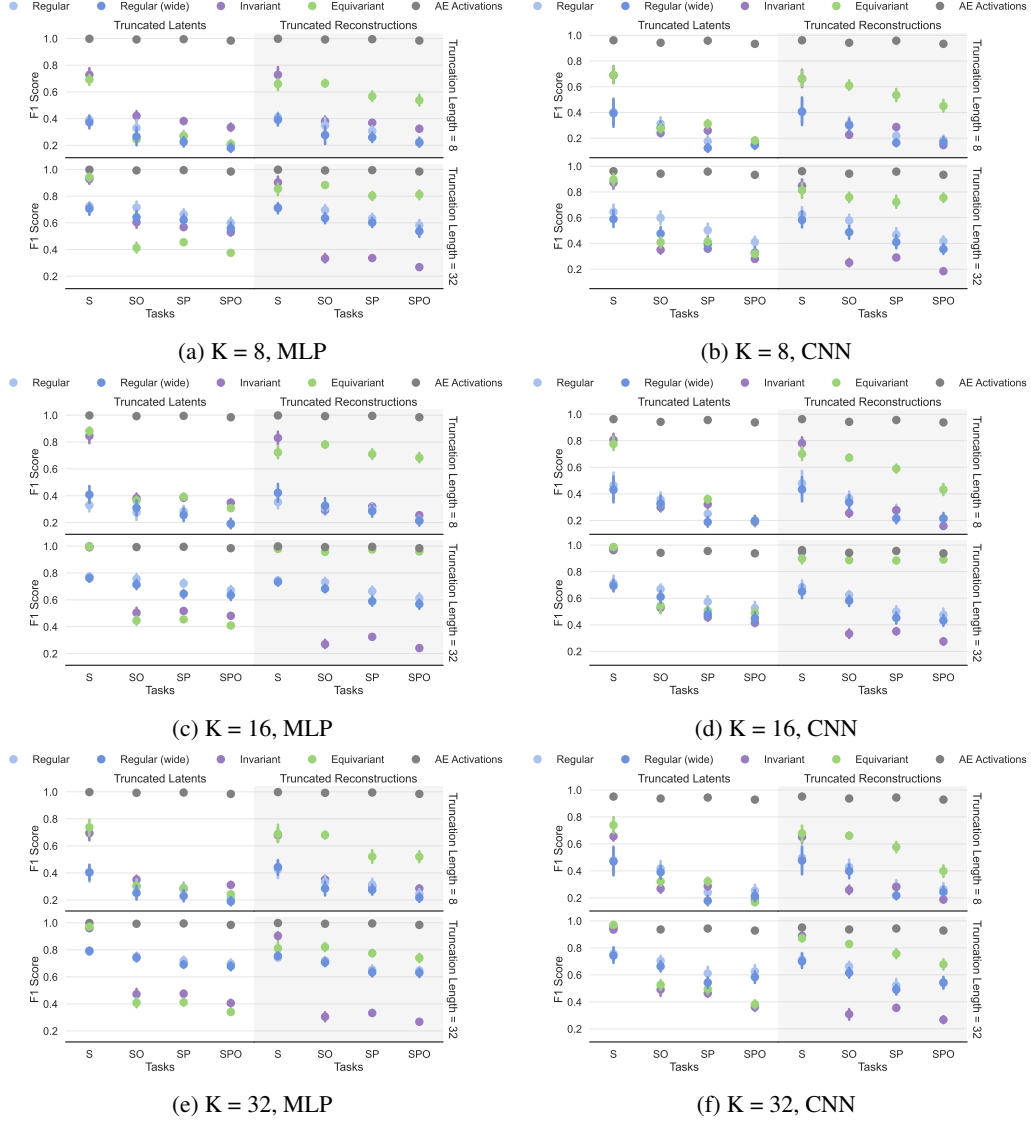


Figure 7: **Further probing results with different TopK values.** Although increasing TopK and the truncation length increases the performance of the probes trained over the regular SAEs activations and reconstructions, results follow a similar trend with those in Figure 4, with the equivariant SAE leading to the best overall probes over its reconstructions.