

---

# Probe Generalization as Subspace Selection for OOD Deception Detection

---

Anonymous Authors<sup>1</sup>

## Abstract

Linear probes can be used to detect behaviors and concepts inside language model activations, but may fail to transfer to out-of-distribution examples. When studying the generalization performance of Llama-3.1-8B-Instruct probes over 3 held-out deception detection datasets, we find that projecting inputs onto a small subset of principal components (PCs) from the training distribution of activations enables cross-domain transfer that nearly matches the performance of probes trained directly on the test distribution. Furthermore, we find that PC interpretations can be used to find a subset of those transferable PCs. By using an LLM judge to score each PC on whether its most/least activating examples imply a transferable deception direction, then probing on the highest-scoring PCs, we close the baseline-to-oracle gap by 78% on Insider Trading Report and by 25% on Sandbagging. The directions a source probe weights heavily appear to encode source-specific surface features, while the directions that actually transfer appear to encode the same contrast more abstractly, in a way natural language descriptions can capture. Broadly, our results suggest that the OOD robustness of probes is largely determined by subspace selection.

## 1. Introduction

Linear probes are standard tooling for reading information from language model activations (Alain & Bengio, 2018; Hewitt & Liang, 2019). Linear probes are cheap to train and expose some semantically meaningful information in the otherwise uninterpretable representation space. However, probes that look strong in-distribution are sometimes unreliable out-of-distribution (Kirch et al., 2026). Existing work shows that key behavioral concepts are linearly represented

---

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

in model activations (Marks & Tegmark, 2024; Gurnee & Tegmark, 2024; Park et al., 2024; Zou et al., 2025), so probes should have the capacity to detect these concepts, yet they appear to learn spurious correlations in practice (Orgad et al., 2025; Wang et al., 2025). This is unacceptable for probe-based monitoring of harmful interactions, setting a clear challenge: *how can we build probes that only learn generalizing features?*

To tackle this, we work in the basis of source-domain principal components (PCs) rather than raw activations, avoiding the assumption that the model’s latent dimensions provide a privileged basis. Principal components offer a simple, data-driven way to study structure in representation spaces, and prior work even suggests that the choice of PCs matters: Mu et al. (2018) shows that removing the top PCs of an embedding space improves downstream performance, and Kantamneni et al. (2025) use top-PC projections as a probe baseline. Thus, the focus of our investigation becomes: *which subset of principal components supports generalizable probe training and how can that subset be identified?*

We use deception probes on Llama-3.1-8B-Instruct (Grattafiori et al., 2024) as our testbed, training them on a single source dataset (Roleplaying deception scenarios) and evaluating on 3 held-out target datasets (Insider Trading Report, Insider Trading Confirm, and Sandbagging) (Goldowsky-Dill et al., 2025). Working in the source-domain PCA basis, we explore strategies for selecting a subset of training dimensions that avoid spurious correlations and transfer to target domains.

This work contributes the following findings:

- **A small transferable subspace exists** (Section 3.2). A target-supervised greedy search over source PCs finds 8–15 components on which a source-trained probe nearly matches a target-trained probe.
- **Source-side and target-side rankings both fail to recover the subspace** (Section 3.3). Ranking PCs by the source or target probe’s weights does not produce a basis that transfers better than the full-representation baseline.
- **An LLM judge partially recovers the subspace without target data** (Section 3.5). An LLM judge can interpret and select PCs that are likely to encode trans-

ferable behavioral contrast rather than surface patterns, enabling probes to recover most of the AUROC improvement available from training on target.

## 2. Setup

**Datasets.** Following Goldowsky-Dill et al. (2025), we use 4 of their deception detection datasets (Roleplaying, Insider Trading Report, Insider Trading Confirm, Sandbagging), with descriptions being provided in Appendix B.

Roleplaying is our source domain and the other 3 are target test sets. Each response is labeled as honest (0), ambiguous (-1), or deceptive (1) by GPT-4o (OpenAI et al., 2024), with the ambiguous responses being removed for our experiments.

**Language models.** Following Ying et al. (2026), we use layer 15 of Llama-3.1-8B-Instruct (Grattafiori et al., 2024), the layer with the highest cross-domain AUROC in their analysis. To extract activations, we pass each prompt through the model and average token activations over just the assistant’s response.

**Probe training and evaluation.** We standardize source activations to zero mean and unit variance, then apply the same source statistics to standardize all target activations. We run PCA on the scaled source activations and retain the top 100 PCs. All probes are logistic regression classifiers ( $\lambda = 10$ ); for experiments using a subset  $S$  of these 100 PCs, activations are projected onto  $S$  before training and evaluation, while the baseline probes train directly on the full standardized activations.

We use the top 100 PCs because alignment with SAE decoder features (Bricken et al., 2023) decreases with PC index (Section D), meaning lower-index PCs are more amenable to natural-language interpretation and hence comparable with interpretability-driven selection methods. This is not a claim that low-index PCs generalize better; Section 3.1 confirms the 100-PC subspace retains enough information to separate deceptive from honest responses on every target.

## 3. Experiments

### 3.1. Setting a baseline

We first test how well probes trained on the full source representations transfer to each target using source data. The results are shown in Figure 1. Insider Trading Confirm transfers reasonably well (0.874 AUROC) and Insider Trading Report performs moderately (0.706 AUROC), while Sandbagging performs poorly (0.606 AUROC).

To check that the 100-PC subspace itself is not the bottleneck, we also train a probe directly on target data within this

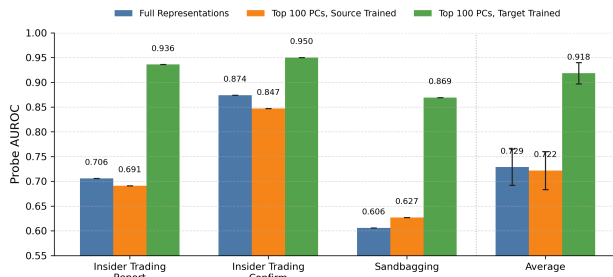


Figure 1. Probe AUROC on target datasets for source-trained probes on full representations (blue) and top-100 PCs (orange), and target-trained probes on top-100 source PCs (green, 5-fold CV). Error bars: 95% CIs

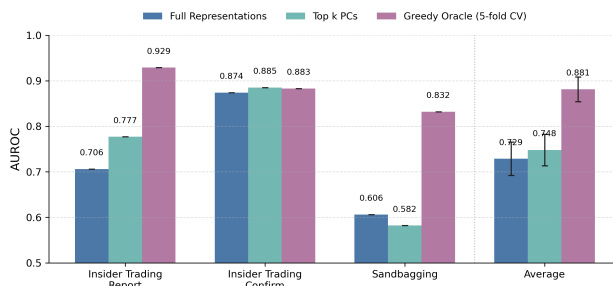


Figure 2. Probe AUROC across the three targets for the full-representation baseline (blue), top-k PCs by variance (teal, k matched to greedy), and greedy oracle PCs (purple). Error bars: 95% CIs.

subspace using 5-fold cross-validation. The probe achieves 0.936 / 0.950 / 0.869 AUROC across the three targets. These results confirm that the 100-PC subspace encodes enough information to separate deceptive from honest responses on every target.

