

# Faster & More Reliable Tuning of Neural Networks:

## *Bayesian Optimization with Importance Sampling*

Setareh Ariaifar, Zelda Mariet, Ehsan Elhamifar, Dana Brooks, Jennifer Dy,  
Jasper Snoek

# Motivation

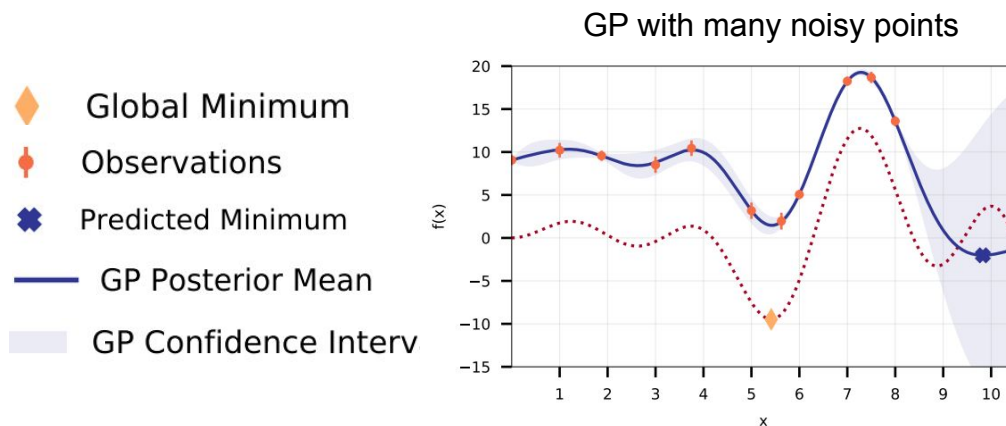
- Training neural nets  $\longrightarrow$  **expensive**
- Bayesian Optimization (BO)  $\longrightarrow$  **limited hyperparameters**
- **Low-fidelity** observations

## Pros

- **Increased #** of explored hyperparameters via:
  - Cheap partially trained models
  - Extrapolate to fully trained models

## Cons

- Adds to the **randomness/noise** of BO
- Challenging extrapolation



# Proposed Solution

- **Decrease randomness** by using Information of each training example
- BO + Importance Sampling (IS)

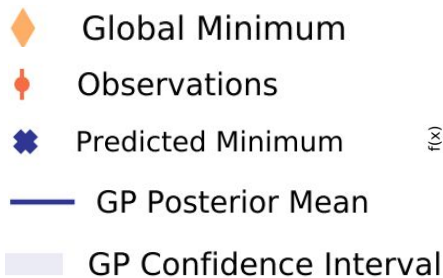
## Pros

- **High-fidelity** observation
- More accurate models
- **Less # of** observations required

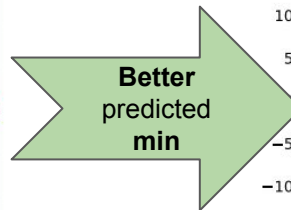
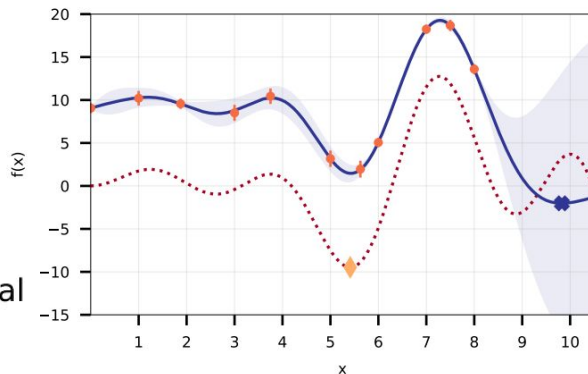
## Cons

