

# A Trajectory-Based Bayesian Approach to Multi-Objective Hyperparameter Optimization with Epoch-Aware Trade-Offs

Wenyu Wang<sup>1</sup>, Zheyi Fan<sup>2,3</sup>, Szu Hui Ng<sup>1,\*</sup>

<sup>1</sup>Industrial Systems Engineering and Management, National University of Singapore, Singapore

<sup>2</sup>Academy of Mathematics and System Science, Chinese Academy of Sciences, China

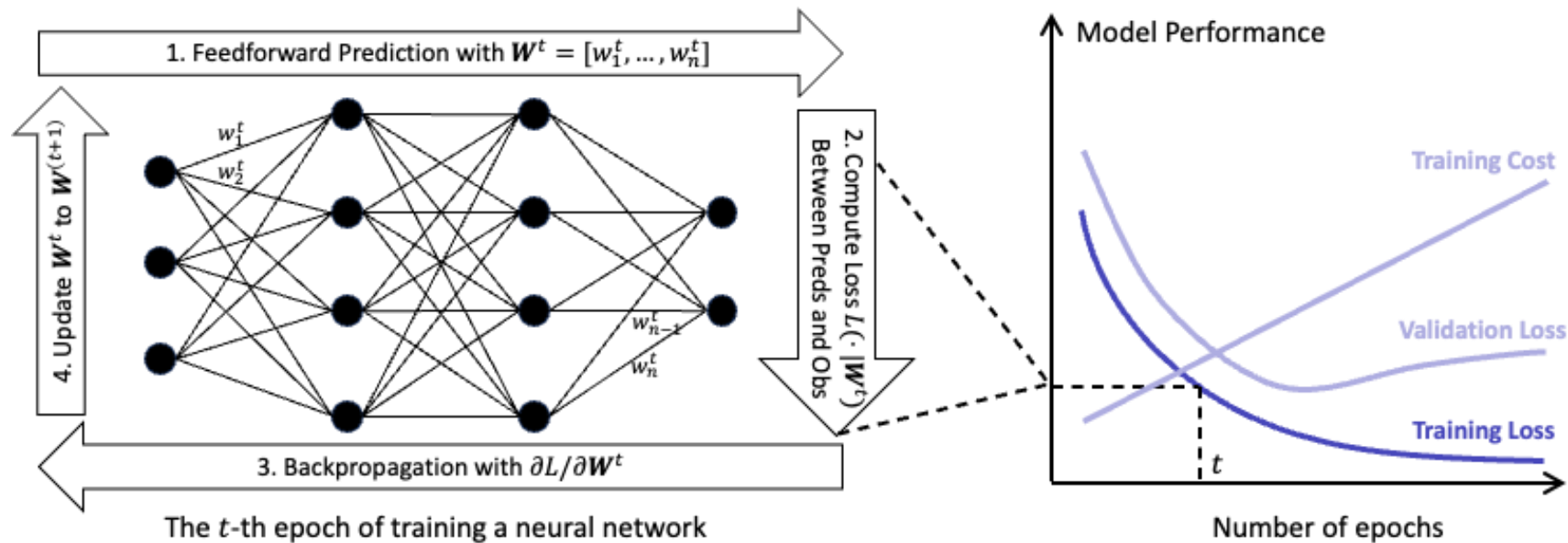
<sup>3</sup>School of Mathematical Sciences, University of Chinese Academy of Sciences, China

# Outline

1. Introduction
2. Problem and Methodology
  - 2.1 Enhanced Multi-Objective Hyperparameter Tuning
  - 2.2 Trajectory-Based Bayesian Optimization Approach
3. Numerical Experiments
4. Conclusions

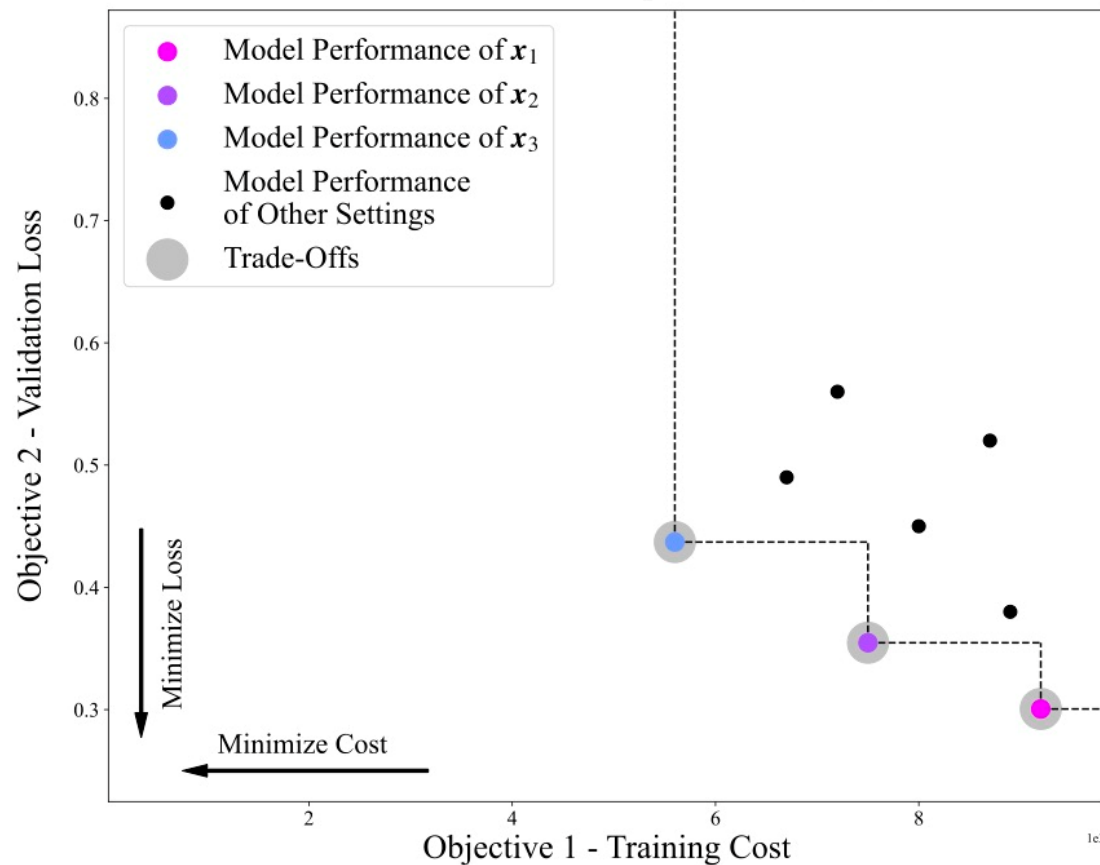
# Hyperparameter Optimization

- Addressing Hyperparameter Optimization (HPO) problem has long been challenging as it involves **resource-intensive model training** that prevents optimizers from exhaustively exploring the hyperparameter space.
- More recently, the surge in the demand for HPO is not only in pursuit of prediction accuracy but also for ensuring the efficiency and robustness of models, which leads to **Multi-Objective HPO (MOHPO)**.



