

---

# Robust Multi-task Modeling for Bayesian Optimization via In-Context Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Bayesian optimization is a sample-efficient optimization technique for black-box  
2 optimization, and leveraging historical information from related tasks can greatly  
3 improve its performance. Gaussian processes (GPs) are commonly used to model  
4 this multi-task data; however, they trade off complexity with expressivity. Jointly  
5 modeling all tasks can be computationally infeasible for GPs, while scalable  
6 approaches may fail to effectively utilize inter-task relationships. Moreover, these  
7 methods are often prone to negative transfer, where the inclusion of unrelated tasks  
8 degrades predictive performance. In this paper, we present Multi-Task Prior-Data  
9 Fitted Networks (MTPFNs), a multi-task model that efficiently and jointly models  
10 all tasks and data points. We show that MTPFNs serve as a compelling surrogate  
11 model that is robust to negative transfer, and their flexibility enables more efficient  
12 exploration. We demonstrate the effectiveness of our approach across a variety of  
13 synthetic and real-world benchmarks including hyperparameter optimization.

## 14 1 Introduction

15 Optimizing expensive black-box functions is critical across scientific and engineering domains,  
16 and important applications include as hyperparameter tuning for machine learning models [29],  
17 chemical reaction optimization [28], and engineering design problems [21]. For these resource-  
18 intensive settings, Bayesian optimization (BO) is a sample-efficient method that aims to find the  
19 global optimum with minimal function evaluations by using probabilistic surrogate models to select  
20 future evaluations [12].

21 Transfer learning can be used to further improve sample efficiency by extracting information from  
22 related tasks. Existing transfer-learning BO approaches often use a surrogate to jointly model the data  
23 from different tasks. However, these models face limitations. The commonly-used Gaussian process  
24 (GP) surrogate models often trade off with data efficiency with robustness: the models which jointly  
25 fit all of the data [4, 30, 35, 26, 18] may make strong assumptions about how different tasks are  
26 correlated and experience *negative transfer*, where unrelated tasks hinder performance; meanwhile,  
27 other approaches which fit separate GPs per task and ensemble their predictions [13, 10, 33, 7, 31]  
28 are more robust to negative transfer but cannot jointly capture cross-task information. Furthermore,  
29 multi-task GP may become computationally infeasible as we scale the number of tasks and data  
30 points, while more flexible variants may fail to efficiently utilize the full dataset.

31 Neural-network based approaches are an attractive alternative to Gaussian Processes. Transformer  
32 neural processes (TNP) [24] and prior fitted networks (PFNs) [22] are transformers trained to  
33 approximate the posterior predictive distribution given a data generating function, and they are  
34 capable of efficiently approximating the posterior for complex and bespoke priors. While TNPs and  
35 PFNs have successfully applied to Bayesian optimization in the single-task setting [23, 25], there  
36 has been no prior work which explores the use of in-context transfer of related tasks to accelerate

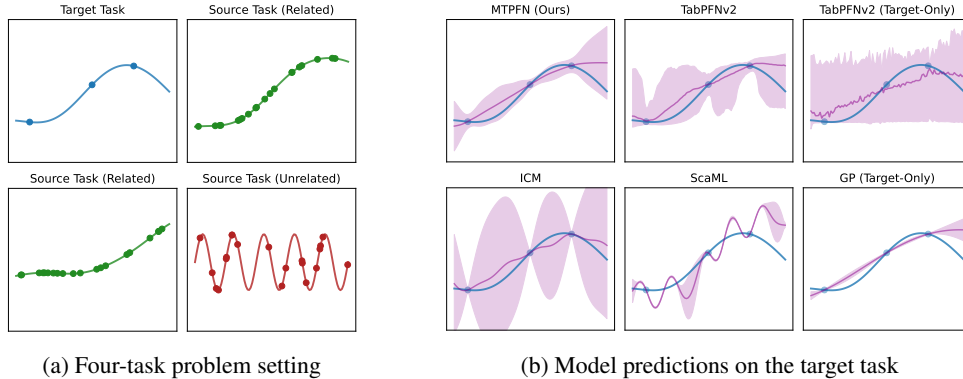


Figure 1: **MTPFNs effectively transfer information from related tasks while remaining robust to unrelated tasks.** Compared to joint models such as ICM, ScaML, and TabPFNv2 with a categorical task variable, our MTPFN demonstrates improved robustness to the unrelated source task (red). The MTPFN is also able to borrow strength from the related source tasks (green) and outperforms models which only consider the target task. We plot the mean and 95% confidence intervals for each model.

37 optimization. For more details, we provide a comprehensive background on Bayesian optimization,  
 38 multi-task GPs, neural-network approaches, and other related work in Appendix A.

39 We propose Multi-Task Prior-data Fitted Networks (MTPFNs), which use a novel data generating  
 40 process that samples data from both related and unrelated tasks. MTPFNs achieve better generalization  
 41 on classical transfer tasks compared to previous GP-based approaches, while being more robust to  
 42 negative transfer from irrelevant auxiliary sources. We demonstrate the efficacy of MTPFNs as a  
 43 surrogate model in Bayesian optimization on synthetic problems and real-world AutoML benchmarks.

## 44 2 Method

45 In this section, we present the Multi-Task Prior-Data Fitted Network (MTPFN), a scalable model that  
 46 transfers relevant knowledge from auxiliary tasks to model a target task via in-context learning.

### 47 2.1 Data Generation Process

48 PFNs are trained to approximate the posterior of a data generation process (DGP), and the design of  
 49 this prior significantly influences predictive performance. While various DGPs have been proposed  
 50 in previous works [e.g. 1, 15, 23], the multi-task setting poses unique challenges: there are complex  
 51 relationships between tasks where information may transfer through shared latent structures, and real-  
 52 world scenarios frequently contain unrelated tasks which should not negatively impact performance.  
 53 To address these challenges, we propose a multi-task DGP that learns complex relationships between  
 54 relevant tasks while mitigating corruption from irrelevant tasks.

55 Our approach combines two key insights: when tasks are truly related, strong transfer can be achieved  
 56 by sharing statistical properties like lengthscales across tasks; however, since we cannot know a  
 57 priori which auxiliary tasks will be helpful, the model must explicitly account for task irrelevance.  
 58 Our DGP models task relationships through a shared covariance structure when tasks are related and  
 59 introduces a probability  $p \in [0, 1]$  that each auxiliary task is instead modeled independently when  
 60 it would be irrelevant or harmful. In Figure 1, existing GP models and PFN priors perform poorly  
 61 in this synthetic setting due to influence from the unrelated task in red. In contrast, our MTPFN is  
 62 robust to the irrelevant task and accurately mirrors the true behavior, demonstrating the benefit of our  
 63 robust prior generation procedure. We present the full algorithm in Algorithm A.1, with additional  
 64 discussion and ablations for alternative DGPs in Appendix B.

### 65 2.2 Hierarchical Attention Mechanism

66 For multi-task regression problems, it is important to consider how the task itself should be encoded  
 67 within the model. Two natural approaches include representing the associated task as a *categorical*  
 68 *feature* or a continuous *task embedding*. However, these methods have distinct limitations: they

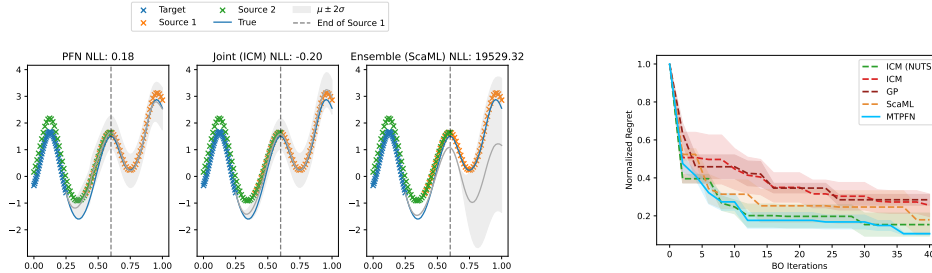


Figure 2: **MTPFNs jointly model the data from the target and all auxiliary tasks and perform fully Bayesian inference.** (Left): MTPFNs perform similarly to other joint models like ICM and outperform ensemble-based models like ScaML. (Right): MTPFNs have comparable performance to fully Bayesian methods like ICM with MCMC (NUTS) sampling.

