

# Transfer Learning for Optimal Configuration of Big Data Software

*Pooyan Jamshidi, Giuliano Casale  
Imperial College London  
p.jamshidi@imperial.ac.uk*

*Department of Computing  
Research Associate Symposium  
14th June 2016*

# Motivation

1- Many different Parameters =>

- large state space
- interactions

2- Defaults are typically used =>

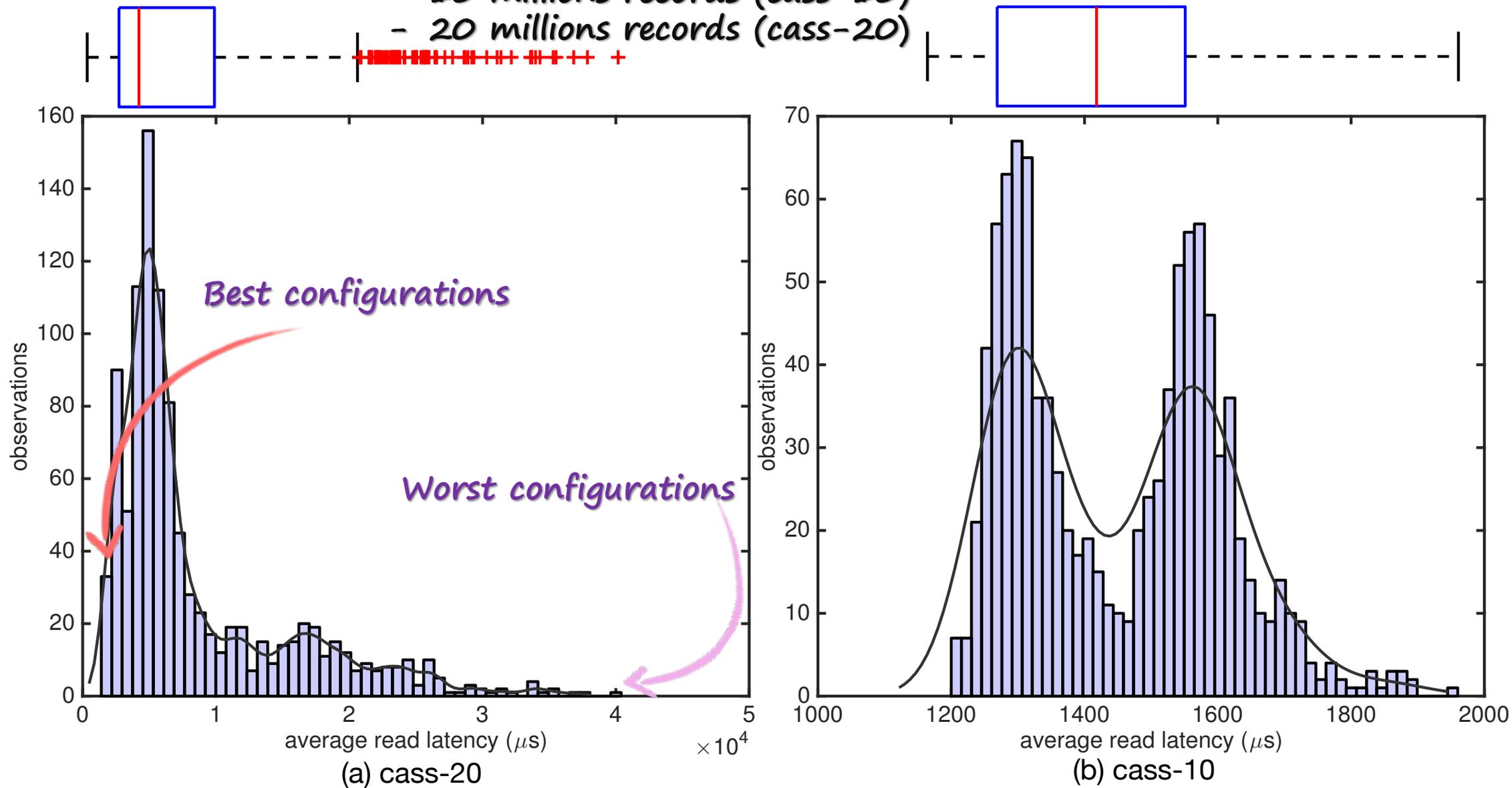
- poor performance

```
102
103 drpc.port: 3772
104 drpc.worker.threads: 64
105 drpc.max_buffer_size: 1048576
106 drpc.queue.size: 128
107 drpc.invocations.port: 3773
108 drpc.invocations.threads: 64
109 drpc.request.timeout.secs: 600
110 drpc.childopts: "-Xmx768m"
111 drpc.http.port: 3774
112 drpc.https.port: -1
113 drpc.https.keystore.password: ""
114 drpc.https.keystore.type: "JKS"
115 drpc.http.creds.plugin: org.apache.storm.security.auth.DefaultHttpCredentialsPlugin
116 drpc.authorizer.acl.filename: "drpc-auth-acl.yaml"
117 drpc.authorizer.acl.strict: false
118
119 transactional.zookeeper.root: "/transactional"
120 transactional.zookeeper.servers: null
121 transactional.zookeeper.port: null
122
123 ## blobstore configs
124 supervisor.blobstore.class: "org.apache.storm.blobstore.NimbusBlobStore"
125 supervisor.blobstore.download.thread.count: 5
126 supervisor.blobstore.download.max_retries: 3
127 supervisor.localizer.cache.target.size.mb: 10240
128 supervisor.localizer.cleanup.interval.ms: 600000
129
```

# Motivation

Experiments on  
Apache Cassandra:

- 6 parameters, 1024 configurations
- Average read latency
- 10 millions records (*cass-10*)
- 20 millions records (*cass-20*)



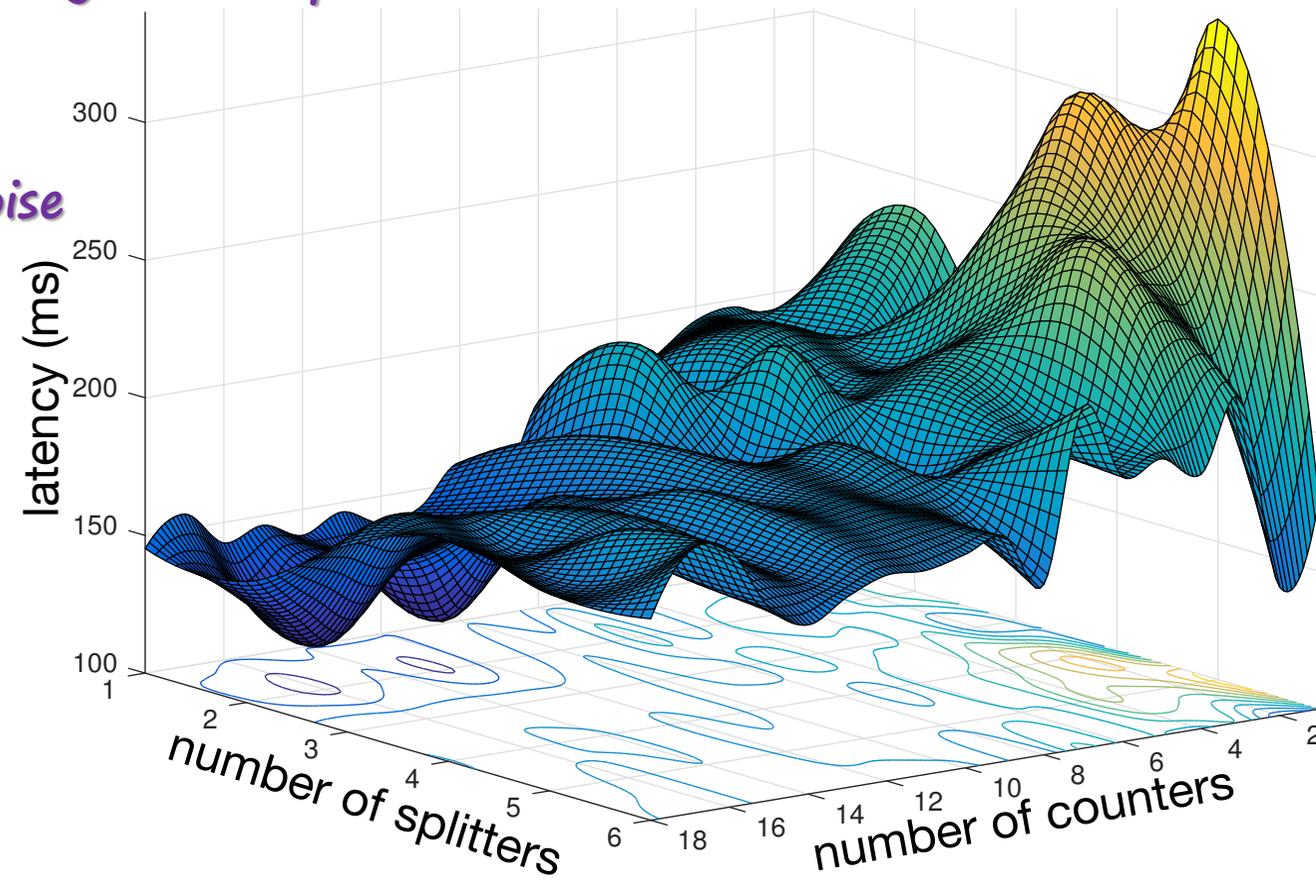
# Goal!

$$\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbb{X}} f(\mathbf{x})$$

$$\mathbb{X} = \text{Dom}(X_1) \times \cdots \times \text{Dom}(X_d) \quad \text{Configuration space}$$

$$y_i = f(\mathbf{x}_i), \mathbf{x}_i \subset \mathbb{X} \quad \text{Partially known}$$

$$y_i = f(\mathbf{x}_i) + \epsilon \quad \text{Measurements contain noise}$$



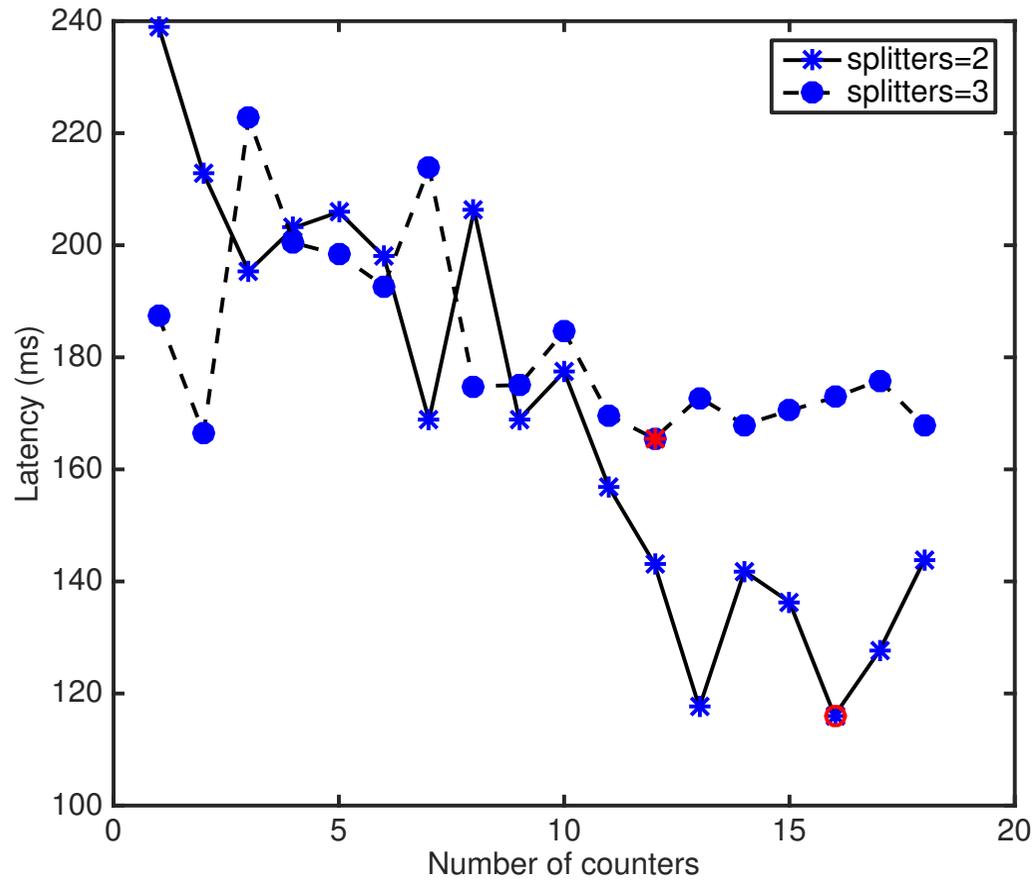
# Finding optimum configuration is difficult!