Thus, the difficulty does not stem from the relevant information being missing, but that the source probe learns a direction in the subspace that does not generalize.

### 3.2. Showing a transferable subspace exists

*Is there any method for training probes on source data that closes the gap in performance between the baseline and the target-trained probes?* To test this, we first introduce and investigate a greedy algorithm for PC selection. Starting from an empty set, we iteratively add PCs from the source dataset. At each step, we search over the top 100 PCs and add the one that most improves AUROC on the target. To evaluate each candidate, we project source activations onto the current subset and train a logistic regression probe on those projections.

Probe weights are always fit on source data, but the choice of PC subsets is guided by performance on labeled examples from the target dataset. Using 5-fold cross-validation over

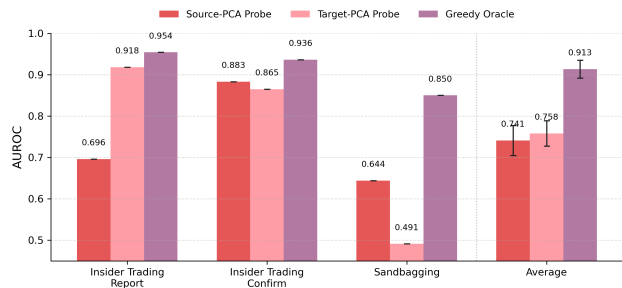


Figure 3. Probe AUROC across the three targets for three PC-selection rules: ranking by source-probe weights (red), by target-probe weights (pink), and greedy oracle selection (purple). Error bars: 95% CIs.

the target dataset, we perform greedy PC selection on the non-held-out target folds and evaluate the selected source-trained probe on the held-out target fold. We stop when no candidate PC improves AUROC by more than 0.001. The search converges on small subsets of PCs — on average, 15 PCs for Insider Trading Report, 8 for Insider Trading Confirm, and 14 for Sandbagging.

The greedy probe provides an upper bound on transfer performance, not an unbiased OOD estimate, because subset selection uses some target examples. To ground our results, we compare against simply selecting the top  $k$  PCs by variance, where  $k$  matches the number of PCs selected by the greedy algorithm. The results are shown in Figure 2.

The greedy method comes close to a probe trained directly on target dataset for Insider Trading Report and Sandbagging, improving on the baseline by 0.22 and 0.23 points of AUROC, respectively. Thus, a transferable subspace does exist in principle. The variance control method matches greedy only for Insider Trading Confirm (the easy case) and falls far behind for Insider Trading Report and Sandbagging, the harder generalization tasks. These findings suggest that the choice of PCs matters and that the amount of variance explained by those PCs does not predict their usefulness.

### 3.3. Testing hypotheses for transferable principal components

*What property of these PCs makes them transferable and can that property be detected from the source-side alone?* To test this, we investigate 2 simple hypotheses for what the greedy subset contains. The first posits that the transferable PCs are simply the ones that separate the target classes well by themselves across the test set. — greedy happens to pick them because it has access to target data. The second posits that they are the PCs the source probe already leans on most heavily — greedy is just rediscovering the source signal in a more compact form.

We compare 2 rankings of PCs, both based on the standard-

ized coefficient (the absolute value of each probe coefficient multiplied by the standard deviation of the corresponding PC) given to each PC by a probe trained on activations projected on the source-PCA basis.

The ‘source-PCA probe’ ranks PCs by the weights learned by a probe trained on source data, while the ‘target-PCA probe’ ranks PCs by the weights learned by a probe trained on target data. In both cases, the activations are projected onto the source-PCA basis. In both cases, the probes are trained on source data that has been projected onto the top  $k$  PCs, selected based to their respective rankings. If either hypothesis were correct, the corresponding probes should recover the greedy method’s performance.

To set  $k$ , we run the greedy method on the full target dataset rather than per fold, giving a single stable PC set:  $k = 13$  for Insider Trading Report, 7 for Insider Trading Confirm, and 14 for Sandbagging. Because we are using the full dataset rather than doing CV, the greedy numbers differ slightly from the CV numbers in the previous section. The results are shown in Figure 3.

Both hypotheses fail, but in different ways. The source-PCA probe transfers about as well as the full representation baseline — the PCs that the source probe relies on are those that do not generalize. The target-PCA probe does well on Insider Trading Report but performs worse than the greedy probe on Insider Trading Confirm and has below random transfer performance on Sandbagging. The PCs that best separate target classes are not the PCs that, when weighted using source data, continue to transfer the best to the target data. For further PC set comparisons, see Appendix E.

### 3.4. Interpreting the transferable subspace

So far we have shown that a small subset of source PCs supports cross-domain transfer, and that this subset is not the same as the directions a source classification or target classification relies on. *Does the transferable subset correspond to anything human-interpretable?* To test this, we adopt the interpretability method of Bills et al. (2023) and Casademunt et al. (2025).

For each of the top 100 source PCs, we find the 8 most and 8 least PC-activating examples from the source data, before prompting an LLM to generate a natural-language description of the positive pole, the negative pole, and the PC direction as a whole. The LLM judge then rates the generalizability of each direction on a 1–10 scale (OOD score), with the full judge prompt and model detail being provided in Appendix C. Since no PC received a score above 6, we set this to be the maximum OOD score.

There were 5 PCs that showed up in the greedy selections of more than one target. 3 out of 5 PCs received the maximum OOD score, suggesting that LLMs can partially recover

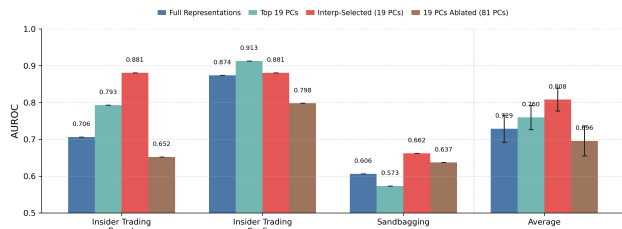


Figure 4. Probe AUROC across the three targets for the interp-selected subspace (red; 19 PCs) versus three controls: full-representation baseline (blue), top-19 PCs by variance (teal), and the remaining 81 PCs after ablating the interp-selected PCs (brown). Error bars: 95% CIs.

transferable directions. For comparison, only 1 of the top 5 PCs selected by the source-PCA probes received the maximum OOD score.

