

---

# fmxcoders: Factorized Masked Crosscoders for Cross-Layer Feature Discovery

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Many features in pretrained Transformers span multiple layers: they emerge  
2 through stages of inference, persist in the residual stream, or are built jointly by  
3 parallel MLPs. Crosscoders (namely, sparse dictionaries trained jointly across  
4 layers) aim to recover these cross-layer features in a single shared latent space. We  
5 show that standard crosscoders largely fail at this purpose. Although their decoder  
6 weight norms spread evenly across layers, a functional coherence metric we in-  
7 troduce reveals that each latent’s activation is effectively driven by only one or two  
8 layers on average. While functionally coherent latents act as human-interpretable  
9 concept detectors (e.g., US states and cities), the layer-localized latents that  
10 crosscoders predominantly learn collapse onto surface-level patterns such as  
11 digit detectors. We trace this failure to two structural limitations: unconstrained  
12 cross-layer parameterization and unregularized cross-layer dependence. We  
13 address both by introducing `fmxcoders`, which (i) replace the encoder and  
14 decoder with low-rank tensor factorizations that draw every latent’s per-layer  
15 weights from a shared cross-layer basis, and (ii) apply stochastic layer masking, a  
16 denoising regularizer along the layer axis that penalizes latents whose contribution  
17 collapses when a single layer is masked. Across GPT2-Small, Pythia-410M,  
18 Pythia-1.4B, and Gemma2-2B, `fmxcoders` lift mean probing F1 by 10–30 points,  
19 surpassing per-layer SAE baselines that standard crosscoders fail to reach, reduce  
20 reconstruction MSE by 25–50%, and roughly double mean functional coherence.  
21 An LLM-as-a-judge evaluation further shows that `fmxcoders` recover 3–13× more  
22 semantically coherent latents than standard crosscoders across all four base LLMs.

## 23 1 Introduction

24 The principal aim of mechanistic interpretability is the decomposition of network activations into  
25 human-interpretable features and the identification of circuits that connect these features [1, 2].  
26 Focusing on the former, the dominant method for feature detection is sparse dictionary learning,  
27 which trains an encoder to map activations into a sparse code and a decoder to reconstruct them [3–9].  
28 Each coordinate of the sparse code is referred to as a *latent*, and corresponds to a learned dictionary  
29 direction whose activation is interpreted as the presence of a feature in the input.

30 A deep network gives rise to two kinds of features that dictionary learning aims to recover: *local*  
31 features, tied to a specific layer or a narrow band of adjacent layers, and *global* features, built jointly  
32 by several layers or persisting across many. Both are well documented in the literature. On the  
33 local side, knowledge neurons store specific factual associations at identifiable layers [10]. At a  
34 coarser scale, transformer inference appears to unfold in universal stages (detokenization, feature  
35 engineering, prediction ensembling, and residual sharpening) each tied to a range of layers [11]. On  
36 the global side, the residual-stream view holds that transformer layers incrementally update a shared  
37 hidden state across the network [12]. Universal neurons recur with similar activation patterns across  
38 different LLMs [13]. Per-layer sparse autoencoder (SAE) analyses reveal features that activate with

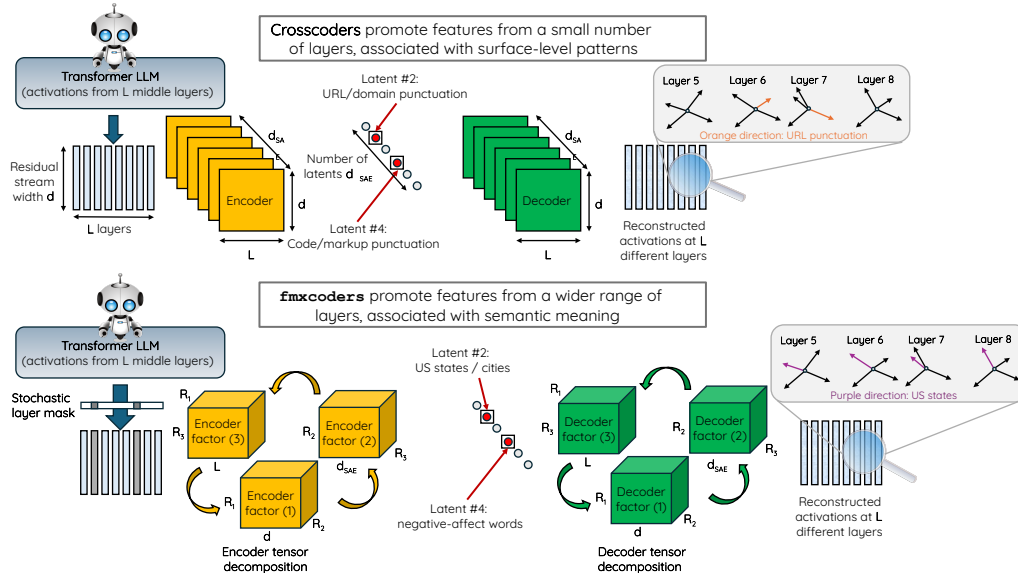


Figure 1: Standard crosscoders (**top**) parameterize the encoder and decoder as  $L$  independent per-layer matrices and recover latents whose functional dependence collapses onto a few layers, capturing surface-level patterns. **fmxcoders** (**bottom**) replace these with low-rank tensor factorizations that share a cross-layer basis across all latents, and add stochastic layer masking as a denoising regularizer along the layer axis, producing latents that span more layers and act as semantic concept detectors. Insets show per-layer decoder directions for one example latent.

39 near-identical directions across many consecutive layers [14, 15], as well as features built by multiple  
 40 MLPs acting in parallel on the residual stream, forming *cross-layer superposition* [16].

41 SAEs are able to properly recover local features [3, 4] by tying a dictionary to a single layer.  
 42 Crosscoders [16, 17] extend this to attain global features by learning a single dictionary whose latent  
 43 neurons read from and write to multiple layers at once. Although crosscoder weight norms spread  
 44 roughly evenly across layers [16, 18], we show that the functional dependence of each latent in  
 45 fact collapses onto one or two on average. To measure this gap, we introduce a *coherence* function,  
 46 defined as the ratio of the sum of per-layer dependence to its maximum per-layer dependence (see  
 47 Section 4.2), bounded between 1 (all dependence on one layer) and  $L$  (equal across all  $L$  layers).  
 48 Two choices of defining this per-layer dependence give two metrics: *norm coherence*, which uses  
 49 the decoder weight norm at each layer (a generalization of typical crosscoder diagnostics [16, 18]),  
 50 and *functional coherence*, which zeros all activations of a specific layer at the crosscoder’s input and  
 51 measures the resulting change in the latent’s activation. We validate that, as already known in the  
 52 literature [16, 18], norm coherence is high ( $> 7$ ), but this is misleading since functional coherence is  
 53 low ( $\sim 2$ ) and thus crosscoder latents model information from just two layers on average.

54 This divergence reflects a structural failure of the standard crosscoder to learn cross-layer features  
 55 which we trace to two concrete limitations:

56 **Limitation 1 (L1): Unconstrained cross-layer structure.** The encoder and decoder are parame-  
 57 terized as  $L$  independent per-layer matrices with no structural link across layers or across the latent  
 58 population. The reconstruction loss (Equation (2)) decomposes additively across layers and, as a  
 59 result, the objective is invariant to reallocations of a latent’s weight mass that preserve per-layer  
 60 reconstructions.

61 **Limitation 2 (L2): Unregularized cross-layer dependence.** Even with a constrained structure, the  
 62 reconstruction objective places no penalty on how a latent’s activation depends on the different layers,  
 63 so a latent can still collapse its functional dependence onto a single dominant layer at no cost.

64 We address these limitations by introducing **fmxcoders**, with two complementary modifications.  
 65 First, **fmxcoders** constrain the cross-layer structure by stacking the per-layer decoder/encoder  
 66 matrices into a 3D tensor and factorizing it [19], directly addressing **L1**. The layer axis of the  
 67 input passes through a shared low-dimensional space rather than  $L$  independent matrices, so each

68 latent’s engagement across layers is drawn from a structured family shared across all latents. Second,  
69 `fmxcoders` regularize the cross-layer dependence with *stochastic layer masking*, directly addressing  
70 **L2**. During training, the activations from specific layers are randomly masked with a probability  $p$ ,  
71 encouraging redundant codes that distribute features across channels, similar to coding theory [20, 21].  
72 In this setting, latents supported by a single layer pay a reconstruction cost growing with  $p$ , while  
73 latents distributed across layers do not, so the optimization pressure favors cross-layer support.

74 Our contributions are summarized below:

75 **C1**. By introducing a functional coherence metric we show that **standard crosscoders concentrate**  
76 **functional dependence on two layers on average**, despite spreading their weight norms uniformly.  
77 This pathology is invisible to the norm-based diagnostics used in prior works.

78 **C2**. We identify two architectural limitations as the sources of the pathology described in **C1**,  
79 **unconstrained cross-layer structure (L1)** and **unregularized cross-layer dependence (L2)**.

80 **C3**. To address these limitations we propose `fmxcoders`, a crosscoder variant combining two  
81 complementary modifications: (i) **low-rank tensor decomposition** of the encoder and decoder  
82 weights, which couple per-layer parameters through a shared cross-layer basis, and (ii) **stochastic**  
83 **layer masking**, which acts as a regularizer along the layer axis and pressures latents to spread their  
84 support across multiple layers.

85 **C4**. Across GPT2-Small, Pythia-410M, Pythia-1.4B, and Gemma2-2B, the combined method **raises**  
86 **probing F1** above the per-layer SAE baseline that the standard crosscoder fails to reach by large  
87 margins, **reduces reconstruction MSE** by 25–50%, and roughly **doubles average functional**  
88 **coherence**.

89 **C5**. A qualitative inspection together with an LLM-as-a-judge evaluation shows that `fmxcoders`  
90 **recover 3–13× more semantically coherent latents** than the standard crosscoders across all LLMs,  
91 with high-coherence latents acting as concept detectors (US states, negative-affect words, etc.) and  
92 low-coherence latents collapsing onto surface-level patterns (digits, punctuation, etc.).

## 93 2 Related Work

94 **Sparse dictionary learning for interpretability.** Sparse autoencoders have become the standard  
95 tool for decomposing transformer activations into monosemantic features [3–5], with variants such  
96 as TopK [6], JumpReLU [7], and BatchTopK [8] refining the sparsification operation. SAEs learn a  
97 feature dictionary at a single location in the transformer, and training multiple SAEs across different  
98 transformer layers forms the basis of circuit discovery [22–24].

99 **Multi-layer and cross-layer methods.** Several approaches extend dictionary learning across layers.  
100 Multi-layer SAEs [14] and feature-matching procedures [15] align per-layer dictionaries post hoc,  
101 while crosscoders [16] train a single shared dictionary over all layers and provide the basis for our  
102 work. They have been used to study feature evolution within a single LLM [17, 25, 26], and to track  
103 features across base-finetune LLM pairs or training checkpoints [27–29]. The crosscoder encoders  
104 and decoders remain unconstrained, which [18] flagged in the model-axis setting as producing  
105 artifactual latents. We identify the layer-axis analogue of this failure, show it persists despite high  
106 dictionary norms across the layers, and address it by applying tensor factorization and layer masking.

