

An Uncertainty-Aware Approach to Optimal Configuration of Stream Processing Systems

Pooyan Jamshidi
(joint work with Giuliano Casale)
Imperial College London
p.jamshidi@imperial.ac.uk

University of Bern
1st Nov 2016

**Imperial College
London**

Motivation

1- Many different Parameters =>
- large state space
- interactions

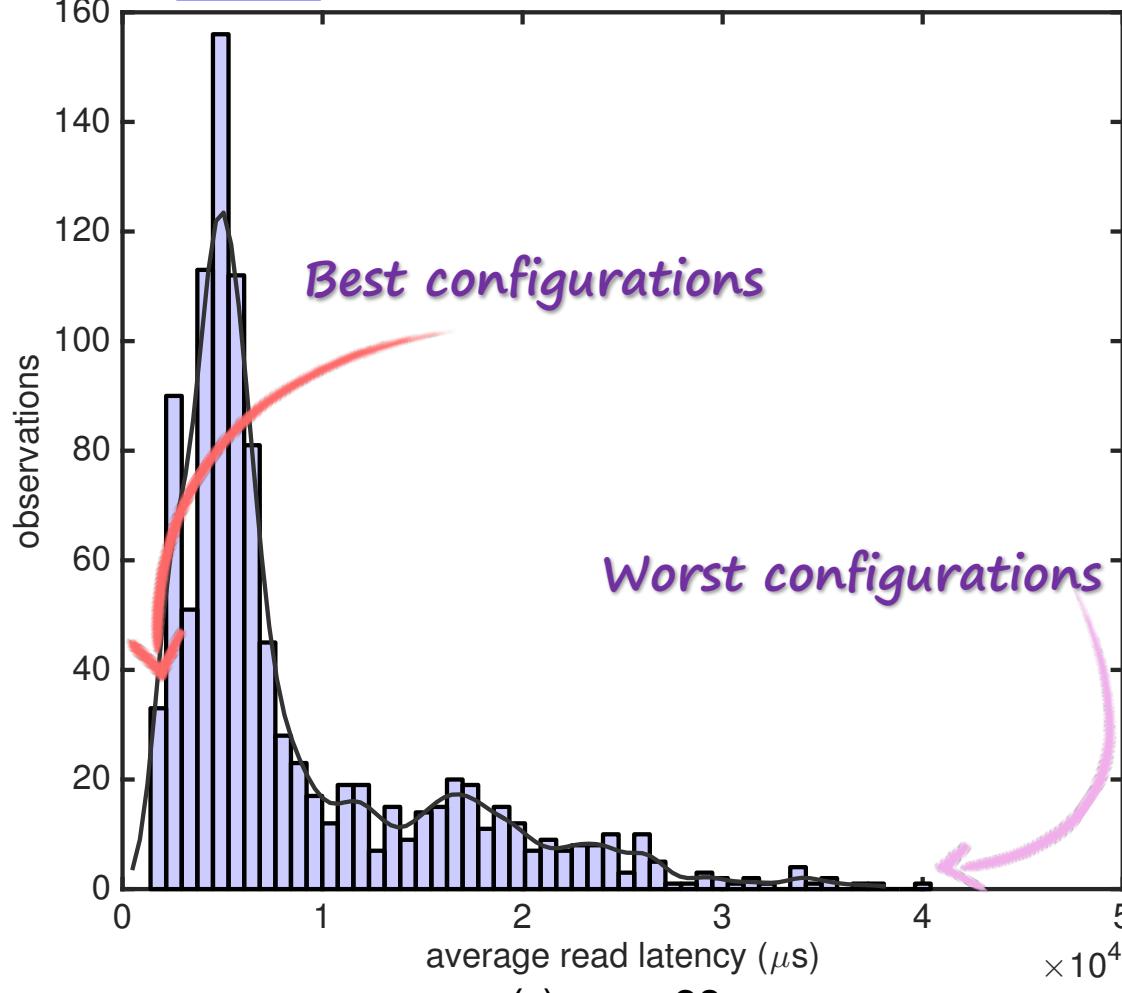
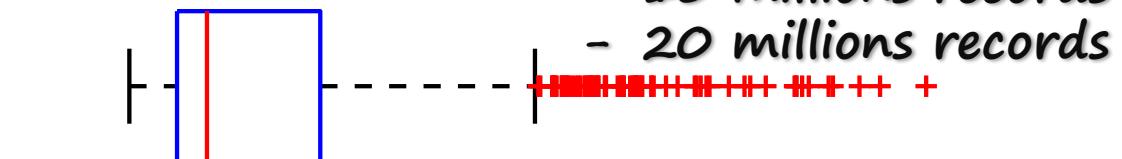
2- Defaults are typically used =>
- poor performance

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

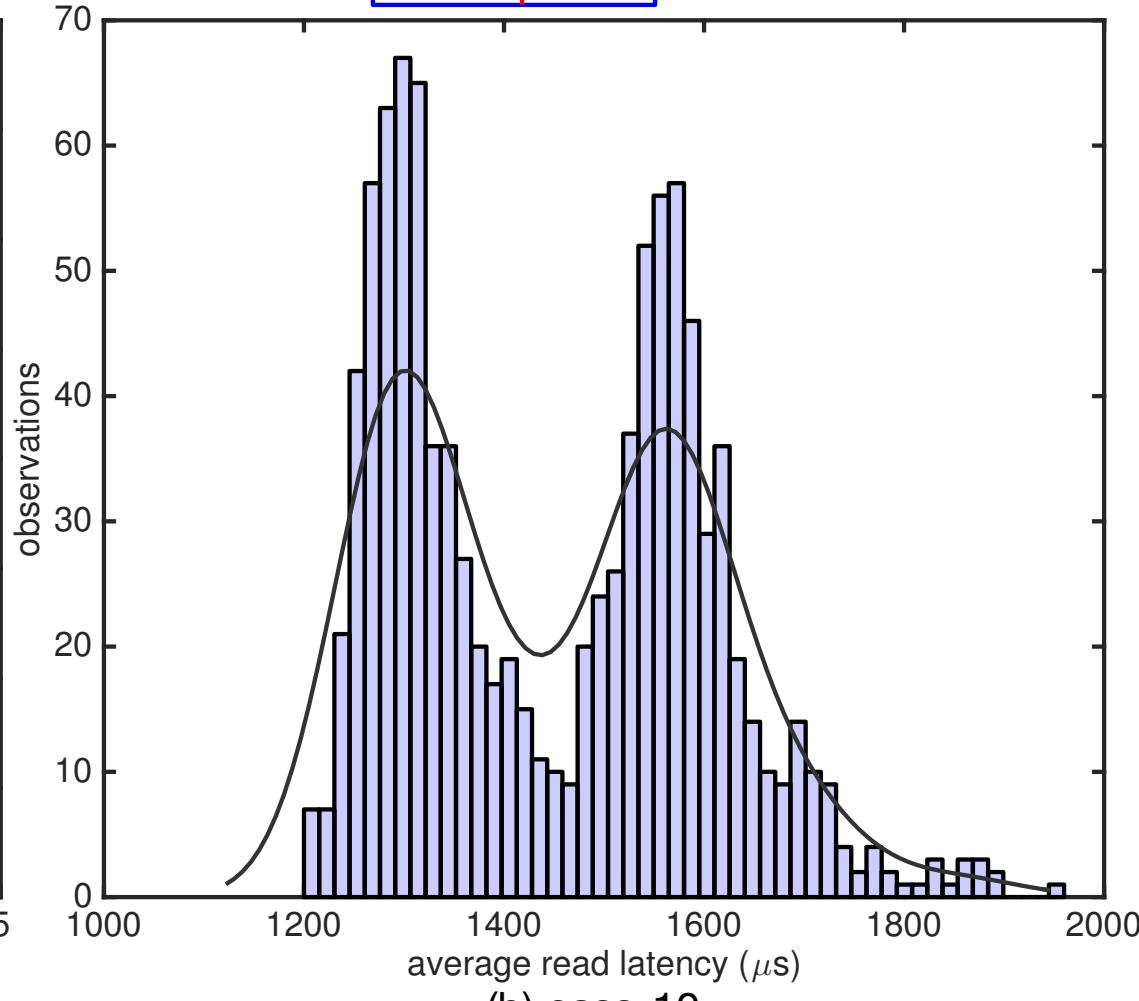
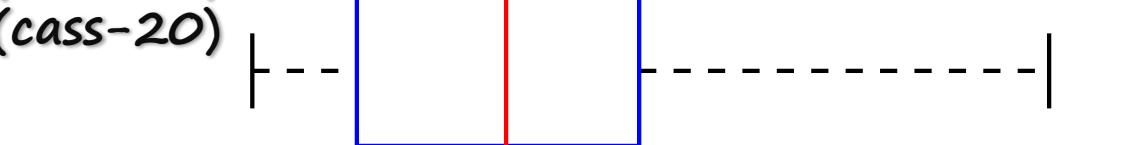
Motivation

Experiments on Apache Cassandra:

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



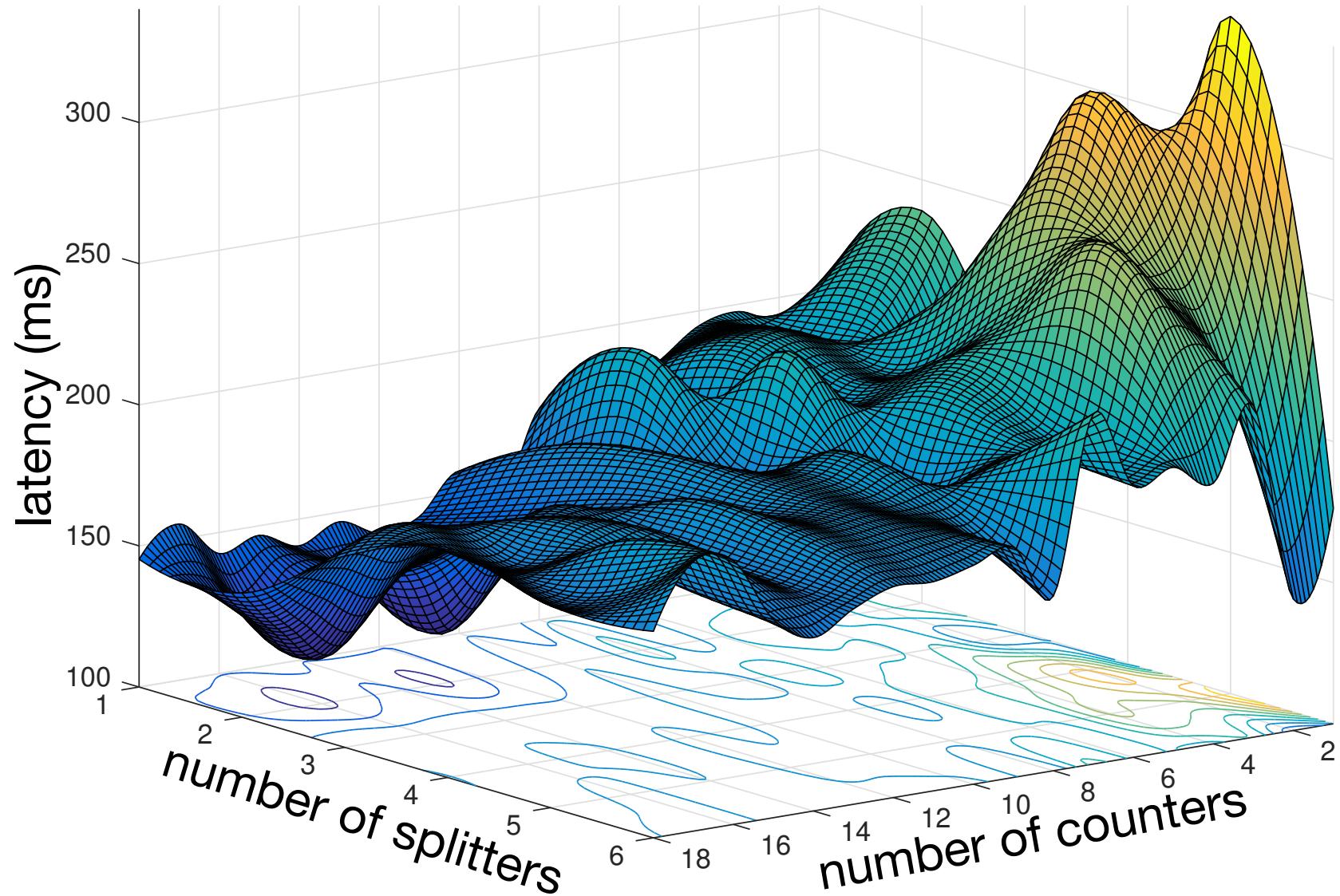
(a) cass-20



(b) cass-10

Motivation (Apache Storm)

In our experiments we observed improvement up to 100%



Goal

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

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

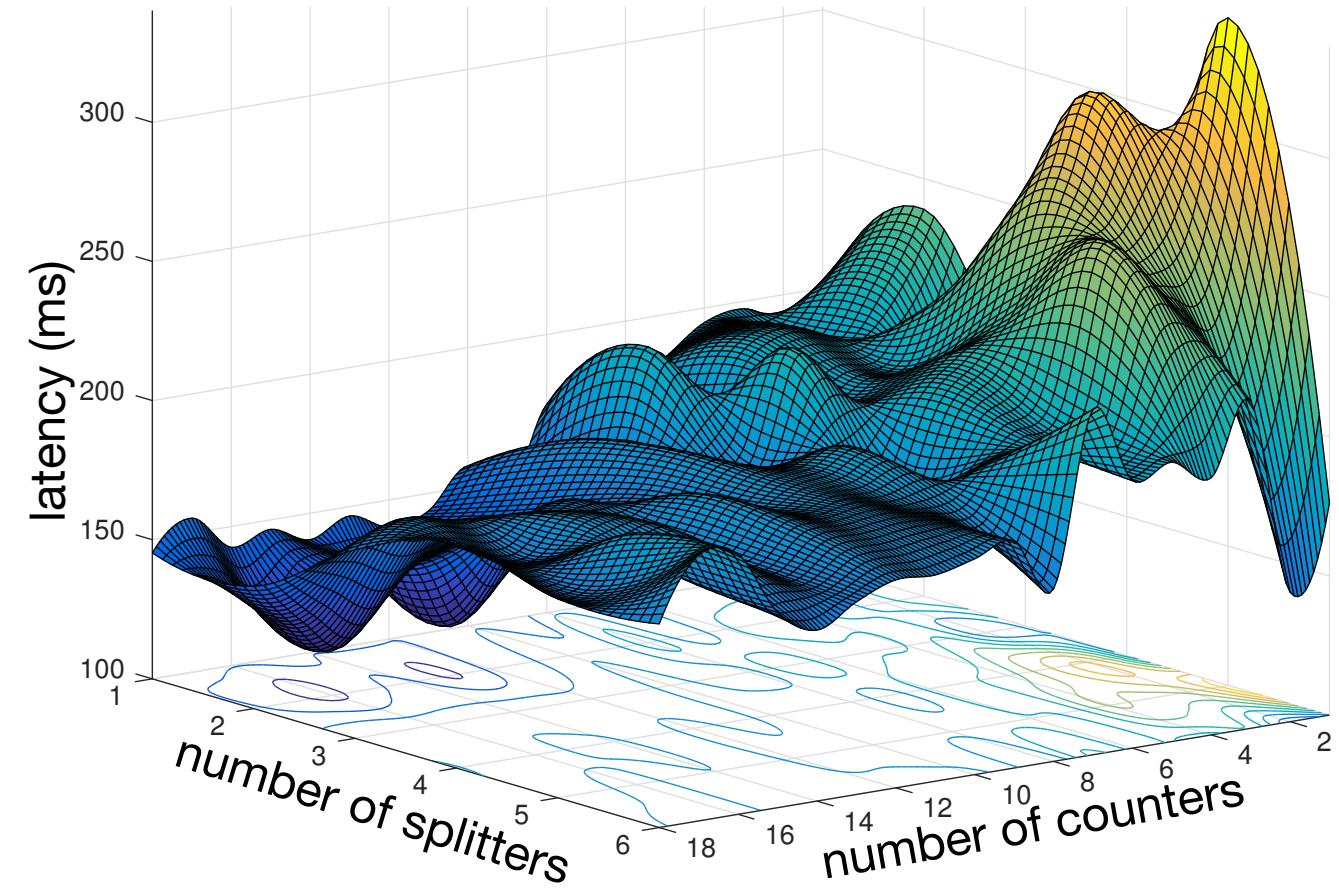
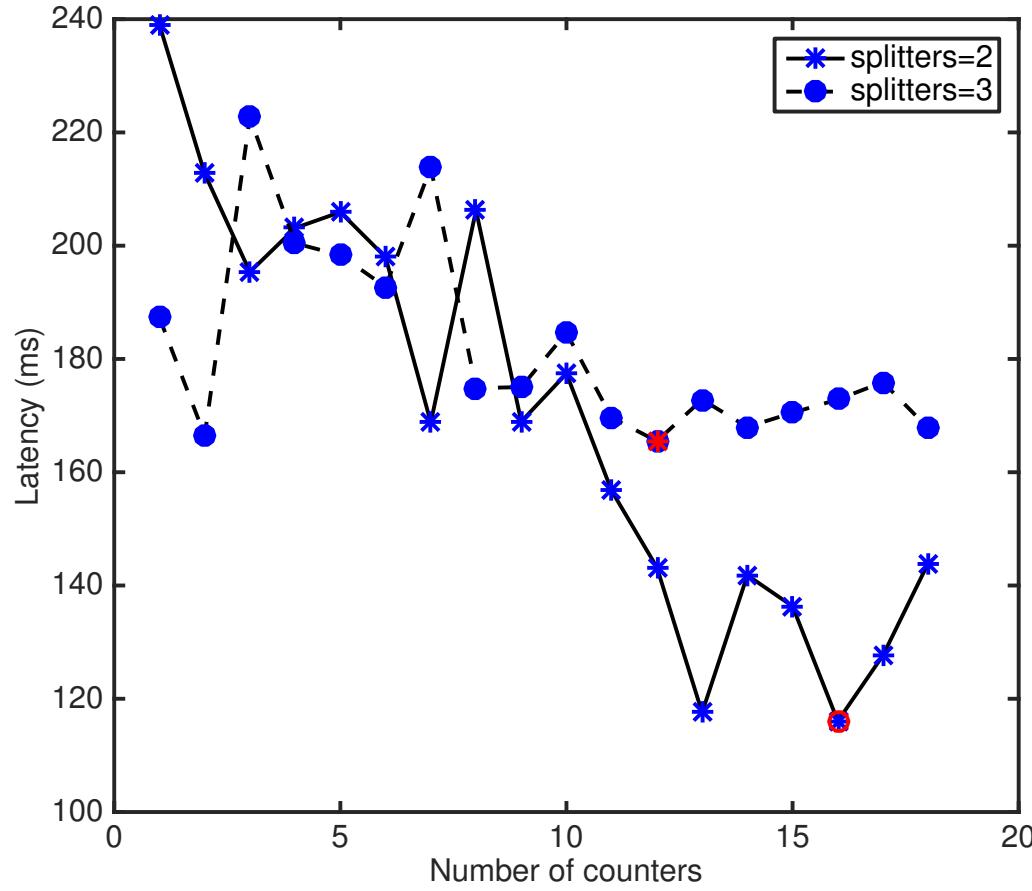
$y_i = f(\boldsymbol{x}_i), \boldsymbol{x}_i \subset \mathbb{X}$ Partially known