$$\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbb{X}} f(\mathbf{x})$$

$$\mathbb{X} = \text{Dom}(X_1) \times \dots \times \text{Dom}(X_d)$$

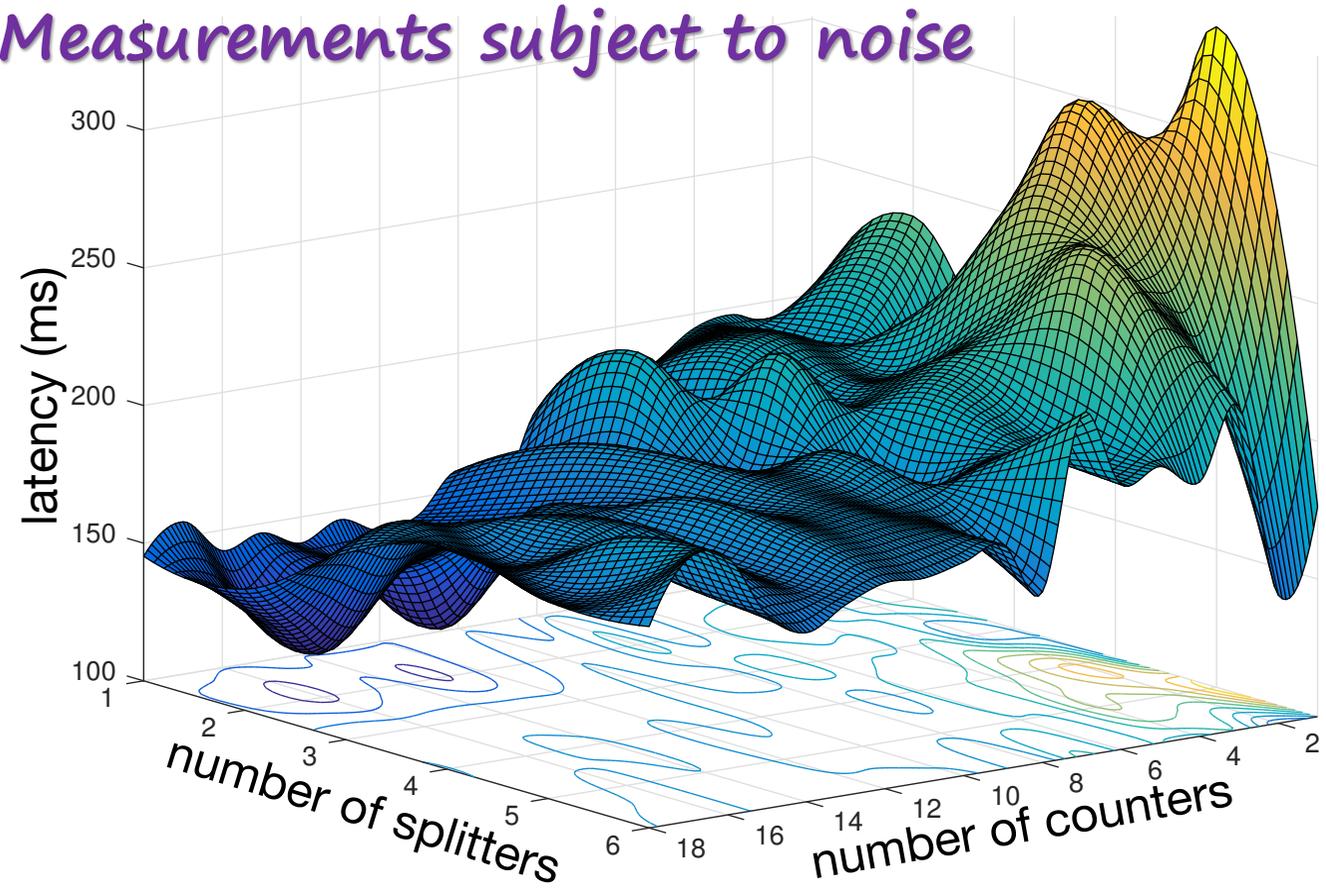
$$y_i = f(\mathbf{x}_i), \mathbf{x}_i \in \mathbb{X}$$

$$y_i = f(\mathbf{x}_i) + \epsilon$$



Response surface is:

- Non-linear
- Non convex
- Multi-modal
- Measurements subject to noise





Part

# Bayesian Optimization for Configuration Optimization (BO4CO)

Code: <https://github.com/dice-project/DICE-Configuration-BO4CO>

# GP for modeling blackbox response function

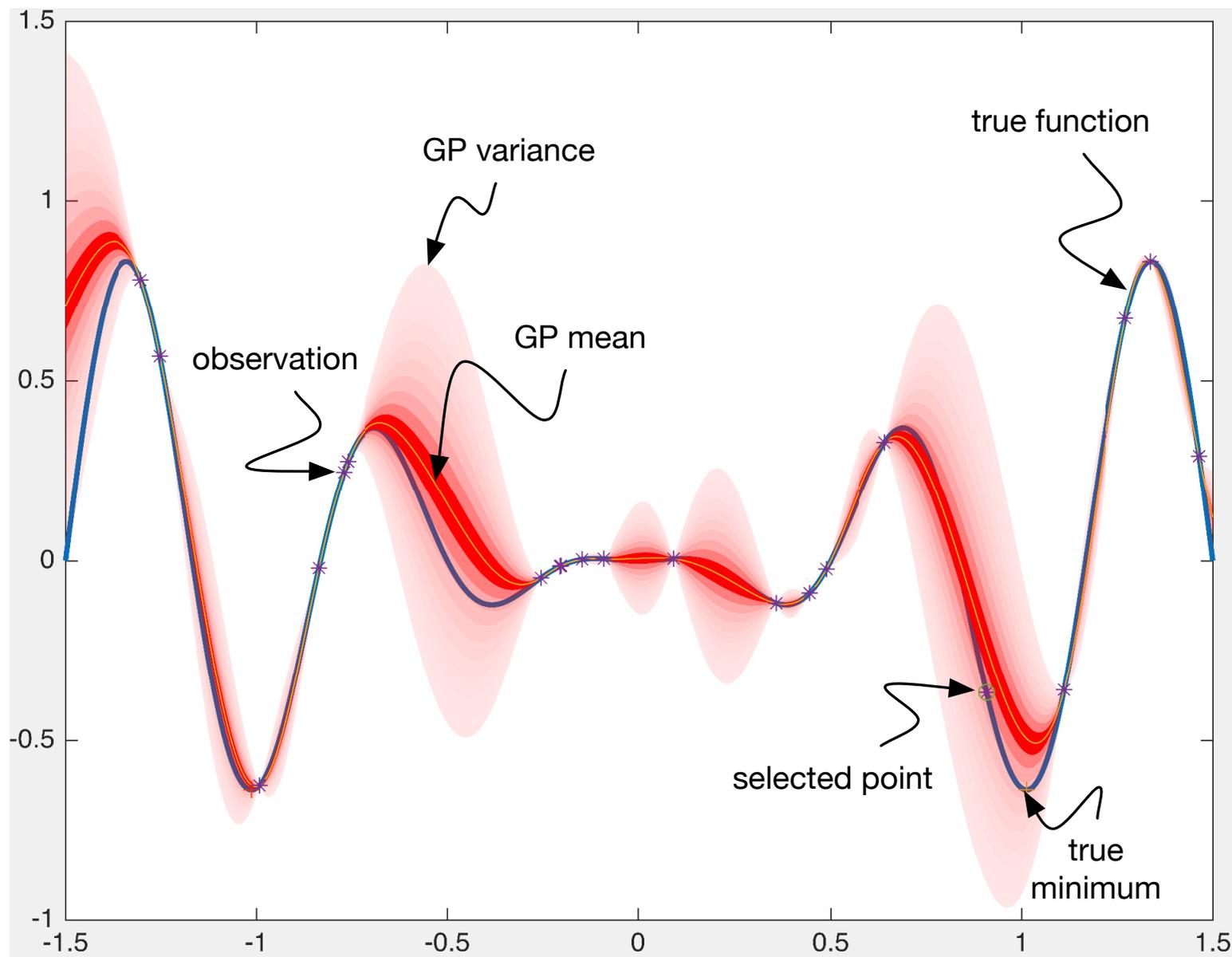
$$y = f(\mathbf{x}) \sim \mathcal{GP}(\mu(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')),$$

$$\mu_t(\mathbf{x}) = \mu(\mathbf{x}) + \mathbf{k}(\mathbf{x})^\top (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} (\mathbf{y} - \boldsymbol{\mu})$$

$$\sigma_t^2(\mathbf{x}) = k(\mathbf{x}, \mathbf{x}) + \sigma^2 \mathbf{I} - \mathbf{k}(\mathbf{x})^\top (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}(\mathbf{x})$$

## Motivations:

- 1- mean estimates + variance
- 2- all computations are linear algebra



## Kernel function:

$$\mathbf{K} := \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \dots & k(\mathbf{x}_1, \mathbf{x}_t) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_t, \mathbf{x}_1) & \dots & k(\mathbf{x}_t, \mathbf{x}_t) \end{bmatrix}$$

$$k_\theta(\mathbf{x}_i, \mathbf{x}_j) = \exp(\sum_{\ell=1}^d (-\theta_\ell \delta(\mathbf{x}_i \neq \mathbf{x}_j))),$$

## Acquisition function:

$$u_{LCB}(\mathbf{x} | \mathcal{M}, \mathbb{S}_{1:n}) = \operatorname{argmin}_{\mathbf{x} \in \mathbb{X}} \mu_t(\mathbf{x}) - \kappa \sigma_t(\mathbf{x}),$$

## Code:

<https://github.com/pooyanjamshidi/BO4CO>

---

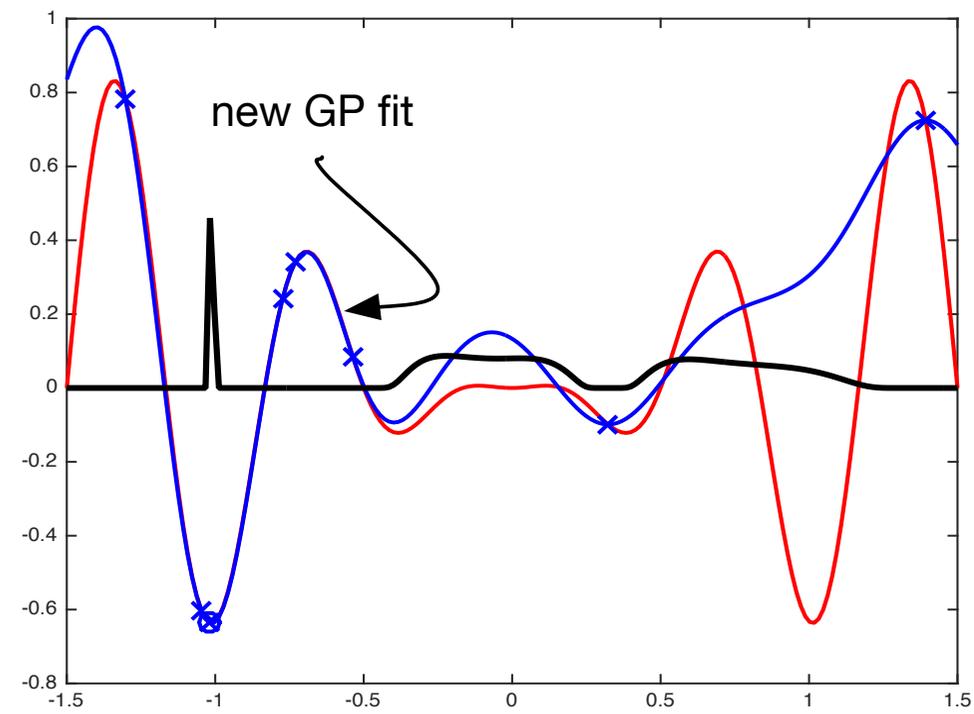
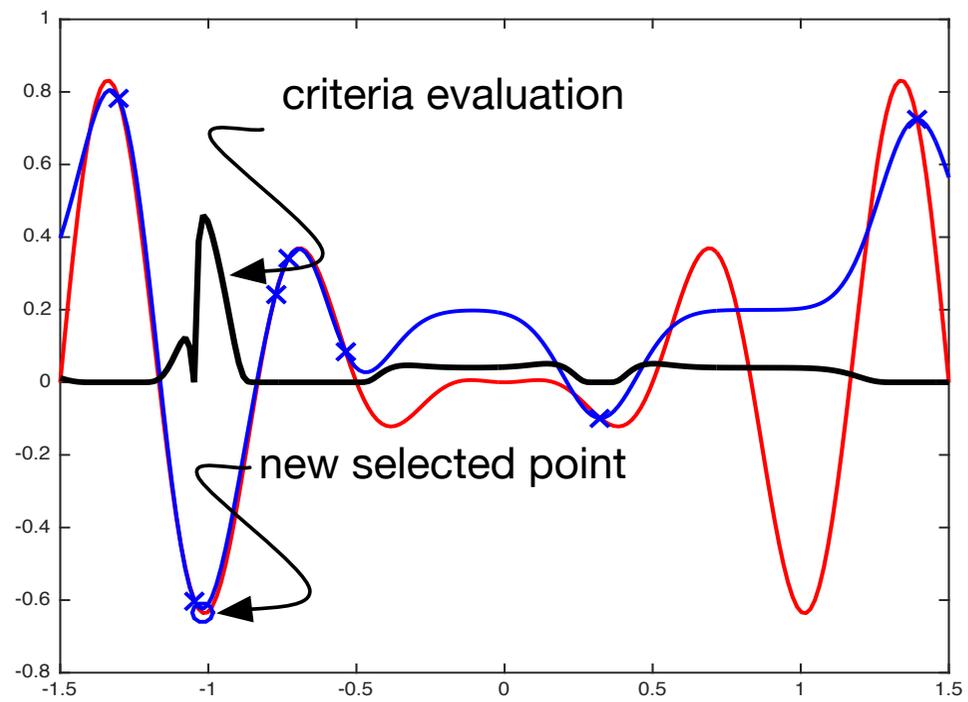
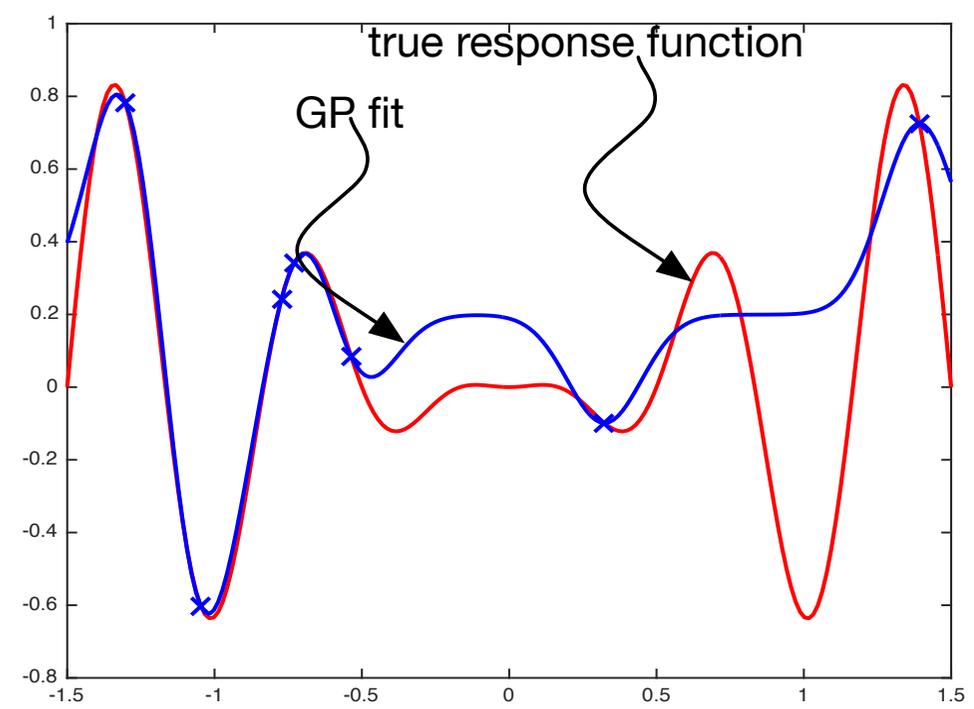
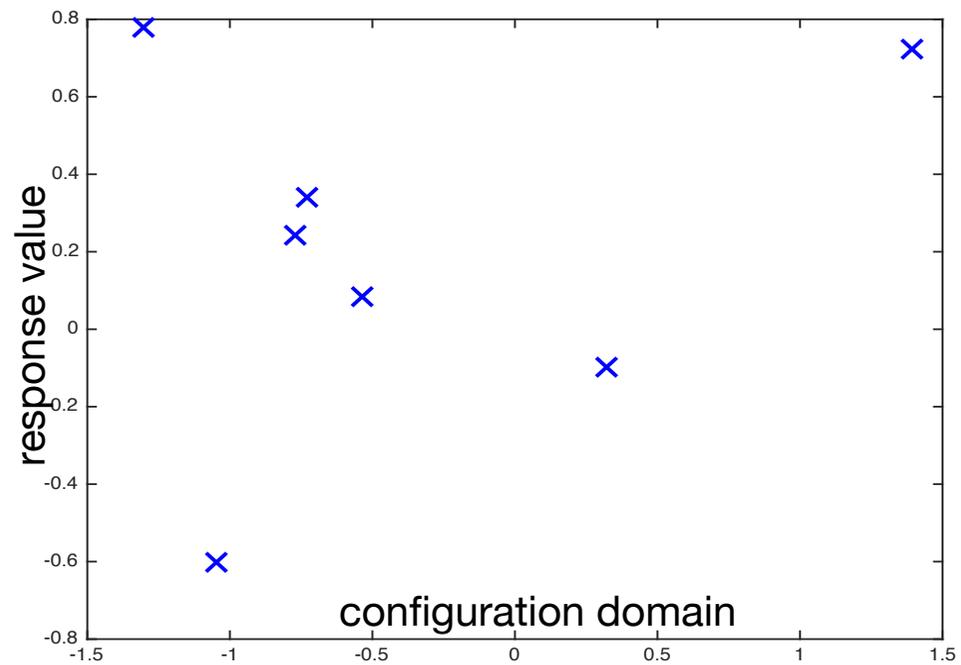
## Algorithm 1 : BO4CO