69 require the maximum number of tasks to be specified during train-time, and they will not generalize to  
 70 larger number of tasks. Furthermore, the model’s representation of a task is fixed and is not influenced  
 71 by the data points associated with that task.

72 To address these limitations, we propose a novel scalable attention mechanism for PFNs that effectively  
 73 leverages the natural hierarchical structure of multi-task data, as shown in Figure A.3. Our  
 74 approach applies hierarchical attention [34] to the multi-task regression setting, and we interleave  
 75 intra-task and inter-task blocks. Our architecture introduces learnable “[Task]” tokens that summarize  
 76 task-specific properties. The intra-task transformer blocks learn relationships between data points  
 77 within each task while updating both point embeddings and the “[Task]” token summary, requiring  
 78  $O(D^2)$  compute per task. The inter-task encoders then learn relationships between tasks by attending  
 79 only to the summary “[Task]” embedding with  $O(T^2)$  complexity. This hierarchical design reduces  
 80 the overall attention complexity from the naive global setting of  $O(D^2T^2)$  to  $O(D^2T + T^2)$ , enabling  
 81 significantly longer contexts while still allowing for every data point to influence others.

82 Our hierarchical attention directly addresses many of the limitations of other task encoders. First, our  
 83 attention mechanism naturally handles inputs of varying lengths, allowing the model to generalize  
 84 to any number of tasks. Furthermore, our approach enables the model to dynamically learn data-  
 85 dependent task representations, so tasks with similar patterns can develop similar representations,  
 86 allowing the model to better capture the potentially complex relationships between tasks. We include  
 87 further details in Appendix C.

### 88 3 Advantages of MTPFNs

89 In this section, we provide explicit demonstrations of the strengths of MTPFNs for multi-task learning.  
 90 We compare against baselines: (1) ICM [14], a joint method which trains a multi-task GP on the  
 91 combined target and auxiliary data; (2) ScaML [31], an ensemble method that fits individual GPs to  
 92 each auxiliary task; and (3) a single-task GP which only uses the target task and ignores the auxiliary  
 93 tasks. See Appendix D for additional details about the empirical evaluation and further results.

94 **MTPFNs are robust to negative transfer** In Appendix D.2, we evaluate the performance of  
 95 multi-task models in settings with varying proportions of unrelated tasks. When only 1 of the 3  
 96 auxiliary tasks is unrelated, we find that all of the multi-task methods perform similarly. However,  
 97 as we increase the number of unrelated tasks to 2 out of 3, we find that the MTPFNs trained on data  
 98 with a higher proportion of corrupted tasks outperform the ICM, which is more sensitive to negative  
 99 transfer. When we increase the number of unrelated auxiliary tasks 3, we find that the MTPFN trained  
 100 with  $p = 0.2$  is comparable to the single-task GP, which is the underlying DGP for this problem.

101 **MTPFNs efficiently model inter-task relationships** We design a regression problem in Figure 2 to  
 102 highlight the importance of joint modeling. In this setting, all of the data points across all source tasks  
 103 are drawn from the same function, and this function is highly correlated with the target task. We see  
 104 that ensemble methods such as ScaML, which do not model the joint interactions between auxiliary  
 105 tasks, are unable to effectively use the relevant information. However, traditional joint modeling  
 106 methods like ICM are unable to scale to a large number of tasks and data points. In Appendix D.3,  
 107 we benchmark the runtimes of multi-task models as we increase the number of tasks and the number

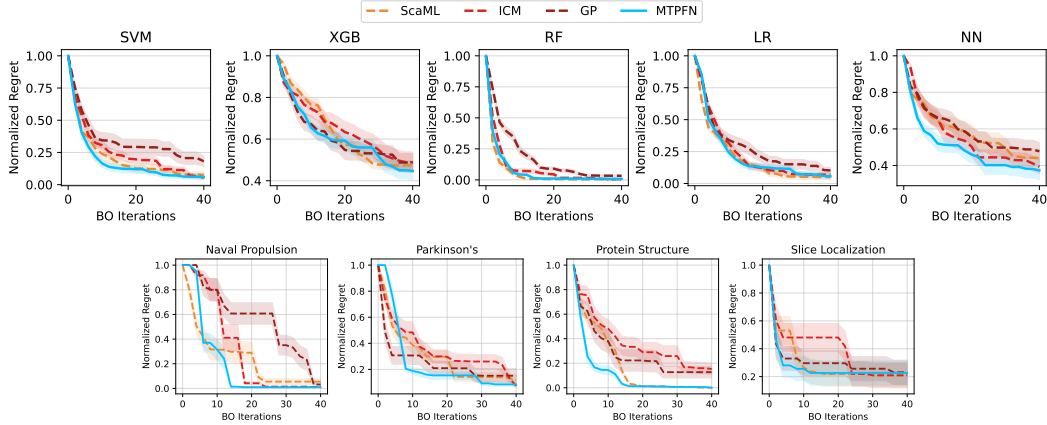


Figure 3: **MTPFNs are competitive across diverse real-world benchmarks.** Each plot shows the normalized regret for Bayesian optimization loop that was initialized with 3 auxiliary tasks, 20 observations from auxiliary task, and 5 observations from the target task.

108 of data points per task, and we find that this problem quickly becomes unmanageable for ICM. In  
 109 contrast, MTPFNs are able to scale to large amounts of data while jointly modeling all interactions.

110 **MTPFNs quickly perform fully Bayesian inference** Müller et al. [22] demonstrate that PFNs  
 111 perform fully Bayesian inference by implicitly learning the posterior predictive distribution that  
 112 marginalizes over all possible samples from the prior which are consistent with the observed data.  
 113 This ability extends to the multi-task setting: in Figure 2 (Right), we initialize the problem with 20  
 114 samples from three auxiliary tasks and 2 samples from the target task. MTPFNs are comparable to  
 115 the fully Bayesian inference using MCMC sampling through NUTS; furthermore, the MTPFN is able  
 116 to make predictions using one forward pass of the model in approximately 0.5 seconds on average,  
 117 while NUTS takes orders of magnitudes longer at 352 seconds per iteration. We showcase further  
 118 demonstrations of the importance of fully Bayesian inference in Appendix D.4.

119 **MTPFNs can leverage domain data** When making predictions with PFNs, domain data can be  
 120 incorporated through fine-tuning or in-context learning. Fine-tuning updates model parameters to  
 121 adapt to specific domains, enabling specialization but at high computational cost and risk of overfitting  
 122 that hurts generalization. In Appendix D.5, we demonstrate these drawbacks on HPOBench’s Logistic  
 123 Regression dataset [8], which contains 25 training tasks with 4 held-out evaluation tasks. We find  
 124 that fine-tuning improves target domain performance but severely hurts generalization on other  
 125 domains like SVM and GP tasks due to overfitting, whereas MTPFNs and in-context learning achieve  
 126 comparable target performance while maintaining strong generalization across all datasets.

127 **MTPFNs are effective across real-world optimization benchmarks** We compare the effectiveness  
 128 of the methods on a set of hyperparameter optimization problems [8] for machine learning model  
 129 described in Appendix D.6. Following Tighineanu et al. [31], we consider the hyperparameter  
 130 optimization for five types of models and we share the results of our benchmark in the top panel of  
 131 Figure 3. We also consider tabular FC-Net benchmarks: Slice Localization, Protein Structure, Naval  
 132 Propulsion, and Parkinson’s Telemonitoring. For each benchmark, we set the benchmark to be the  
 133 target task, and we use the three other tabular dataset as the auxiliary tasks. MTPFNs are competitive  
 134 across all of the settings. Specifically, we find that in instances where the meta-tasks contain helpful  
 135 information (ScaML and ICM outperform GP), the MTPFNs are also able to effectively utilize this  
 136 data. Furthermore, in cases like XGB where there is negative transfer for the ICM model, we find  
 137 that MTPFNs are more robust and perform similarly to the standard single-task GP.

## 138 4 Discussion

139 In this work, we present MTPFNs, a scalable and robust surrogate model for Bayesian optimiza-  
 140 tion. Our results highlight the effectiveness of leveraging domain data through in-context learning.  
 141 MTPFNs are able to successfully capture the complex relationships between the information sources  
 142 and thus can leverage auxiliary information without expensive model-fitting or fine-tuning procedures.

