

# Priors in probabilistic numerics

Zi Wang

ProbNum @ Dagstuhl

Google Research



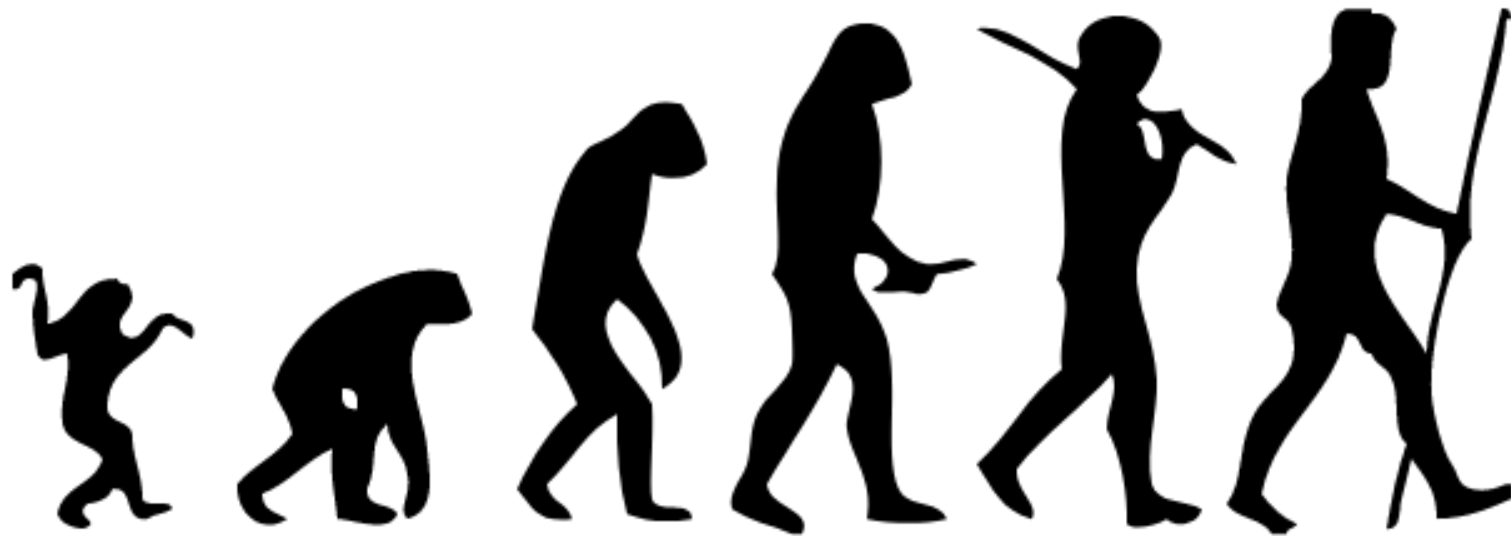
Prior, engineering prior and data prior

Where does the prior  
come from?