---

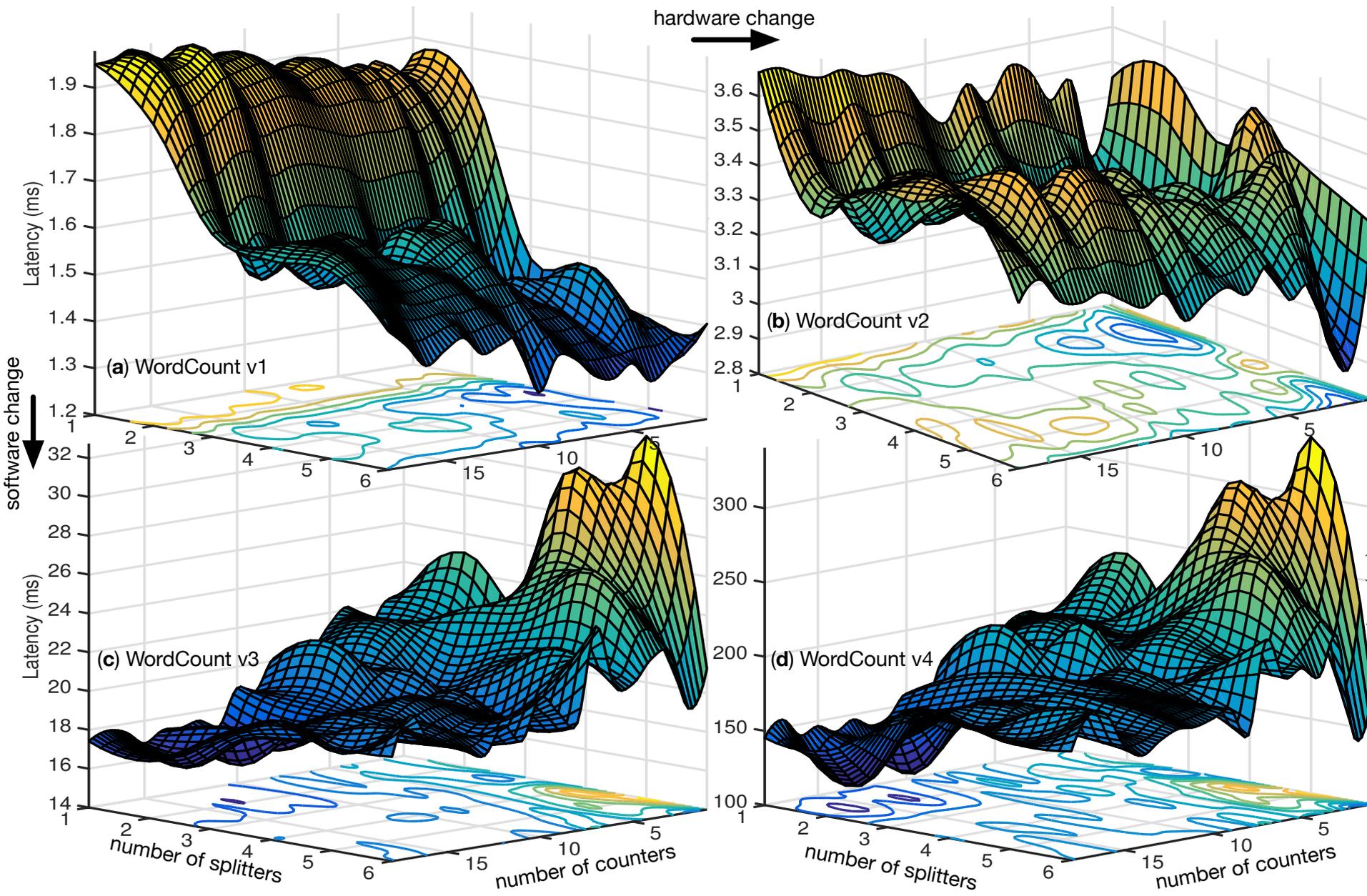
**Input:** Configuration space  $\mathbb{X}$ , Maximum budget  $N_{max}$ , Response function  $f$ , Kernel function  $\mathbf{K}_\theta$ , Hyper-parameters  $\theta$ , Design sample size  $n$ , learning cycle  $N_l$

**Output:** Optimal configurations  $\mathbf{x}^*$  and learned model  $\mathcal{M}$

- 1: choose an initial sparse design (*lhd*) to find an initial design samples  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$
  - 2: obtain *performance measurements* of the initial design,  $y_i \leftarrow f(\mathbf{x}_i) + \epsilon_i, \forall \mathbf{x}_i \in \mathcal{D}$
  - 3:  $\mathbb{S}_{1:n} \leftarrow \{(\mathbf{x}_i, y_i)\}_{i=1}^n; t \leftarrow n + 1$
  - 4:  $\mathcal{M}(\mathbf{x} | \mathbb{S}_{1:n}, \theta) \leftarrow$  fit a  $\mathcal{GP}$  model to the design  $\triangleright$  Eq.(3)
  - 5: **while**  $t \leq N_{max}$  **do**
  - 6:     if  $(t \bmod N_l = 0)$   $\theta \leftarrow$  *learn* the kernel hyper-parameters by maximizing the likelihood
  - 7:     find *next configuration*  $\mathbf{x}_t$  by optimizing the selection criteria over the estimated response surface given the data,  $\mathbf{x}_t \leftarrow \operatorname{argmax}_{\mathbf{x}} u(\mathbf{x} | \mathcal{M}, \mathbb{S}_{1:t-1})$   $\triangleright$  Eq.(9)
  - 8:     obtain performance for the *new configuration*  $\mathbf{x}_t, y_t \leftarrow f(\mathbf{x}_t) + \epsilon_t$
  - 9:     Augment the configuration  $\mathbb{S}_{1:t} = \{\mathbb{S}_{1:t-1}, (\mathbf{x}_t, y_t)\}$
  - 10:      $\mathcal{M}(\mathbf{x} | \mathbb{S}_{1:t}, \theta) \leftarrow$  *re-fit* a new GP model  $\triangleright$  Eq.(7)
  - 11:      $t \leftarrow t + 1$
  - 12: **end while**
  - 13:  $(\mathbf{x}^*, y^*) = \min_{\mathbb{S}_{1:N_{max}}}$
  - 14:  $\mathcal{M}(\mathbf{x})$
-



# Correlations across different versions



correlation coefficient →

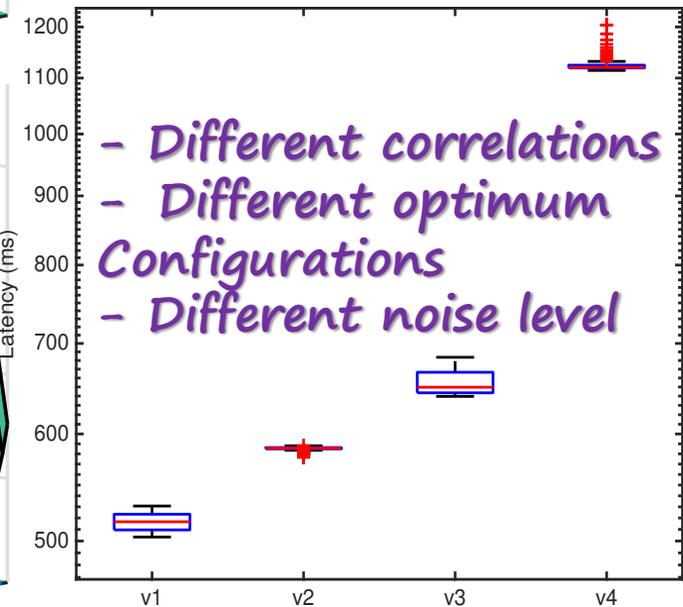
p-value ↓

	v1	v2	v3	v4
v1	1	0.41	-0.46	-0.50
v2	7.36E-06	1	-0.20	-0.18
v3	6.92E-07	0.04	1	0.94
v4	2.54E-08	0.07	1.16E-52	1

(e) Pearson correlation coefficients

	v1	v2	v3	v4
v1	1	0.49	-0.51	-0.51
v2	5.50E-08	1	-0.2793	-0.24
v3	1.30E-08	0.003	1	0.88
v4	1.40E-08	0.01	8.30E-36	1

(f) Spearman correlation coefficients

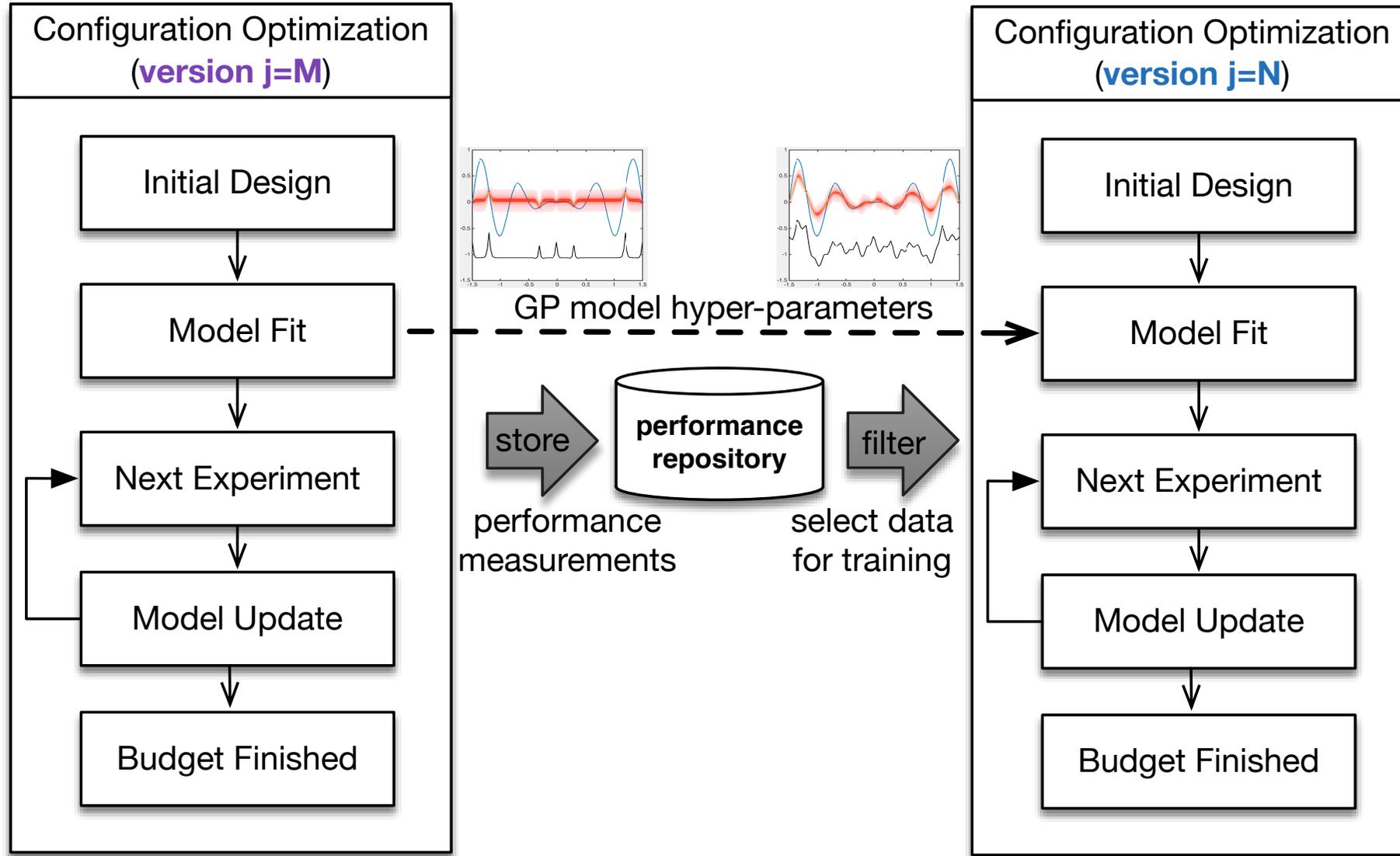


(g) Measurement noise across WordCount versions

# DevOps

- Different versions are continuously delivered (daily basis).
- Big Data systems are developed using similar frameworks (Apache Storm, Spark, Hadoop, Kafka, etc).
- Different versions share similar business logics.