143 **References**

- 144 [1] Adriaensen, S., Rakotoarison, H., Müller, S., and Hutter, F. Efficient bayesian learning curve  
145 extrapolation using prior-data fitted networks, 2023. URL [https://arxiv.org/abs/2310.](https://arxiv.org/abs/2310.20447)  
146 20447.
- 147 [2] Balandat, M., Karrer, B., Jiang, D., Daulton, S., Letham, B., Wilson, A. G., and Bakshy, E.  
148 Botorch: A framework for efficient monte-carlo bayesian optimization. *Advances in neural*  
149 *information processing systems*, 33:21524–21538, 2020.
- 150 [3] Beltagy, I., Peters, M. E., and Cohan, A. Longformer: The long-document transformer. *arXiv*  
151 *preprint arXiv:2004.05150*, 2020.
- 152 [4] Bonilla, E. V., Chai, K., and Williams, C. Multi-task gaussian process prediction. In Platt, J.,  
153 Koller, D., Singer, Y., and Roweis, S. (eds.), *Advances in Neural Information Processing Sys-*  
154 *tems*, volume 20. Curran Associates, Inc., 2007. URL [https://proceedings.neurips.cc/](https://proceedings.neurips.cc/paper_files/paper/2007/file/66368270ffd51418ec58bd793f2d9b1b-Paper.pdf)  
155 [paper\\_files/paper/2007/file/66368270ffd51418ec58bd793f2d9b1b-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2007/file/66368270ffd51418ec58bd793f2d9b1b-Paper.pdf).
- 156 [5] Chalkidis, I., Dai, X., Fergadiotis, M., Malakasiotis, P., and Elliott, D. An exploration of  
157 hierarchical attention transformers for efficient long document classification. *arXiv preprint*  
158 *arXiv:2210.05529*, 2022.
- 159 [6] Chen, Y., Song, X., Lee, C., Wang, Z., Zhang, R., Dohan, D., Kawakami, K., Kochanski, G.,  
160 Doucet, A., Ranzato, M., et al. Towards learning universal hyperparameter optimizers with  
161 transformers. *Advances in Neural Information Processing Systems*, 35:32053–32068, 2022.
- 162 [7] Dai, Z., Chen, Y., Yu, H., Low, B. K. H., and Jaillet, P. On provably robust meta-bayesian  
163 optimization. In *Uncertainty in Artificial Intelligence*, pp. 475–485. PMLR, 2022.
- 164 [8] Eggenesperger, K., Müller, P., Mallik, N., Feurer, M., Sass, R., Klein, A., Awad, N., Lindauer, M.,  
165 and Hutter, F. HPOBench: A collection of reproducible multi-fidelity benchmark problems for  
166 HPO. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Bench-*  
167 *marks Track (Round 2)*, 2021. URL <https://openreview.net/forum?id=1k4rJYEWda->.
- 168 [9] Fan, Z., Han, X., and Wang, Z. Hyperbo+: Pre-training a universal prior for bayesian optimiza-  
169 tion with hierarchical gaussian processes. *arXiv preprint arXiv:2212.10538*, 2022.
- 170 [10] Feurer, M., Letham, B., and Bakshy, E. Scalable meta-learning for bayesian optimization using  
171 ranking-weighted gaussian process ensembles. In *AutoML Workshop at ICML*, volume 7, pp. 5,  
172 2018.
- 173 [11] Gardner, J. R., Pleiss, G., Bindel, D., Weinberger, K. Q., and Wilson, A. G. Gpytorch:  
174 Blackbox matrix-matrix gaussian process inference with gpu acceleration. In *Advances in*  
175 *Neural Information Processing Systems*, 2018.
- 176 [12] Garnett, R. *Bayesian Optimization*. Cambridge University Press, 2023.
- 177 [13] Golovin, D., Solnik, B., Moitra, S., Kochanski, G., Karro, J., and Sculley, D. Google vizier: A  
178 service for black-box optimization. In *Proceedings of the 23rd ACM SIGKDD international*  
179 *conference on knowledge discovery and data mining*, pp. 1487–1495, 2017.
- 180 [14] Goovaerts, P. *Geostatistics for natural resources evaluation*, volume 483. Oxford University  
181 Press, 1997.
- 182 [15] Hollmann, N., Müller, S., Purucker, L., Krishnakumar, A., Körfer, M., Hoo, S. B., Schirrmester,  
183 R. T., and Hutter, F. Accurate predictions on small data with a tabular foundation model. *Nature*,  
184 637(8045):319–326, 2025.
- 185 [16] Hvarfner, C., Hellsten, E. O., and Nardi, L. Vanilla Bayesian optimization performs great in  
186 high dimensions. In Salakhutdinov, R., Kolter, Z., Heller, K., Weller, A., Oliver, N., Scarlett,  
187 J., and Berkenkamp, F. (eds.), *Proceedings of the 41st International Conference on Machine*  
188 *Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 20793–20817. PMLR,  
189 21–27 Jul 2024. URL <https://proceedings.mlr.press/v235/hvarfner24a.html>.

- 190 [17] Jones, D. R., Schonlau, M., and Welch, W. J. Efficient global optimization of expensive  
191 black-box functions. *Journal of Global Optimization*, 13:455–492, 1998.
- 192 [18] Joy, T. T., Rana, S., Gupta, S., and Venkatesh, S. A flexible transfer learning framework for  
193 bayesian optimization with convergence guarantee. *Expert Systems with Applications*, 115:  
194 656–672, 2019.
- 195 [19] Katharopoulos, A., Vyas, A., Pappas, N., and Fleuret, F. Transformers are rnns: Fast autoregres-  
196 sive transformers with linear attention. In *International conference on machine learning*, pp.  
197 5156–5165. PMLR, 2020.
- 198 [20] Kitaev, N., Kaiser, Ł., and Levskaya, A. Reformer: The efficient transformer. *arXiv preprint*  
199 *arXiv:2001.04451*, 2020.
- 200 [21] Liao, X., Li, Q., Yang, X., Zhang, W., and Li, W. Multiobjective optimization for crash  
201 safety design of vehicles using stepwise regression model. *Structural and Multidisciplinary*  
202 *Optimization*, 35:561–569, 06 2008. doi: 10.1007/s00158-007-0163-x.
- 203 [22] Müller, S., Hollmann, N., Arango, S. P., Grabocka, J., and Hutter, F. Transformers can do  
204 bayesian inference. *arXiv preprint arXiv:2112.10510*, 2021.
- 205 [23] Müller, S., Feurer, M., Hollmann, N., and Hutter, F. Pfns4bo: In-context learning for bayesian  
206 optimization. In *International Conference on Machine Learning*, pp. 25444–25470. PMLR,  
207 2023.
- 208 [24] Nguyen, T. and Grover, A. Transformer neural processes: Uncertainty-aware meta learning  
209 via sequence modeling. In *International Conference on Machine Learning*, pp. 16569–16594.  
210 PMLR, 2022.
- 211 [25] Nguyen, T., Zhang, Q., Yang, B., Lee, C., Bornschein, J., Miao, Y., Perel, S., Chen, Y., and  
212 Song, X. Predicting from strings: Language model embeddings for bayesian optimization, 2024.  
213 URL <https://arxiv.org/abs/2410.10190>.
- 214 [26] Poloczek, M., Wang, J., and Frazier, P. Multi-information source optimization. *Advances in*  
215 *neural information processing systems*, 30, 2017.
- 216 [27] Rasmussen, C. E. *Gaussian Processes in Machine Learning*, pp. 63–71. Springer Berlin  
217 Heidelberg, Berlin, Heidelberg, 2004.
- 218 [28] Shields, B. J., Stevens, J., Li, J., Parasram, M., Damani, F., Alvarado, J. I. M., Janey, J. M.,  
219 Adams, R. P., and Doyle, A. G. Bayesian reaction optimization as a tool for chemical synthesis.  
220 *Nature*, 590(7844):89–96, 2021.
- 221 [29] Snoek, J., Larochelle, H., and Adams, R. P. Practical bayesian optimization of machine learning  
222 algorithms. In *Advances in neural information processing systems*, pp. 2951–2959, 2012.
- 223 [30] Swersky, K., Snoek, J., and Adams, R. P. Multi-task bayesian optimization. *Advances in neural*  
224 *information processing systems*, 26, 2013.
- 225 [31] Tighineanu, P., Grossberger, L., Baireuther, P., Skubch, K., Falkner, S., Vinogradska, J., and  
226 Berkenkamp, F. Scalable meta-learning with gaussian processes. In *International Conference*  
227 *on Artificial Intelligence and Statistics*, pp. 1981–1989. PMLR, 2024.
- 228 [32] Wang, Z., Dahl, G. E., Swersky, K., Lee, C., Nado, Z., Gilmer, J., Snoek, J., and Ghahramani,  
229 Z. Pre-trained gaussian processes for bayesian optimization. *Journal of Machine Learning*  
230 *Research*, 25(212):1–83, 2024.
- 231 [33] Wistuba, M., Schilling, N., and Schmidt-Thieme, L. Scalable gaussian process-based transfer  
232 surrogates for hyperparameter optimization. *Machine Learning*, 107(1):43–78, 2018.
- 233 [34] Wu, C., Wu, F., Qi, T., and Huang, Y. Hi-transformer: Hierarchical interactive transformer for  
234 efficient and effective long document modeling. *arXiv preprint arXiv:2106.01040*, 2021.
- 235 [35] Yogatama, D. and Mann, G. Efficient transfer learning method for automatic hyperparameter  
236 tuning. In *Artificial intelligence and statistics*, pp. 1077–1085. PMLR, 2014.