Interpretations of the greedy-selected PCs describe the contrast at a level of abstraction that survives the move to other domains. A full breakdown of PCs selected by each method is provided in Section F.

### 3.5. Selecting a subspace by interpreting PCs

*Does selecting PCs based on the interpretability-based OOD scores recover a meaningful fraction of the greedy probe performance, without using any target data?* To test this, we select the PCs which were given the maximum OOD score by the judge model and use them as the probe subspace. Because LLM judge scores are noisy, we run the scoring twice and only keep PCs that received the maximum score in both runs, giving us 19 interp-selected PCs in total. As a control, we evaluate against a probe trained on the ablation of the selected PCs (the 81 PCs not used) as well as a probe trained on the top 19 source PCs by variance. The results are shown in Figure 4.

The interp-selected PC probe outperforms the baseline on all 3 targets and outperforms the source-variance-selected control, by wide margins on Insider Trading Report and Sandbagging. As expected, the interp-selected PC probe also outperforms the interp-ablated PC probe across all 3 targets. The size of the recovery varies sharply by dataset. On Insider Trading Report, interp-selection closes 78% of the gap between baseline (0.706) and the greedy oracle (0.929). On Insider Trading Confirm, the baseline transfer performance is already near the ceiling so there is little headroom for comparisons, thus the result is uninformative — each selection method (aside from interp-ablated) does roughly as well as one another. On Sandbagging, interp-selection closes 25% of the gap between the greedy oracle (0.832).

Overall, selecting PCs purely by the natural language description of what they encode, with no use of target data,

provides a meaningful improvement over baseline performance on the generalization tasks.

## 4. Limitations

Our strongest claims are still quite narrow and there are several limitations that should be considered by the reader.

**Single model and source dataset.** All experiments use a single model (Llama-3.1-8B-Instruct, layer 15) and a single source dataset (Roleplaying (Goldowsky-Dill et al., 2025)). Moreover, Insider Trading Confirm results are uninformative, as baseline transfer on Insider Trading Confirm already reaches 0.874. Whether our results generalize to other model families or other datasets is an open question we leave to future work.

**Sandbagging margins are small.** The weakest results are for Sandbagging, where interpretability selection raises AUROC from 0.606 to 0.662, recovering only 25% of the gap to the greedy upper bound (0.832). Sandbagging is also the most out-of-distribution target, and the modest gain likely reflects this distribution shift. The result for Insider Trading Report (0.706 to 0.881, closing 78% of the gap to greedy at 0.929) is the more robust.

**No causal validation.** We make structural claims about what each PC “encodes” based on top-activating examples and an LLM judge’s verbalization of the contrast. These are descriptive claims, not causal ones. We have not ablated the high-scoring PCs during inference to verify that they carry the contrast they appear to encode.

## 5. Conclusion

In probe-based deception monitoring, a small subset of source-domain principal components (PCs) enables cross-domain transfer that nearly matches target-trained probes, where this subset cannot be found with source-side or target-side separability alone. An interpretability score over each PC’s top-activating examples, produced by an LLM judge with no access to target data, partially recovers the subspace.

A key implication is that robust probing requires identifying not only whether a behavioral signal is linearly accessible, but where in activation space that signal is represented most robustly. Linear probes trained over full representations may combine general behavioral features with dataset-specific surface cues. These heuristics can support high in-distribution accuracy while hurting transfer. In contrast, sparse subspaces of PCs may better isolate directions corresponding to more general behavioral structure. This reframes OOD probing as a representation-selection problem: the goal is not only to fit a classifier, but also choose a subspace that preserves the intended behavioral feature and filters out spurious correlation.

## Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

## References

Alain, G. and Bengio, Y. Understanding intermediate layers using linear classifier probes, 2018. URL <https://arxiv.org/abs/1610.01644>.

andyrdt. saes-llama-3.1-8b-instruct. <https://huggingface.co/andyrdt/saes-llama-3.1-8b-instruct>, 2024. Residual stream sparse autoencoders for Llama-3.1-8B-Instruct. Accessed 2026-05-04.

Arditi, A., Obeso, O., Syed, A., Paleka, D., Panickssery, N., Gurnee, W., and Nanda, N. Refusal in language models is mediated by a single direction, 2024. URL <https://arxiv.org/abs/2406.11717>.

Bills, S., Cammarata, N., Mossing, D., Tillman, H., Gao, L., Goh, G., Sutskever, I., Leike, J., Wu, J., and Saunders, W. Language models can explain neurons in language models. <https://openaipublic.blob.core.windows.net/neuron-explainer/paper/index.html>, 2023.

Bricken, T., Templeton, A., Batson, J., Chen, B., Jermyn, A., Conerly, T., Turner, N., Anil, C., Denison, C., Askell, A., Lasenby, R., Wu, Y., Kravec, S., Schiefer, N., Maxwell, T., Joseph, N., Hatfield-Dodds, Z., Tamkin, A., Nguyen, K., McLean, B., Burke, J. E., Hume, T., Carter, S., Henighan, T., and Olah, C. Towards monosemanticity: Decomposing language models with dictionary learning. *Transformer Circuits Thread*, 2023. <https://transformer-circuits.pub/2023/monosemantic-features/index.html>.

Burns, C., Ye, H., Klein, D., and Steinhardt, J. Discovering latent knowledge in language models without supervision, 2024. URL <https://arxiv.org/abs/2212.03827>.

Casademunt, H., Juang, C., Karvonen, A., Marks, S., Rajamanoharan, S., and Nanda, N. Steering out-of-distribution generalization with concept ablation fine-tuning, 2025. URL <https://arxiv.org/abs/2507.16795>.

Goldowsky-Dill, N., Chughtai, B., Heimersheim, S., and Hobbhahn, M. Detecting strategic deception using linear probes, 2025. URL <https://arxiv.org/abs/2502.03407>.

Grattafiori, A., Dubey, A., and Jauhri, A. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.

Gurnee, W. and Tegmark, M. Language models represent space and time, 2024. URL <https://arxiv.org/abs/2310.02207>.

Hewitt, J. and Liang, P. Designing and interpreting probes with control tasks, 2019. URL <https://arxiv.org/abs/1909.03368>.

Kantamneni, S., Engels, J., Rajamanoharan, S., Tegmark, M., and Nanda, N. Are sparse autoencoders useful? a case study in sparse probing, 2025. URL <https://arxiv.org/abs/2502.16681>.

