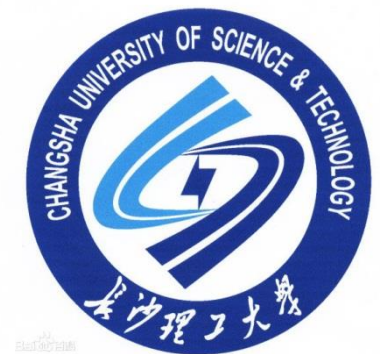


Learning-based Adaptive Data Placement for Low Latency in Data Center Networks

Kaiyang Liu, Jingrong Wang, Zhuofan Liao,
Boyang Yu, Jianping Pan



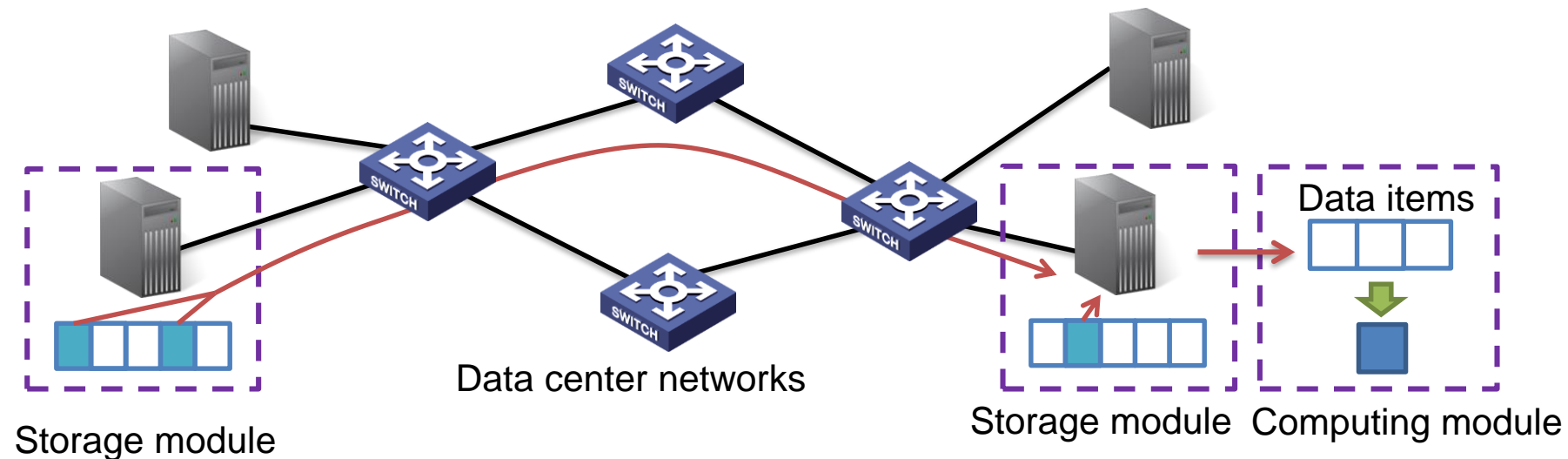
**University
of Victoria**



Data Analytics Services

➤ Data-intensive applications

- Data items need to be moved **frequently** between storage nodes
- This introduces **increased** and **changing** data access latency



➤ Data placement problem arises

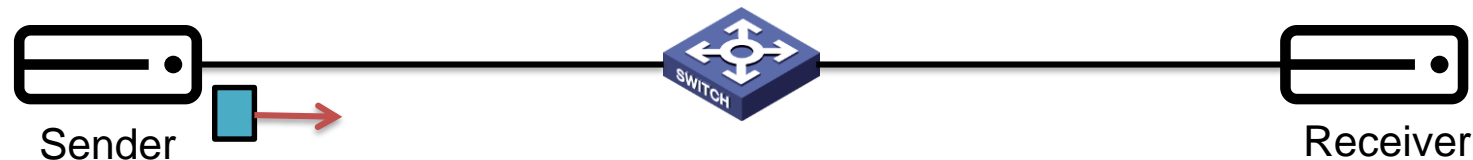
Data Placement

➤ Data storage

- Storage locations affect the finish time of the distributed computation tasks
- Main bottleneck: data movement latency [1]
- Amazon and Google reported that a slight increase in delay will lead to observable fewer user accesses

➤ Data movement latency

- Network latency = Processing + Queuing + Transmission + Propagation latency



- Different factors contribute to the latency

Related Work

- **Existing research efforts** [2]—[5]
 - Analyzing the factors that may affect the network latency
 - Hand-crafted design of optimization models
- **Limitations:** not flexible enough to deal with a dynamic environment
 - Different latency factors could be time-variant
 - Many uncertainties
 - Unreliable network links, variable user request patterns, and evolving system configurations

[2] Y. Xiang, et al., “Joint latency and cost optimization for erasure-coded data center storage,” IEEE/ACM Trans. Netw., 2016.

[3] X. Ren, et al., “Datum: Managing data purchasing and data placement in a geo-distributed data market,” IEEE/ACM Trans. Netw., 2018.

[4] B. Yu, et al., “A framework of hypergraph-based data placement among geo-distributed datacenters,” IEEE Trans. Serv. Comput., 2017.

[5] Y. Hu, et al., “Latency reduction and load balancing in coded storage systems,” in Proc. of ACM SoCC, 2017.

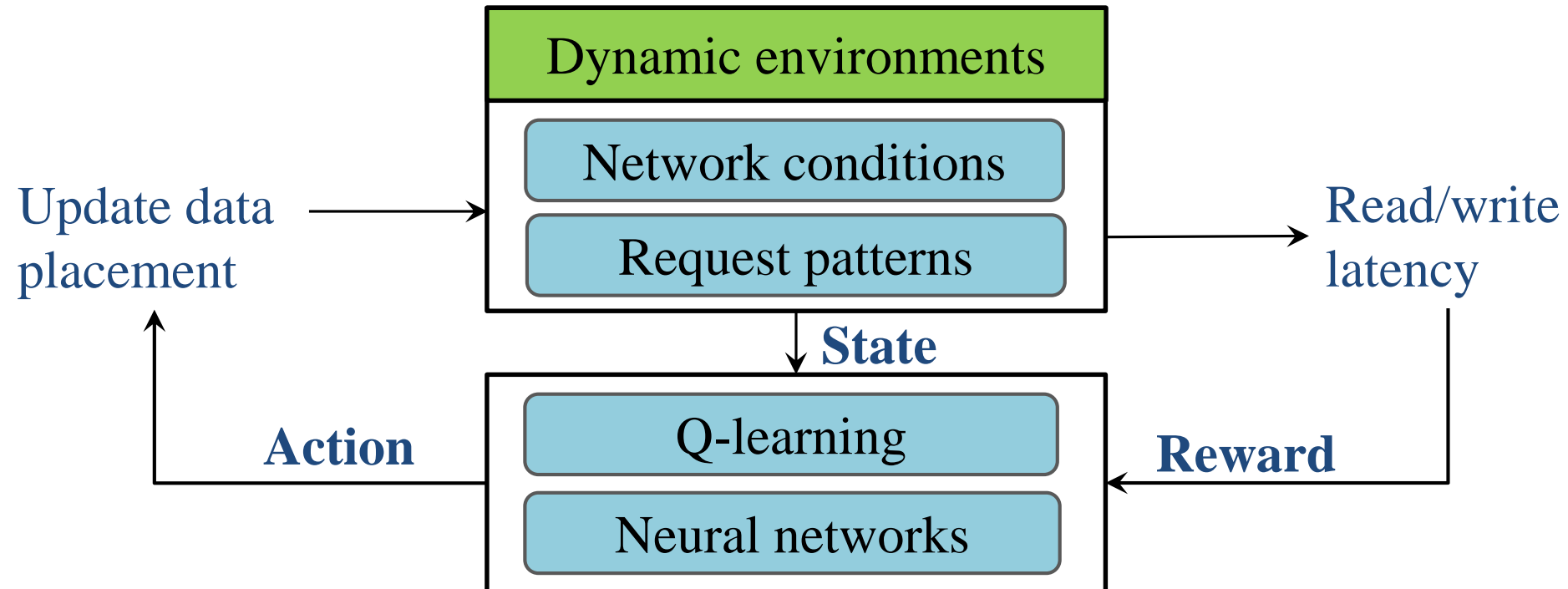
Key Questions

- **Data placement problem:** how to choose the storage locations of data items for low latency?