- 237 [36] Zaheer, M., Guruganesh, G., Dubey, K. A., Ainslie, J., Alberti, C., Ontanon, S., Pham, P.,  
238 Ravula, A., Wang, Q., Yang, L., et al. Big bird: Transformers for longer sequences. *Advances*  
239 *in neural information processing systems*, 33:17283–17297, 2020.
- 240 [37] Zhuang, B., Liu, J., Pan, Z., He, H., Weng, Y., and Shen, C. A survey on efficient training of  
241 transformers. *arXiv preprint arXiv:2302.01107*, 2023.

## 242 A Background and Related Work

### 243 A.1 Bayesian Optimization

244 Bayesian optimization (BO) [12] is a sample-efficient black-box optimization method. Given a  
245 function  $f$  and a compact search space  $\mathcal{X} \subset \mathbb{R}^D$ , BO aims to find the global maximum<sup>1</sup>  $x^* =$   
246  $\arg \max_{x \in \mathcal{X}} f(x)$  by iteratively selecting points to evaluate, conditional on observations made  
247 previously. Typically, the observations are corrupted by noise  $y = f(x) + \epsilon(x)$  where  $\epsilon(x)$  is a noise  
248 process. At each step  $t$ , a probabilistic model (often a Gaussian process [27]) parameterized by  $\theta$   
249 is fit to the data collected so far  $\mathcal{D}_t = \{(x_i, y_i)\}_{i=1}^t$ . An acquisition function  $a(x; \theta)$  (e.g. expected  
250 improvement [17]), which balances exploration and exploitation, then uses the model’s posterior  
251 distribution  $p(f|\mathcal{D}_t, \theta)$  to determine the next point to evaluate  $x_{t+1} = \arg \max_x a(x; \theta)$ . Finally,  
252 we evaluate the function for the selected point  $x_{t+1}$ , add the new observation to the set of function  
253 evaluations  $\mathcal{D}_{t+1} = \mathcal{D}_t \cup \{(x_{t+1}, y_{t+1})\}$ , and proceed to the next iteration. The surrogate model is  
254 the key determinant of how successful BO will be, and therefore improvements in the model typically  
255 lead to improvements in performance [16].

256 In the multi-task setting, we have access to auxiliary data from functions which may be similar to the  
257 black-box function we wish to optimize. Formally, we aim to find the global optimum of a target task  
258 (denoted by task ID 0)  $x^* = \arg \max_{x \in \mathcal{X}} f_0(x)$  while having access to evaluations from auxiliary  
259 tasks  $\{f_k\}_{k=1}^K$ , which we assumed are defined over the same search space. Compared to single-task  
260 BO, each observation consists of an additional task index  $k$ , and the dataset is  $\mathcal{D}_t = \{(x_i, y_i, k_i)\}_{i=1}^{n_t}$ .  
261 The acquisition function  $a(x; \theta)$  is then used to select the next point for the target task 0. Successful  
262 multi-task BO requires a surrogate model that is able to correctly infer the relationship between the  
263 tasks and use this information effectively.

### 264 A.2 Multi-Task Surrogate Models

265 One approach to model multi-task data for BO is to jointly model the full collection of target and  
266 auxiliary data using a single GP with a multi-task kernel [4, 30, 35, 26, 18]. We will refer to these  
267 models as multi-task GPs (MTGPs). While these joint modeling approaches are effective in the  
268 low-data regime, they become computationally infeasible as we scale the number of tasks and data  
269 points due to their cubic complexity in the total number of data evaluations.

The intrinsic coregionalization model ICM, [14] is a common choice of MTGP due to its simplicity  
and was proposed in Swersky et al. [30] for multi-task BO. The ICM models the functions with a  
kernel that decomposes into two components

$$k((x, t), (x', t')) = k_{\text{inputs}}(x, x') \cdot k_{\text{tasks}}(t, t'),$$

270 where  $k_{\text{inputs}}$  is a kernel (often RBF or Matérn) that represents the covariance between inputs and  $k_{\text{tasks}}$   
271 captures the covariance between tasks. Because this model assumes a single shared latent function  
272 across all tasks, it can efficiently transfer knowledge between similar tasks; however, the ICM model  
273 may perform poorly when this assumption does not hold and the tasks have distinct characteristics,  
274 e.g., when they should be modelled with different lengthscales.

275 There are other multi-task models which require weaker assumptions. For example, the linear model  
276 of coregionalization LMC, [14] is a generalization of the ICM, which uses multiple latent functions that  
277 are linearly combined for each task, instead of assuming a single latent function for all tasks. However,  
278 although this model is more flexible than the ICM, it has increased computational complexity and  
279 may lead to overfitting.

280 Alternative approaches which focus on scaling to large multi-task datasets have also been proposed.  
281 Many methods fit separate GPs to each auxiliary task and ensemble their predictions to inform the  
282 target task [13, 10, 33, 7]. Although these approaches are scalable, they are not able to jointly capture  
283 information across the related tasks and instead rely on heuristics to determine the relevance of each  
284 GP. Other methods use the auxiliary data to learn a better prior over GP hyperparameters for the  
285 target task [32, 9]; however, these methods are unable to utilize the specific information within each  
286 task. Tighineanu et al. [31] propose a scalable joint modeling approach between the target task and  
287 auxiliary tasks; however, this does not model the correlations between the auxiliary tasks, thereby not

---

<sup>1</sup>Or equivalently the global minimum. Without loss of generality we consider maximization.

288 taking advantage of the entire dataset. In contrast, we propose a scalable method that jointly models  
289 the full interaction between all data points and tasks.

### 290 A.3 Bayesian Optimization with Transformers

291 For single-task Bayesian optimization, there has been a growing interest in using neural-network  
292 based approaches. OptFormer [6] leverages a transformer directly trained on data collected from  
293 BO loops. In this setup the model does not only model  $y$ , like in classical BO, but directly predicts  
294 proposals for the next  $x$ . This method requires access to large amounts of domain data during training,  
295 which may not be possible in many settings.

296 Transformer neural processes (TNP) [24] and prior fitted networks (PFNs) [22] are transformers  
297 trained to approximate the posterior predictive distribution given a data generating function (prior).  
298 These models approximate posterior predictive distribution for a prior specified over a hypothesis  
299 space  $\mathcal{H}$  where each hypothesis  $h \in \mathcal{H}$  defines a relationship between inputs  $x$  and outputs  $y$ . This  
300 work builds off of the PFN framework due to PFN’s strong empirical performance in Bayesian  
301 optimization tasks [23] and prediction for tabular data [15].