107 **Tensor factorization in deep networks.** Tensor decompositions including Canonical Polyadic [30]  
108 (CP) and Tensor Ring [31] (TR) have a long history of compressing weight tensors and inducing  
109 structured priors in deep models [19]. In interpretability research, tensor decompositions have been  
110 applied on the mixture-of-experts paradigm to provide scalable expert specialization [32] and faithful  
111 approximation to multilayer perceptron (MLP) components [33], as well as to model quadratic and  
112 cubic feature interactions in polynomial SAEs [9]. In the present study, tensor decompositions are  
113 repurposed to provide a shared representation space across the latent population and across layers,  
114 addressing the structural limitation identified in standard crosscoders.

115 **Denosing and structured masking.** Denoising autoencoders [34, 35] corrupt inputs and require  
116 reconstruction of clean targets to learn robust and distributed representations; an argument that was  
117 formalized via coding theory [20]. The result is redundancy in information encoding so that masking  
118 some still leaves recoverable signal. The principle has scaled to masked language modeling [36],  
119 masked image modeling [37], and feature-level dropout [21], all of which corrupt at the input or

120 feature level. In our work, we apply the same principles along the layer axis where the corruption  
 121 pattern correlates with the failure mode being targeted, namely layer-specialized latents.

### 122 3 Factorized Masked Crosscoders

#### 123 3.1 Preliminaries

124 **Notation.** Vectors and matrices are denoted by lowercase and uppercase bold letters respectively,  $\mathbf{u}$ ,  
 125  $\mathbf{U}$  and tensors are represented by uppercase bold upright letters  $\mathbf{U}$ . An element (scalar) of a matrix or  
 126 tensor is accessed by subscript, e.g.  $\mathbf{U}_{i,j,k}$ . Mode- $n$  fibers (the generalization of rows and columns  
 127 from matrices to tensors) are obtained by fixing all indices except the one that corresponds to the  $n^{\text{th}}$   
 128 mode, e.g. a mode-1 fiber of  $\mathbf{U}$  is  $\mathbf{U}_{:,j,k}$ . Similarly, tensor slices are obtained by fixing a single index,  
 129 e.g.  $\mathbf{U}_{i,:,:}$ . The  $i$ -th column of  $\mathbf{M}$  is  $\mathbf{m}_i$ , and  $\mathbf{M}_{:,1:r}$  denotes its first  $r$  columns. For depth stacking,  
 130 i.e. combining matrices as slices of a tensor along a new dimension, we use  $[\mathbf{X}, \mathbf{Y}] \in \mathbb{R}^{d_1 \times d_2 \times 2}$  for  
 131  $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{d_1 \times d_2}$ . The symbol  $\circ$  denotes the vector outer product, so  $\mathbf{a} \circ \mathbf{b} \circ \mathbf{c}$  is the rank-one tensor  
 132 with entries  $a_i b_j c_k$ .  $\mathbb{R}^d$  and  $\mathbb{R}^{d_{\text{sae}}}$  denote the activation and sparse-code spaces, with  $d_{\text{sae}} \gg d$ .

133 **Sparse Autoencoders.** SAEs build on overcomplete dictionary learning [38, 6] to decompose neural  
 134 activations into a sparse set of latent features. Given activations  $\mathbf{x} \in \mathbb{R}^d$  from an intermediate  
 135 layer of a pretrained network, an SAE learns a sparse code  $\mathbf{z} \in \mathbb{R}^{d_{\text{sae}}}$  via  $\hat{\mathbf{x}} = \mathbf{b}^{\text{dec}} + \mathbf{D}\mathbf{z}$ ,  $\mathbf{z} =$   
 136  $\mathcal{S}(\text{ReLU}(\mathbf{E}^\top \mathbf{x} + \mathbf{b}^{\text{enc}}))$ , where  $\mathbf{E} \in \mathbb{R}^{d \times d_{\text{sae}}}$  is a linear encoder,  $\mathbf{D} \in \mathbb{R}^{d \times d_{\text{sae}}}$  is the decoder  
 137 dictionary,  $\mathbf{b}^{\text{enc}} \in \mathbb{R}^{d_{\text{sae}}}$  and  $\mathbf{b}^{\text{dec}} \in \mathbb{R}^d$  are the encoder and decoder bias terms,  $\hat{\mathbf{x}} \in \mathbb{R}^d$  are the  
 138 reconstructed activations, and  $\mathcal{S}$  is a sparsification operator such as TopK [6] or BatchTopK [8].  
 139 Feature  $i$  activation depends on the extent to which  $\mathbf{x}$  aligns with the encoder direction  $\mathbf{e}_i$  and, in the  
 140 linear regime, contributes by  $z_i \mathbf{d}_i$  to the reconstruction.

141 **Crosscoders.** A crosscoder [16] applies sparse dictionary learning to jointly encode activations from  
 142 *multiple* layers into a shared sparse latent space and use it to reconstruct all layers. For a given input  
 143 token, let  $\mathbf{x}_\ell \in \mathbb{R}^d$  denote the activation at layer  $\ell \in [1, L]$ , where  $L$  is the number of transformer  
 144 layers considered. The crosscoder assigns encoders  $\mathbf{E}_\ell \in \mathbb{R}^{d \times d_{\text{sae}}}$  and decoders  $\mathbf{D}_\ell \in \mathbb{R}^{d \times d_{\text{sae}}}$  to  
 145 each layer to reconstruct the activations of each layer by computing,

$$146 \quad \mathbf{z} = \mathcal{S} \left( \text{ReLU} \left( \sum_{\ell \in [1, L]} \mathbf{E}_\ell^\top \mathbf{x}_\ell + \mathbf{b}^{\text{enc}} \right) \right), \quad \hat{\mathbf{x}}_\ell = \mathbf{D}_\ell \mathbf{z} + \mathbf{b}_\ell^{\text{dec}}, \quad (1)$$

146 where  $\mathbf{b}^{\text{enc}} \in \mathbb{R}^{d_{\text{sae}}}$  is a pre-activation bias shared across layers,  $\mathbf{b}_\ell^{\text{dec}} \in \mathbb{R}^d$  is a per-layer decoder  
 147 bias,  $\mathcal{S}$  is a sparsification operator (e.g. TopK or BatchTopK), and  $\hat{\mathbf{x}}_\ell$  is the reconstruction at layer  $\ell$   
 148 produced by the encoder–decoder map. The crosscoder is trained to minimize the expected mean  
 149 squared reconstruction error across all layers,

$$150 \quad \mathcal{L}_{\text{recon}} = \frac{1}{L} \sum_{\ell \in [1, L]} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[ \|\mathbf{x}_\ell - \hat{\mathbf{x}}_\ell(\mathbf{x})\|_2^2 \right], \quad (2)$$

150 where  $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_L)$  is the tuple of per-layer activations for a token, and  $\mathcal{D}$  is the joint distribution  
 151 over such tuples induced by passing tokens from the training corpus through the base LLM. In practice  
 152 the expectation is approximated by the empirical mean over a minibatch of tokens.

#### 153 3.2 Tensor Factorization of Crosscoder Weights

154 Depth-stacking the per-layer encoders expresses the total encoder as a tensor  $\mathbf{E} = [\mathbf{E}_1, \dots, \mathbf{E}_L] \in$   
 155  $\mathbb{R}^{d \times d_{\text{sae}} \times L}$ , and the decoder as  $\mathbf{D} = [\mathbf{D}_1, \dots, \mathbf{D}_L] \in \mathbb{R}^{d \times d_{\text{sae}} \times L}$ . In this setting, the mode-1 fiber  
 156  $\mathbf{E}_{:,i,\ell} \in \mathbb{R}^d$  is the encoder direction for feature  $i$  at layer  $\ell$ , and  $\mathbf{D}_{:,i,\ell}$  its decoder atom.

157 This tensor view makes the standard crosscoder’s structural assumption explicit: the  $L$  different  
 158 frontal slices  $\mathbf{E}_{:,i,\ell}$  are parameterized independently, with no mechanism coupling a feature’s behavior  
 159 across layers. This independence is precisely the unconstrained cross-layer structure flagged as **L1**  
 160 in Section 1, and it manifests in two ways. First, *representational fragmentation*: together with  
 161 a reconstruction loss (Equation (2)) that scores each layer separately, nothing pushes latent  $i$  to  
 162 represent the same feature at layer  $\ell$  and layer  $\ell'$ , so the shared latent space does not yield global

163 features but only places layer-local features in a common coordinate system. Second, *parameter*  
 164 *redundancy*: at  $\mathcal{O}(d_{\text{sae}} \cdot d \cdot L)$  parameters, the crosscoder costs as much as  $L$  stacked SAEs and pays  
 165 to re-discover cross-layer features at every layer.

166 We observe that both issues admit a shared fix. Inducing structure across layers resolves the fragmen-  
 167 tation, and constraining how the three modes interact eliminates the redundancy. This choice has both  
 168 empirical and theoretical precedent. Empirically, structural intervention has been shown to resolve a  
 169 different crosscoder pathology, artifactual latents along the model axis in model-diffing settings [18]  
 170 between a base and a fine-tuned model. While the underlying inductive bias of low-rank coupling is  
 171 shared, the layer-axis failure we identify is distinct: it concerns functional dependence collapsing  
 172 within a single model with many layers rather than spurious latents arising between two models, and  
 173 it requires complementary regularization (Section 3.3) that does not apply in the model-axis case.  
 174 Theoretically, low-rank tensor factorizations are a standard inductive bias for multi-task settings with  
 175 shared latent structure across tasks [39, 40]: independent per-task parameterization fails to exploit  
 176 this structure even when it is present, while factorizing along the task axis biases learning toward  
 177 recovering the shared representation. In our setting, each layer must reconstruct its own activations  
 178 from a representation shared with the other layers, making this inductive bias directly applicable.

179 Concretely, we replace  $\mathbf{E}$  and  $\mathbf{D}$  with low-rank tensor decompositions that jointly factorize all three  
 180 modes, forcing the input through a low-dimensional bottleneck. The main results in this study are  
 181 presented using the TR decomposition because it is suited for tensors with large differences between  
 182 dimensions. We also report results using the CP decomposition in Appendix B.

183 **Tensor Ring Decomposition.** The Tensor Ring (TR) decomposition [31] represents each element of  
 184 the tensor as the trace of a product of matrix slices drawn from three factor tensors:

$$\mathbf{E}_{j,i,\ell} = \text{Tr}(\mathbf{G}_{:,j,:}^{(1)} \mathbf{G}_{:,i,:}^{(2)} \mathbf{G}_{:, \ell, :}^{(3)}), \quad (3)$$

185 where  $\mathbf{G}^{(1)} \in \mathbb{R}^{R_1 \times d \times R_2}$ ,  $\mathbf{G}^{(2)} \in \mathbb{R}^{R_2 \times d_{\text{sae}} \times R_3}$ ,  $\mathbf{G}^{(3)} \in \mathbb{R}^{R_3 \times L \times R_1}$ .

186 TR couples the three modes through three independent dimensions, the decomposition ranks  
 187  $(R_1, R_2, R_3)$ , one per mode boundary in the ring. Each mode contributes a matrix and these matrices  
 188 compose multiplicatively: the contribution of layer  $\ell$  to  $\mathbf{E}_{j,i,\ell}$  is also affected by the “feature” slice  
 189  $\mathbf{G}_{:,i,:}^{(2)} \in \mathbb{R}^{R_2 \times R_3}$ , so the same layer pattern can act differently depending on which feature it repre-  
 190 sents. The inductive bias is therefore *context-dependent coupling*: features are no longer confined to  
 191 a single shared basis of layer patterns but participate in subspaces that vary with feature identity, and  
 192 decomposition ranks  $(R_1, R_2, R_3)$  control expressiveness separately at each mode boundary. Tensor  
 193 factorization restricts the cross-layer parameterizations of  $\mathbf{E}$  and  $\mathbf{D}$  (addressing **L1**) without biasing  
 194 the optimizer toward solutions that spread features across layers (**L2**).

