

---

# ROUTERINTERP: UNDERSTANDING SUPERPOSED SPECIALISATION IN MOE ROUTING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Sparse Mixture of Experts (MoE) models scale more efficiently than dense models by routing tokens to modular expert networks that are only active when relevant to the task. A leading hypothesis for the performance of MoE models is that each expert specialises in a single, coherent domain. However, interpretability efforts that assume this hypothesis have generally been unsuccessful. We propose and present evidence for an alternative account that we call the *Superposed Specialisation Hypothesis* (SSH): experts specialise in a disjoint union of fine-grained features rather than one broad domain. Leveraging the SSH, we introduce *RouterInterp*, a method for interpreting expert routing that identifies Sparse Autoencoder features most predictive of routing decisions and produces unified natural language explanations. On gpt-oss-20b, explanations from RouterInterp predict expert routing with 77% higher accuracy than prior methods. This work provides a scalable method for generating concise and more accurate explanations of expert routing and increases our understanding of a previously uninterpretable component of foundation models.

## 1 INTRODUCTION

Sparse Mixture-of-Experts (MoE) transformers have emerged as a promising approach for scaling frontier language models (Cai et al., 2025; Fedus et al., 2022b; Du et al., 2022). The strong performance of Sparse MoE models has often been attributed to expert specialisation (Lewis et al., 2021): if each expert learns to handle a subset of the input data distribution or perform only a subset of computations, then using only a subset of the parameters for each input can be both *effective* and *efficient*. We call this explanation for the success of MoEs the **Specialisation Hypothesis**.

Prior work has struggled to recover interpretable patterns of expert specialisation (Jiang et al., 2024; Lewis et al., 2021; Zoph et al., 2022). We formalise two distinct forms of the Specialisation Hypothesis. Under the **Domain Specialisation Hypothesis** (DSH), each expert specialises in a single coherent domain: semantically similar fine-grained categories (micro-domains) cluster within an expert. Under the **Superposed Specialisation Hypothesis** (SSH), each expert specialises in a disjoint collection of features spanning multiple unrelated micro-domains, mirroring the superposition of features in shared neurons (Elhage et al., 2022).

Inspired by SSH, we develop **RouterInterp**, a method that interprets MoE routing via sparse autoencoder (SAE) latents (Makhzani & Frey, 2013; Cunningham et al., 2024; Bricken et al., 2023): it identifies SAE latents most predictive of each expert’s activation, generates natural language explanations for those latents, and aggregates them into expert-level explanations (Section 3). RouterInterp achieves 89% and 81% accuracy in predicting expert routing from SAE features on OLMoE-1B-7B and gpt-oss-20b, and attains a 0.62 explanation score on gpt-oss-20b, outperforming token-based explanations (0.35) (Section 4).

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

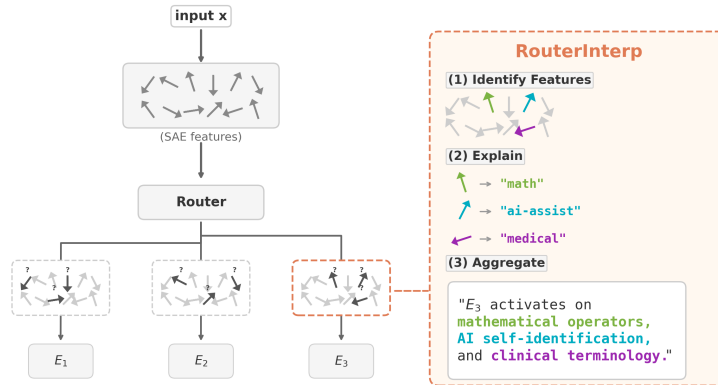


Figure 1: RouterInterp explains routing decisions as a combination of interpretable features. Sparse autoencoder (SAE) latents, sparse feature representations extracted from model activations, pick out directions in activation space that correspond to interpretable features. RouterInterp *identifies* SAE latents that are most predictive of routing, *explains* these features in concise natural language strings, and *aggregates* these explanations into an expert-level explanation. While no single monosemantic concept captures an expert’s behavior, aggregating multiple features makes routing understandable.

## 2 SPECIALISATION HYPOTHESES

Here, we formalise the two hypotheses (DSH and SSH) introduced in Section 1, state their predictions, and give theoretical motivations for why SSH might hold in practice. Suppose that we have a set of  $D$  micro-domains occurring in a corpus  $\mathcal{C}$ , denoted  $\mathcal{D} = \{d_1, \dots, d_D\}$ . Suppose also that we have an MoE layer with  $E$  experts. Firstly, note that when  $D = E$ , then we should expect the experts to specialise in a single micro-domain. Secondly, when  $D < E$ , then we should expect multiple experts to specialise in the same micro-domain or for there to be redundant experts which are never routed to. Both of these cases are consistent with both the DSH and SSH. However, in case of  $D > E$ , the two hypotheses make different predictions. Under the DSH, similar micro-domains are routed to the same expert and so an expert specialises in a semantically coherent domain consisting of adjacent micro-domains - expert specialisation can be explained by a single (monosemantic) concept or domain. Under the SSH, dissimilar micro-domains are routed to the same expert and so an expert specialises in a collection of semantically disjoint domains - expert specialisation can only be explained by multiple (polysemantic) concepts or domains.

There are two core arguments for why we might expect the Superposed Specialisation Hypothesis to be a more accurate description of expert specialisation. First, we argue that the router is incentivised to assign dissimilar domains to the same expert so that the expert can perform Computation in Superposition (Hänni et al., 2024; Linsefors & Bushnaq, 2025a;b; Newgas, 2025), i.e., a single expert can have multiple disjoint transforms that it applies depending on the input domain. This mirrors Elhage et al. (2022): just as sparse features that are not frequently co-activated can share the same neurons, dissimilar domains that are not frequently co-activated can share an expert. In both cases, low co-activation reduces the *interference* cost of multiple features or domains sharing the same capacity, which is why superposition is effective for models: it allows representing more features (micro-domains) than the number of available neurons (experts). Second, MoE training commonly includes a *load-balancing loss* that encourages the router to distribute tokens evenly across experts (Shazeer et al., 2017a; Fedus et al., 2022a). With relatively small batch sizes, this can unintentionally encourage routing dissimilar inputs to the same expert: a single sequence typically contains tokens from the same macro-domain, so a small batch contains few macro-domains and the load-balancing loss distributes tokens from these few across all experts. This necessarily requires that some tokens from the same macro-domain will be routed to different experts - in opposition to the DSH, where we would expect all tokens from the same macro-domain to be routed to the same expert.