# Solution: Transfer Learning for Configuration Optimization





Part

# Transfer Learning for Configuration Optimization (TL4CO)

Code: <https://github.com/dice-project/DICE-Configuration-TL4CO>

$$\mu_t(\mathbf{x}) = \mu(\mathbf{x}) + \mathbf{k}^\top (\mathbf{K}(X, X) + \Sigma)^{-1} (\mathbf{y} - \boldsymbol{\mu})$$

$$\sigma_t^2(\mathbf{x}) = k(\mathbf{x}, \mathbf{x}) + \sigma^2 - \mathbf{k}^\top (\mathbf{K}(X, X) + \Sigma)^{-1} \mathbf{k},$$

$$k_{TL4CO}(l, l', \mathbf{x}, \mathbf{x}') = k_t(l, l') \times k_{xx}(\mathbf{x}, \mathbf{x}'),$$

correlation between  
different versions

covariance functions

---

### Algorithm 1 : TL4CO

---

**Input:** Configuration space  $\mathbb{X}$ , Number of historical configuration optimization datasets  $T$ , Maximum budget  $N_{max}$ , Response function  $f$ , Kernel function  $\mathbf{K}$ , Hyperparameters  $\boldsymbol{\theta}^j$ , The current dataset  $\mathbb{S}^{j=1}$ , Datasets belonging to other versions  $\mathbb{S}^{j=2:T}$ , Diagonal noise matrix  $\Sigma$ , Design sample size  $n$ , Learning cycle  $N_l$

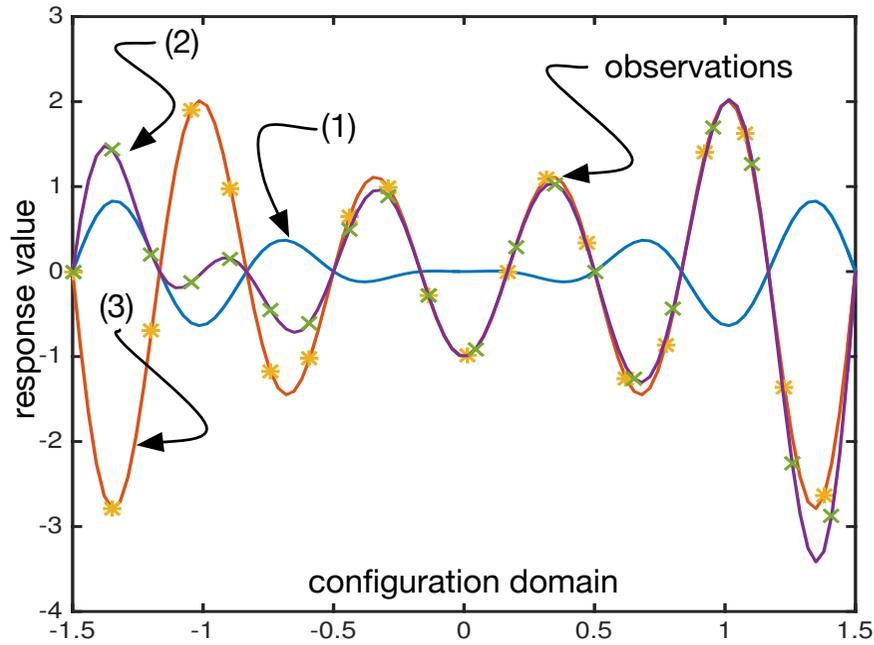
**Output:** Optimal configurations  $\mathbf{x}^*$  and learned model  $\mathcal{M}$

- 1:  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  create an initial sparse design (*lhd*)
  - 2: Obtain the *performance*  $\forall x_i \in \mathcal{D}, y_i \leftarrow f(\mathbf{x}_i) + \epsilon_i$
  - 3:  $\mathbb{S}^{j=2:T} \leftarrow$  select  $m$  points from other versions ( $j = 2 : T$ ) that reduce entropy  $\triangleright$  Eq.(8)
  - 4:  $\mathbb{S}_{1:n}^1 \leftarrow \mathbb{S}^{j=2:T} \cup \{(\mathbf{x}_i, y_i)\}_{i=1}^n; t \leftarrow n + 1$
  - 5:  $\mathcal{M}(\mathbf{x} | \mathbb{S}_{1:n}^1, \boldsymbol{\theta}^j) \leftarrow$  fit a multi-task GP to  $\mathcal{D}$   $\triangleright$  Eq. (6)
  - 6: **while**  $t \leq N_{max}$  **do**
  - 7:     If  $(t \bmod N_l = 0)$  then [ $\boldsymbol{\theta} \leftarrow$  learn the kernel hyperparameters by maximizing the likelihood]
  - 8:     Determine the *next configuration*,  $\mathbf{x}_t$ , by optimizing LCB over  $\mathcal{M}$ ,  $\mathbf{x}_t \leftarrow \arg \max_{\mathbf{x}} u(\mathbf{x} | \mathcal{M}, \mathbb{S}_{1:t-1}^1)$   $\triangleright$  Eq.(11)
  - 9:     Obtain the *performance* of  $\mathbf{x}_t$ ,  $y_t \leftarrow f(\mathbf{x}_t) + \epsilon_t$
  - 10:     Augment the configuration  $\mathbb{S}_{1:t}^1 = \mathbb{S}_{1:t-1}^1 \cup \{(\mathbf{x}_t, y_t)\}$
  - 11:      $\mathcal{M}(\mathbf{x} | \mathbb{S}_{1:t}^1, \boldsymbol{\theta}) \leftarrow$  re-fit a new GP model  $\triangleright$  Eq.(6)
  - 12:      $t \leftarrow t + 1$
  - 13: **end while**
  - 14:  $(\mathbf{x}^*, y^*) = \min[\mathbb{S}_{1:N_{max}}^1 - \mathbb{S}^{j=2:T}]$  locating the optimum configuration by discarding data for other system versions
  - 15:  $\mathcal{M}(\mathbf{x})$  storing the learned model
- 

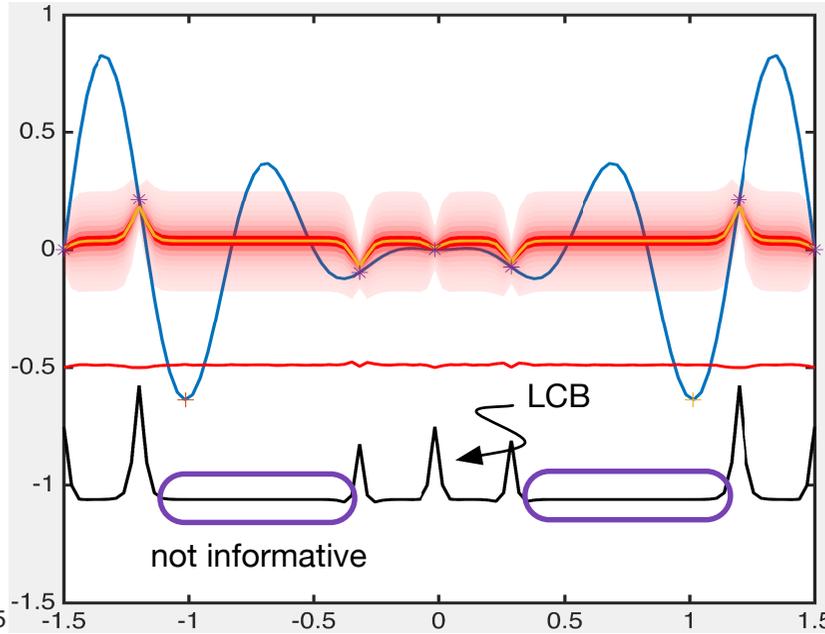
**Code:**

<https://github.com/pooyanjamshidi/TL4CO>

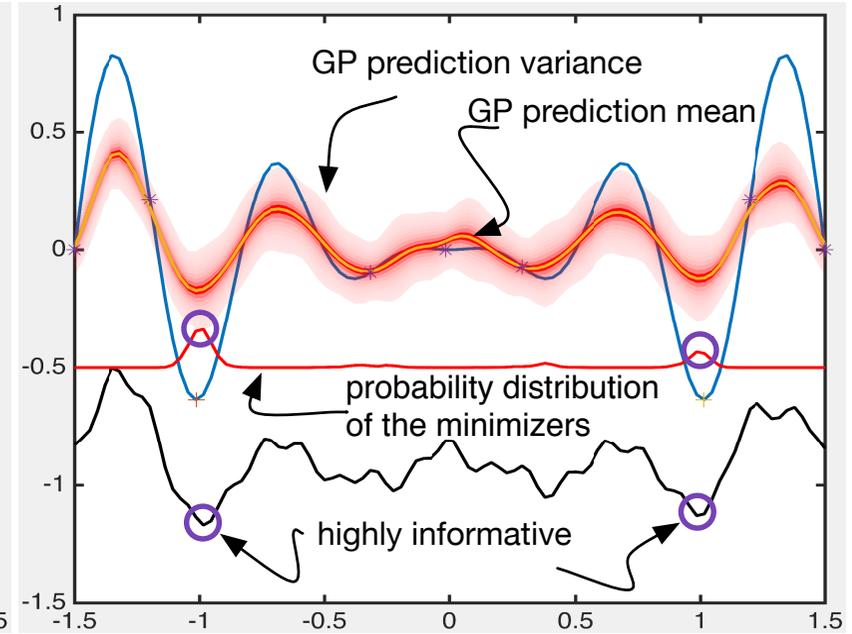
# Multi-task GP vs single-task GP



(a) 3 sample response functions



(b) GP fit for (1) ignoring observations for (2),(3)

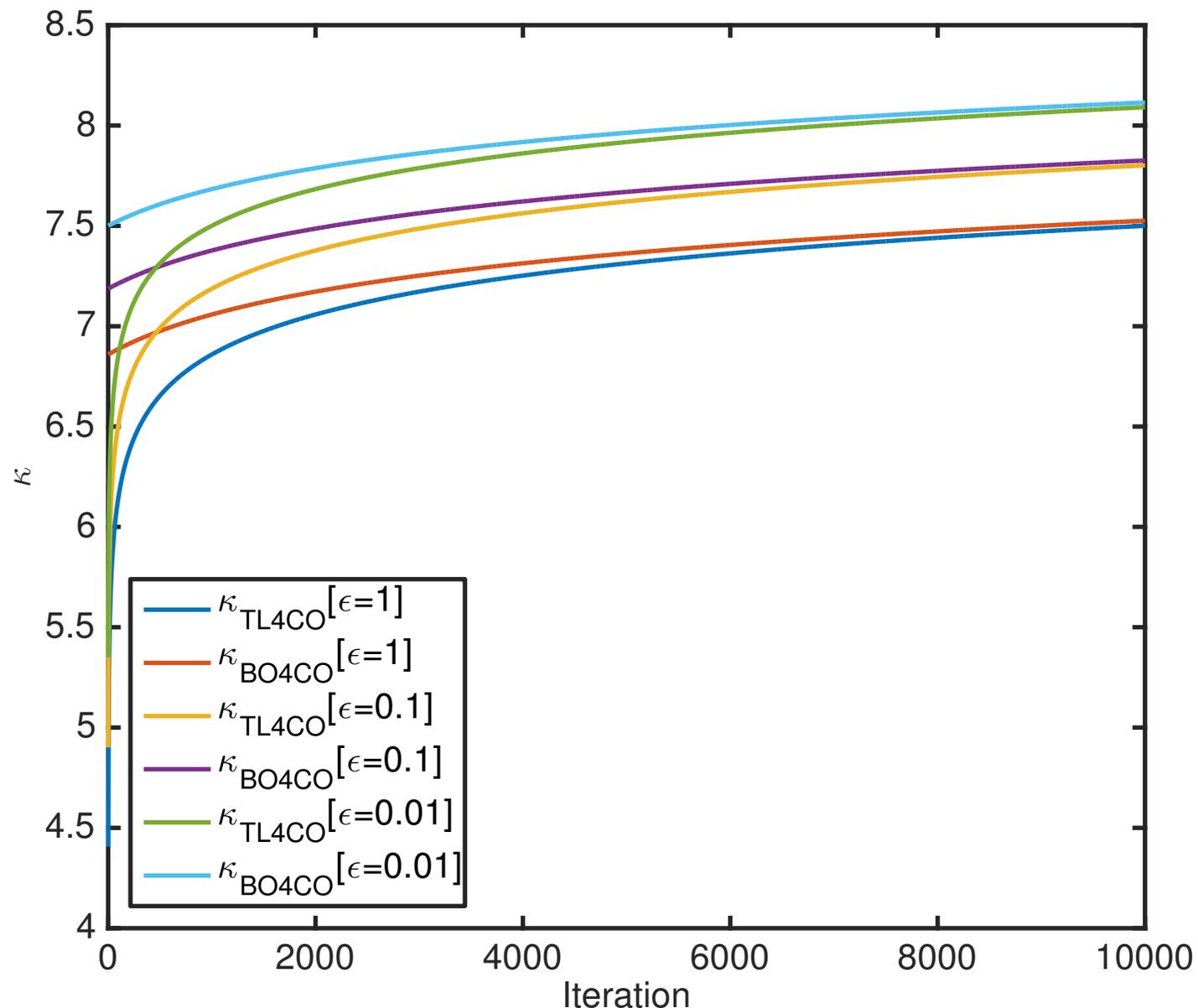


(c) multi-task GP fit for (1) by transfer learning from (2),(3)