### 195 3.3 Stochastic Layer Masking

196 Within the constrained encoder/decoder parametrizations, a latent can still concentrate all its functional  
 197 dependence at a single layer, as nothing in the reconstruction objective of Equation (2) penalizes such  
 198 a configuration. To bias learning toward the cross-layer regime, we complement the architectural  
 199 constraint with an input-level perturbation, *stochastic layer masking*.

200 For each input token and each layer  $\ell$ , we independently draw a Bernoulli mask  $m_\ell \in \{0, 1\}$   
 201 with  $\text{Pr}(m_\ell = 0) = p$ , and the encoder receives the corrupted activations  $\tilde{\mathbf{x}}_\ell = m_\ell \cdot \mathbf{x}_\ell$ . The  
 202 reconstruction loss targets the unmasked activations  $\mathbf{x}_\ell$ , meaning that the crosscoder is trained to  
 203 reconstruct every layer, including the masked ones. This places our objective in the denoising  
 204 autoencoder framework [34, 35] (corrupting inputs and require reconstruction of clean targets),  
 205 instantiated with layer-structured masking noise, a setting in which input noise acts as an implicit  
 206 Tikhonov regularizer on the encoder-decoder mapping [41]. The masking probability  $p$  directly  
 207 controls how strongly the regularizer suppresses localized features. Larger  $p$  makes the cross-layer  
 208 pressure dominate while smaller  $p$  lets localized features survive.

209 A complementary perspective from coding theory shows that information traversing a sequence of  
 210 channels retains capacity only when the code is redundant, and that networks trained under stochastic  
 211 neuron dropout learn precisely such codes, distributing each feature across multiple units rather  
 212 than localizing it [20]. Stochastic layer masking applies the same principle along the layer axis: a  
 213 single-layer latent pays a reconstruction cost that grows with  $p$ , mirroring the redundancy induced by

214 unit-level dropout [21]. This contrasts with the windowed “convolutional” crosscoder of [16], which  
 215 architecturally restricts each latent to  $K < L$  contiguous layers, imposing a hard locality constraint,  
 216 opposite in effect to our stochastic regularizer.

## 217 4 Experiments

### 218 4.1 Experimental Setup

219 The training procedure for `fmxcoders` was implemented in `SAELens` [42], with all implementation  
 220 details provided in Appendix A. The results presented in this study refer to crosscoders with 16,384  
 221 latents trained on four pretrained base LLMs of different scales: GPT2-Small [43], Pythia-410M  
 222 and Pythia-1.4B [44], and Gemma2-2B [45]. Activations were collected from the residual stream at  
 223 eight post-MLP locations spanning the middle layers of each base LLM. All crosscoders were trained  
 224 using a ReLU pre-activation, followed by a BatchTopK [8] sparsifier with  $K = 64$ . Gemma2-2B  
 225 and GPT2-Small activations were generated using OpenWebText [46], while the activations from the  
 226 Pythia models were generated using an uncopyrighted variant of the deduplicated Pile [47]. Training  
 227 runs used 500M tokens (300M for GPT2-Small) with a context length of 128. The evaluations  
 228 that follow compare `fmxcoders` with and without layer masking against base SAEs and standard  
 229 crosscoders trained with the same hyperparameters. In all main-text experiments, `fmxcoders` are  
 230 matched to the standard crosscoder’s parameter count by appropriate choice of decomposition ranks.  
 231 Rank and masking probability ablations are presented in Appendix C.

### 232 4.2 Coherence Diagnostics

233 To assess whether a crosscoder latent neuron represents a genuinely multi-layer feature or merely  
 234 responds to information from a single dominant layer, we introduce two complementary coherence  
 235 metrics: *norm coherence* and *functional coherence*.

236 **Norm Coherence.** The simplest approach to measuring layer distribution is to examine the cross-  
 237 coder’s weights directly, which is a generalized version of the *relative decoder norms* shown in [16],  
 238 suitable for any number of layers. For a crosscoder with decoder weights  $\mathbf{D} \in \mathbb{R}^{d \times d_{\text{sac}} \times L}$ , we compute  
 239 the L2 norm for each latent’s dictionary entry corresponding to each layer  $\ell \in [1, L]$  via the decoder’s  
 240 mode-1 fibers  $\mathbf{N}_{i,\ell} = \|\mathbf{D}_{:,i,\ell}\|_2$ . The *norm coherence* of latent  $i$  is then

$$c_i^n = \frac{\sum_{\ell \in [1,L]} \mathbf{N}_{i,\ell}}{\max_{\ell \in [1,L]} \mathbf{N}_{i,\ell}}, \quad (4)$$

241 which ranges from 1 (all weight mass concentrated in one layer) to  $L$  (equal norms across all layers).

242 **Functional Coherence.** A more informative measure examines functional dependence, focusing  
 243 on how a latent’s activation changes when information from a specific layer is removed. For an  
 244 input  $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_L) \sim \mathcal{D}$ , let  $M_\ell$  denote the operator that zeros the  $\ell$ -th component of  $\mathbf{x}$ ,  
 245  $M_\ell(\mathbf{x}) = (\mathbf{x}_1, \dots, \mathbf{x}_{\ell-1}, \mathbf{0}, \mathbf{x}_{\ell+1}, \dots, \mathbf{x}_L)$ , so that  $z_i(M_\ell(\mathbf{x}))$  is the latent activation when layer  $\ell$   
 246 is masked while all other layers remain intact. We define latent  $i$ ’s functional coherence  $c_i^f$  as

$$c_i^f = \frac{\sum_{\ell \in [1,L]} \mathbf{S}_i^\ell}{\max_{\ell \in [1,L]} \mathbf{S}_i^\ell}, \quad \text{where } \mathbf{S}_i^\ell = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[ \left| \frac{z_i(\mathbf{x}) - z_i(M_\ell(\mathbf{x}))}{z_i(\mathbf{x}) + \epsilon} \right| \right]. \quad (5)$$

247  $\mathbf{S}_i^\ell$  represents the relative change of latent  $i$ ’s activation with respect to masking layer  $\ell$ , with  
 248  $0 < \epsilon \ll 1$  a small constant for numerical stability. A high value of  $\mathbf{S}_i^\ell$  indicates that latent  $i$  is  
 249 sensitive to information from layer  $\ell$ , while a low value indicates that layer  $\ell$  contributes little to that  
 250 latent’s activation.  $c_i^f$  therefore, measures the ratio of the total relative change across layers to the  
 251 largest single-layer relative change. Analogously to norm coherence,  $c_i^f$  ranges from 1 to  $L$ .

252 Figure 2 shows the latent distributions of norm and functional coherence for crosscoder and `fmxcoder`  
 253 variants trained on Pythia-410M. Norm coherence remains high (avg.  $c^n > 7$ ) in all cases, while  
 254 functional coherence is much lower (avg.  $c^f \approx 2$ ), except for `fmxcoder` with  $p=0.05$ , where it  
 255 doubles. Factorization alone therefore does not increase functional coherence; *the gain emerges only*  
 256 *when factorization and masking are combined*. This interaction reflects the complementary roles

Table 1: Examples of high functional coherence  $c^f$  latents (top half of the table) and low  $c^f$  latents (bottom half of the table). The corresponding norm coherence  $c^n$  is also reported. The results correspond to an `fmxcoder` with  $p=0.05$  trained on the activations of Gemma2-2B.

| Latent firing pattern   | Description                | $c^f$ | $c^n$ |
|---|----------------------------|-------|-------|
| [ <i>lobbying, campaigned, pleaded, pitch, ...</i> ]          | verbs of advocacy          | 7.1   | 6.1   |
| [ <i>regret, despise, failure, gruesome, ...</i> ]            | negative-affect words      | 7.0   | 6.9   |
| [ <i>Philadelphia, Pennsylvania, Maryland, Florida, ...</i> ] | US states and cities       | 6.9   | 6.7   |
| [ <i>colleagues, friends, accomplice, partner, ...</i> ]      | companion relationship     | 6.8   | 6.3   |
| [ <i>vampire, werewolf, supernatural, superpowers, ...</i> ]  | supernatural beings/powers | 6.7   | 7.4   |
| [1, 2, 3, 4, ...]   | digit detector             | 1.0   | 5.8   |
| [ <i>def, set, define, -&gt;, ...</i> ]                       | code/markup punctuation    | 1.0   | 6.2   |
| [ <i>d#, 3, :#, bbbb, ...</i> ]                               | hex/colour-code characters | 1.0   | 6.0   |
| [ <i>G, TT, BB, QL, ...</i> ]                                 | short sub-word fragments   | 1.0   | 6.5   |
| [ <i>], ;, ), ...</i> ]                                       | closing-brace family       | 1.0   | 5.5   |

257 of the two interventions. In a standard (unfactorized) crosscoder, the per-layer dictionary entries  
 258 (encoder/decoder mode-1 fibers) of each latent are independent parameters. Layer masking only  
 259 adjusts the weights of layers a feature already occupies, with no mechanism to extend its support  
 260 to new ones, despite the correlations between residual-stream activations across layers. Tensor  
 261 factorization removes this barrier by tying the per-layer dictionary entries through shared factors, so  
 262 weight updates propagate jointly across layers and latents. Layer masking, applied on top, supplies  
 263 the optimization pressure that turns this coupling into genuinely cross-layer features.

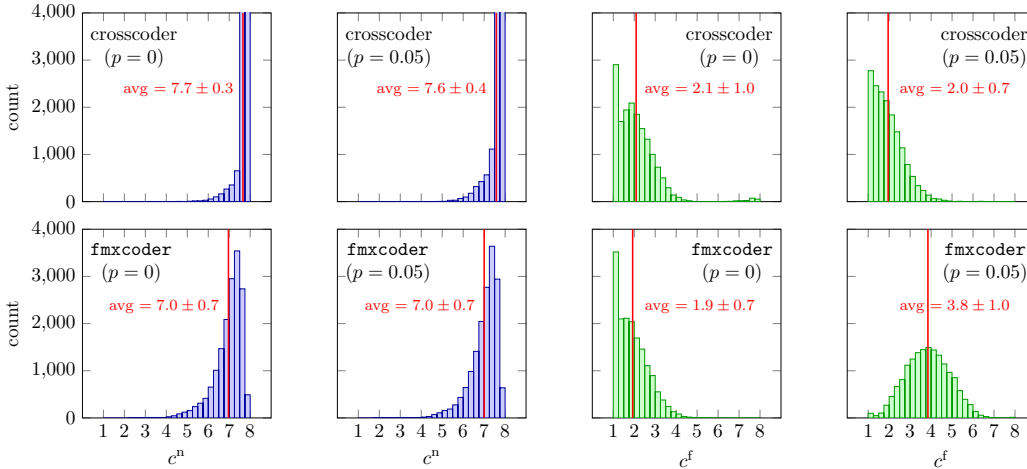


Figure 2: Norm coherence  $c^n$  (blue) and functional coherence  $c^f$  (green) latent distributions for crosscoders (top row) and `fmxcoder` (bottom row) variants trained on Pythia-410M. The average value and standard deviation of each coherence distribution are also shown in red color.