Kirch, N., Dower, S., Skapars, A., Yannakoudakis, H., Lubana, E. S., and Krasheninnikov, D. The impact of off-policy training data on probe generalisation, 2026. URL <https://arxiv.org/abs/2511.17408>.

Kramár, J., Engels, J., Wang, Z., Chughtai, B., Shah, R., Nanda, N., and Conmy, A. Building production-ready probes for gemini, 2026. URL <https://arxiv.org/abs/2601.11516>.

Lu, C., Gallagher, J., Michala, J., Fish, K., and Lindsey, J. The assistant axis: Situating and stabilizing the default persona of language models, 2026. URL <https://arxiv.org/abs/2601.10387>.

Marks, S. and Tegmark, M. The geometry of truth: Emergent linear structure in large language model representations of true/false datasets, 2024. URL <https://arxiv.org/abs/2310.06824>.

Mu, J., Bhat, S., and Viswanath, P. All-but-the-top: Simple and effective postprocessing for word representations, 2018. URL <https://arxiv.org/abs/1702.01417>.

OpenAI, :, Hurst, A., Lerer, A., Goucher, A. P., Perelman, A., Ramesh, A., Clark, A., Ostrow, A., Welihinda, A., Hayes, A., Radford, A., Madry, A., Baker-Whitcomb, A., Beutel, A., Borzunov, A., Carney, A., Chow, A., Kirillov, A., Nichol, A., Paino, A., Renzin, A., Passos, A. T., Kirillov, A., Christakis, A., Conneau, A., Kamali, A., Jabri, A., Moyer, A., Tam, A., Crookes, A., Tootoochian, A., Tootoonchian, A., Kumar, A., Vallone, A., Karpathy, A., Braunstein, A., Cann, A., Codispoti, A., Galu, A., Kondrich, A., Tulloch, A., Mishchenko, A., Baek, A., Jiang, A., Pelisse, A., Woodford, A., Gosalia, A., Dhar, A., Pantuliano, A., Nayak, A., Oliver, A., Zoph, B., Ghorbani, B., Leimberger, B., Rossen, B., Sokolowsky, B., Wang, B., Zweig, B., Hoover, B., Samic, B., McGrew, B., Spero, B., Gierler, B., Cheng, B., Lightcap, B., Walkin, B., Quinn,

