
Inference-Time Alignment via Hypothesis Reweighting

Anonymous Authors¹

Abstract

Chat assistants must handle diverse and often conflicting user preferences, requiring adaptability to various user needs. We propose a lightweight framework to address the general challenge of aligning models to user intent at inference time. Our approach involves training an efficient ensemble, i.e., a single neural network with multiple prediction heads, each representing a different function consistent with the training data. Our main contribution is HYRE, a simple adaptation technique that dynamically reweights ensemble members at test time using a small set of labeled examples from the target distribution, which can be labeled in advance or actively queried from a larger unlabeled pool. The computational cost of our training procedure is comparable to fine-tuning a single model, and thus scales to large pretrained backbones. We empirically validate HYRE in several target evaluation distributions. With as few as five preference pairs from each target distribution, adaptation via HYRE surpasses state-of-the-art reward models on RewardBench at both the 2B and 8B parameter scales.

1 Introduction

Task specification—describing precisely what a machine learning model should do—is inherently iterative and fundamentally incomplete under any finite set of instructions or training examples. As models grow more powerful and are applied to increasingly complex and nuanced tasks, this problem arises in many forms, from spurious correlations in the data to conflicting user preferences. Consider a chatbot trained via Reinforcement Learning from Human Feedback (RLHF) (Siththaranjan et al., 2023) on a broad distribution of user preferences. Such models often perform ad-

equately in aggregate but systemically fail to address specific users’ needs, since different individuals have distinct, sometimes contradictory, notions of desirable responses. Meeting these user-specific requirements necessitates rapid model adaptation with minimal supervision. However, existing adaptation strategies, such as prompt-based methods (Gao et al., 2020; Khattab et al., 2023; Yuksekgonul et al., 2024) and fine-tuning (Houlsby et al., 2019; Hu et al., 2021; Liu et al., 2024; Wu et al., 2024), can be computationally heavy, typically requiring multiple forward-backward passes or large-scale gradient updates. This renders them unsuitable for on-the-fly adaptation at test time.

To efficiently resolve ambiguity at test time, we draw on recent progress in efficient ensemble architectures (Osband et al., 2023). These methods let a single backbone network represent a broad range of plausible functions at a small overhead, capturing the different ways the model can interpret the training set. While prior work focuses largely on using ensembles for uncertainty estimation, we propose using them to disambiguate tasks in real time: by quickly assessing which members of the ensemble best match a new distribution, we can “pick the right interpretation” for that scenario.

We introduce Hypothesis Reweighting (HYRE), a two-step approach that scales to large models. First, we train an ensemble of function heads on top of a shared backbone, ensuring each head individually fits the training data. Next, at inference time, we gather a few labeled examples from the target distribution—either proactively queried or provided in advance—and measure each head’s performance. We then reweight the ensemble using a generalized Bayesian update that favors the heads performing best on the adaptation set. Crucially, this update supports non-differentiable metrics like 0-1 error, and it requires only a single forward pass over the adaptation set, making it far more efficient than conventional fine-tuning.

We evaluate HYRE across over 20 target distributions including preference personalization tasks and benchmarks for response safety and usefulness. With just 1-5 adaptation examples, HYRE reweights a 100-head ensemble at negligible (< 1%) overhead, improving the state-of-the-art reward model by an average of 20% absolute accuracy across 32 tasks. Adaptation via HYRE also outperforms the state-

¹Anonymous 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.

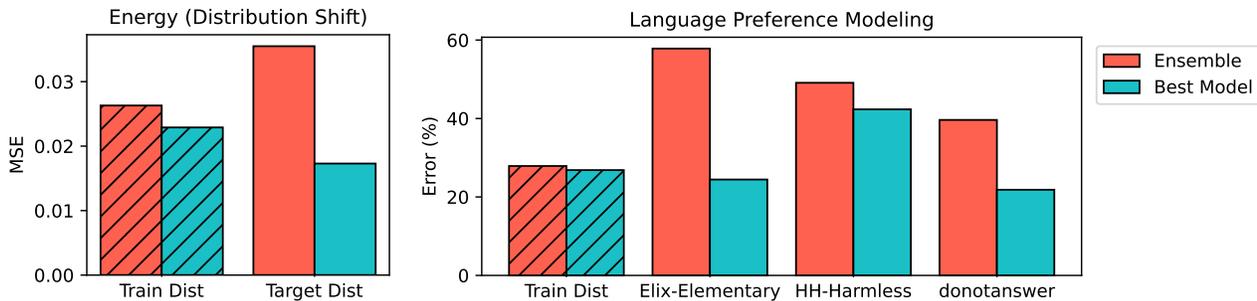


Figure 1: Motivating observation: the ensemble average often performs worse than a single well-chosen member. This tendency is particularly pronounced further away from the training distribution. **HYRE goes further than selecting the best head—it finds a continuous weighting of all heads.**

Algorithm 1 HYRE (Inference Time)

Require: Ensemble members $f_{1..K}$, unlabeled dataset $x_{1..N}$, query budget B

- 1: Initialize weights $w \leftarrow [\frac{1}{K}, \dots, \frac{1}{K}]$, query set $Q \leftarrow \emptyset$
- 2: **for** $i \leftarrow 1$ to B **do**
- 3: (Optional) Query label y_n for $\arg \max_n c(x_n)$ and add (x_n, y_n) to Q (Appendix F)
- 4: Compute accuracy $\text{acc}_k = \sum_{n \in Q} \text{acc}(f_k, x_n, y_n)$ for each k
- 5: Update ensemble weight $w_k \propto \exp(\text{acc}_k + p)$ (Section 2)
- 6: **end for**
- 7: **Return** final weighted ensemble function $f_w : x \mapsto \sum_{k=1}^K w_k f_k(x)$

of-the-art models on RewardBench (Lambert et al., 2024) at both the 2B and 8B parameter scales. These findings demonstrate that fast inference-time reweighting of a well-chosen ensemble can effectively adapt to new tasks with minimal supervision.

2 HYRE: Fast Inference-Time Ensemble Reweighting

In this section, we motivate and describe Hypothesis Reweighting (HYRE), a simple and computationally efficient method for few-shot adaptation to new tasks. HYRE dynamically adjusts the weights assigned to different ensemble members at test time based on a few labeled samples from the new task. HYRE leverages the ensemble’s diversity—each member representing a different function that fits the training data—to efficiently adapt without re-training any model parameters.

Our method is motivated by Figure 1, which shows that the ensemble average often performs substantially worse than the best weighted ensemble. This tendency is particularly pronounced further away from the training distribution.

Given an ensemble of K models f_1, \dots, f_K , we aim to dynamically update their weights based on adaptation data. As a practical test-time assumption in settings where we cannot further train neural networks, we can think of the “best” model as being one of the K ensemble particles that performs best on the evaluation distribution. Starting with uniform weights $w_k = \frac{1}{K}$, we update them as new labeled data from P_{eval} becomes available.

The weighted ensemble prediction is $f_w(x) = \sum_{i=1}^K w_i f_i(x)$, where each $w_i \geq 0$ and $\sum_{i=1}^K w_i = 1$. We measure each member’s performance using a loss function $l(f_k, x, y)$ and compute their cumulative loss on adaptation data $\mathcal{L}(f_k, \mathcal{D}_{\text{adapt}}) = \sum_{(x,y) \in \mathcal{D}_{\text{adapt}}} l(f_k, x, y)$. The weights are updated using a softmax on negative cumulative loss:

$$w_k = \frac{\exp(-\mathcal{L}(f_k, \mathcal{D}_{\text{adapt}}))}{\sum_{i=1}^K \exp(-\mathcal{L}(f_i, \mathcal{D}_{\text{adapt}}))}. \quad (1)$$

As the loss $l(f_k, x, y)$, we use 0-1 error for classification and mean squared error for regression, though HYRE supports any performance metric since the weight update remains valid for non-differentiable functions. The complete adaptation procedure is summarized in Algorithm 1.

3 Experiments

We now empirically validate HYRE. We focus on three key questions: (1) Can HYRE effectively handle mild covariate shift? (2) Does HYRE scale to large models? (3) How robust and computationally efficient is HYRE? We describe the detailed setup for each experiment in the appendix.

3.1 Regression Data with Mild Covariate Shift

We evaluate HYRE on three UCI regression datasets (Kelly et al.)—Energy Efficiency, Kin8nm, and CCPP—using the protocol of Sharma et al. (2023): the top and bottom 5% of the data (sorted by mean input features) form an OOD target set, while the central 90% is split into train and validation sets. All methods employ 100 two-layer MLPs with

Model	Type	Overall	Chat	Chat Hard	Safety	Reasoning
Tulu-2-DPO-70B	DPO	79.1	97.5	60.5	84.5	74.1
StableLM-2-12B-Chat	DPO	79.9	96.6	55.5	78.1	89.4
Claude-3 Sonnet (June 2024)	Gen	84.2	96.4	74.0	81.6	84.7
GPT-4 (May 2024)	Gen	84.6	96.6	70.4	86.5	84.9
GPT-4 (Aug 2024)	Gen	86.7	96.1	76.1	88.1	86.6
Gemini-1.5-Pro-0924	Gen	86.8	94.1	77.0	85.8	90.2
Skywork-Reward-Gemma-2-27B	Seq	94.3	96.1	89.9	93.0	98.1
INF-ORM-Llama3.1-70B	Seq	95.1	96.6	91.0	93.6	99.1
GRM-Gemma2-2B	Seq	88.4	93.0	77.2	92.2	91.2
+ Ours (uniform)	Seq	87.1	96.4	73.1	87.4	89.8
+ Ours (N=1)	Seq + HYRE	86.5	92.4	71.5	85.1	92.5
+ Ours (N=5)	Seq + HYRE	88.5	95.0	72.5	90.3	93.1
+ Ours (N=10)	Seq + HYRE	89.7	96.4	74.7	92.4	93.5
+ Ours (best head oracle)*	Seq + Oracle	91.8	97.2	80.0	96.2	94.2
+ Ours (best weight oracle)*	Seq + Oracle	93.1	98.3	83.4	96.7	94.9
Skywork-Llama-3.1-8B	Seq	94.0	94.7	88.6	92.7	96.7
+ Ours (uniform)	Seq	94.0	95.0	87.2	93.0	96.8
+ Ours (N=1)	Seq + HYRE	94.3	95.2	87.8	93.0	97.5
+ Ours (N=5)	Seq + HYRE	94.7	95.5	88.6	93.2	97.8
+ Ours (N=10)	Seq + HYRE	95.0	95.9	89.3	93.5	97.9
+ Ours (best head oracle)*	Seq + Oracle	96.4	98.3	91.2	95.7	98.4
+ Ours (best weight oracle)*	Seq + Oracle	97.2	99.2	93.0	96.5	98.8

* Oracle methods show an upper bound on performance, using the test set.

Table 1: Accuracy across tasks in RewardBench. We report overall performance and breakdowns by task category for all models. **HYRE improves upon the state-of-the-art models at the 2B and 8B parameter scales with as few as 1-5 labeled samples per distribution.**

50 units each. As baselines, we consider a vanilla ensemble of independently trained models and MC Dropout (Gal & Ghahramani, 2016). We report the best-performing MC Dropout results across all architectures. Results in Table 4 demonstrate that uniform ensembles perform strongly in these OOD generalization settings and that HYRE consistently improves over the uniform ensemble.

3.2 Scalable Personalization of Preference Models

Experimental setup. We evaluate personalization using four sets of human preference benchmarks: Elix (Singh et al., 2025), RewardBench (Lambert et al., 2024), PERSONA (Castricato et al., 2024), and Anthropic HH (Bai et al., 2022). Together, these benchmarks contain 32 datasets, each encoding a different aspect of human preferences. To train HYRE on preference data, we attach Shared-Base ensemble heads to a pretrained 2B reward model and fine-tune it on the UltraFeedback (Cui et al., 2023) dataset, a standard dataset for reward model training. We use two public finetune checkpoints of Gemma 2B models, which achieve state-of-the-art performance on RewardBench at the 2B parameter scale, even outperforming GPT-4o (Achiam et al., 2023). Refer to Appendix G for our detailed setup.