108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161

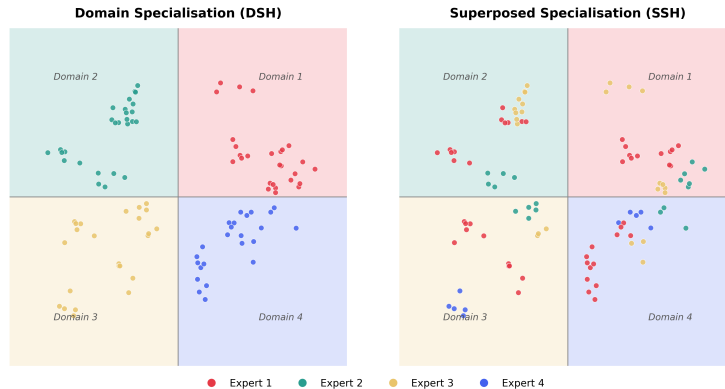


Figure 2: Two hypotheses for expert specialisation in MoE models: (a) Domain Specialisation (DSH) and (b) Superposed Specialisation (SSH). (a) Under the DSH, examples from the same semantic domain (coloured region) are routed to the same expert (same colour); experts specialise in one coherent domain. (b) Under the SSH, examples routed to the same expert (same colour) are scattered across disjoint domains (different coloured regions); clusters of points with the same expert are *micro-domains* (fine-grained features), and experts specialise in a disjoint union of micro-domains. We suggest the SSH better describes expert specialisation in practice.

### 3 ROUTERINTERP

RouterInterp operates in three stages: (1) identifying SAE features most relevant to each expert’s routing behaviour, (2) generating natural language explanations for each feature, and (3) aggregating these into expert-level explanations.

**Feature Selection.** We identify which SAE features are most predictive of each expert’s activation using the ‘ $\rho$ -usefulness’ metric (Ilyas et al., 2019), a data-driven selection criterion that measures how discriminative a feature is for expert selection. A feature  $f$  is  $\rho$ -useful for expert  $E_i$  if the expected activation of  $f$  is higher when expert  $E_i$  is selected than when it is not:

$$\rho(f, E_i) = \mathbb{E}[f(\mathbf{x}) \mid E_i \in \mathcal{T}(\mathbf{x})] - \mathbb{E}[f(\mathbf{x}) \mid E_i \notin \mathcal{T}(\mathbf{x})] \quad (1)$$

where  $\mathcal{T}(\mathbf{x})$  denotes the set of selected experts for input  $\mathbf{x}$ . Features with high  $\rho$ -usefulness scores are specifically predictive of a particular expert’s activation. For each expert  $E_i$ , we select the top- $n$  features ranked by  $\rho$ -usefulness score, denoting this set of features  $\mathbb{F}_i$ <sup>1</sup>.

**Feature Explanations.** Given the selected feature set  $\mathbb{F}_i$  for each expert  $E_i$ , we generate natural language explanations for each feature following the automated interpretability framework of (Paulo et al., 2025). For each feature  $f \in \mathbb{F}_i$ , we present a language model (the *explainer*) with *positive examples* where  $f$  activates strongly and *contrastive negatives* where  $f$  does not activate despite being semantically similar (in embedding space). The explainer infers what concept or pattern the feature detects, producing a concise natural language description.

**Aggregation.** To obtain an expert-level explanation, we prompt a language model with all  $n$  feature explanations for expert  $E_i$  and ask it to synthesise a single coherent paragraph describing when the expert activates. The goal is an explanation that enables predicting when the expert activates; we evaluate this using an LLM as a proxy for human judgment.

<sup>1</sup>We evaluate the effect of feature set size  $|\mathbb{F}_i|$  on explanation quality and compare against an alternative selection method based on cosine similarity between router weight vectors and SAE decoder directions in Section F.

	Layer 4	Layer 8	Layer 12	Layer 16	Layer 20
RouterInterp	<b>0.60</b>	<b>0.64</b>	<b>0.61</b>	<b>0.62</b>	<b>0.60</b>
Bigram Baseline	0.23	0.47	0.42	0.35	0.40

Table 1: RouterInterp’s natural language explanations predict expert routing with much higher explanation scores than explanations generated from high-frequency token statistics. Explanation scores measure how accurately a language model can predict expert activation given only the explanation. Across different depths of gpt-oss-20b, RouterInterp outperforms the bigram based explanations by 36–161%. Each expert explanation aggregates the 10 most  $\rho$ -useful SAE features.

**Evaluation.** To evaluate the quality of the explanations, we adapt the AutoInterp scoring framework (Paulo et al., 2025). First, we construct a balanced evaluation set for each expert containing *positive examples* (contexts where a token was routed to the expert) and *hard negatives* (semantically similar contexts with no tokens routing to that expert). Then, an LLM (the *scorer*) receives the expert explanation and the example from the evaluation set and predicts whether the expert would be activated. This yields a binary prediction per example, which we compare against the ground truth routing decision. Final explanation score is reported as F1 score of this binary classification task. We compare RouterInterp’s explanations against an n-gram baseline: we take token n-grams (unigrams or bigrams) that most frequently co-occur with each expert and prompt an LLM to summarise them into a natural language description.

## 4 RESULTS

Table 1 shows results for gpt-oss-20b, where each expert’s explanation aggregates explanations of the 10 SAE features with highest  $\rho$ -usefulness scores for that expert. RouterInterp achieves 0.60–0.64 across layers, substantially outperforming the explanations generated from bigram-expert co-occurrence frequencies (bigram baseline: 0.23–0.47). We illustrate examples of RouterInterp’s explanations in Section G. To test the SSH, we also evaluate whether SAE latents can predict expert routing well. We extract model activations and encode them with an SAE encoder to produce sparse feature representations. We train a linear classifier on these sparse representations to predict whether a given expert is activated for some context. On OLMoE-1B-7B and gpt-oss-20b, the SAE-based predictor achieves 0.89 and 0.81 Recall respectively, substantially outperforming unigram and bigram baselines (0.52–0.74). These results support our hypothesis that routing is determined by features rather than surface-level token statistics. Setup, metric definition, and results across layers are in Section C.

## 5 CONCLUSION

We introduced *RouterInterp*, a method that interprets MoE routing decisions by identifying sparse autoencoder features most predictive of expert selection and aggregating their explanations into unified natural language descriptions. RouterInterp’s explanations achieve 0.62 explanation score compared to 0.35 for token-expert co-occurrence methods. RouterInterp emerged from formalizing two alternative specialisation hypotheses: the *Domain Specialisation Hypothesis* (DSH), which posits that experts specialize in semantically coherent domains, and the *Superposed Specialisation Hypothesis* (SSH), which models experts as responding to a superposition of distinct features spanning multiple unrelated micro-domains. Our experiments provide initial evidence for the SSH: High quality (81%/89% recall) of SAE-based routing prediction shows that experts are highly polysemantic, with the set of features most aligned with any single expert often spanning disjoint concepts. We acknowledge the following limitations of our method: (1) RouterInterp is downstream of SAE quality (whether features are monosemantic) and of the quality of feature-level explanations, and (2) textual summaries struggle to be both concise and complete as more features are included, so interactive presentations may better serve expert interpretation. We are excited to improve upon these issues in future work.