- 275 B., Guarraci, B., Hsu, B., Kellogg, B., Eastman, B., Lu-  
 276 garesi, C., Wainwright, C., Bassin, C., Hudson, C., Chu,  
 277 C., Nelson, C., Li, C., Shern, C. J., Conger, C., Barette,  
 278 C., Voss, C., Ding, C., Lu, C., Zhang, C., Beaumont, C.,  
 279 Hallacy, C., Koch, C., Gibson, C., Kim, C., Choi, C.,  
 280 McLeavey, C., Hesse, C., Fischer, C., Winter, C., Czar-  
 281 necki, C., Jarvis, C., Wei, C., Koumouzelis, C., Sherburn,  
 282 D., Kappler, D., Levin, D., Levy, D., Carr, D., Farhi, D.,  
 283 Mely, D., Robinson, D., Sasaki, D., Jin, D., Valladares,  
 284 D., Tsipras, D., Li, D., Nguyen, D. P., Findlay, D., Oiwoh,  
 285 E., Wong, E., Asdar, E., Proehl, E., Yang, E., Antonow, E.,  
 286 Kramer, E., Peterson, E., Sigler, E., Wallace, E., Brevdo,  
 287 E., Mays, E., Khorasani, F., Such, F. P., Raso, F., Zhang,  
 288 F., von Lohmann, F., Sulit, F., Goh, G., Oden, G., Salmon,  
 289 G., Starace, G., Brockman, G., Salman, H., Bao, H.,  
 290 Hu, H., Wong, H., Wang, H., Schmidt, H., Whitney, H.,  
 291 Jun, H., Kirchner, H., de Oliveira Pinto, H. P., Ren, H.,  
 292 Chang, H., Chung, H. W., Kivlichan, I., O’Connell, I.,  
 293 O’Connell, I., Osband, I., Silber, I., Sohl, I., Okuyucu,  
 294 I., Lan, I., Kostrikov, I., Sutskever, I., Kanitscheider, I.,  
 295 Gulrajani, I., Coxon, J., Menick, J., Pachocki, J., Aung, J.,  
 296 Betker, J., Crooks, J., Lennon, J., Kiros, J., Leike, J., Park,  
 297 J., Kwon, J., Phang, J., Teplitz, J., Wei, J., Wolfe, J., Chen,  
 298 J., Harris, J., Varavva, J., Lee, J. G., Shieh, J., Lin, J., Yu,  
 299 J., Weng, J., Tang, J., Yu, J., Jang, J., Candela, J. Q., Beut-  
 300 ler, J., Landers, J., Parish, J., Heidecke, J., Schulman, J.,  
 301 Lachman, J., McKay, J., Uesato, J., Ward, J., Kim, J. W.,  
 302 Huizinga, J., Sitkin, J., Kraaijeveld, J., Gross, J., Kap-  
 303 lan, J., Snyder, J., Achiam, J., Jiao, J., Lee, J., Zhuang,  
 304 J., Harriman, J., Fricke, K., Hayashi, K., Singhal, K.,  
 305 Shi, K., Karthik, K., Wood, K., Rimbach, K., Hsu, K.,  
 306 Nguyen, K., Gu-Lemberg, K., Button, K., Liu, K., Howe,  
 307 K., Muthukumar, K., Luther, K., Ahmad, L., Kai, L., Itow,  
 308 L., Workman, L., Pathak, L., Chen, L., Jing, L., Guy, L.,  
 309 Fedus, L., Zhou, L., Mamitsuka, L., Weng, L., McCal-  
 310 lum, L., Held, L., Ouyang, L., Feuvrier, L., Zhang, L.,  
 311 Kondraciuk, L., Kaiser, L., Hewitt, L., Metz, L., Doshi,  
 312 L., Aflak, M., Simens, M., Boyd, M., Thompson, M.,  
 313 Dukhan, M., Chen, M., Gray, M., Hudnall, M., Zhang, M.,  
 314 Aljube, M., Litwin, M., Zeng, M., Johnson, M., Shetty,  
 315 M., Gupta, M., Shah, M., Yatbaz, M., Yang, M. J., Zhong,  
 316 M., Glaese, M., Chen, M., Janner, M., Lampe, M., Petrov,  
 317 M., Wu, M., Wang, M., Fradin, M., Pokrass, M., Castro,  
 318 M., de Castro, M. O. T., Pavlov, M., Brundage, M., Wang,  
 319 M., Khan, M., Murati, M., Bavarian, M., Lin, M., Yesil-  
 320 dal, M., Soto, N., Gimelshein, N., Cone, N., Staudacher,  
 321 N., Summers, N., LaFontaine, N., Chowdhury, N., Ryder,  
 322 N., Stathas, N., Turley, N., Tezak, N., Felix, N., Kudige,  
 323 N., Keskar, N., Deutsch, N., Bundick, N., Puckett, N.,  
 324 Nachum, O., Okelola, O., Boiko, O., Murk, O., Jaffe, O.,  
 325 Watkins, O., Godement, O., Campbell-Moore, O., Chao,  
 326 P., McMillan, P., Belov, P., Su, P., Bak, P., Bakkum, P.,  
 327 Deng, P., Dolan, P., Hoeschele, P., Welinder, P., Tillet,  
 328 P., Pronin, P., Tillet, P., Dhariwal, P., Yuan, Q., Dias,  
 R., Lim, R., Arora, R., Troll, R., Lin, R., Lopes, R. G.,  
 Puri, R., Miyara, R., Leike, R., Gaubert, R., Zamani, R.,  
 Wang, R., Donnelly, R., Honsby, R., Smith, R., Sahai, R.,  
 Ramchandani, R., Huet, R., Carmichael, R., Zellers, R.,  
 Chen, R., Chen, R., Nigmatullin, R., Cheu, R., Jain, S.,  
 Altman, S., Schoenholz, S., Toizer, S., Miserendino, S.,  
 Agarwal, S., Culver, S., Ethersmith, S., Gray, S., Grove,  
 S., Metzger, S., Hermani, S., Jain, S., Zhao, S., Wu, S.,  
 Jomoto, S., Wu, S., Shuaiqi, Xia, Phene, S., Papay, S.,  
 Narayanan, S., Coffey, S., Lee, S., Hall, S., Balaji, S.,  
 Broda, T., Stramer, T., Xu, T., Gogineni, T., Christian-  
 son, T., Sanders, T., Patwardhan, T., Cunningham, T.,  
 Degry, T., Dimson, T., Raoux, T., Shadwell, T., Zheng,  
 T., Underwood, T., Markov, T., Sherbakov, T., Rubin, T.,  
 Stasi, T., Kaftan, T., Heywood, T., Peterson, T., Walters,  
 T., Eloundou, T., Qi, V., Moeller, V., Monaco, V., Kuo,  
 V., Fomenko, V., Chang, W., Zheng, W., Zhou, W., Man-  
 assra, W., Sheu, W., Zaremba, W., Patil, Y., Qian, Y.,  
 Kim, Y., Cheng, Y., Zhang, Y., He, Y., Zhang, Y., Jin, Y.,  
 Dai, Y., and Malkov, Y. Gpt-4o system card, 2024. URL  
<https://arxiv.org/abs/2410.21276>.
- Orgad, H., Toker, M., Gekhman, Z., Reichart, R., Szpektor,  
 I., Kotek, H., and Belinkov, Y. Llm know more than  
 they show: On the intrinsic representation of llm hallu-  
 cinations, 2025. URL <https://arxiv.org/abs/2410.02707>.
- Park, K., Choe, Y. J., and Veitch, V. The linear representa-  
 tion hypothesis and the geometry of large language mod-  
 els, 2024. URL <https://arxiv.org/abs/2311.03658>.
- Paulo, G., Mallen, A., Juang, C., and Belrose, N. Auto-  
 matically interpreting millions of features in large lan-  
 guage models, 2025. URL <https://arxiv.org/abs/2410.13928>.
- Ravichander, A., Belinkov, Y., and Hovy, E. Probing the  
 probing paradigm: Does probing accuracy entail task  
 relevance?, 2021. URL <https://arxiv.org/abs/2005.00719>.
- Singh, A., Fry, A., Perelman, A., Tart, A., Ganesh, A.,  
 El-Kishky, A., McLaughlin, A., Low, A., Ostrow, A.,  
 Ananthram, A., Nathan, A., Luo, A., Helyar, A., Madry,  
 A., Efremov, A., Spyra, A., Baker-Whitcomb, A., Beutel,  
 A., Karpenko, A., Makelov, A., Neitz, A., Wei, A., Barr,  
 A., Kirchmeyer, A., Ivanov, A., Christakis, A., Gillespie,  
 A., Tam, A., Bennett, A., Wan, A., Huang, A., Sandjideh,  
 A. M., Yang, A., Kumar, A., Saraiva, A., Vallone, A.,  
 Gheorghe, A., Garcia, A. G., Braunstein, A., Liu, A.,  
 Schmidt, A., Mereskin, A., Mishchenko, A., Applebaum,  
 A., Rogerson, A., Rajan, A., Wei, A., Kotha, A., Srivas-  
 tava, A., Agrawal, A., Vijayvergiya, A., Tyra, A., Nair,  
 A., Nayak, A., Eggers, B., Ji, B., Hoover, B., Chen, B.,