### 264 4.3 Qualitative and LLM-Judge Analysis of Features

265 Table 1 contrasts examples of highest- and lowest- $c^f$  latents recovered by an `fmxcoder` with  $p=0.05$   
 266 trained on the activations of the Gemma2-2B model. *High- $c^f$  latents act as concept detectors* that fire  
 267 on diverse tokens unified by meaning, including morphological variants and named-entity instances.  
 268 *Low- $c^f$  latents instead collapse onto character classes*, single punctuation symbols, or short subword  
 269 fragments with no shared semantics. Across both halves of the table,  $c^n$  remains in the narrow band  
 270 5.5–7.4 while  $c^f$  spans the full 1.0–7.1 range, confirming that weight magnitudes spread evenly across  
 271 layers even when functional dependence does not, and motivating  $c^f$  as a more informative diagnostic.

272 To scale this analysis, we perform an *LLM-as-a-judge* evaluation, asking GPT-4o-mini as the LLM  
 273 judge to score each latent on two independent axes: (i) the latent has discovered a semantic concept  
 274 (fires share a meaning, concept, topic, sentiment, entity, etc.), and (ii) a surface-level pattern

275 (grammatical role, character class, formatting, etc.). The latents are then labeled accordingly as  
 276 “semantic” or “surface”. Full procedure and prompts are provided in Appendix E. The total number of  
 277 semantic and surface latents for all base LLMs is presented in Table 2, comparing `fmxcoder` ( $p=0.05$ )  
 278 results against crosscoder results. Across all four base LLMs, `fmxcoder` *recovers substantially more*  
 279 *semantic latents* than the standard crosscoder, lifting the absolute count from at most a few hundred  
 280 to between 1.2k and 3.3k. The effect is most pronounced on Gemma2-2B, the only setting in which  
 281 semantic latents approach parity with surface ones. Surface-latent counts also rise on three of the  
 282 four base LLMs and fall only on Gemma2-2B, so the recovered dictionary shifts toward semantic  
 283 latents in composition, while the total number of clearly categorized latents increases overall.

Table 2: Number of semantic vs. surface latent counts per base LLM, using an LLM-as-a-judge procedure. Full procedure and prompts are provided in Appendix E.

| Method                               | GPT2-Small |      | Pythia-410M |      | Pythia-1.4B |      | Gemma2-2B |      |
|--------------------------------------|------------|------|-------------|------|-------------|------|-----------|------|
|                                      | Sem.       | Sur. | Sem.        | Sur. | Sem.        | Sur. | Sem.      | Sur. |
| Crosscoder                           | 160        | 2948 | 128         | 3184 | 399         | 8693 | 405       | 6624 |
| <code>fmxcoder</code> ( $p = 0.05$ ) | 2045       | 5934 | 1324        | 9074 | 1195        | 8925 | 3319      | 4153 |

#### 284 4.4 Reconstruction and Sparse Probing Results

285 For the evaluation of the trained models we report reconstruction metrics on held-out data following  
 286 the SAEBench library [48]. Table 3 reports the mean squared error (MSE), explained variance  
 287 (EV), and cosine similarity (CS) across LLMs. We report all three metrics because they capture  
 288 complementary aspects of reconstruction quality: MSE gives absolute error in native units, EV  
 289 normalizes the error to enable cross-LLM comparison, and CS isolates directional fidelity, which  
 290 downstream layers depend on. The standard crosscoder’s high EV (0.986) but low CS (0.870) on  
 291 Pythia-410M illustrate why all three metrics are needed. Across all four base LLMs, the standard  
 292 crosscoder is the weakest reconstructor on every metric, with MSE up to roughly 100% higher than  
 293 the per-layer SAE baseline and consistent drops in CS. The `fmxcoder` variants reverse this pattern,  
 294 reducing MSE by 25–50% relative to the standard crosscoder, recovering CS that, in some cases,  
 295 exceeds the SAE baseline, and matching the SAE explained variance. The masking variant pays  
 296 a small reconstruction cost relative to the unmasked variant, which is the expected trade-off of a  
 297 denoising objective and is offset by the functional coherence gains reported in Section 4.2.

Table 3: Reconstruction metrics and probing-based results for different base LLMs, for SAE, cross-  
 coder, and `fmxcoders` variants. We report average F1 across six classification tasks and the detailed  
 results are presented in Appendix D. SAE results were reproduced using the settings from [9].

| LLM         | Variant                              | MSE ↓       | EV ↑          | CS ↑          | Avg. F1 ↑   | Avg. Wass. ↑<br>( $\times 10^{-3}$ ) |
|-------------|--------------------------------------|-------------|---------------|---------------|-------------|--------------------------------------|
| GPT2-Small  | SAE                                  | 0.53        | 0.9256        | 0.9623        | 65.7        | 8.4                                  |
|             | Crosscoder                           | 0.64        | 0.8704        | 0.9346        | 67.6        | 11.6                                 |
|             | <code>fmxcoder</code> ( $p = 0$ )    | <b>0.32</b> | <b>0.9359</b> | <b>0.9682</b> | <b>78.3</b> | <b>120.2</b>                         |
|             | <code>fmxcoder</code> ( $p = 0.05$ ) | 0.33        | 0.9332        | 0.9670        | 73.5        | 68.4                                 |
| Pythia-410m | SAE                                  | <b>0.03</b> | <b>0.9918</b> | <b>0.9364</b> | 65.0        | 0.5                                  |
|             | Crosscoder                           | 0.05        | 0.9856        | 0.8696        | 47.4        | <0.1                                 |
|             | <code>fmxcoder</code> ( $p = 0$ )    | <b>0.03</b> | 0.9917        | 0.9275        | <b>74.3</b> | <b>10.2</b>                          |
|             | <code>fmxcoder</code> ( $p = 0.05$ ) | <b>0.03</b> | 0.9914        | 0.9257        | 68.1        | 4.5                                  |
| Pythia-1.4b | SAE                                  | <b>0.22</b> | <b>0.9784</b> | <b>0.9289</b> | 64.6        | 0.8                                  |
|             | Crosscoder                           | 0.31        | 0.9691        | 0.8976        | 47.0        | 0.2                                  |
|             | <code>fmxcoder</code> ( $p = 0$ )    | <b>0.22</b> | 0.9780        | 0.9273        | <b>78.4</b> | <b>36.1</b>                          |
|             | <code>fmxcoder</code> ( $p = 0.05$ ) | 0.23        | 0.9774        | 0.9252        | 71.2        | 13.1                                 |
| Gemma2-2b   | SAE                                  | <b>1.58</b> | <b>0.8702</b> | 0.9215        | 64.8        | 2.7                                  |
|             | Crosscoder                           | 3.29        | 0.7809        | 0.8737        | 51.2        | 0.2                                  |
|             | <code>fmxcoder</code> ( $p = 0$ )    | 2.03        | 0.8632        | <b>0.9229</b> | <b>81.9</b> | <b>198.7</b>                         |
|             | <code>fmxcoder</code> ( $p = 0.05$ ) | 2.07        | 0.8608        | 0.9218        | 73.0        | 39.1                                 |

298 Furthermore, Table 3 also reports sparse probing performance on six classification tasks: Bias in  
 299 Bios [49], AG News [50], EuroParl [51], GitHub programming languages [52], Amazon Sentiment,

300 and Amazon 15 [53]. While the F1 score measures whether a latent fires on the correct class, it is  
301 invariant to how strongly the firing distributions for positive and negative examples are separated. We  
302 complement it with the 1-Wasserstein distance between the activation distributions of the two classes,  
303 computed per latent and aggregated across the probing dataset [9]. Wasserstein distance measures the  
304 minimum “transport cost” to convert one distribution into the other, capturing both the gap between  
305 their means and differences in their shapes. A latent with high F1 but low Wasserstein distance  
306 discriminates the classes only marginally, whereas a latent with both high F1 and high Wasserstein  
307 distance separates the classes with a wide margin, indicating a more robust feature. Averaged across  
308 all six classification tasks, the standard crosscoder underperforms the per-layer SAE baseline by  
309 13–18 F1 points on Pythia-410M, Pythia-1.4B, and Gemma2-2B, and shows near-zero Wasserstein  
310 distance across all four base LLMs, indicating latents that barely separate the probing classes. As with  
311 reconstruction metrics, `fmxcoders` reverse this gap entirely: without masking, F1 rises to 74–82%,  
312 exceeding the SAE baseline by 9–17 points on every model, and Wasserstein distance increases  
313 by one to three orders of magnitude, confirming that the gains reflect well-separated activation  
314 distributions rather than thin-margin classifications. The masking variants incur a small F1 cost (4–9  
315 points) and lower Wasserstein distance, but they nonetheless remain well above the SAE and standard  
316 crosscoder baselines on F1 and Wasserstein distance. The individual sparse probing results across all  
317 the classification datasets considered in this study are presented in Appendix D.

## 318 5 Limitations

319 Several aspects of the present study limit the scope of these conclusions. First, our experiments are  
320 restricted to eight post-MLP residual-stream locations spanning the middle layers of each model. The  
321 behavior of `fmxcoders` when applied to the full depth of a network, or jointly to attention-output  
322 and MLP-output streams, remains untested, and the case  $L = 8$  is small relative to the depth of  
323 frontier models. Second, `fmxcoders` introduce two new hyperparameters, the masking probability  $p$   
324 and the factorization ranks, whose optimal settings appear to be model- and dataset-dependent (see  
325 Appendix C). Our sweeps explored only a coarse grid and we do not provide an *a priori* procedure for  
326 selecting  $p$  when the downstream task is unknown. Third, all main-table results use a single sparsifier  
327 (BatchTopK at  $K=64$ ) and a fixed dictionary width of 16,384 latents. How `fmxcoders` interact with  
328 alternative sparsity mechanisms (JumpReLU, L1) and with significantly larger dictionary widths is  
329 left to future work. Fourth, distinguishing semantic from surface latents is itself a subtle judgment,  
330 often hinging on whether a shared grammatical role also carries shared meaning, and affordable  
331 judges such as *GPT-4o-mini* are not always reliable at this distinction. Fifth, `fmxcoders` incur a slight  
332 additional training-time cost relative to standard crosscoders (<10%) due to the tensor-factorization  
333 forward/backward passes, with compute costs reported in Appendix A. Finally, all main-table results  
334 come from a single training run per LLM/`fmxcoder` variant due to compute constraints.

## 335 6 Conclusion

336 We introduced the `fmxcoder` architecture, that uses low-rank tensor decomposition for the en-  
337 coder/decoder tensors and applies stochastic layer masking to act as a remedy to standard crosscoders’  
338 tendency to learn primarily layer-localized features. Together, the two interventions raise mean prob-  
339 ing F1 by 10–30 points across GPT2-Small, Pythia-410M, Pythia-1.4B, and Gemma2-2B, surpassing  
340 per-layer SAE baselines that the standard crosscoder fails to reach, reduce reconstruction MSE by  
341 25–50%, and roughly double the average functional coherence of the recovered latents. A qualitative  
342 inspection further indicates that high-coherence latents act as semantic concept detectors (US states  
343 and cities, negative-affect words, verbs of advocacy), while low-coherence latents collapse onto  
344 surface-level character classes and syntactic tokens. An LLM-as-a-judge evaluation corroborates this  
345 shift quantitatively: `fmxcoders` recover 1.2k–3.3k semantic concept latents per model, roughly an  
346 order of magnitude above the standard crosscoder, with semantic latents approaching parity with  
347 surface ones on Gemma2-2B.