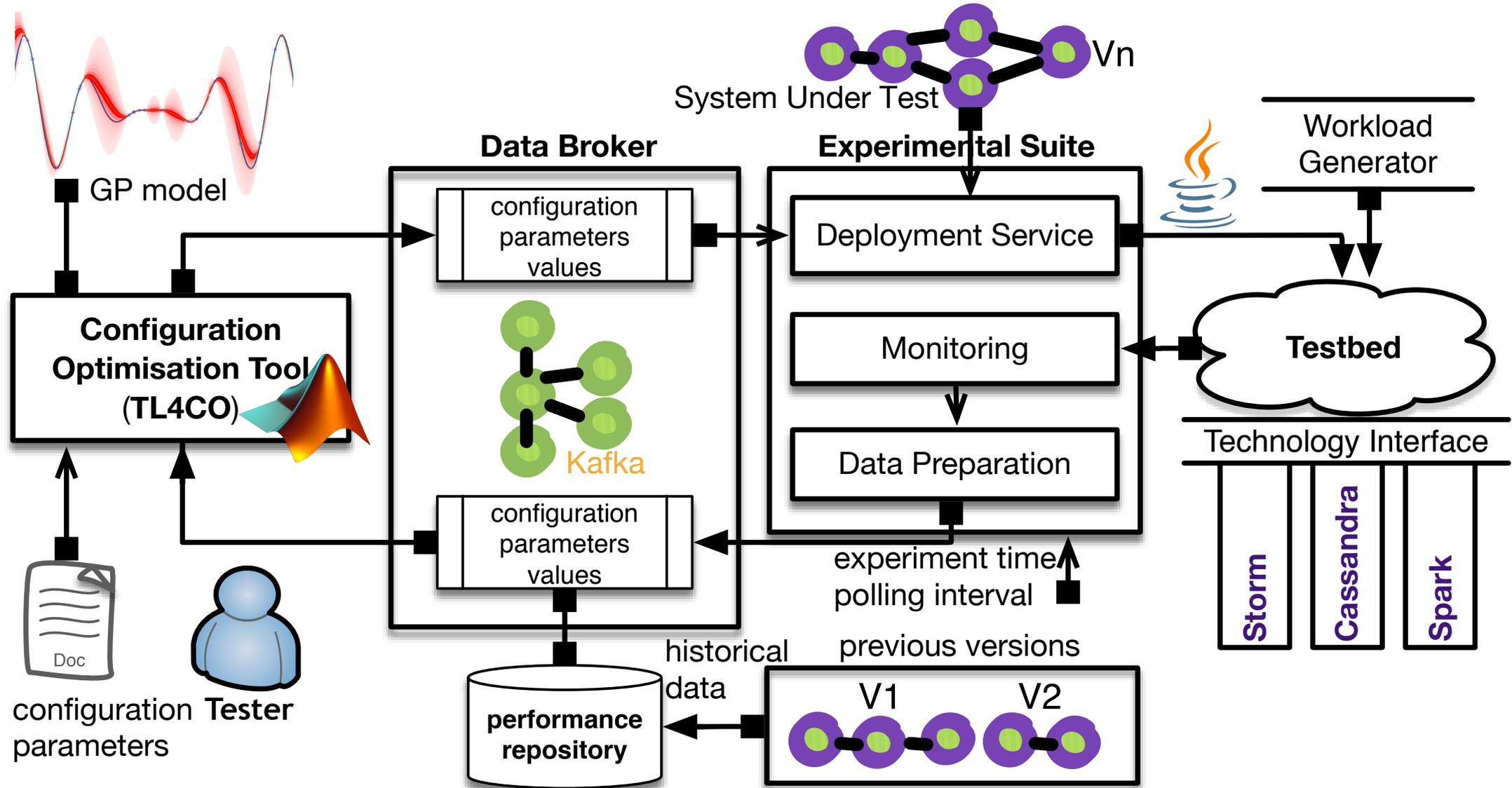
# Exploitation vs exploration

$$\mathbf{x}_{t+1} = \operatorname{argmax}_{\mathbf{x} \in \mathbb{X}} u(\mathbf{x} | \mathcal{M}, \mathbb{S}_{1:t}^1)$$

$$u_{LCB}(\mathbf{x} | \mathcal{M}, \mathbb{S}_{1:t}^1) = \operatorname{argmin}_{\mathbf{x} \in \mathbb{X}} \mu_t(\mathbf{x}) - \kappa \sigma_t(\mathbf{x}),$$

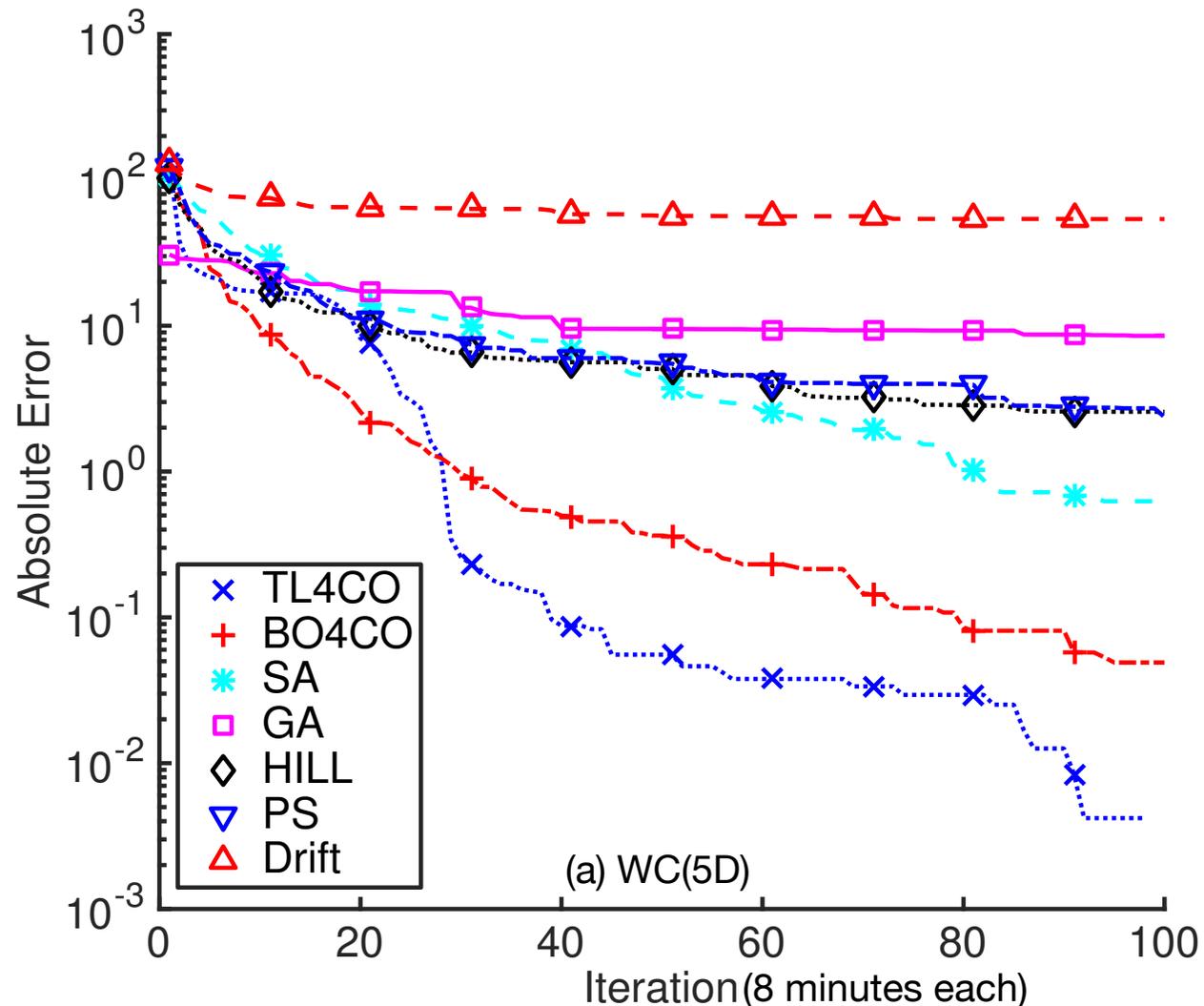
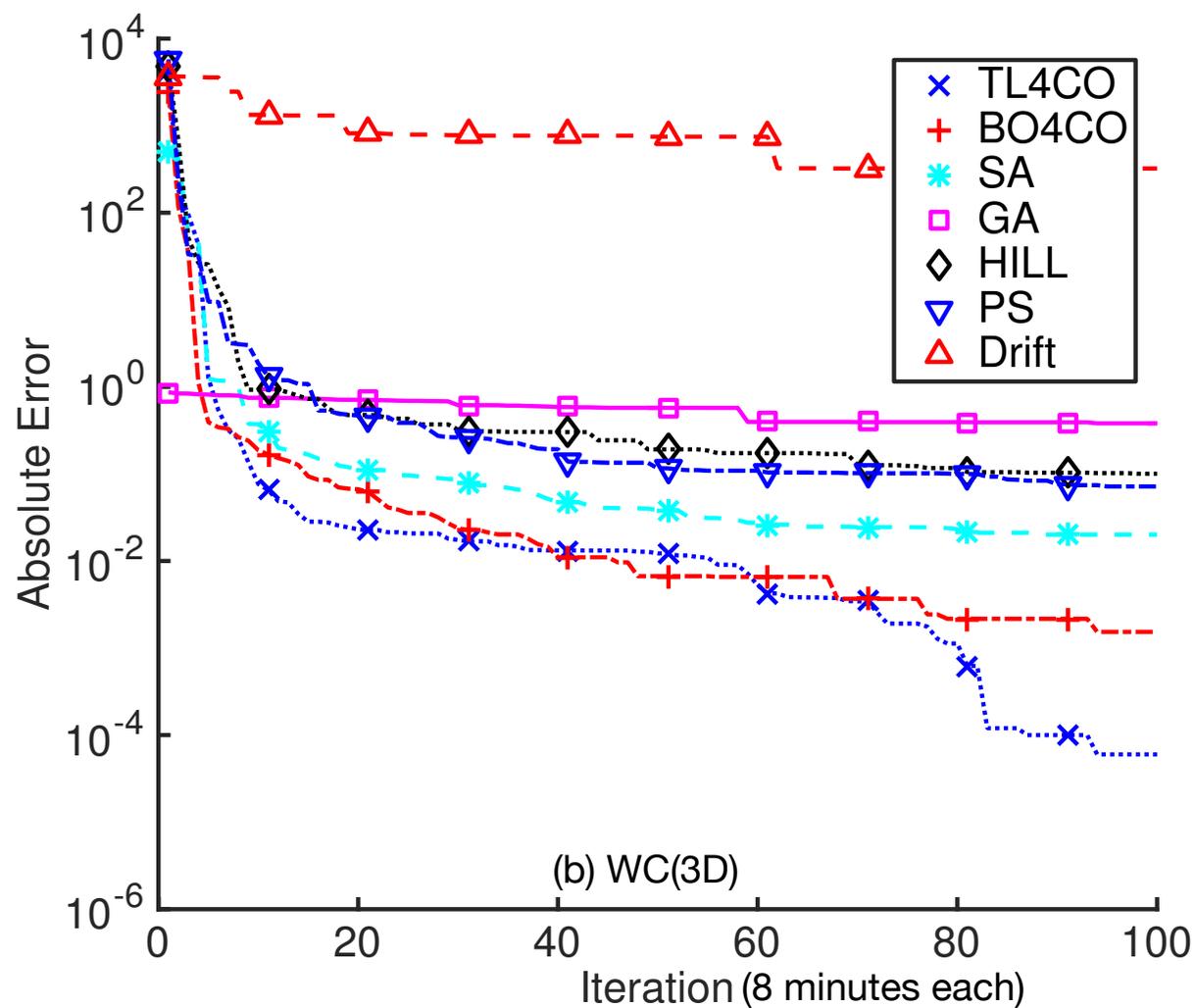