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269

---

## REFERENCES

- Samira Abnar, Harshay Shah, Dan Busbridge, Alaaeldin Mohamed Elnouby Ali, Josh Susskind, and Vimal Thilak. Parameters vs flops: Scaling laws for optimal sparsity for mixture-of-experts language models. In *ICLR Workshop, ICML*, 2025. URL <https://arxiv.org/abs/2501.12370>.
- Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, Andy Applebaum, Edwin Arbus, Rahul K Arora, Yu Bai, Bowen Baker, Haiming Bao, et al. gpt-oss-120b & gpt-oss-20b model card. *arXiv preprint arXiv:2508.10925*, 2025.
- Anthropic. System card: Claude haiku 4.5, October 2025. URL <https://www.anthropic.com/claude-haiku-4-5-system-card>.
- Kola Ayonrinde. Awesome adaptive computations, 2023. URL <https://github.com/koayon/awesome-adaptive-computation/>.
- Kola Ayonrinde, Michael T Pearce, and Lee Sharkey. Interpretability as compression: Reconsidering SAE explanations of neural activations. In *NeurIPS 2024 Workshop on Scientific Methods for Understanding Deep Learning*, 2024. URL <https://openreview.net/forum?id=hAqeEZRVSd>.
- Andrea Banino, Jan Balaguer, and Charles Blundell. Pondernet: Learning to ponder. *arXiv preprint arXiv:2107.05407*, 2021.
- Nora Belrose, Zach Furman, Logan Smith, Danny Halawi, Igor Ostrovsky, Lev McKinney, Stella Biderman, and Jacob Steinhardt. Eliciting latent predictions from transformers with the tuned lens. *arXiv preprint arXiv:2303.08112*, 2023.
- Leonard Bereska and Stratis Gavves. Mechanistic interpretability for AI safety - a review. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL <https://openreview.net/forum?id=ePUVetPKu6>. Survey Certification, Expert Certification.
- Dan Braun, Lucius Bushnaq, and Lee Sharkey. Stochastic parameter decomposition. In *Submitted to ILLIAD 2: ODYSSEY*, 2025. URL <https://openreview.net/forum?id=dEdS9ao8gN>. under review.
- Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick Turner, Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu, Shauna Kravec, Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex Tamkin, Karina Nguyen, Brayden McLean, Josiah E Burke, Tristan Hume, Shan Carter, Tom Henighan, and Christopher Olah. Towards monosemanticity: Decomposing language models with dictionary learning. *Transformer Circuits Thread*, 2023. <https://transformer-circuits.pub/2023/monosemantic-features/index.html>.
- Lucius Bushnaq, Stefan Heimersheim, Nicholas Goldowsky-Dill, Dan Braun, Jake Mendel, Kaarel Hänni, Avery Griffin, Jörn Stöhler, Magdalena Wache, and Marius Hobbhahn. The local interaction basis: Identifying computationally-relevant and sparsely interacting features in neural networks. *arXiv preprint arXiv:2405.10928*, 2024.
- Weilin Cai, Juyong Jiang, Fan Wang, Jing Tang, Sunghun Kim, and Jiayi Huang. A survey on mixture of experts in large language models. *IEEE Transactions on Knowledge and Data Engineering*, 2025.
- Marmik Chaudhari, Jeremi Nuer, and Rome Thorstenson. Superposition in mixture of experts. In *Mechanistic Interpretability Workshop at NeurIPS 2025*, 2025. URL <https://openreview.net/forum?id=bZqopmfZDE>.
- Jeff Clune, Jean-Baptiste Mouret, and Hod Lipson. The evolutionary origins of modularity. *Proceedings of the Royal Society b: Biological sciences*, 280(1755):20122863, 2013.
- Valérie Costa, Thomas Fel, Ekdeep Singh Lubana, Bahareh Tolooshams, and Demba Ba. From flat to hierarchical: Extracting sparse representations with matching pursuit. *arXiv preprint arXiv:2506.03093*, 2025.

---

270 Róbert Csordás, Christopher Potts, Christopher D Manning, and Atticus Geiger. Recurrent  
271 neural networks learn to store and generate sequences using non-linear representations.  
272 In Yonatan Belinkov, Najoung Kim, Jaap Jumelet, Hosein Mohebbi, Aaron Mueller,  
273 and Hanjie Chen (eds.), *Proceedings of the 7th BlackboxNLP Workshop: Analyzing and*  
274 *Interpreting Neural Networks for NLP*, pp. 248–262, Miami, Florida, US, November 2024.  
275 Association for Computational Linguistics. doi:10.18653/v1/2024.blackboxnlp-1.17. URL  
276 <https://aclanthology.org/2024.blackboxnlp-1.17/>.

277 Hoagy Cunningham, Aidan Ewart, Logan Riggs Smith, Robert Huben, and Lee Sharkey.  
278 Sparse autoencoders find highly interpretable features in language models. In *The Twelfth*  
279 *International Conference on Learning Representations*, 2024. URL [https://openreview](https://openreview.net/forum?id=F76bwRSLek)  
280 [.net/forum?id=F76bwRSLek](https://openreview.net/forum?id=F76bwRSLek).

282 Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Łukasz Kaiser.  
283 Universal transformers. *arXiv preprint arXiv:1807.03819*, 2018.

284 Nan Du, Yanping Huang, Andrew M Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu,  
285 Maxim Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, et al. Glam: Efficient scaling of  
286 language models with mixture-of-experts. In *International conference on machine learning*,  
287 pp. 5547–5569. PMLR, 2022.

289 Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna  
290 Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse,  
291 Sam McCandlish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher  
292 Olah. Toy models of superposition. *Transformer Circuits Thread*, 2022. [https://](https://transformer-circuits.pub/2022/toy_model/index.html)  
293 [transformer-circuits.pub/2022/toy\\_model/index.html](https://transformer-circuits.pub/2022/toy_model/index.html).

294 Joshua Engels, Eric J Michaud, Isaac Liao, Wes Gurnee, and Max Tegmark. Not all language  
295 model features are one-dimensionally linear. In *The Thirteenth International Conference on*  
296 *Learning Representations*, 2025. URL [https://openreview](https://openreview.net/forum?id=d63a4AM4hb)  
297 [.net/forum?id=d63a4AM4hb](https://openreview.net/forum?id=d63a4AM4hb).

298 William Fedus, Jeff Dean, and Barret Zoph. A review of sparse expert models in deep  
299 learning. *arXiv preprint arXiv:2209.01667*, 2022a.

300 William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion  
301 parameter models with simple and efficient sparsity. *Journal of Machine Learning Research*,  
302 23(120):1–39, 2022b.

304 Daniel Filan, Stephen Casper, Shlomi Hod, Cody Wild, Andrew Critch, and Stuart Russell.  
305 Clusterability in neural networks. *arXiv preprint arXiv:2103.03386*, 2021.

306 Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster,  
307 Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor  
308 Leahy. The pile: An 800gb dataset of diverse text for language modeling, 2020. URL  
309 <https://arxiv.org/abs/2101.00027>.

311 Leo Gao, Tom Dupre la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford,  
312 Ilya Sutskever, Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders.  
313 In *The Thirteenth International Conference on Learning Representations*, 2025. URL  
314 <https://openreview.net/forum?id=tcsZt9ZNKD>.

315 Alex Graves. Adaptive computation time for recurrent neural networks. *arXiv preprint*  
316 *arXiv:1603.08983*, 2016.

318 Yizeng Han, Gao Huang, Shiji Song, Le Yang, Honghui Wang, and Yulin Wang. Dynamic  
319 neural networks: A survey. *IEEE transactions on pattern analysis and machine intelligence*,  
320 44(11):7436–7456, 2021.

322 Kaarel Hänni, Jake Mendel, Dmitry Vaintrob, and Lawrence Chan. Mathematical models of  
323 computation in superposition. In *ICML 2024 Workshop on Mechanistic Interpretability*,  
2024.