# MOHPO with Iterative Learning Procedure

An Illustrative Example for MOHPO

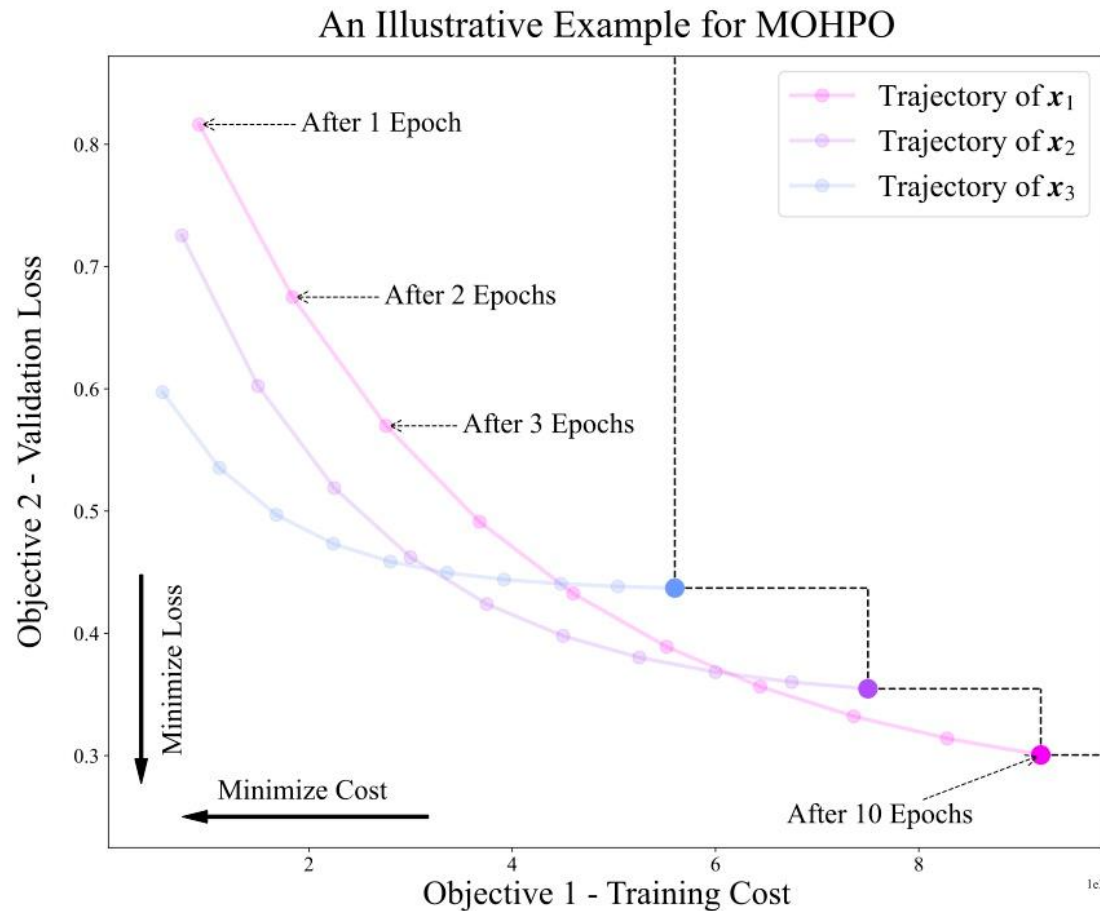


- Minimizing an MOHPO is equivalent to finding **trade-offs** between performances at the end of training.

$$\min_{x \in \mathbb{X}} f(x) = [f_1(x), f_2(x)]$$

- A solution **dominates** another if it is no worse in all objectives and strictly better in at least one.
- The **Pareto-optimal front** consists of all non-dominated solutions

# MOHPO with Iterative Learning Procedure



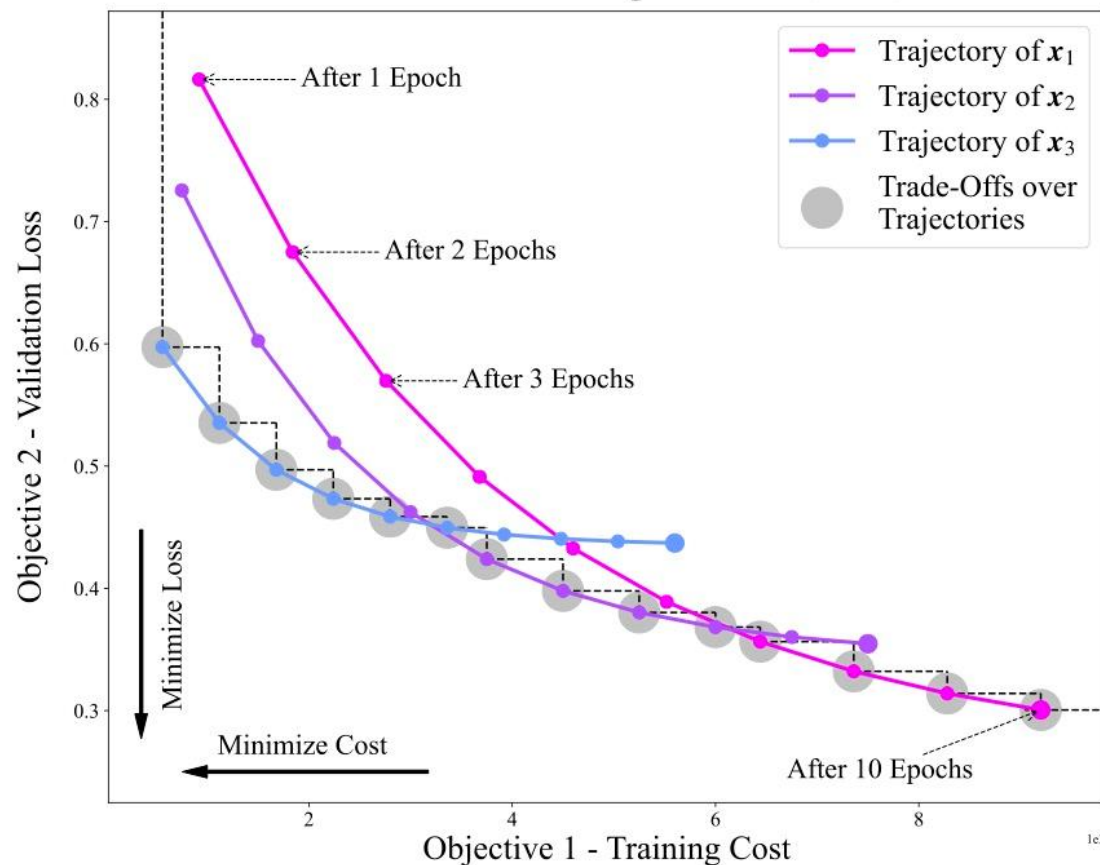
- Minimizing an MOHPO is equivalent to finding **trade-offs** between performances at the end of training.

$$\min_{\mathbf{x} \in \mathbb{X}} \mathbf{f}(\mathbf{x}) = [f_1(\mathbf{x}), f_2(\mathbf{x})]$$

- ❑ A solution **dominates** another if it is no worse in all objectives and strictly better in at least one.
- ❑ The **Pareto-optimal front** consists of all non-dominated solutions
- Training ML model is an **iterative learning procedure**, allowing epoch-wise tracking on model performances.
  - ❑ Does a trade-off emerge when the number of training epochs is fewer than the maximum allowed?

# MOHPO with Iterative Learning Procedure

An Illustrative Example for EMOHPO



- Minimizing an MOHPO is equivalent to finding **trade-offs** between performances at the end of training.

$$\min_{\mathbf{x} \in \mathbb{X}} \mathbf{f}(\mathbf{x}) = [f_1(\mathbf{x}), f_2(\mathbf{x})]$$

- A solution **dominates** another if it is no worse in all objectives and strictly better in at least one.
- The **Pareto-optimal front** consists of all non-dominated solutions

- Training ML model is an **iterative learning procedure**, allowing epoch-wise tracking on model performances.