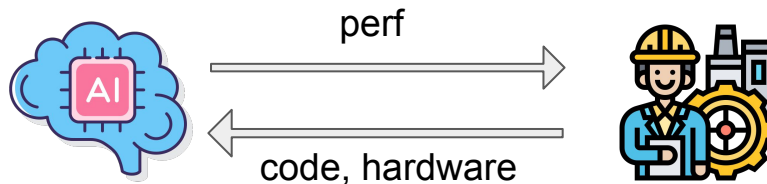
# Priors in nature



# Priors in nature through evolution



# Priors in software: engineering prior



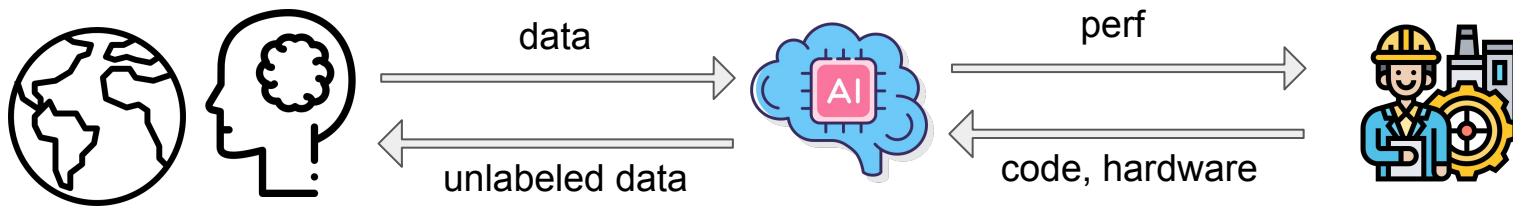
- engineering prior  $\approx$  (a sample of) the posterior of engineers
- posterior of engineers builds on human knowledge as a whole (education, books, journals, blogs + engineer experience with models)

# Example of engineering prior: A\*



- Very generalizable;
- Wide applications in robot planning and path planning in games;
- Algorithm entirely built-in by expert knowledge and abstraction of how humans solve path planning problem.

# Software: engineering prior + data



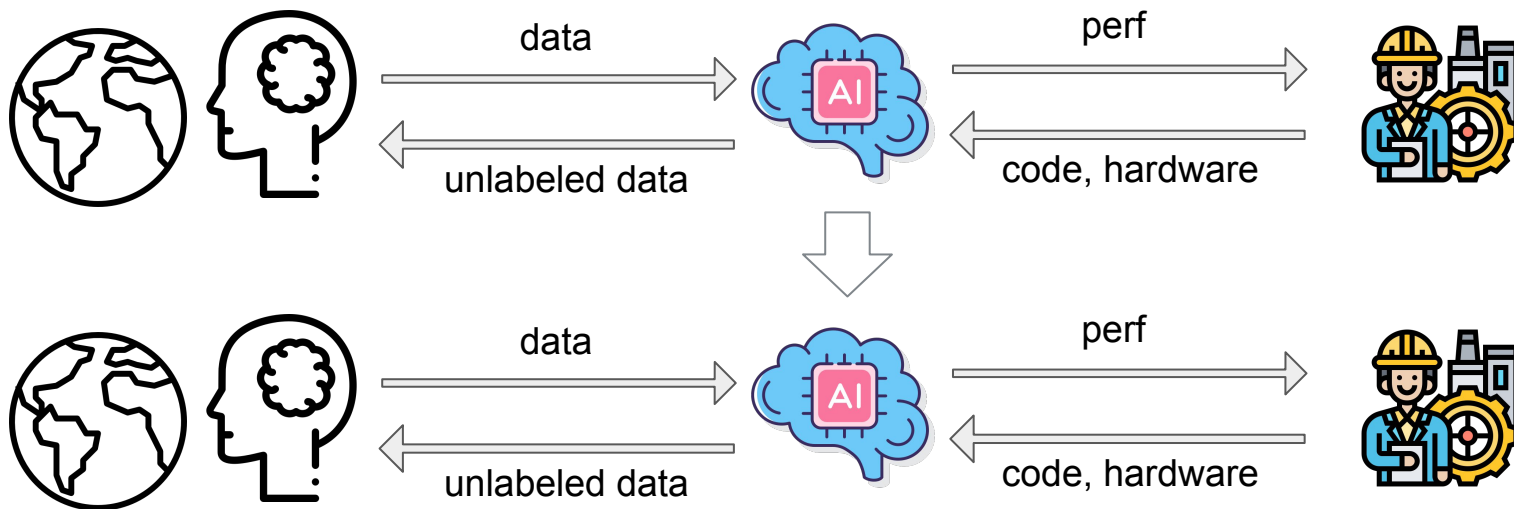
- engineering prior  $\approx$  (a sample of) the posterior of engineers
- data = a set of data points returned by the world or human annotators

# Example: robot learning with strong engineering priors





# Data prior: priors in the evolution of software



- data prior  $\approx$  pretraining

# Obtaining a data prior through meta learning

Task 1	$(x_{11}, y_{11})$	$(x_{12}, y_{12})$	.....	$(x_{1M}, y_{1M})$
Task 2	$(x_{21}, y_{21})$	$(x_{22}, y_{22})$	.....	$(x_{2M}, y_{2M})$
.....	.....	.....	.....	.....
Task N	$(x_{N1}, y_{N1})$	$(x_{N2}, y_{N2})$	.....	$(x_{NM}, y_{NM})$
New Task	?	?	.....	?