- 330 Chen, B., Barak, B., Minaiev, B., Hao, B., Baker, B.,  
 331 Lightcap, B., McKinzie, B., Wang, B., Quinn, B., Fioca,  
 332 B., Hsu, B., Yang, B., Yu, B., Zhang, B., Brenner, B.,  
 333 Zetino, C. R., Raymond, C., Lugaresi, C., Paz, C., Hud-  
 334 son, C., Whitney, C., Li, C., Chen, C., Cole, C., Voss,  
 335 C., Ding, C., Shen, C., Huang, C., Colby, C., Hallacy, C.,  
 336 Koch, C., Lu, C., Kaplan, C., Kim, C., Minott-Henriques,  
 337 C., Frey, C., Yu, C., Czarnecki, C., Reid, C., Wei, C.,  
 338 Decareaux, C., Scheau, C., Zhang, C., Forbes, C., Tang,  
 339 D., Goldberg, D., Roberts, D., Palmie, D., Kappler, D.,  
 340 Levine, D., Wright, D., Leo, D., Lin, D., Robinson, D.,  
 341 Grabb, D., Chen, D., Lim, D., Salama, D., Bhattacharjee,  
 342 D., Tsipras, D., Li, D., Yu, D., Strouse, D., Williams, D.,  
 343 Hunn, D., Bayes, E., Arbus, E., Akyurek, E., Le, E. Y.,  
 344 Widmann, E., Yani, E., Proehl, E., Sert, E., Cheung, E.,  
 345 Schwartz, E., Han, E., Jiang, E., Mitchell, E., Sigler, E.,  
 346 Wallace, E., Ritter, E., Kavanaugh, E., Mays, E., Nikishin,  
 347 E., Li, F., Such, F. P., de Avila Belbute Peres, F., Raso,  
 348 F., Bekerman, F., Tsimpourlas, F., Chantzis, F., Song, F.,  
 349 Zhang, F., Raila, G., McGrath, G., Briggs, G., Yang, G.,  
 350 Parascandolo, G., Chabot, G., Kim, G., Zhao, G., Valiant,  
 351 G., Leclerc, G., Salman, H., Wang, H., Sheng, H., Jiang,  
 352 H., Wang, H., Jin, H., Sikchi, H., Schmidt, H., Aspegren,  
 353 H., Chen, H., Qiu, H., Lightman, H., Covert, I., Kivlichan,  
 354 I., Silber, I., Sohl, I., Hammoud, I., Clavera, I., Lan, I.,  
 355 Akkaya, I., Kostrikov, I., Kofman, I., Etinger, I., Singal,  
 356 I., Hehir, J., Huh, J., Pan, J., Wilczynski, J., Pachocki, J.,  
 357 Lee, J., Quinn, J., Kiros, J., Kalra, J., Samaroo, J., Wang,  
 358 J., Wolfe, J., Chen, J., Wang, J., Harb, J., Han, J., Wang,  
 359 J., Zhao, J., Chen, J., Yang, J., Tworek, J., Chand, J., Lan-  
 360 don, J., Liang, J., Lin, J., Liu, J., Wang, J., Tang, J., Yin,  
 361 J., Jang, J., Morris, J., Flynn, J., Ferstad, J., Heidecke, J.,  
 362 Fishbein, J., Hallman, J., Grant, J., Chien, J., Gordon, J.,  
 363 Park, J., Liss, J., Kraaijeveld, J., Guay, J., Mo, J., Lawson,  
 364 J., McGrath, J., Vendrow, J., Jiao, J., Lee, J., Steele, J.,  
 365 Wang, J., Mao, J., Chen, K., Hayashi, K., Xiao, K., Salahi,  
 366 K., Wu, K., Sekhri, K., Sharma, K., Singhal, K., Li, K.,  
 367 Nguyen, K., Gu-Lemberg, K., King, K., Liu, K., Stone,  
 368 K., Yu, K., Ying, K., Georgiev, K., Lim, K., Tirumala,  
 369 K., Miller, K., Ahmad, L., Lv, L., Clare, L., Fauconnet,  
 370 L., Itow, L., Yang, L., Romaniuk, L., Anise, L., Byron,  
 371 L., Pathak, L., Maksin, L., Lo, L., Ho, L., Jing, L., Wu,  
 372 L., Xiong, L., Mamitsuka, L., Yang, L., McCallum, L.,  
 373 Held, L., Bourgeois, L., Engstrom, L., Kuhn, L., Fevrier,  
 374 L., Zhang, L., Switzer, L., Kondraciuk, L., Kaiser, L.,  
 375 Joglekar, M., Singh, M., Shah, M., Stratta, M., Williams,  
 376 M., Chen, M., Sun, M., Cayton, M., Li, M., Zhang, M.,  
 377 Aljube, M., Nichols, M., Haines, M., Schwarzer, M.,  
 378 Gupta, M., Shah, M., Guan, M. Y., Huang, M., Dong, M.,  
 379 Wang, M., Glaese, M., Carroll, M., Lampe, M., Malek,  
 380 M., Sharman, M., Zhang, M., Wang, M., Pokrass, M.,  
 381 Florian, M., Pavlov, M., Wang, M., Chen, M., Wang,  
 382 M., Feng, M., Bavarian, M., Lin, M., Abdool, M., Ro-  
 383 haninejad, M., Soto, N., Staudacher, N., LaFontaine, N.,  
 384 Marwell, N., Liu, N., Preston, N., Turley, N., Ansman, N.,  
 Blades, N., Pancha, N., Mikhaylin, N., Felix, N., Handa,  
 N., Rai, N., Keskar, N., Brown, N., Nachum, O., Boiko,  
 O., Murk, O., Watkins, O., Gleeson, O., Mishkin, P.,  
 Lesiewicz, P., Baltescu, P., Belov, P., Zhokhov, P., Pronin,  
 P., Guo, P., Thacker, P., Liu, Q., Yuan, Q., Liu, Q., Dias,  
 R., Puckett, R., Arora, R., Mullapudi, R. T., Gaon, R., Mi-  
 yara, R., Song, R., Aggarwal, R., Marsan, R., Yemiru, R.,  
 Xiong, R., Kshirsagar, R., Nuttall, R., Tsiupa, R., Eldan,  
 R., Wang, R., James, R., Ziv, R., Shu, R., Nigmatullin,  
 R., Jain, S., Talaie, S., Altman, S., Arnesen, S., Toizer, S.,  
 Toyer, S., Miserendino, S., Agarwal, S., Yoo, S., Heon, S.,  
 Ethersmith, S., Grove, S., Taylor, S., Bubeck, S., Baniesiu,  
 S., Amdo, S., Zhao, S., Wu, S., Santurkar, S., Zhao, S.,  
 Chaudhuri, S. R., Krishnaswamy, S., Shuaiqi, Xia, Cheng,  
 S., Anadkat, S., Fishman, S. P., Tobin, S., Fu, S., Jain,  
 S., Mei, S., Egoian, S., Kim, S., Golden, S., Mah, S.,  
 Lin, S., Imm, S., Sharpe, S., Yadlowsky, S., Choudhry,  
 S., Eum, S., Sanjeev, S., Khan, T., Stramer, T., Wang,  
 T., Xin, T., Gogineni, T., Christianson, T., Sanders, T.,  
 Patwardhan, T., Degry, T., Shadwell, T., Fu, T., Gao, T.,  
 Garipov, T., Sriskandarajah, T., Sherbakov, T., Korbak,  
 T., Kaftan, T., Hiratsuka, T., Wang, T., Song, T., Zhao, T.,  
 Peterson, T., Kharitonov, V., Chernova, V., Kosaraju, V.,  
 Kuo, V., Pong, V., Verma, V., Petrov, V., Jiang, W., Zhang,  
 W., Zhou, W., Xie, W., Zhan, W., McCabe, W., DePue,  
 W., Ellsworth, W., Bain, W., Thompson, W., Chen, X.,  
 Qi, X., Xiang, X., Shi, X., Dubois, Y., Yu, Y., Khakbaz,  
 Y., Wu, Y., Qian, Y., Lee, Y. T., Chen, Y., Zhang, Y.,  
 Xiong, Y., Tian, Y., Cha, Y., Bai, Y., Yang, Y., Yuan,  
 Y., Li, Y., Zhang, Y., Yang, Y., Jin, Y., Jiang, Y., Wang,  
 Y., Wang, Y., Liu, Y., Stubenvoll, Z., Dou, Z., Wu, Z.,  
 and Wang, Z. Openai gpt-5 system card, 2026. URL <https://arxiv.org/abs/2601.03267>.
- Wang, C., Wei, Z., Liu, Q., and Chen, M. False sense of security: Why probing-based malicious input detection fails to generalize, 2025. URL <https://arxiv.org/abs/2509.03888>.
- Ying, Z. J., Ravfogel, S., Kriegeskorte, N., and Hase, P. The truthfulness spectrum hypothesis, 2026. URL <https://arxiv.org/abs/2602.20273>.
- Zou, A., Phan, L., Chen, S., Campbell, J., Guo, P., Ren, R., Pan, A., Yin, X., Mazeika, M., Dombrowski, A.-K., Goel, S., Li, N., Byun, M. J., Wang, Z., Mallen, A., Basart, S., Koyejo, S., Song, D., Fredrikson, M., Kolter, J. Z., and Hendrycks, D. Representation engineering: A top-down approach to ai transparency, 2025. URL <https://arxiv.org/abs/2310.01405>.