We first evaluate the effectiveness of HYRE in adapting our reward model ensemble to new distributions at test time, comparing its performance to that of the original reward model. As shown in Figure 2, a simple uniform ensemble initially underperforms the original model, indicating that naive ensembling alone cannot ensure broad generalization. Nevertheless, HYRE quickly surpasses the baseline with just a few labeled examples per distribution. We show detailed dataset-level results in the appendix (Figure 9).

We compare HYRE against state-of-the-art reward models on the RewardBench leaderboard at both the 2B and 8B parameter scales. As shown in Table 8, HYRE—with only 1-5 labeled examples per distribution—exceeds the performance of many much larger reward models. We note that these reward models outperform strong generative reward models including Claude 3.5 Sonnet, GPT-4, and Gemini-1.5-Pro (Achiam et al., 2023; Anthropic, 2024; Team et al., 2024). This indicates that inference-time alignment can be a powerful alternative to naively scaling up reward models.

3.3 Comparison with Alternative Adaptation Methods

Few-shot prompting. We compare HYRE with few-shot prompting using GPT-4o-mini on two datasets from Re-

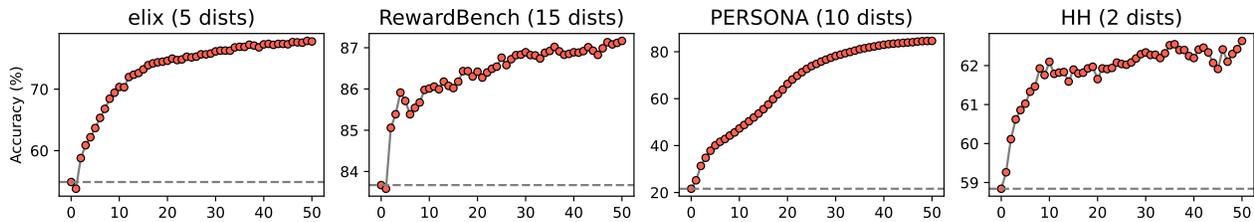


Figure 2: Average reward model accuracy as a function of adaptation set size. The dashed line shows the best available static 2B reward model for each dataset group. **HYRE consistently outperforms the state-of-the-art reward model with as few as 1-5 examples per distribution.**

Dataset	N=0	N=1	N=5	N=10	N=20	N=40	N=80
GPT-4o-mini N-Shot Prompting							
donotanswer	44.4	50.3	60.2	64.7	<u>68.7</u>	66.4	67.0
refusals	79.4	<u>82.7</u>	80.5	82.2	82.1	78.0	82.1
Llama-3.1-8B N-Shot Prompting							
donotanswer	46.6	52.8	59.6	<u>63.4</u>	41.2	62.6	*
refusals	61.6	<u>82.4</u>	80.0	72.2	44.4	79.0	*
Llama-3.1-8B + HYRE							
donotanswer	58.6	60.8	69.1	71.3			
refusals	88.9	90.0	94.0	95.2			

Table 2: Comparison with few-shot prompting on two datasets from RewardBench. (*) exceeds Together AI API token limit. We see a degradation in performance for both GPT-4o-mini and Llama-3.1-8B as we increase the number of examples, whereas HYRE consistently outperforms both across all sample sizes. **HYRE provides reliable test-time alignment, unlike few-shot prompting, which can degrade with too much context.**

wardBench. As shown in Table 2, HYRE consistently outperforms GPT-4o-mini across all sample sizes. Even in the zero-shot setting, specialized reward models like HYRE achieve higher performance than general-purpose language models like GPT-4o-mini. Notably, we observe that too many few-shot examples can actually harm performance, as seen with GPT-4o-mini’s performance drop after N=10 for donotanswer and N=20 for refusals. These results demonstrate that specialized reward models with inference-time adaptation can more efficiently leverage few-shot examples than general-purpose language models.

Fine-tuning on target data. We compare HYRE against models fine-tuned on the helpful-base and harmless-base training sets in the Anthropic-HH dataset. Results in Table 5 indicate that while targeted fine-tuning models achieve higher performance in their respective target metrics, they significantly reduce performance in the other. In contrast, our HYRE-adapted ensemble not only increases performance across each data distribution but also retains or slightly improves performance in the other split. We emphasize that we show fine-tuning performance only as a point of comparison; **fine-tuning a model for a target distribution is usually too computationally expensive to be done at inference time, and is thus not a practical**

	donotanswer	xstest-sr	refusals
HYRE w/ Cross-Entropy			
N=0	58.60 ± 4.93	82.80 ± 2.43	88.90 ± 3.25
N=1	62.53 ± 4.00	85.78 ± 3.06	92.49 ± 3.28
N=5	62.57 ± 1.88	87.38 ± 1.26	93.17 ± 2.19
N=10	62.25 ± 1.71	87.51 ± 1.19	93.21 ± 1.77
HYRE w/ Accuracy			
N=0	58.60 ± 4.93	82.80 ± 2.43	88.90 ± 3.25
N=1	60.81 ± 5.29	85.80 ± 2.45	90.00 ± 3.63
N=5	69.12 ± 5.81	89.28 ± 2.86	94.00 ± 2.45
N=10	71.32 ± 6.33	90.32 ± 3.20	95.20 ± 2.32
Oracle	76.54 ± 2.35	90.32 ± 1.91	99.50 ± 0.87

Table 3: Cross-entropy ablation experiment. We report average and std of accuracy (%) with varying numbers of adaptation examples (N) on three datasets. **Using accuracy as the adaptation objective for HYRE significantly improves post-adaptation performance.**

solution for inference-time alignment.

3.4 Ablation Studies and Analysis

Ablation on reweighting criteria. We investigate the impact of using binary cross entropy instead of accuracy for reweighting. As shown in Table 3, using cross-entropy loss for reweighting significantly degrades performance across three representative RewardBench datasets. This is because cross-entropy loss is sensitive to outliers as it is unbounded from above, and thus can quickly overfit to a head.

Computational overhead. HYRE is designed to be efficient enough to be used at inference time. We use a single pre-trained backbone with K prediction heads, where each head is a very small MLP compared to the backbone. The parameter overhead is negligible: in our reward model experiments, 100 ensemble heads (5.5×10^5 parameters) add less than 0.03% to the parameter count of the Gemma-2B backbone (2.0×10^9 parameters). At inference time, reweighting requires only a single forward pass through the backbone and heads, with the subsequent weight calculation being minimal. The total cost increase in time and memory for using HYRE compared to the single base model is less than 1%.

References

- Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Ahmed, A. M., Rafailov, R., Sharkov, S., Li, X., and Koyejo, S. Scalable ensembling for mitigating reward overoptimisation. *arXiv preprint arXiv:2406.01013*, 2024.
- Anthropic. Claude 3.5 sonnet. Accessed via Claude.ai, API, and cloud platforms, 2024. URL <https://www.anthropic.com>. Enhanced reasoning, state-of-the-art coding skills, computer use, and 200K context window. Available on Anthropic API, Amazon Bedrock, and Google Cloud’s Vertex AI.
- Arjovsky, M., Bottou, L., Gulrajani, I., and Lopez-Paz, D. Invariant risk minimization. *arXiv preprint arXiv:1907.02893*, 2019.
- Bai, Y., Jones, A., Ndousse, K., Askell, A., Chen, A., Das-Sarma, N., Drain, D., Fort, S., Ganguli, D., Henighan, T., et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- Barreto, A., Dumoulin, V., Mao, Y., Perez-Nieves, N., Shahriari, B., Dauphin, Y., Precup, D., and Larochelle, H. Capturing individual human preferences with reward features. *arXiv preprint arXiv:2503.17338*, 2025.
- Bissiri, P. G., Holmes, C. C., and Walker, S. G. A General Framework for Updating Belief Distributions. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 78(5):1103–1130, 02 2016. ISSN 1369-7412. doi: 10.1111/rssb.12158.
- Brown, B., Juravsky, J., Ehrlich, R., Clark, R., Le, Q. V., Ré, C., and Mirhoseini, A. Large language monkeys: Scaling inference compute with repeated sampling. *arXiv preprint arXiv:2407.21787*, 2024.
- Burda, Y., Edwards, H., Storkey, A., and Klimov, O. Exploration by random network distillation. *arXiv preprint arXiv:1810.12894*, 2018.
- Castricato, L., Lile, N., Rafailov, R., Fränken, J.-P., and Finn, C. Persona: A reproducible testbed for pluralistic alignment. *arXiv preprint arXiv:2407.17387*, 2024.
- Chen, D., Chen, Y., Rege, A., and Vinayak, R. K. Modeling the plurality of human preferences via ideal points. In *ICML 2024 Workshop on Models of Human Feedback for AI Alignment*, 2024. URL <https://openreview.net/forum?id=Ykd0xVxy6a>.
- Christiano, P. F., Leike, J., Brown, T., Martic, M., Legg, S., and Amodei, D. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.
- Coste, T., Anwar, U., Kirk, R., and Krueger, D. Reward model ensembles help mitigate overoptimization. *arXiv preprint arXiv:2310.02743*, 2023.
- Cui, G., Yuan, L., Ding, N., Yao, G., Zhu, W., Ni, Y., Xie, G., Liu, Z., and Sun, M. Ultrafeedback: Boosting language models with high-quality feedback. *arXiv preprint arXiv:2310.01377*, 2023.
- D’Amour, A., Heller, K., Moldovan, D., Adlam, B., Alipanahi, B., Beutel, A., Chen, C., Deaton, J., Eisenstein, J., Hoffman, M. D., et al. Underspecification presents challenges for credibility in modern machine learning. *Journal of Machine Learning Research*, 23(226):1–61, 2022.
- Fedus, W., Zoph, B., and Shazeer, N. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120):1–39, 2022.
- Gal, Y. and Ghahramani, Z. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pp. 1050–1059. PMLR, 2016.
- Gal, Y., Islam, R., and Ghahramani, Z. Deep bayesian active learning with image data. In *International conference on machine learning*, pp. 1183–1192. PMLR, 2017.
- Gao, L., Schulman, J., and Hilton, J. Scaling laws for reward model overoptimization. In *International Conference on Machine Learning*, pp. 10835–10866. PMLR, 2023.
- Gao, T., Fisch, A., and Chen, D. Making pre-trained language models better few-shot learners. *arXiv preprint arXiv:2012.15723*, 2020.
- Geirhos, R., Jacobsen, J.-H., Michaelis, C., Zemel, R., Brendel, W., Bethge, M., and Wichmann, F. A. Short-cut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673, 2020.
- Hadfield-Menell, D., Russell, S. J., Abbeel, P., and Dragan, A. Cooperative inverse reinforcement learning. *Advances in neural information processing systems*, 29, 2016.
- Houlsby, N., Huszár, F., Ghahramani, Z., and Lengyel, M. Bayesian active learning for classification and preference learning. *arXiv preprint arXiv:1112.5745*, 2011.