# TL4CO architecture

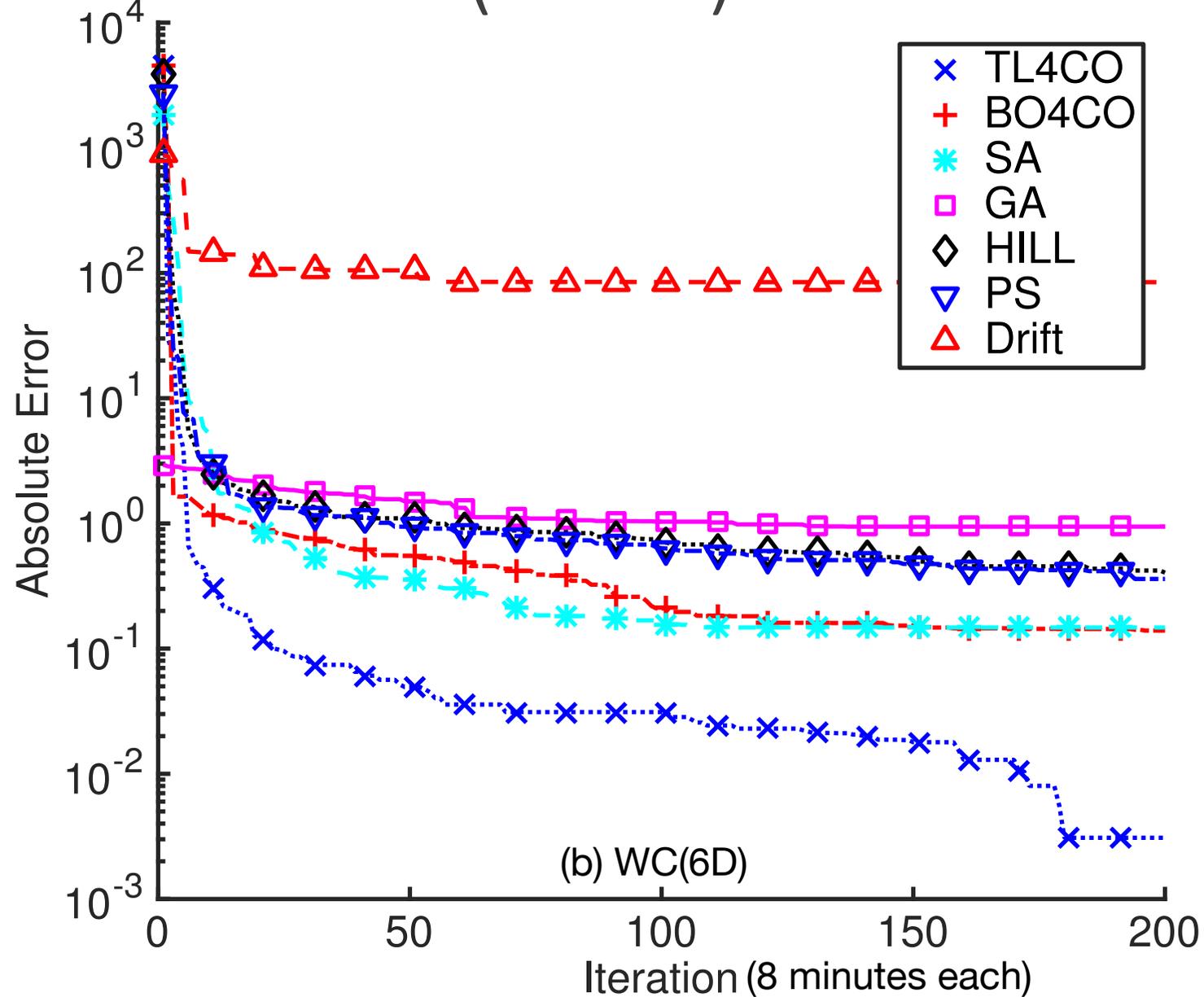


# Experimental results

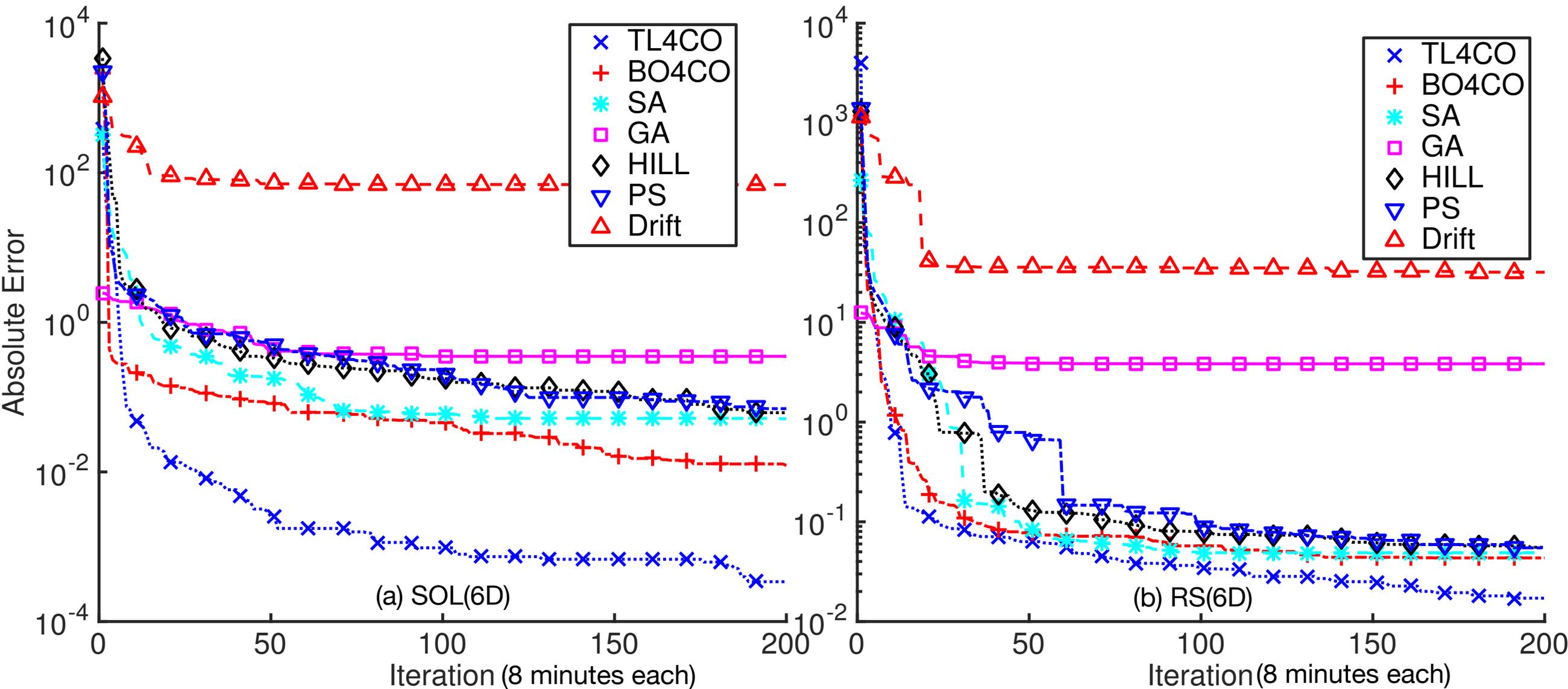
- 30 runs, report average performance  
- Yes, we did full factorial measurements and we know where global min is... 😊



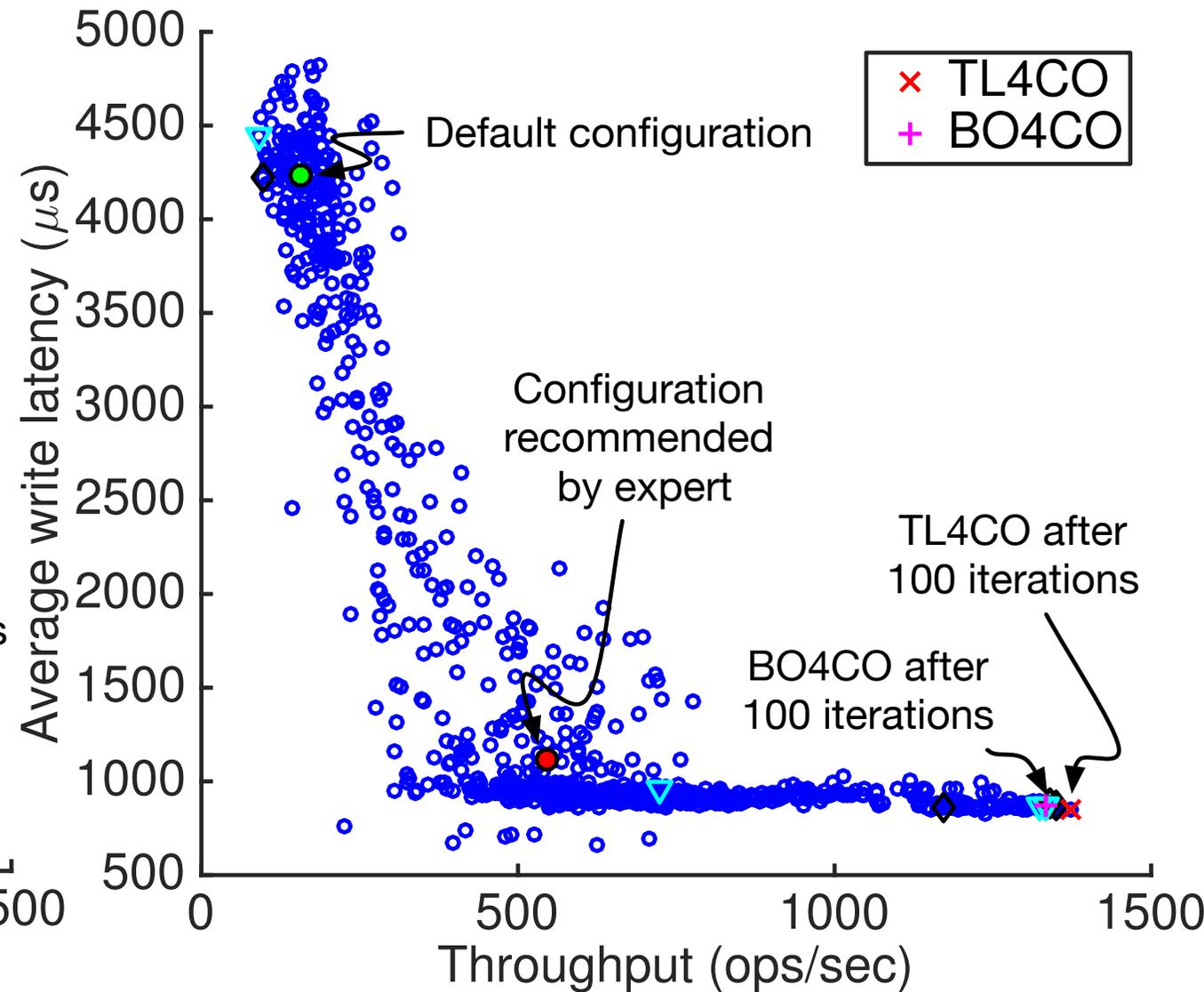
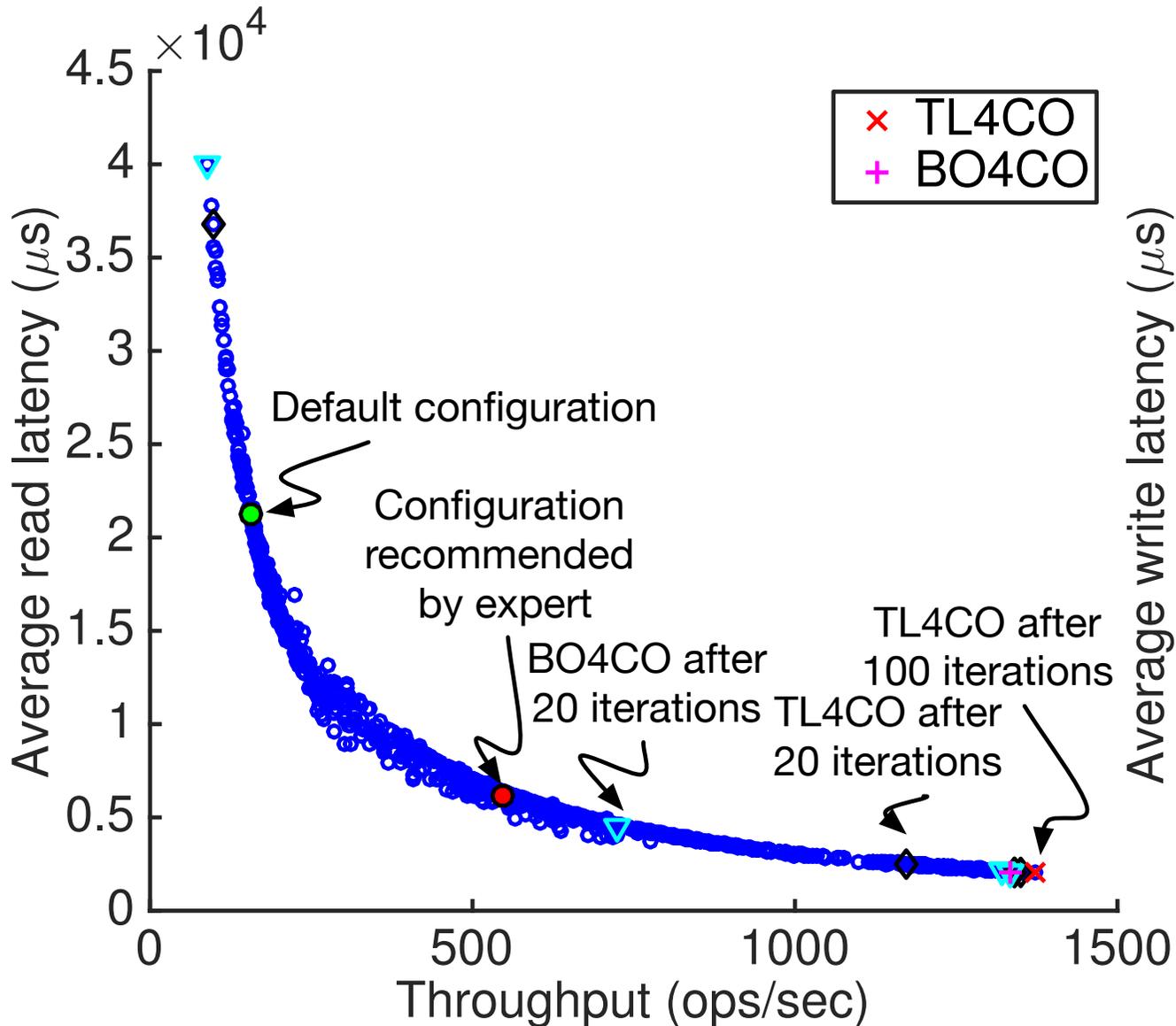
# Experimental results (cont.)



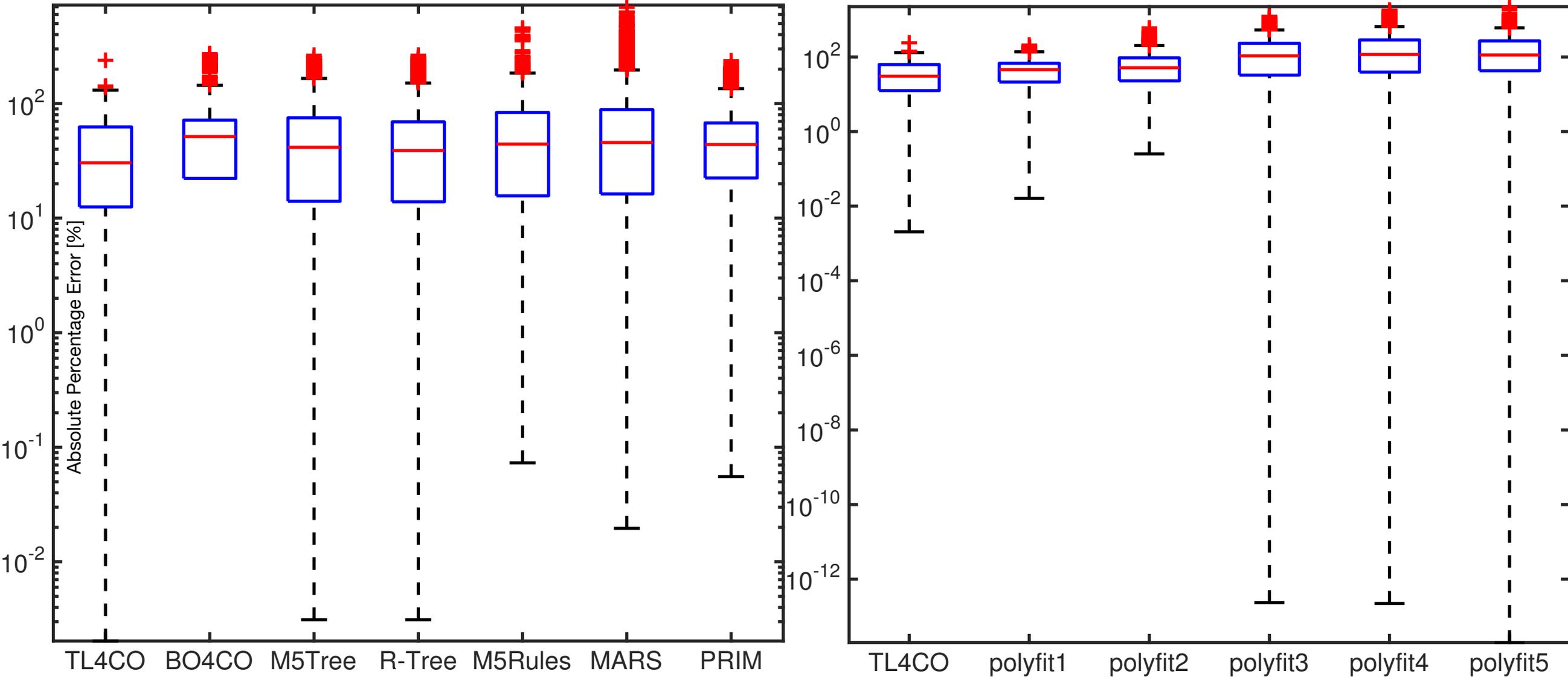
# Experimental results (cont.)