- **Challenges**
 - **Adaptability:** online schemes to deal with the network uncertainties
 - **Easy Implementation**
 - Low overhead
 - No need to modify the existing storage architecture

DataBot: A Learning-based Solution

➤ Design overview



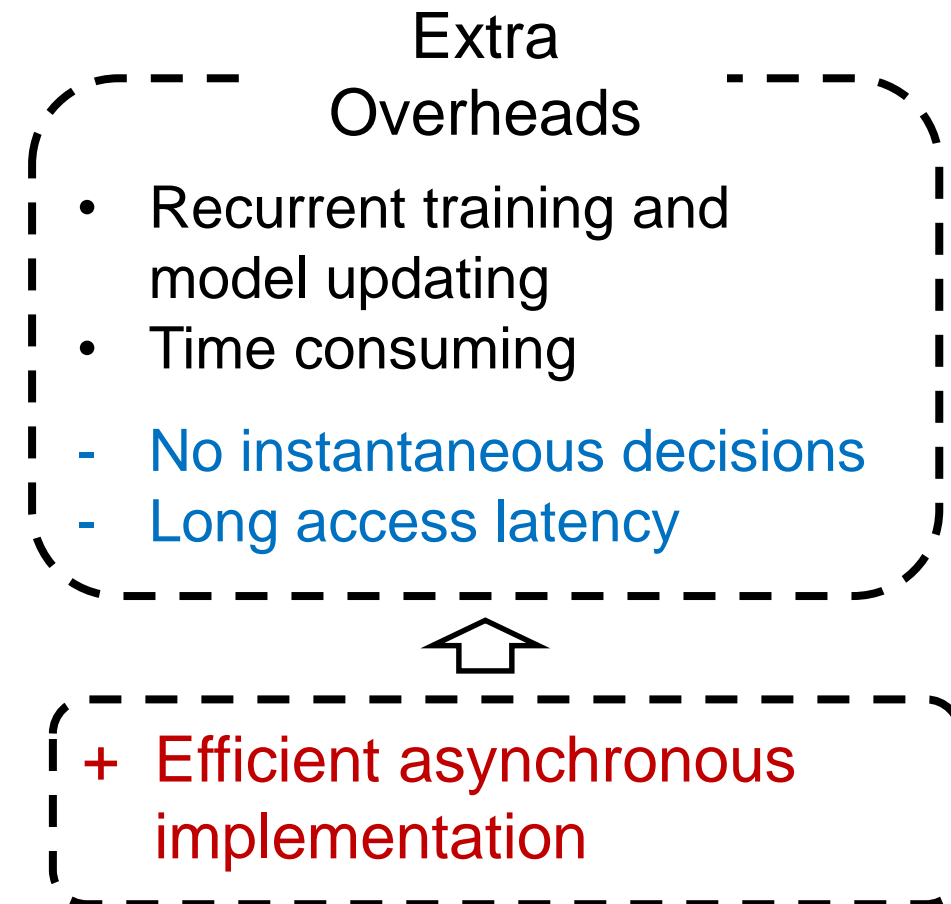
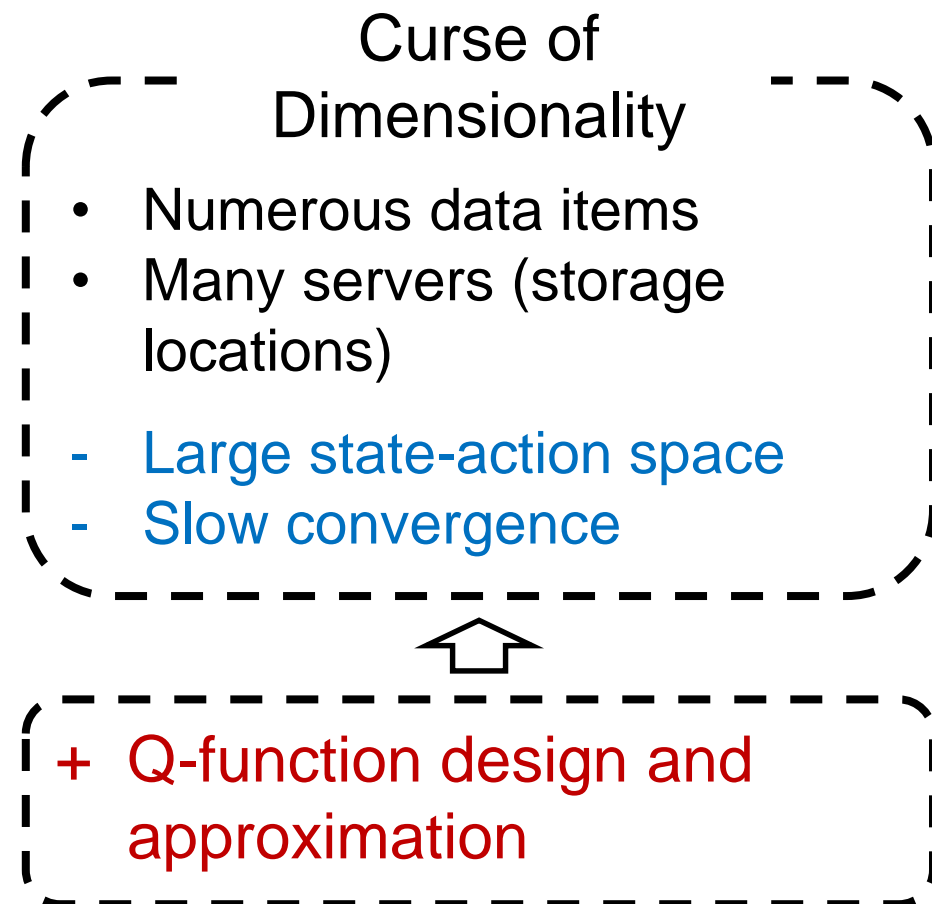
$$Q: \text{States}(S) \times \text{Action}(\mathcal{A}) = \text{Reward}(\mathcal{R})$$

Q-learning based Data Placement

- Data placement can be treated as a finite Markov Decision Process (**FMDP**)
 - The number of storage nodes is finite
 - Each action of data placement is independent
 - The performance of placement only depends on the current states and decisions
- The model-free Q-learning can find an optimal action selection policy for any given FMDP [6]

Our Contributions

➤ Solutions to address the limitations of conventional Q-learning



Q-Function Design (1)

➤ States (\mathcal{S})

a) Network conditions: $\{L_{ij}^{[R]}, L_{ij}^{[W]}, \forall i, j \in \mathcal{N}\}$