- The software to solve a new task corresponds to a new generation of models.
- A data prior can be obtained through all previous generations.

# Data prior example on BayesOpt

“HyperBO does not assume the knowledge of any GP parameters; instead, we learn the GP mean function, kernel function, and possible observation noise from data in the form of point sets, i.e. *i.i.d.* sets of correlated points.”

— Automatic prior selection for meta Bayesian optimization with a case study on tuning deep neural network optimizers.

Joint work with G. E. Dahl, K. Swersky, C. Lee, Z. Mariet, Z. Nado, J. Gilmer, J. Snoek, Z. Ghahramani. <https://arxiv.org/abs/2109.08215>

Assumption: evaluation functions on hyperparameters for all tasks are *i.i.d.* function samples from a GP

Task f_1	(x_11, y_11)	(x_12, y_12)	.....	(x_1M, y_1M)
.....	.....	.....	.....	.....
Task f_i	(x_i1, y_i1)	(x_i2, y_i2)	.....	(x_iM, y_iM)
.....	.....	.....	.....	.....
Task f_N	(x_N1, y_N1)	(x_N2, y_N2)	.....	(x_NM, y_NM)
New Task	?	?	.....	?

$D_{f_i}$

$D_N = \{D_{f_i}\}_{i=1}^N$

$$f_i \sim \mathcal{GP}(\mu, k)$$

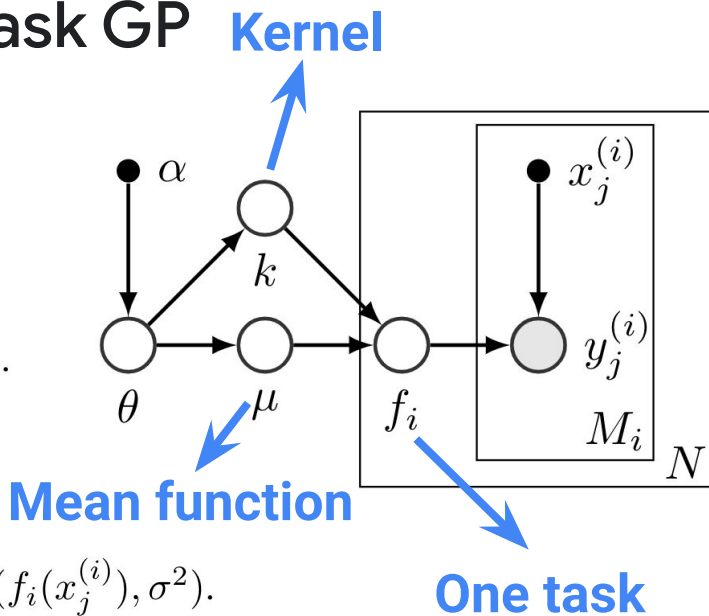
$$f \sim \mathcal{GP}(\mu, k)$$

$$\text{solve } \max_{x \in \mathfrak{X}} f(x)$$

# Hierarchical Gaussian process: a different viewpoint of multi-task GP

All tasks are IID samples from a GP

- Draw parameter  $\theta$  from  $p(\theta; \alpha)$ .
- Draw mean function  $\mu$  and kernel function  $k$  from  $p(\mu, k \mid \theta)$ .
- For each outer iteration  $i$  from 1 to  $N$ ,
  - Draw a function  $f_i$  from  $\mathcal{GP}(\mu, k)$ .
  - For each inner loop iteration from 1 to  $M_i$ ,
    - \* Given input  $x_j^{(i)}$ , we draw the observation  $y_j^{(i)} \sim \mathcal{N}(f_i(x_j^{(i)}), \sigma^2)$ .