$y_i = f(\boldsymbol{x}_i) + \epsilon$ Measurements subject to noise

Non-linear interactions

Response surface is:

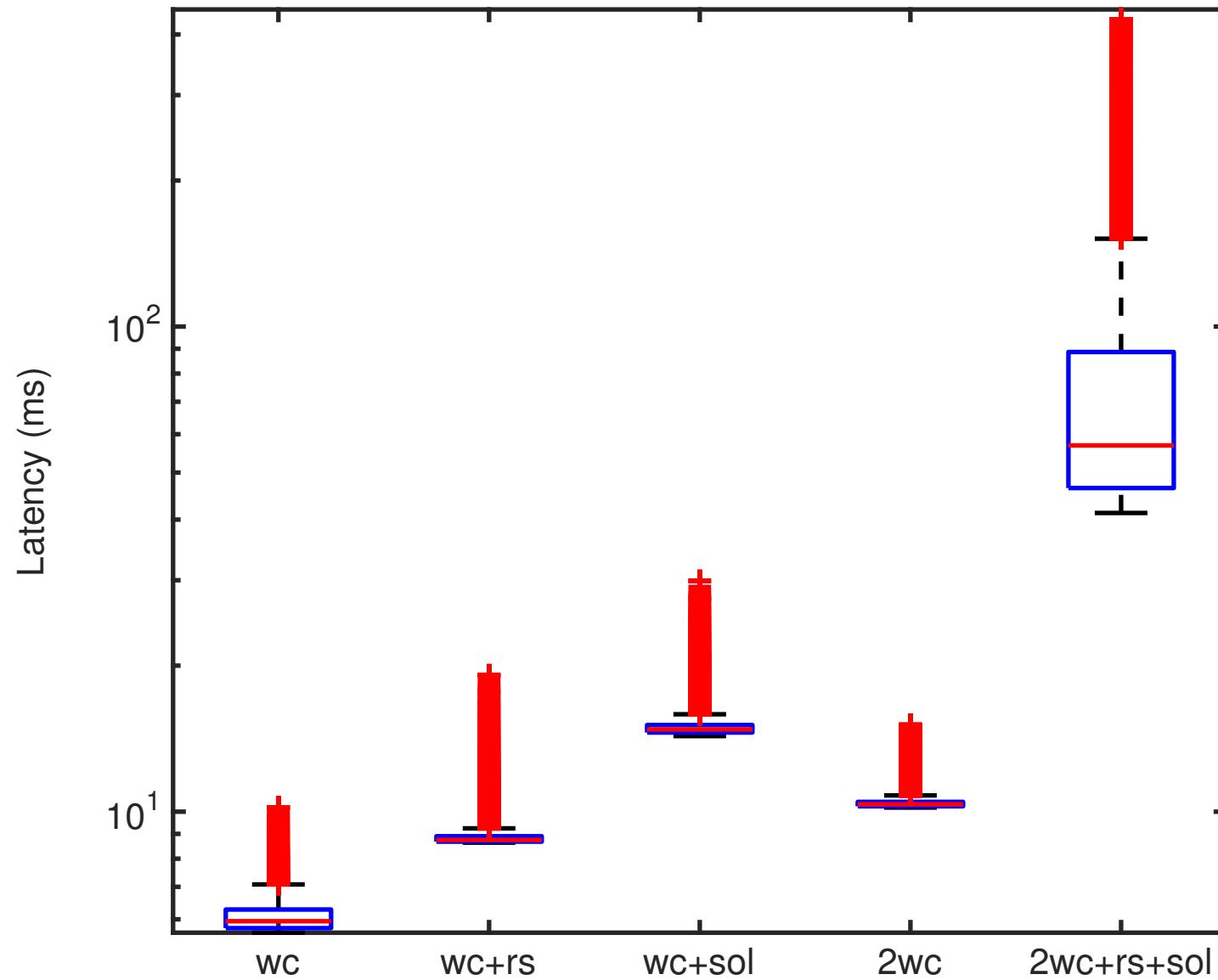
- Non-linear
- Non convex
- Multi-modal



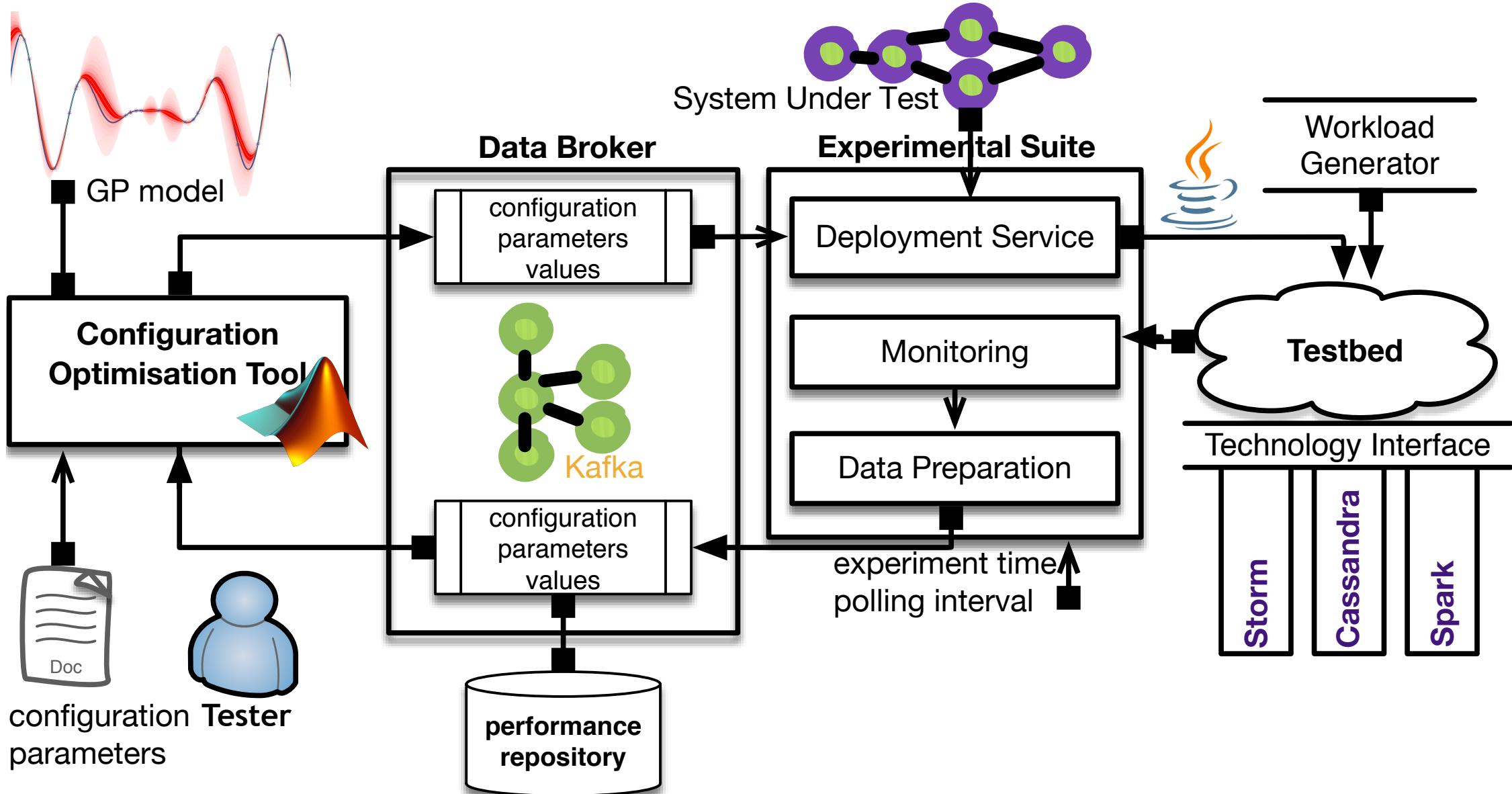
The measurements are subject to variability

The scale of measurement variability is different in different deployments (heteroscedastic noise)

$$y_i = f(x_i) + \epsilon_i$$



BO4CO architecture



GP for modeling blackbox response function

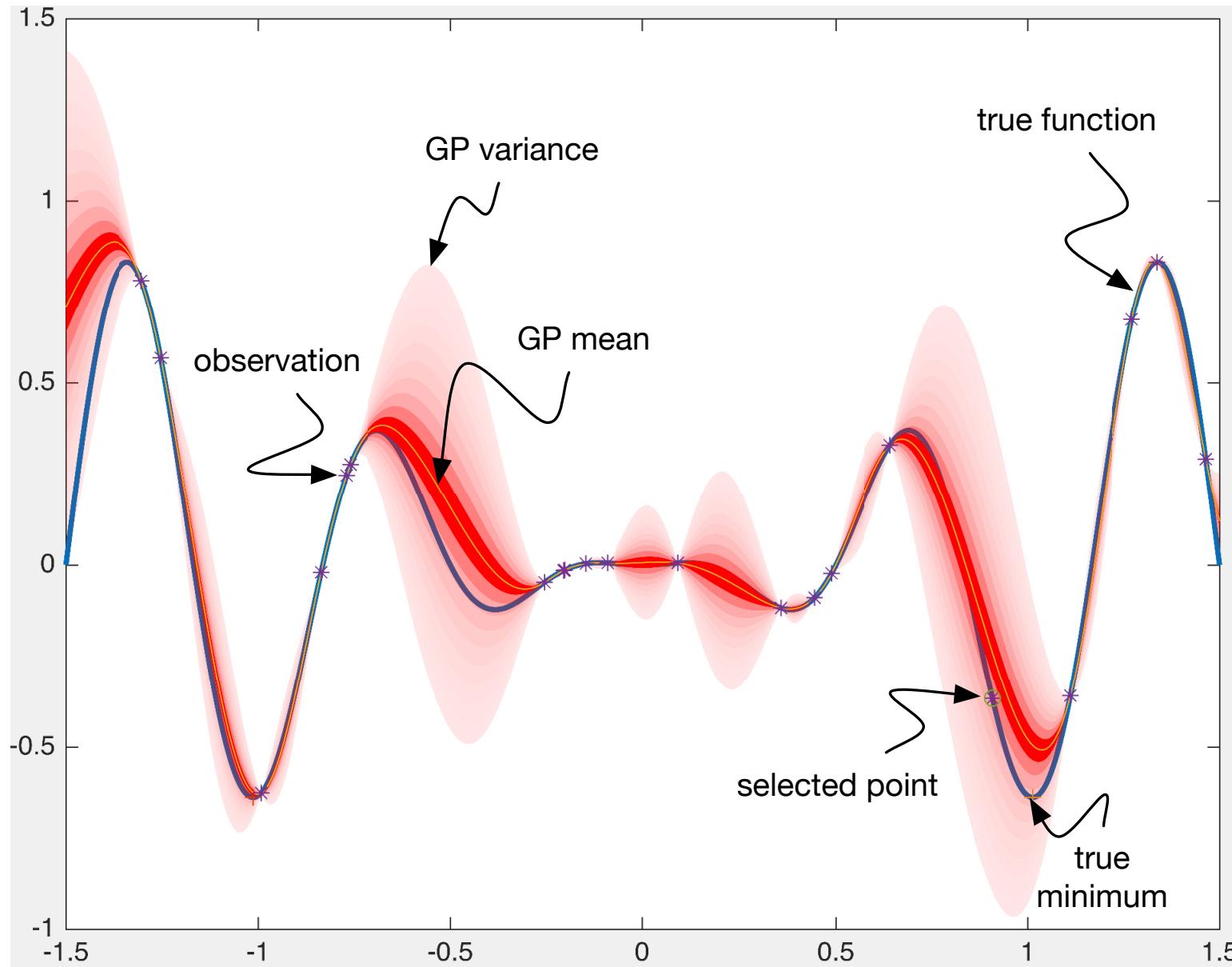
$$y = f(\mathbf{x}) \sim \mathcal{GP}(\mu(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')),$$

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

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

Motivations:

- 1- mean estimates + variance
- 2- all computations are linear algebra
- 3- good estimations when few data



Sparsity of Effects

- Correlation-based feature selector
- Merit is used to select subsets that are highly correlated with the response variable
- At most 2-3 parameters were strongly interacting with each other

	Topol.	Parameters	Main factors	Merit	Size	Testbed
1	wc(6D)	1-spouts, 2-max_spout, 3-spout_wait, 4-splitters, 5-counters, 6-netty_min_wait	{1, 2, 5}	0.787	2880	C1
2	sol(6D)	1-spouts, 2-max_spout, 3-top_level, 4-netty_min_wait, 5-message_size, 6-bolts	{1, 2, 3}	0.447	2866	C2
3	rs(6D)	1-spouts, 2-max_spout, 3-sorters, 4-emit_freq, 5-chunk_size, 6-message_size	{3}	0.385	3840	C3
4	wc(3D)	1-max_spout, 2-splitters, 3-counters	{1, 2}	0.480	756	C4
5	wc(5D)	1-spouts, 2-splitters, 3-counters, 4-buffer-size, 5-heap	{1}	0.851	1080	C5

Experiments on:

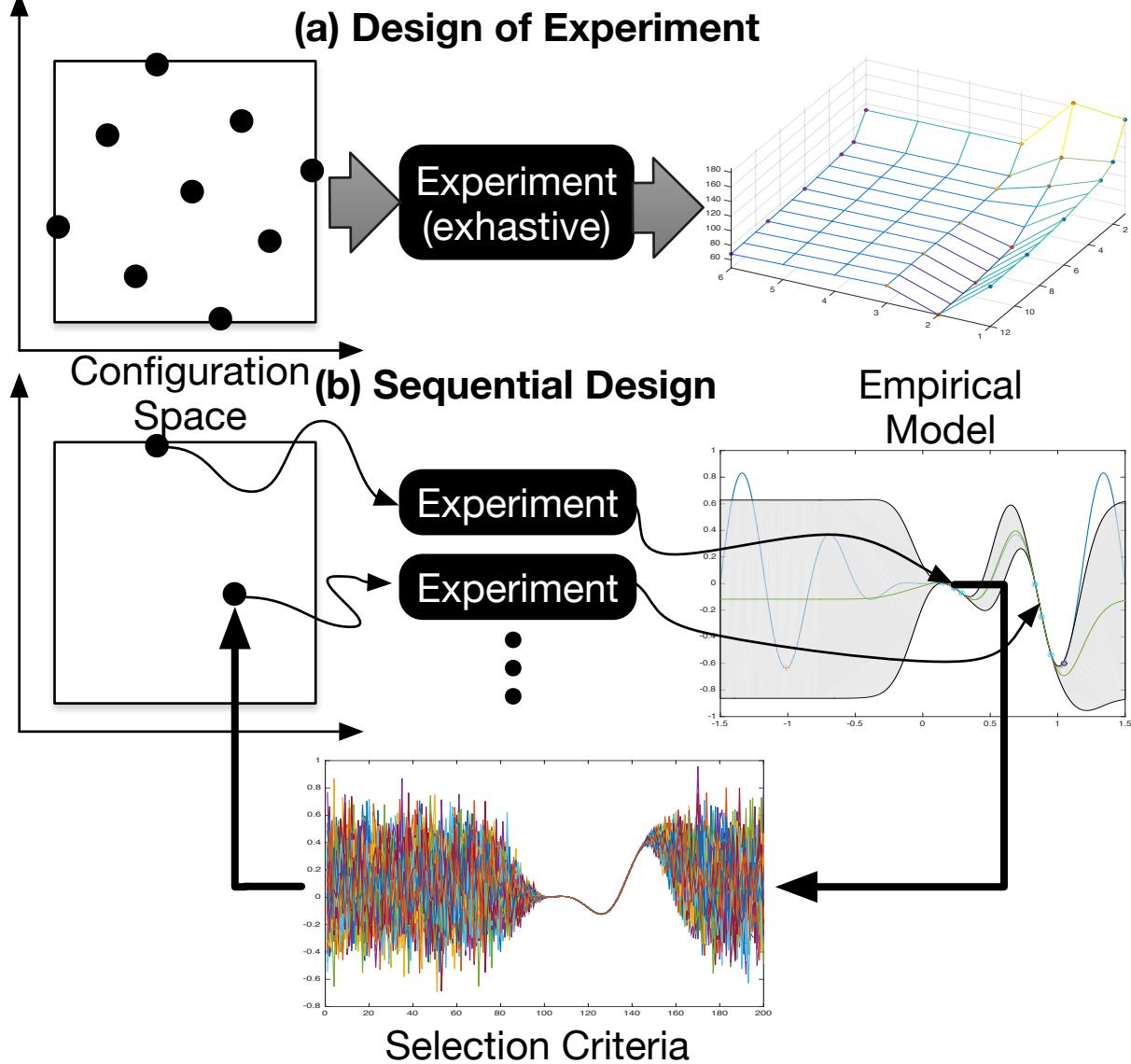
1. C1: OpenNebula (X)
2. C2: Amazon EC2 (Y)
3. C3: OpenNebula (3X)
4. C4: Amazon EC2 (2Y)
5. C5: Microsoft Azure (X)

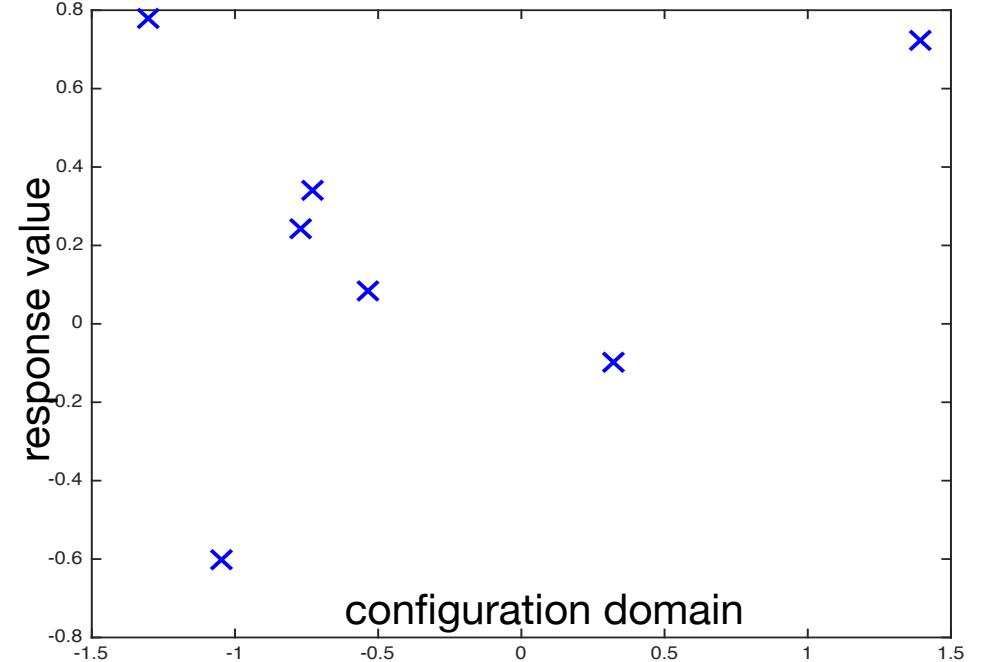
Algorithm 1 : BO4CO

Input: Configuration space \mathbb{X} , Maximum budget N_{max} , Response function f , Kernel function K_θ , Hyper-parameters θ , 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, $\mathbf{x}_t \leftarrow \arg \max_{\mathbf{x}} u(\mathbf{x}|\mathcal{M}, \mathbb{S}_{1:t-1})$ \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: $\mathcal{M}(\mathbf{x})$

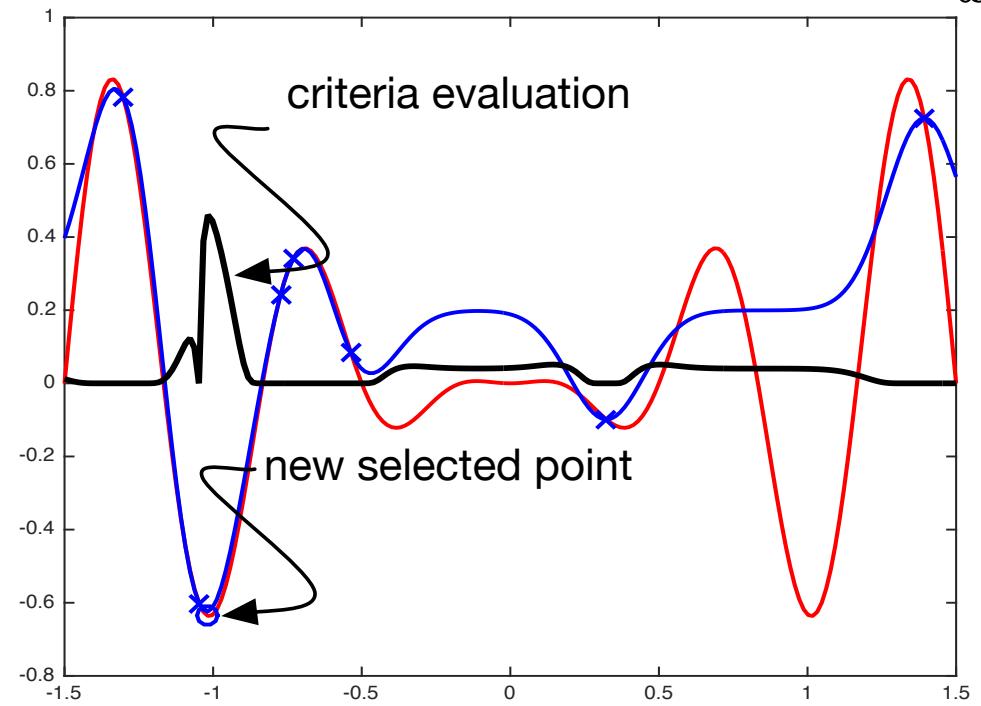




configuration domain

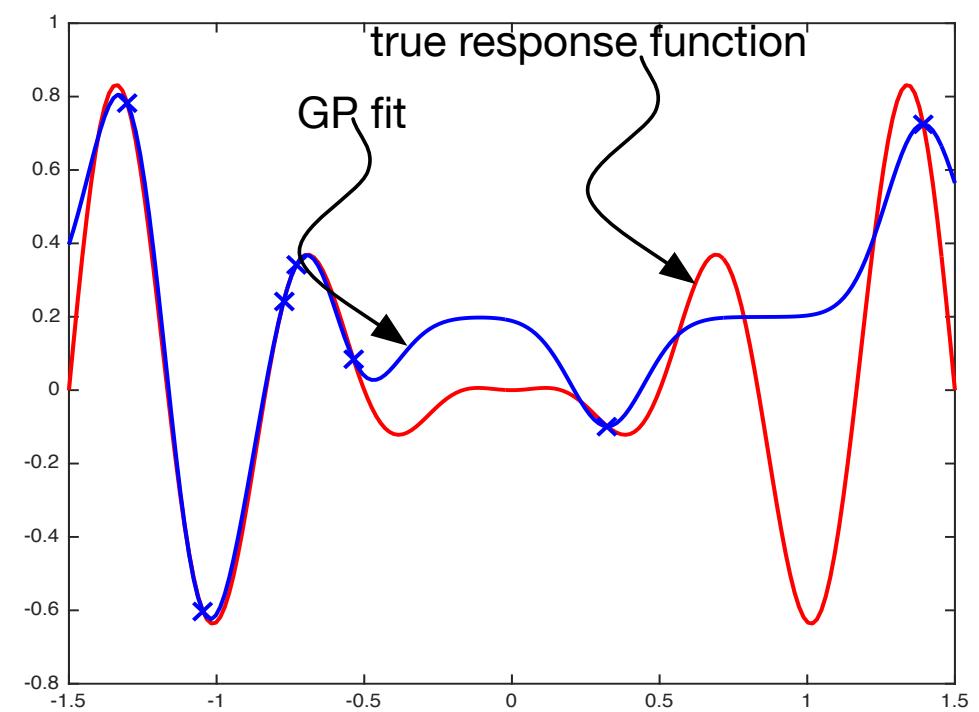
Acquisition function:

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



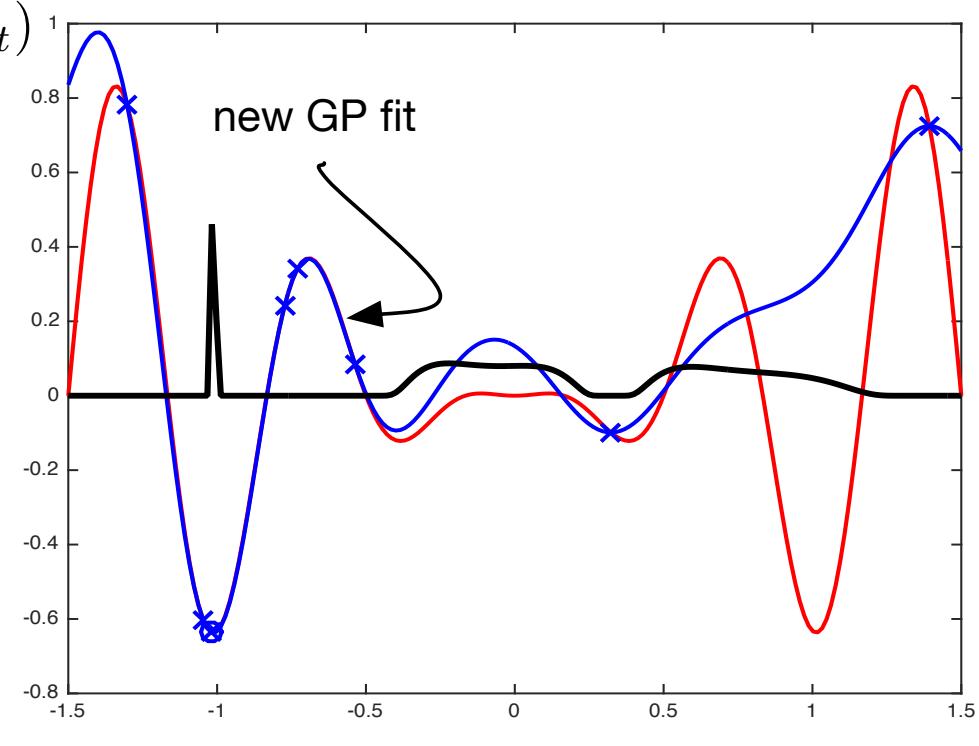
criteria evaluation

new selected point



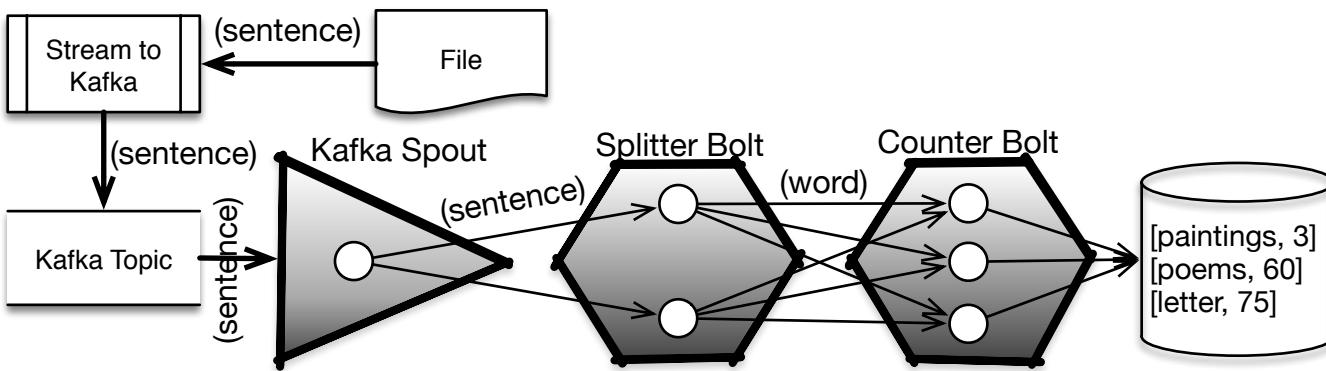
'true response function'