Average read/write latencies

Set of servers

Exponentially Weighted Moving Average (EWMA) mechanism [7]:

$$L_{ij}^{[R/W]} = \alpha_l l + (1 - \alpha_l) L_{ij}^{[R/W]}$$

Measured latency for each data movement

Discount factor

Benefits of EWMA: it only needs $O(1)$ space for latency estimation

Q-Function Design (2)

➤ States (\mathcal{S})

b) Request Patterns: $\{F_{im}^{[R]}, F_{im}^{[W]}, \tilde{F}_i^{[R]}, \tilde{F}_i^{[W]}, \forall i \in \mathcal{N}\}$



Read/write rates to data m from i



Total read/write rates from server i

Discounting Rate Estimator (DRE) method [8]

- Maintains a counter for each item
- Increases with every read/write operation
- Decreases periodically

Benefits of DRE: 1) it reacts quickly to the changes, and 2) only needs $O(1)$ space

Q-Function Design (3)

➤ States (\mathcal{S})

c) Source location: 0-1 vector

The size of state s will be: $|s| = 2N^2 + 5N = O(N^2) \rightarrow$ Number of servers

- The **number of data items** will not affect the deployment complexity

➤ Actions (\mathcal{A}): 0-1 vector (storage locations)

➤ Rewards (\mathcal{R}):
$$r_t = \omega \cdot \frac{1}{l^{[W]}} + (1 - \omega) \cdot \frac{1}{|\mathcal{P}|} \cdot \sum_{p \in \mathcal{P}} \frac{1}{l_p^{[R]}}$$

↓
Tradeoff parameter

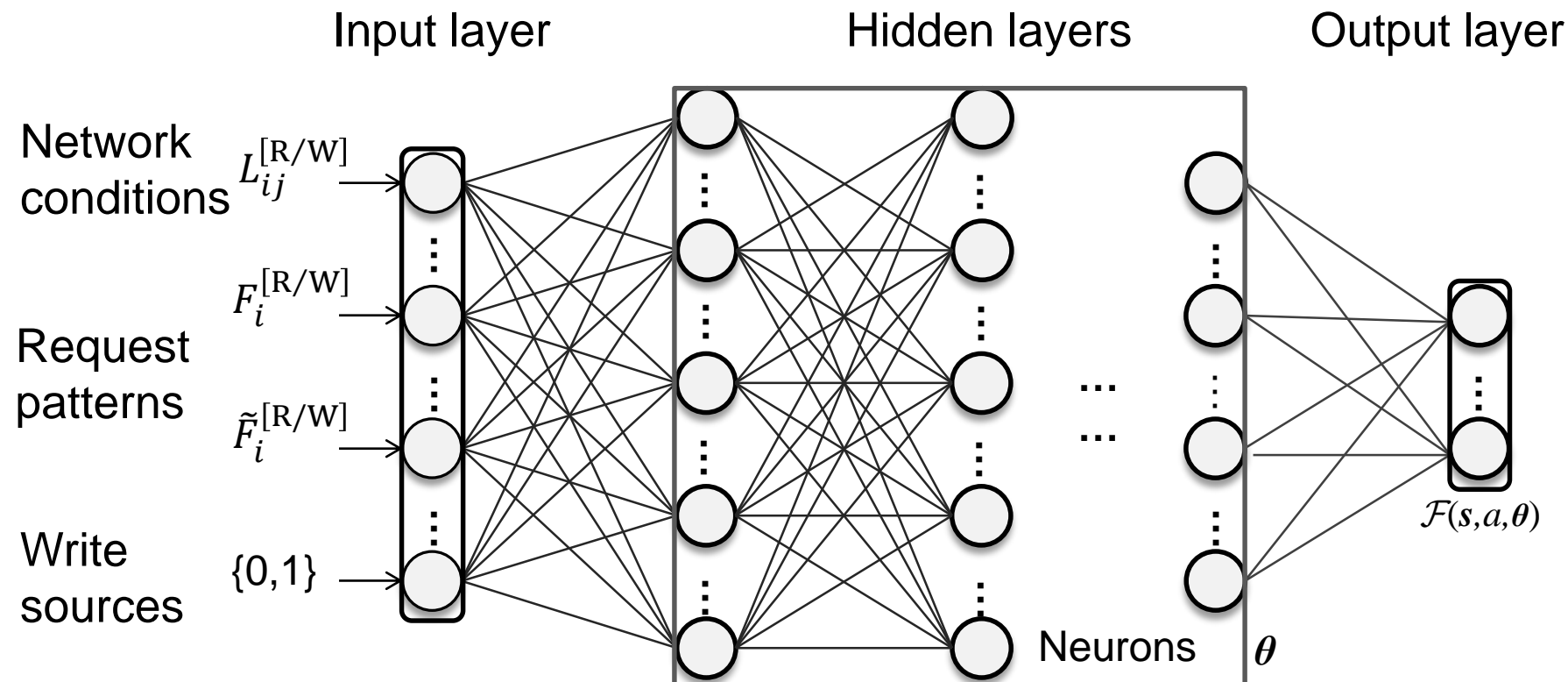
↓
of read operations between two write operations

The **measured** read/write latencies are used as the reward

Q-function Approximation

➤ Neural networks (NN)

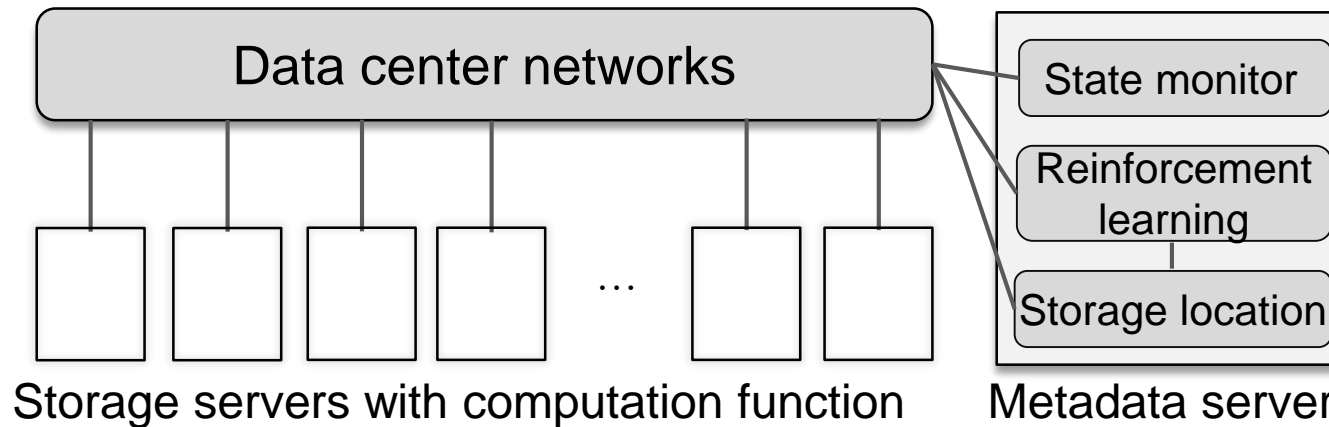
- Learn to output the expected rewards of data placement actions
- Lower the scale of the state space (**number of servers**)