Instead of learning correlations among tasks,  
we learn the GP that generated all tasks

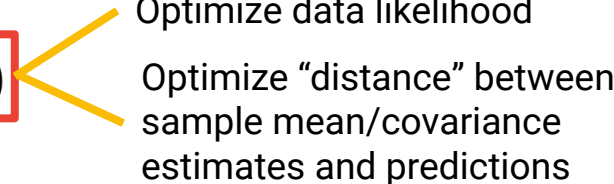
# HyperBO: a practical meta Bayesian optimization method

---

**Algorithm 1** HyperBO with acquisition function  $\alpha(\cdot)$ .

---

```
1: function HYPERBO ( $f, D_N$ )
2:    $\mathcal{GP}(\hat{\mu}, \hat{k}) \leftarrow \text{TRAINGP}(D_N)$ 
3:    $D_f \leftarrow \emptyset$ 
4:   for  $t = 1, \dots, T$  do
5:      $x_t \leftarrow \arg \max_{x \in \mathcal{X}} \alpha \left( x; \mathcal{GP}(\hat{\mu}, \hat{k} \mid D_f) \right)$ 
6:      $y_t \leftarrow \text{OBSERVE}(f(x_t))$ 
7:      $D_f \leftarrow D_f \cup \{(x_t, y_t)\}$ 
8:   end for
9:   return  $D_f$ 
10: end function
```



Optimize data likelihood

Optimize “distance” between sample mean/covariance estimates and predictions

# Optimize data likelihood

$$\begin{aligned}\log p(D_N \mid \mu, k, \sigma^2) &= \sum_{i=1}^N \log p(D_{f_i} \mid \mu, k, \sigma^2) \\ &= \sum_{i=1}^N \left( -\frac{1}{2} \bar{\mathbf{y}}_{(i)}^\top K^{-1} \bar{\mathbf{y}}_{(i)} - \frac{1}{2} \log |K| - \frac{M_i}{2} \log 2\pi \right)\end{aligned}$$

$$\bar{\mathbf{y}}_{(i)} = \mathbf{y}^{(i)} - \mu(\mathbf{x}^{(i)}), K = k(\mathbf{x}^{(i)}) + \sigma^2 \mathbf{I}, \mathbf{x}^{(i)} = [x_j^{(i)}]_{j=1}^{M_i} \text{ and } \mathbf{y}^{(i)} = [y_j^{(i)}]_{j=1}^{M_i}$$

Also possible to search over different mean/kernel architectures

# Optimize distance between sample mean/covariance estimates and predictions (for same inputs across tasks)

Task 1	(x <sub>1</sub> , y <sub>11</sub> )	(x <sub>2</sub> , y <sub>12</sub> )	.....	(x <sub>M</sub> , y <sub>1M</sub> )
Task 2	(x <sub>1</sub> , y <sub>21</sub> )	(x <sub>2</sub> , y <sub>22</sub> )	.....	(x <sub>M</sub> , y <sub>2M</sub> )
.....	.....	.....	.....	.....
Task N	(x <sub>1</sub> , y <sub>N1</sub> )	(x <sub>2</sub> , y <sub>N2</sub> )	.....	(x <sub>M</sub> , y <sub>NM</sub> )
New Task	?	?	.....	?

$$\mathbf{y} = [\mathbf{y}_j]_{j=1}^M \in \mathbb{R}^{M \times N}$$

$$\hat{\boldsymbol{\mu}} = \frac{1}{N} \mathbf{y} \mathbf{1}_N \in \mathbb{R}^M$$

$$\hat{K} = \frac{1}{N} (\mathbf{y} - \hat{\boldsymbol{\mu}} \mathbf{1}_N^\top) (\mathbf{y} - \hat{\boldsymbol{\mu}} \mathbf{1}_N^\top)^\top \in \mathbb{R}^{M \times M}$$

$$\mathcal{D}_{\text{KL}}(\hat{\boldsymbol{\mu}}, \hat{K}, \mu(\mathbf{x}), k(\mathbf{x}) + \mathbf{I}\sigma^2) = \frac{1}{2} \left( \text{tr}(K^{-1} \hat{K}) + (\mu(\mathbf{x}) - \hat{\boldsymbol{\mu}})^\top K^{-1} (\mu(\mathbf{x}) - \hat{\boldsymbol{\mu}}) + \ln \frac{|K|}{|\hat{K}|} - M \right)$$