- 275 Houlsby, N., Giurgiu, A., Jastrzebski, S., Morrone, B.,
 276 De Laroussilhe, Q., Gesmundo, A., Attariyan, M., and
 277 Gelly, S. Parameter-efficient transfer learning for nlp. In
 278 *International conference on machine learning*, pp. 2790–
 279 2799. PMLR, 2019.
- 280 Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y.,
 281 Wang, S., Wang, L., and Chen, W. Lora: Low-rank
 282 adaptation of large language models. *arXiv preprint*
 283 *arXiv:2106.09685*, 2021.
- 284
 285 Izmailov, P., Nicholson, P., Lotfi, S., and Wilson, A. G.
 286 Dangers of bayesian model averaging under covariate
 287 shift. In Ranzato, M., Beygelzimer, A., Dauphin, Y.,
 288 Liang, P., and Vaughan, J. W. (eds.), *Advances in Neural*
 289 *Information Processing Systems*, volume 34, pp. 3309–
 290 3322. Curran Associates, Inc., 2021.
- 291
 292 Jacobs, R. A., Jordan, M. I., Nowlan, S. J., and Hinton,
 293 G. E. Adaptive mixtures of local experts. *Neural compu-*
 294 *tation*, 3(1):79–87, 1991.
- 295
 296 Jang, J., Kim, S., Lin, B. Y., Wang, Y., Hessel, J., Zettle-
 297 moyer, L., Hajishirzi, H., Choi, Y., and Ammanabrolu,
 298 P. Personalized soups: Personalized large language
 299 model alignment via post-hoc parameter merging. *arXiv*
 300 *preprint arXiv:2310.11564*, 2023.
- 301
 302 Jiang, A. Q., Sablayrolles, A., Roux, A., Mensch, A.,
 303 Savary, B., Bamford, C., Chaplot, D. S., Casas, D. d. l.,
 304 Hanna, E. B., Bressand, F., et al. Mixtral of experts.
 305 *arXiv preprint arXiv:2401.04088*, 2024.
- 306
 307 Jimenez, D. Dynamically weighted ensemble neural
 308 networks for classification. In *1998 IEEE Internation-*
 309 *al Joint Conference on Neural Networks Proceed-*
 310 *ings. IEEE World Congress on Computational Intelli-*
 311 *gence (Cat. No.98CH36227)*, volume 1, pp. 753–756
 312 vol.1, 1998. doi: 10.1109/IJCNN.1998.682375.
- 313
 314 Jordan, M. I. and Jacobs, R. A. Hierarchical mixtures of
 315 experts and the em algorithm. *Neural computation*, 6(2):
 181–214, 1994.
- 316
 317 Kelly, M., Longjohn, R., and Nottingham, K. Uci machine
 318 learning repository. URL [https://archive.ics.](https://archive.ics.uci.edu)
 319 [uci.edu](https://archive.ics.uci.edu). Accessed October 2024.
- 320
 321 Khattab, O., Singhvi, A., Maheshwari, P., Zhang, Z., San-
 322 thanam, K., Vardhamanan, S., Haq, S., Sharma, A.,
 323 Joshi, T. T., Moazam, H., Miller, H., Zaharia, M.,
 324 and Potts, C. Dspy: Compiling declarative language
 325 model calls into self-improving pipelines. *arXiv preprint*
 326 *arXiv:2310.03714*, 2023.
- 327
 328 Koh, P. W., Sagawa, S., Marklund, H., Xie, S. M., Zhang,
 329 M., Balsubramani, A., Hu, W., Yasunaga, M., Phillips,
 R. L., Gao, I., et al. Wilds: A benchmark of in-the-wild
 distribution shifts. In *International conference on ma-*
chine learning, pp. 5637–5664. PMLR, 2021.
- Krogh, A. and Vedelsby, J. Neural network ensembles,
 cross validation, and active learning. In Tesauro, G.,
 Touretzky, D., and Leen, T. (eds.), *Advances in Neural*
Information Processing Systems, volume 7. MIT Press,
 1994.
- Lakshminarayanan, B., Pritzel, A., and Blundell, C. Simple
 and scalable predictive uncertainty estimation using deep
 ensembles. *Advances in neural information processing*
systems, 30, 2017.
- Lambert, N., Pyatkin, V., Morrison, J., Miranda, L., Lin,
 B. Y., Chandu, K., Dziri, N., Kumar, S., Zick, T., Choi,
 Y., Smith, N. A., and Hajishirzi, H. Rewardbench: Eval-
 uating reward models for language modeling, 2024.
- Laskin, M., Yarats, D., Liu, H., Lee, K., Zhan, A., Lu,
 K., Cang, C., Pinto, L., and Abbeel, P. Urlb: Unsuper-
 vised reinforcement learning benchmark. *arXiv preprint*
arXiv:2110.15191, 2021.
- Lee, Y., Yao, H., and Finn, C. Diversify and disambiguate:
 Learning from underspecified data. *International Con-*
ference on Learning Representations, 2023.
- Lepikhin, D., Lee, H., Xu, Y., Chen, D., Firat, O., Huang,
 Y., Krikun, M., Shazeer, N., and Chen, Z. Gshard: Scal-
 ing giant models with conditional computation and auto-
 matic sharding. *arXiv preprint arXiv:2006.16668*, 2020.
- Li, X., Lipton, Z. C., and Leqi, L. Personalized language
 modeling from personalized human feedback. *arXiv*
preprint arXiv:2402.05133, 2024.
- Liu, S.-Y., Wang, C.-Y., Yin, H., Molchanov, P., Wang,
 Y.-C. F., Cheng, K.-T., and Chen, M.-H. Dora:
 Weight-decomposed low-rank adaptation. *arXiv preprint*
arXiv:2402.09353, 2024.
- Osband, I., Wen, Z., Asghari, S. M., Dwaracherla, V.,
 Ibrahimi, M., Lu, X., and Van Roy, B. Epistemic neural
 networks. *Advances in Neural Information Processing*
Systems, 36, 2023.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C.,
 Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A.,
 et al. Training language models to follow instructions
 with human feedback. *Advances in neural information*
processing systems, 35:27730–27744, 2022.
- Pan, A., Bhatia, K., and Steinhardt, J. The effects of
 reward misspecification: Mapping and mitigating mis-
 aligned models. *arXiv preprint arXiv:2201.03544*, 2022.

- 330 Poddar, S., Wan, Y., Ivison, H., Gupta, A., and Jaques, N.
 331 Personalizing reinforcement learning from human feed-
 332 back with variational preference learning. *arXiv preprint*
 333 *arXiv:2408.10075*, 2024.
- 334 Rafailov, R., Sharma, A., Mitchell, E., Manning, C. D., Er-
 335 mon, S., and Finn, C. Direct preference optimization:
 336 Your language model is secretly a reward model. *Ad-*
 337 *vances in Neural Information Processing Systems*, 36,
 338 2024.
- 340 Sagawa, S., Koh, P. W., Hashimoto, T. B., and Liang, P.
 341 Distributionally robust neural networks for group shifts:
 342 On the importance of regularization for worst-case gen-
 343 eralization. *arXiv preprint arXiv:1911.08731*, 2019.
- 345 Shahhosseini, M., Hu, G., and Pham, H. Optimizing en-
 346 semble weights and hyperparameters of machine learn-
 347 ing models for regression problems. *Machine Learning*
 348 *with Applications*, 7:100251, 2022.
- 349 Sharma, M., Farquhar, S., Nalisnick, E., and Rainforth, T.
 350 Do bayesian neural networks need to be fully stochastic?,
 351 2023.
- 353 Shazeer, N., Mirhoseini, A., Maziarz, K., Davis, A., Le,
 354 Q., Hinton, G., and Dean, J. Outrageously large neural
 355 networks: The sparsely-gated mixture-of-experts layer.
 356 *arXiv preprint arXiv:1701.06538*, 2017.
- 357 Shi, Y., Seely, J., Torr, P. H., Siddharth, N., Hannun, A.,
 358 Usunier, N., and Synnaeve, G. Gradient matching for do-
 359 main generalization. *arXiv preprint arXiv:2104.09937*,
 360 2021.
- 362 Singh, A., Hsu, S., Hsu, K., Mitchell, E., Ermon, S.,
 363 Hashimoto, T., Sharma, A., and Finn, C. Fspo: Few-shot
 364 preference optimization of synthetic preference data in
 365 llms elicits effective personalization to real users. *arXiv*
 366 *preprint arXiv:2502.19312*, 2025.
- 368 Siththaranjan, A., Laidlaw, C., and Hadfield-Menell, D.
 369 Distributional preference learning: Understanding and
 370 accounting for hidden context in rlhf. *arXiv preprint*
 371 *arXiv:2312.08358*, 2023.
- 372 Skalse, J., Howe, N., Krashennnikov, D., and Krueger, D.
 373 Defining and characterizing reward gaming. *Advances in*
 374 *Neural Information Processing Systems*, 35:9460–9471,
 375 2022.
- 377 Snell, C., Lee, J., Xu, K., and Kumar, A. Scaling llm test-
 378 time compute optimally can be more effective than scal-
 379 ing model parameters. *arXiv preprint arXiv:2408.03314*,
 380 2024.
- 382 Sorensen, T., Moore, J., Fisher, J., Gordon, M., Miresghal-
 383 lah, N., Rytting, C. M., Ye, A., Jiang, L., Lu, X., Dziri,
 384 N., et al. A roadmap to pluralistic alignment. *arXiv*
preprint arXiv:2402.05070, 2024.
- Sun, B., Feng, J., and Saenko, K. Correlation alignment for
 unsupervised domain adaptation. *Domain adaptation in*
computer vision applications, pp. 153–171, 2017.
- Team, G., Georgiev, P., Lei, V. I., Burnell, R., Bai, L.,
 Gulati, A., Tanzer, G., Vincent, D., Pan, Z., Wang, S.,
 et al. Gemini 1.5: Unlocking multimodal understand-
 ing across millions of tokens of context. *arXiv preprint*
arXiv:2403.05530, 2024.
- Teney, D., Abbasnejad, E., Lucey, S., and van den Hengel,
 A. Evading the simplicity bias: Training a diverse set
 of models discovers solutions with superior ood general-
 ization. In *Proceedings of the IEEE/CVF Conference on*
Computer Vision and Pattern Recognition (CVPR), pp.
 16761–16772, June 2022.
- Tipping, M. E. and Bishop, C. M. Probabilistic principal
 component analysis. *Journal of the Royal Statistical So-*
ciety Series B: Statistical Methodology, 61(3):611–622,
 1999.
- von Werra, L., Belkada, Y., Tunstall, L., Beeching, E.,
 Thrush, T., Lambert, N., Huang, S., Rasul, K., and
 Gallowédec, Q. Trl: Transformer reinforcement learn-
 ing. <https://github.com/huggingface/trl>,
 2020.
- Wirth, C., Akrou, R., Neumann, G., and Fürnkranz, J. A
 survey of preference-based reinforcement learning meth-
 ods. *Journal of Machine Learning Research*, 18(136):
 1–46, 2017.
- Wu, Z., Arora, A., Wang, Z., Geiger, A., Jurafsky, D.,
 Manning, C. D., and Potts, C. Reft: Representa-
 tion finetuning for language models. *arXiv preprint*
arXiv:2404.03592, 2024.
- Yang, R., Ding, R., Lin, Y., Zhang, H., and Zhang, T. Reg-
 ularizing hidden states enables learning generalizable re-
 ward model for llms. *arXiv preprint arXiv:2406.10216*,
 2024.
- Yao, H., Wang, Y., Li, S., Zhang, L., Liang, W., Zou, J.,
 and Finn, C. Improving out-of-distribution robustness
 via selective augmentation. In *International Conference*
on Machine Learning, pp. 25407–25437. PMLR, 2022.
- Yuksekgonul, M., Bianchi, F., Boen, J., Liu, S., Huang, Z.,
 Guestrin, C., and Zou, J. Textgrad: Automatic "differen-
 tiation" via text. *arXiv preprint arXiv:2406.07496*, 2024.
- Yuksel, S. E., Wilson, J. N., and Gader, P. D. Twenty years
 of mixture of experts. *IEEE transactions on neural net-*
works and learning systems, 23(8):1177–1193, 2012.

385 Zhang, S., Chen, Z., Chen, S., Shen, Y., Sun, Z., and
386 Gan, C. Improving reinforcement learning from human
387 feedback with efficient reward model ensemble. *arXiv*
388 *preprint arXiv:2401.16635*, 2024.

389 Zhuang, S. and Hadfield-Menell, D. Consequences of mis-
390 aligned ai. In Larochelle, H., Ranzato, M., Hadsell, R.,
391 Balcan, M., and Lin, H. (eds.), *Advances in Neural In-*
392 *formation Processing Systems*, volume 33, pp. 15763–
393 15773. Curran Associates, Inc., 2020.
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439

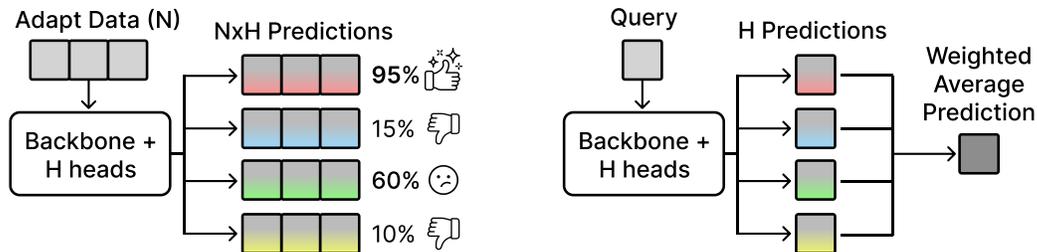


Figure 3: **Overview of HYRE.** We train multiple prediction heads on a shared backbone. (Left) At inference time, we evaluate each head on a small labeled adaptation set drawn from the target distribution. (Right) We reweight the heads according to the sum of their accuracies on the adaptation set, and use the weighted ensemble to make predictions on new inputs.