- Does a trade-off emerge when the number of training epochs is fewer than the maximum allowed?
- E.g., (1) Partially-trained model; (2) Overfitting.

$$\min_{(\mathbf{x}, t) \in \mathbb{X} \times \mathbb{T}} \mathbf{f}(\mathbf{x}, t) = [f_1(\mathbf{x}, t), f_2(\mathbf{x}, t)]$$

# Outline

1. Introduction
2. Problem and Methodology
  - 2.1 Enhanced Multi-Objective Hyperparameter Tuning
  - 2.2 Trajectory-Based Bayesian Optimization Approach
3. Numerical Experiments
4. Conclusions

# Enhanced Multi-Objective Hyperparameter Optimization Problem

Consider the sequential minimization of an EMOHPO in the following form:

$$\min_{(\mathbf{x}, t) \in \mathbb{X} \times \mathbb{T}} \mathbf{f}(\mathbf{x}, t) = [f_1(\mathbf{x}, t), \dots, f_k(\mathbf{x}, t)], \quad (2)$$

- $\mathbf{x}$  denotes a  $d$ -dimensional hyperparameter setting with  $\mathbf{x} \in \mathbb{X} \subset \mathbb{R}^d$ .
- $t$  denotes the number of training epochs with  $t \in \mathbb{T} = \{1, \dots, t_{max}\}$ .
- $\mathbf{f}: \mathbb{X} \times \mathbb{T} \mapsto \mathbb{R}^k$  comprises  $k$  objectives, each of which represents a specific performance measure of the ML model after training with setting  $\mathbf{x}$  for  $t$  epochs.



# Enhanced Multi-Objective Hyperparameter Optimization Problem

Consider the sequential minimization of an EMOHPO in the following form:

$$\min_{(\mathbf{x}, t) \in \mathbb{X} \times \mathbb{T}} \mathbf{f}(\mathbf{x}, t) = [f_1(\mathbf{x}, t), \dots, f_k(\mathbf{x}, t)], \quad (2)$$

- $\mathbf{x}$  denotes a  $d$ -dimensional hyperparameter setting with  $\mathbf{x} \in \mathbb{X} \subset \mathbb{R}^d$ .
- $t$  denotes the number of training epochs with  $t \in \mathbb{T} = \{1, \dots, t_{max}\}$ .
- $\mathbf{f}: \mathbb{X} \times \mathbb{T} \mapsto \mathbb{R}^k$  comprises  $k$  objectives, each of which represents a specific performance measure of the ML model after training with setting  $\mathbf{x}$  for  $t$  epochs.
- Assume that when querying at any feasible pair  $\mathbf{z} = (\mathbf{x}, t) \in \mathbb{X} \times \mathbb{T}$ ,

- **Noisy Observation:** Each observed model performance is noisy, i.e., for any  $i = 1, \dots, k$ ,

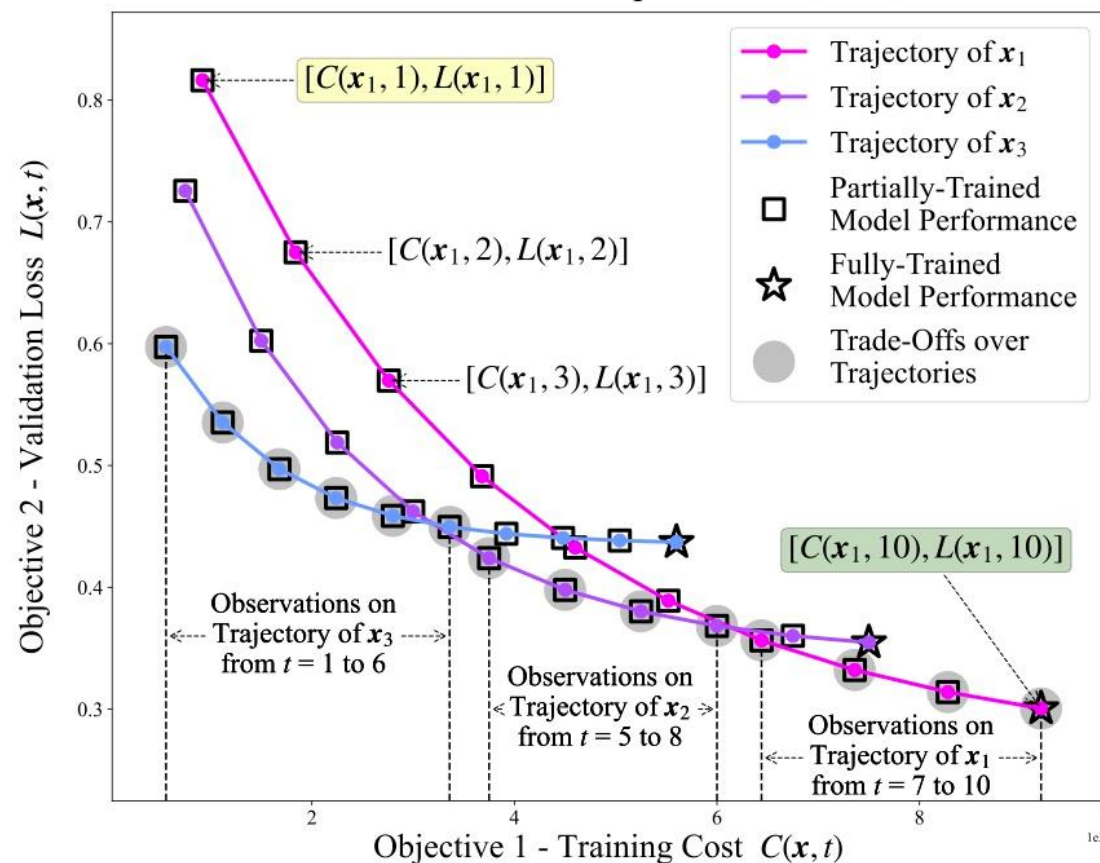
$$y_i(\mathbf{x}, t) = f_i(\mathbf{x}, t) + \varepsilon_i \quad \text{and} \quad \varepsilon_i \sim \mathcal{N}(0, \sigma_i^2).$$

- **Iterative Learning:** A sequence of multi-objective model performances are observed, i.e.,  $\{\mathbf{y}(\mathbf{x}, 1), \dots, \mathbf{y}(\mathbf{x}, t)\}$  with

$$\mathbf{y}(\mathbf{x}, \tau) = [y_1(\mathbf{x}, \tau), \dots, y_k(\mathbf{x}, \tau)], \quad \forall \tau = 1, \dots, t.$$

# Challenges in Solving EMOHPO

An Illustrative Example for EMOHPO



In the objective space of EMOHPO,

## 1. How to make prediction on trajectory?

The model should be able to capture the characteristics of the trajectory as the epoch changes.

# Gaussian Process for Trajectory Prediction

## ➤ The Prior:

- Consider a function  $f(\mathbf{z})$  to be sampled from a **Gaussian Process (GP)** with kernel  $K(\mathbf{z}, \mathbf{z}')$  and let  $K(Z, Z) \in \mathbb{R}^{n \times n}$  with  $[K(Z, Z)]_{i,j} = K(\mathbf{z}_i, \mathbf{z}_j)$ .