# Preliminary Results

Reduced complexity in #Tasks

$$O(M^3 N)$$

## Theoretical results: near-zero regret with unknown GP priors

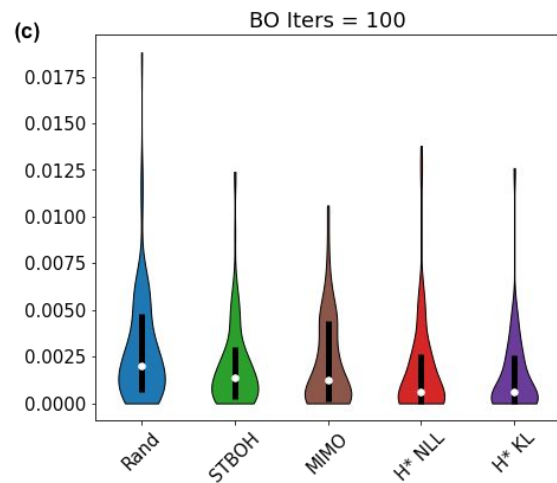
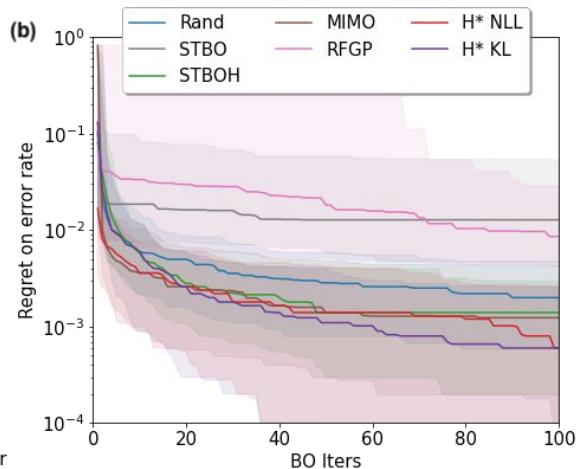
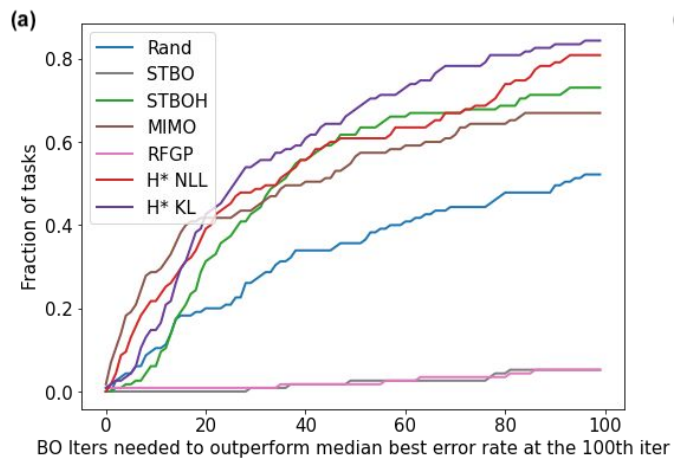
**Proposition 1.** *For any  $M, d, N \in \mathbb{Z}^+$ ,  $\mathbf{x} \in \mathbb{R}^{M \times d}$ ,  $\boldsymbol{\mu} \in \mathbb{R}^M$ ,  $V \in \mathbb{R}^{N \times M}$  and  $K = V^\top V$ , there exists a Gaussian process  $\mathcal{GP}(\hat{\mu}, \hat{k})$  such that  $\mathcal{D}_{EUC}(\hat{\mu}(\mathbf{x}), \hat{k}(\mathbf{x}), \boldsymbol{\mu}, K) \equiv 0$ .*