## A. Related Work

**Linear concept directions in activation space.** A related line of work studies whether high-level behavioral or semantic properties are represented as directions or low-dimensional subspaces in activation space. Marks & Tegmark (2024) report linear structure in true/false representations across datasets, while Park et al. (2024) formalize the linear representation hypothesis and connect it to both probing and steering. Similarly, Zou et al. (2025) model high-level properties such as honesty and emotion as population-level directions that can be read from and intervened on in model activations.

This perspective has also been used to explain and control model behavior. (Arditi et al., 2024) show that a single direction can mediate refusal behavior, suggesting that some high-level behaviors may be localized in simple linear features. Lu et al. (2026) similarly identify an “assistant” axis corresponding to assistant-like behavior. Together, these works suggest that linear structure can support both mechanistic understanding and activation-space interventions, while also raising questions about when such directions capture stable concepts rather than context-dependent correlations.

**Linear probes and their limitations.** Linear and lightweight probes have a long history in interpretability as tools for analyzing what information is represented in neural networks (Alain & Bengio, 2018; Hewitt & Liang, 2019). Recent work has applied such probes to safety-relevant signals in LLM activations, including latent truthfulness without direct supervision (Burns et al., 2024) and strategic deception (Goldowsky-Dill et al., 2025).

At the same time, methodological skepticism is central to the probing literature. Hewitt & Liang (2019) argue that high probe accuracy alone is insufficient evidence that a model already represents the target property, since the probe may itself learn the task or exploit surface-level correlates. Similarly, Ravichander et al. (2021) question whether probing reliably establishes what information is encoded or merely recoverable. Recent work also finds that probes with strong in-distribution performance can generalize unevenly out of distribution (Kramár et al., 2026). These concerns motivate careful evaluation of whether safety probes identify robust internal signals rather than dataset-specific artifacts.

**Interpreting features with natural-language explanations.** Another approach aims to interpret internal representations by assigning natural-language descriptions to neurons or learned features. Bills et al. (2023) introduced the use of LLMs to generate explanations of neurons from high-activation examples. This approach has since been extended to the automatic labeling of features in sparse autoencoders ((Paulo et al., 2025), (Bricken et al., 2023)), enabling large-scale interpretation of learned feature dictionaries.

Recent work has further explored whether natural-language feature descriptions can be used not only for interpretation but also for intervention. For example, Casademunt et al. (2025) identify concepts through automated natural-language interpretation and ablate them during fine-tuning to reduce misaligned behavior.

## B. Dataset Details

**Roleplaying:** the model is placed in a simple scenario with an opportunity to lie.

**Insider Trading Report:** acting as a trader, the model decides whether to honestly report that it used insider information.

**Insider Trading Confirm:** after hiding earlier use of insider information, the model is asked to confirm it did not use any.

**Sandbagging:** the model is told to deliberately under-perform on some tasks while behaving normally on others.

Dataset	Total Examples	Label Counts (-1 / 0 / 1)
Roleplaying	371	140 / 59 / 172
Insider Trading Report	473	218 / 194 / 61
Insider Trading Confirm	151	29 / 38 / 84
Sandbagging	1000	379 / 375 / 246

Table 1. Dataset label distributions.

### C. LLM Judge Prompt

We use GPT-5-mini (Singh et al., 2026) to interpret directions.

You are interpreting a principal component direction in a neural representation space. This direction is bipolar:

- Positive activations mean movement toward one pole.
- Negative activations mean movement toward the opposite pole.
- Larger magnitudes mean a stronger match to that pole.

Each line is formatted as token<TAB>activation, with activations normalized to integers from -10 to +10.

The examples below were selected from a dataset with this setup: The model is instructed to Roleplay in a simple scenario. The model is asked a question which gives a strategic chance to lie.

High-positive examples:

{positive\_examples}

High-negative examples:

{negative\_examples}

Infer the dominant contrast captured by the direction using both the positive and negative examples.

Important instructions:

- Do not assume the direction is a broad sentiment direction unless the examples clearly support that.
- Prefer the narrowest accurate description over a broad generic one.
- Do not overgeneralize beyond what is supported by the examples.
- Rate how likely this direction is to be useful for out-of-distribution deception on a scale from 1 to 10, and briefly explain why.
- Give low scores to directions that are mainly surface cues, lexical habits, or task-specific phrasing, even if they correlate with deception in these examples.
- A direction that is noisy, brittle, context-dependent, or likely to fail under different wording or domains should usually be scored around 1 to 4, not 5 or above.
- Reserve scores of 5 or higher for directions that seem plausibly tied to a more general deception-relevant strategy or behavior, rather than local wording.

Write your answer in exactly this format:

Positive pole: ...

Negative pole: ...

Shared direction: ...

OOD deception usefulness (1-10): ...