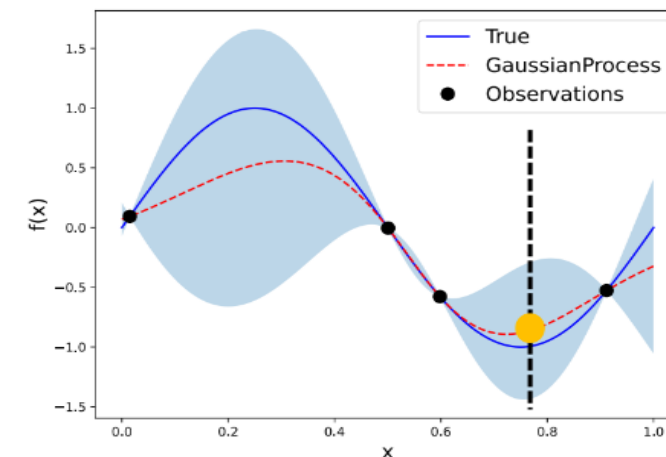
$$f(Z) \sim \mathcal{N}(0, K(Z, Z)), \quad (3)$$

## ➤ The Posterior:

- Conditioning on the observations  $Y = \{y_i\}_{i=1}^n$  at  $Z$ , the posterior predictive distribution at any input  $\mathbf{z} \in Z$  is given by,

$$f(\mathbf{z}) \mid Z, Y \sim \mathcal{N}(\mu(\mathbf{z}), \Sigma(\mathbf{z})), \quad (4)$$

with  $\mu(\mathbf{z}) = K(\mathbf{z}, Z)[K(Z, Z) + \sigma^2 I]^{-1}Y$  and  $\Sigma(\mathbf{z}) = K(\mathbf{z}, \mathbf{z}) - K(\mathbf{z}, Z)[K(Z, Z) + \sigma^2 I]^{-1}K(Z, \mathbf{z})$ .



# Gaussian Process for Trajectory Prediction

## ➤ The Prior:

- Consider a function  $f(\mathbf{z})$  to be sampled from a **Gaussian Process (GP)** with kernel  $K(\mathbf{z}, \mathbf{z}')$  and let  $K(Z, Z) \in \mathbb{R}^{n \times n}$  with  $[K(Z, Z)]_{i,j} = K(\mathbf{z}_i, \mathbf{z}_j)$ .

$$f(Z) \sim \mathcal{N}(0, K(Z, Z)), \quad (3)$$

## ➤ The Posterior:

- Conditioning on the observations  $Y = \{y_i\}_{i=1}^n$  at  $Z$ , the posterior predictive distribution at any input  $\mathbf{z} \in Z$  is given by,

$$f(\mathbf{z}) | Z, Y \sim \mathcal{N}(\mu(\mathbf{z}), \Sigma(\mathbf{z})), \quad (4)$$

with  $\mu(\mathbf{z}) = K(\mathbf{z}, Z)[K(Z, Z) + \sigma^2 I]^{-1}Y$  and  $\Sigma(\mathbf{z}) = K(\mathbf{z}, \mathbf{z}) - K(\mathbf{z}, Z)[K(Z, Z) + \sigma^2 I]^{-1}K(Z, \mathbf{z})$ .

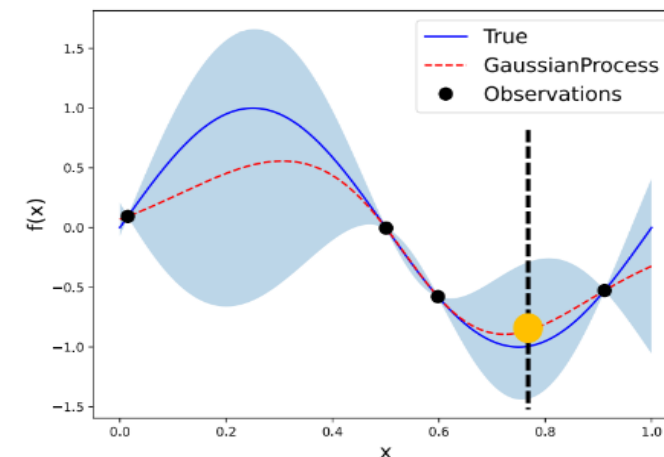
## ➤ Product Kernel:

- As a pair  $\mathbf{z} = (\mathbf{x}, t) \in \mathbb{X} \times \mathbb{T}$ , a kernel can be decomposed into two parts to capture the iterative learning characteristics

$$K(\mathbf{z}, \mathbf{z}') = \boxed{K_1(\mathbf{x}, \mathbf{x}')} \times \boxed{K_2(t, t')}$$

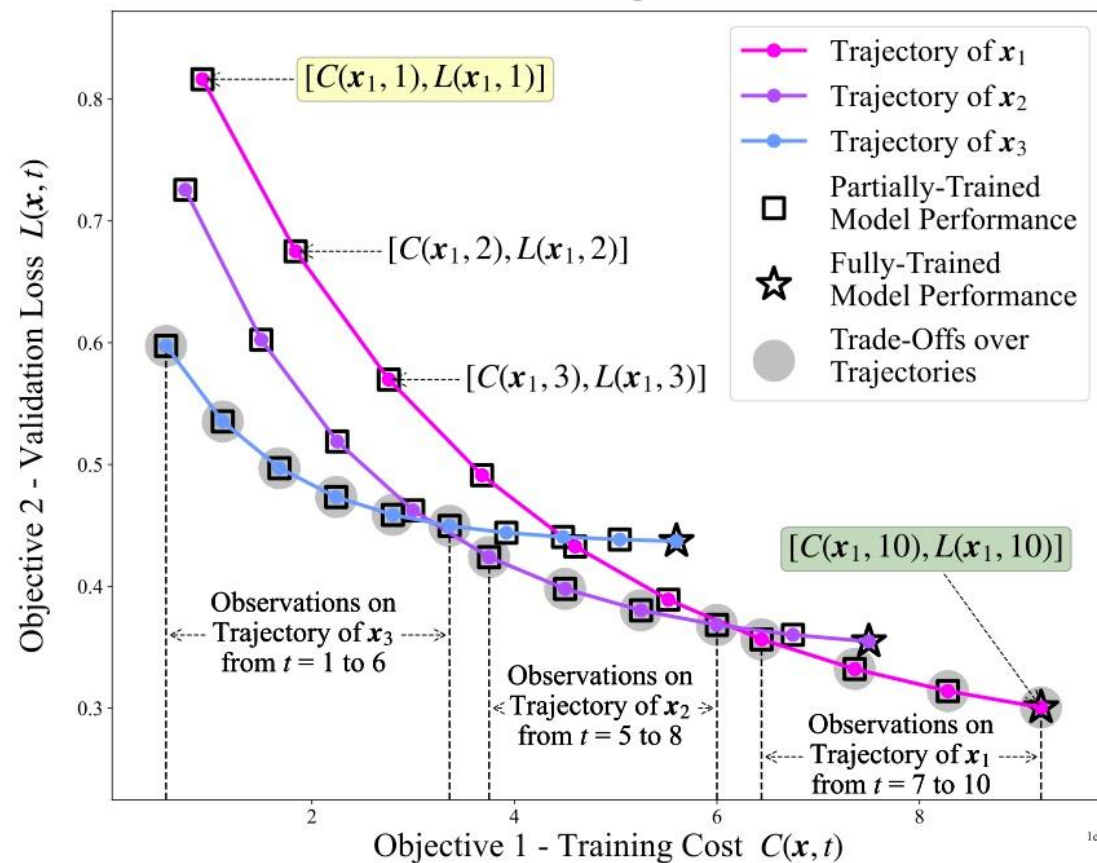
Kernel over hyperparameter setting  
e.g., Matérn kernel

Kernel over epoch  
e.g., Exponential decay or linear kernel



# Challenges in Solving EMOHPO

An Illustrative Example for EMOHPO



In the objective space of EMOHPO,

## 1. How to make prediction on trajectory?

The model should be able to capture the characteristics of the trajectory as the epoch changes.

