# Explainable Meta Bayesian Optimization with Human Feedback for Scientific Applications like Fusion Energy

**Ricardo Luna Gutierrez**[1†], **Sahand Ghorbanpour**[1†], **Vineet Gundecha**[1†], **Rahman Ejaz**[2†],
**Varchas Gopalaswamy**[2], **Riccardo Betti**[2], **Avisek Naug**[1], **Desik Rengarajan**[1],
**Ashwin Ramesh Babu**[1], **Paolo Faraboschi**[1], **Soumyendu Sarkar**[1†*]

[1]Hewlett Packard Enterprise, [2]University of Rochester

{vineet.gundecha, rluna, sahand.ghorbanpour, avisek.naug, desik.rengarajan,
ashwin.ramesh-babu, paolo.faraboschi, soumyendu.sarkar}@hpe.com
{reja, vgopalas, betti}@lle.rochester.edu

## Abstract

We introduce Meta Bayesian Optimization with Human Feedback (MBO-HF), which integrates Meta-Learning and expert preferences to enhance BO. MBO-HF employs Transformer Neural Processes (TNPs) to create a meta-learned surrogate model and a human-informed acquisition function (AF) to suggest and explain proposed candidate experiments. MBO-HF outperforms current methods in optimizing various scientific experiments and benchmarks in simulation, including the energy yield of the inertial confinement fusion (ICF), practical molecular optimization (PMO), and critical temperature maximization for superconducting materials.

## 1 Introduction

Due to the high costs and limited experimentation opportunities of some scientific applications, efficient optimization techniques incorporating expert knowledge are crucial. Bayesian Optimization (BO) has shown potential in optimizing expensive black-box functions (Feurer et al., 2015; Snoek et al., 2012; Wang et al., 2023; Shmakov et al., 2023; Gundecha et al., 2024; Ghorbanpour et al., 2024) but faces challenges in real-world applications due to trustworthiness concerns (Bouthillier and Varoquaux, 2020). Recent advancements in human-AI collaborative approaches have improved BO by incorporating expert input (Colella et al., 2020; A V et al., 2022; Gupta et al., 2023; Hvarfner et al., 2022; Adachi et al., 2024), though these methods typically focus on single-task scenarios, failing to leverage knowledge from previous tasks.

Meta-Learning within the BO framework (Meta-BO) (Volpp et al., 2020) addresses sample efficiency by enabling models to adapt quickly to new tasks using knowledge from related tasks, enhancing sample efficiency. However, trustworthiness and explainability, critical for the adoption of black-box optimization by scientists, are under-explored in Meta-BO. Moreover, previous Meta-BO studies have not investigated the impact of integrating expert knowledge through methods like preference learning on optimization performance. To address this, we introduce Meta Bayesian Optimization with Human Feedback (MBO-HF), a framework that integrates expert knowledge into Meta-BO using Transformer Neural Processes (TNPs) to build an adaptable surrogate model. MBO-HF provides an explainable system that helps experts evaluate optimization recommendations, improving both

---

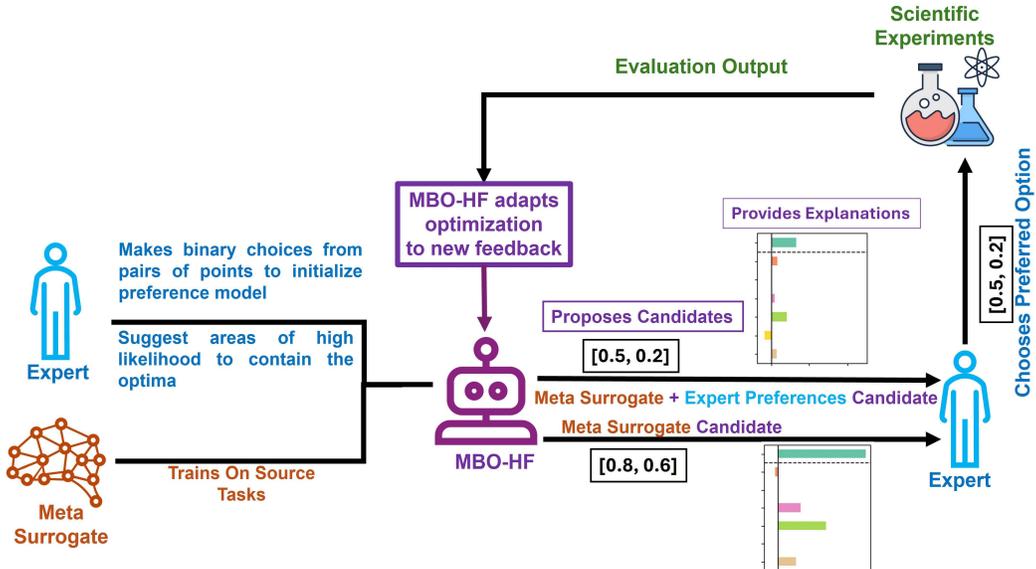[*]Corresponding author. [†]These authors contributed equally.