System Architecture

➤ Data storage architecture

- DataBot is implemented in the metadata server



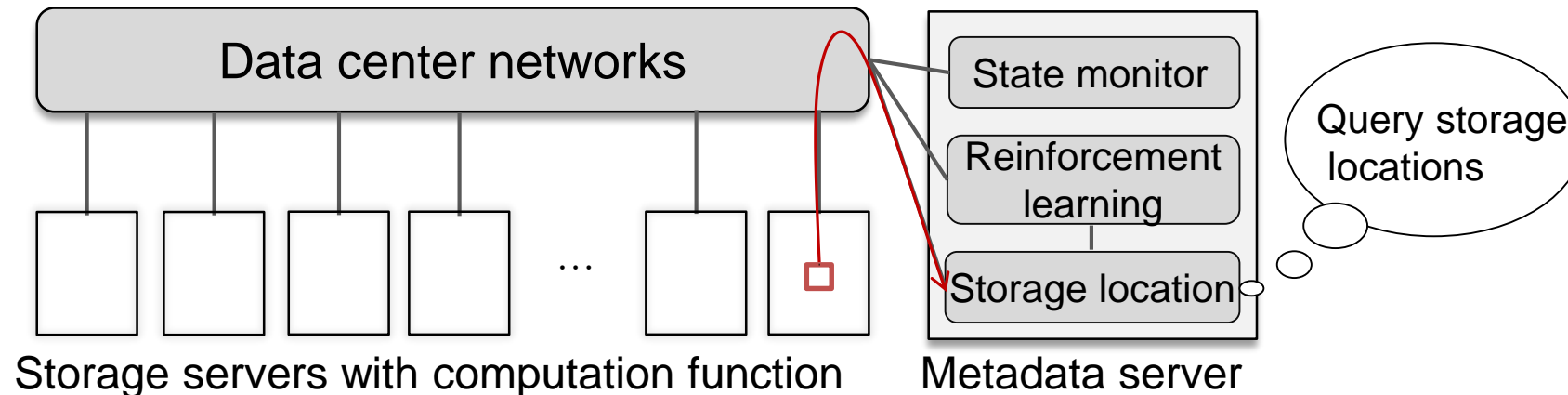
➤ Metadata server

- Manages the storage locations of data items, e.g., using hashtag
- Captures the **logs** of the read/write requests: (**TS**, **R/W**, **Src**, **Dst**, **Lat**)

System Architecture

➤ Data storage architecture

- DataBot is implemented in the metadata server



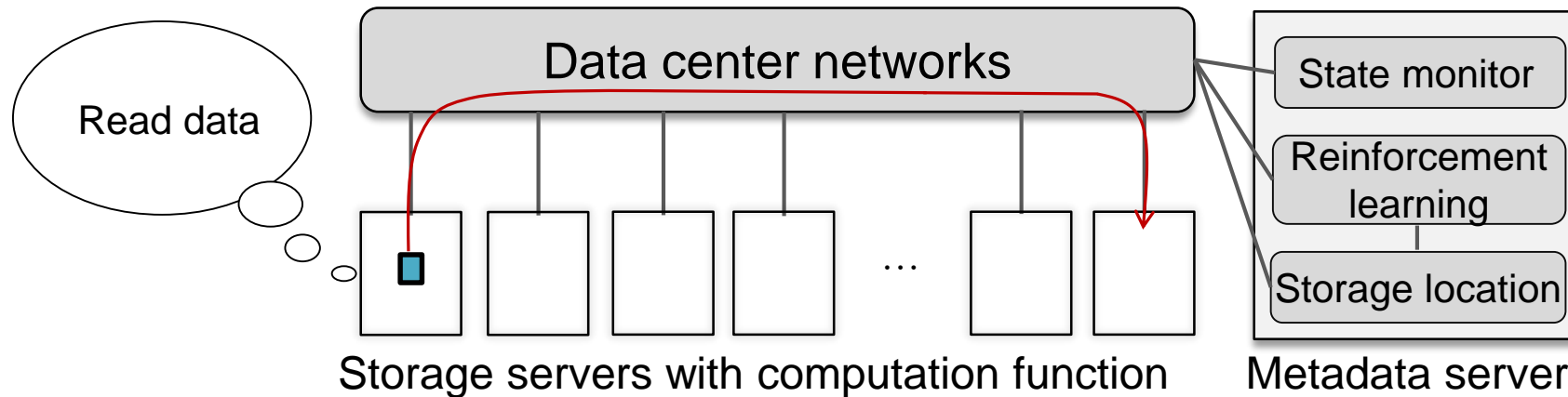
➤ Metadata server

- Manages the storage locations of data items, e.g., using hashtag
- Captures the **logs** of the read/write requests: (**TS**, **R/W**, **Src**, **Dst**, **Lat**)

System Architecture

➤ Data storage architecture

- DataBot is implemented in the metadata server



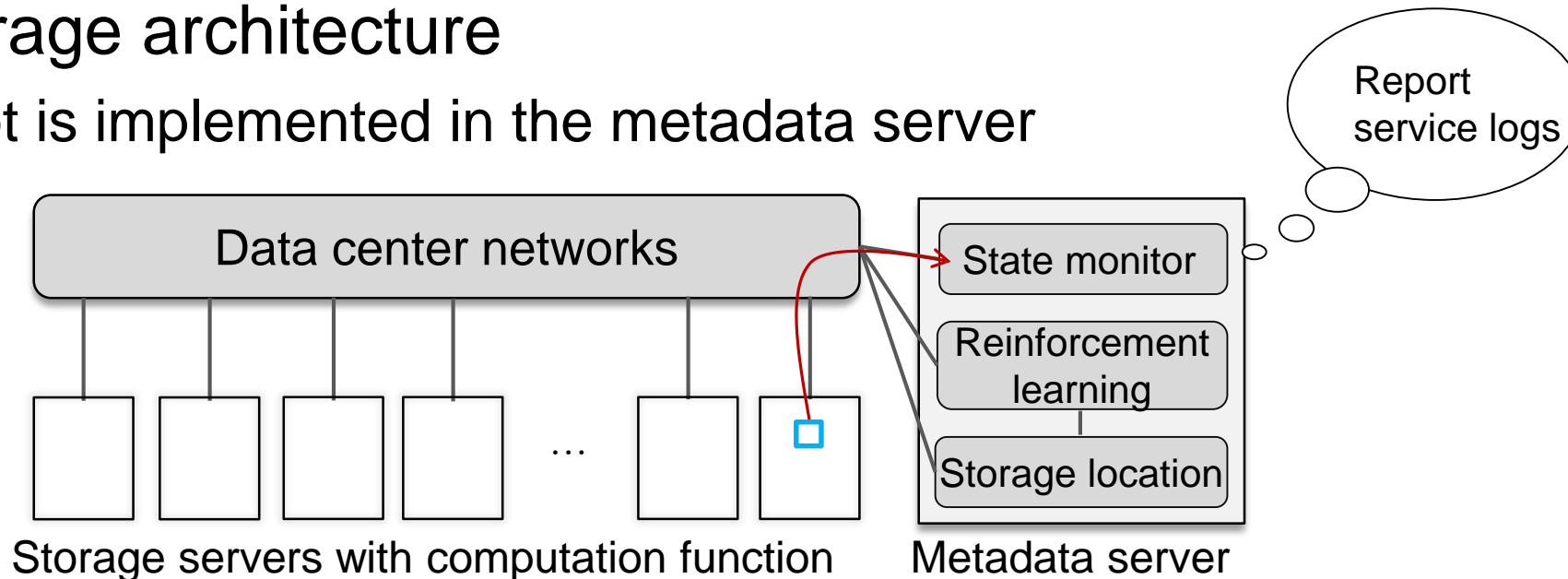
➤ Metadata server

- Manages the storage locations of data items, e.g., using hashtag
- Captures the **logs** of the read/write requests: (**TS**, **R/W**, **Src**, **Dst**, **Lat**)