## 2. How to sequentially determine next hyperparameter setting with trajectory prediction? (i.e., new $x'$ )

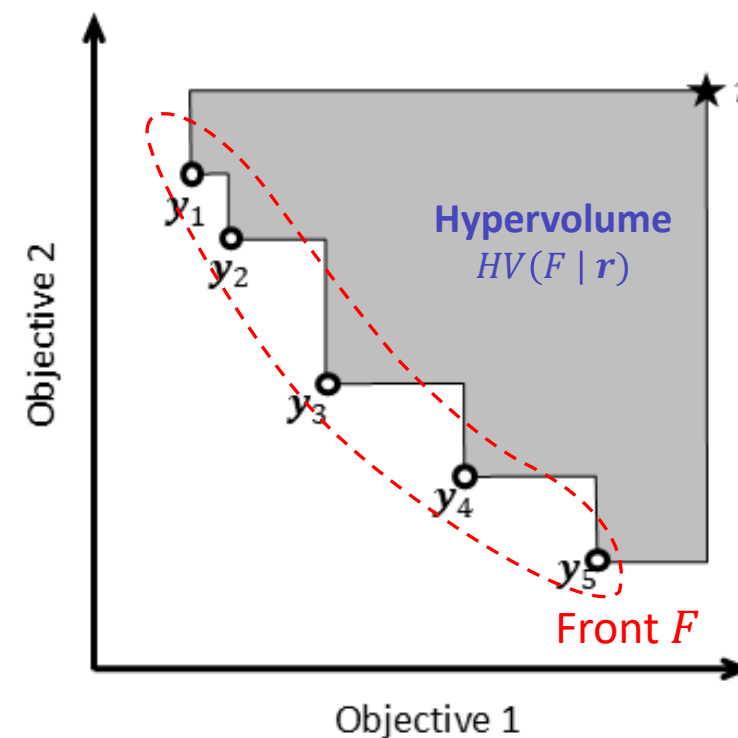
Assessing the quality of a setting requires consideration of its entire trajectory instead of a single position.

# Choose Setting $\mathbf{x}'$ - Trajectory-Based Acquisition Function

## Definition 1: Hypervolume Improvement (HVI)

Given a front  $F$  and a fixed point  $\mathbf{r}$ , the HVI of an objective vector  $\mathbf{y}'$  is the change in Hypervolume before and after including  $\mathbf{y}'$  into the front  $F$ , i.e.,

$$HVI(\mathbf{y}' | F, \mathbf{r}) = HV(F \cup \{\mathbf{y}'\} | \mathbf{r}) - HV(F | \mathbf{r}).$$

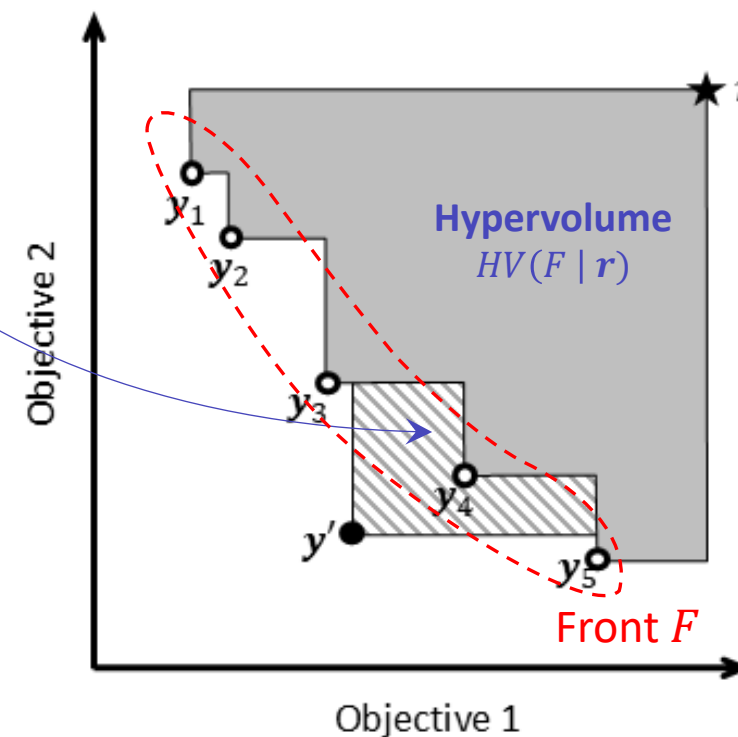


# Choose Setting $\mathbf{x}'$ - Trajectory-Based Acquisition Function

## Definition 1: Hypervolume Improvement (HVI)

Given a front  $F$  and a fixed point  $\mathbf{r}$ , the HVI of an objective vector  $\mathbf{y}'$  is the change in Hypervolume before and after including  $\mathbf{y}'$  into the front  $F$ , i.e.,

$$HVI(\mathbf{y}' | F, \mathbf{r}) = HV(F \cup \{\mathbf{y}'\} | \mathbf{r}) - HV(F | \mathbf{r}).$$



# Choose Setting $\mathbf{x}'$ - Trajectory-Based Acquisition Function

## Definition 1: Hypervolume Improvement (HVI)

Given a front  $F$  and a fixed point  $\mathbf{r}$ , the HVI of an objective vector  $\mathbf{y}'$  is the change in Hypervolume before and after including  $\mathbf{y}'$  into the front  $F$ , i.e.,

$$HVI(\mathbf{y}' | F, \mathbf{r}) = HV(F \cup \{\mathbf{y}'\} | \mathbf{r}) - HV(F | \mathbf{r}).$$

## Definition 2: Trajectory

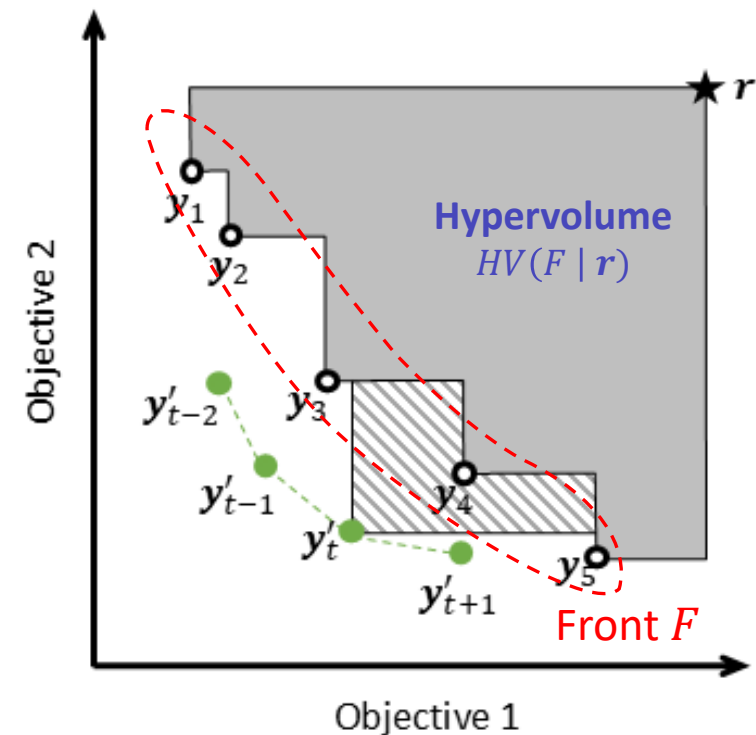
The trajectory of a hyperparameter setting  $\mathbf{x}$  is defined as the collection of all model performances observed during the entire training with  $\mathbf{x}$ , i.e.,

$$Trj(\mathbf{x}) := \{\mathbf{f}(\mathbf{x}, t)\}_{t=1}^{t_{max}} = \{f_1(\mathbf{x}, t), \dots, f_k(\mathbf{x}, t)\}_{t=1}^{t_{max}}.$$