Figure 1: MBO-HF combines meta-learning and experts' knowledge for a highly efficient optimization process, with the expert-in-the-loop.

the effectiveness and the trustworthiness of the optimization process, as shown in Figure 1. We summarize our main contributions as follows:

- We introduce Meta Bayesian Optimization with Human Feedback (MBO-HF), an explainable framework that enhances Meta-BO by incorporating human expertise. We demonstrate that MBO-HF outperforms SOTA in both human-aided BO and standard Meta-BO.

- We present a novel expert-informed acquisition function (AF) that integrates preference learning with expert hypotheses, and evaluate the impact of human expertise on optimization.

- We implement an explainable framework for Meta-BO, utilizing Shapley values (Shapley, 1952) to provide the human-in-the-loop with insights into the optimization decisions.

## 2   MBO-HF

Unlike single-task approaches, MBO-HF trains a meta-surrogate using Transformer Neural Processes (TNPs) on prior tasks and incorporates expert preferences to optimize the target function effectively.

### 2.1   Preference Learning

To fit the preference model, binary preference learning (Bradley and Terry, 1952; Lun Chau et al., 2022; Adachi et al., 2024) is used, where users compare pairs of candidate points $\{x_1, x_2\}$, to indicate which is more favorable, generating a dataset that informs the preference function. This process is repeated $M$ times, resulting in a dataset $D_{\text{pref}} = \{x_1^i, x_2^i, y_{\text{pref}}^i\}_{i=1}^M$, where $y_{\text{pref}}$ is 1 if $x_1$ is favored over $x_2$, and 0 otherwise. $D_{\text{pref}}$ is then used to build a binary preference function $g$, which follows a likelihood model:

$$\mathbb{P}(y_{\text{pref}}|x_1, x_2) = S(y_{\text{pref}}; g(x_1, x_2)) \tag{1}$$

where $S(y_{\text{pref}}; z) := z^{y_{\text{pref}}}(1-z)^{1-y_{\text{pref}}}$ represents a Bernoulli likelihood and $g(x_1, x_2)$ represents the user's preference level for $x_1$ compared to $x_2$. This function is modeled using Dirichlet-based Gaussian Processes (GPs) (Adachi et al., 2024) with skew-symmetric data augmentation, estimating a user-aligned posterior distribution for target locations. Following this approach allows us to estimate a preference model $\pi$ as GPs that generate a distribution $p_\pi(\cdot|X_T, \mathcal{D}_c)$ that approximates a user aligned posterior given a set of target locations $X_T$ and context samples $\mathcal{D}_c$. Instead of randomly selecting

candidate pairs from the entire search space, we explore the use of experts' hypotheses (Cisse et al., 2024), which define promising regions for pair selection during the creation of the dataset. In our framework, preference learning is applied only during the testing phase for the target function, not during meta-training on source functions.