---

324 Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and  
325 Aleksander Madry. Adversarial examples are not bugs, they are features. In H. Wal-  
326 lach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.),  
327 *Advances in Neural Information Processing Systems*, volume 32. Curran Associates,  
328 Inc., 2019. URL [https://proceedings.neurips.cc/paper\\_files/paper/2019/file/  
329 e2c420d928d4bf8ce0ff2ec19b371514-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/e2c420d928d4bf8ce0ff2ec19b371514-Paper.pdf).

330 Robert A. Jacobs, Michael I. Jordan, Steven J. Nowlan, and Geoffrey E. Hinton. Adaptive mix-  
331 tures of local experts. *Neural Computation*, 3(1):79–87, 1991. doi:10.1162/neco.1991.3.1.79.  
332

333 Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary,  
334 Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian  
335 Bressand, et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.

336 Jakub Krajewski, Jan Ludziejewski, Kamil Adamczewski, Maciej Pióro, Michal Krutul,  
337 Szymon Antoniak, Kamil Ciebiera, Krystian Król, Tomasz Odrzygózd, Piotr Sankowski,  
338 Marek Cygan, and Sebastian Jaszczur. Scaling laws for fine-grained mixture of experts.  
339 *CoRR*, abs/2402.07871, 2024. URL <https://doi.org/10.48550/arXiv.2402.07871>.

340 Edmund Lau, Zach Furman, George Wang, Daniel Mufet, and Susan Wei. The local  
341 learning coefficient: A singularity-aware complexity measure. In Yingzhen Li, Stephan  
342 Mandt, Shipra Agrawal, and Emtiyaz Khan (eds.), *Proceedings of The 28th International  
343 Conference on Artificial Intelligence and Statistics*, volume 258 of *Proceedings of Machine  
344 Learning Research*, pp. 244–252. PMLR, 03–05 May 2025. URL [https://proceedings.  
345 mlr.press/v258/lau25a.html](https://proceedings.mlr.press/v258/lau25a.html).

346 Mike Lewis, Shruti Bhosale, Tim Dettmers, Naman Goyal, and Luke Zettlemoyer. Base layers:  
347 Simplifying training of large, sparse models. In *International Conference on Machine  
348 Learning*, pp. 6265–6274. PMLR, 2021.

349 Johnny Lin. Llama 3.3 70b, temporal feature analysis, featured research, gpt-oss-20b, and \*a  
350 lot more\*, November 2025. URL <https://www.neuronpedia.org/blog/fall-update>.

351 Linda Linsefors and Lucius Bushnaq. Circuits in superposition 2: Now with less  
352 wrong math, 2025a. URL [https://www.lesswrong.com/posts/FWkZYQceEzL84tNej/  
353 circuits-in-superposition-2-now-with-less-wrong-math](https://www.lesswrong.com/posts/FWkZYQceEzL84tNej/circuits-in-superposition-2-now-with-less-wrong-math).

354 Linda Linsefors and Lucius Bushnaq. Rotations in superposition, 2025b. URL <https://www.lesswrong.com/posts/LZ7YMPJueB6qjL24n/rotations-in-superposition>.

355 Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao,  
356 Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv  
357 preprint arXiv:2412.19437*, 2024.

358 Ekdeep Singh Lubana, Can Rager, Sai Sumedh R Hindupur, Valerie Costa, Greta Tuckute,  
359 Oam Patel, Sonia Krishna Murthy, Thomas Fel, Daniel Wurgaft, Eric J Bigelow, et al.  
360 Priors in time: Missing inductive biases for language model interpretability. *arXiv preprint  
361 arXiv:2511.01836*, 2025.

362 Alireza Makhzani and Brendan Frey. K-sparse autoencoders. *arXiv preprint arXiv:1312.5663*,  
363 2013.

364 Samuel Marks, Can Rager, Eric J Michaud, Yonatan Belinkov, David Bau, and Aaron  
365 Mueller. Sparse feature circuits: Discovering and editing interpretable causal graphs in  
366 language models. In *The Thirteenth International Conference on Learning Representations*,  
367 2025. URL <https://openreview.net/forum?id=I4e82CIDxv>.

368 Callum McDougall, Avery, and Lucius Bushnaq. Theories of modularity in the biological  
369 literature, 2022. URL [https://www.lesswrong.com/posts/JzTfKrgC7Lfz3zcwM/  
370 theories-of-modularity-in-the-biological-literature](https://www.lesswrong.com/posts/JzTfKrgC7Lfz3zcwM/theories-of-modularity-in-the-biological-literature).

371 Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word  
372 representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.

---

378 Niklas Muennighoff, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Jacob Morrison, Sewon  
379 Min, Weijia Shi, Evan Pete Walsh, Oyvind Tafford, Nathan Lambert, Yuling Gu, Shane  
380 Arora, Akshita Bhagia, Dustin Schwenk, David Wadden, Alexander Wettig, Binyuan Hui,  
381 Tim Dettmers, Douwe Kiela, Ali Farhadi, Noah A. Smith, Pang Wei Koh, Amanpreet  
382 Singh, and Hannaneh Hajishirzi. OLMoe: Open mixture-of-experts language models.  
383 In *The Thirteenth International Conference on Learning Representations*, 2025. URL  
384 <https://openreview.net/forum?id=xXTkbTBmqq>.

385 Neel Nanda, Andrew Lee, and Martin Wattenberg. Emergent linear representations in  
386 world models of self-supervised sequence models. In *Proceedings of the 6th BlackboxNLP  
387 Workshop: Analyzing and Interpreting Neural Networks for NLP*, pp. 16–30, 2023.

388 Adam Newgas. Compressed computation: Dense circuits in a toy model of the universal-  
389 AND problem. In *Mechanistic Interpretability Workshop at NeurIPS 2025*, 2025. URL  
390 <https://openreview.net/forum?id=J3Kds2Rxov>.

391 Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan  
392 Carter. Zoom in: An introduction to circuits. *Distill*, 2020. doi:10.23915/distill.00024.001.  
393 <https://distill.pub/2020/circuits/zoom-in>.

394 Kiho Park, Yo Joong Choe, and Victor Veitch. The linear representation hypothesis and the  
395 geometry of large language models. In *Forty-first International Conference on Machine  
396 Learning*, 2024. URL <https://openreview.net/forum?id=UGpGkLzwpP>.

397 Gonalo Santos Paulo, Alex Troy Mullen, Caden Juang, and Nora Belrose. Automatically  
398 interpreting millions of features in large language models. In *Forty-second Interna-  
399 tional Conference on Machine Learning*, 2025. URL [https://openreview.net/forum?  
400 id=EemtbbhJ0Xc](https://openreview.net/forum?id=EemtbbhJ0Xc).

401 Guilherme Penedo, Hynek Kydlıek, Loubna Ben allal, Anton Lozhkov, Margaret Mitchell,  
402 Colin Raffel, Leandro Von Werra, and Thomas Wolf. The fineweb datasets: Decanting  
403 the web for the finest text data at scale. In *The Thirty-eight Conference on Neural  
404 Information Processing Systems Datasets and Benchmarks Track*, 2024. URL [https:  
405 //openreview.net/forum?id=n6SCKn2QaG](https://openreview.net/forum?id=n6SCKn2QaG).

406 Jonas Pfeiffer, Sebastian Ruder, Ivan Vuli, and Edoardo Maria Ponti. Modular deep learning.  
407 *arXiv preprint arXiv:2302.11529*, 2023.