302 A PFN, denoted by  $f_\theta$ , takes as input a dataset  $\mathcal{D}$  and test point  $x_{\text{test}}$  and outputs a distribution over  
303 the target variable  $p(y_{\text{test}}|x_{\text{test}}, \mathcal{D})$ . To train  $f_\theta$  to approximate the posterior predictive distribution, we  
304 repeatedly sample datasets by first sampling a hypothesis  $h \sim p(h)$  which defines a datasets’ input-  
305 output relationship, and then sampling a dataset  $\mathcal{D} \sim p(\mathcal{D}|h)$ . The PFN parameters  $\theta$  are optimized  
306 by minimizing the negative log-likelihood on held-out test examples across datasets, expressed as  
307  $\mathcal{L}_{\text{NLL}} = \mathbb{E}_{\mathcal{D} \sim p(\mathcal{D}|h)}[-\log f_\theta(y_{\text{test}}|x_{\text{test}}, \mathcal{D}_{\text{train}})]$ , where  $\mathcal{D}$  is split into  $\mathcal{D}_{\text{train}} \cup \{(x_{\text{test}}, y_{\text{test}})\}$ .

308 While TNPs and PFNs have successfully applied to Bayesian optimization in the single-task setting  
309 [23, 25], there has been no prior work which explores the use of in-context transfer of related tasks to  
310 accelerate optimization. While the PFN’s pre-training can already be interpreted as meta-learning, we  
311 take it a step further by training them to do Bayesian inference over several related tasks in-context;  
312 in this sense, the multi-task PFN acts as a *meta-meta-learning* model.

### 313 A.4 Long Contexts

314 Because of the increased number of tasks and data-points required for multi-task Bayesian optimiza-  
315 tion, the underlying architecture needs to support significantly longer context windows compared to  
316 the single-task setting. Various approaches have been proposed to extend the attention mechanisms in  
317 transformers to longer contexts, such as sparse attention [3, 36], hierarchical attention [34, 5], and  
318 others [19, 20]. See Zhuang et al. [37] for a survey of efficient methods.

## 319 B Data Generation Processes for PFNs

320 MTPFNs are flexible and have the capacity to incorporate various data generation processes during  
321 training. In this section, we explore the impacts of various data generation processes for MTPFNs,  
322 each designed to capture distinct inductive biases which improve model performance across different  
323 types of tasks.

### 324 B.1 Robust Isotropic Full-Rank ICM

---

**Algorithm A.1** Data Generation Using a Robust Isotropic Full-Rank ICM

---

**Require:** Sequence length  $n$ , number of tasks  $T$ , unrelated task probability  $p$

- ▷ **Input Sampling and Task Assignment**
  - 1: Sample inputs  $\{x_i\}_{i=1}^n \sim \text{Uniform}([0, 1]^d)$
  - 2: Sample task proportions  $\pi \sim \text{Dirichlet}(\alpha)$
  - 3: **for**  $i = 1$  to  $n$  **do**
  - 4:     Sample task ID  $t_i \sim \text{Categorical}(\pi)$
  - 5: **end for**
  - ▷ **Isotropic ICM Covariance Structure**
  - 6: Sample task covariance matrix  $K_T \sim \text{LKJ}(\eta = 1)$
  - 7: Sample input lengthscale  $\ell \sim \text{Gamma}(3, 6)$
  - 8: Define input covariance  $K_X$  on  $\{x_i\}_{i=1}^n$  as RBF kernel with lengthscale  $\ell$
  - 9: Compute full ICM kernel  $K = K_T \otimes K_X$
  - 10: Sample  $y \sim \mathcal{N}(0, K)$
  - ▷ **Sampling Unrelated Tasks**
  - 11: **for** each source task  $j$  **do**
  - 12:     With probability  $p$ :
  - 13:         Sample new lengthscale  $\ell_j \sim \text{Gamma}(3, 6)$
  - 14:         Define RBF kernel  $K_X^{(j)}$  on  $\{x_i : t_i = j\}$  using  $\ell_j$
  - 15:         Resample  $y^{(j)} \sim \mathcal{N}(0, K_X^{(j)})$
  - 16: **end for**
- 

325 In the main text, we train PFNs with the data generation process described in Algorithm A.1: we  
326 sample datapoints across tasks from a full-rank isotropic ICM model. This approach assumes that  
327 all the input dimensions share identical lengthscales; this assumption imposes a strong prior on the  
328 relationship between tasks, and enables effective information transfer across related tasks when the  
329 assumption is met. This data generation process enables us to learn information across related tasks,  
330 since the full-rank isotropic model makes strong assumptions.

331 To improve our model’s robustness to negative transfer, we also incorporate an additional hyperpa-  
332 rameter  $p \in [0, 1]$  which dictates the relatedness of the tasks during training. Specifically,  $p$  is the  
333 probability that any given source task is drawn independently from the target task, and thus may have  
334 completely different behaviors and lengthscales. Our data generation procedure enables the model to  
335 see a diverse group of datasets which consist of a mix of related and unrelated source tasks.

336 This  $p$  hyperparameter plays a crucial part in the robustness of the model against negative transfer:  
337 because the model is able to see many examples of unrelated tasks during training, it becomes more  
338 robust to seeing unrelated tasks during inference time and is less likely to be negatively impacted  
339 from irrelevant information.

340 We first introduce inter-task relationships by sampling from an ICM MTGP, where the ICM’s  
341 assumption of a shared lengthscale across tasks enables strong transfer when the tasks are related.  
342 Specifically, we sample an inter-task covariance matrix from an LKJ prior with a concentration of 1.0,  
343 which provides us with a diverse set of relationships between tasks, and we sample the shared RBF  
344 kernel lengthscale from a Gamma (3, 6) prior following the default lengthscale prior in BoTorch v1.11  
345 [2]. To prevent negative transfer, our DGP explicitly encodes the belief that each source task may be  
346 irrelevant to the target task by introducing a probability  $p \in [0, 1]$  that the task is instead modeled  
347 independently using a separate RBF GP with its own lengthscale. In the following sections, we  
348 present results under a simple and transparent DGP, but different priors and more sophisticated DGPs  
349 can easily be accommodated within the PFN framework. See Appendix B for additional discussions.

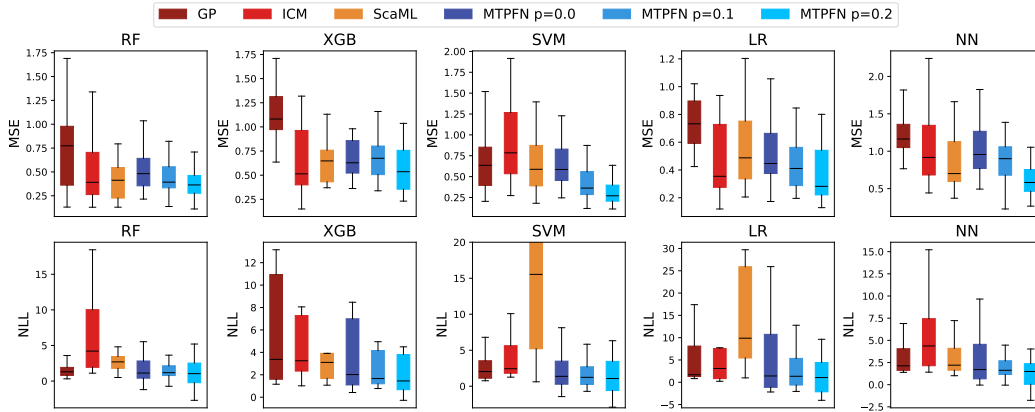


Figure A.1: As we increase  $p$  (the probability that each source task is unrelated to the target task during data generation), the model becomes more robust to negative transfer and achieves better performance on real-world benchmarks. We visualize the model’s predictive performance on the HPOBench dataset, where we sample 5 data points from the target task and 20 data points each from three source tasks. We plot the average MSE and NLL on holdout data from the target task across 25 trials.