GP fit



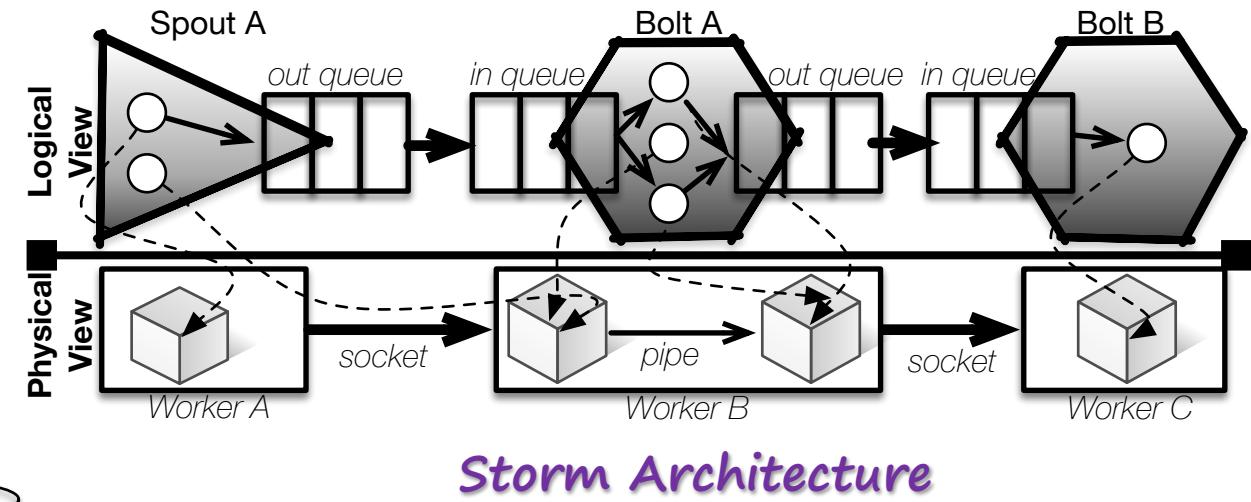
new GP fit

Stream Processing Systems

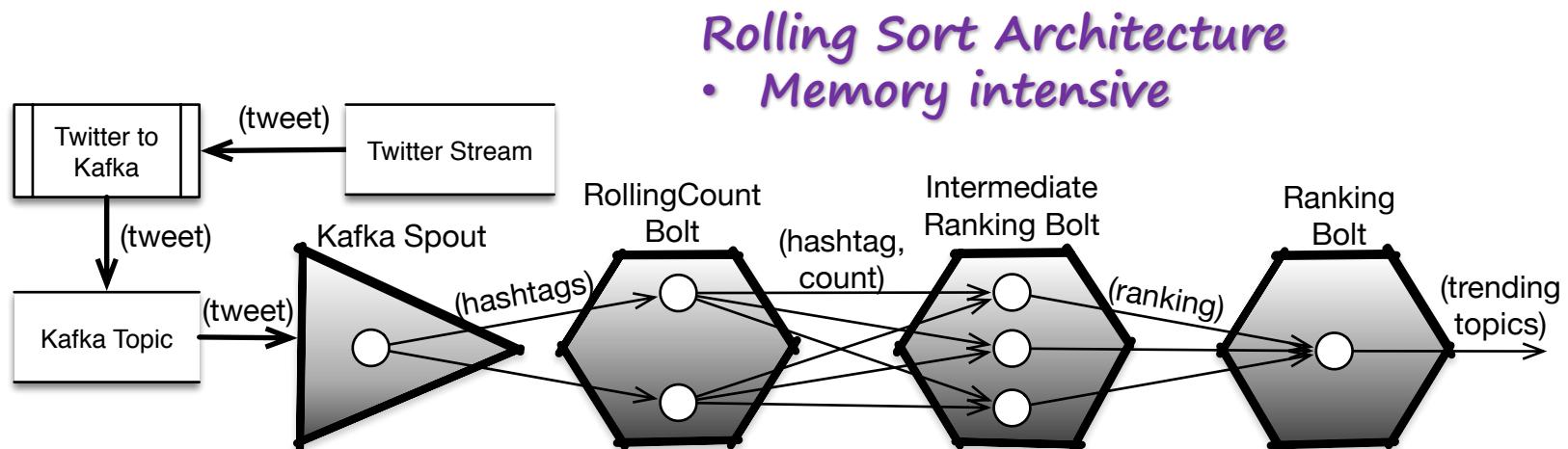


Word Count Architecture
• CPU intensive

Applications:
• Fraud detection
• Trending topics



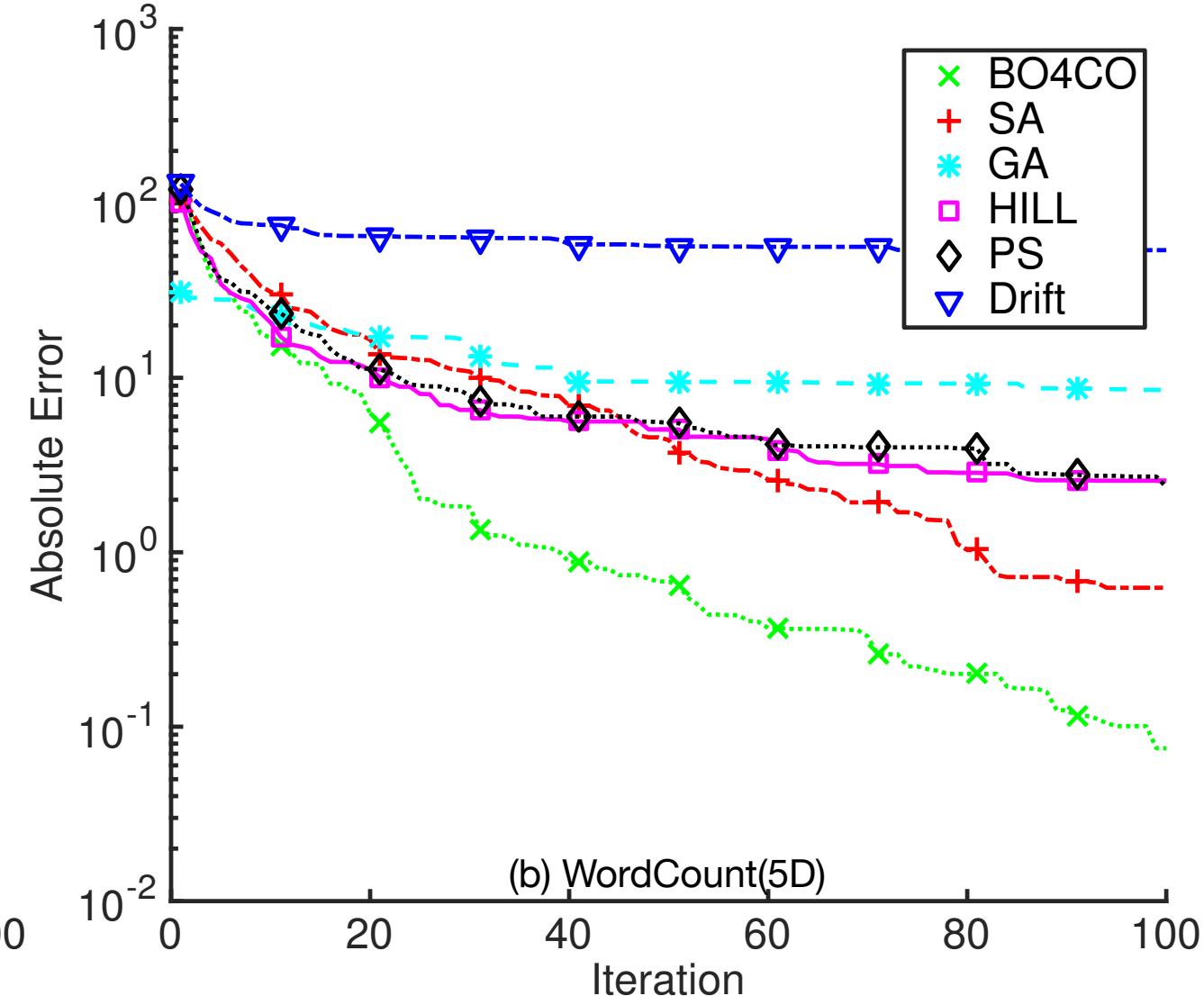
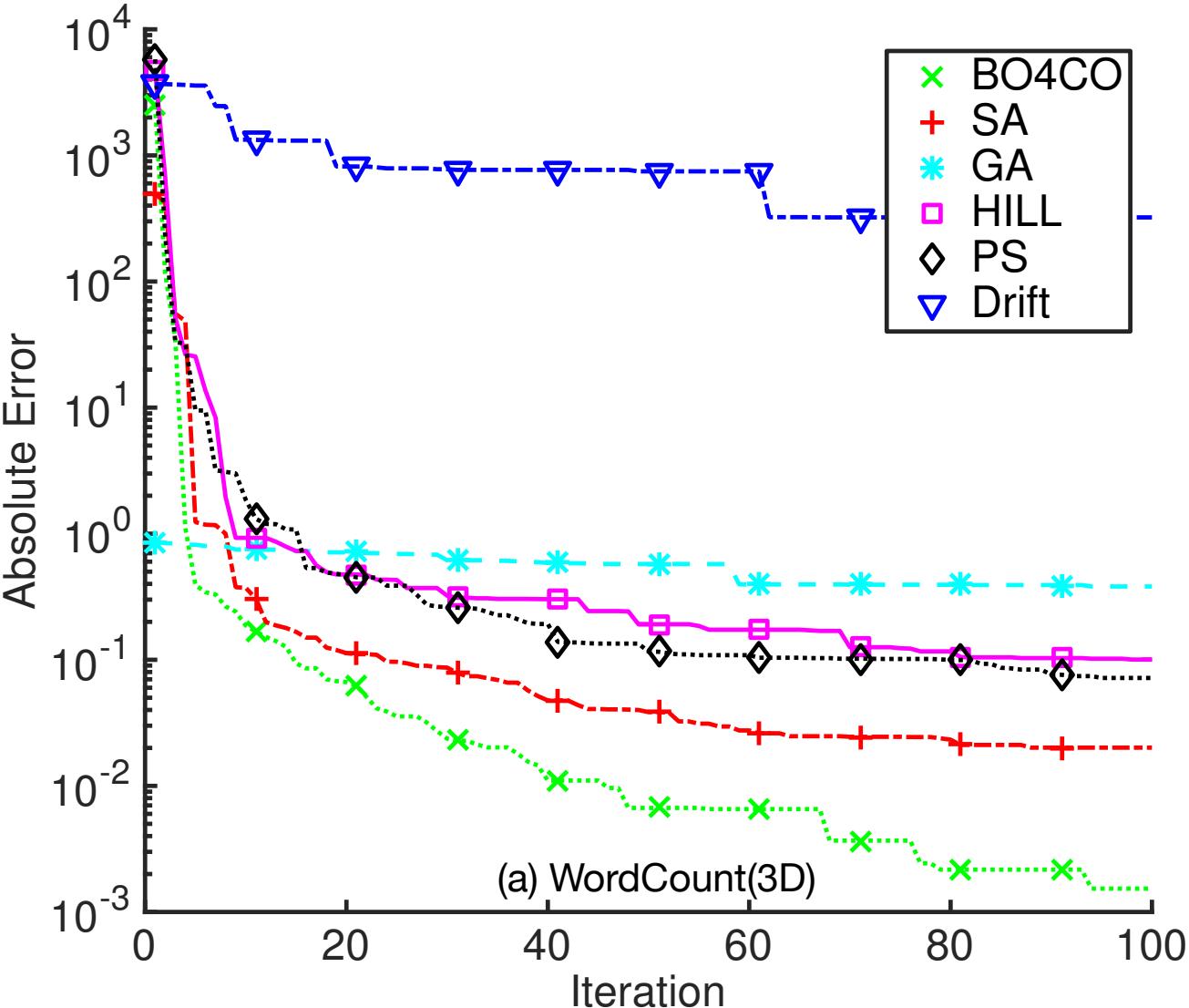
Storm Architecture



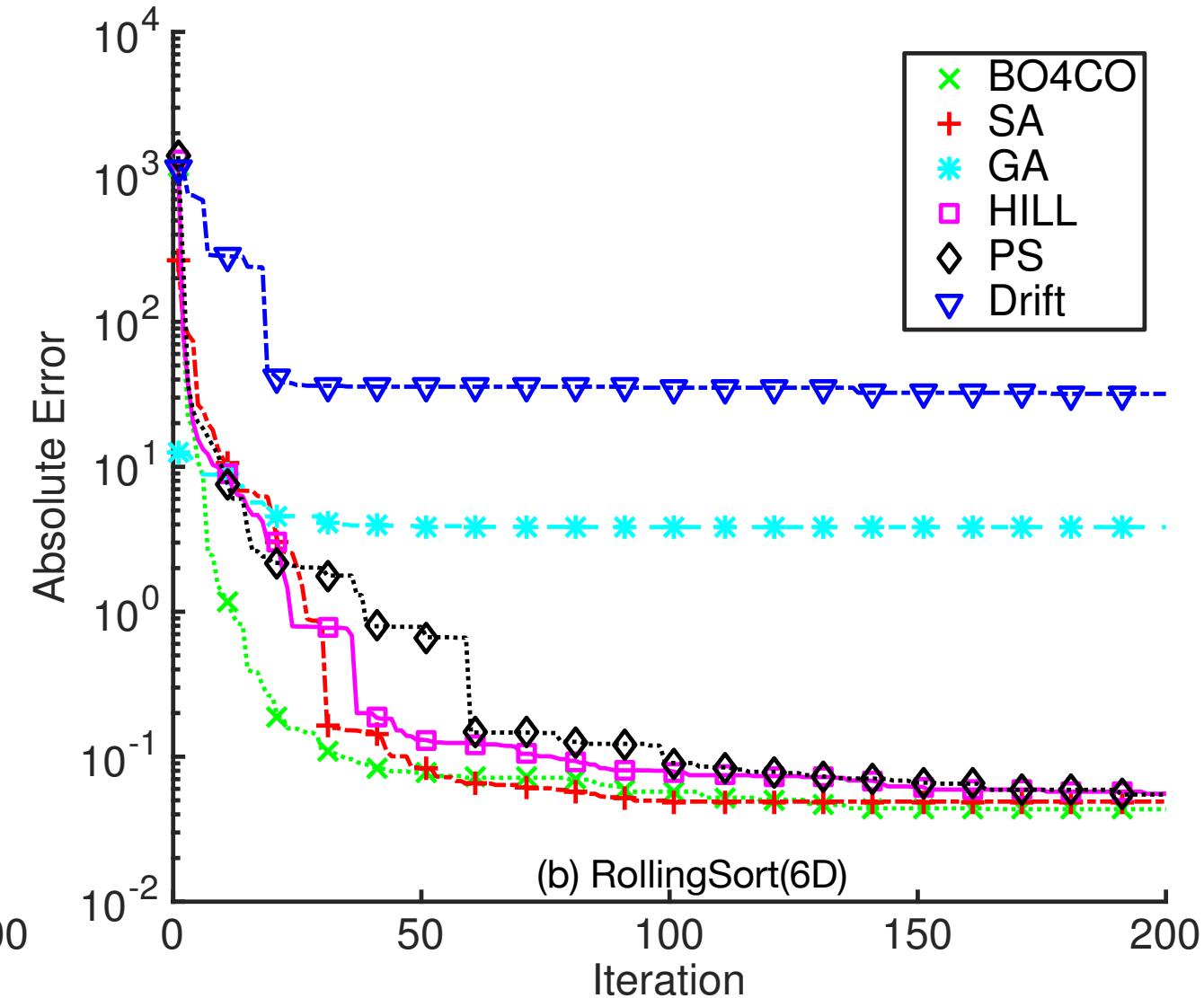
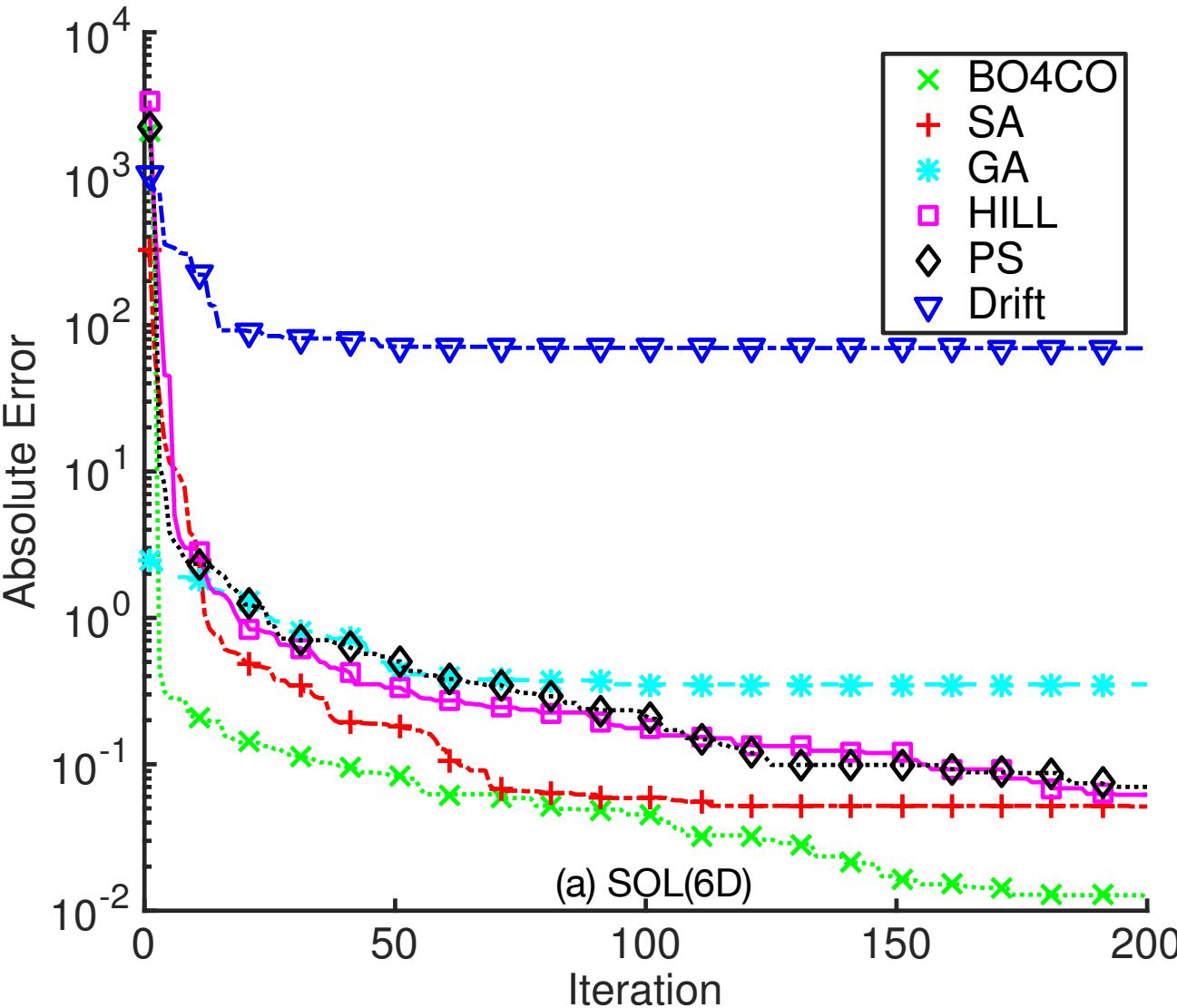
Rolling Sort Architecture
• Memory intensive