## 2.2 Acquisition Function

In expert-aided Bayesian Optimization (BO), CoExBO (Adachi et al., 2024) combines a surrogate model and a preference model, using Gaussian Processes (GPs), to create an acquisition function (AF) that balances model predictions with expert preferences. However, CoExBO is limited to single-task cases and uses an Upper Confidence Bound (UCB) based AF, which can be outperformed by Expected Improvement (EI) (Merrill et al., 2021). To refine these, we propose replacing the GP-based surrogate with a Transformer Neural Process (TNP) $\mathcal{M}$ for improved optimization by leveraging past data and reformulating the AF to an EI-based approach, while retaining the preference-surrogate trade-off as:

$$\alpha_{\mathcal{M},\pi}(x) = (\mu_{\mathcal{M},\pi}(x) - f(x^+) - \xi)\Phi(Z) + \sigma_{\mathcal{M},\pi}(x)\phi(Z)$$

$$Z = \frac{\mu_{\mathcal{M},\pi}(x) - f(x^+) - \xi}{\sigma_{\mathcal{M},\pi}(x)}$$

$$\mu_{\mathcal{M},\pi}(x) := \frac{\sigma^2_{\mathcal{M},\pi}(x)}{\mathscr{S}^2_\pi(x)}\mu_\pi(x) + \frac{\sigma^2_{\mathcal{M},\pi}(x)}{\sigma^2_{\mathcal{M}}(x)}\mu_{\mathcal{M}}(x)$$

$$\sigma^2_{\mathcal{M},\pi}(x) := \frac{\mathscr{S}^2_\pi(x)\sigma^2_{\mathcal{M}}(x)}{\mathscr{S}^2_\pi(x) + \sigma^2_{\mathcal{M}}(x)}$$

$$\mathscr{S}^2_\pi(x) := \sigma^2_\pi(x) + \gamma^{t^2}\sigma^2_{\mathcal{M}}(x)$$

where $f(x^+)$ is the best observed value so far, the predictive mean $\mu_{\mathcal{M}}$ and variance $\sigma^2_{\mathcal{M}}$ of $\mathcal{M}$, $\Phi(\cdot)$ is the cumulative distribution function (CDF) and $\phi(\cdot)$ is the probability density function (PDF) of $Z$ (Tacq, 2010). The parameter $\xi$ is a small positive value to encourage exploration. $\gamma$ is a decay factor and $t$ represents the current time step in the optimization process. $\gamma$ ensures that the information provided by $\pi$ decays, so the impact of $\pi$ is high at the start of the optimization but allows $\mathcal{M}$ to take over at later stages. Furthermore, we evaluated AF $\alpha_{\mathcal{M}}(x)$, which applies the standard EI over $\mu_{\mathcal{M}}$ and $\sigma^2_{\mathcal{M}}$, without preference information.

## 2.3 Explainability

To support human-experts in making informed decisions during optimization, we extend the CoExBO framework to Meta-BO and incorporate Shapley values for explainability. In this work, to calculate Shapley values, we adapt KernelExplainer from SHAP (Lundberg and Lee, 2017). This quantify the contribution of each feature to the model's output, helping users understand the decisions made by the acquisition functions during the optimization process.

---

**Algorithm 1** MBO-HF algorithm

---

**Input:** trained meta-surrogate model $\mathcal{M}$, preference model $\pi$, target black-box function $f_{\text{target}}$, target locations $X_T$, $\mathcal{I}$ initial samples, optimization budget $B$.
**Output:** Points to evaluate.

1: Collect $\mathcal{I}$ input-output pairs $X_{\mathcal{I}}$.
2: Initialize $D_c \leftarrow X_{\mathcal{I}}$.
3: Set $i \leftarrow 0$.
4: **while** $i < B$ **do**
5:     Forward pass on $\mathcal{M}$ to obtain posterior distribution $p(\cdot|X_T, \mathcal{D}_c)$.
6:     Run acquisition functions $\alpha_{\mathcal{M}}$ and $\alpha_{\mathcal{M},\pi}$ to propose candidate points $\{x_1, x_2\}$.
7:     Run explainability framework on candidates $\{x_1, x_2\}$ and present to user.
8:     Expert selects candidate $x_p \in \{x_1, x_2\}$.
9:     Evaluate $f_{\text{target}}$ on $x_p$ to obtain $y_p$.
10:     Update $D_c \leftarrow D_c \cup \{x_p, y_p\}$.
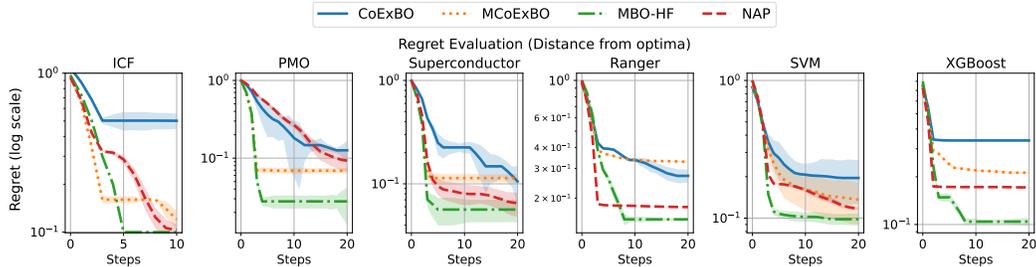11:     Increment $i$ by 1.
12: **end while**