**Theorem 2.** *Let  $N \geq 4 \log \frac{6}{\delta} + T + 2$ . With probability at least  $1 - \delta$ , simple regret in  $T$  iterations of HyperBO with special cases of either GP-UCB or PI satisfies*

$$R_T < O \left( \sqrt{\frac{1}{N - T}} + \left( \log \frac{1}{\delta} \right)^{\frac{1}{2}} \right) O(\rho_T / T + \sigma), \quad (1)$$

where  $\rho_T = \max_{A \subset \mathfrak{X}, |A|=T} \frac{1}{2} \log |\mathbf{I} + \sigma^{-2} k(A)|$ .

# Hyper BO achieves better empirical results on offline optimizer hyperparameter tuning tasks

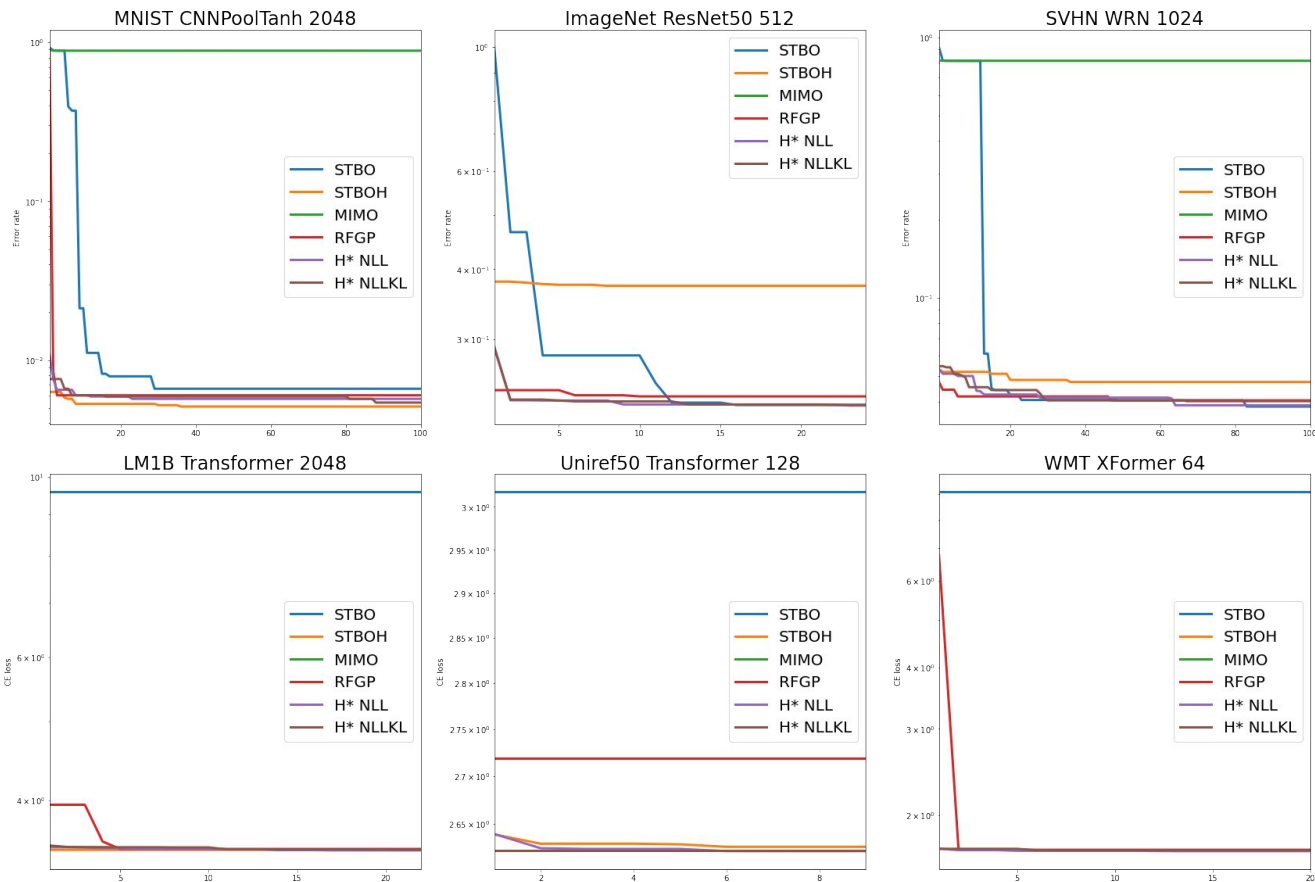


hyperparameter tuning dataset generated by [github.com/google/init2winit](https://github.com/google/init2winit)