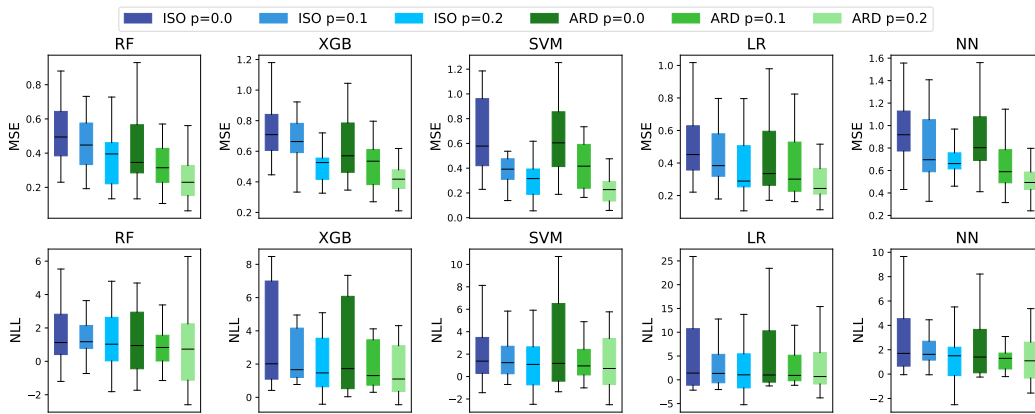


Figure A.2: MTPFNs trained with the ARD data generation process tend to outperform MTPFNs trained with the isotropic process (ISO) and achieve lower MSE and NLLs on HPOBench problems.

350 In Figure A.1, we study the impact of  $p$  on the model’s ability to accurately predict the empirical data  
 351 from HPOBench. Specifically, for each model type (SVM, LR, XGB, NN, and RF), we randomly  
 352 sample one task to be the target task, and we sample 3 auxiliary tasks from the metadata. The target  
 353 task is randomly initialized with 5 samples, and we also sample 20 points for each of the auxiliary  
 354 tasks. We measure the mean squared error (MSE) and the negative log-likelihood (NLL) of each  
 355 surrogate model on heldout examples from the target task, and we repeat this procedure 25 times and  
 356 plot the average MSE and NLL for each trial.

357 We find that increasing  $p$ , which increases the diversity of the data that the model sees during training,  
 358 leads to improved model performance on real-world benchmarks. We see that the MTPFN trained  
 359 with  $p = 0.2$  consistently outperforms other MTPFNs trained with lower values of  $p$ , and this MTPFN  
 360 also outperforms baselines such as the standard ICM model, which assumes that all tasks share the  
 361 same lengthscale.

## 362 B.2 Full-Rank ICM with Automatic Relevance Determination

363 We can also relax the assumption that all of the input dimensions share the same lengthscale, and  
 364 instead sample datapoints from an ICM model with Automatic Relevance Determination (ARD),  
 365 where we assume that each input dimensions has an independent lengthscale. This enables the PFNs

---

**Algorithm A.2** Data Generation Using a Robust ARD Full-Rank ICM

---

**Require:** Sequence length  $n$ , number of tasks  $T$ , unrelated task probability  $p$

- ▷ **Input Sampling and Task Assignment**
- 1: Sample inputs  $\{x_i\}_{i=1}^n \sim \text{Uniform}([0, 1]^d)$
- 2: Sample task proportions  $\pi \sim \text{Dirichlet}(\alpha)$
- 3: **for**  $i = 1$  to  $n$  **do**
- 4:     Sample task ID  $t_i \sim \text{Categorical}(\pi)$
- 5: **end for**
- ▷ **ARD ICM Covariance Structure**
- 6: Sample task covariance matrix  $K_T \sim \text{LKJ}(\eta = 1)$
- 7: **Sample independent input lengthscales:**  $\ell = (\ell_1, \dots, \ell_d) \sim \text{Gamma}(3, 6)^d$
- 8: **Define input covariance**  $K_X$  on  $\{x_i\}_{i=1}^n$  as an RBF kernel with ARD lengthscales  $\ell$
- 9: Compute full ICM kernel  $K = K_T \otimes K_X$
- 10: Sample  $y \sim \mathcal{N}(0, K)$
- ▷ **Sampling Unrelated Tasks**
- 11: **for** each source task  $j$  **do**
- 12:     With probability  $p$ :
- 13:         Sample new lengthscale  $\ell_j \sim \text{Gamma}(3, 6)$
- 14:         Define RBF kernel  $K_X^{(j)}$  on  $\{x_i : t_i = j\}$  using  $\ell_j$
- 15:         Resample  $y^{(j)} \sim \mathcal{N}(0, K_X^{(j)})$
- 16: **end for**

---

366 to have more flexibility and fit more complex problems; however, this weaker assumption may reduce  
367 the model’s ability to effective transfer information compared to the isotropic settings. We describe  
368 this data generation process in Algorithm A.2 and highlight the differences from the isotropic data  
369 generation in green.

370 In Figure A.2, we compare the performance of the MTPFN trained with isotropic lengthscales (ISO)  
371 to the performance of the MTPFNs trained with the ARD lengthscales. This experiment follows  
372 an identical setup to Figure A.1, where we sample 5 points from a target task and 20 points each  
373 from 3 source tasks, and evaluate the MTPFNs on held-out data from the target task. We plot the  
374 experiments across 25 trials.

375 We find that the improved flexibility of the ARD lengthscale generally enables the model to have  
376 better performance on the testing data, with the ARD outperforming ISO across many datasets.  
377 However, in some settings such as SVM, we find that the model performance of the isotropic ICM  
378 and the ARD ICM are comparable. This similar performance may be because the assumption of the  
379 shared lengthscale across input dimension is satisfied in this setting, so the additional flexibility of  
380 the ARD is unnecessary.

## 381 C Hierarchical Attention Details

382 For multi-task regression problems, the encoding of the task can influence how the model integrates  
383 information from the various sources and impact its ability to learn helpful relationships and differen-  
384 tiate between tasks with irrelevant characteristics. In this section, we describe two straight-forward  
385 methods to represent the task information, and we discuss their implementations and also describe  
386 their limitations for our problem setting.

387 **Categorical Feature** For a data point  $(\mathbf{x}_i, y_i)$ , its associated task can be represented as a categorical  
388 feature  $t_i \in \{1, 2, \dots, T\}$ , where  $T$  is the number of distinct tasks. We can represent  $t_i$  using a  
389 one-hot encoding  $\mathbf{1}_{t_i} \in \{0, 1\}^T$ , and the final input is formed by concatenating the original feature  
390 vector with the one-hot encoding  $\mathbf{x}'_i = [\mathbf{x}_i; \mathbf{1}_{t_i}]$ . Although simple, this approach does not provide the  
391 model with information related to the task itself. Furthermore, the maximum number of tasks must  
392 be specified at train-time, and the model is also unable to generalize to a larger number of tasks at  
393 test-time since the categorical feature has a fixed number of dimensions.

394 **Task Embedding** Rather than directly using the one-hot encoding of the task, we can use a task  
395 encoder to map the task  $t_i$  to a continuous embedding vector  $\mathbf{e}_{t_i} \in \mathbb{R}^d$ , where  $d$  is the embedding  
396 dimension of the model. This embedding is jointly learned with the model parameters, allowing  
397 the task representation to adapt to task-specific characteristics. The original input  $(\mathbf{x}_i, y_i)$  is first  
398 transformed into a feature  $\mathbf{z}_i = \phi(\mathbf{x}_i, y_i)$ , and then this feature is combined with the task embedding  
399  $\mathbf{z}'_i = \mathbf{z}_i + \mathbf{e}_{t_i}$ . This approach integrates the task information directly into the feature space; however,  
400 the representation for the task is still learned independently of the information within each task, and  
401 the model remains unable to generalize to more tasks at test-time.

### 402 C.1 Hierarchical Attention

403 To address these limitations, we propose a novel scalable attention mechanism for PFNs that effec-  
404 tively leverages the natural hierarchical structure of multi-task data, as shown in Figure A.3. Our  
405 approach applies hierarchical attention [34] to the multi-task regression setting and uses specialized  
406 transformer blocks to separately model intra-task and inter-task relationships.