Asynchronous Implementation (1)

➤ Data storage architecture

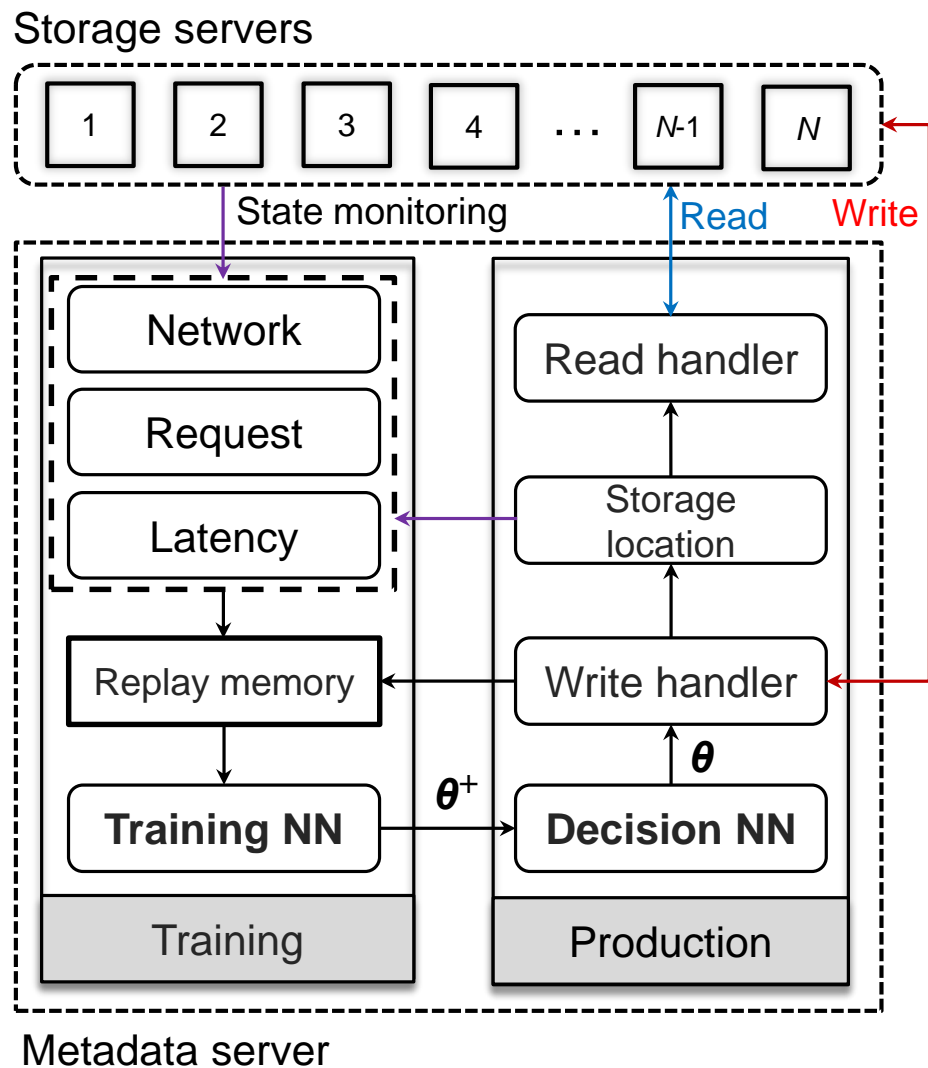
- DataBot is implemented in the metadata server



➤ Metadata server

- Manages the storage locations of data items, e.g., using hashtag
- Captures the **logs** of the read/write requests: (**TS**, **R/W**, **Src**, **Dst**, **Lat**)

Asynchronous Implementation



➤ Production system: decision NN

- **Input:** state s_t ; **Output:** expected reward $\mathcal{F}(s_t, a_t, \theta)$
- **ϵ -greedy method:** with probability ϵ to select the action a_t that maximizes the output value
- A tuple $\tau = (s_t, a_t, s_{t+1}, r_t)$ is stored for each request

➤ Training system: training NN

- Tuples for a period \rightarrow Replay memory \mathcal{R} for training
- Mini-batch stochastic gradient descent (SGD) [9]: training weight vector θ^+
 - **Batch:** all tuples in \mathcal{R} are partitioned into mini-batches
 - **Epoch:** mini-batches are trained with multiple iterations
- Weight update: $\theta \leftarrow \theta^+$

Performance Evaluation

Evaluation Setup: Traces

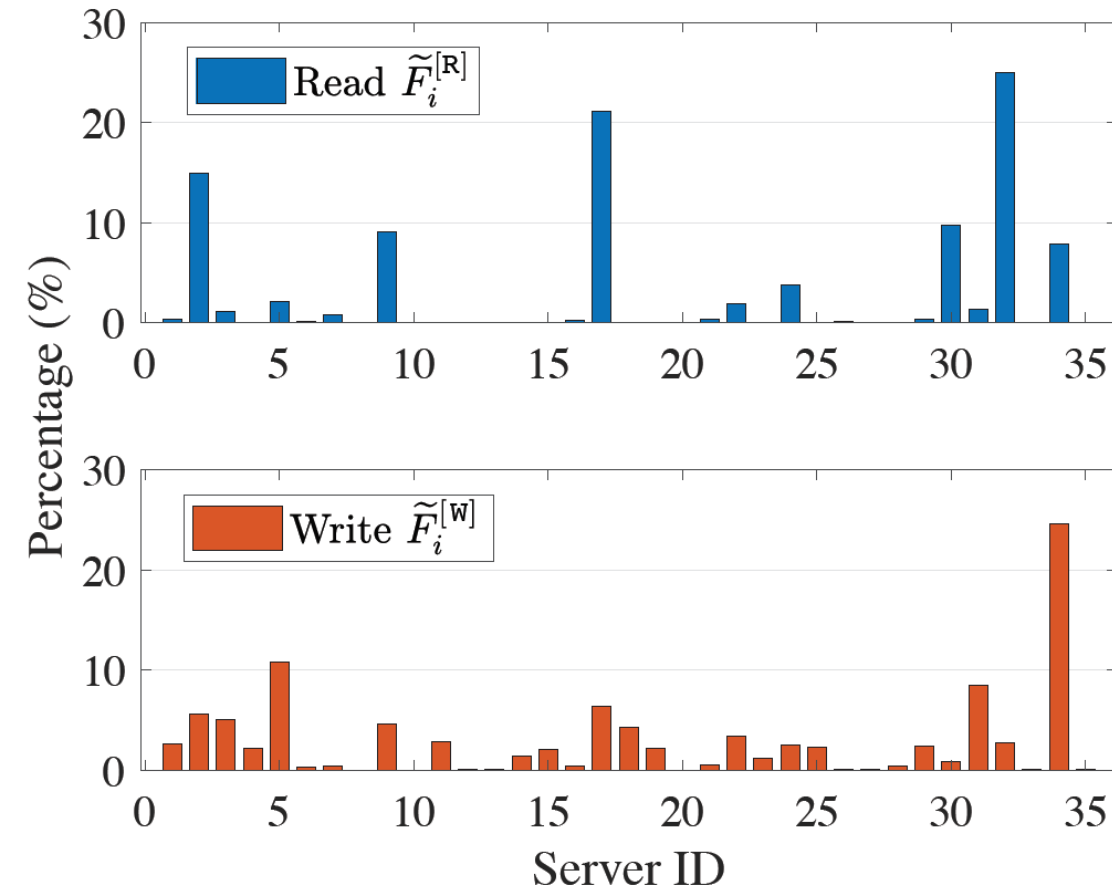
➤ MSR Cambridge Traces [10]