---

Figure 2: Regret minimization evaluation, ower the better. MBO-HF outperforms the baselines on all 6 benchmarks. The shaded areas represent a $\pm 1$ SD. Regrets are normalized.
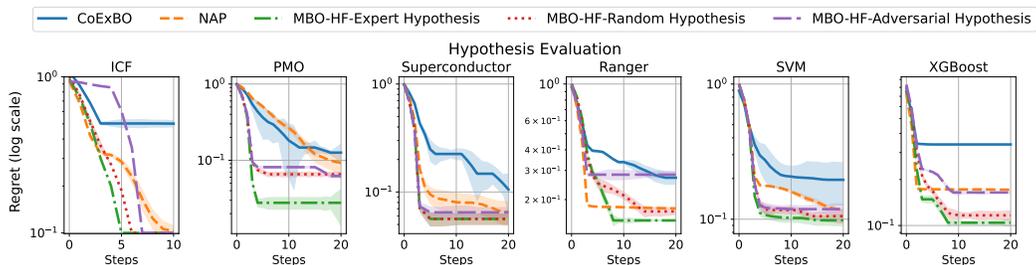


Figure 3: Comparison between different hypothesis: "Expert" best-case, "Adversarial" worst-case, and Random sampling in the preference learning initialization. The shaded areas represent a $\pm 1$ SD. Regrets are normalized.

# 3 Experiments

To conduct our evaluations, we make use of the HPO-B benchmark (Pineda Arango et al., 2021). HPO-B is a popular hyperparameter optimization benchmark used to evaluate Meta-BO approaches. HPO-B encompasses a collection of datasets, featuring hyperparameters of different classification models alongside their respective accuracy, spanning various categories and search spaces. For each different search space, HPO-B provides a training (source), validation and test (target) set, which we used as is. We selected 3 different datasets for this evaluation; Ranger (6D), SVM (8D), and XGBoost (16D), covering problems from 6 to 16 input dimensions (hyperparameters to optimize).

Moreover, we evaluate MBO-HF in the real-world tasks of energy yield optimization for inertial confinement fusion **(ICF)**, (Betti and Hurricane, 2016; Lees et al., 2021; Gopalaswamy et al., 2019, 2024), practical molecular optimization **(PMO)** (Gao et al., 2022), and critical temperature maximization for **(Superconductor)** materials (Trabucco et al., 2022). The input dimensions are 5, 86 and 2048 for ICF, superconductor and PMO respectively. Furthermore, we evaluated our approach using the HPO-B benchmark (Pineda Arango et al., 2021), which includes datasets with hyperparameters and accuracy metrics for various classification models across different search spaces. We selected Ranger 6D, SVM 8D and XGBoost 16D datasets to cover problems with 6 to 16 input dimensions.

# 4 Results

MBO-HF outperforms previous BO approaches in terms of simple regret while optimizing unseen target tasks as seen in Figure 2. We compare MBO-HF against CoExBO (Adachi et al., 2024), the SOTA approach for human-in-the-loop BO, and NAP (Maraval et al., 2023), the SOTA for Meta-BO. Moreover, to evaluate the isolated impact of our proposed AF we design a new baseline, which we call MCoExBO. MCoExBO uses the same AF proposed in CoExBO. However, the single-task GPs are replaced by TNPs. Both MBO-HF and MCoExBO were meta-trained on the same training sets. Similarly, for fairness of comparison, the GP used by CoExBO was initialized with the training data.