# Best validation error rates for each task in the offline leave-one-out experiments (after 100 BO iterations)

					HyperBO w/ different objectives	
	Rand	STBOH	MIMO	MAF	H* NLL	H* KL
WMT XFormer 64	34.27 $\pm$ 0.16	34.15 $\pm$ 0.15	34.40 $\pm$ 0.13	34.09 $\pm$ 0.09	<b>33.91 <math>\pm</math> 0.01</b>	33.97 $\pm$ 0.02
Uniref50 Transformer 128	79.06 $\pm$ 0.04	78.92 $\pm$ 0.12	79.17 $\pm$ 0.13	79.34 $\pm$ 0.27	78.71 $\pm$ 0.06	<b>78.64 <math>\pm</math> 0.00</b>
LM1B Transformer 2048	61.96 $\pm$ 0.03	61.95 $\pm$ 0.04	61.96 $\pm$ 0.05	62.02 $\pm$ 0.10	<b>61.81 <math>\pm</math> 0.01</b>	<b>61.81 <math>\pm</math> 0.01</b>
SVHN WRN 1024	3.99 $\pm$ 0.04	4.05 $\pm$ 0.10	<b>3.83 <math>\pm</math> 0.04</b>	4.10 $\pm$ 0.09	4.10 $\pm$ 0.02	4.08 $\pm$ 0.01
SVHN WRN 256	3.71 $\pm$ 0.01	3.72 $\pm$ 0.02	<b>3.65 <math>\pm</math> 0.01</b>	3.69 $\pm$ 0.03	3.78 $\pm$ 0.01	3.72 $\pm$ 0.03
ImageNet ResNet50 256	23.03 $\pm$ 0.07	22.66 $\pm$ 0.07	22.73 $\pm$ 0.07	26.44 $\pm$ 1.98	<b>22.53 <math>\pm</math> 0.00</b>	22.58 $\pm$ 0.04
ImageNet ResNet50 512	23.02 $\pm$ 0.11	22.74 $\pm$ 0.05	23.01 $\pm$ 0.05	25.46 $\pm$ 1.41	<b>22.65 <math>\pm</math> 0.02</b>	22.79 $\pm$ 0.03
MNIST CNNPoolTanh 2048	0.55 $\pm$ 0.01	<b>0.53 <math>\pm</math> 0.01</b>	<b>0.53 <math>\pm</math> 0.01</b>	<b>0.52 <math>\pm</math> 0.01</b>	0.59 $\pm$ 0.02	0.54 $\pm$ 0.00
MNIST CNNPoolTanh 256	0.51 $\pm$ 0.01	0.48 $\pm$ 0.01	<b>0.47 <math>\pm</math> 0.00</b>	<b>0.47 <math>\pm</math> 0.01</b>	<b>0.46 <math>\pm</math> 0.01</b>	<b>0.47 <math>\pm</math> 0.01</b>
MNIST CNNPoolReLU 2048	0.69 $\pm$ 0.01	0.73 $\pm$ 0.02	0.67 $\pm$ 0.02	0.68 $\pm$ 0.01	<b>0.64 <math>\pm</math> 0.00</b>	0.70 $\pm$ 0.03
MNIST CNNPoolReLU 256	0.51 $\pm$ 0.01	0.55 $\pm$ 0.03	0.50 $\pm$ 0.01	0.51 $\pm$ 0.01	<b>0.49 <math>\pm</math> 0.00</b>	<b>0.49 <math>\pm</math> 0.00</b>