Experimental results

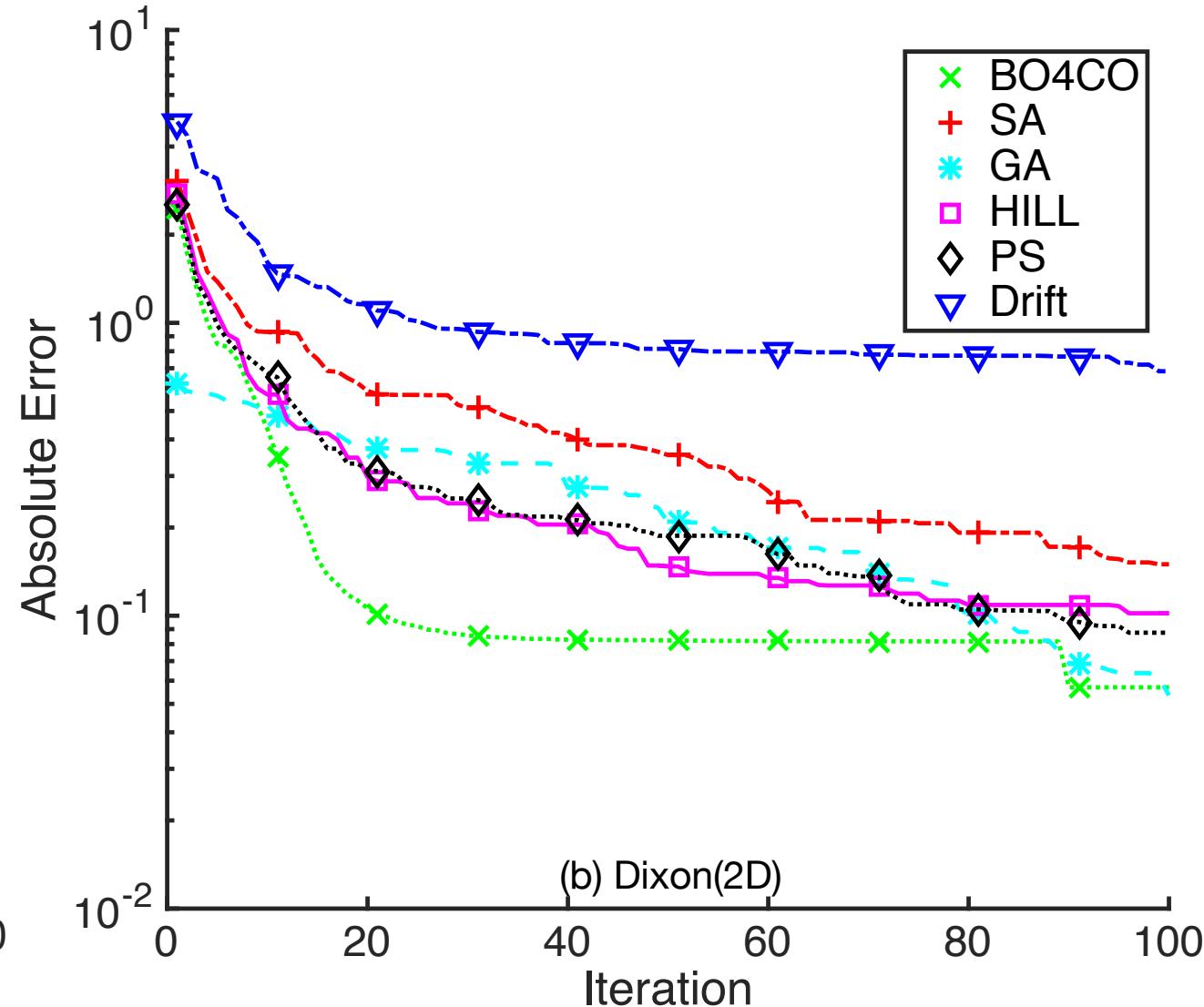
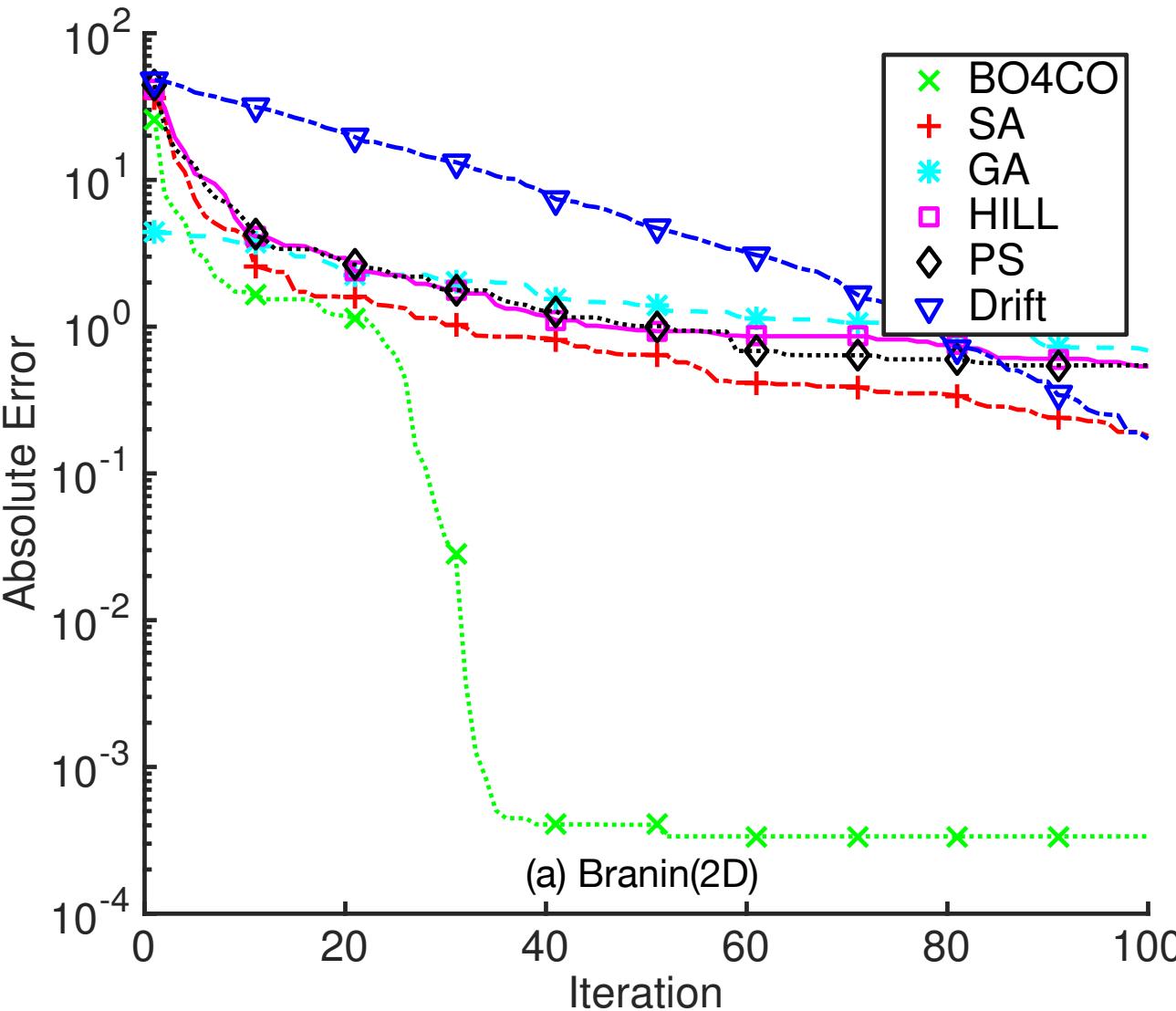
- 30 runs, report average performance
- Yes, we did full factorial measurements and we know where global min is...



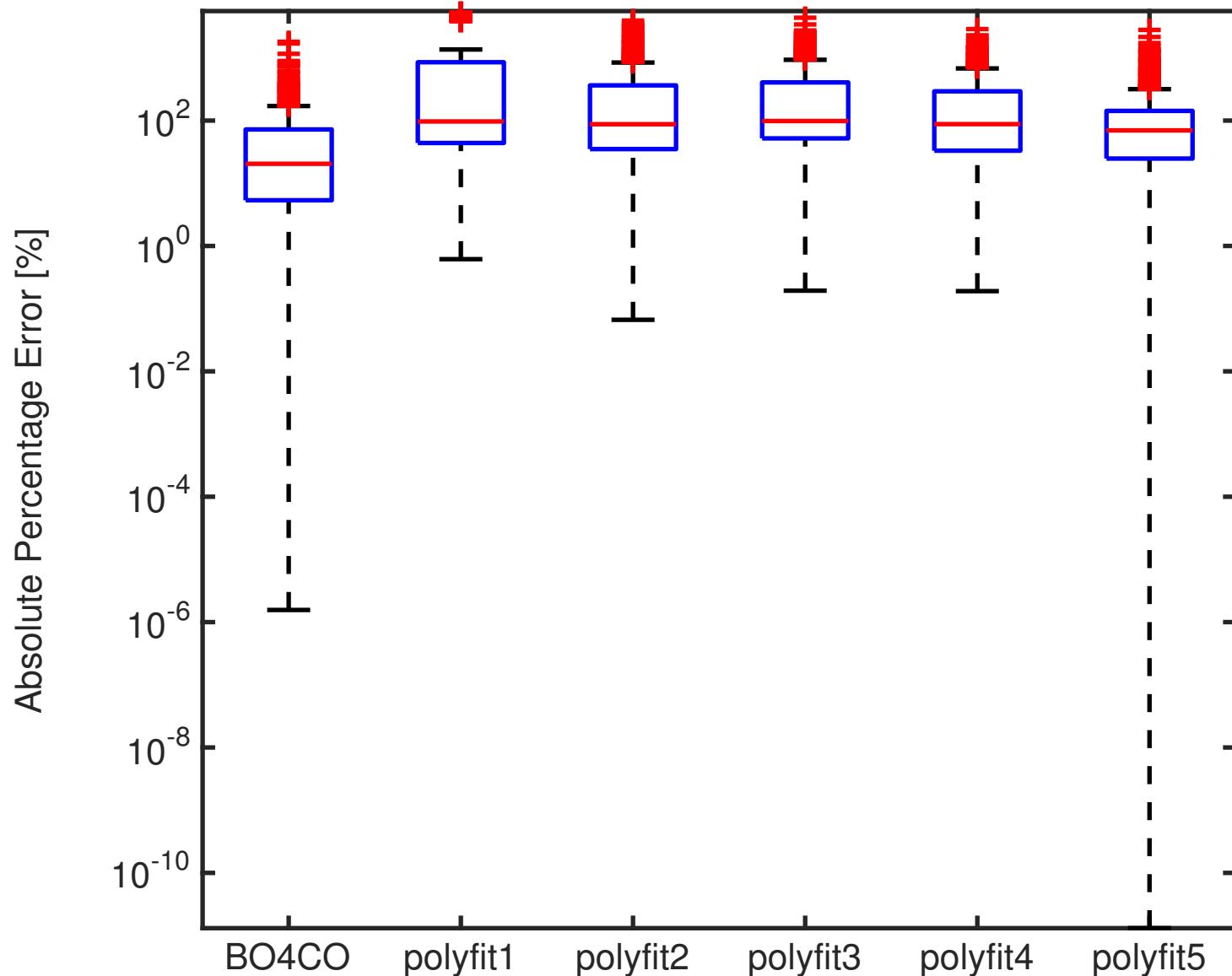
Experimental results



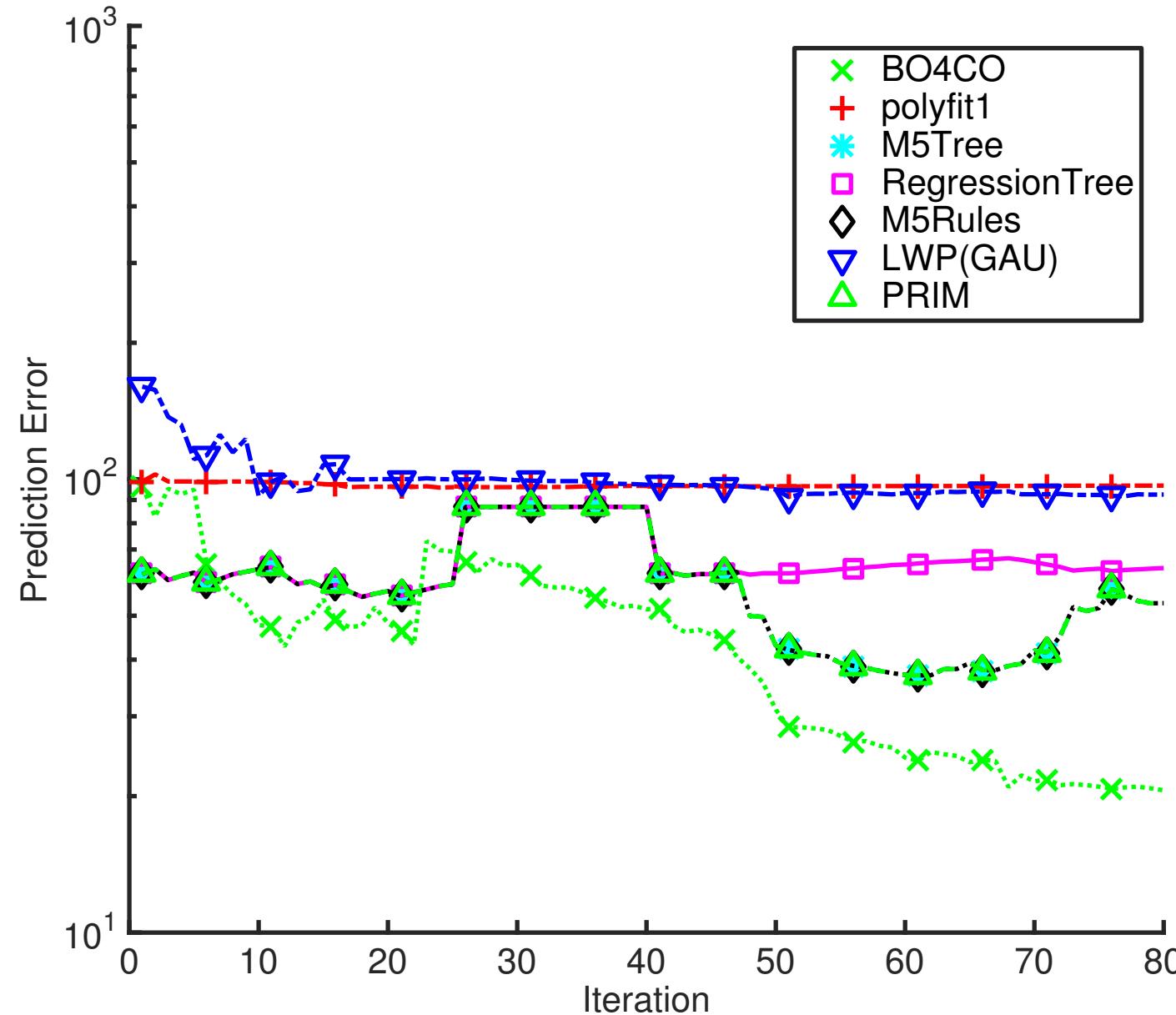
Experimental results



Model accuracy (comparison with polynomial regression models)