408 David Raposo, Sam Ritter, Blake Richards, Timothy Lillicrap, Peter Conway Humphreys, and  
409 Adam Santoro. Mixture-of-depths: Dynamically allocating compute in transformer-based  
410 language models. *arXiv preprint arXiv:2404.02258*, 2024.

411 Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese  
412 bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural  
413 Language Processing*. Association for Computational Linguistics, 11 2019. URL [https:  
414 //arxiv.org/abs/1908.10084](https://arxiv.org/abs/1908.10084).

415 L Sharkey. Sparsify: A mechanistic interpretability research agenda. In *AI Alignment Forum*,  
416 2024.

417 Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton,  
418 and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts  
419 layer. *arXiv preprint arXiv:1701.06538*, 2017a.

420 Noam Shazeer, \*Azalia Mirhoseini, \*Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey  
421 Hinton, and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-  
422 of-experts layer. In *International Conference on Learning Representations*, 2017b. URL  
423 <https://openreview.net/forum?id=B1ckMDq1g>.

424 Lewis Smith. The ‘strong’feature hypothesis could be wrong. In *AI Alignment Forum*, 2024.

425 Curt Tigges. Thread on 2025 goodfire moe interpretability hackathon. [https://x.com/  
426 CurtTigges/status/1953877787552755890](https://x.com/CurtTigges/status/1953877787552755890), Aug 2025. X (formerly Twitter) thread by  
427 @CurtTigges; accessed 2025-12-17.

---

432 George Wang, Jesse Hoogland, Stan van Wingerden, Zach Furman, and Daniel Mufet.  
433 Differentiation and specialization of attention heads via the refined local learning coefficient.  
434 In *The Thirteenth International Conference on Learning Representations*, 2025. URL  
435 <https://openreview.net/forum?id=SUc1U0Wndp>.  
436

437 Martin Wattenberg and Fernanda Viégas. Relational composition in neural networks: A  
438 survey and call to action. In *ICML 2024 Workshop on Mechanistic Interpretability*, 2024.  
439 URL <https://openreview.net/forum?id=zzCEiUIPk9>.

440 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu,  
441 Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint*  
442 *arXiv:2505.09388*, 2025a. URL <https://arxiv.org/abs/2505.09388>.  
443

444 Liu Yang, Kangwook Lee, Robert D Nowak, and Dimitris Papailiopoulos. Looped transformers  
445 are better at learning learning algorithms. In *The Twelfth International Conference on*  
446 *Learning Representations*, 2024. URL <https://openreview.net/forum?id=HHbRxoDTxE>.

447 Xingyi Yang, Constantin Venhoff, Ashkan Khakzar, Christian Schroeder de Witt, Puneet K.  
448 Dokania, Adel Bibi, and Philip Torr. Mixture of experts made intrinsically interpretable.  
449 In *Forty-second International Conference on Machine Learning*, 2025b. URL <https://openreview.net/forum?id=6QERrXMLP2>.  
450

451 Barret Zoph, Irwan Bello, Sameer Kumar, Nan Du, Yanping Huang, Jeff Dean, Noam  
452 Shazeer, and William Fedus. St-moe: Designing stable and transferable sparse expert  
453 models. *arXiv preprint arXiv:2202.08906*, 2022.  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

---

## 486 A BACKGROUND

### 487 A.1 SPARSE AUTOENCODERS (SAE)

488 A fundamental challenge in interpreting LLM activations is *superposition*: models encode  
489 more features than the number of available dimensions, resulting in *polysemantic neurons*  
490 that activate for multiple unrelated concepts (Elhage et al., 2022). Sparse Autoencoders  
491 (SAEs) address this by mapping an activation vector  $\mathbf{x} \in \mathbb{R}^N$  to a corresponding sparse  
492 latent representation  $\mathbf{z} \in \mathbb{R}^F$ .<sup>2</sup> Ideally,  $\mathbf{z}$  should represent an “unpacking” of the compressed  
493 representations into monosemantic, single-concept latents referred to as *features* (Bricken  
494 et al., 2023; Cunningham et al., 2024).

495 The SAE architecture consists of an encoder that projects activations into a sparse latent  
496 space, and a decoder that reconstructs the original activation from the sparse latents:

$$497 \mathbf{z} = \sigma_s(\mathbf{W}_{\text{enc}}(\mathbf{x} - \mathbf{b}_{\text{pre}}) + \mathbf{b}_{\text{enc}}) \quad (2)$$

$$498 \hat{\mathbf{x}} = \mathbf{W}_{\text{dec}}\mathbf{z} + \mathbf{b}_{\text{pre}} \quad (3)$$

499 where  $\mathbf{W}_{\text{enc}} \in \mathbb{R}^{N \times F}$ ,  $\mathbf{W}_{\text{dec}} \in \mathbb{R}^{F \times N}$ , and  $\sigma_s$  is a sparsifying activation function. We follow  
500 Gao et al. (2025)’s **Top-K SAE** approach, where  $\sigma_s = \text{TopK}(\cdot, k)$  retains only the  $k$  largest  
501 activations per input and zeros out the rest, ensuring that the latent is sparse. The SAE is  
502 trained to optimise reconstruction fidelity  $\mathcal{L} = \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2$ , ensuring that the original activation  
503 can be accurately reconstructed from the sparse latent representation  $\mathbf{z}$ .

### 504 A.2 MIXTURE OF EXPERTS (MoE)

505 In a standard MoE layer, the dense feed-forward network (FFN) is replaced by  $E$  parallel  
506 expert networks  $\{\mathbf{E}_i\}_{i=1}^E$ , each typically an FFN itself. A learned router network ( $\mathbf{W}_r$ )  
507 determines which experts process each token. Given an input token representation  $\mathbf{x}$ , the  
508 router computes routing logits  $h(\mathbf{x}) = \mathbf{W}_r \cdot \mathbf{x}$  and converts them to a probability distribution  
509 over experts via softmax:

$$510 p_i(\mathbf{x}) = \frac{\exp(h(\mathbf{x})_i)}{\sum_{j=1}^E \exp(h(\mathbf{x})_j)} \quad (4)$$

511 To enforce sparsity, only the top- $k$  experts with the highest probabilities are selected. Letting  
512  $\mathcal{T}$  denote the set of selected expert indices, the layer output is computed as the weighted  
513 sum of expert outputs:

$$514 \mathbf{y} = \sum_{i \in \mathcal{T}} p_i(\mathbf{x}) \cdot \mathbf{E}_i(\mathbf{x}) \quad (5)$$

515 A key challenge in MoE training is *load balancing*: without intervention, models tend to  
516 collapse to using only a few experts, leaving others undertrained (Shazeer et al., 2017b).  
517 This is addressed through auxiliary losses that encourage uniform expert utilization.

## 518 B RELATED WORK

### 519 B.1 SPARSE FEATURE DECOMPOSITIONS FOR MODEL INTERPRETABILITY

520 SAEs rely on the Linear Representation Hypothesis (Mikolov et al., 2013; Bricken et al.,  
521 2023; Olah et al., 2020; Nanda et al., 2023; Park et al., 2024): the hypothesis that features  
522 in LMs are represented linearly (Bereska & Gavves, 2024). The LRH is often a controversial  
523 assumption within interpretability (Smith, 2024; Csordás et al., 2024). However, in our case  
524 of interpreting expert routing, routing is usually performed by a linear map and so only  
525 linearly accessible information can be important for routing. This makes SAEs uniquely well  
526 suited to our problem<sup>3</sup>.