## Definition 3: Trajectory-Based Expected HVI (TEHVI)

TEHVI estimates the gain of an out-of-sample setting  $\mathbf{x}$  by taking the expectation of HVI over the predictive distribution of its trajectory,

$$TEHVI(\mathbf{x} | F, \mathbf{r}) := \mathbb{E}[HVI(Trj(\mathbf{x}) | F, \mathbf{r})] = \mathbb{E}[HVI(\{\mathbf{f}(\mathbf{x}, t)\}_{t=1}^{t_{max}} | F, \mathbf{r})].$$



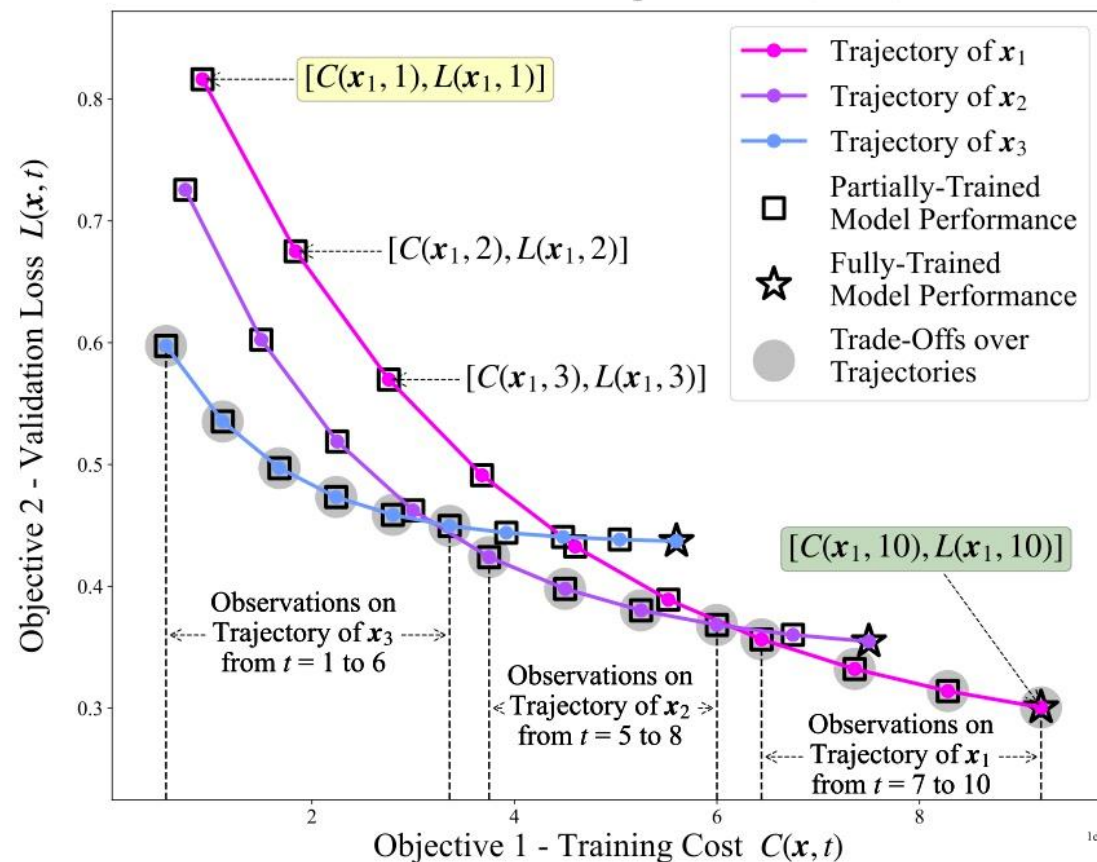
$$\approx \operatorname{argmax}_{\mathbf{x} \in \mathbb{X}} \frac{1}{M} \sum_{m=1}^M HVI(\widehat{Trj}_m(\mathbf{x}) | F, \mathbf{r})$$

by Monte Carlo integration



# Challenges in Solving EMOHPO

An Illustrative Example for EMOHPO



In the objective space of EMOHPO,

## 1. How to make prediction on trajectory?

The model should be able to capture the characteristics of the trajectory as the epoch changes.

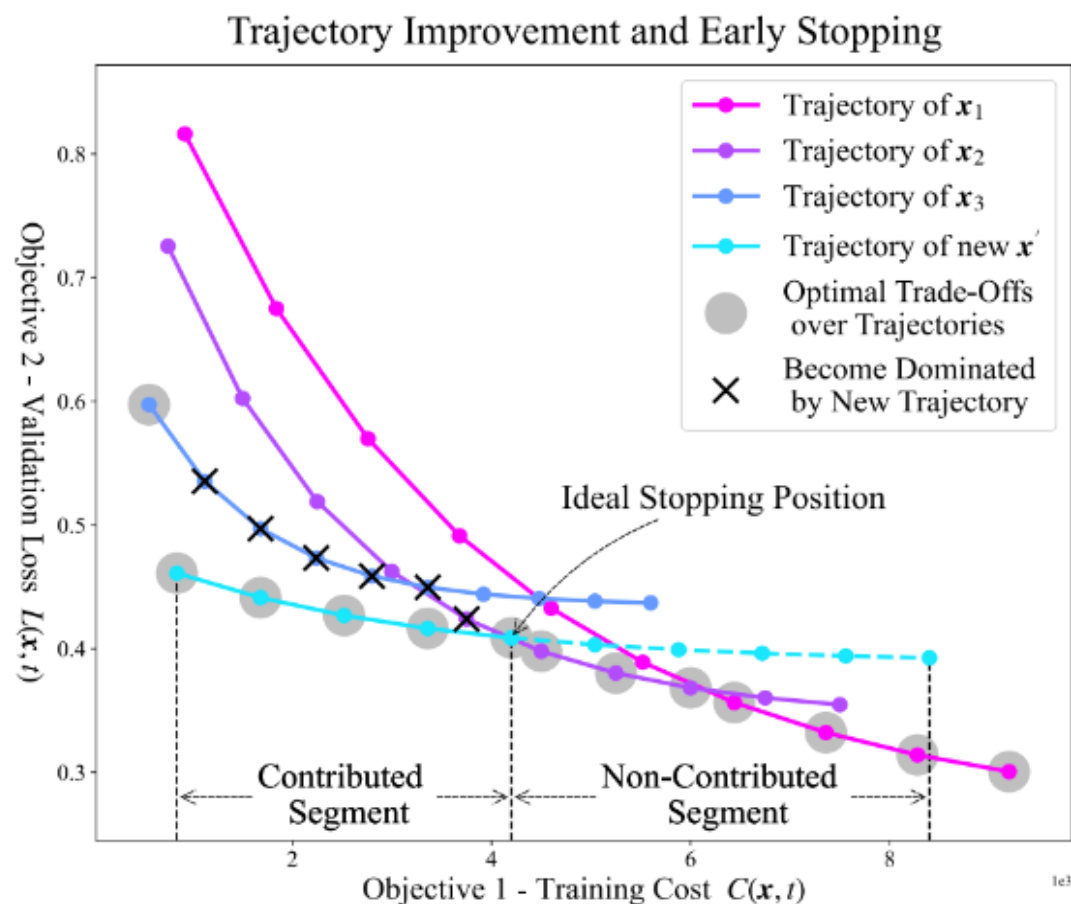
## 2. How to sequentially determine next hyperparameter setting with trajectory prediction? (i.e., new $x'$ )

Assessing the quality of a setting requires consideration of its entire trajectory instead of a single position.

## 3. How to execute early stopping without compromising optimization results? (i.e., new $t'$ )

A training procedure should only be stopped after as many trade-offs as possible have been observed along trajectory.