- Large **overhead cost** challenging

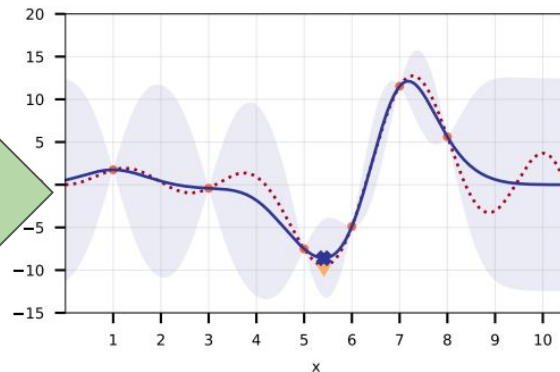
Solve via **multi-task BO over importance sampling design**  
Learn when high-fidelity is worth the cost



GP with many noisy points



GP with few noiseless points



# Proposed Solution

- Importance dist  $\longrightarrow$  **Initially similar to uniform sampling & expensive**

IS-SGD [1]

Start from uniform sampling  
Track variance reduction  
Switch to IS if variance reduction large

Select a random super-batch of size **B**  
Select mini-batches with IS from super-batch

[1] Katharopoulos &amp; Fleuret., ICML 2018

To learn the **trade-off** parameter **B**  $\longrightarrow$  Maximize  $\alpha_n(x, B)$

$$\alpha_n(x, B) = \frac{1}{\mu(c_n(x | B))} \left[ H(\mathbb{P}[x^* | B = |\mathbf{D}|, \mathcal{D}_n]) - \mathbb{E}_y [H(\mathbb{P}[x^* | B = |\mathbf{D}|, \mathcal{D}_n \cup \{x, B, y\}])] \right],$$

Expected training cost  
for  $x, \mathbf{B}$

Expected entropy reduction from training on hyperparameter  $x$  via  
IS-SGD routine with super-batch size **B**