## 348 References

349 [1] N. Elhage, T. Hume, C. Olsson, N. Schiefer, T. Henighan, S. Kravec, Z. Hatfield-Dodds,  
350 R. Lasenby, D. Drain, C. Chen, R. Grosse, S. McCandlish, J. Kaplan, D. Amodei,

- 351 M. Wattenberg, and C. Olah, “Toy models of superposition,” *Transformer Circuits Thread*,  
352 2022. [Online]. Available: [https://transformer-circuits.pub/2022/toy\\_model/index.html](https://transformer-circuits.pub/2022/toy_model/index.html)
- 353 [2] L. Sharkey, B. Chughtai, J. Batson, J. Lindsey, J. Wu, L. Bushnaq, N. Goldowsky-Dill, S. Heimer-  
354 sheim, A. Ortega, J. Bloom, S. Biderman, A. Garriga-Alonso, A. Conmy, N. Nanda, J. Rumbel-  
355 low, M. Wattenberg, N. Schoots, J. Miller, E. J. Michaud, S. Casper, M. Tegmark, W. Saunders,  
356 D. Bau, E. Todd, A. Geiger, M. Geva, J. Hoogland, D. Murfet, and T. McGrath, “Open problems  
357 in mechanistic interpretability,” *arXiv preprint arXiv:2501.16496*, 2025.
- 358 [3] T. Bricken, A. Templeton, J. Batson, B. Chen, A. Jermyn, T. Conerly, N. Turner,  
359 C. Anil, C. Denison, A. Askell, R. Lasenby, Y. Wu, S. Kravec, N. Schiefer, T. Maxwell,  
360 N. Joseph, Z. Hatfield-Dodds, A. Tamkin, K. Nguyen, B. McLean, J. E. Burke, T. Hume,  
361 S. Carter, T. Henighan, and C. Olah, “Towards monosemanticity: Decomposing language  
362 models with dictionary learning,” *Transformer Circuits Thread*, 2023. [Online]. Available:  
363 <https://transformer-circuits.pub/2023/monosemantic-features/index.html>
- 364 [4] H. Cunningham, A. Ewart, L. Riggs, R. Huben, and L. Sharkey, “Sparse autoencoders find  
365 highly interpretable features in language models,” in *International Conference on Learning  
366 Representations (ICLR)*, 2024.
- 367 [5] A. Templeton, T. Conerly, J. Marcus, J. Lindsey, T. Bricken, B. Chen, A. Pearce,  
368 C. Citro, E. Ameisen, A. Jones, H. Cunningham, N. L. Turner, C. McDougall,  
369 M. MacDiarmid, C. D. Freeman, T. R. Sumers, E. Rees, J. Batson, A. Jermyn,  
370 S. Carter, C. Olah, and T. Henighan, “Scaling monosemanticity: Extracting interpretable  
371 features from Claude 3 Sonnet,” *Transformer Circuits Thread*, 2024. [Online]. Available:  
372 <https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html>
- 373 [6] L. Gao, T. Dupré la Tour, H. Tillman, G. Goh, R. Troll, A. Radford, I. Sutskever, J. Leike, and  
374 J. Wu, “Scaling and evaluating sparse autoencoders,” in *International Conference on Learning  
375 Representations (ICLR)*, 2025, oral presentation.
- 376 [7] S. Rajamanoharan, T. Lieberum, N. Sonnerat, A. Conmy, V. Varma, J. Kramár, and N. Nanda,  
377 “Jumping ahead: Improving reconstruction fidelity with JumpReLU sparse autoencoders,” *arXiv  
378 preprint arXiv:2407.14435*, 2024.
- 379 [8] B. Bussmann, P. Leask, and N. Nanda, “BatchTopK sparse autoencoders,” in *NeurIPS 2024  
380 Workshop on Scientific Methods for Understanding Neural Networks*, 2024.
- 381 [9] P. Koromilas, A. D. Demou, J. Oldfield, Y. Panagakis, and M. Nicolaou, “Polysae: Mod-  
382 eling feature interactions in sparse autoencoders via polynomial decoding,” *arXiv preprint  
383 arXiv:2602.01322*, 2026.
- 384 [10] D. Dai, L. Dong, Y. Hao, Z. Sui, B. Chang, and F. Wei, “Knowledge neurons in pretrained  
385 transformers,” in *Proceedings of the 60th Annual Meeting of the Association for Computational  
386 Linguistics (Volume 1: Long Papers)*, 2022, pp. 8493–8502.
- 387 [11] V. Lad, J. H. Lee, W. Gurnee, and M. Tegmark, “Remarkable robustness of LLMs: Stages of  
388 inference?” in *The Thirty-ninth Annual Conference on Neural Information Processing Systems*,  
389 2025. [Online]. Available: <https://openreview.net/forum?id=Wxh5Xz7NpJ>
- 390 [12] M. Geva, A. Caciularu, K. Wang, and Y. Goldberg, “Transformer feed-forward layers build pre-  
391 dictions by promoting concepts in the vocabulary space,” in *Proceedings of the 2022 conference  
392 on empirical methods in natural language processing*, 2022, pp. 30–45.
- 393 [13] W. Gurnee, T. Horsley, Z. C. Guo, T. R. Kheirkhah, Q. Sun, W. Hathaway, N. Nanda, and  
394 D. Bertsimas, “Universal neurons in gpt2 language models,” *arXiv preprint arXiv:2401.12181*,  
395 2024.
- 396 [14] T. Lawson, L. Farnik, C. Houghton, and L. Aitchison, “Residual stream analysis with multi-layer  
397 SAEs,” in *International Conference on Learning Representations (ICLR)*, 2025.
- 398 [15] N. Balagansky, I. Maksimov, and D. Gavrilov, “Mechanistic permutability: Match features  
399 across layers,” in *International Conference on Learning Representations (ICLR)*, 2025.

- 400 [16] J. Lindsey, A. Templeton, J. Marcus, T. Conerly, J. Batson, and C. Olah, “Sparse crosscoders  
401 for cross-layer features and model diffing,” *Transformer Circuits Thread*, 2024. [Online].  
402 Available: <https://transformer-circuits.pub/2024/crosscoders/index.html>
- 403 [17] E. Ameisen, J. Lindsey, A. Pearce, W. Gurnee, N. L. Turner, B. Chen, C. Citro, D. Abrahams,  
404 S. Carter, B. Hosmer, J. Marcus, M. Sklar, A. Templeton, T. Bricken, C. McDougall,  
405 H. Cunningham, T. Henighan, A. Jermyn, A. Jones, A. Persic, Z. Qi, T. B. Thompson,  
406 S. Zimmerman, K. Rivoire, T. Conerly, C. Olah, and J. Batson, “Circuit tracing: Revealing  
407 computational graphs in language models,” *Transformer Circuits Thread*, 2025. [Online].  
408 Available: <https://transformer-circuits.pub/2025/attribution-graphs/methods.html>
- 409 [18] C. Dumas, J. Minder, C. Juang, B. Chughtai, and N. Nanda, “Overcoming sparsity artifacts  
410 in crosscoders to interpret chat-tuning,” in *Mechanistic Interpretability Workshop at NeurIPS*  
411 *2025*, 2025.
- 412 [19] Y. Panagakis, J. Kossaiji, G. G. Chrysos, J. Oldfield, M. A. Nicolaou, A. Anandkumar, and  
413 S. Zafeiriou, “Tensor methods in computer vision and deep learning,” *Proceedings of the IEEE*,  
414 vol. 109, no. 5, pp. 863–890, 2021.
- 415 [20] S. C. Marshall and J. H. Kirchner, “Understanding polysemanticity in neural networks through  
416 coding theory,” *arXiv preprint arXiv:2401.17975*, 2024.
- 417 [21] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, “Dropout: A  
418 simple way to prevent neural networks from overfitting,” *Journal of Machine Learning Research*,  
419 vol. 15, no. 56, pp. 1929–1958, 2014.
- 420 [22] A. Conmy, A. Mavor-Parker, A. Lynch, S. Heimersheim, and A. Garriga-Alonso, “Towards  
421 automated circuit discovery for mechanistic interpretability,” *Advances in Neural Information*  
422 *Processing Systems*, vol. 36, pp. 16 318–16 352, 2023.
- 423 [23] M. Hanna, O. Liu, and A. Variengien, “How does gpt-2 compute greater-than?: Interpreting  
424 mathematical abilities in a pre-trained language model,” *Advances in Neural Information*  
425 *Processing Systems*, vol. 36, pp. 76 033–76 060, 2023.
- 426 [24] S. Marks, C. Rager, E. J. Michaud, Y. Belinkov, D. Bau, and A. Mueller, “Sparse feature  
427 circuits: Discovering and editing interpretable causal graphs in language models,” *arXiv preprint*  
428 *arXiv:2403.19647*, 2024.
- 429 [25] D. Troitskii, K. Pal, C. Wendler, C. S. McDougall, and N. Nanda, “Internal states before wait  
430 modulate reasoning patterns,” *Proceedings of the Findings of the Association for Computational*  
431 *Linguistics: EMNLP*, 2025.
- 432 [26] A. Wadell, A. Bhutani, V. Azumah, A. R. Ellis-Mohr, C. Kelly, H. Zhao, A. K. Nayak, K. Hegazy,  
433 A. Brace, and H. Lin, “Foundation models for discovery and exploration in chemical space,”  
434 *arXiv preprint arXiv:2510.18900*, 2025.
- 435 [27] D. D. Baek and M. Tegmark, “Towards understanding distilled reasoning models: A representa-  
436 tional approach,” *arXiv preprint arXiv:2503.03730*, 2025.
- 437 [28] X. Ge, W. Shu, J. Wu, Y. Zhou, Z. He, and X. Qiu, “Evolution of concepts in language model  
438 pre-training,” in *The Fourteenth International Conference on Learning Representations*, 2026.
- 439 [29] D. Bayazit, A. Mueller, and A. Bosselut, “Crosscoding through time: Tracking emergence  
440 & consolidation of linguistic representations throughout LLM pretraining,” *arXiv preprint*  
441 *arXiv:2509.05291*, 2025.
- 442 [30] F. L. Hitchcock, “The expression of a tensor or a polyadic as a sum of products,” *Journal of*  
443 *Mathematics and Physics*, vol. 6, no. 1-4, pp. 164–189, 1927.
- 444 [31] Q. Zhao, G. Zhou, S. Xie, L. Zhang, and A. Cichocki, “Tensor Ring Decomposition,” *arXiv*  
445 *preprint arXiv:1606.05535*, 2016.