407 For intra-task encoding, we introduce a learnable “[Task]” token to each task that summarizes task-  
408 specific properties. The intra-task transformer blocks are responsible for learning the relationships of  
409 the data points within each task. By performing attention over these points, the intra-task block updates  
410 the embeddings for each data point and also updates the “[Task]” token with a summary embedding  
411 for the task, requiring  $O(D^2)$  total compute per task. Then, the inter-task encoders, responsible for  
412 learning the relationship between tasks, attend to these summary “[Task]” embeddings, with  $O(T^2)$   
413 complexity. This hierarchical design reduces the overall attention complexity from the naive global  
414 setting of  $O(D^2T^2)$  to  $O(D^2T + T^2)$ , enabling significantly longer contexts while still allowing  
415 for every data point to influence others. We interleave the intra-task and inter-task blocks in our  
416 architecture, although Chalkidis et al. [5] show that other topologies may also be effective

417 Our hierarchical attention directly addresses many of the limitations of other task encoders. First, our  
418 attention mechanism naturally handles inputs of varying lengths, allowing the model to generalize to  
419 any number of tasks. This flexibility ensures that even if the model encounters more tasks at test time  
420 than it did during training, it can still meaningfully integrate new task representations. Furthermore,  
421 our approach enables the model to dynamically learn task representations which depend on the data  
422 from the task, and its representation of each task evolves through the many layers of attention. This  
423 enables tasks with similar patterns to develop similar representations, allowing the model to better  
424 capture the potentially complex relationships between tasks.

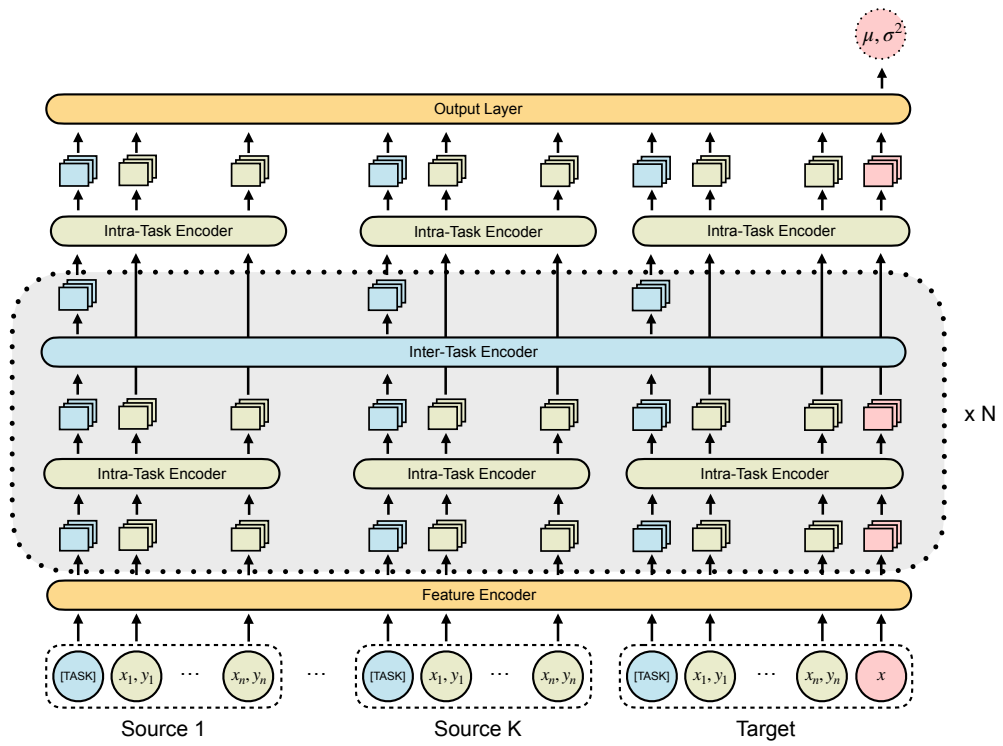


Figure A.3: MTPFNs use hierarchical attention and jointly model data across information sources.

## 425 D Empirical Results and Details

### 426 D.1 Setup Details

427 For our empirical results, we use a transformer backbone with 23 attention layers, where twelve  
428 intra-task attention layers are interwoven between eleven inter-task layers. Each attention layer  
429 has 4 attention heads with a hidden size of 512. The model is trained on approximately 50 million  
430 synthetically generated datasets as described in Section 2.1, with a batch size of 16 and AdamW with  
431 a learning rate of  $1e-4$  and cosine annealing.

432 We compare our method, MTPFN, to several baselines: (1) ICM [14], a joint method which trains a  
433 multi-task GP on the combined target and auxiliary data; (2) ScaML [31], an ensemble method that  
434 fits individual GPs to each auxiliary task; and (3) a single-task GP which only uses the target task and  
435 ignores the auxiliary tasks. Our Bayesian optimization results were implemented using BoTorch [2]  
436 and GPyTorch [11], and we will provide access to our code upon acceptance.

### 437 D.2 Robust to Negative Transfer

438 We explore the impact of  $p$ , the probability that any task is unrelated to target task from the data  
439 generation procedure. We train three MTPFNs with varying values of  $p$  to understand its effect on  
440 the robustness of the model to negative transfer. See Appendix B for further ablations.

441 For our test settings, we use an ICM to generate 5 different multi-task datasets with 3 input dimensions,  
442 each with 2 samples from the target task and 20 samples from each of the 3 auxiliary tasks. The  
443 auxiliary tasks have varying amounts of correlations with the target task. We then perform 5 runs of  
444 Bayesian inference for each multi-task dataset, and plot the mean and standard error of the mean.  
445 We find that MTPFNs trained with higher proportions of unrelated tasks are more robust to negative  
446 transfer and outperform other baselines.

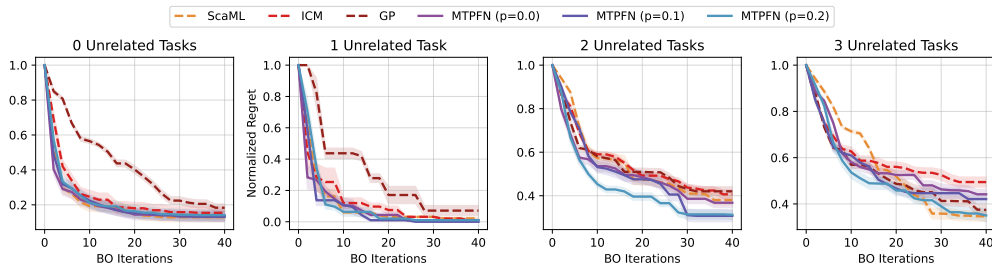


Figure A.4: **MTPFNs are robust to negative transfer from unrelated tasks.** We evaluate BO across four multi-task settings, where the target task is related to  $\{0, \dots, 3\}$  out of 3 auxiliary tasks. We compare the performance of MTPFNs trained with different  $p$ , where  $p$  represents the probability that an auxiliary task was drawn independently from the target task during training. As we increase the number of unrelated tasks during evaluation, the MTPFNs which were exposed to unrelated tasks during training ( $p > 0$ ) outperform the ICM model, which suffers from negative transfer. We plot the mean and standard error of the mean over 5 trials.

447 **D.3 Efficient Modeling of Inter-Task Relationships**

448 MTPFNs are able to do Bayesian inference with a single forward pass. Furthermore, our proposed  
 449 hierarchical attention mechanism enables the MTPFN to scale in  $O(TD^2 + T^2)$ , where  $D$  is the  
 450 number of data points per task and  $T$  is the number of tasks. We compare the runtime of MTPFNs to  
 451 joint-modeling methods such as ICM and ensemble-based methods such as ScaML in Figure A.5.  
 452 We see that MTPFNs are able to perform inference on an order of magnitude more data points and  
 453 tasks compared to traditional GP methods.

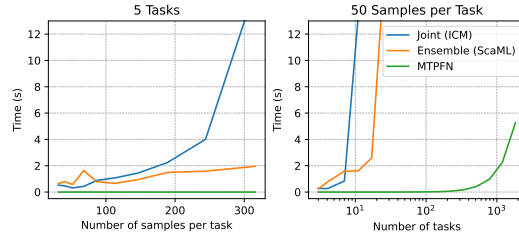


Figure A.5: MTPFNs are significantly faster than alternative GP-based methods.