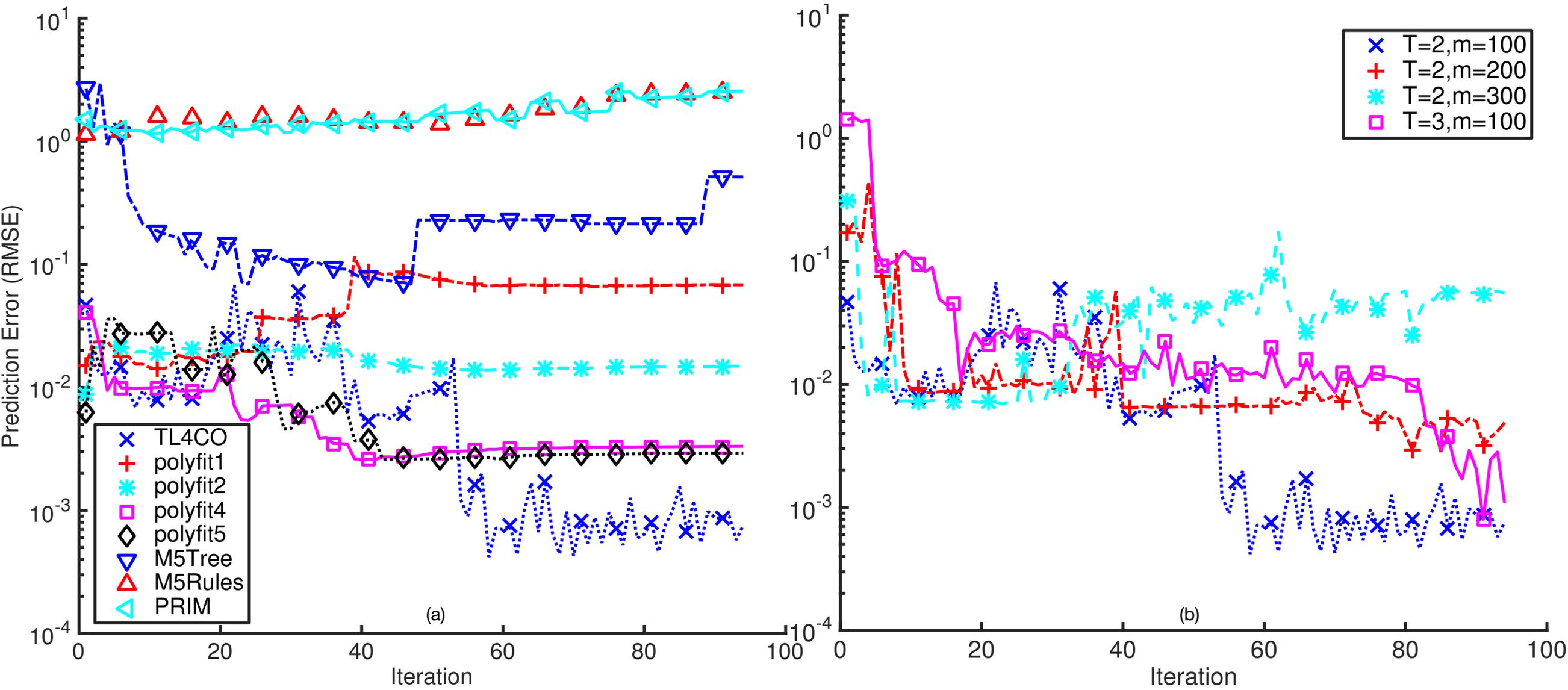
# Comparison with default and expert prescription



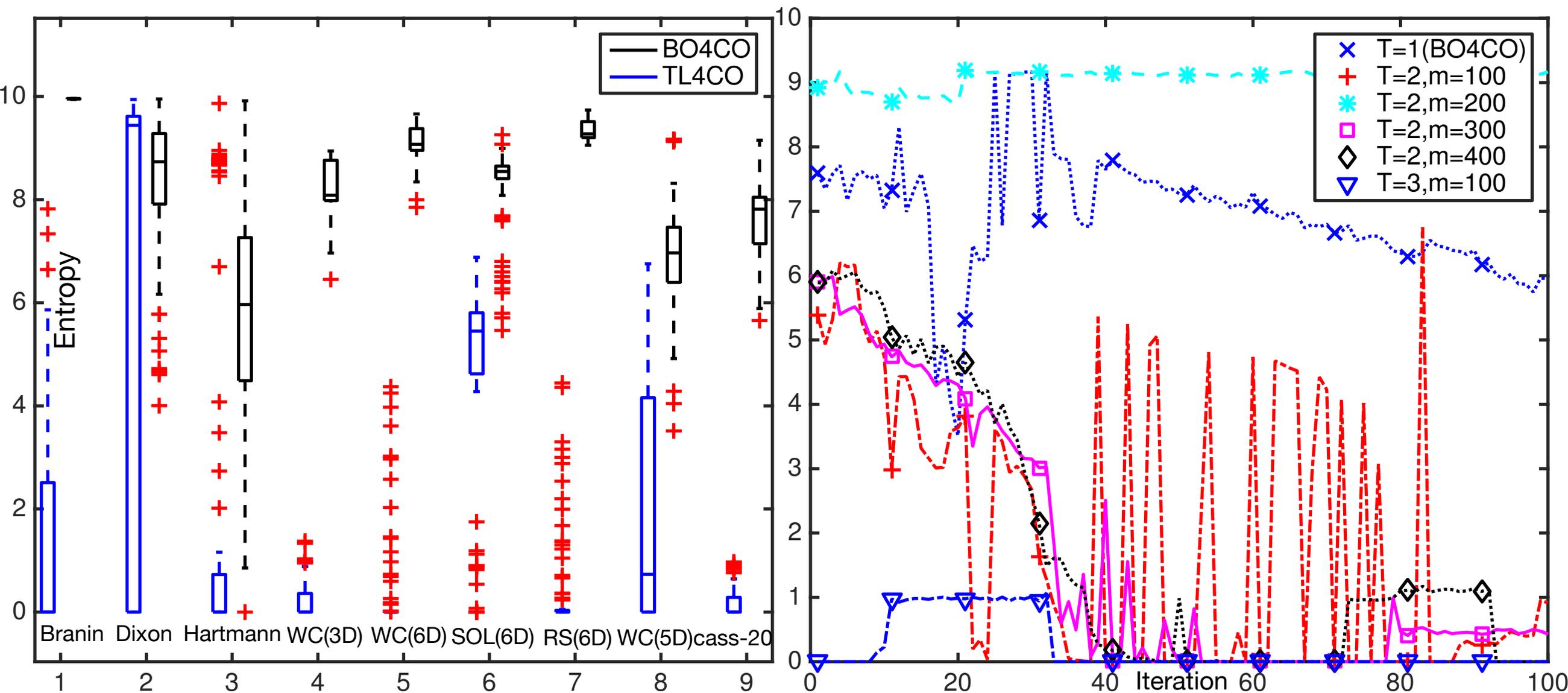
# Model accuracy



# Prediction accuracy over time

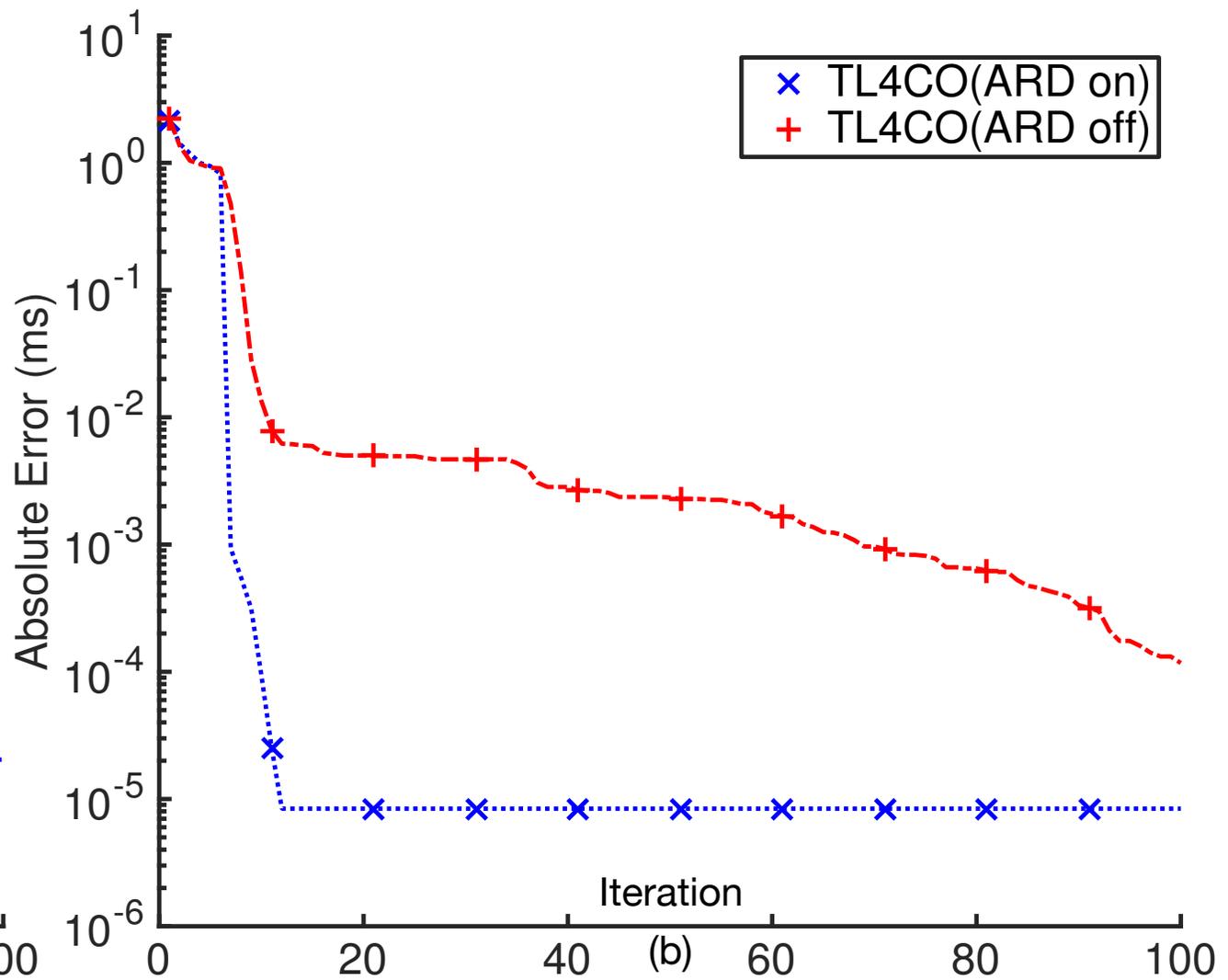
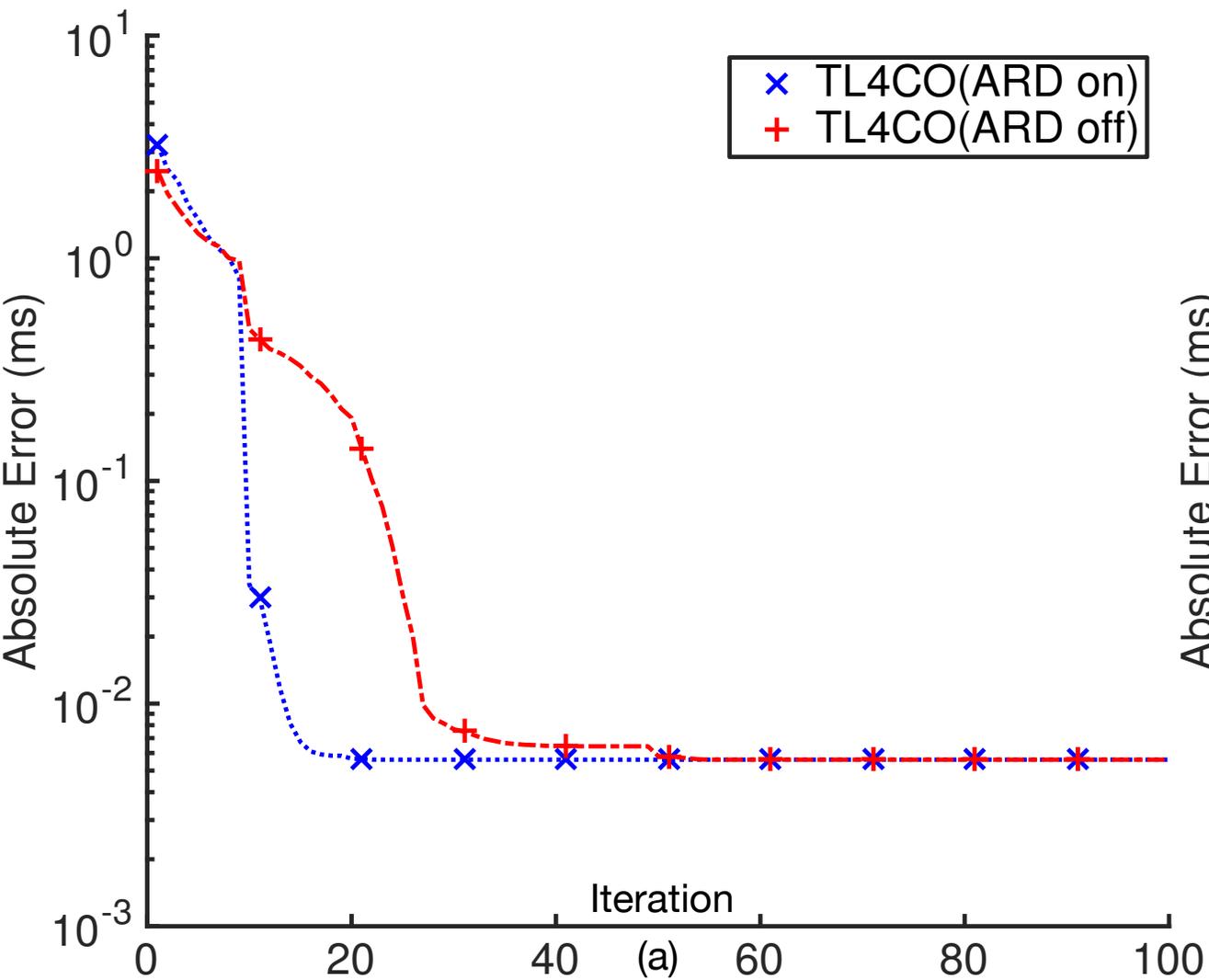