- 446 [32] J. Oldfield, M. Georgopoulos, G. G. Chrysos, C. Tzelepis, Y. Panagakis, M. A. Nicolaou,  
447 J. Deng, and I. Patras, “Multilinear mixture of experts: Scalable expert specialization through  
448 factorization,” *Advances in Neural Information Processing Systems*, vol. 37, pp. 53 022–53 063,  
449 2024.
- 450 [33] J. Oldfield, S. Im, S. Li, M. A. Nicolaou, I. Patras, and G. G. Chrysos, “Towards interpretability  
451 without sacrifice: Faithful dense layer decomposition with mixture of decoders,” *arXiv preprint*  
452 *arXiv:2505.21364*, 2025.
- 453 [34] P. Vincent, H. Larochelle, Y. Bengio, and P.-A. Manzagol, “Extracting and composing robust  
454 features with denoising autoencoders,” in *Proceedings of the 25th International Conference on*  
455 *Machine Learning (ICML)*. ACM, 2008, pp. 1096–1103.
- 456 [35] P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, and P.-A. Manzagol, “Stacked denoising  
457 autoencoders: Learning useful representations in a deep network with a local denoising criterion,”  
458 *Journal of Machine Learning Research*, vol. 11, pp. 3371–3408, 2010.
- 459 [36] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “BERT: Pre-training of deep bidirec-  
460 tional transformers for language understanding,” in *Proceedings of the 2019 Conference of the*  
461 *North American Chapter of the Association for Computational Linguistics: Human Language*  
462 *Technologies (NAACL-HLT)*, 2019, pp. 4171–4186.
- 463 [37] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick, “Masked autoencoders are scalable  
464 vision learners,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern*  
465 *recognition*, 2022, pp. 16 000–16 009.
- 466 [38] S. G. Mallat and Z. Zhang, “Matching pursuits with time-frequency dictionaries,” *IEEE Trans-*  
467 *actions on signal processing*, vol. 41, no. 12, pp. 3397–3415, 1993.
- 468 [39] B. Romera-Paredes, H. Aung, N. Bianchi-Berthouze, and M. Pontil, “Multilinear multitask  
469 learning,” in *International Conference on Machine Learning*. PMLR, 2013, pp. 1444–1452.
- 470 [40] Y. Yang and T. Hospedales, “Deep multi-task representation learning: A tensor factorisation  
471 approach,” *arXiv preprint arXiv:1605.06391*, 2016.
- 472 [41] C. M. Bishop, “Training with noise is equivalent to Tikhonov regularization,” *Neural Computa-*  
473 *tion*, vol. 7, no. 1, pp. 108–116, 1995.
- 474 [42] J. Bloom, C. Tigges, A. Duong, and D. Chanin, “SAELens,” [https://github.com/jbloomAus/](https://github.com/jbloomAus/SAELens)  
475 [SAELens](https://github.com/jbloomAus/SAELens), 2024, gitHub repository.
- 476 [43] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, “Language models are unsu-  
477 pervised multitask learners,” OpenAI, Tech. Rep., 2019. [Online]. Available: [https://cdn.openai.](https://cdn.openai.com/better-language-models/language_models_are_unsupervised_multitask_learners.pdf)  
478 [com/better-language-models/language\\_models\\_are\\_unsupervised\\_multitask\\_learners.pdf](https://cdn.openai.com/better-language-models/language_models_are_unsupervised_multitask_learners.pdf)
- 479 [44] S. Biderman, H. Schoelkopf, Q. G. Anthony, H. Bradley, K. O’Brien, E. Hallahan, M. A.  
480 Khan, S. Purohit, U. S. Prashanth, E. Raff, A. Skowron, L. Sutawika, and O. Van Der Wal,  
481 “Pythia: A suite for analyzing large language models across training and scaling,” in  
482 *Proceedings of the 40th International Conference on Machine Learning*, ser. Proceedings of  
483 Machine Learning Research, vol. 202. PMLR, 2023, pp. 2397–2430. [Online]. Available:  
484 <https://proceedings.mlr.press/v202/biderman23a.html>
- 485 [45] Gemma Team, “Gemma 2: Improving open language models at a practical size,”  
486 Google DeepMind, Tech. Rep., 2024. [Online]. Available: [https://storage.googleapis.com/](https://storage.googleapis.com/deepmind-media/gemma/gemma-2-report.pdf)  
487 [deepmind-media/gemma/gemma-2-report.pdf](https://storage.googleapis.com/deepmind-media/gemma/gemma-2-report.pdf)
- 488 [46] A. Gokaslan, V. Cohen, E. Pavlick, and S. Tellex, “OpenWebText corpus,” Zenodo, 2019.  
489 [Online]. Available: <https://zenodo.org/records/3834942>
- 490 [47] L. Gao, S. Biderman, S. Black, L. Golding, T. Hoppe, C. Foster, J. Phang, H. He, A. Thite,  
491 N. Nabeshima, S. Presser, and C. Leahy, “The Pile: An 800GB dataset of diverse text for  
492 language modeling,” *arXiv preprint arXiv:2101.00027*, 2021.

- 493 [48] A. Karvonen, C. Rager, J. Lin, C. Tigges, J. I. Bloom, D. Chanin, Y.-T. Lau, E. Farrell,  
 494 C. S. McDougall, K. Ayonrinde, D. Till, M. Wearden, A. Conmy, S. Marks, and N. Nanda,  
 495 “SAEBench: A comprehensive benchmark for sparse autoencoders in language model  
 496 interpretability,” in *Proceedings of the 42nd International Conference on Machine Learning*,  
 497 ser. Proceedings of Machine Learning Research, vol. 267. PMLR, 2025, pp. 29 223–29 264.  
 498 [Online]. Available: <https://proceedings.mlr.press/v267/karvonen25a.html>
- 499 [49] M. De-Arteaga, A. Romanov, H. Wallach, J. Chayes, C. Borgs, A. Chouldechova, S. C. Geyik,  
 500 K. Kenthapadi, and A. T. Kalai, “Bias in bios: A case study of semantic representation bias  
 501 in a high-stakes setting,” in *Proceedings of the Conference on Fairness, Accountability, and*  
 502 *Transparency (FAT\* ’19)*. ACM, 2019, pp. 120–128.
- 503 [50] X. Zhang, J. Zhao, and Y. LeCun, “Character-level convolutional networks for text classification,”  
 504 in *Advances in Neural Information Processing Systems*, vol. 28, 2015. [Online]. Available: <https://proceedings.neurips.cc/paper/2015/file/250cf8b51c773f3f8dc8b4be867a9a02-Paper.pdf>
- 506 [51] P. Koehn, “Europarl: A parallel corpus for statistical machine translation,” in *Proceedings of*  
 507 *Machine Translation Summit X: Papers*, Phuket, Thailand, Sep. 13-15 2005, pp. 79–86.
- 508 [52] CodeParrot, “Github code dataset,” <https://huggingface.co/datasets/codeparrot/github-code>,  
 509 2022.
- 510 [53] Y. Hou, J. Li, Z. He, A. Yan, X. Chen, and J. McAuley, “Bridging language and items for  
 511 retrieval and recommendation,” *arXiv preprint arXiv:2403.03952*, 2024.
- 512 [54] R. A. Harshman, “Foundations of the PARAFAC procedure: Models and conditions for an  
 513 “explanatory” multi-modal factor analysis,” *UCLA working papers in phonetics*, vol. 16, no. 1,  
 514 p. 84, 1970.

## 515 A Implementation Details

516 This appendix documents the training setup, hyperparameters, and infrastructure used for all experiments in the  
 517 main paper. All crosscoder variants (standard crosscoder, TR `fmxcoder`, and CP `fmxcoder`) were trained with  
 518 identical optimizer, schedule, and data settings; the only differences are the encoder/decoder parameterization  
 519 and the presence of stochastic layer masking.

### 520 A.1 Architectures and Hyperparameters

521 **Dictionary configuration.** All crosscoders and `fmxcoders` use a dictionary width of  $d_{\text{sae}} = 16,384$  latents.  
 522 The sparsity operator is BatchTopK [8] with  $K = 64$  active latents per token. A pre-activation ReLU is applied  
 523 before BatchTopK selection. The encoder bias  $\mathbf{b}^{\text{enc}} \in \mathbb{R}^{d_{\text{sae}}}$  is shared across layers; per-layer decoder biases  
 524  $\mathbf{b}_\ell^{\text{dec}} \in \mathbb{R}^d$  are learned independently for each  $\ell \in [1, L]$ .

525 **Layer selection.** For each base LLM, we use  $L = 8$  residual-stream activations collected at post-MLP positions  
 526 spanning the middle layers. The specific layer indices are (zero-indexed):

- 527 • GPT2-Small (12 layers): layers 2–9.
- 528 • Pythia-410M (24 layers): layers 8–15.
- 529 • Pythia-1.4B (24 layers): layers 8–15.
- 530 • Gemma2-2B (26 layers): layers 9–16.

531 **Factorization ranks.** For TR `fmxcoders` matching the standard crosscoder parameter count, the ranks  
 532  $(R_1, R_2, R_3)$  were chosen subject to  $R_2/R_1 = \sqrt{d/L}$  and  $R_3/R_1 = \sqrt{d_{\text{sae}}/d}$ . The specific values per  
 533 base LLM are:

- 534 • GPT2-Small ( $d = 768$ ): (5, 244, 25).
- 535 • Pythia-410M ( $d = 1024$ ): (7, 302, 27).
- 536 • Pythia-1.4B ( $d = 2048$ ): (11, 501, 31).
- 537 • Gemma2-2B ( $d = 2304$ ): (12, 545, 32).

538 For CP `fmxcoders`, the rank  $R$  was chosen to match the standard crosscoder parameter count; specific values  
539 per base LLM are:

- 540 • GPT2-Small:  $R = 5866$ .
- 541 • Pythia-410M:  $R = 7707$ .
- 542 • Pythia-1.4B:  $R = 14557$ .
- 543 • Gemma2-2B:  $R = 16153$ .

544 **Masking probability.** The default masking probability for `fmxcoders` reported in the main results is  $p = 0.05$ .  
545 The unmasked variant uses  $p = 0$ . Bernoulli masks are drawn independently per token and per layer at every  
546 training step. Ablations of the masking probability are presented in Appendix C.

## 547 A.2 Optimization

548 **Optimizer.** We use Adam with a learning rate  $3 \times 10^{-4}$ ,  $(\beta_1, \beta_2) = (0.9, 0.999)$ , with no warmup or decay  
549 schedules. For training stabilization, we use gradient clipping with a maximum norm of 1.0.

550 **Batch size.** 4096.

551 **Training tokens.** All training runs use 500M tokens (300M for GPT2-Small) at a context length of 128.

552 **Initialization.** Encoder and decoder tensors / decomposition factors follow a variance-matched i.i.d. Gaussian  
553 initialization, chosen so that the materialized 3-D weight tensor has the same per-element variance as a Kaiming-  
554 uniform-initialized 2-D slice. Biases are initialized as zero.

555 **Decoder normalization.** No decoder normalization is applied (e.g. similar to [3]).

## 556 A.3 Data

557 **Activation generation.** For Gemma2-2B and GPT2-Small, activations were generated using OpenWebText [46].  
558 For Pythia-410M and Pythia-1.4B, activations were generated using an uncopyrighted variant of the deduplicated  
559 Pile [47]. All activations are cached at the post-MLP residual-stream position for each of the  $L = 8$  selected  
560 layers.

561 **Held-out evaluation set.** Reconstruction metrics are computed on a held-out set of approximately 410k tokens  
562 not seen during training. The same held-out set is used for all variants of a given base LLM to ensure direct  
563 comparability.

## 564 A.4 Sparse Probing Setup

565 We follow the protocol of [9] for sparse probing. For each classification task, we select the top- $K = 1$  latent  
566 that best separates classes by F1 on a training split, then report F1 and Wasserstein distance on a held-out  
567 evaluation split. Class labels and dataset splits follow the original sources for Bias in Bios [49], AG News [50],  
568 EuroParl [51], GitHub programming languages [52], Amazon Sentiment, and Amazon 15 [53].