---

527 <sup>2</sup>Typically  $F > N$ , allowing the SAE to represent more features than the dimensionality of the  
528 original activation vector  $\mathbf{x}$ .

529 <sup>3</sup>Also note that the Tuned Lens (Belrose et al., 2023) also validates that linear probes on  
530 intermediate residual stream states capture meaningful structure.

---

540 SAEs have drawbacks, however. SAEs can struggle to capture logical hierarchical structure  
541 (Costa et al., 2025); they do not efficiently capture structures of multi-dimensional features  
542 (Engels et al., 2025); they do not generally capture relations across multiple tokens (Lubana  
543 et al., 2025); and the structure of Relational Composition (Wattenberg & Viégas, 2024)  
544 remains difficult for SAEs.

545 Most work interpreting representations to date has focused on the impact of intermediate  
546 representations on the output logits of a forward pass (Bricken et al., 2023) or on representa-  
547 tions in a subsequent layer (Marks et al., 2025). We instead focus on how representations  
548 impact the behaviour of a subsequent routing layer.

549 Inspired by the apparent success of SAE-based explanations, Sharkey (2024) describes a  
550 framework for the continual improvement of MI explanations. The framework’s three stages  
551 are: (1) *mathematical description* (breaking down the neural network into functional parts  
552 <sup>4</sup>), (2) *semantic description* (labelling each functional part in a way that is understandable  
553 to humans <sup>5</sup>) and (3) *validation* (using the semantic description to make predictions about  
554 model behaviour and evaluating these predictions). Here our combination of the trained  
555 SAEs and our Ansatz (conjecture) that router weights are linear combinations of weights  
556 provide our mathematical description of routing. We leverage AutoInterp (Paulo et al., 2025)  
557 to aid with semantic description and evaluation of the functional parts and the empirical  
558 success of this approach provides us with evidence for the Ansatz.

## 560 B.2 MIXTURE OF EXPERTS INTERPRETABILITY

562 Shazeer et al. (2017b) designed the modern Sparse MoE layer to improve the efficiency  
563 and scalability of very large neural networks. Other researchers, however, hoped that the  
564 specialisation routing provided might also extend some interpretability benefits (Jacobs  
565 et al., 1991; Fedus et al., 2022a). However, Jiang et al. (2024) find it difficult to obtain clean  
566 interpretability results with their MoE model Mixtral, remarking “Surprisingly, we do not  
567 observe obvious patterns in the assignment of experts based on the topic [of the input text].”  
568 Lewis et al. (2021) and Zoph et al. (2022) analyse the unigram patterns of routing and find  
569 that the observed level of specialisation varies dramatically making interpretability difficult.  
570 Tigges (2025) report some success with unigram analysis of a few experts such as a ‘business’  
571 expert. We show, however, that by using SAE features as a basis for explanation rather than  
572 tokens, we are able to achieve more accurate and concise explanations.

573 Yang et al. (2025b) suggest that with very large numbers of experts, routing can be made  
574 somewhat interpretable. However, their setting requires more experts than is compute  
575 optimal given scaling laws for sparse models (Abnar et al., 2025; Krajewski et al., 2024) or  
576 more experts than is typical in high performing open source MoE models (Liu et al., 2024;  
577 Yang et al., 2025a; Agarwal et al., 2025; Muennighoff et al., 2025).

578 Chaudhari et al. (2025) provide an alternative model for understanding the phenomena  
579 of superposition (Elhage et al., 2022) in MoE models. Chaudhari et al. (2025) find that  
580 individual experts exhibit greater monosemanticity than equivalent dense models (in a toy  
581 setting). Crucially, their analysis measures monosemanticity at the level of *expert weight*  
582 *matrices*, not at the level of *routing decisions*, which we focus on. We argue that even if  
583 individual experts represent their assigned features monosemantically, the routing decision  
584 itself may be polysemantic. The router partitions the input space such that features assigned  
585 to the same expert only compete with each other for representational capacity (Chaudhari  
586 et al., 2025). This creates an incentive to route *dissimilar* domains to the same expert:  
587 micro-domains that rarely co-occur can share an expert with little interference, enabling the  
588 expert to effectively perform Computation in Superposition (Hänni et al., 2024).

---

590 <sup>4</sup>where here the parts could be in terms of representations (Bricken et al., 2023) or computations  
591 (Braun et al., 2025)

592 <sup>5</sup>note that implicitly this formulation requires the Independent Additivity Principle (Ayonrinde  
593 et al., 2024) - if multiple parts are relevant at one time, then it must be the case that we can  
understand the whole simply from understanding the constituent parts

---

### B.3 ADAPTIVE COMPUTATION

Sparse MoE models are a type of Adaptive Computation model (Graves, 2016; Ayonrinde, 2023)<sup>6</sup>. Adaptive Computation models are typically either more FLOP or parameter efficient than dense models as they either use a subset of their parameters (for example sparse MoE models or Early Exit models (Banino et al., 2021)) or reuse parameters (for example looping models (Dehghani et al., 2018; Yang et al., 2024)) respectively. There has been little work on the interpretability of Adaptive Computation models and while our work focuses on interpreting routing in the sparse MoE layer, we would be excited about work generalising this approach to routing in early-exit models like Mixture of Depths (Raposo et al., 2024).

### B.4 MODULARITY IN NEURAL NETWORKS

Many brain-inspired neural network architectures employ modularity (the organisation of a system into functional, sparsely connected subunits) as a core component (Pfeiffer et al., 2023), as modularity is believed to be one of the key properties that makes brains efficient and effective (Clune et al., 2013).

We might hypothesise that because neural networks generalise so well and have much lower effective dimensionality than the model dimension would suggest (Lau et al., 2025), we should expect the models to contain intrinsic modularity. That is to say that even though NNs look fully connected and highly entangled, perhaps there is some way of viewing computation such that the computation is in fact highly modularised (Bushnaq et al., 2024). However, efforts to find such modularity have proven difficult (Filan et al., 2021), partially because it is not clear in what form we should expect such modularity to appear.

One way to understand the modularity that neural networks naturally and implicitly form, is to enforce some modular computation and see how the neural networks react. Sparse MoE models have this kind of enforced modularity and we hope that in studying these models we can develop better tools for uncovering possible latent modularity in dense neural networks. In particular, one promising path for understanding modularity may be through the optimisation pressures that resulted in such modularity. For example specialisation and reducing connection costs are two evolutionary pressures for modularity to develop in biological models (Clune et al., 2013; McDougall et al., 2022). Analogously, we might expect that specialisation pressures in the training processes for AI systems also results in the development of modularity as in Wang et al. (2025).

## C SAE FEATURES PREDICT MOE ROUTING

Beside generating explanations, we also evaluate whether SAE features can linearly predict routing patterns. We extract activations  $\mathbf{x}$  from the residual stream before the routing function and use an SAE encoder to produce sparse feature representations  $\mathbf{z}$ . We then train a logistic regression classifier on these sparse feature representations, one per expert  $E_i$ , to predict whether a given expert is activated for some context. We compare to n-gram baselines: for each token in the vocabulary, we count how often it co-occurs with each expert’s activation on a training set. At inference, we predict that a token routes to the experts it most frequently co-occurred with.