Figure 3 shows the ablation study where we explore the influence of experts' hypotheses on the performance of MBO-HF. Similar to Cisse et al. (2024), we assume access to the target function

$f_{target}$ and consider three experimental settings: Expert Hypothesis (EH), Random Hypothesis (RH), and Adversarial Hypothesis (AH).

## 5 Conclusion

We presented MBO-HF, an explainable approach for Meta Bayesian Optimization with human feedback. We evaluated our method on energy yield optimization in ICF and across five different hyperparameter optimization tasks. Our results showed that MBO-HF surpasses current SOTA for human-AI collaboration and Meta-BO. Moreover, we assessed the impact of experts' knowledge on the optimization process and showed that its integration with Meta-BO is beneficial. Moreover enabling human-in-the-loop is essential for ICF given its experimental complexities and variabilities.

## 6 Acknowledgement

## References

M. Feurer, A. Klein, K. Eggensperger, J. Springenberg, M. Blum, F. Hutter, Efficient and robust automated machine learning, in: C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, R. Garnett (Eds.), Advances in Neural Information Processing Systems, volume 28, Curran Associates, Inc., 2015. URL: `https://proceedings.neurips.cc/paper_files/paper/2015/file/11d0e6287202fced83f79975ec59a3a6-Paper.pdf`.

J. Snoek, H. Larochelle, R. P. Adams, Practical bayesian optimization of machine learning algorithms, in: F. Pereira, C. Burges, L. Bottou, K. Weinberger (Eds.), Advances in Neural Information Processing Systems, volume 25, Curran Associates, Inc., 2012. URL: `https://proceedings.neurips.cc/paper_files/paper/2012/file/05311655a15b75fab86956663e1819cd-Paper.pdf`.

X. Wang, Y. Jin, S. Schmitt, M. Olhofer, Recent advances in bayesian optimization, ACM Comput. Surv. 55 (2023). URL: `https://doi.org/10.1145/3582078`. doi:10.1145/3582078.

A. Shmakov, A. Naug, V. Gundecha, S. Ghorbanpour, R. L. Gutierrez, A. R. Babu, A. Guillen, S. Sarkar, Rtdk-bo: High dimensional bayesian optimization with reinforced transformer deep kernels, in: 2023 IEEE 19th International Conference on Automation Science and Engineering (CASE), IEEE, 2023, pp. 1–8.

V. Gundecha, R. Gutierrez Luna, S. Ghorbanpour, R. Ejaz, V. Gopalaswamy, R. Betti, A. Naug, P. Faraboschi, S. Sarkar, Meta-learned bayesian optimization for energy yield in inertial confinement fusion, in: NeurIPS 2024 Workshop on Tackling Climate Change with Machine Learning, 2024.

S. Ghorbanpour, R. L. Gutierrez, V. Gundecha, D. Rengarajan, A. R. Babu, S. Sarkar, Llm enhanced bayesian optimization for scientific applications like fusion, in: NeurIPS Workshop on Machine Learning and the Physical Sciences, 2024.

X. Bouthillier, G. Varoquaux, Survey of machine-learning experimental methods at NeurIPS2019 and ICLR2020, Research Report, Inria Saclay Ile de France, 2020. URL: `https://hal.science/hal-02447823`.

F. Colella, P. Daee, J. Jokinen, A. Oulasvirta, S. Kaski, Human strategic steering improves performance of interactive optimization, in: Proceedings of the 28th ACM Conference on User Modeling, Adaptation and Personalization, UMAP '20, ACM, 2020. URL: `http://dx.doi.org/10.1145/3340631.3394883`. doi:10.1145/3340631.3394883.

A. K. A V, S. Rana, A. Shilton, S. Venkatesh, Human-ai collaborative bayesian optimisation, in: S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, A. Oh (Eds.), Advances in Neural Information Processing Systems, volume 35, Curran Associates, Inc., 2022, pp. 16233–16245. URL: `https://proceedings.neurips.cc/paper_files/paper/2022/file/6751611b394a3464cea53eed91cf163c-Paper-Conference.pdf`.