# Choose Epoch $t'$ - Trajectory-Based Early Stopping

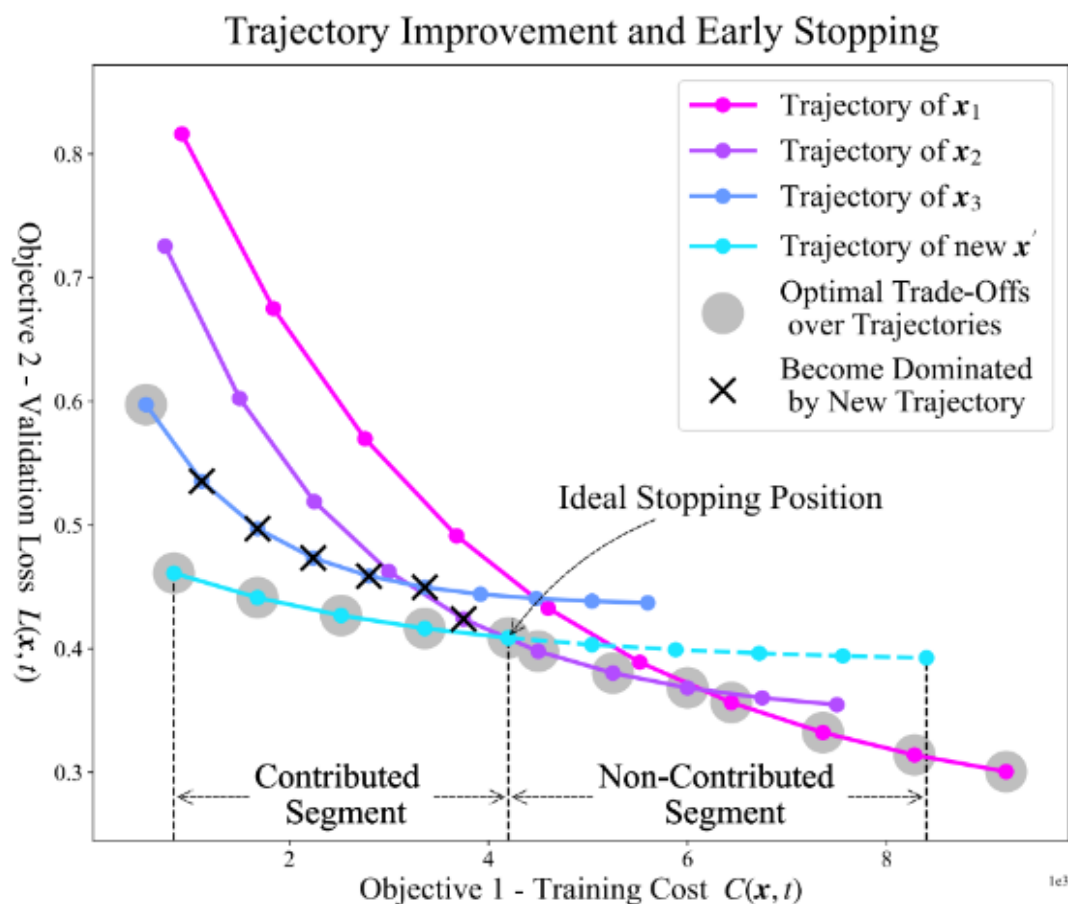


## ➤ Conservative stopping epoch

$$t^* = \sup\{t \in \mathbb{T} \mid \mu(x', t) - \beta \Sigma(x', t) < y, \exists y \in F\}$$

- $\mu(x', t) - \beta \Sigma(x', t)$  denotes the lower bound of the performance at  $t$ , with  $\beta$  controls the confidence level.
- Intuitively,  $t^*$  is the number of epochs after which future training results are unlikely to improve front  $F$ .

# Choose Epoch $t'$ - Trajectory-Based Early Stopping



## ➤ Conservative stopping epoch

$$t^* = \sup\{t \in \mathbb{T} \mid \mu(x', t) - \beta \Sigma(x', t) < y, \exists y \in F\}$$

- $\mu(x', t) - \beta \Sigma(x', t)$  denotes the lower bound of the performance at  $t$ , with  $\beta$  controls the confidence level.
- Intuitively,  $t^*$  is the number of epochs after which future training results are unlikely to improve front  $F$ .

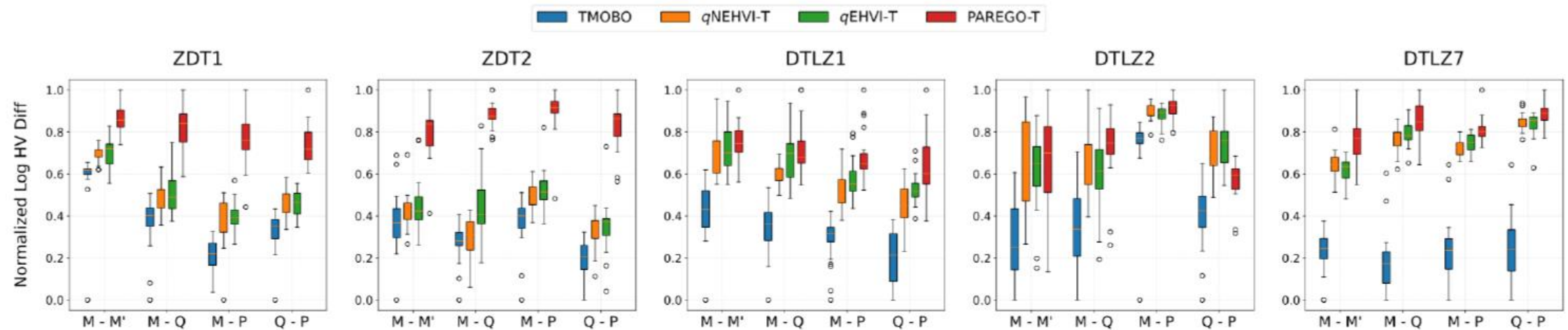
## ➤ Early stopping strategy ( $t'$ increases from 1 to $t_{max}$ ):

- if  $t' \leq t^*$ , continue training with  $x'$  for one epoch and let  $t' = t' + 1$  and recompute  $t^*$ ;
- if  $t' > t^*$ , terminate the training with  $x'$  immediately.

# Outline

1. Introduction
2. Problem and Methodology
  - 2.1 Enhanced Multi-Objective Hyperparameter Tuning
  - 2.2 Trajectory-Based Bayesian Optimization Approach
3. Numerical Experiments
4. Conclusions

# Results on Synthetic Simulations



- TMOBO consistently achieves the lowest HV difference over  $5 \times 4$  synthetic problems, which are modeled by

$$\min_{(x,t) \in \mathbb{X} \times \mathbb{T}} \mathbf{f}(x, t) = [f_1(x) \cdot g_1(t), \dots, f_k(x) \cdot g_k(t)].$$