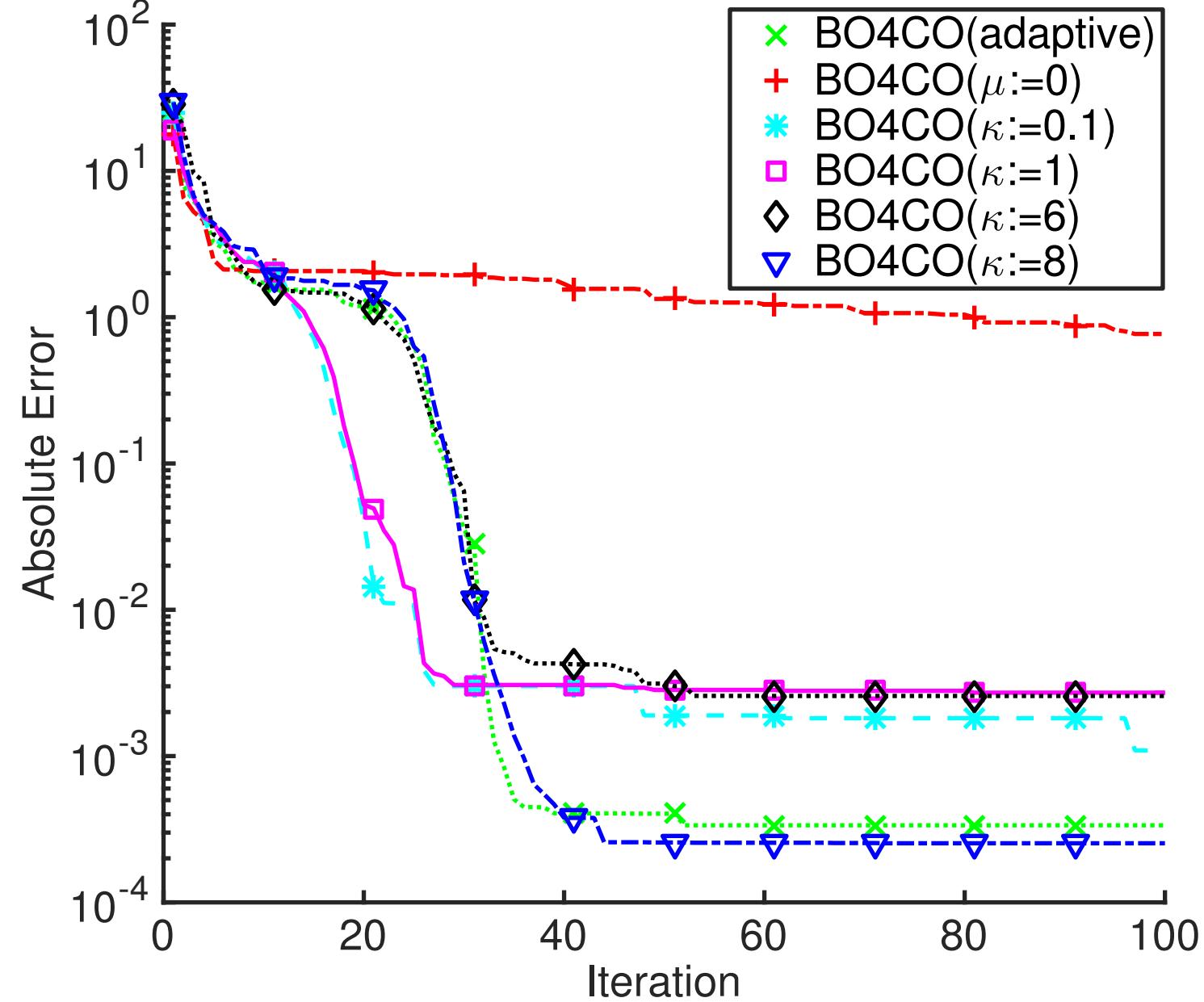
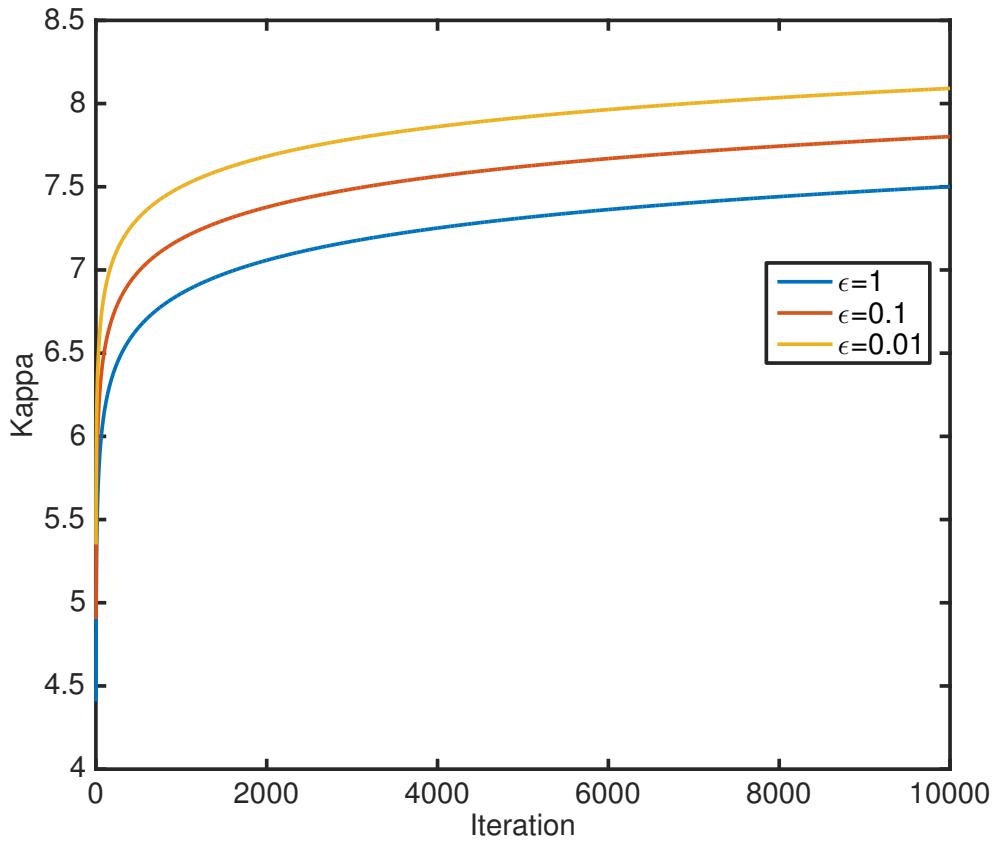
Prediction accuracy over time



Exploitation vs exploration

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

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



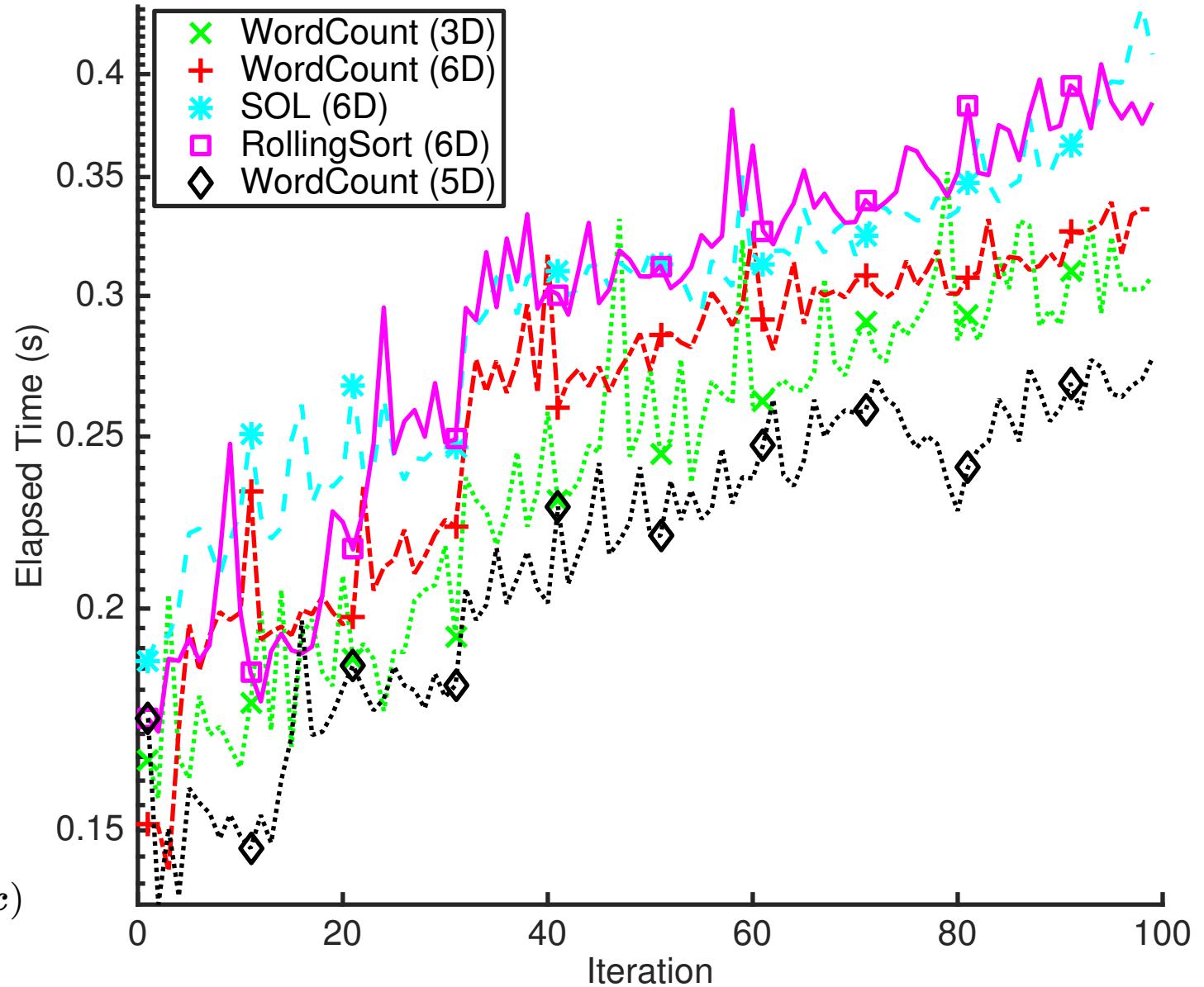
Runtime overhead

- The computation time in larger datasets is higher than those with less data and lower.
- The computation time increases over time since the matrix size for Cholesky inversion gets larger.

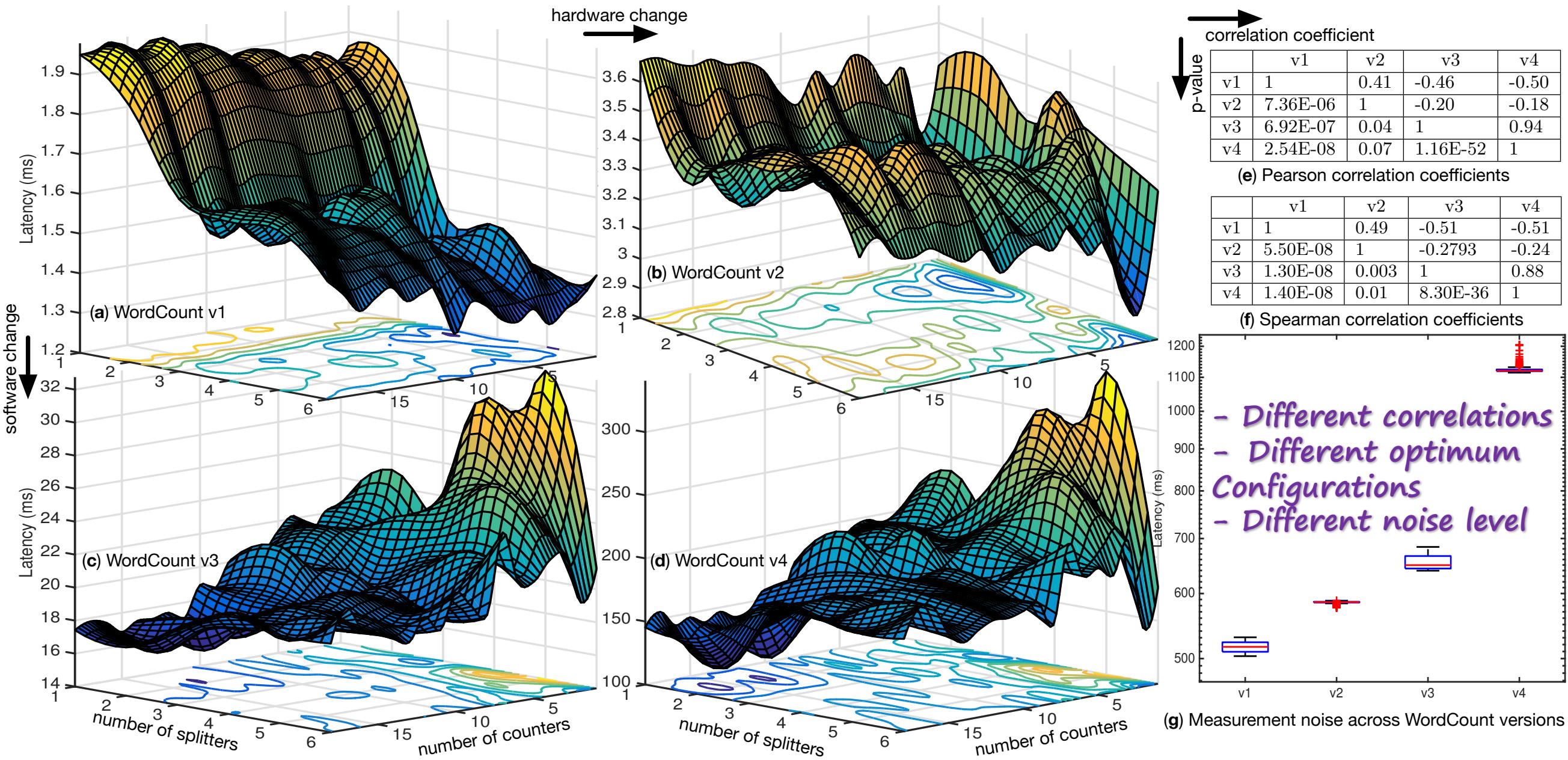
$$y = f(\mathbf{x}) \sim \mathcal{GP}(\mu(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')),$$