We report *Recall* as our main evaluation metric, measuring the proportion of experts that are correctly predicted as being routed to.<sup>7</sup> For a token  $t$  with actual expert set  $E_{\text{actual}}(t)$  and predicted expert set  $E_{\text{pred}}(t)$ , both of size  $k$ :

$$\text{Recall} = \frac{1}{N} \sum_{t=1}^N \frac{|E_{\text{pred}}(t) \cap E_{\text{actual}}(t)|}{|E_{\text{actual}}(t)|} \quad (6)$$

where  $N$  is the total number of tokens.

<sup>6</sup>Also known as Conditional or Dynamic Computation models (Han et al., 2021)

<sup>7</sup>Since both the predicted set  $E_{\text{pred}}(t)$  and actual set  $E_{\text{actual}}(t)$  have cardinality  $k$  by construction (due to top- $k$  routing), Recall equals Precision and F1 in this setting (equal set sizes imply  $\text{FP} = k - |E_{\text{pred}} \cap E_{\text{actual}}| = \text{FN}$ ).

648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

	OLMoE-1B-7B	gpt-oss-20b
Unigram Baseline	0.70	0.52
Bigram Baseline	0.74	0.55
SAE Predictor (Ours)	<b>0.89</b>	<b>0.81</b>

Table 2: Token-level statistics are insufficient to predict MoE routing: richer representations are required. N-gram baselines, which predict routing from token co-occurrence frequencies, achieve only 0.52–0.74 Recall, while our SAE-based predictor achieves 0.81–0.89 by leveraging learned feature representations. The performance gap widens on the larger model (gpt-oss-20b), suggesting routing becomes increasingly context-dependent at scale. This motivates using SAE features, rather than tokens, as the basis for interpreting expert routing. Results reported as Recall for layer 11 of OLMoE-1B-7B and layer 16 of gpt-oss-20b.

	Layer 4	Layer 8	Layer 12	Layer 16	Layer 20
Unigram Baseline	0.57	0.60	0.52	0.57	0.55
Bigram Baseline	0.63	0.64	0.56	0.59	0.59
SAE Predictor (Ours)	<b>0.82</b>	<b>0.85</b>	<b>0.79</b>	<b>0.82</b>	<b>0.83</b>

Table 3: The advantage of SAE-based prediction over token co-occurrence frequencies holds across all network depths. SAE-based prediction achieves 0.79–0.85 Recall across layers 4–20 of gpt-oss-20b, while n-gram baselines plateau at 0.52–0.64. The narrow range of SAE predictor performance suggests routing complexity is similar across layers.

We evaluate whether SAE latents can predict MoE routing decisions across different model architectures. Table 2 compares our SAE-based predictor against n-gram baselines on both OLMoE-1B-7B and gpt-oss-20b. The SAE predictor achieves 0.89 and 0.81 Recall respectively, substantially outperforming both baselines across both architectures. N-gram baselines perform markedly worse on the larger model (0.74 for OLMoE vs 0.55 for gpt-oss), suggesting that routing decisions in larger models are more context-dependent and cannot be predicted from surface-level token statistics alone. In contrast, the SAE predictor maintains strong performance across both architectures, demonstrating that routing information is consistently encoded in SAE features regardless of model size, architecture, or training distribution.

Table 3 shows that the advantage of SAE-based prediction holds across all layers of gpt-oss-20b. The consistent advantage at both early and late layers indicates that routing is better explained by SAE latents than by token statistics, supporting our hypothesis that routing operates at the feature level. We provide a detailed ablation on the number of active SAE latents  $m$  needed to outperform n-gram baselines in Section D.

## D ROUTING DEPENDS ON MULTIPLE FEATURES

In Section C, we reported routing prediction results using  $m = 32$  active features. Here we ablate the effect of the number of active features  $m$  on routing prediction Recall. Table 4 shows Recall for varying  $m \in \{1, 2, 4, 8, 16, 32, 64\}$  across five layers of gpt-oss-20b.

The results reveal two key patterns. First, Recall continues to improve with more features across all layers, indicating that routing decisions depend on multiple distinct concepts—consistent with polysemantic expert behaviour. Second, even a single feature ( $m = 1$ ) outperforms n-gram baselines at most layers (underlined), demonstrating that individual SAE features carry substantial routing information. This aligns with our finding that  $\rho$ -useful features achieve strong explanation performance with just 1–2 features (Appendix F).

	Layer 4	Layer 8	Layer 12	Layer 16	Layer 20
Unigram Baseline	0.570	0.602	0.523	0.573	0.553
Bigram Baseline	0.631	0.639	0.560	0.593	0.589
SAE Predictor ( $m = 1$ )	<u>0.690</u>	0.462	<u>0.597</u>	<u>0.687</u>	<u>0.748</u>
SAE Predictor ( $m = 4$ )	0.790	<u>0.809</u>	0.754	0.811	0.847
SAE Predictor ( $m = 8$ )	0.803	0.824	<b>0.811</b>	0.800	0.844
SAE Predictor ( $m = 16$ )	0.817	0.830	0.810	<b>0.830</b>	0.861
SAE Predictor ( $m = 32$ )	<b>0.824</b>	<b>0.865</b>	0.801	0.820	0.866
SAE Predictor ( $m = 64$ )	0.822	0.858	0.793	0.825	<b>0.873</b>

Table 4: Routing prediction improves with more SAE features, consistent with polysemantic expert behaviour predicted by the SSH. Underlined values outperform n-gram baselines; **bold** marks best Recall with fewest features. Notably, even a single feature ( $m = 1$ ) outperforms n-gram baselines at 4 of 5 layers, while performance continues improving up to  $m = 32$ – $64$ , indicating that routing depends on multiple distinct concepts rather than a single domain.

## E EXPERIMENTAL SETUP

**Models** We evaluate on two MoE architectures: OLMoE-1B-7B (Muennighoff et al., 2025) and gpt-oss-20b (Agarwal et al., 2025).

**Datasets.** For experiments with OLMoE, we use OLMoE-mix-0924 (Muennighoff et al., 2025) for SAE training and routing prediction experiments. For gpt-oss, we use a following mixture: 75% gpt-oss-generated continuations from The Pile (Gao et al., 2020) prompts<sup>8</sup> and 25% FineWeb (Penedo et al., 2024). We use The Pile (Gao et al., 2020) as a dataset for explanation scoring.

**SAEs.** For OLMoE, we train Top-K SAEs for 100M tokens on activations sampled from layers 3, 7, 11, and 15 (out of OLMoE’s 16 layers). Our SAEs have  $k = 32$  and 32,768 features (an expansion of 16x compared to the residual stream size). For gpt-oss (24 layers), we use trained BatchTopK SAEs<sup>9</sup> (Lin, 2025) on layers 4, 8, 12, 16, and 20, with 131,072 features and  $k = 64$ .

**Routing Prediction.** Linear classifiers are trained on SAE features from 1M tokens. We vary the number of active features  $m$  to assess how routing information distributes across features. Specifically, for OLMoE we use  $m \in \{1, 2, 4, 8, 16, 32\}$ , while for gpt-oss we use  $m \in \{1, 2, 4, 8, 16, 32, 64\}$ .