S. Gupta, A. Shilton, A. K. A. V. au2, S. Ryan, M. Abdolshah, H. Le, S. Rana, J. Berk, M. Rashid, S. Venkatesh, Bo-muse: A human expert and ai teaming framework for accelerated experimental design, 2023. `arXiv:2303.01684`.

C. Hvarfner, D. Stoll, A. Souza, M. Lindauer, F. Hutter, L. Nardi, $\pi$bo: Augmenting acquisition functions with user beliefs for bayesian optimization, 2022. `arXiv:2204.11051`.

M. Adachi, B. Planden, D. A. Howey, M. A. Osborne, S. Orbell, N. Ares, K. Muandet, S. L. Chau, Looping in the human collaborative and explainable bayesian optimization, 2024. `arXiv:2310.17273`.

M. Volpp, L. P. Fröhlich, K. Fischer, A. Doerr, S. Falkner, F. Hutter, C. Daniel, Meta-learning acquisition functions for transfer learning in bayesian optimization, 2020. URL: `https://arxiv.org/abs/1904.02642`. `arXiv:1904.02642`.

L. S. Shapley, A Value for N-Person Games, RAND Corporation, Santa Monica, CA, 1952. doi:`10.7249/P0295`.

R. A. Bradley, M. E. Terry, Rank analysis of incomplete block desings: The method of pairder comparisons, Biometrika 39 (1952) 324–345. URL: `https://doi.org/10.1093/biomet/39.3-4.324`. doi:`10.1093/biomet/39.3-4.324`.

S. Lun Chau, J. Gonzalez, D. Sejdinovic, Learning inconsistent preferences with gaussian processes, in: G. Camps-Valls, F. J. R. Ruiz, I. Valera (Eds.), Proceedings of The 25th International Conference on Artificial Intelligence and Statistics, volume 151 of *Proceedings of Machine Learning Research*, PMLR, 2022, pp. 2266–2281. URL: `https://proceedings.mlr.press/v151/lun-chau22a.html`.

A. Cisse, X. Evangelopoulos, S. Carruthers, V. V. Gusev, A. I. Cooper, Hypbo: Accelerating black-box scientific experiments using experts' hypotheses, 2024. URL: `https://arxiv.org/abs/2308.11787`. `arXiv:2308.11787`.

E. Merrill, A. Fern, X. Fern, N. Dolatnia, An empirical study of bayesian optimization: Acquisition versus partition, Journal of Machine Learning Research 22 (2021) 1–25. URL: `http://jmlr.org/papers/v22/18-220.html`.

J. Tacq, The normal distribution and its applications, in: P. Peterson, E. Baker, B. McGaw (Eds.), International Encyclopedia of Education (Third Edition), third edition ed., Elsevier, Oxford, 2010, pp. 467–473. URL: `https://www.sciencedirect.com/science/article/pii/B9780080448947015633`. doi:`https://doi.org/10.1016/B978-0-08-044894-7.01563-3`.

S. M. Lundberg, S.-I. Lee, A unified approach to interpreting model predictions, in: I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, R. Garnett (Eds.), Advances in Neural Information Processing Systems 30, Curran Associates, Inc., 2017, pp. 4765–4774. URL: `http://papers.nips.cc/paper/7062-a-unified-approach-to-interpreting-model-predictions.pdf`.

R. Betti, O. A. Hurricane, Inertial-confinement fusion with lasers, Nature Physics 12 (2016) 435–448. URL: `https://doi.org/10.1038/nphys3736`. doi:`10.1038/nphys3736`.