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

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



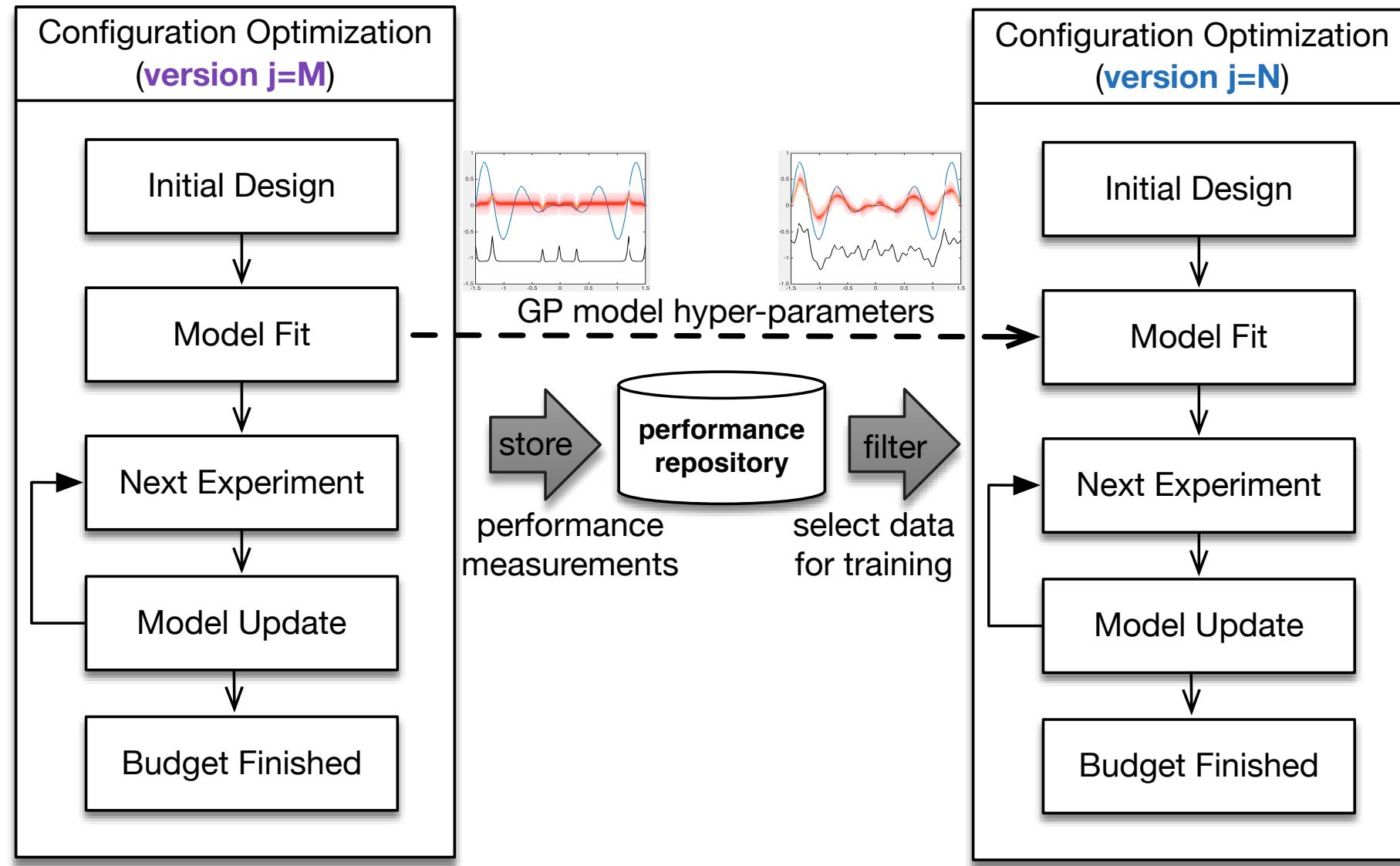
Correlations: SPS experiments



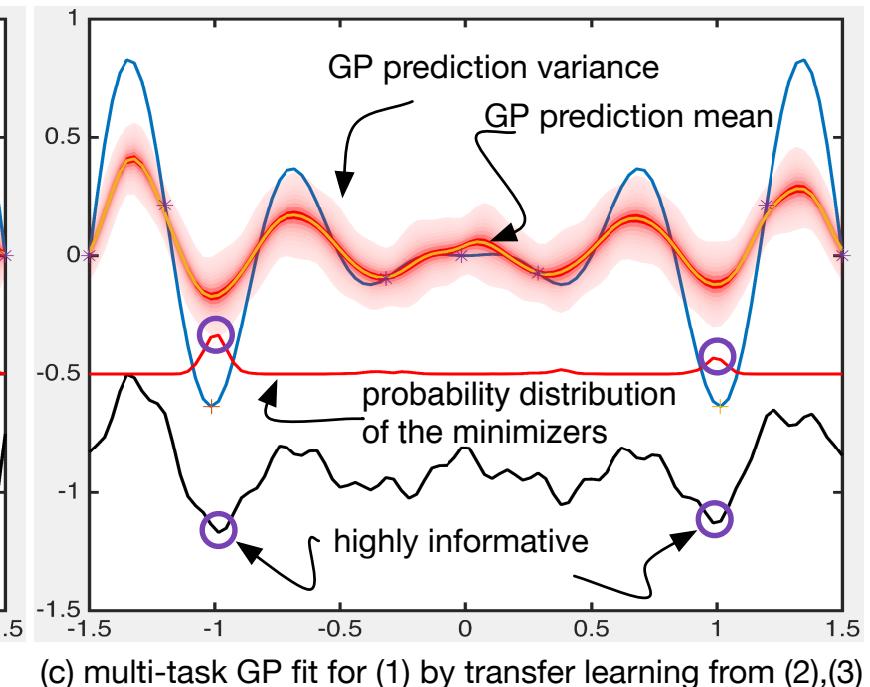
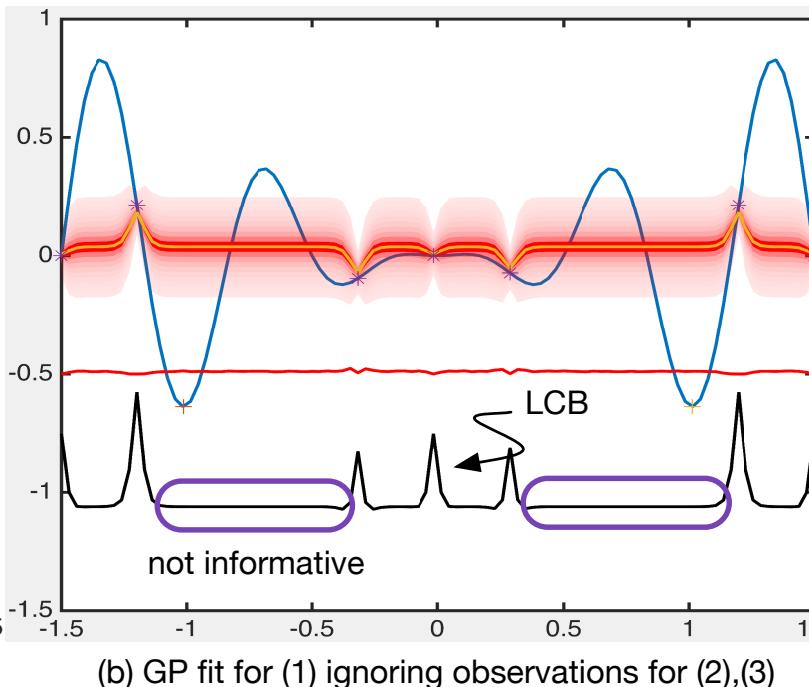
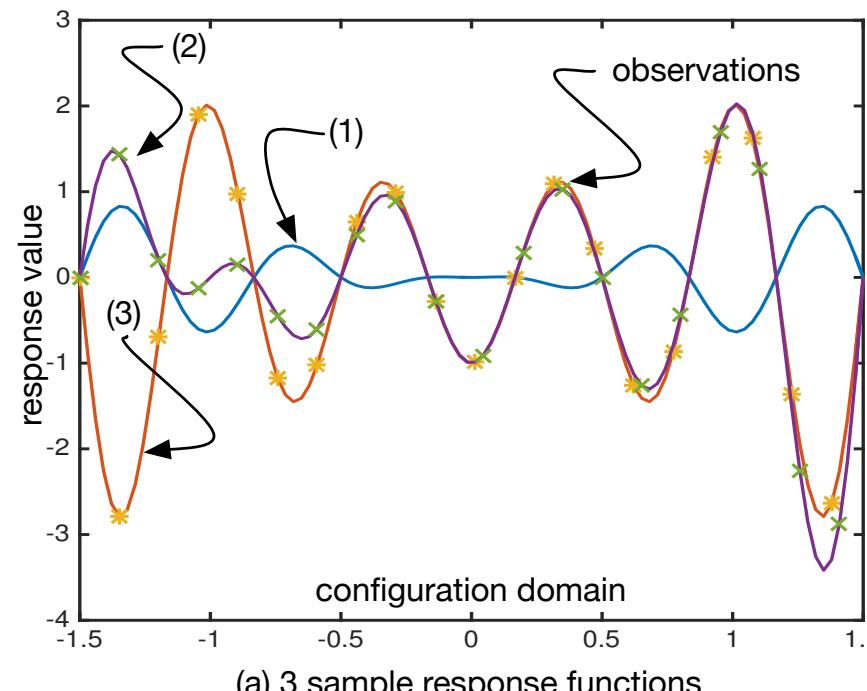
DevOps

- 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.

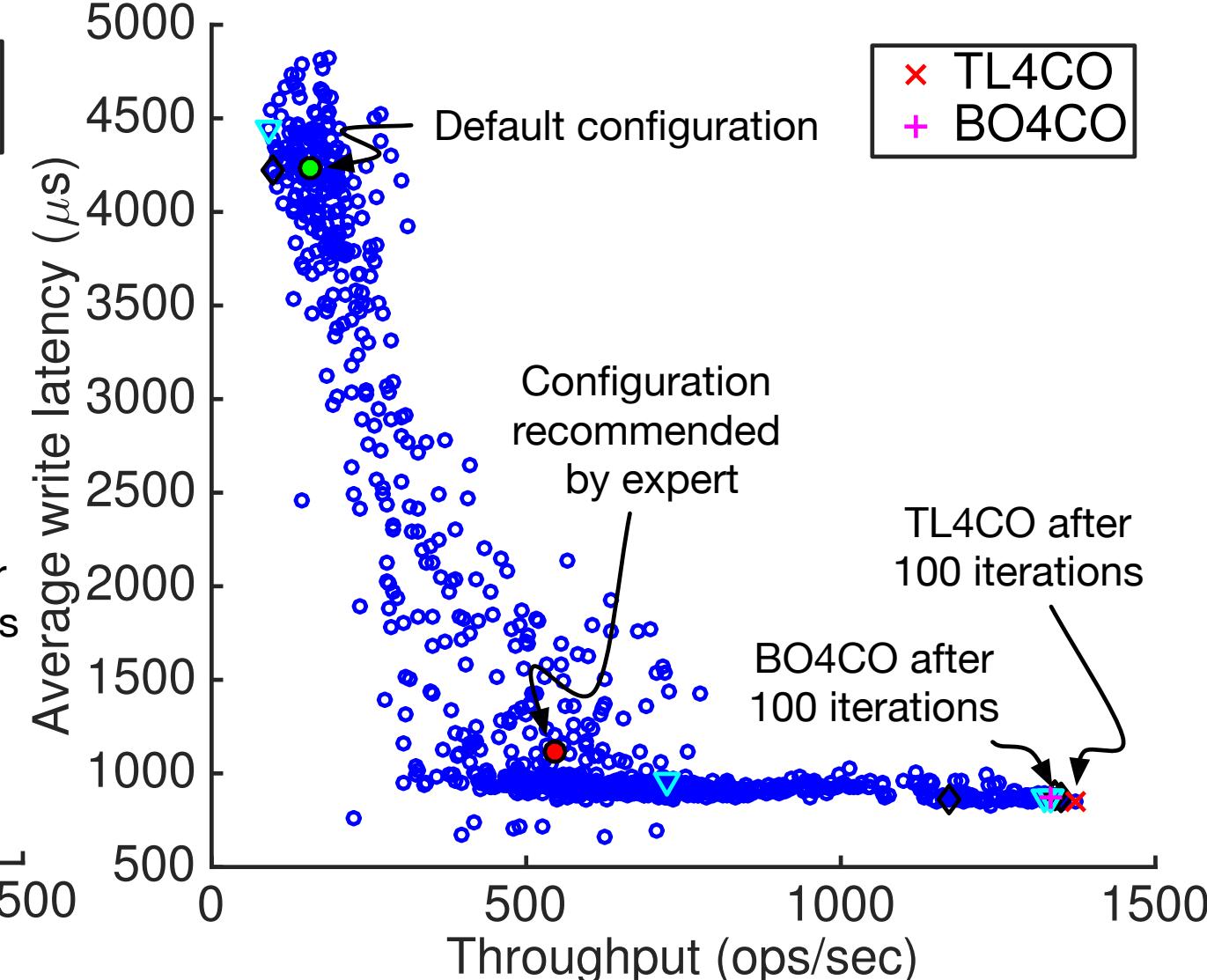
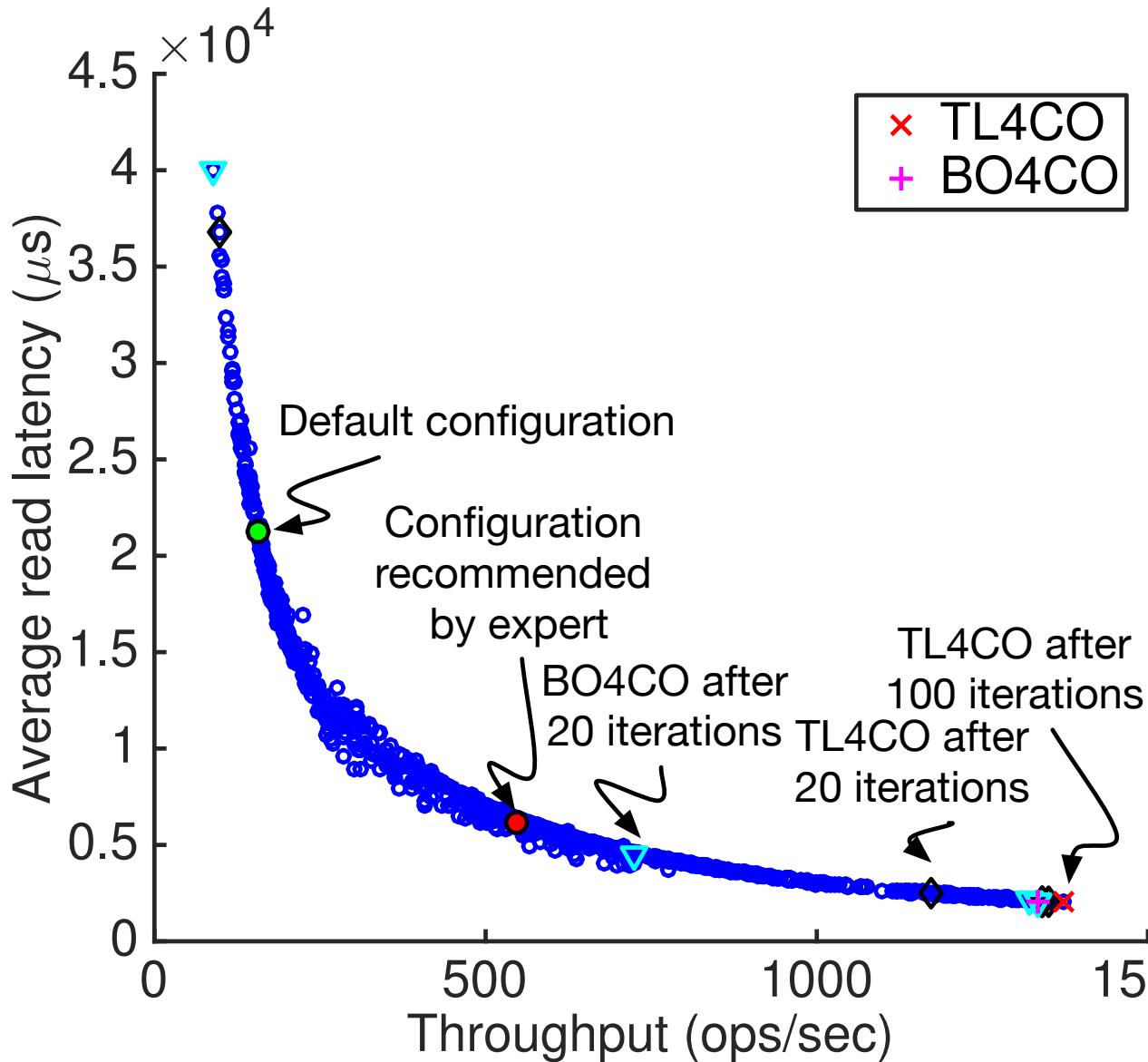
Solution: Transfer Learning for Configuration Optimization



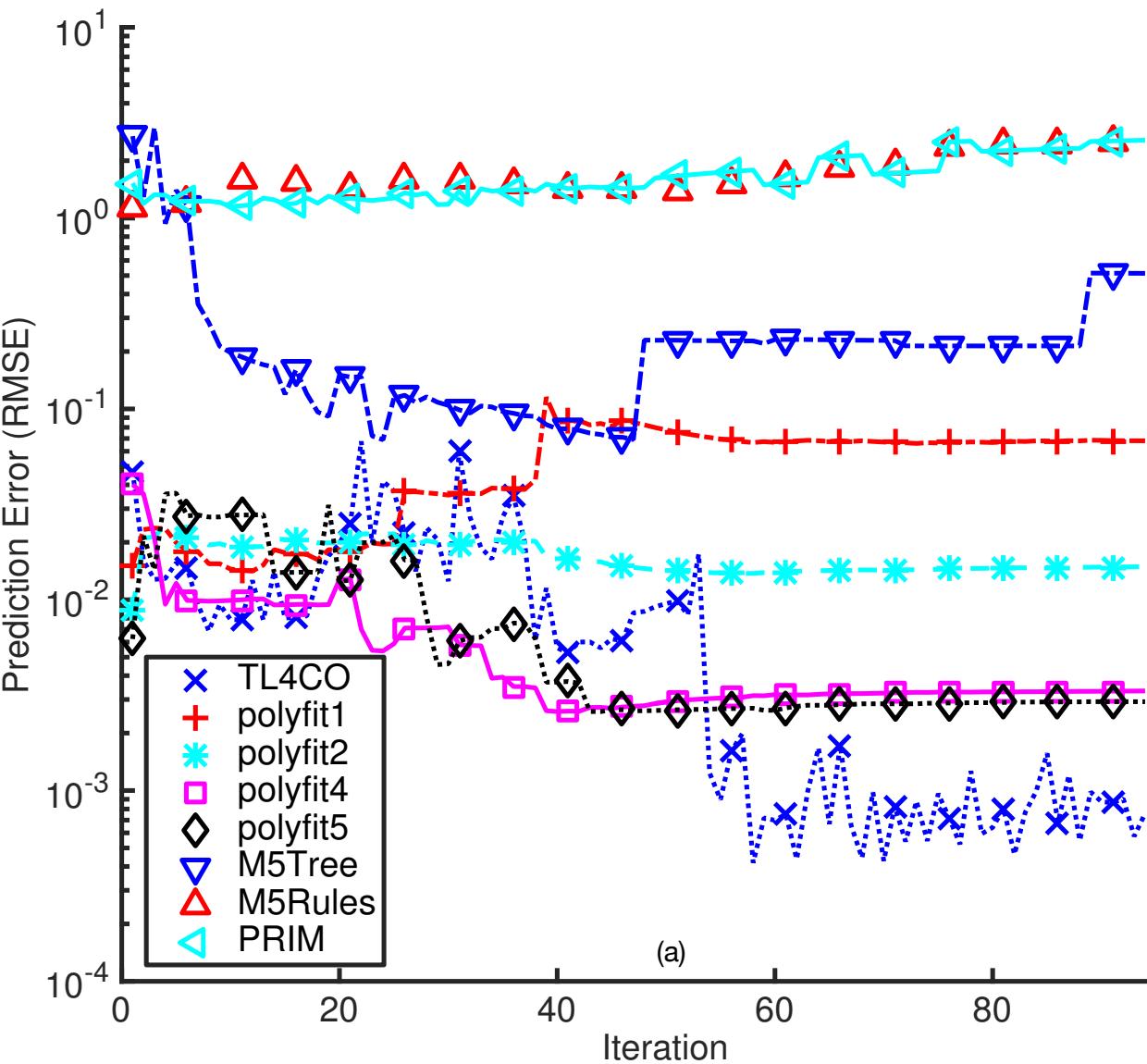
The case where we learn from correlated responses