### D. Lower PCs Show Stronger Decoder Alignment

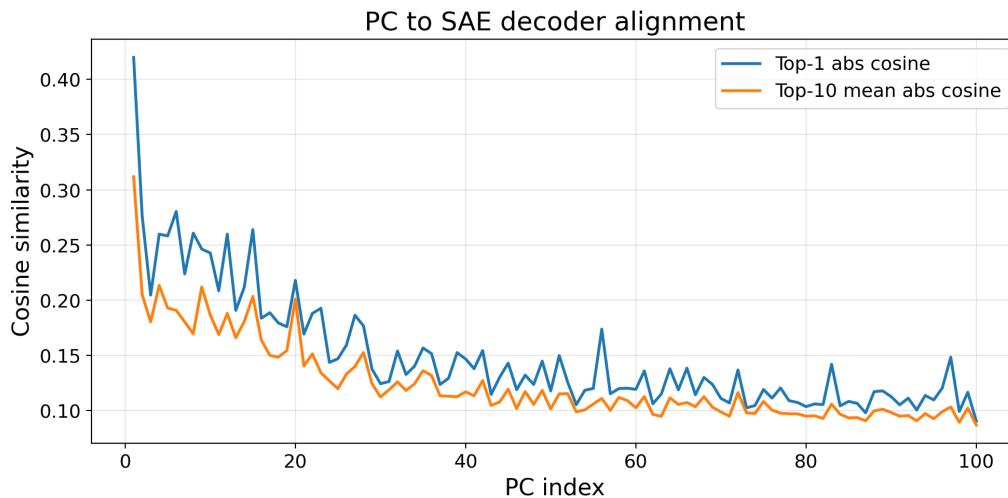


Figure 5. Lower PCs show stronger decoder alignment, both for the best-matching SAE feature and for the mean of the top 10 features (Spearman  $\rho = -0.857$  and  $-0.917$ , respectively). We used andyrdt/saes-llama-3.1-8b-instruct (andyrdt, 2024)

**E. PC Set Comparisons**

Jaccard overlap between the PC sets of 2 methods (greedy and target-pca) is 0.182 for Insider Trading Report, 0.077 for Insider Trading Confirm, and 0.167 for Sandbagging. The 2 methods agree on less than a fifth of their selections.

Together these results suggest that the transferable subspace is not the part of representation space that is most discriminative on the source, and not the part most discriminative on the target.

Selection method	Dataset / target	Selected PCs
Source probe weights	Source	PC7, PC8, PC11, PC12, PC4, PC5, PC2, PC3, PC6, PC27, PC22, PC21, PC1, PC41, PC18
Greedy selection	Insider Trading Report	PC6, PC89, PC62, PC66, PC26, PC14, PC57, PC25, PC93, PC9, PC70, PC97, PC78
Greedy selection	Insider Trading Confirm	PC71, PC53, PC28, PC86, PC32, PC81, PC77
Greedy selection	Sandbagging	PC37, PC62, PC54, PC41, PC46, PC9, PC57, PC44, PC51, PC43, PC19, PC48, PC86, PC81, PC33
Target separability ranking	Insider Trading Report	PC5, PC6, PC89, PC18, PC79, PC46, PC87, PC4, PC14, PC34, PC29, PC1, PC66, PC67, PC19
Target separability ranking	Insider Trading Confirm	PC4, PC1, PC8, PC53, PC30, PC7, PC13, PC18, PC5, PC20, PC26, PC49, PC28, PC31, PC97
Target separability ranking	Sandbagging	PC51, PC91, PC97, PC13, PC39, PC77, PC4, PC66, PC53, PC86, PC62, PC23, PC41, PC10, PC61

Table 2. Comparison of PCs selected by source-only probe weights, greedy transfer selection, and target separability ranking.

F. PC Interpretability

Table 3. Principal components selected for likely transferability, based on interpretation with an LLM judge.

PC	Selected On	OOD Score	Full explanation
PC62	Insider Trading Report + Sandbagging	6	OOD deception usefulness (1-10): 6 — This direction maps onto a general deception-relevant behavior (deny/deflect vs provide details), so it can be useful across settings, but it is noisy: true speakers also deny or obfuscate for legitimate reasons and liars sometimes produce elaborate fabricated detail, and the signal may pick up simple lexical negation rather than semantic strategy.
PC57	Insider Trading Report + Sandbagging	6	This direction aligns with a general deception-relevant strategy (deny/minimize or present innocuous explanations vs give specific/incriminating details), so it can help flag rhetorical moves often used when hiding or deflecting. However it is partly lexical and task-specific (relies on surface denials, hedges, and domain terms), so it will be brittle and context-dependent out of distribution.
PC9	Insider Trading Report + Sandbagging	3	OOD deception usefulness (1-10): 3 — This direction mainly picks up domain-specific lexical cues (medical/leave language vs project/status language) and stylistic framing rather than a generalizable deception strategy. It may flag genuine illness claims and is brittle across domains and rephrasings.
PC81	Insider Trading Confirm + Sandbagging	3	OOD deception usefulness (1-10): 3 — The direction picks up a plausible deception-relevant contrast (admit vs deny/hedge), but it mainly relies on surface/formal cues (negations, hedges, résumé style, enumerations) that are task- and wording-dependent and likely brittle across domains and phrasings.
PC86	Insider Trading Confirm + Sandbagging	6	OOD deception usefulness (1-10): 6 — Hedging and defensive PR-style wording is a general tactic associated with deceptive or self-protective behavior, so the direction could transfer across domains. However it is noisy (used legitimately for privacy, politeness, or negotiation), sensitive to phrasing and role, and will produce false positives, so it’s only moderately reliable out-of-distribution.

Table 4. Principal components and their OOD scores with interpretations.

PC	OOD Score	Full explanation
PC7	3	OOD deception usefulness (1-10): 3 — This direction mostly picks up a pragmatic contrast (deny/deflect vs admit) which can correlate with deception in some settings, but it relies heavily on surface cues (explicit “No”, word order, presence of “I” + admission) and task-specific phrasing. It is therefore brittle and likely to fail or mislead in different domains or with varied wording.
PC8	4	OOD deception usefulness (1-10): 4 — The direction captures a useful narrow contrast (denial/excuse vs admission) that can correlate with evasive or deceptive responses in similar QA/roleplay tasks, but it largely picks up surface cues (explicit “No”/“Yes” and common denial phrasing) and is likely brittle under paraphrase, different domains, or subtler deception strategies.
PC11	3	OOD deception usefulness (1-10): 3 — The direction mainly picks up a stylistic/pragmatic contrast (assertive first-person role/defense vs problem-reporting/negation). That can correlate with evasive or face-saving lies in some scenarios, but it’s surfacey and context dependent and likely to fail across different domains or phrasings.
PC12	3	OOD deception usefulness (1-10): 3 — This direction mainly captures discourse style/forcefulness (a surface cue). While overconfident or over-assertive phrasing can sometimes correlate with deceptive defensive behavior, it is brittle and domain-dependent and will produce many false positives/negatives when applied across different speakers, topics, or cultural styles.
PC4	6	OOD deception usefulness (1-10): 6 — This captures a general communicative strategy relevant to deception (deny vs admit), so it can transfer beyond the specific prompts. However it also leans on surface cues (explicit “Yes”/“No” turns and common defensive phrasings) and may be brittle when denials or lies use different wording or more subtle tactics.