MNIST CNNReLU 2048	1.14 $\pm$ 0.03	1.20 $\pm$ 0.09	1.10 $\pm$ 0.01	1.17 $\pm$ 0.02	<b>1.06 <math>\pm</math> 0.00</b>	1.11 $\pm$ 0.02
MNIST CNNReLU 256	1.09 $\pm$ 0.02	1.06 $\pm$ 0.01	1.08 $\pm$ 0.02	1.07 $\pm$ 0.02	<b>1.03 <math>\pm</math> 0.00</b>	1.04 $\pm$ 0.01
Fashion CNNPoolTanh 2048	7.14 $\pm$ 0.06	7.10 $\pm$ 0.05	<b>7.01 <math>\pm</math> 0.04</b>	7.12 $\pm$ 0.04	<b>7.00 <math>\pm</math> 0.04</b>	<b>7.02 <math>\pm</math> 0.07</b>
Fashion CNNPoolTanh 256	6.51 $\pm$ 0.03	6.67 $\pm$ 0.18	6.40 $\pm$ 0.05	6.47 $\pm$ 0.03	6.40 $\pm$ 0.04	<b>6.34 <math>\pm</math> 0.04</b>
Fashion CNNPoolReLU 2048	<b>7.47 <math>\pm</math> 0.02</b>	<b>7.48 <math>\pm</math> 0.04</b>	7.54 $\pm$ 0.06	7.63 $\pm$ 0.04	<b>7.47 <math>\pm</math> 0.03</b>	<b>7.47 <math>\pm</math> 0.02</b>
Fashion CNNPoolReLU 256	6.78 $\pm$ 0.04	<b>6.74 <math>\pm</math> 0.01</b>	7.03 $\pm$ 0.07	6.84 $\pm$ 0.05	<b>6.74 <math>\pm</math> 0.03</b>	6.81 $\pm$ 0.05
Fashion CNNReLU 2048	7.70 $\pm$ 0.03	<b>7.47 <math>\pm</math> 0.09</b>	7.60 $\pm$ 0.04	40.40 $\pm$ 17.80	<b>7.54 <math>\pm</math> 0.01</b>	7.57 $\pm$ 0.02
Fashion CNNReLU 256	7.70 $\pm$ 0.04	7.46 $\pm$ 0.11	7.84 $\pm$ 0.06	24.13 $\pm$ 14.54	<b>7.29 <math>\pm</math> 0.05</b>	<b>7.25 <math>\pm</math> 0.05</b>
CIFAR100 WRN 2048	21.28 $\pm$ 0.27	<b>20.78 <math>\pm</math> 0.19</b>	21.75 $\pm$ 0.15	50.70 $\pm$ 15.44	21.22 $\pm$ 0.23	<b>20.82 <math>\pm</math> 0.19</b>
CIFAR100 WRN 256	19.17 $\pm$ 0.19	19.02 $\pm$ 0.03	19.12 $\pm$ 0.04	19.84 $\pm$ 0.13	<b>19.00 <math>\pm</math> 0.00</b>	19.04 $\pm$ 0.05
CIFAR10 WRN 2048	3.73 $\pm$ 0.05	<b>3.43 <math>\pm</math> 0.07</b>	<b>3.46 <math>\pm</math> 0.05</b>	<b>3.40 <math>\pm</math> 0.06</b>	3.55 $\pm$ 0.10	<b>3.43 <math>\pm</math> 0.05</b>
CIFAR10 WRN 256	2.84 $\pm$ 0.04	2.88 $\pm$ 0.06	2.89 $\pm$ 0.06	3.04 $\pm$ 0.05	2.82 $\pm$ 0.03	<b>2.74 <math>\pm</math> 0.01</b>

# Online experiment results: HyperBO is most stable



# Thank you!

zi-wang.com