- I/O traces of an enterprise data center
- Hostname, request type (read/write), and timestamp
- Request distribution is **biased** among 36 storage servers

➤ **Limitation:** do not specify the detailed data item for each read/write request

➤ **Assumption**

- Number of data items: 10,000
- The request rates of data items follow a Zipf distribution among servers

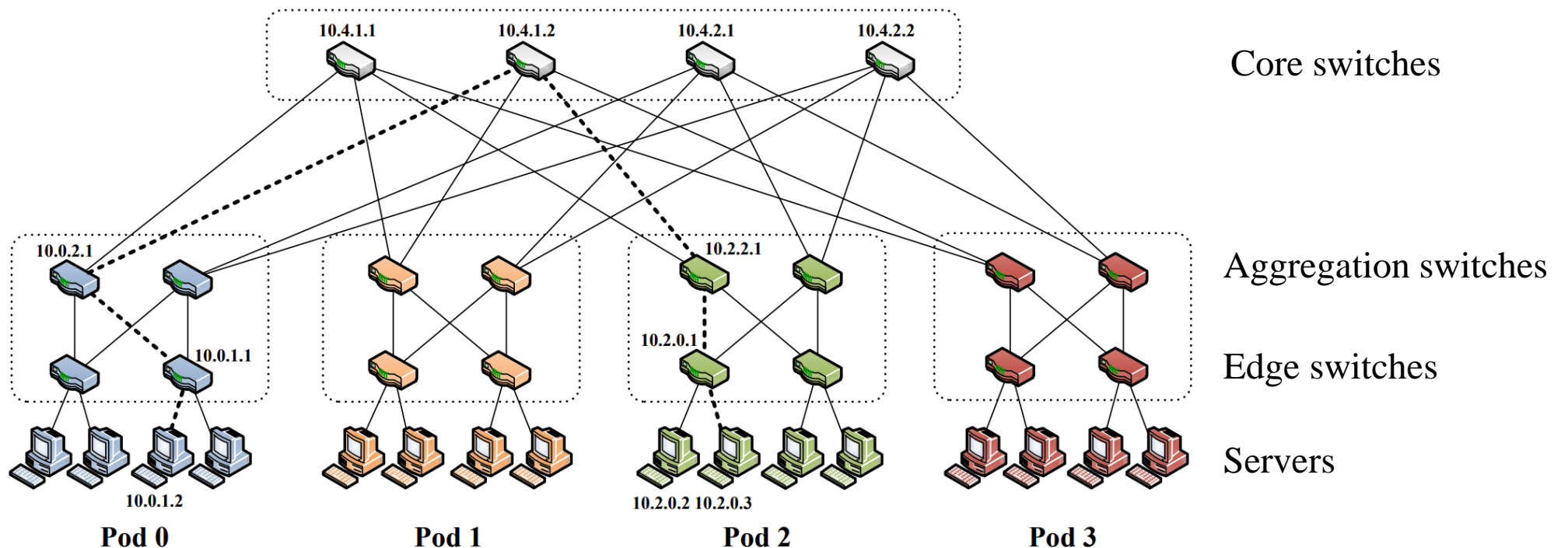


Arrival rates of the read/write requests

Evaluation Setup: Scenarios (1)

➤ Data center network emulation: **Mininet**

- Representative network topology: Fat-Tree [11];
 - Default data block size: 64 MB [12]
- Link capacity: 1 Gbps



[11] M. Al-Fares, et al., “A scalable, commodity data center network architecture,” in Proc. of ACM SIGCOMM, 2008.

[12] “HDFS Architecture Guide.” [Online]: <https://hadoop.apache.org/>

Evaluation Setup: Scenarios (2)

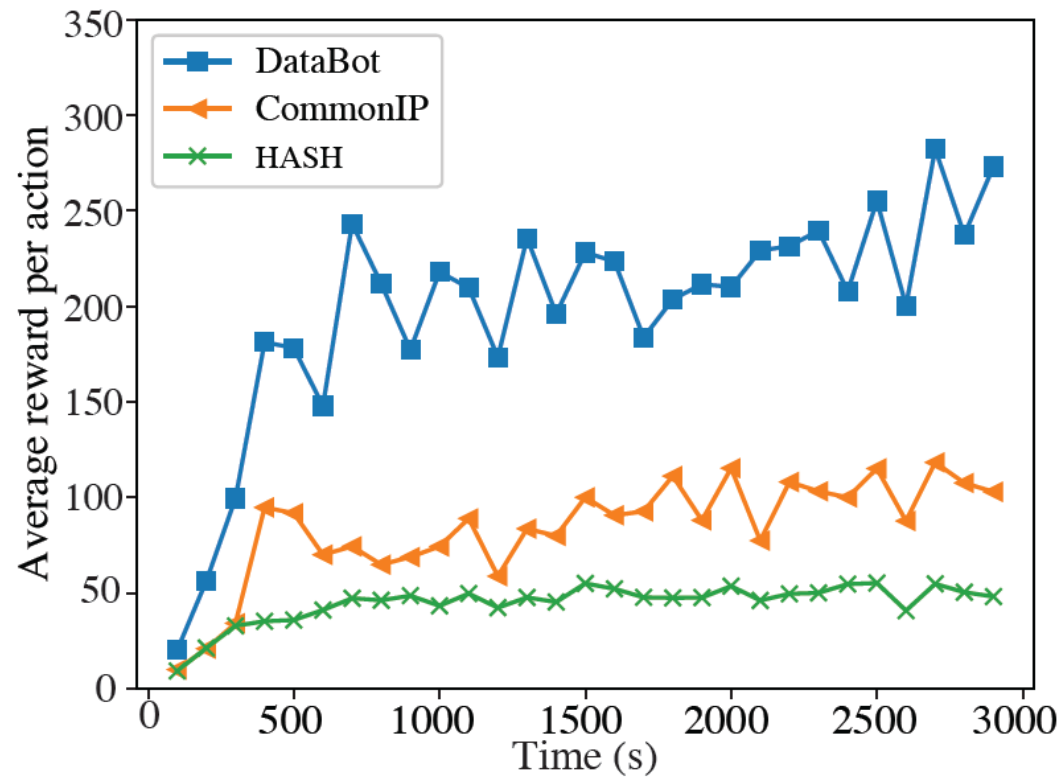
- Client program: Initiates the read/write requests
 - Memcached: the end of data flows for data caching in RAM
- Metadata server program
 - State monitoring
 - Write destination decision
 - NN training: Multilayer perceptron with one kernel
- Performance baselines
 - **HASH** [11] hashes data to servers for load-balancing
 - **CommonIP** [12] places data as close as possible to the IP address that most commonly accesses the data under the constraint of storage capacity