A Preliminaries

Problem setup. We consider a general supervised learning setting that includes classification, preference learning, and regression tasks. Let \mathcal{X} represent the input space and \mathcal{Y} the output space, with training distribution P_{train} and evaluation distribution P_{eval} defined over $\mathcal{X} \times \mathcal{Y}$. The training dataset $\mathcal{D}_{\text{train}} = \{(x_i, y_i)\}_{i=1}^N$ consists of N examples drawn from P_{train} . We explore few-shot adaptation settings such as chatbot personalization, where a small adaptation set $\mathcal{D}_{\text{adapt}} \sim P_{\text{eval}}$ only partially informs model performance under P_{eval} . The adaptation set $\mathcal{D}_{\text{adapt}}$ can be labeled in advance or actively queried, and is much smaller than the training set ($|\mathcal{D}_{\text{adapt}}| \ll |\mathcal{D}_{\text{train}}|$). For instance, in our main experiment, $|\mathcal{D}_{\text{adapt}}| = 16$ compared to $|\mathcal{D}_{\text{train}}| > 300,000$, with adaptation occurring near-instantly after a single forward pass through the network.

Ensemble architectures. We train an ensemble of K models f_1, \dots, f_K on the training data $\mathcal{D}_{\text{train}}$. We consider parameterizations of the ensemble that aim to represent a distribution over functions by training multiple models on the same dataset $\mathcal{D}_{\text{train}}$, ensuring diversity without computational overhead beyond training a single model. To achieve this, we employ *prior networks* (Osband et al., 2023): fixed, randomly initialized models whose outputs are added to each ensemble member’s output. This mechanism preserves diversity among ensemble members during training, even as individual models converge. We consider two computationally efficient ensemble architectures:

1. **Shared-Base Ensemble:** A single neural network that parameterizes both the prior and ensemble components by sharing a common base.
2. **Epinet:** A base network augmented by a small auxiliary network that introduces diversity via a learned index.

We train all ensemble members jointly by minimizing $\sum_{k=1}^K \mathcal{L}(f_k, \mathcal{D}_{\text{train}})$ using SGD. These architectures have negligible overhead—in our reward model experiments, 100 ensemble heads add only 550K parameters (0.03%) to the 2B-parameter Gemma backbone. Please refer to [Appendix H](#) for architectural details.

B When is Ensemble Reweighting Effective, and Why?

This section explores *the conditions under which ensemble reweighting is effective* through three illustrative examples: analyzing ensemble diversity through PCA, examining decision boundaries in classification, and comparing adaptation strategies.

Ensemble diversity reflects task ambiguity. We visualize how an ensemble’s diversity reflects the axes of task ambiguity. We consider a synthetic regression task where the training data is sampled from a Gaussian Process (GP) prior. For target inputs x_1, \dots, x_M , each ensemble member f_k produces predictions $v_k \in \mathbb{R}^M$. We perform Principal Component Analysis (PCA) on the prediction matrix $V = (v_1, \dots, v_K) \in \mathbb{R}^{M \times K}$ to yield components $u_1, \dots, u_m \in \mathbb{R}^M$ that capture the main variations between ensemble members.

Using an ensemble of 100 models trained on 7 inputs and evaluated on 1000 test inputs, we visualize the first three principal components in [Figure 4](#). Each component represents a distinct mode of variation while preserving smoothness and fit to training data. Like wavelets, these components are localized in input space and form a basis for approximating the ensemble. See [Appendix J](#) for further analysis of PCA applied to ensemble predictions.

Ensembles as diverse sharp decision boundaries. We build on an alternative interpretation of the Bradley-Terry model, where the model can be seen as representing a population of deterministic decision-makers. For items i and j with param-

Inference-Time Alignment via Hypothesis Reweighting

Method	Energy	Kin8nm	CCPP
MC Dropout	0.3033	0.6494	0.3761
Vanilla	0.1664	0.4514	0.2920
+ HYRE	0.1572 (-0.0092)	0.4498 (-0.0016)	0.2902 (-0.0018)
Epinet	0.1396	0.4823	0.3068
+ HYRE	0.1345 (-0.0051)	0.4814 (-0.0009)	0.3036 (-0.0032)
Shared-Base	0.1508	0.5316	0.2976
+ HYRE	0.1431 (-0.0077)	0.5314 (-0.0002)	0.2955 (-0.0021)

Table 4: RMSE (lower is better) on test data with distribution shifts across three UCI datasets. We compare the performance of various ensemble architectures with test-time adaptation using HYRE. We find that **for all three ensemble architectures, HYRE is consistently able to adapt to the distribution shift between training and test data.**

Model	Helpful	Harmless
Fine-Tune (Helpful)	73.03	32.59
Fine-Tune (Harmless)	32.06	73.30
Pretrained RM	68.01	52.16
Ensemble	66.34	50.90
+ HYRE (Helpful)	68.44	51.21
+ HYRE (Harmless)	64.24	57.66

Table 5: Helpful vs harmless tradeoff. To establish an upper bound on performance, we fine-tune the reward model on the helpful and harmless datasets separately. **Reweighting an ensemble model with HYRE allows us to flexibly trade off between the two desiderata.**

eters $\theta_i, \theta_j \in \mathbb{R}$, the preference probability under the Bradley-Terry model is:

$$P(i \succ j) = \frac{e^{\theta_i}}{e^{\theta_i} + e^{\theta_j}} = P(\theta_i + \epsilon_i > \theta_j + \epsilon_j), \quad (2)$$

where $\epsilon_i, \epsilon_j \sim \text{Gumbel}(0, 1)$. Rather than a single stochastic decision-maker, the model can be seen as representing a population of deterministic decision-makers. Each decision-maker is characterized by a pair (ϵ_i, ϵ_j) , and makes sharp choices based on which among $\theta_i + \epsilon_i$ and $\theta_j + \epsilon_j$ is larger. The model’s probabilistic behavior emerges from averaging across this population.

We hypothesize that diverse ensembles can learn such sharp decision boundaries from aggregate data across a population of annotators. To test this, we construct a synthetic preference learning task with conflicting labelers. We sample inputs (x_1, x_2) from $[0, 1]^2$ and generate diverse linear decision boundaries $w_1 x_1 + w_2 x_2 > 0$, with $w_1, w_2 \sim N(0, 1)$. As shown in Figure 6, our ensemble quickly adapts to new decision boundaries, outperforming single models. The average ensemble prediction matches the “average” decision-maker, while individual members capture distinct boundaries. In particular, higher diversity coefficients for the prior network yields sharper boundaries per ensemble member. In Section 3, we show this enables rapid personalization in real-world preference tasks.

HYRE outperforms fine-tuning in low-data regimes. We compare HYRE to model fine-tuning on a synthetic binary classification task. The training set contains inputs from $[0, 1]^5$ labeled as 1 and inputs from $[-1, 0]^5$ labeled as 0. The target distribution is uniform over $[-1, 1]^5$ with a random linear decision boundary. Results in Figure 5 show that HYRE outperforms fine-tuning in the low-data regime, achieving high accuracy with few queries. Fine-tuning eventually surpasses reweighting with more data due to its higher capacity. This illustrates a bias-variance tradeoff: reweighting reduces variance by restricting solutions to the ensemble’s span, providing an advantage with limited data. Additionally, HYRE requires only a single forward pass and negligible weight computation cost (1), making it especially suitable for large models and resource-constrained settings.

Interpreting HYRE as generalized Bayesian inference. The weight update in (1) can be interpreted as a form of generalized Bayesian inference (Bissiri et al., 2016). Given an initial belief state $\pi(w)$, the updated belief after observing $\mathcal{D}_{\text{adapt}}$ is:

$$\pi(w | \mathcal{D}_{\text{adapt}}) \propto \exp(-\mathcal{L}(w, \mathcal{D}_{\text{adapt}})) \pi(w), \quad (3)$$

which generalizes classical Bayesian inference by allowing arbitrary loss functions. Standard Bayes is recovered when $l(w, x)$ is the negative log-likelihood. Under mild conditions like i.i.d. sampling, these updates are consistent and coherent: they converge to the optimal weighting and yield identical posteriors whether applied incrementally or in batches. For classification tasks, using 0-1 loss instead of log-likelihood provides more stable updates by avoiding outlier dominance (Izmailov et al., 2021). This makes HYRE particularly suitable for robust adaptation with non-differentiable metrics.

C Related Work

Ensembles and mixture-of-experts. A long-standing theme in machine learning is using ensembles to improve predictive performance and uncertainty estimates when different members make independent mistakes (Krogh & Vedelsby, 1994;

Lakshminarayanan et al., 2017). This principle underlies Mixture-of-Experts (MoE) models, where a gating mechanism dynamically selects experts (Jacobs et al., 1991; Jordan & Jacobs, 1994; Yuksel et al., 2012), recently scaled to large neural networks via conditional activation (Fedus et al., 2022; Jiang et al., 2024; Lepikhin et al., 2020; Shazeer et al., 2017). Our approach diverges from these methods in a fundamental way: rather than learning a routing function during training, we perform adaptive reweighting of ensemble members at test time. Building on efficient ensemble methods with shared backbones (Osband et al., 2023), we extend prior work on dynamic ensemble weighting (Jimenez, 1998; Shahhosseini et al., 2022), though these typically focus on differentiable loss-based objectives. In contrast, our method dynamically adjusts ensemble weights based on non-differentiable evaluation metrics, allowing for more effective inference-time alignment.

Task underspecification and scalable alignment. In many machine learning tasks, the training data fails to fully define desired model behavior (D’Amour et al., 2022; Geirhos et al., 2020). This challenge intensifies under limited data or distribution shifts, where multiple hypotheses remain consistent with observations. Reinforcement learning faces similar issues: reward specification is difficult in open-ended environments, and optimizing misspecified objectives can lead to unintended behaviors (Gao et al., 2023; Pan et al., 2022; Skalse et al., 2022; Zhuang & Hadfield-Menell, 2020). Instead of fully defining a task upfront, one can collect human demonstrations or pairwise preferences, framing task specification as a cooperative game between agents and humans (Hadfield-Menell et al., 2016). Reinforcement Learning from Human Feedback (RLHF) operationalizes this idea by using user preferences to guide post-training (Christiano et al., 2017; Ouyang et al., 2022; Rafailov et al., 2024; Wirth et al., 2017), with some using ensembles (Ahmed et al., 2024; Coste et al., 2023; Zhang et al., 2024). Recent work on pluralistic alignment (Sorensen et al., 2024) uses explicit domain labels or per-user data to improve personalization (Barreto et al., 2025; Chen et al., 2024; Jang et al., 2023; Li et al., 2024; Poddar et al., 2024). However, these methods require explicit domain labels or per-user data. HYRE demonstrates that this additional information is not necessary during training: a diverse ensemble trained on aggregate data can capture ambiguity, which we can use to directly adapt to new users. Our experiments show this insight generalizes across several problem settings.

D Discussion

Our results demonstrate how efficient ensemble architectures can offer a practical path to inference-time alignment in large models. By attaching lightweight ensemble heads to a shared backbone, we can capture multiple plausible interpretations of the training distribution at negligible extra cost. Then, through a simple reweighting step that leverages just a handful of target-domain examples, the ensemble can effectively pick out the functions that align best with a new task. Our findings complement recent efforts that flexibly distribute compute at inference time (Brown et al., 2024; Snell et al., 2024). A natural next step is to close the loop by pairing our approach with a parameterization of the reward model that allows for direct behavior adjustments (Rafailov et al., 2024).

Our method currently relies on a small batch of labeled examples from the target distribution, and does not address single-sample or online streaming adaptation. Furthermore, while relying on minimal data, our reweighting still assumes that the ensemble’s functional diversity covers the new domain’s core behaviors. Extending our framework to dynamically expand or augment the ensemble as new tasks emerge is an exciting direction. Nevertheless, our results demonstrate that lightweight ensembles with inference-time reweighting offer a promising and practical approach for aligning large models at inference time.

Limitations. On the WILDS (Koh et al., 2021) benchmark, we observe limited gains over the uniform ensemble in four out of the five datasets we tested (see Table 9). We attribute this to insufficient functional diversity relevant to these specific natural distribution shifts. Thus, while HYRE significantly improves performance on personalization tasks, its effectiveness is limited on settings with more severe distribution shifts. We leave the exploration of how to best approach inference-time adaptation in such settings to future work.