# Entropy of the density function of the minimizers

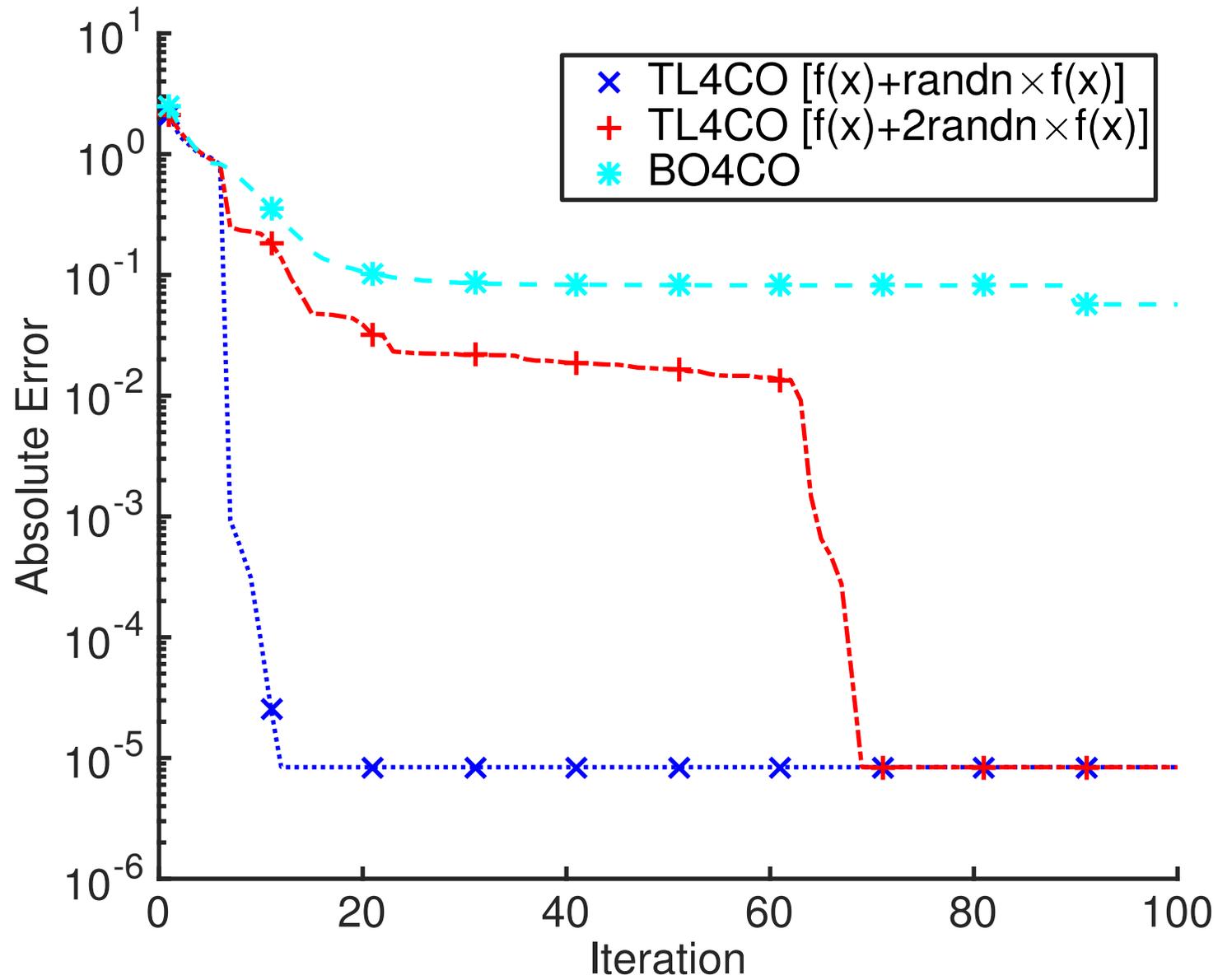


# Effect of ARD

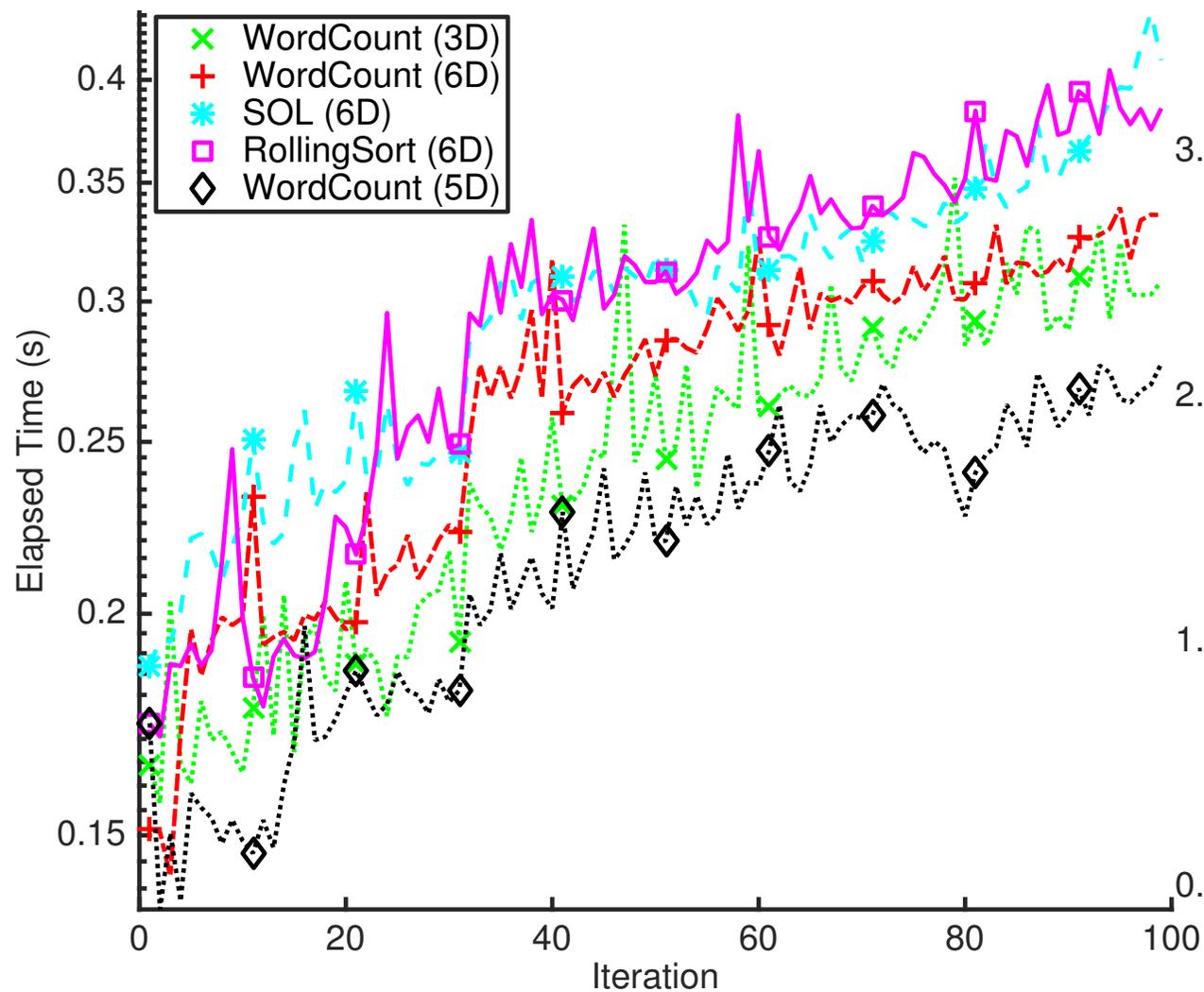
$$k_{xx}(\mathbf{x}_i, \mathbf{x}_j) = \exp(\sum_{\ell=1}^d (-\theta_{\ell} \delta(\mathbf{x}_i \neq \mathbf{x}_j))),$$



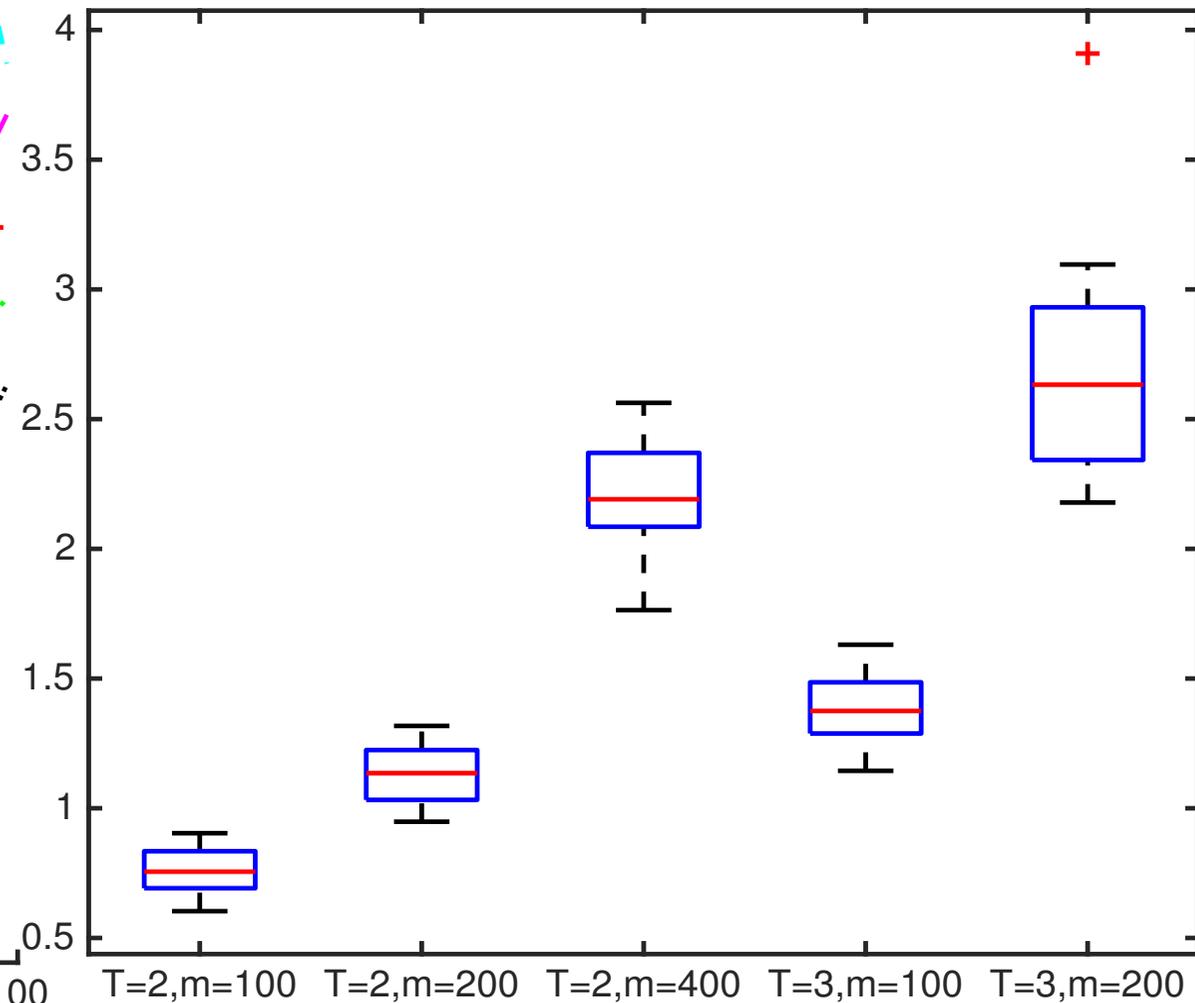
# Effect of correlation between tasks



# Runtime overhead



(a) TL4CO runtime for different datasets

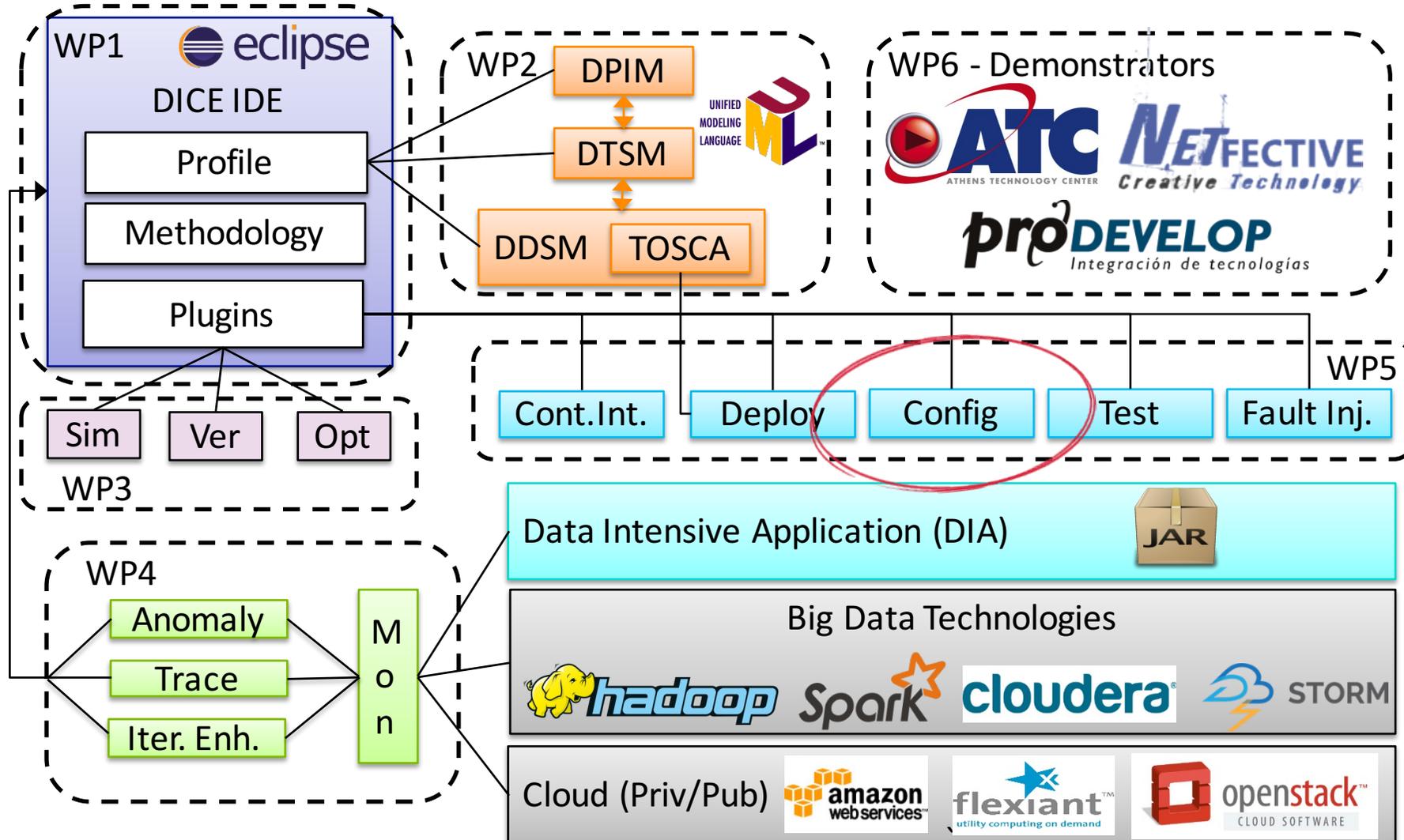


(b) TL4CO runtime for different size of observations

# Key takeaways

- A principled way to leverage prior knowledge gained from searches over previous versions of the system.
- Multi-task GPs can be used in order to capture correlation between related tuning tasks.
- MTGPs are more accurate than STGP and other regression models.
- Application to SPSs, Batch and NoSQL
- Lead to a better performance in practice

Acknowledgement: BO4CO and TL4CO are now integrated with other DevOps tools in the delivery pipeline for Big Data in H2020 DICE project (<http://www.dice-h2020.eu/>)



Tool Demo: <http://www.slideshare.net/pooyanjamshidi/configuration-optimization-tool>