[11] “HDFS Architecture Guide.” [Online]: <https://hadoop.apache.org/>

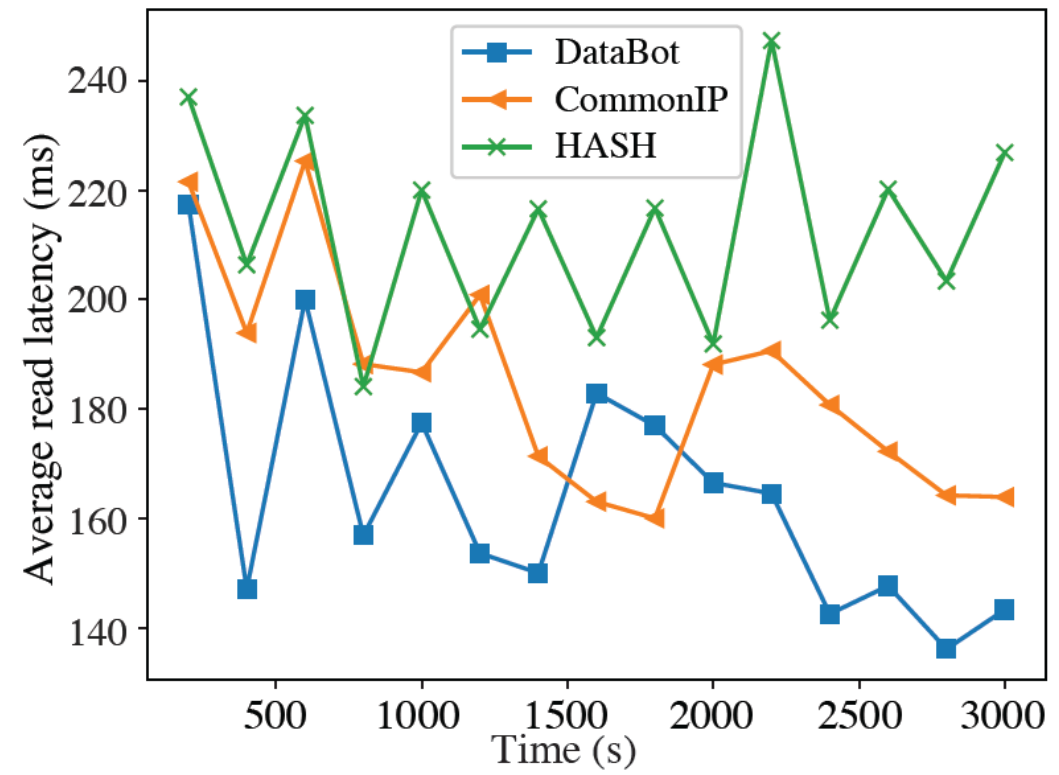
[12] S. Agarwal, et al., “Volley: Automated data placement for geo-distributed cloud services,” in Proc. of USENIX NSDI, 2010.

Results: Read optimized

➤ Read optimized: write weight $\omega = 0.2$



(a) Average reward per action r_t



(b) Average read latency

Results: Other Factors (1)

➤ Parameter impacts: write weight and number of replicas

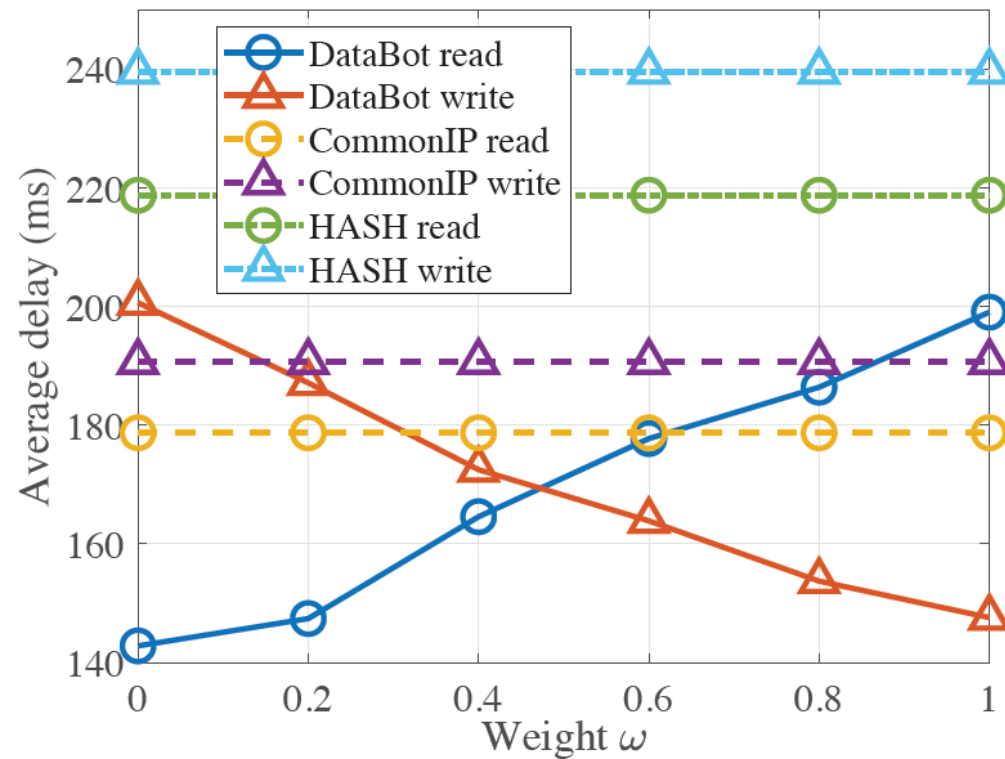


Fig. 7. Impact of weight ω .

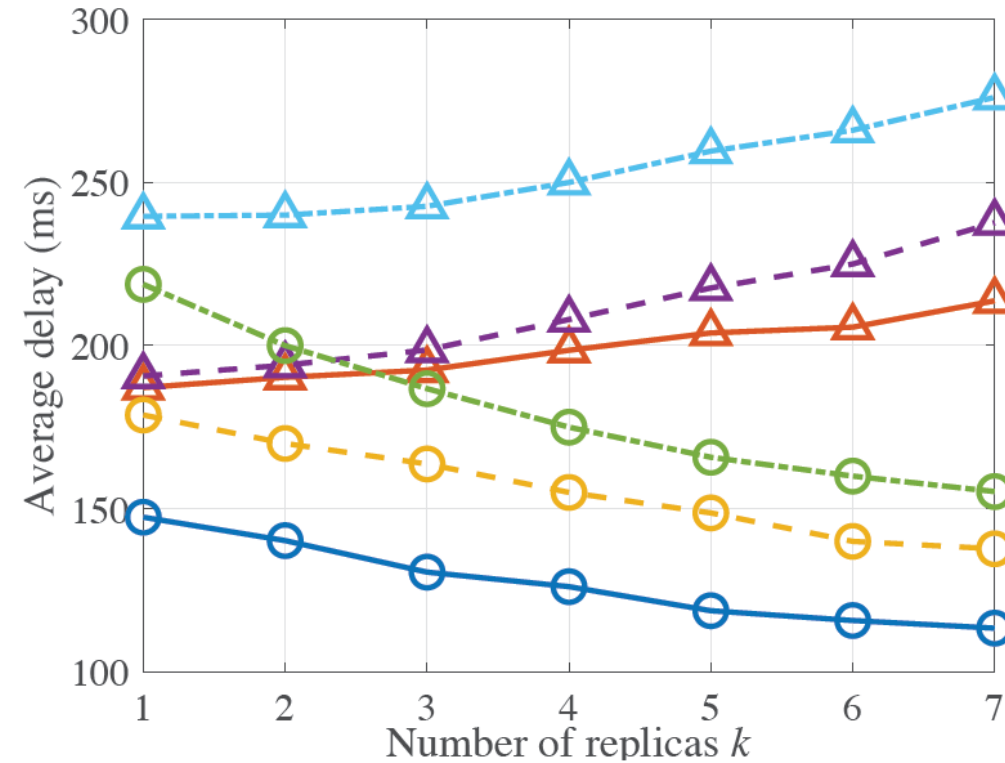


Fig. 8. Impact of replica number k .