E Additional Experiments

Effect of sampling strategy. In Table 6, we compare the performance of different active learning criteria for selecting adaptation data points. We consider random sampling, BALD, and entropy, measuring their performance over 0 to 40 target examples. Across the acquisition of 40 examples, active learning methods (BALD and entropy) demonstrated slightly better performance compared to random sampling. Even random sampling consistently improves performance, indicating that HYRE can be used with data collected before inference without sacrificing performance.

WILDS experiments. We evaluate a trained Shared-Base ensemble, both with and without HYRE on the WILDS-

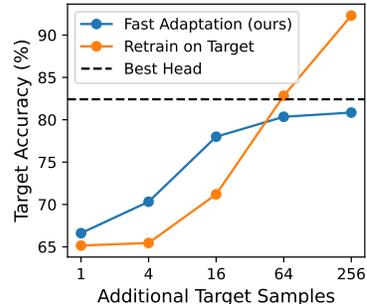
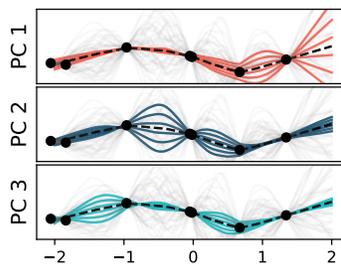
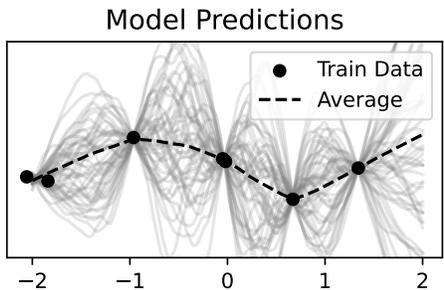
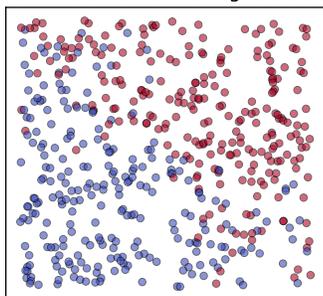


Figure 4: Principal component analysis of an ensemble of regression models. Left: Each gray line is the prediction of an ensemble member; the dashed line shows the ensemble mean. Right: The top three principal components of the ensemble’s predictions reveal distinct axes of variation in predictive behavior. **Searching among ensemble weights like HYRE acts as a strong inductive bias towards simple functions consistent with the training data.**

Figure 5: HYRE vs. fine-tuning with different amounts of adaptation data. Despite using only a single forward pass, **HYRE outperforms fine-tuning in low-data regimes.**

Train Data w/ Conflicting Labels



Ensemble Predictions

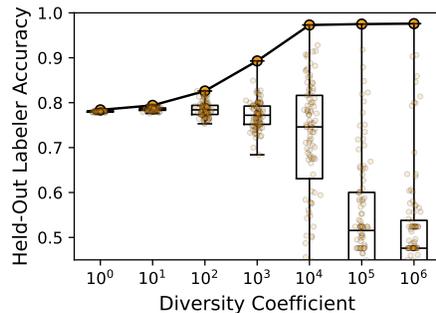
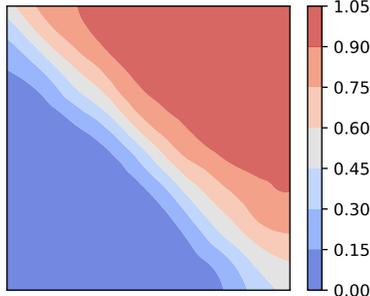


Figure 6: Ensemble behavior under label ambiguity. (Left) We simulate conflicting preferences between labelers on a synthetic dataset. (Center) Ensemble-averaged predictions approximate the consensus, smoothing over disagreements. (Right) We measure the maximum agreement between an ensemble and a held-out labeler: **increasing model diversity improves alignment with individual labelers.**

Camelyon17 dataset (Koh et al., 2021), comparing against several representative methods for OOD generalization from the official WILDS benchmark. As shown in Table 7, test-time adaptation with HYRE consistently outperforms other methods that do not use domain labels and remains competitive with LISA (Yao et al., 2022), a strong method that leverages domain labels for targeted data augmentation. We also test Shared-Base ensembles on four additional WILDS datasets (CivilComments, Amazon, FMoW, iWildCam), but did not observe further improvements from ensemble reweighting via HYRE, as detailed in Table 9. Nonetheless, training a diverse ensemble consistently improved OOD generalization in these datasets. We attribute the limited benefit of ensemble reweighting in these cases to some natural distribution shifts behaving similarly to in-distribution data in terms of task underspecification. For further discussion on the conditions that can make a single model outperform the ensemble, see Appendix B.

We further compare the performance of HYRE with few-shot fine-tuning with the same amount of adaptation data. We evaluate both HYRE and fine-tuning with $\{4, 8, 16, 32\}$ datapoints from the OOD test set. Our results in Figure 7 show that ensemble reweighting outperforms fine-tuning in the low-data regime (4 and 8) examples, and fine-tuning eventually surpasses the performance of ensemble reweighting.

F Active Learning Details

We also consider an active learning setup in which the N datapoints to label for HYRE are chosen at test time from a larger unlabeled pool of data. Rather than choosing all datapoints at once, we choose one datapoint at the time based on one of the following three criteria:

Method	N=0	N=1	N=5	N=10	N=20	N=40
HYRE + Random	84.40	85.33	86.97	87.34	88.01	88.83
HYRE + Entropy	84.40	84.25	86.73	87.54	88.60	89.76
HYRE + BALD	84.40	84.28	87.13	87.78	88.60	88.99

Table 6: Accuracies on RewardBench with different datapoint selection strategies. While active sampling methods perform slightly better, **even random sampling consistently improves performance with the HYRE reweighting process.**

Algorithm	DL	Test Acc
IRM	O	64.2 (8.1)
CORAL	O	59.5 (7.7)
Group DRO	O	68.4 (7.3)
Fish	O	74.7 (7.1)
LISA	O	77.1 (6.9)
ERM	X	70.3 (6.4)
Evading	X	73.6 (3.7)
Ensemble	X	71.5 (3.4)
Ensemble + HYRE	X	75.2 (5.3)

Table 7: Test set accuracy on Camelyon17. HYRE achieves competitive performance without using domain labels (DL).

- **Entropy** (classification): $H\left(\sum_{h=1}^H w_h f_h(x)\right)$. This criterion selects datapoints where the weighted ensemble is most uncertain, promoting the exploration of ambiguous regions.
- **BALD** (classification): $H\left(\sum_{i=1}^H w_i f_i(x)\right) - \sum_{i=1}^H w_i H(f_i(x))$. BALD considers both ensemble uncertainty and disagreement among members, balancing exploration and exploitation (Gal et al., 2017; Hounsby et al., 2011).
- **Variance** (regression): $\sum_{i=1}^H w_i (f_i(x) - \bar{f}(x))^2$, where $\bar{f}(x) = \sum_{i=1}^H w_i f_i(x)$. This criterion focuses on points where ensemble predictions have the highest variance, which is a good indicator of uncertainty in regression tasks.

Each of these criteria can be computed quickly. Because the belief states w has a closed-form update that can be computed very quickly, we can efficiently recompute the next best data point after each active label query.

We note that the first criterion (Entropy) does not distinguish between so-called aleatoric uncertainty and epistemic uncertainty. Therefore, this criterion is susceptible to the “noisy TV problem”, where an agent fixates on a source of uncertainty that cannot be resolved (Burda et al., 2018; Laskin et al., 2021). In practice, we find that HYRE is robust to the choice of active learning criterion, and even random selection is effective at adapting to the target distribution.

G Experimental Details

Unless specified otherwise, we use the following configuration for the ensemble networks. We use an ensemble of 100 models. The learnable and prior networks are each a one-hidden-layer MLP with 128 units. For the epinet, the epistemic index is 10-dimensional. For ensemble reweighting via HYRE, we use 32 examples from the target dataset, actively queried based on the BALD (classification) or Variance (regression) criterion. We found that final performance is not very sensitive to the choice of active learning criterion, and even random sampling resulted in consistent benefits.

WILDS. We closely follow the reference WILDS implementation for each dataset (Koh et al., 2021), including the choice of backbone, learning rate, and weight decay. We briefly describe the baseline methods used in our experiments:

- **CORrelation ALignment** (Sun et al., 2017, CORAL): CORAL is an unsupervised domain adaptation method that aligns the second-order statistics (covariances) of source and target feature distributions.

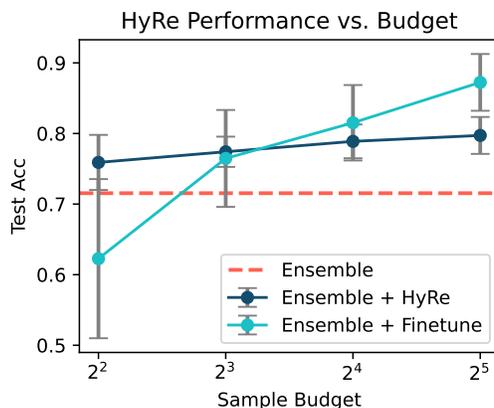


Figure 7: Comparison of HYRE and few-shot fine-tuning on the Camelyon17 OOD test set. HYRE outperforms fine-tuning in the low-data regime despite requiring significantly less computational cost.

- **Invariant Risk Minimization** (Arjovsky et al., 2019, IRM): IRM aims to learn data representations that capture invariant correlations across multiple training distributions.
- **Group Distributionally Robust Optimization** (Sagawa et al., 2019, Group DRO): Group DRO seeks to minimize the worst-case training loss over predefined groups within the data.
- **Fish** (Shi et al., 2021): Fish is a domain generalization technique that approximates inter-domain gradient matching by maximizing the inner product between gradients from different domains.
- **LISA** (Yao et al., 2022): LISA builds on MixUp and selectively interpolates data samples to achieve domain invariance.

LLM Preference Learning We finetune three reward model checkpoints (Yang et al., 2024):

- <https://huggingface.co/Ray2333/GRM-Gemma-2B-rewardmodel-ft>
- <https://huggingface.co/Ray2333/GRM-Gemma2-2B-rewardmodel-ft>
- <https://huggingface.co/Skywork/Skywork-Reward-Llama-3.1-8B-v0.2>

Our ensemble architecture uses these networks as the backbone, and small MLPs for the learnable and prior networks which take the backbone’s final embedding as input. We use the TRL codebase for reward model training (von Werra et al., 2020). We train with bfloat16 mixed precision. We use a learning rate of 0.0001, no weight decay, a batch size of 16, and train for 5000 steps. We consider four collections of preference datasets:

- **Elix** (Singh et al., 2025) is inspired by the “Explain like I’m 5” subreddit. It consists of questions answered at five educational levels: elementary, middle, high, college, and expert. Preference pairs are created by scoring how different pairs of GPT-4 generated responses meet the expected comprehension at each level.
- **RewardBench** (Lambert et al., 2024) is a suite of 27 preference datasets designed to test reward models on a broad spectrum of tasks, including chat quality, safety, reasoning, coding, and refusal handling. In our aggregate results Figure 2, we drop datasets with less than 100 examples. In our RewardBench experiments Table 8, we use all datasets to ensure a fair comparison with existing methods.
- **PERSONA** (Castricato et al., 2024) contains preference data derived from a collection of synthetic personas with diverse demographic attributes and values. We sample 10 personas and treat each as a target distribution. Further details are in Appendix K.
- **Anthropic HH** (Bai et al., 2022) contains human-labeled preferences focused on helpfulness and harmlessness. We use the helpfulness-base and harmlessness-base splits as evaluation distributions to measure the tradeoff between the two objectives.