# Results- ResNet on CIFAR10

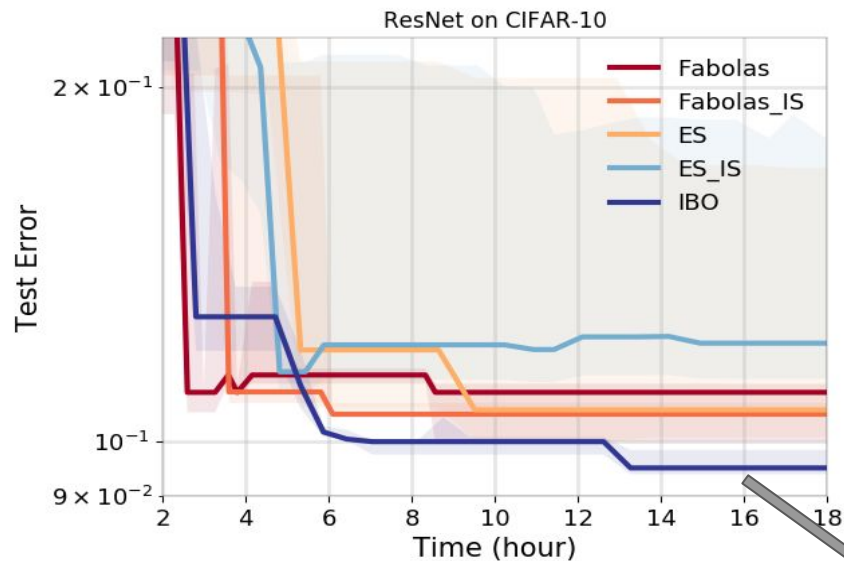


Google AI

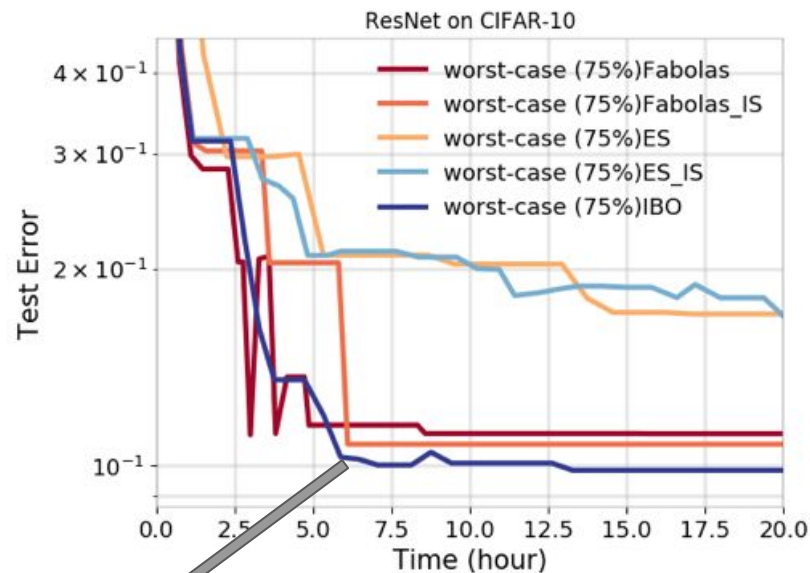


Northeastern University

- Improved worst-case performance



Average Performance



Worst-case Performance

Ours

# Results- ResNet on CIFAR100

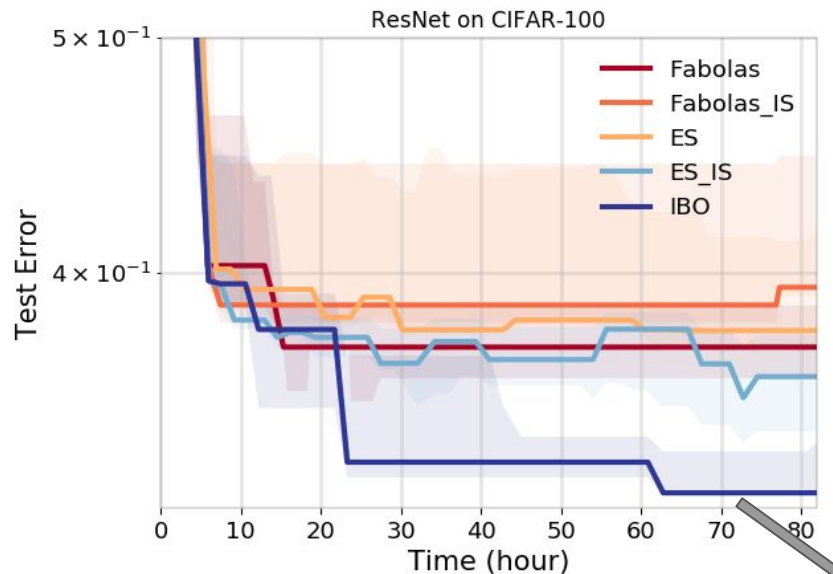


Google AI

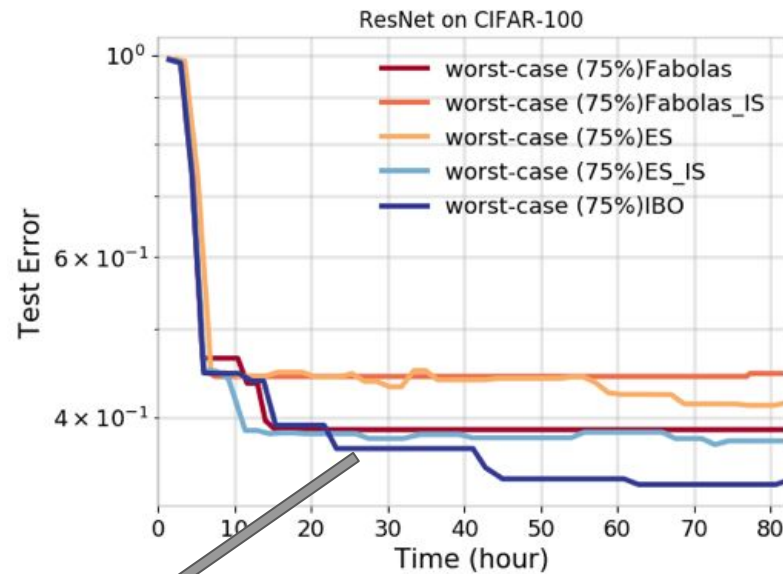


Northeastern University

- Improved worst-case performance



Average Performance



Worst-case Performance

Ours