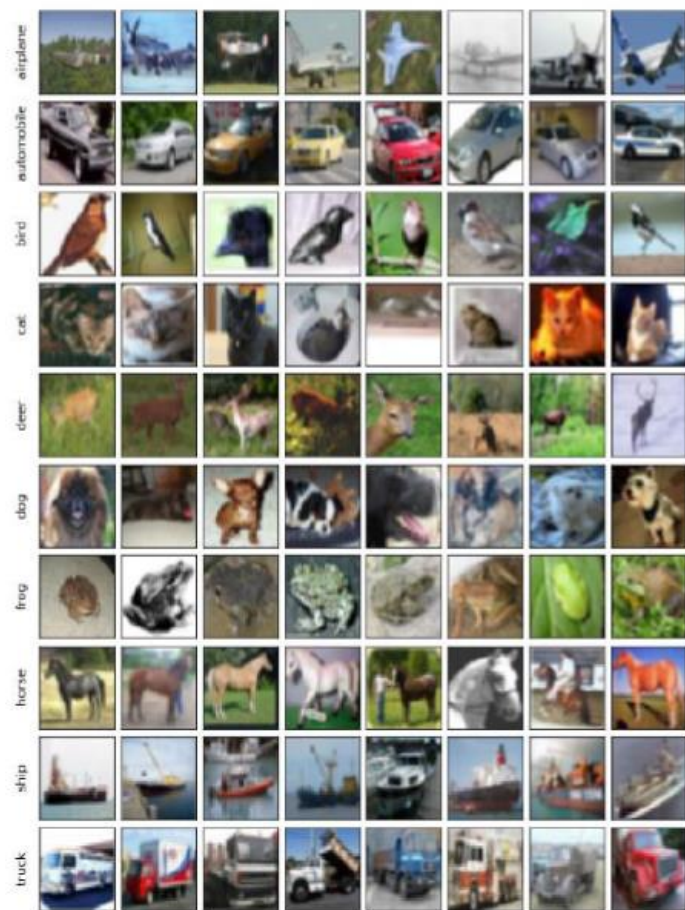
MO benchmark's functions

Functions to simulate iterative learning

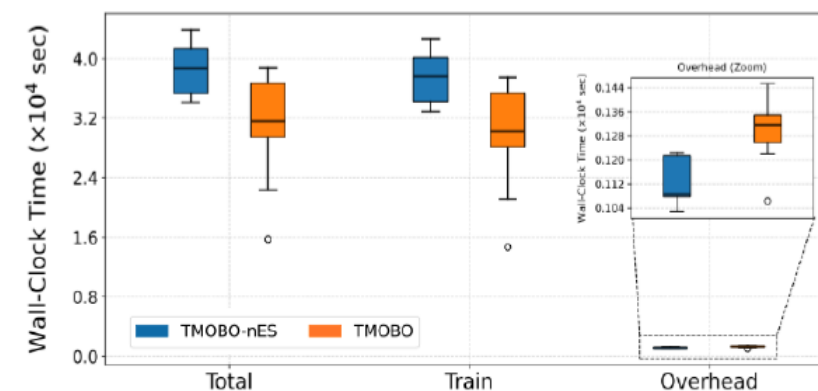
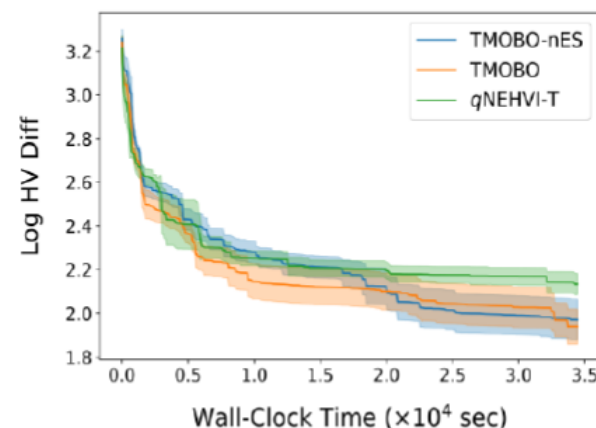
- Over 20 independent trials, the solutions obtained by TMOBO generally dominate a large proportion of those obtained by three alternative enhanced multi-objective optimizers.

□ E.g.,  $q$ NEHVI-T denotes the enhanced  $q$ NEHVI<sup>[3]</sup> by collecting all trajectory observations, similarly for  $q$ EHVI-T and ParEGO-T.

# Results on Real-World Benchmarks



- Tuning a more complex MobileNetV2 model<sup>[11]</sup> with iterative learning on CIFAR-10 image datasets:



- ❑ [Left] TMOBO and its variant TMOBO-nES outperform *q*NEHVI-T, with TMOBO demonstrating faster early convergence.
- ❑ [Right] TMOBO reduces model training time more than TMOBO-nES, though it incurs slightly higher (but negligible) computation overhead.

# Outline

1. Introduction
2. Problem and Methodology
  - 2.1 Enhanced Multi-Objective Hyperparameter Tuning
  - 2.2 Trajectory-Based Bayesian Optimization Approach
3. Numerical Experiments
4. Conclusions



# Conclusions

- Considering MOHPO with iterative learning, our interest centers on (1) how trajectory information affects the distribution of trade-offs and (2) how to leverage this information to search trade-offs.

**Problem Definition:** Introduce EMOHPO problem by including the number of training epoch as an explicit decision variable to reveal the trade-offs that may occur along trajectories.

**Methodology:** Propose TMOBO method that iteratively samples setting based on trajectory-based contribution and decides when to stop training based on trajectory predictions.

**Numerical Study:** Show the advantage of TMOBO over alternative methods in locating trade-offs for EMOHPO through synthetic and real-world benchmarks.

- For future research of this study
  - ❑ Development of the analytical form or more efficient approximation for the computation of TEHVI.
  - ❑ Application of scalable GP and more advanced data augmentation strategy in large-scale applications.
  - ❑ Extend the EMOHPO framework to other scenarios, such as drug design and material engineering where an iterative procedure typically exists.



# References

1. Knowles, J. (2006). ParEGO: A hybrid algorithm with on-line landscape approximation for expensive multiobjective optimization problems. *IEEE transactions on evolutionary computation*, 10(1), 50-66.
2. Daulton, S., Balandat, M., & Bakshy, E. (2020). Differentiable expected hypervolume improvement for parallel multi-objective Bayesian optimization. *Advances in Neural Information Processing Systems*, 33, 9851-9864.
3. Daulton, S., Balandat, M., & Bakshy, E. (2021). Parallel bayesian optimization of multiple noisy objectives with expected hypervolume improvement. *Advances in Neural Information Processing Systems*, 34, 2187-2200.
4. Swersky, K., Snoek, J., & Adams, R. P. (2014). Freeze-thaw Bayesian optimization. *arXiv preprint arXiv:1406.3896*.
5. Falkner, S., Klein, A., & Hutter, F. (2018, July). BOHB: Robust and efficient hyperparameter optimization at scale. In *International conference on machine learning* (pp. 1437-1446). PMLR.
6. Nguyen, V., Schulze, S., & Osborne, M. (2020). Bayesian optimization for iterative learning. *Advances in Neural Information Processing Systems*, 33, 9361-9371.
7. Dai, Z., Yu, H., Low, B. K. H., & Jaillet, P. (2019, May). Bayesian optimization meets Bayesian optimal stopping. In *International conference on machine learning* (pp. 1496-1506). PMLR.
8. Belakaria, S., Deshwal, A., & Doppa, J. R. (2020, April). Multi-fidelity multi-objective Bayesian optimization: An output space entropy search approach. In *Proceedings of the AAAI Conference on artificial intelligence* (Vol. 34, No. 06, pp. 10035-10043).
9. Schmucker, R., Donini, M., Zafar, M. B., Salinas, D., & Archambeau, C. (2021). Multi-objective asynchronous successive halving. *arXiv preprint arXiv:2106.12639*.
10. Zimmer, L., Lindauer, M., & Hutter, F. (2021). Auto-pytorch: Multi-fidelity metalearning for efficient and robust autodl. *IEEE transactions on pattern analysis and machine intelligence*, 43(9), 3079-3090.
11. Sandler, M., Howard, A., Zhu, M., Zhmoginov, A., & Chen, L. C. (2018). Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 4510-4520).

# Thanks for Your Attention!

## Q & A

Wenyu Wang<sup>1</sup>, Zheyi Fan<sup>2,3</sup>, Szu Hui Ng<sup>1,\*</sup>

<sup>1</sup>Industrial Systems Engineering and Management, National University of Singapore, Singapore

<sup>2</sup>Academy of Mathematics and System Science, Chinese Academy of Sciences, China

<sup>3</sup>School of Mathematical Sciences, University of Chinese Academy of Sciences, China