Results: Other Factors (2)

➤ Parameter impacts: number of training epochs and batch size

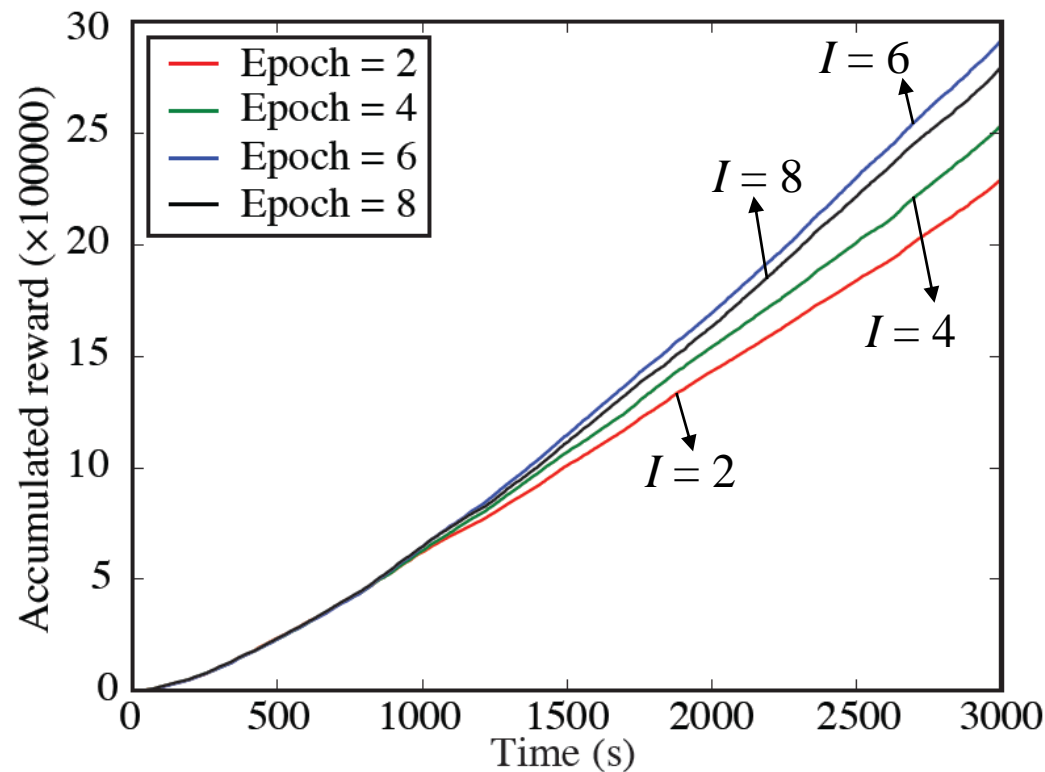


Fig. 9. Impact of epoch number I .

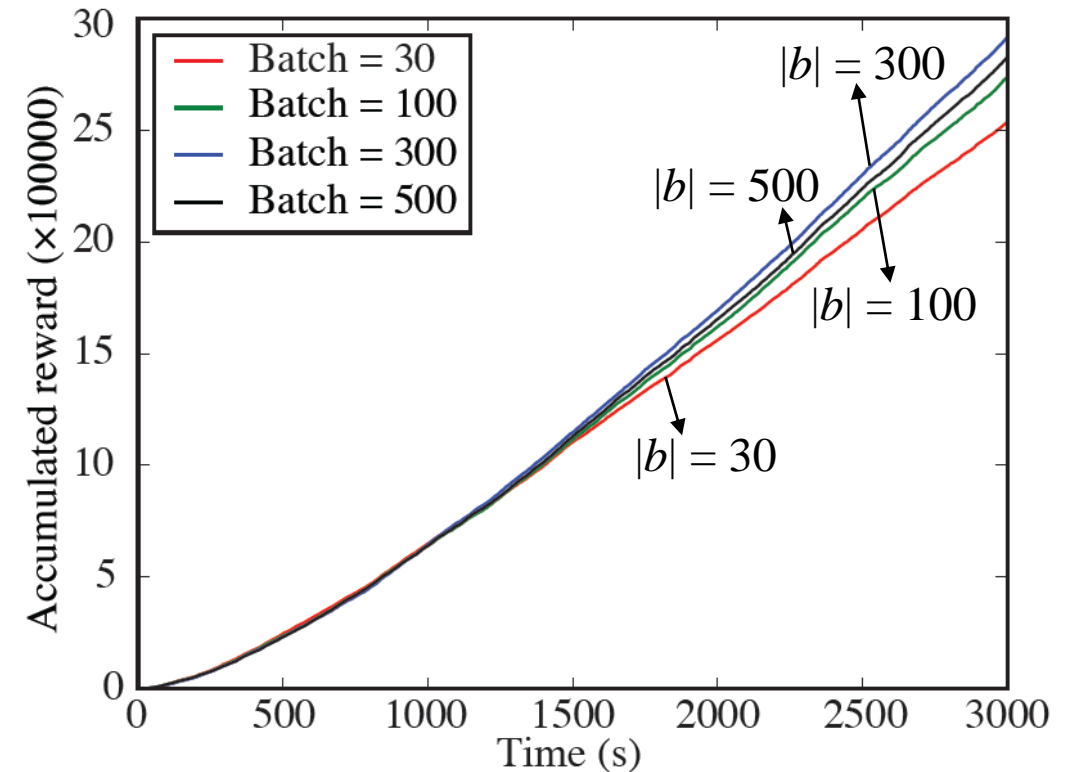


Fig. 10. Impact of batch size $|b|$.

Conclusions

- **DataBot** automatically learns the optimal data placement policies to handle the **system uncertainties**
 - With no future information about the **dynamics**
- Neural networks achieve a quick approximation when combined with Q-learning
- Asynchronous implementation
 - Online decision making and offline training

Thanks

Happy to answer your questions

Backup Slides

Neural Networks

➤ Structure of NN

- Multilayer perceptron (MLP) with one kernel
- **Input layer:** 1,476 features; **Output layer:** 36 features
- **Two hidden layers:** 1,000 and 600 features

➤ Weight vector training

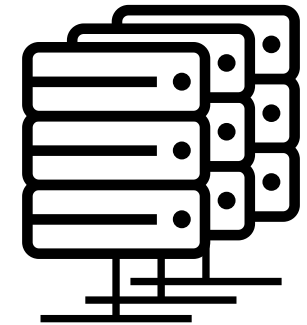
- Traditional back propagation method

➤ Implementation based on popular learning frameworks

- Keras deep learning library [14] (with TensorFlow as backend)

Scalability

- How to improve the scalability when the serves are deployed on a large scale?
 - Hundreds or thousands of servers
- Our solution in the future work



Distributed learning mechanism

- Multiple workers run in parallel to train the partitions of the input dataset
- Works update shared model parameters for training
- The learning process can be sped up with no need of aggregating raw data to a centralized metadata server.