Inference-Time Alignment via Hypothesis Reweighting

Model	Type	Overall	Chat	Chat Hard	Safety	Reasoning
Mixtral-8x7B-Instruct-v0.1	DPO	77.6	95.0	64.0	72.6	78.7
LLaMA-3-Tulu-2-DPO-70B	DPO	77.2	96.4	57.5	74.9	80.2
Tulu-2-DPO-13B	DPO	76.7	95.8	58.3	79.5	73.2
Tulu-2-DPO-70B	DPO	79.1	97.5	60.5	84.5	74.1
StableLM-2-12B-Chat	DPO	79.9	96.6	55.5	78.1	89.4
Claude-3 Sonnet (June 2024)	Gen	84.2	96.4	74.0	81.6	84.7
GPT-4 (May 2024)	Gen	84.6	96.6	70.4	86.5	84.9
GPT-4 (Aug 2024)	Gen	86.7	96.1	76.1	88.1	86.6
Gemini-1.5-Pro-0924	Gen	86.8	94.1	77.0	85.8	90.2
Skywork-Reward-Gemma-2-27B	Seq	94.3	96.1	89.9	93.0	98.1
INF-ORM-Llama3.1-70B	Seq	95.1	96.6	91.0	93.6	99.1
GRM-Gemma-2B	Seq	84.5	89.4	75.2	84.5	88.8
+ Ours (uniform)	Seq	84.5	88.6	72.9	83.7	89.8
+ Ours (N=1)	Seq + HYRE	85.3	88.5	72.7	85.5	91.4
+ Ours (N=5)	Seq + HYRE	86.4	90.3	72.6	89.1	91.4
+ Ours (N=10)	Seq + HYRE	87.2	90.4	72.5	90.0	92.3
+ Ours (best head oracle)*	Seq + Oracle	88.6	91.1	78.1	91.9	92.3
+ Ours (best weight oracle)*	Seq + Oracle	90.0	92.3	81.8	92.5	93.1
GRM-Gemma2-2B	Seq	88.4	93.0	77.2	92.2	91.2
+ Ours (uniform)	Seq	87.1	96.4	73.1	87.4	89.8
+ Ours (N=1)	Seq + HYRE	86.5	92.4	71.5	85.1	92.5
+ Ours (N=5)	Seq + HYRE	88.5	95.0	72.5	90.3	93.1
+ Ours (N=10)	Seq + HYRE	89.7	96.4	74.7	92.4	93.5
+ Ours (best head oracle)*	Seq + Oracle	91.8	97.2	80.0	96.2	94.2
+ Ours (best weight oracle)*	Seq + Oracle	93.1	98.3	83.4	96.7	94.9
Skywork-Llama-3.1-8B	Seq	94.0	94.7	88.6	92.7	96.7
+ Ours (uniform)	Seq	94.0	95.0	87.2	93.0	96.8
+ Ours (N=1)	Seq + HYRE	94.3	95.2	87.8	93.0	97.5
+ Ours (N=5)	Seq + HYRE	94.7	95.5	88.6	93.2	97.8
+ Ours (N=10)	Seq + HYRE	95.0	95.9	89.3	93.5	97.9
+ Ours (best head oracle)*	Seq + Oracle	96.4	98.3	91.2	95.7	98.4
+ Ours (best weight oracle)*	Seq + Oracle	97.2	99.2	93.0	96.5	98.8

* Oracle methods show an upper bound on performance, using the test set.

Table 8: Accuracy across tasks in RewardBench. We report overall performance and breakdowns by task category for all models. **HYRE improves upon the state-of-the-art models at the 2B and 8B parameter scales with as few as 1-5 labeled samples per distribution.**

For the few-shot prompting experiments, we use GPT-4o-mini. For each number of “shots” $N \in \{0, 1, 5, 10, 20, 40, 80\}$, we sample 1000 examples from the target distribution and use them to prompt GPT-4o-mini.

H Diverse Ensemble Architectures

We describe the diverse ensemble architectures used in our experiments. Each architecture is designed to parameterize an ensemble of H models, whose outputs are later combined to form an ensemble prediction. The key goal of these architectures is to produce diverse predictions across the ensemble at a low computational cost.

All architectures are trained end-to-end by minimizing the sum of a standard loss function (cross-entropy for classification,

Algorithm	DL	CivilComments	Amazon	FMoW	iWildCam
		Worst-Group Acc	10% Acc	Worst-Reg Acc	Macro F1
IRM	O	66.3 (2.1)	52.4 (0.8)	32.8 (2.09)	15.1 (4.9)
IRMX	O	73.4 (1.4)	-	33.7 (0.95)	26.7 (1.1)
IRMX (PAIR)	O	74.2 (1.4)	-	35.4 (1.3)	27.9 (0.9)
CORAL	O	65.6 (1.3)	52.9 (0.8)	32.8 (0.66)	32.7 (0.2)
Group DRO	O	70.0 (2.0)	53.3 (0.0)	31.1 (1.66)	23.8 (2.0)
DFR	O	72.5 (0.9)	-	42.8 (0.42)	-
Fish	O	75.3 (0.6)	53.3 (0.0)	34.6 (0.18)	22.0 (1.8)
LISA	O	72.9 (1.0)	54.7 (0.0)	35.5 (0.81)	-
ERM	X	56.0 (3.6)	53.8 (0.8)	31.3 (0.17)	30.8 (1.3)
Shared-Base	X	58.1 (2.2)	54.2 (0.6)	32.8 (0.4)	30.9 (0.8)
Shared-Base + HYRE	X	58.1 (0.2)	54.2 (0.6)	32.8 (0.4)	31.0 (0.8)

Table 9: Performance on additional WILDS benchmark datasets. The DL column indicates whether the algorithm uses domain labels. Using a Shared-Base ensemble consistently results in gains in OOD generalization metrics over prior methods. However, we observe no further benefits from reweighting the ensemble via HYRE on these datasets.

MSE for regression) over all ensemble members:

$$\sum_{h=1}^H \mathcal{L}(f_h(x), y). \tag{4}$$

Here, x is an input example, y is the true label, and f^i is the i -th ensemble member. While each individual model minimizes the training loss, we want the ensemble members to extrapolate to unseen data in diverse ways. The specific ensemble parameterizations, which we describe below, are designed to achieve this goal.

H.1 Vanilla Ensemble

A vanilla ensemble consists of H independently initialized and trained neural networks with identical architectures. Each network f_h takes an input x and produces an output $f_h(x)$. No parameters are shared. While simple to implement, this approach scales poorly as H increases since both memory and computation scale linearly with H .

H.2 Shared-Base Ensemble

We propose a scalable neural network architecture that can represent thousands of diverse ensemble members. The network outputs H real-valued predictions in parallel, with the output space being \mathbb{R}^H . The architecture comprises a frozen prior network f_p and a learnable network f_θ , both of which produce outputs of shape \mathbb{R}^H . Although the architectures of f_p and f_θ are identical in our experiments, this is not a requirement.

For a given input x , the network output is

$$f^p(z) + f^\theta(z) = \begin{bmatrix} f_1^p(z) + f_1^\theta(z) \\ f_2^p(z) + f_2^\theta(z) \\ \vdots \\ f_H^p(z) + f_H^\theta(z) \end{bmatrix} \in \mathbb{R}^H \tag{5}$$

where each prediction $f_i^p(z) + f_i^\theta(z)$ is compared against the ground-truth label y . The parameters of f^p are fixed at initialization and do not change during training; the parameters of f^θ are learnable.

Using the frozen prior network f^p is crucial to the diversity in this architecture. If we were to only train f^θ , the ensemble of the H predictions would have low diversity due to co-adaptation. To understand why this architecture produces a diverse ensemble, note that each learnable head solves a shifted task determined by the corresponding prior network head. Since we undo this shifting when producing the final prediction, we can view the different learnable heads as solving a different yet equivalent task.

H.3 Epinet

The epinet architecture combines a base model $f^{\text{base}} : \mathcal{X} \rightarrow \mathbb{R}^K$ with an epistemic network $f^{\text{epi}} : \mathcal{Z} \times \mathbb{R}^{d_{\text{firs}}} \times \mathcal{X} \rightarrow \mathbb{R}^K$. The base model can be any regular neural network, including a large pretrained model, and is used to extract features through a feature extractor $\phi : \mathcal{X} \rightarrow \mathbb{R}^{d_{\text{firs}}}$. Here, d_{firs} is the dimension of the extracted intermediate representations.

The epistemic network (epinet) is composed of two parts:

- A frozen prior network $f^{\text{epi-frozen}} : \mathcal{X} \rightarrow \mathbb{R}^{1, \dots, d_{\text{index}} \times K}$. The parameters of this network are fixed at initialization and do not change during training.
- A trainable network $f^{\text{epi-trainable}} : \mathcal{Z} \times \mathbb{R}^{d_{\text{firs}}} \times \mathcal{X} \rightarrow \mathbb{R}^K$.

Given an epistemic index $z \in \mathbb{R}^d$ and input $x \in \mathcal{X}$, we compute the model output as:

$$f(z, x) = f^{\text{base}}(x) + v f^{\text{epi-frozen}}(x) \cdot z + f^{\text{epi-trainable}}(z, \phi(x), x) \cdot z \quad (6)$$

where \cdot is the dot product and $v \in (0, \infty)$ is the so-called prior scale. At each step, we sample multiple epistemic indices z to form an ensemble, i.e., $f_1(x), \dots, f_H(x) = f(z_1, x), \dots, f(z_H, x)$. This architecture efficiently generates diverse predictions by sampling different epistemic indices z while leveraging a potentially large pretrained base model.

I Repulsion vs Random Priors for Diversity

A line of prior work use repulsion for enforcing diversity between ensemble members. The high-level idea is to add a regularization term to the loss function that is minimized when the ensemble members are sufficiently “different” according to some distance metric. For example, [Teney et al. \(2022\)](#) uses a repulsion term that maximizes the cosine distance between the gradient of each ensemble member, and [Lee et al. \(2023\)](#) maximizes the mutual information of ensemble predictions on OOD inputs. While these techniques have seen success in certain settings, our early experiments indicate that such explicit regularization often results in a suboptimal ensemble. The repulsion term can overpower the learning signal in the training data, leading to ensemble members that are diverse but inaccurate.

In contrast, diversification via random priors ([Osband et al., 2023](#)) provides a more balanced approach. The key idea is to initialize each ensemble member with a different random prior function which is fixed throughout training. This introduces diversity from the start without explicitly optimizing for it during training. This approach maintains diversity without sacrificing accuracy on the training data, and the degree of diversification is easily controlled by scaling the prior functions.

J Function-Space Dimensionality Reduction

Here, we expand on the idea of PCA on ensemble predictions. A central challenge with large model ensembles is understanding the commonalities and differences among the individual models. The high-level idea is that PCA applied to ensemble predictions reveals the major direction of variation within an ensemble of models. This dimensionality reduction allows us to clearly interpret model behaviors and identify groups of related datapoints. Additionally, PCA enables the generation of new functions with similar statistical properties by parameterizing a low-rank Gaussian distribution in the joint prediction space, which we can sample from.

J.1 Motivating Example

Consider three models f_1, \dots, f_3 and five inputs z_1, \dots, z_5 . Denoting each model’s predicted probability for an input as $p_{nh} = \sigma(f_h(z_n)) \in [0, 1]$, assume that the matrix of predictions is

$$\begin{pmatrix} p_{11} & p_{12} & p_{13} & p_{14} & p_{15} \\ p_{21} & p_{22} & p_{23} & p_{24} & p_{25} \\ p_{31} & p_{32} & p_{33} & p_{34} & p_{35} \end{pmatrix} = \begin{pmatrix} 1 & 0 & 1 & 0 & 1/2 \\ 0 & 1 & 1/2 & 1/2 & 1/2 \\ 1/2 & 1/2 & 0 & 1 & 1/2 \end{pmatrix}. \quad (7)$$

Each row of this matrix shows one model’s prediction on the entire pool of inputs, and each column shows every model’s prediction on a single input. We can analyze such a matrix of predictions on three levels, each revealing increasing amounts of structure within the ensemble:

Level 1: Per-sample ensemble uncertainty. We can first compute the average prediction $\bar{p}(x) = \frac{1}{H} \sum_h p_{nh}$ for each datapoint. For the predictions in (7), the average prediction is $\bar{p}(x) = 1/2$ for every input x , and thus the collection of models may be viewed as equally uncertain about each of the 5 inputs. This is the measure of ensemble uncertainty commonly used for ensembles (Lakshminarayanan et al., 2017).

Level 2: Per-sample disagreement. We can further account for the amount of disagreement among ensemble members for each datapoint. Note that for the four inputs z_1, z_2, z_3, z_4 , there is strong disagreement between two functions where one predicts 0 and the other predicts 1. This is not true of z_5 , where all functions predict $1/2$. Uncertainty metrics that take disagreement into account, such as the BALD criterion (Houlsby et al., 2011), will reveal that the ensemble is more uncertain about z_1, z_2, z_3, z_4 than it is about z_5 .