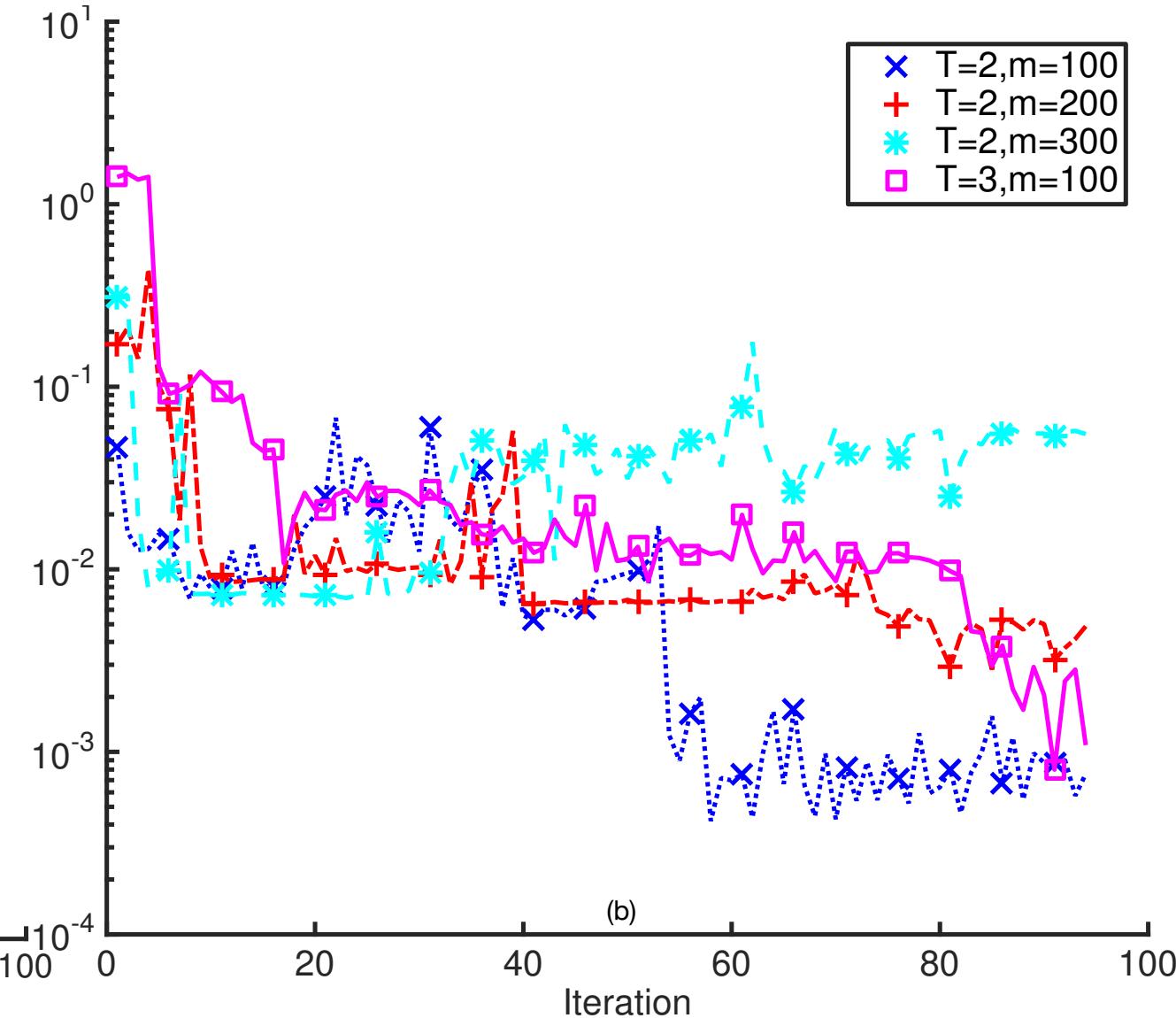
Comparison with default and expert prescription



Prediction accuracy over time



(a)



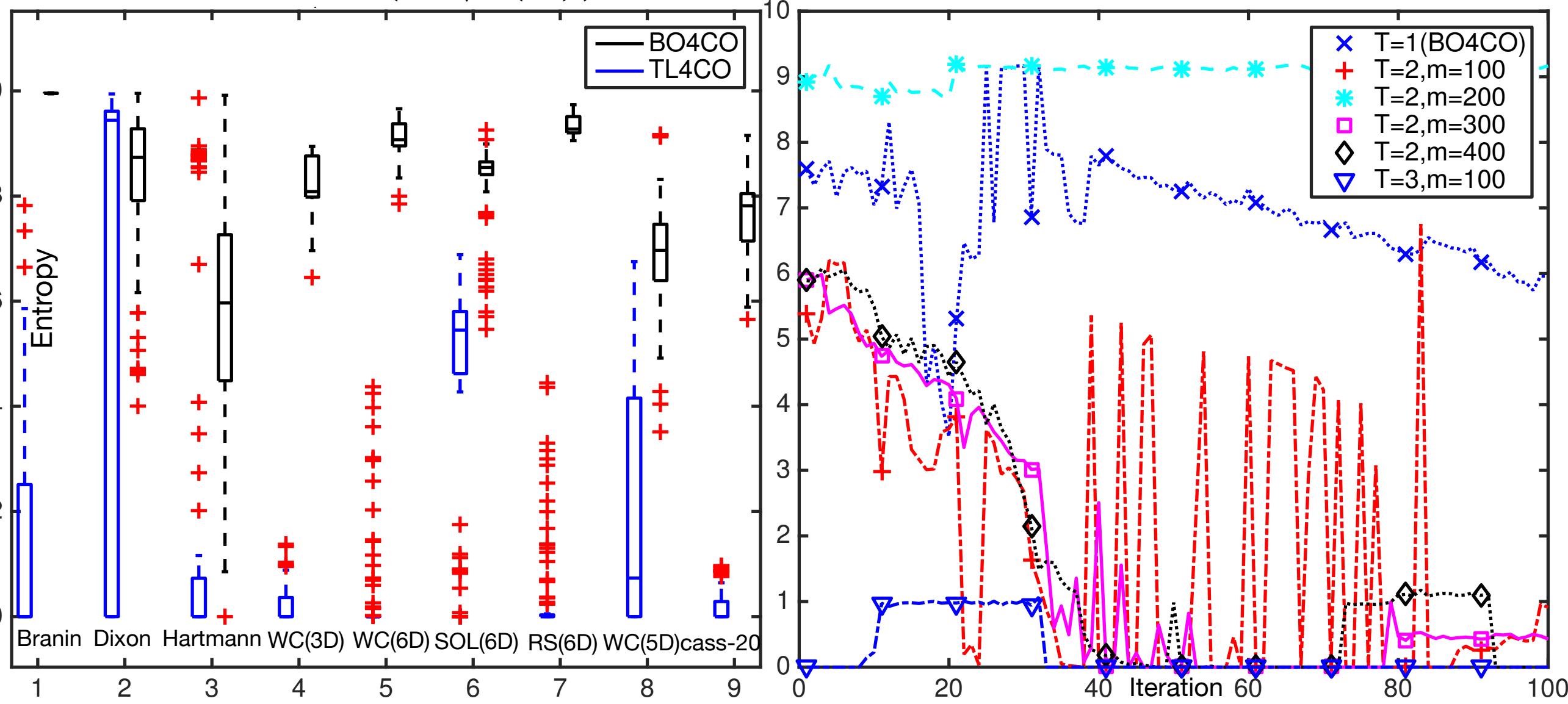
(b)

\times T=2,m=100
 $+$ T=2,m=200
 $*$ T=2,m=300
 \square T=3,m=100

Entropy of the density function of the minimizers

$$X^* = \Pr(x^* | f(x))$$

Knowledge about the location of the minimizer

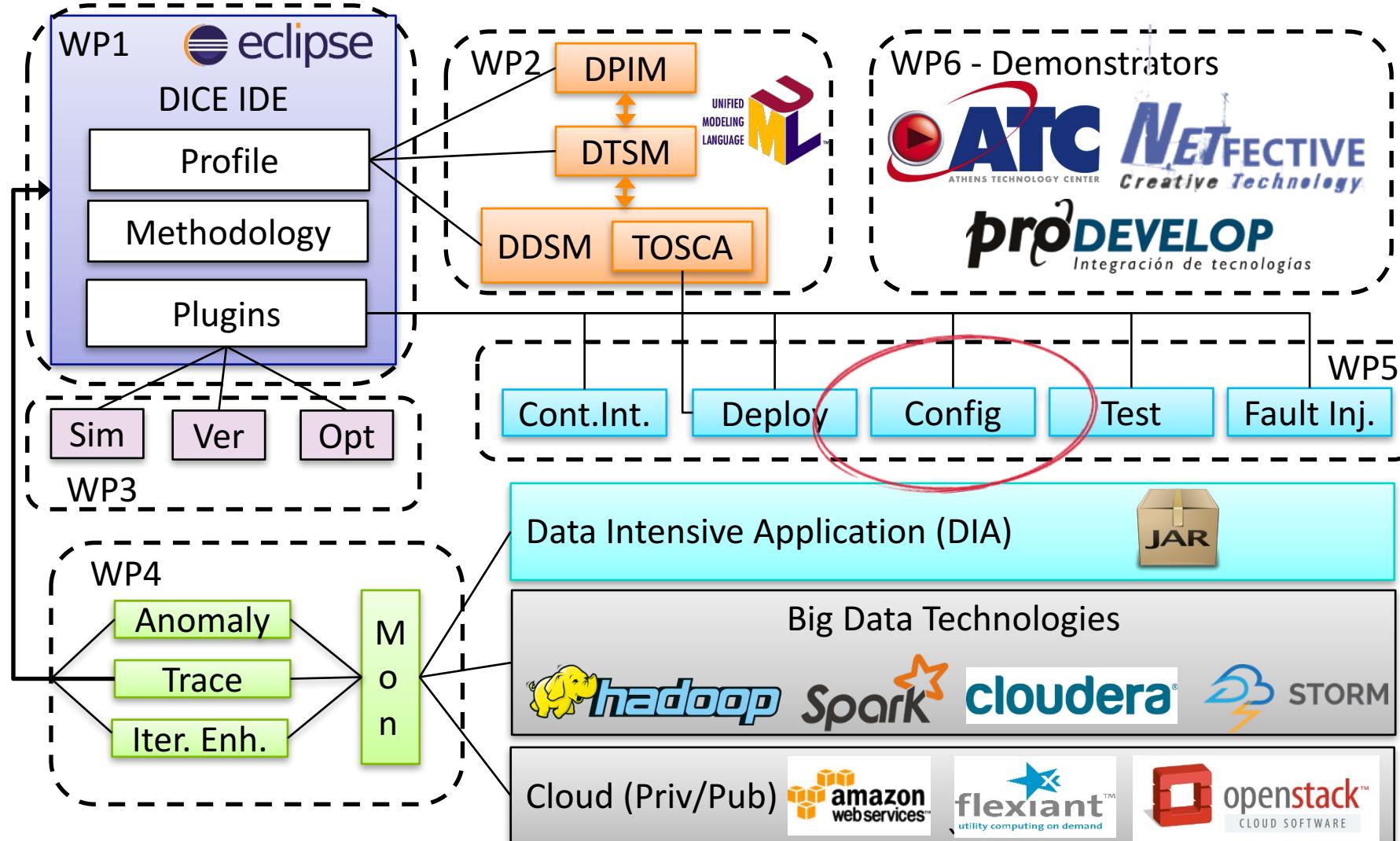


Takeaways

- Be aware of Uncertainty
 - By **quantifying** the uncertainty
 - Make decisions taking into account the right level of uncertainty (homoscedastic vs heteroscedastic)
 - Uncertainty sometimes helps (models that provide an estimation of the uncertainty are typically more **informative**)
 - By exploiting this knowledge you can only explore **interesting zones** rather than learning the whole performance function
- You can learn from operational data
 - Not only from the current version, but from **previous measurements** as well
 - Use the learning from past measurements as **prior knowledge**
 - **Too much data** can be also harmful, it would slow down or blur the proper learning (**negative transfer**)

Acknowledgement:

- BO4CO as a part of DevOps pipeline in H2020 DICE
- BO4CO is being acquired by TATA (TCS)



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

Submit to SEAMS 2017

- Any work on Self-*
- Abstract Submission: 6 Jan, 2017 (firm)
- Paper Submission: 13 Jan, 2017 (firm)
- Page limit:
 - Long: 10+2,
 - New ideas and tools: 6+1
- More info: <https://wp.doc.ic.ac.uk/seams2017/>
- Symposium: 22-23 May, 2017
- We accept artifacts submissions (tool, data, model)

12th International Symposium on Software Engineering for Adaptive and Self-Managing Systems
Buenos Aires, Argentina, May 22-23, 2017, <http://wp.doc.ic.ac.uk/seams2017>



Co-located with



General Chair

David Garlan, USA

Program Chair

Bashar Nuseibeh, UK & Ireland

Artifacts Chair

Javier Cámará, USA

Publicity Chair

Pooyan Jamshidi, UK

Local Chair

Nicolas D'Ippolito, Argentina

Program Committee

Dalal Alraieh, UK

Jesper Andersson, Sweden

Rami Bahsoon, UK

Arosa Bandara, UK

Luciano Baresi, Italy

Jacob Beal, USA

Nelly Bencomo, UK

Amel Bennaceur, UK

Victor Gómezberumen, Argentina

Tomas Bures, Czech Republic

Radu Calinescu, UK

Javier Cámará, USA

Betty Cheng, USA

Siobhán Clarke, Ireland

Rogério de Lemos, UK

Elisabetta Di Nitto, Italy

Nicolás D'Ippolito, Argentina

Ada Diaconescu, France

Gregor Engels, Germany

Antonio Filleri, UK

Erik Fredericks, USA

Holger Giese, Germany

Hassan Gomaa, USA

Joel Greenyer, Germany

Mark Harman, UK

Valérie Issarny, France

Pooyan Jamshidi, UK

Jean-Marc Jézéquel, France

Samuel Kounev, Germany

Philippe Lalanda, France

Seok-Won Lee, South Korea

Mari Litou, Canada

Xiaoming Ma, China

Martina Maggio, Sweden

Sam Malek, USA

Nenad Medvidovic, USA

Hausi Seebach, Canada

Henry Muccini, Italy

John Mylopoulos, Canada

Ingrid Nunes, Brazil

Liliana Pasquale, Ireland

Patrizio Pelliccione, Sweden

Xin Peng, China

David Rosenblum, Singapore

Bradley Schmerl, USA

Hella Seebach, Germany

Amir Molzahn Sharifloo, Germany

Vitor Silva Sousa, Brazil

Jan-Philipp Steghöfer, Sweden

Ladan Tahvildari, Canada

Kenji Tei, Japan

Axel van Lamsweerde, Belgium

Giuseppe Valetto, Italy

Mirko Viroli, Italy

Danny Weijns, Belgium

Yijun Yu, UK

Important Dates:

Abstract Submission: 6 Jan, 2017 (AoE,firm)

Paper Submission: 13 Jan, 2017 (AoE,firm)

Notification: 21 February, 2017

Camera ready: 6 Mar, 2017

*There will be a specific session to be dedicated to artifacts that may be useful for the community as a whole. Please see <http://wp.doc.ic.ac.uk/seams2017/call-for-artifacts/> for more details.

Selected papers will be invited to submit to the ACM Transactions on Autonomous and Adaptive Systems (TAAS).

Artifact Evaluation Committee
Konstantinos Angelopoulos, UK
Nuno Antunes, Portugal
Amel Bennaceur, UK
Javier Cámará, USA
Ilias Gerostathopoulos, Germany
Mahmoud Hammad, USA
Muhammad Usman Iftikhar, Sweden
Ashutosh Pandey, USA
Rojkong Sukerd, USA
Christos Tsigkanos, Italy