We compare our SAE-based classifiers to unigram and bigram models that predict routing based on  $n$ -gram frequencies in training data. Concretely, for each  $n$ -gram (unigram or bigram) in the train set, we collect the number of times it activates each expert. We then predict the activated experts by taking the top- $k$  experts with the highest activation frequency for a given  $n$ -gram.

**Expert Interpretation.** For gpt-oss, we generate explanations across layers using the top-10  $\rho$ -useful features per expert. We use the Delphi library (Paulo et al., 2025) with Claude Haiku 4.5 (Anthropic, 2025) for explanation generation, summarisation, and calculating explanation scores. We use the all-MiniLM-L6-v2 model from Sentence-Transformers library (Reimers & Gurevych, 2019), for calculating embeddings for samples in the evaluation set for explanation scoring,

<sup>8</sup><https://huggingface.co/datasets/andyrdt/gpt-oss-20b-rollouts>

<sup>9</sup><https://huggingface.co/andyrdt/saes-gpt-oss-20b>

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

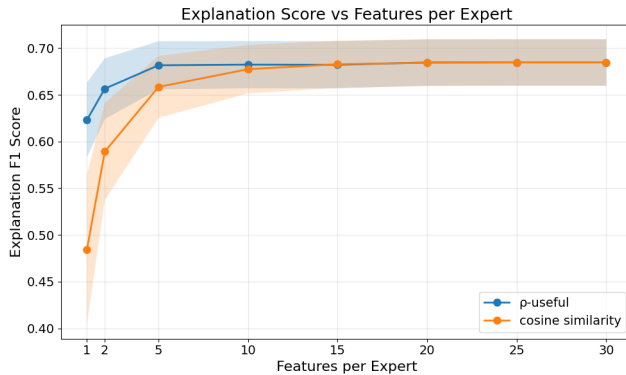


Figure 3: The most useful features (ranked by the  $\rho$ -usefulness metric) are individually *more* predictive of expert routing than features with highest cosine similarity with router weights. Despite selecting largely non-overlapping feature sets, both methods converge to similar performance at 10 features, suggesting routing information is redundantly encoded across many SAE features. Results are reported for layer 16 of gpt-oss-20b. Shaded regions show standard deviation across experts.

## F FEATURE SELECTION METHOD AND SET SIZE

In Section 3, we describe selecting the top- $n$  features by  $\rho$ -usefulness for each expert. Here we justify this choice by evaluating (1) how the feature set size  $n$  affects explanation quality, and (2) how  $\rho$ -usefulness compares to an alternative selection method based on cosine similarity.

As an alternative to  $\rho$ -usefulness, we consider selecting features based on the geometric alignment between SAE decoder directions  $\mathbf{d}_f \in W_{dec}$  and router weight vectors  $\mathbf{w}_k \in \mathbf{W}_r$ . For each expert  $E_k$ , we rank features  $f$  by  $\cos(\mathbf{d}_f, \mathbf{w}_k)$  and select the top- $n$ .

To isolate the effect of feature selection, we replace the LLM summarisation with a simpler aggregation: we evaluate each feature’s explanation independently and label an example as positive if *at least one* feature’s explanation predicts expert activation. This setup directly measures whether adding more features increases predictive coverage: as the feature set grows, we track whether additional features contribute new predictive information.

Figure 3 shows that individual  $\rho$ -useful features carry strong predictive signal— even a single feature achieves  $F1 \approx 0.62$ . In contrast, features with high cosine similarity to router weights are only predictive in aggregate ( $F1 \approx 0.49$  at  $n = 1$ ), suggesting that geometric alignment with router weights does not guarantee a feature fires when that expert is selected. Both methods converge to approximately the same explanation score ( $\approx 0.68$ ) at 10 features.

The convergence at 10 features is interesting given that the two criteria select largely non-overlapping feature sets. This suggests that routing-relevant information is encoded across many SAE features, and with enough features, both methods eventually capture sufficient signal to predict routing.

Based on these results, we use  $n = 10$  features per expert throughout our experiments.

## G EXAMPLE EXPERT EXPLANATIONS

Table 5 shows different levels of RouterInterp explanations for a randomly selected expert (Expert 8) from layer 16 of gpt-oss-20b. It illustrates how expert-level explanations aggregate multiple feature-level explanations, which are in turn acquired by explaining examples where the feature activates.

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

Level	Explanation
<b>Expert 8</b>	This expert activates across five distinct contexts: (1) Structural and grammatical elements—function words, articles, prepositions, conjunctions, and punctuation that organize text and parse grammatical relationships. (2) Informational and technical content—medical/clinical terminology, dates, measurements, quantitative data, multilingual technical terms, and metadata markers, with stronger activation in content-dense, semantically specific text. (3) Logical and reasoning anchors—punctuation marking boundaries, conjunctions signaling relationships, mathematical operators, and discourse markers in technical or argumentative contexts. (4) Self-identification phrase—high activation on tokens within "a large language model trained by OpenAI" in system messages. (5) Mathematical and technical notation—numerical values, operators, variables, LaTeX notation, and domain-specific keywords encoding mathematical constraints and relationships.
<i>Feature 46328</i>	<i>Self-identification phrase—tokens within "a large language model trained by OpenAI" in system messages, particularly "large," "model," "trained," and "AI."</i> <b>Activating Examples:</b> [1] Chat GPT , a large language model trained by Open AI .
<i>Feature 70455</i>	<i>Logical and reasoning anchors—punctuation marking boundaries, conjunctions signaling relationships, mathematical operators, and discourse markers in technical or argumentative contexts.</i> <b>Activating Examples:</b> [1] a GUI ( like Chess Base , Arena , Stockfish 's own GUI , or
<i>Feature 73639</i>	<i>Structural and grammatical elements—articles, prepositions, conjunctions, punctuation, and mathematical delimiters that organize and parse text.</i> <b>Activating Examples:</b> [1] and \$ y \$ . Let \$ n \$ be the number of possible values of [2] s 393 56 ash Maybe it's best to ignore
<i>Feature 81948</i>	<i>Informational and technical content—medical/clinical terms, dates, measurements, quantitative data, multilingual technical terms, and metadata markers. Higher activation in content-dense, semantically specific text.</i> <b>Activating Examples:</b> [1] though he termed it "speculative ." He acknowledged [2] I 've never been through anything like this before .
<i>Feature 91358</i>	<i>Mathematical and technical notation—numerical values, operators, variables, LaTeX notation, and domain-specific keywords encoding mathematical constraints and relationships.</i> <b>Activating Examples:</b> [1] But it's given that he secured 110 marks . So 170 - T [2] $2 - y - 1$ } + \{ \frac { y - 2 ^ 2 } { }

Table 5: RouterInterp first explains features by looking at examples where these features activate strongly, and then aggregates these explanations into an expert-level explanation. For example, Expert 8 from gpt-oss-20b layer 16 specialises in AI self-identification (*Feature 46328*), medical terms (*Feature 81948*), mathematical notation (*Feature 91358* and *Feature 73639*), and reasoning anchors (*Feature 70455*). Here we show explanations for the top 5  $\rho$ -useful SAE features (rather than 10) for brevity.