Level 3: Joint predictions. First, note that the two approaches above discard all information about which ensemble member made which individual prediction for a given input, by (1) averaging all predictions or (2) considering only the unordered set of predictions. There is additional structure to the differences among ensemble members that we can extract by considering the joint predictions, i.e., viewing each column of (7) as an object in itself. The pair of inputs (z_1, z_2) are closely related since they deviate from the ensemble prediction in the same “direction” in the joint prediction space (\mathbb{R}^H). We can make the same observation about the pair (z_3, z_4) . To see this structure more clearly, consider the matrix of deviations from the ensemble prediction $\delta_{nh} = p_{nh} - \frac{1}{H} \sum_h p_{nh}$:

$$\begin{pmatrix} \delta_{11} & \delta_{12} & \delta_{13} & \delta_{14} & \delta_{15} \\ \delta_{21} & \delta_{22} & \delta_{23} & \delta_{24} & \delta_{25} \\ \delta_{31} & \delta_{32} & \delta_{33} & \delta_{34} & \delta_{35} \end{pmatrix} = \frac{1}{2} \begin{pmatrix} 1 & -1 & 1 & -1 & 0 \\ -1 & 1 & 0 & 0 & 0 \\ 0 & 0 & -1 & 1 & 0 \end{pmatrix}. \quad (8)$$

This clearly shows that the vector of joint deviations $(\delta_{11}, \delta_{12}, \delta_{13})$ is the negative of that of $(\delta_{21}, \delta_{22}, \delta_{23})$. More generally, we can view the vector of deviations $(\delta_{1n}, \delta_{2n}, \delta_{3n})$ as a representation of the datapoint z_n in the joint prediction space. In this sense, the matrix of predictions $\{p_{nh}\}$ can be explained by the mean prediction 0.5 for each datapoint, together with two factors of variation $(1, -1, 0)$ and $(1, 0, -1)$ appropriately applied to each input. We next describe how to automatically extract such consistent high-level factors in an ensemble from the matrix of predictions.

J.2 PCA on Ensemble Predictions

We propose to apply PCA to the $H \times N$ matrix of residual predictions to obtain P principal components. Each principle component is a vector of size H that captures the orthogonal factors of variation in how ensemble members extrapolated from the training data. Given a set of weights w_1, \dots, w_P over principal components, we can “reconstruct” a set of joint predictions as

$$p(x) = \bar{p}(x) + (w_1 \quad \dots \quad w_P) \begin{pmatrix} c_{11} & \dots & c_{1H} \\ c_{21} & \dots & c_{2H} \\ \vdots & \ddots & \vdots \\ c_{P1} & \dots & c_{PH} \end{pmatrix} \begin{pmatrix} p_1(x) - \bar{p}(x) \\ p_2(x) - \bar{p}(x) \\ \vdots \\ p_H(x) - \bar{p}(x) \end{pmatrix}, \quad (9)$$

where we denote the mean prediction as $\bar{p}(x) = \frac{1}{H} \sum_h p_{nh}$ and the P principal components as $C \in \mathbb{R}^{P \times H}$.

We highlight two known interpretations of PCA that have interesting implications for our goal of summarizing ensemble predictions:

Maximum mutual information / variance after projection. PCA finds the linear projection $y = w^\top x$ with unit vector w that achieves maximum mutual information $I(x; y)$, or equivalently, maximum variance $\text{Var}(y)$. Each principal component finds the linear combination of ensemble members that preserves the most information about the set of joint ensemble predictions. This is closely related to the disagreement term in Bayesian active learning (Houlsby et al., 2011).

Factor model. The principal components are maximum likelihood parameters under a linear Gaussian factor model of the data (Tipping & Bishop, 1999). Indeed, we can view our principal components as orthogonal modifications to the mean prediction $\bar{p}(x)$. The distribution of ensemble members is closely approximated by “reconstructed predictions” (9), where $z_{1:P} \sim \mathcal{N}(0, I^P)$. We can view each principal component as a consistent high-level direction of functional variation in which the training data provided insufficient information.

K PERSONA Dataset Details

Below, we list the personas used in our PERSONA (Castricato et al., 2024) experiments. The dataset includes 1000 personas in total, each with 200 preference pairs. We subsampled 10 personas from the original dataset of 1000, ensuring a diverse set of backgrounds, ages, and lifestyles.

Persona 1. Age: 1. Sex: Male. Race: White alone. Ancestry: Irish. Household language: English only. Education: Not applicable. Employment status: Not applicable. Class of worker: Not applicable. Industry category: Not applicable. Occupation category: Not applicable. Detailed job description: Not applicable. Income: Not applicable. Marital status: Too young to be married. Household type: Cohabiting couple household with children of the householder less than 18. Family presence and age: With related children under 5 years only. Place of birth: Missouri/MO. Citizenship: Born in the United States. Veteran status: Not applicable. Disability: None. Health insurance: With health insurance coverage. Fertility: Not applicable. Hearing difficulty: None. Vision difficulty: None. Cognitive difficulty: None. Ability to speak english: Not applicable. Big five scores: Openness: High, Conscientiousness: High, Extraversion: Low, Agreeableness: Extremely High, Neuroticism: Extremely Low. Defining quirks: Loves to play with his food. Mannerisms: Waves hands when excited. Personal time: Spends most of his time playing, sleeping, and learning to walk. Lifestyle: Lives a carefree and playful lifestyle. Ideology: Not applicable. Political views: Not applicable. Religion: Other Christian.

Persona 2. Age: 11. Sex: Male. Race: White alone. Ancestry: Irish. Household language: English only. Education: Grade 4. Employment status: Unemployed. Class of worker: Not applicable. Industry category: Not applicable. Occupation category: Not applicable. Detailed job description: Student. Income: 0. Marital status: Never married or under 15 years old. Household type: Cohabiting couple household with children of the householder less than 18. Family presence and age: With related children 5 to 17 years only. Place of birth: Louisiana/LA. Citizenship: Born in the United States. Veteran status: Not applicable. Disability: None. Health insurance: With health insurance coverage. Big five scores: Openness: Low, Conscientiousness: Low, Extraversion: High, Agreeableness: High, Neuroticism: Average. defining quirks: Loves to draw and create stories. Mannerisms: Often seen doodling or daydreaming. Personal time: Spends free time drawing or playing video games. Lifestyle: Active and playful, enjoys school and spending time with friends. Ideology: Undeveloped. Political views: Undeveloped. Religion: Religiously Unaffiliated.

Persona 3. Age: 19. Sex: Male. Race: Asian Indian alone. Ancestry: Indian. Household language: Hindi. Education: 1 or more years of college credit, no degree. Employment status: Not in labor force. Class of worker: Not Applicable. Industry category: Not Applicable. Occupation category: Not Applicable. Detailed job description: Not Applicable. Income: -60000.0. Marital status: Never married or under 15 years old. Household type: Living with parents. Family presence and age: Living with two parents. Place of birth: India. Citizenship: Not a U.S. citizen. Veteran status: Non-Veteran. Disability: None. Health insurance: With health insurance coverage. Big five scores: Openness: Average, Conscientiousness: High, Extraversion: Extremely Low, Agreeableness: Extremely High, Neuroticism: Extremely Low. defining quirks: Passionate about music. Mannerisms: Expressive hand gestures when speaking. Personal time: Practicing music or studying. Lifestyle: Student and Music Enthusiast. Ideology: Liberal. Political views: Liberal. Religion: Other Christian.

Persona 4. Age: 29. Sex: Female. Race: Laotian alone. Ancestry: Laotian. Household language: Asian and Pacific Island languages. Education: Some college, but less than 1 year. Employment status: Armed forces, at work. Class of worker: Federal government employee. Industry category: MIL-U.S. Navy. Occupation category: MIL-Military Enlisted Tactical Operations And Air/Weapons Specialists And Crew Members. Detailed job description: Maintains and operates tactical weapons systems. Income: 81000.0. Marital status: Married. Household type: Married couple household with children of the householder less than 18. Family presence and age: With related children 5 to 17 years only. Place of birth: California/CA. Citizenship: Born in the United States. Veteran status: Now on active duty. Disability: None. Health insurance: With health insurance coverage. Big five scores: Openness: Average, Conscientiousness: High, Extraversion: Average, Agreeableness: High, Neuroticism: Average. Defining quirks: Collects military memorabilia. Mannerisms: Frequently uses military jargon. Personal time: Spends time with family and collecting military memorabilia. Lifestyle: Disciplined and active. Ideology: Conservative. Political views: Republican. Religion: Protestant.

Persona 5. Age: 36. Sex: Female. Race: Some Other Race alone. Ancestry: Hispanic. Household language: English. Education: Regular high school diploma. Employment status: Civilian employed, at work. Class of worker: Employee of a private for-profit company or business, or of an individual, for wages, salary, or commissions. Industry category: FIN-Insurance Carriers. Occupation category: OFF-Insurance Claims And Policy Processing Clerks. Detailed job description: Processes insurance claims and policies. Income: 182000.0. Marital status: Married. Household type: Married couple household with children of the householder less than 18. Family presence and age: With related children under 5 years

1045 only. Place of birth: New Mexico/NM. Citizenship: Born in the United States. veteran status: Non-Veteran Disability:
 1046 None. Health insurance: With health insurance coverage. Big five scores: Openness: Extremely Low, Conscientiousness:
 1047 Extremely High, Extraversion: Extremely High, Agreeableness: High, Neuroticism: Average. Defining quirks: Enjoys
 1048 bird-watching. Mannerisms: Often taps foot when thinking. Personal time: Spends free time with family or in nature.
 1049 Lifestyle: Active and family-oriented. Ideology: Conservative. Political views: Republican. Religion: Other Christian.

1050 **Persona 6.** Age: 44. Sex: Female. Race: Black or African American alone. Ancestry: Haitian. household language:
 1051 Other Indo-European languages education: Associate's degree Employment status: Civilian employed, at work. Class
 1052 of worker: Employee of a private not-for-profit, tax-exempt, or charitable organization. Industry category: FIN-Banking
 1053 And Related Activities. Occupation category: OFF-Tellers. Detailed job description: Handles customer transactions at the
 1054 bank, including deposits, withdrawals, and loan payments. Income: 40000.0. Marital status: Separated. Household type:
 1055 Female householder, no spouse/partner present, with children of the householder less than 18. Family presence and age:
 1056 With related children 5 to 17 years only. Place of birth: Haiti. Citizenship: Not a U.S. citizen. Veteran status: Non-Veteran.
 1057 Disability: None. Health insurance: With health insurance coverage. Big five scores: Openness: High, Conscientiousness:
 1058 Extremely Low, Extraversion: Average, Agreeableness: Average, Neuroticism: Extremely Low. Defining quirks: Loves
 1059 to cook Haitian cuisine. Mannerisms: Often taps her foot when stressed. Personal time: Taking care of her children,
 1060 Pursuing further education. Lifestyle: Busy, Family-oriented. Ideology: Egalitarian. Political views: Democrat. Religion:
 1061 Protestant.

1063 **Persona 7.** Age: 52. Sex: Female. Race: Korean alone. Ancestry: Korean. Household language: Asian and Pacific Island
 1064 languages. Education: Regular high school diploma. Employment status: Civilian employed, at work. Class of worker:
 1065 State government employee. Industry category: ENT-Restaurants And Other Food Services. Occupation category: EAT-
 1066 First-Line Supervisors Of Food Preparation And Serving Workers. Detailed job description: Supervises food preparation
 1067 and serving workers in a state government facility. Income: 133900.0. Marital status: Married. Household type: Married
 1068 couple household, no children of the householder less than 18. Family presence and age: No related children. Place of
 1069 birth: Korea. Citizenship: U.S. citizen by naturalization. Veteran status: Non-Veteran. Disability: None. Health insurance:
 1070 With health insurance coverage. big five scores: Openness: Average, Conscientiousness: Extremely High, Extraversion:
 1071 Extremely Low, Agreeableness: Extremely Low, Neuroticism: Average defining quirks: Deep love for literature and
 1072 reading Mannerisms: Constantly adjusts her glasses. Personal time: Spends free time reading or engaging in community
 1073 activism. Lifestyle: Quiet and community-oriented. Ideology: Liberal. Political views: Democratic. Religion: Protestant.