A. Lees, R. Betti, J. P. Knauer, V. Gopalaswamy, D. Patel, K. M. Woo, K. S. Anderson, E. M. Campbell, D. Cao, J. Carroll-Nellenback, R. Epstein, C. Forrest, V. N. Goncharov, D. R. Harding, S. X. Hu, I. V. Igumenshchev, R. T. Janezic, O. M. Mannion, P. B. Radha, S. P. Regan, A. Shvydky, R. C. Shah, W. T. Shmayda, C. Stoeckl, W. Theobald, C. Thomas, Experimentally inferred fusion yield dependencies of omega inertial confinement fusion implosions, Phys. Rev. Lett. 127 (2021) 105001. URL: `https://link.aps.org/doi/10.1103/PhysRevLett.127.105001`. doi:`10.1103/PhysRevLett.127.105001`.

V. Gopalaswamy, R. Betti, J. P. Knauer, N. Luciani, D. Patel, K. M. Woo, A. Bose, I. V. Igumenshchev, E. M. Campbell, K. S. Anderson, K. A. Bauer, M. J. Bonino, D. Cao, A. R. Christopherson, G. W. Collins, T. J. B. Collins, J. R. Davies, J. A. Delettrez, D. H. Edgell, R. Epstein, C. J. Forrest, D. H. Froula, V. Y. Glebov, V. N. Goncharov, D. R. Harding, S. X. Hu, D. W. Jacobs-Perkins, R. T.

Janezic, J. H. Kelly, O. M. Mannion, A. Maximov, F. J. Marshall, D. T. Michel, S. Miller, S. F. B. Morse, J. Palastro, J. Peebles, P. B. Radha, S. P. Regan, S. Sampat, T. C. Sangster, A. B. Sefkow, W. Seka, R. C. Shah, W. T. Shmyada, A. Shvydky, C. Stoeckl, A. A. Solodov, W. Theobald, J. D. Zuegel, M. G. Johnson, R. D. Petrasso, C. K. Li, J. A. Frenje, Tripled yield in direct-drive laser fusion through statistical modelling, Nature 565 (2019) 581–586.

V. Gopalaswamy, C. A. Williams, R. Betti, D. Patel, J. P. Knauer, A. Lees, D. Cao, E. M. Campbell, P. Farmakis, R. Ejaz, K. S. Anderson, R. Epstein, J. Carroll-Nellenbeck, I. V. Igumenshchev, J. A. Marozas, P. B. Radha, A. A. Solodov, C. A. Thomas, K. M. Woo, T. J. B. Collins, S. X. Hu, W. Scullin, D. Turnbull, V. N. Goncharov, K. Churnetski, C. J. Forrest, V. Y. Glebov, P. V. Heuer, H. McClow, R. C. Shah, C. Stoeckl, W. Theobald, D. H. Edgell, S. Ivancic, M. J. Rosenberg, S. P. Regan, D. Bredesen, C. Fella, M. Koch, R. T. Janezic, M. J. Bonino, D. R. Harding, K. A. Bauer, S. Sampat, L. J. Waxer, M. Labuzeta, S. F. B. Morse, M. Gatu-Johnson, R. D. Petrasso, J. A. Frenje, J. Murray, B. Serrato, D. Guzman, C. Shuldberg, M. Farrell, C. Deeney, Demonstration of a hydrodynamically equivalent burning plasma in direct-drive inertial confinement fusion, Nature Physics 20 (2024) 751–757.

W. Gao, T. Fu, J. Sun, C. Coley, Sample efficiency matters: a benchmark for practical molecular optimization, Advances in Neural Information Processing Systems 35 (2022) 21342–21357.

B. Trabucco, X. Geng, A. Kumar, S. Levine, Design-bench: Benchmarks for data-driven offline model-based optimization, CoRR abs/2202.08450 (2022). URL: `https://arxiv.org/abs/2202.08450`. arXiv:2202.08450.

S. Pineda Arango, H. Jomaa, M. Wistuba, J. Grabocka, Hpo-b: A large-scale reproducible benchmark for black-box hpo based on openml, in: J. Vanschoren, S. Yeung (Eds.), Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks, volume 1, 2021. URL: `https://datasets-benchmarks-proceedings.neurips.cc/paper_files/paper/2021/file/ec8956637a99787bd197eacd77acce5e-Paper-round2.pdf`.

A. Maraval, M. Zimmer, A. Grosnit, H. B. Ammar, End-to-end meta-bayesian optimisation with transformer neural processes, 2023. arXiv:2305.15930.