454 **D.4 Fully Bayesian Inference**

455 When trained on a data-generation process that draws samples from a multi-task GP with an ICM  
 456 kernel, we see in Figure A.6 that the MTPFN and the MTGP (ICM kernel with MAP estimation)  
 457 have comparable behavior across varying levels of correlations. Furthermore, in low-data settings  
 458 demonstrated by Figure A.7, we find that the MTPFN outperforms the MTGP because it considers  
 459 the uncertainty over the task covariance matrix. This demonstrates that fully Bayesian inference may  
 460 be preferable to MAP estimation.

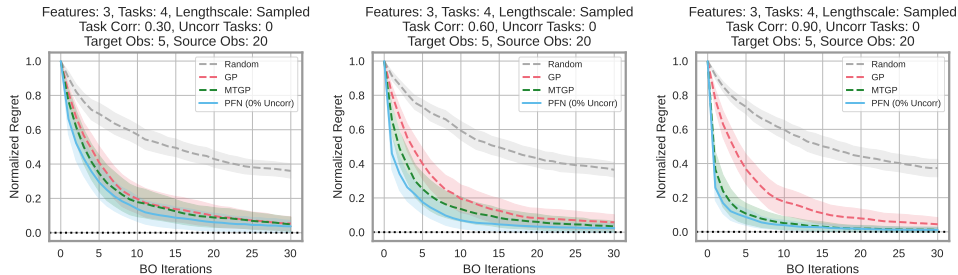


Figure A.6: ICM PFNs are comparable to MTGPs across varying levels of correlations between tasks.

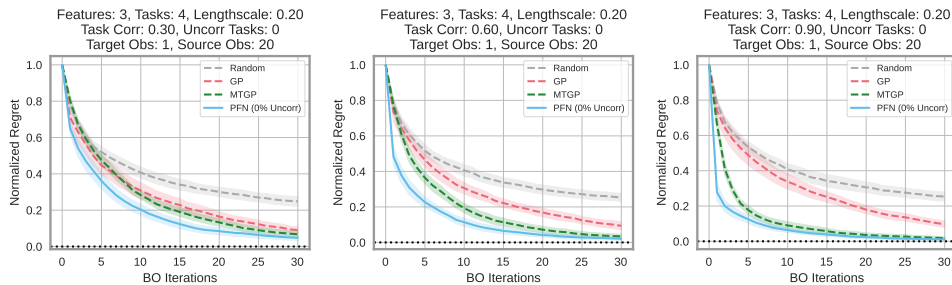


Figure A.7: In low-data settings, ICM PFNs, which approximate fully Bayesian inference, outperform MTGPs with MAP estimation. ICM PFNs are comparable to MTGPs across varying levels of correlations between tasks.

461 **D.5 Leverage Domain Data**

462 In Figure A.8, we demonstrate the benefits of in-context learning compared to fine-tuning on the  
 463 HPOBench dataset [8] for Logistic Regression (LR), which contains 25 tasks with 4 held out for  
 464 evaluation. We compare three approaches: (1) the “Original Single-Task”, a general-purpose single-  
 465 task model trained on Gaussian process draws with an RBF kernel; (2) the “Fine-Tuned Single-Task”,  
 466 the same base model after fine-tuning on data from the 25 LR training tasks; and (3) “MTPFNs”,  
 467 our method that uses only the 4 hold-out tasks in-context during inference, without any fine-tuning  
 468 on the 25 training tasks. We evaluate all models by measuring negative log likelihood on the LR  
 469 evaluation set as well as on other domains (SVM hyperparameter optimization and GP draws) to  
 470 assess generalization.

471 The results clearly demonstrate the benefits of MTPFNs over fine-tuned approaches. While fine-tuning  
 472 does improve performance on the target LR domain, it comes at a severe cost to generalization: as we  
 473 fine-tune on more LR samples, the performance on SVM and GP domains deteriorates significantly  
 474 due to overfitting. In contrast, MTPFNs, which use in-context learning, achieve comparable perfor-  
 475 mance on the target domain while maintaining strong generalization across all evaluated datasets.  
 476 This approach is also computationally efficient, requiring no parameter updates.

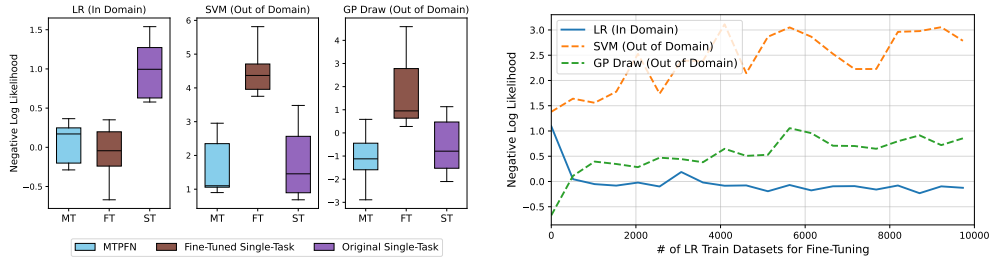


Figure A.8: **MTPFNs, which use domain data through in-context learning, match the performance of Fine-Tuned PFNs on in-domain data while generalizing better to other domains.** (Left): MTPFNs have comparable NLLs to Fine-Tuned PFNs on in-domain data (LR) and outperform Fine-Tuned PFNs on other domains (SVM and GP Draws). (Right): As we fine-tune on more in-domain data, the NLL for Fine-Tuned PFNs significantly worsens for other domains.

477 **Original Single-Task PFN Details** The data generation process for the single-task PFN randomly  
 478 samples inputs  $x$  from the unit cube, and then samples the corresponding outputs  $y$  by drawing a  
 479 sample from a GP with an RBF kernel with a lengthscale sampled from  $\text{Gamma}(3, 6)$ . For this  
 480 experiment, we use a fixed feature size of 2. We train an 8-layer standard transformer (not hierarchical  
 481 attention) with an embedding size of 256 on this data generation process for 4 million sampled  
 482 datasets, with a batch size of 16, and AdamW with a learning rate of  $1e-4$  and cosine annealing.

483 **Fine-Tuned Single-Task PFN Details** To fine-tune on the LR dataset, we develop a subsampling  
 484 data-generation procedure: On the 20 training tasks, we subsample within one task to get 50  $x, y$ . We  
 485 uniformly select some number of them to be used as ICL training, and the remaining to be used as  
 486 the test. We fine-tune the base model described above with a batch size of 16, and AdamW with a  
 487 learning rate of  $1e-4$  and cosine annealing.

488 **D.6 Real-World Optimization Problems**

489 We consider the hyperparameter optimization for five types of models: support vector machines  
 490 (SVM), logistic regression (LR), XGBoost (XGB), neural networks (NN), and random forest (RF).  
 491 For each setting, we randomly sample one task to be the target function, and we sample 3 auxiliary  
 492 tasks from the meta-data. We randomly sample 5 points from the target task and 20 points from each  
 493 of the auxiliary tasks to use as the initialization for Bayesian inference. We measure the normalized  
 494 regret  $(f^* - f_{\text{best}})/(f^* - f_0)$  where  $f^*$  is the optimal value,  $f_{\text{best}}$  is the best value so far, and  $f_0$  is  
 495 the initial value. We run 100 replicates, each with a different combination of target task and auxiliary  
 496 task initializations, and we plot the mean and one standard error.

497 We share the results of our benchmark in the top panel of Figure 3, and MTPFNs are competitive  
 498 across all of the model types. Specifically, we find that in instances where the meta-tasks contain

499 helpful information (ScaML and ICM outperform GP), the MTPFNs are also able to effectively  
500 utilize this data. Furthermore, in cases like XGB where there is negative transfer for the ICM model,  
501 we find that MTPFNs are more robust and perform similarly to the standard single-task GP.

502 We also consider tabular FC-Net benchmarks from: Slice Localization, Protein Structure, Naval  
503 Propulsion, and Parkinson’s Telemonitoring. For each benchmark, we set the benchmark to be the  
504 target task, and we use the three other tabular dataset as the auxiliary tasks. For instance, the results  
505 for Slice Localization use Protein Structure, Naval Propulsion, and Parkinson’s Telemonitoring as  
506 auxiliary data sources. We initialize our Bayesian optimization problem with a random sample of 5  
507 points from the target task and 20 points from each auxiliary task. We report the average normalized  
508 regret over 20 trials in the bottom panel of Figure 3. MTPFNs work well on these tabular datasets,  
509 often outperforming the other baselines.