1074 **Persona 8.** Age: 58. Sex: Male. Race: White. Ancestry: Scottish. Household language: English. Education: Bachelor's
 1075 Degree. Employment status: Employed. Class of worker: Private. industry category: Investigation And Security Ser-
 1076 vices Occupation category: Sales Manager. Detailed job description: Oversees sales teams, sets sales goals, and develops
 1077 strategies to achieve these goals. Income: 198200. Marital status: Married. Household type: Married couple household,
 1078 no children under 18. Family presence and age: No related children. Place of birth: Florida. Citizenship: US Citizen.
 1079 veteran status: Non-Veteran Disability: With a disability. Health insurance: With health insurance coverage. Big five
 1080 scores: Openness: High, Conscientiousness: Extremely High, Extraversion: Average, Agreeableness: Average, Neuroti-
 1081 cism: Average. Defining quirks: Keen interest in security technology and crime novels. mannerisms: Constantly checks
 1082 his surroundings Personal time: Researching the latest security technologies or enjoying a round of golf. Lifestyle: Active
 1083 and health-conscious. Ideology: Conservative. Political views: Republican. Religion: Catholic.

1085 **Persona 9.** Age: 65. Sex: Female. Race: White alone. Ancestry: Italian. Household language: Other Indo-European
 1086 languages. Education: Master's degree. Employment status: Civilian employed, at work. Class of worker: Self-employed
 1087 in own incorporated business, professional practice or farm. Industry category: ENT-Traveler Accommodation. Occu-
 1088 pation category: FIN-Accountants And Auditors. Detailed job description: Manages financial records and tax data for
 1089 her own travel accommodation business. Income: 188600.0. Marital status: Married. Household type: Married couple
 1090 household, no children of the householder less than 18. Family presence and age: No related children. Place of birth:
 1091 Delaware/DE. Citizenship: Born in the United States. Veteran status: Non-veteran. Disability: None. Health insurance:
 1092 With health insurance coverage. ability to speak english: Well. Big five scores: Openness: Average, Conscientiousness:
 1093 Low, Extraversion: Low, Agreeableness: Average, Neuroticism: Extremely High. Defining quirks: Has an extensive col-
 1094 lection of vintage travel posters. Mannerisms: Tends to use Italian phrases in conversation. Personal time: Spends her
 1095 free time exploring new places, trying new cuisines, and learning about different cultures. Lifestyle: Leads a busy lifestyle
 1096 managing her business, but always finds time for her passion for travel and culture. Ideology: Believes in the importance
 1097 of understanding and appreciating different cultures. Political views: Liberal. Religion: Protestant.

1098
 1099

Inference-Time Alignment via Hypothesis Reweighting

1100 **Persona 10.** Age: 75. Sex: Female. Race: White alone. ancestry: Scottish Household language: English only. Education:
1101 Professional degree beyond a bachelor's degree. Employment status: Not in labor force. Class of worker: Retired. Industry
1102 category: Healthcare. Occupation category: Doctor. Detailed job description: Retired pediatrician. Income: 98000.0.
1103 Marital status: Never married. Household type: Female householder, no spouse/partner present, living alone. Family
1104 presence and age: No family. Place of birth: Massachusetts/MA. citizenship: Born in the United States veteran status:
1105 Non-Veteran Disability: None. Health insurance: With health insurance coverage. Big five scores: Openness: Average,
1106 Conscientiousness: Average, Extraversion: High, Agreeableness: Extremely High, Neuroticism: Average. Defining quirks:
1107 Enjoys cooking traditional Scottish meals. Mannerisms: Often hums traditional Scottish tunes. Personal time: Spends free
1108 time volunteering at the local church and community center. Lifestyle: Active but relaxed, with a focus on maintaining
1109 health and staying involved in the community. Ideology: Conservative. Political views: Republican. Religion: Catholic.

1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154

1155
 1156
 1157
 1158
 1159
 1160
 1161
 1162
 1163
 1164
 1165
 1166
 1167
 1168
 1169
 1170
 1171
 1172
 1173
 1174
 1175
 1176
 1177
 1178
 1179
 1180
 1181
 1182
 1183
 1184
 1185
 1186
 1187
 1188
 1189
 1190
 1191
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202
 1203
 1204
 1205
 1206
 1207
 1208
 1209

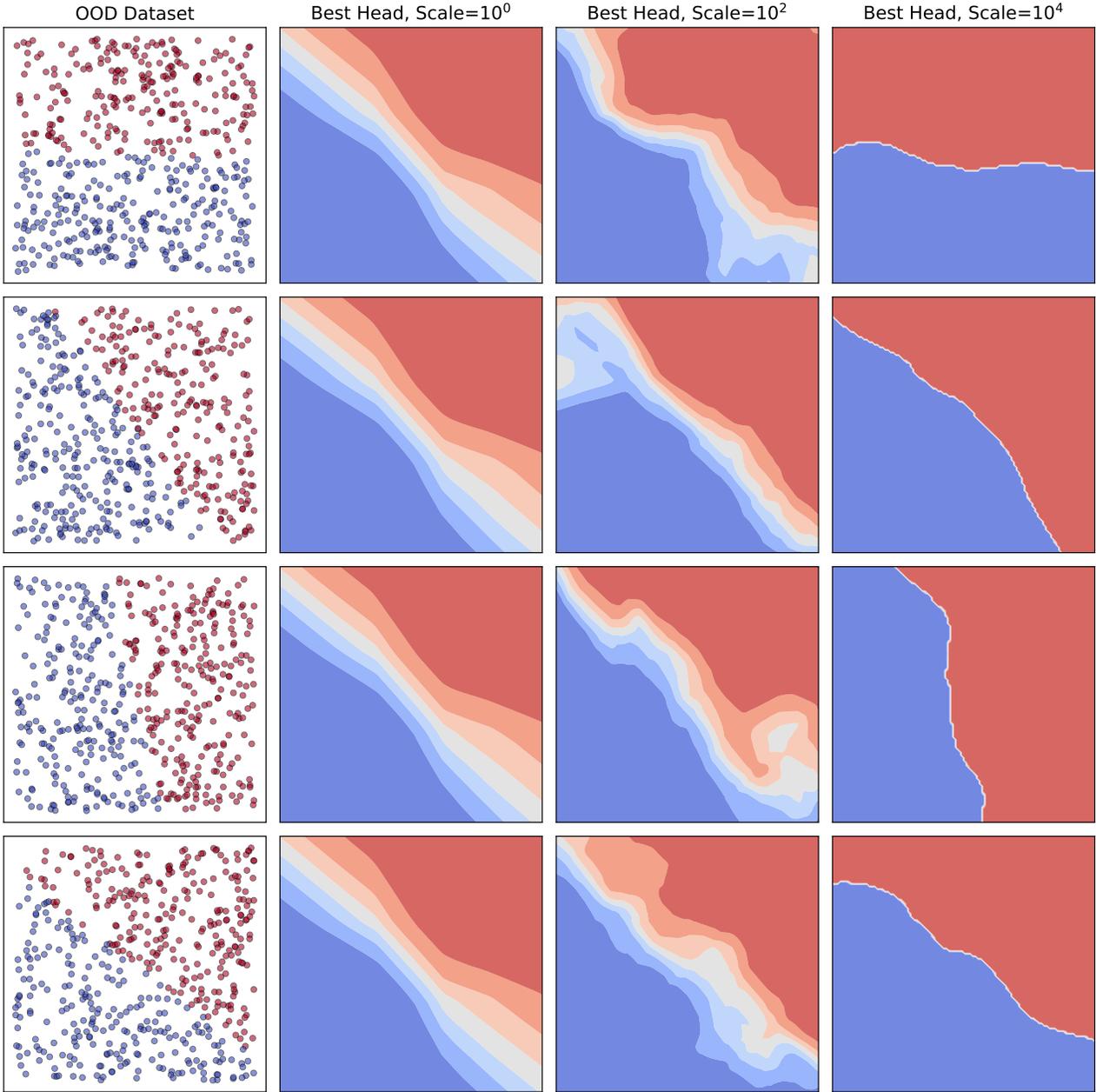


Figure 8: Additional visualizations for the toy conflicting classification example. Increasing the scale hyperparameter results produces heads with sharper decision boundaries.

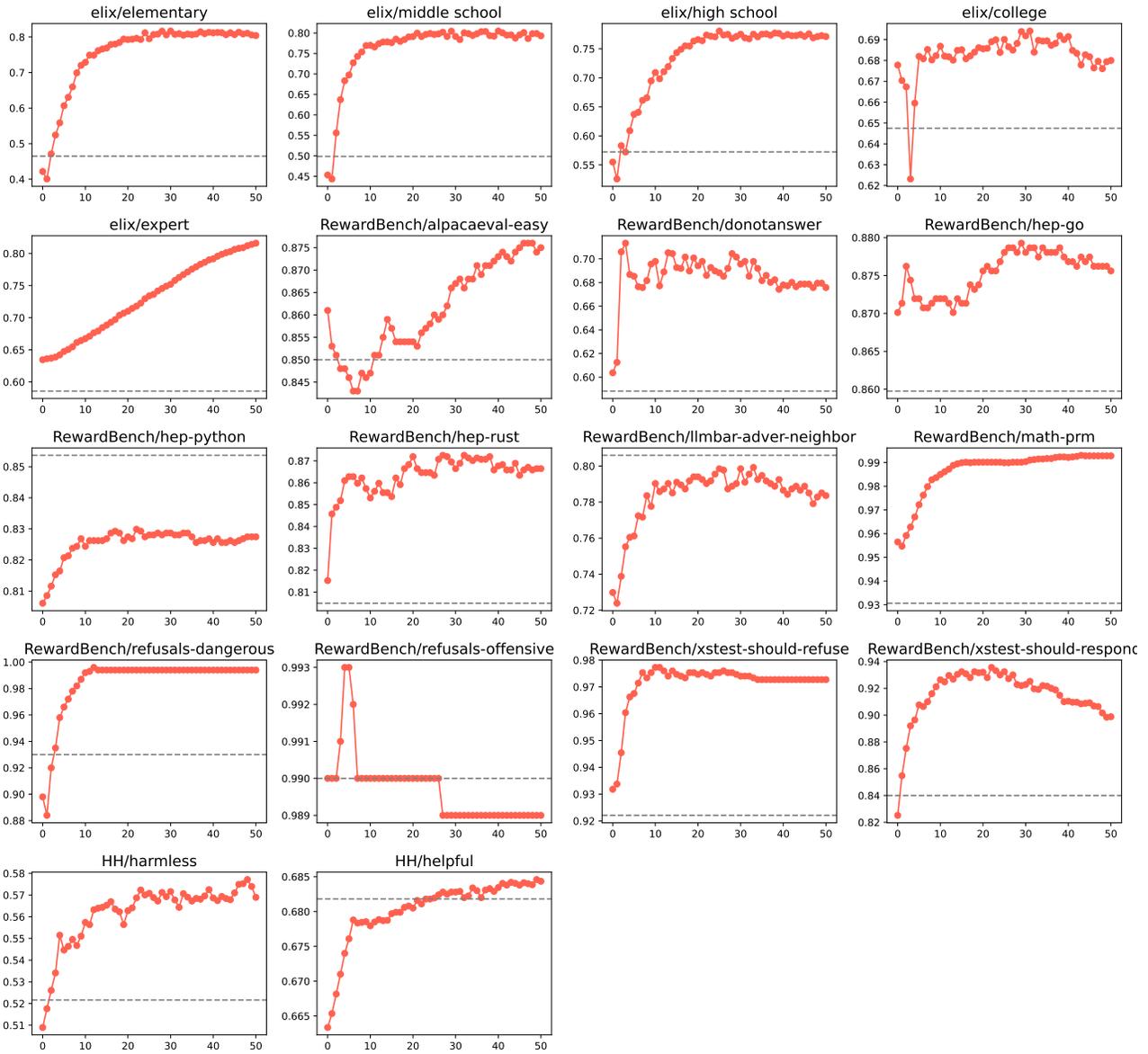


Figure 9: Detailed results for the personalizing preference reward models experiment in Figure 2. Target dataset accuracy (y-axis) after observing different numbers of adaptation samples (x-axis). The dashed line represents the performance of the pretrained reward model.