## 569 A.5 Functional Coherence Computation

570 The functional coherence metric (Section 4.2) is computed on a held-out set of 10M tokens, in order to obtain  
571 enough sample size for individual latent evaluations. For each layer  $\ell \in [1, L]$ , we run an additional forward  
572 pass through the crosscoder with  $\mathbf{x}_\ell$  set to zero while leaving the other layers intact, and record the resulting  
573 per-latent activations. The numerical-stability constant in Equation (5) is set to  $\epsilon = 10^{-8}$ .

## 574 A.6 Compute and Reproducibility

575 **Hardware.** All experiments were run on NVIDIA H100 64GB GPUs, with a single GPU per training run.

576 **Wall-clock time.** A single `fmxcoder` run takes approximately 12 hours for GPT2-Small, 16 hours for Pythia-  
577 410M, 30 hours for Pythia-1.4B, and 42 hours for Gemma2-2B.

578 **Total compute.** The full set of experiments reported in the main paper and appendices required approximately  
579 1500 GPU hours of compute across all base LLMs, factorization variants, and ablation settings.

580 **Software.** Crosscoder training is implemented in SAELens [42], a PyTorch-based SAE training library in  
581 Python. Tensor Ring and CP factorizations are implemented as native PyTorch modules. Reconstruction metrics  
582 (MSE, EV, CS) are computed following the protocols of SAEBench [48]. Probing F1 and Wasserstein distance  
583 follow [9]. The LLM-as-a-judge evaluation calls GPT-4o-mini through the OpenRouter API.

584 **B CP fmxcoders**

585 **B.1 Mathematical definition**

586 The Canonical Polyadic (CP) decomposition [30, 54] approximates a tensor as a sum of rank-one tensors formed  
 587 from vector outer products. For the encoder tensor with rank  $R$ ,

$$\mathbf{E} \approx \sum_{r \in [1, R]} \mathbf{w}_r \circ \mathbf{u}_r \circ \mathbf{v}_r, \quad \mathbf{E}_{j,i,\ell} \approx \sum_{r \in [1, R]} W_{jr} U_{ir} V_{\ell r}, \quad (6)$$

588 where  $\mathbf{w}_r \in \mathbb{R}^d$ ,  $\mathbf{u}_r \in \mathbb{R}^{d_{sae}}$ ,  $\mathbf{v}_r \in \mathbb{R}^L$  are the factor vectors, collected into factor matrices  $\mathbf{W} \in \mathbb{R}^{d \times R}$ ,  
 589  $\mathbf{U} \in \mathbb{R}^{d_{sae} \times R}$ ,  $\mathbf{V} \in \mathbb{R}^{L \times R}$ . The decoder factorizes independently with factors  $\tilde{\mathbf{W}}, \tilde{\mathbf{U}}, \tilde{\mathbf{V}}$  of matching shapes.

590 Each rank-one term is an outer product of an activation-space direction  $\mathbf{w}_r$ , a feature-weight vector  $\mathbf{u}_r$ , and a  
 591 layer pattern  $\mathbf{v}_r$ , making the full encoder a sum of  $R$  atoms shared across all features. This results in a *uniform*  
 592 *coupling* inductive bias: a single rank  $R$  ties the three modes together at once, and every feature expresses its  
 593 cross-layer behavior in the same  $R$ -dimensional basis of layer patterns (the rows of  $\mathbf{V}$ ). Crucially, a feature that  
 594 appears across layers *is represented once* through its  $R$ -dimensional coefficients rather than re-discovered at  
 595 each layer, and the parameter count drops from  $\mathcal{O}(L \cdot d_{sae} \cdot d)$  to  $\mathcal{O}(R \cdot (d + d_{sae} + L))$  per tensor.

596 **B.2 Results**

597 This section presents reconstruction and sparse-probing results for CP fmxcoders, reported in Table 4. The  
 598 qualitative picture matches the TR results in Table 3: CP fmxcoders improve over the standard crosscoder on  
 599 every metric and every LLM, reducing MSE by up to 50%, recovering cosine similarity to within 0.01 of the  
 600 per-layer SAE baseline, and lifting average probing F1 by up to 30 points above the crosscoder. As in the TR  
 601 case, the masking variant trades a small amount of reconstruction quality and F1 for a denoising effect that is  
 602 reflected in the functional coherence metric (see Appendix C). Comparing the two factorizations directly, CP  
 603 fmxcoders match TR fmxcoders closely on MSE, explained variance, and cosine similarity, but under-perform  
 604 on average F1 on 3 out of 4 LLMs, and on Wasserstein distance (by roughly a factor of 2–5  $\times$ ). They nonetheless  
 605 remain well above the crosscoder baseline on Wasserstein, by one to two orders of magnitude. We attribute  
 606 this gap to the structural difference between the two decompositions discussed in Section 3.2: TR couples the  
 607 three modes through three independent ranks ( $R_1, R_2, R_3$ ) with context-dependent interaction across features  
 608 and layers, while CP ties all three modes through a single shared rank  $R$  and a uniform basis of layer patterns.  
 609 The TR parameterization therefore admits richer per-feature layer profiles, more suitable for our setting where  
 610  $L \ll d_{sae}$ , which appears to translate into more separable latent activations on the probing tasks.

Table 4: Reconstruction metrics and probing-based results for different base LLMs, for SAE, cross-  
 coder, and fmxcoder variants. fmxcoders adopt the CP decomposition, with ranks chosen to match  
 the parameters of the crosscoders in each case. We report average F1 across six classification tasks  
 and the detailed results are presented in Appendix D. SAE results were reproduced using the settings  
 from [9].

| LLM         | Variant                 | MSE         | EV $\uparrow$ | CS $\uparrow$ | Avg. F1 $\uparrow$ | Avg. Wass. $\uparrow$<br>( $\times 10^{-3}$ ) |
|-------------|-------------------------|-------------|---------------|---------------|--------------------|---|
| GPT2-Small  | SAE                     | 0.53        | 0.9256        | 0.9623        | 65.7               | 8.4   |
|             | Crosscoder              | 0.64        | 0.8704        | 0.9346        | 67.6               | 11.6  |
|             | fmxcoder ( $p = 0$ )    | <b>0.32</b> | <b>0.9355</b> | <b>0.9679</b> | <b>74.3</b>        | <b>25.8</b>                                   |
|             | fmxcoder ( $p = 0.05$ ) | 0.33        | 0.9330        | 0.9668        | 69.7               | 18.5  |
| Pythia-410m | SAE                     | <b>0.03</b> | <b>0.9918</b> | <b>0.9364</b> | 65.0               | 0.5   |
|             | Crosscoder              | 0.05        | 0.9856        | 0.8696        | 47.4               | <0.1  |
|             | fmxcoder ( $p = 0$ )    | <b>0.03</b> | 0.9917        | 0.9273        | <b>77.5</b>        | <b>4.7</b>                                    |
|             | fmxcoder ( $p = 0.05$ ) | <b>0.03</b> | 0.9915        | 0.9263        | 65.6               | 2.1   |
| Pythia-1.4b | SAE                     | <b>0.22</b> | <b>0.9784</b> | <b>0.9289</b> | 64.6               | 0.8   |
|             | Crosscoder              | 0.31        | 0.9691        | 0.8976        | 47.0               | 0.2   |
|             | fmxcoder ( $p = 0$ )    | <b>0.22</b> | 0.9781        | 0.9275        | <b>76.2</b>        | <b>10.8</b>                                   |
|             | fmxcoder ( $p = 0.05$ ) | 0.23        | 0.9777        | 0.9263        | 65.8               | 5.5   |
| Gemma2-2b   | SAE                     | <b>1.58</b> | <b>0.8702</b> | 0.9215        | 64.8               | 2.7   |
|             | Crosscoder              | 3.29        | 0.7809        | 0.8737        | 51.2               | 0.2   |
|             | fmxcoder ( $p = 0$ )    | 2.05        | 0.8622        | <b>0.9226</b> | <b>79.0</b>        | <b>38.3</b>                                   |
|             | fmxcoder ( $p = 0.05$ ) | 2.05        | 0.8621        | 0.9223        | 65.4               | 16.9  |

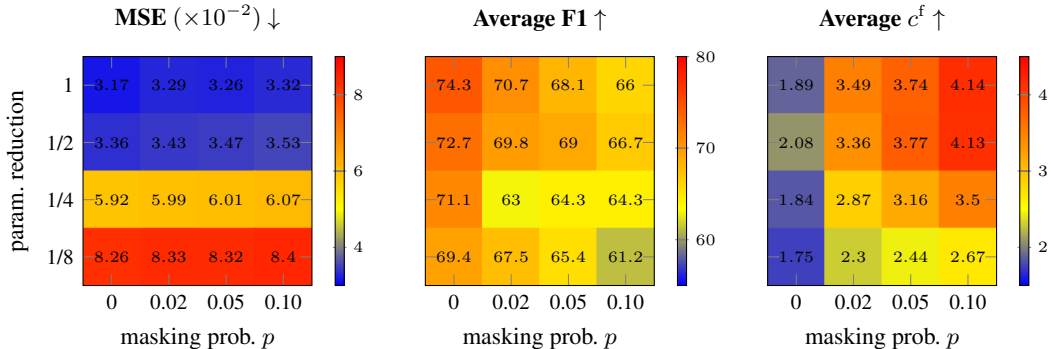


Figure 3: Heat maps for MSE, mean F1 and mean functional coherence  $c^f$ , for different values of the masking probability  $p$  and parameter reduction associated with decreasing factorization ranks for TR fmxcoders.  $(R_1, R_2, R_3)$  take the values  $(7, 302, 27)$  for matching the unfactorized crosscoder number of parameters,  $(5, 214, 19)$  for a 1/2 reduction,  $(3, 151, 13)$  for a 1/4 reduction, and  $(2, 107, 9)$  for a 1/8 reduction. In all cases, the activations were collected from Pythia-410m.

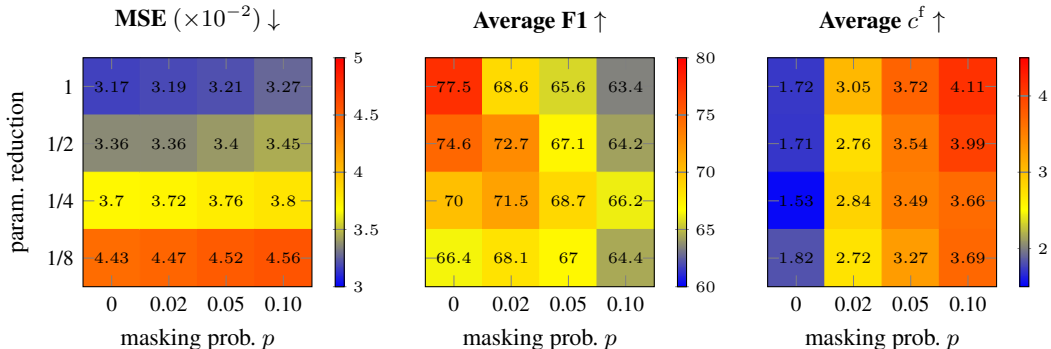


Figure 4: Similar to Figure 3, for CP fmxcoders.  $R=7707$  for matching the unfactorized crosscoder number of parameters,  $R=3853$  for a 1/2 reduction,  $R=1926$  for a 1/4 reduction, and  $R=963$  for a 1/8 reduction.

## 611 C Mask Probability and Rank Ablations

612 This appendix reports ablations of the two hyperparameters introduced by fmxcoders: the masking probability  
 613  $p$  and the factorization ranks. Figure 3 shows results for the TR variant and Figure 4 for the CP variant, both  
 614 trained on Pythia-410m. Each figure reports MSE, mean probing F1, and average functional coherence  $c^f$  as  $p$   
 615 varies in  $\{0, 0.02, 0.05, 0.10\}$  and the parameter count is reduced by factors of 1, 1/2, 1/4, and 1/8 relative to  
 616 the standard crosscoder. The ranks were chosen to match the target parameter count subject to the constraints  
 617  $R_2/R_1 = \sqrt{d/L}$  and  $R_3/R_1 = \sqrt{d_{\text{sac}}/d}$ , which balance the contribution of each mode to the factorization.  
 618 The corresponding rank settings are listed in the figure captions.

619 **Effect of masking probability.** Across both factorizations,  $p$  controls a clean trade-off between reconstruc-  
 620 tion quality and functional coherence. Functional coherence increases monotonically with  $p$ , roughly doubling  
 621 from  $c^f \approx 2$  at  $p=0$  to  $c^f \approx 4$  at  $p=0.1$  for both TR and CP at full rank. Reconstruction MSE is essentially flat in  
 622  $p$  at any fixed rank (e.g.  $3.17 \rightarrow 3.32$  for TR at full rank), confirming that masking exerts pressure on *which*  
 623 *latents are learned* rather than on overall reconstruction capacity. Average F1 declines moderately with  $p$  (up to  
 624 14 points), consistent with the denoising objective penalizing layer-localized features in favor of more distributed  
 625 ones. The choice of  $p$  is therefore a knob the user can set based on whether downstream interpretability or  
 626 sparse-probing performance is the priority: small  $p$  ( $\leq 0.02$ ) preserves probing F1 with modest coherence gains,  
 627 while  $p \geq 0.05$  produces substantially more cross-layer latents at a small reconstruction and F1 cost.

628 **Effect of factorization rank.** Reducing the rank reduces parameter count and degrades reconstruction  
 629 and F1, but does not affect functional coherence. At  $p=0$ ,  $c^f$  remains in the narrow band 1.5–2.1 across the  
 630 full  $8\times$  compression sweep, indicating that capacity controls how well features are reconstructed but not how  
 631 distributed they are across layers. Functional coherence is evidently governed by the optimization pressure from  
 632 masking, not by the size of the parameter budget. The two factorizations differ markedly in how gracefully they

633 tolerate compression: CP’s MSE rises only from approximately 3.2 to 4.5 at 1/8 reduction, while TR’s rises  
 634 from approximately 3.2 to 8.3, more than doubling.

635 **Joint behavior.** The two hyperparameters act largely independently:  $p$  shifts coherence with weak effect on  
 636 reconstruction, while rank shifts reconstruction and F1 with no effect on coherence. The best F1 is consistently  
 637 obtained at full rank with  $p=0$  (77.5 for CP, 74.3 for TR), and the best coherence at full rank with  $p=0.1$   
 638 ( $c^f \approx 4.1$  for both). No setting in the sweep simultaneously maximizes both, reflecting the fundamental trade-off  
 639 introduced by the denoising objective; intermediate values such as  $p=0.05$  at full rank offer a reasonable balance  
 640 and are the settings used in the main results.

## 641 D Full probing results

642 Table 5 presents the per-dataset probing results underlying the averages reported in Table 3 and Table 4. Three  
 643 patterns are worth noting. First, the standard crosscoder’s under-performance relative to the per-layer SAE is  
 644 broad rather than driven by any single dataset, with losses of 10–30 F1 points across most LLM and dataset  
 645 combinations on the three larger models. Second, the gains of TR `fmxcoders` are similarly systematic, beating  
 646 the SAE on 22 of the 24 LLM and dataset combinations. Third, gaps in Wasserstein distances are larger than gaps  
 647 in F1 in relative terms, with `fmxcoder` values exceeding the crosscoder by one to three orders of magnitude even  
 648 on datasets where F1 improvements are modest, supporting the main-text observation that F1 alone understates  
 649 the separation gain. The comparison between TR and CP is dataset-dependent: TR achieves the higher average  
 650 F1 on 3 out of 4 LLMs, but CP wins on individual datasets in several cases.

Table 5: Probing experiments across datasets at  $K=1$ . Format: F1 scores (%) / Wasserstein distance ( $\times 10^{-3}$ ). The last two columns show average F1 and average Wasserstein distance. `fmxcoders` correspond to  $p=0$ .

| LLM         | Variant                  | EuroParl     | Bios         | Amazon Sentiment | GitHub       | AG News      | Amazon 15   | Avg. F1 (%) $\uparrow$ | Avg. Wass. $\uparrow$ ( $\times 10^{-3}$ ) |
|-------------|--------------------------|--------------|--------------|------------------|--------------|--------------|-------------|------------------------|--|
| GPT2-Small  | SAE                      | 67.4 / 19.0  | 59.6 / 7.3   | 68.8 / 4.4       | 68.1 / 8.8   | 65.3 / 8.2   | 65.1 / 2.9  | 65.7                   | 8.4  |
|             | Crosscoder               | 85.0 / 23.0  | 55.7 / 6.2   | 67.9 / 7.5       | 66.6 / 20.4  | 62.5 / 6.0   | 68.1 / 6.2  | 67.6                   | 11.6                                       |
|             | CP <code>fmxcoder</code> | 84.8 / 55.7  | 72.3 / 24.3  | 81.9 / 11.3      | 68.9 / 28.0  | 66.8 / 26.8  | 71.1 / 8.5  | 74.3                   | 25.8                                       |
|             | TR <code>fmxcoder</code> | 93.2 / 256.6 | 78.8 / 114.7 | 82.1 / 53.8      | 67.8 / 126.1 | 75.6 / 130.9 | 72.5 / 38.9 | 78.3                   | 120.2                                      |
| Pythia-410m | SAE                      | 90.9 / 0.8   | 60.5 / 0.3   | 63.9 / 0.4       | 59.7 / 1.1   | 58.7 / 0.3   | 56.6 / 0.3  | 65.0                   | 0.5  |
|             | Crosscoder               | 58.1 / 0.2   | 31.4 / 0.0   | 55.6 / 0.0       | 58.7 / 0.0   | 33.9 / 0.0   | 46.7 / 0.0  | 47.4                   | 0.0  |
|             | CP <code>fmxcoder</code> | 91.5 / 7.0   | 74.9 / 4.4   | 79.7 / 2.7       | 72.0 / 7.8   | 73.6 / 4.3   | 73.0 / 1.9  | 77.5                   | 4.7  |
|             | TR <code>fmxcoder</code> | 87.3 / 13.7  | 65.1 / 8.9   | 86.6 / 6.8       | 75.6 / 18.6  | 69.3 / 8.8   | 61.9 / 4.5  | 74.3                   | 10.2                                       |
| Pythia-1.4b | SAE                      | 74.0 / 0.6   | 65.0 / 0.5   | 57.2 / 0.6       | 63.3 / 2.3   | 65.2 / 0.4   | 63.2 / 0.4  | 64.6                   | 0.8  |
|             | Crosscoder               | 40.1 / 0.0   | 40.2 / 0.0   | 50.7 / 0.2       | 59.3 / 0.9   | 33.5 / 0.0   | 58.5 / 0.2  | 47.0                   | 0.2  |
|             | CP <code>fmxcoder</code> | 98.7 / 14.1  | 70.4 / 9.6   | 88.1 / 6.4       | 60.4 / 20.4  | 69.2 / 9.6   | 70.6 / 4.6  | 76.2                   | 10.8                                       |
|             | TR <code>fmxcoder</code> | 94.3 / 46.6  | 72.2 / 34.0  | 89.2 / 22.9      | 78.4 / 63.1  | 66.6 / 34.3  | 69.6 / 15.8 | 78.4                   | 36.1                                       |
| Gemma2-2b   | SAE                      | 68.3 / 1.9   | 67.6 / 2.6   | 71.1 / 2.8       | 64.4 / 4.5   | 59.9 / 2.7   | 57.6 / 1.9  | 64.8                   | 2.7  |
|             | Crosscoder               | 69.9 / 1.0   | 55.6 / 0.2   | 36.6 / 0.0       | 51.4 / 0.0   | 50.7 / 0.1   | 43.0 / 0.0  | 51.2                   | 0.2  |
|             | CP <code>fmxcoder</code> | 94.8 / 38.3  | 81.5 / 47.4  | 93.1 / 28.7      | 75.3 / 46.7  | 59.9 / 49.1  | 69.6 / 19.6 | 79.0                   | 38.3                                       |
|             | TR <code>fmxcoder</code> | 97.3 / 190.5 | 82.0 / 252.1 | 93.3 / 150.9     | 86.1 / 231.8 | 63.8 / 268.0 | 68.7 / 99.2 | 81.9                   | 198.7                                      |

## 651 E LLM-as-a-Judge Procedure

652 For each latent we collect its top-activating tokens, sorted by maximum activation across the residual-stream  
 653 activations of the base LLM on a held-out subset of its training corpus. We additionally include up to two context  
 654 windows of 80 characters on each side of the activating token. The resulting prompt is sent to *GPT-4o-mini* via  
 655 the OpenRouter API, with a system message instructing it to evaluate the latents. The returned `semantic_score`  
 656 and `surface_score` are each clipped to  $[0, 1]$ , and a latent is labeled *semantic* when `semantic_score`  $> 0.7$   
 657 and `surface_score`  $< 0.3$ , and *surface* with these conditions reversed. The full user-message template is  
 658 shown below:

You are evaluating a latent feature from a sparse autoencoder trained on a language model. You are given the tokens that most strongly activate this feature.

Score the feature on two independent dimensions:

659

**\*\*semantic\_score\*\*** (0 to 1, continuous): How strongly do these tokens collectively represent a coherent, high-level concept that is interpretable by humans?  
- 1.0 = tokens clearly belong to a unified semantic category (e.g. US states and cities, negative-emotion words, cooking verbs, animal species, financial terms)  
- 0.5 = tokens share some thematic connection but it's loose or partial  
- 0.0 = no discernible high-level concept; tokens seem unrelated or random

**\*\*surface\_score\*\*** (0 to 1, continuous): How strongly do these tokens collectively describe low-level linguistic or surface patterns rather than meaning?  
- 1.0 = tokens are unified by syntax, morphology, punctuation, character class, or formatting (e.g. closing braces, digits, -ing suffixes, markup tags, subword fragments)  
- 0.5 = tokens partially share surface-level properties but also carry some semantic content  
- 0.0 = tokens are not unified by any surface-level pattern

These two scores are INDEPENDENT and can vary continuously from 0 to 1. A feature can be high on both (rare), low on both (noise/random), or high on one and low on the other (typical).

Think carefully about whether the unifying pattern is semantic (about meaning) or surface-level (about form/syntax/characters).

**\*\*Feature ID\*\***: {feature\_id}  
**\*\*Top-activating tokens\*\*** (ordered by activation strength): {tokens}

**\*\*Example contexts\*\*** (activating token wrapped in <<>>):  
- token "{tok}":  
  - {context\_window}  
  ...

Return ONLY valid JSON with keys "semantic\_score" and "surface\_score", both floats between 0 and 1.

Example: {"semantic\_score": 0.85, "surface\_score": 0.1}

660