

Orpheus: Efficient Distributed Machine Learning via System and Algorithm Co-design

Pengtao Xie (Petuum Inc)

Jin Kyu Kim (CMU), Qirong Ho (Petuum Inc), Yaoliang Yu (University of Waterloo), Eric P. Xing (Petuum Inc)

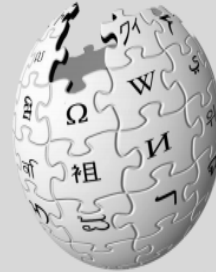
Massive Data



facebook

1B+ USERS

30+ PETABYTES



WIKIPEDIA
The Free Encyclopedia

32 million
pages



You Tube

100+ hours video
uploaded every minute



twitter

645 million users
500 million tweets / day

Distributed ML Systems

Parameter Server
Systems

Yahoo LDA DistBelief
Li & Smola PS Project Adam
Bosen GeePS

Graph Processing
Systems



Pregel GraphX

Dataflow Systems

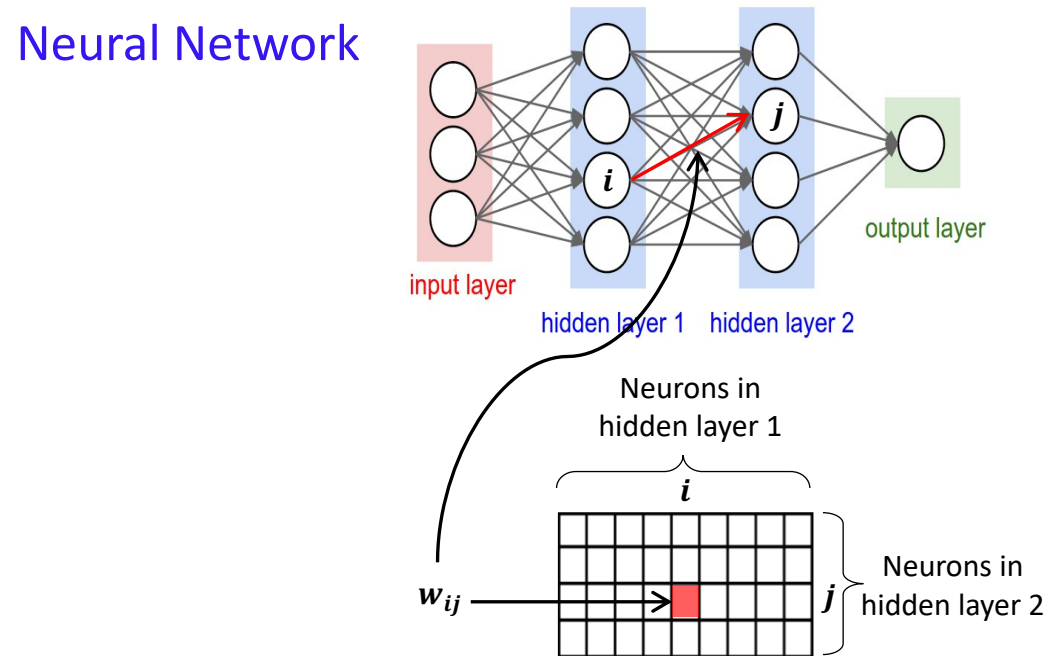


Hybrid Systems



Matrix-Parameterized Models (MPMs)

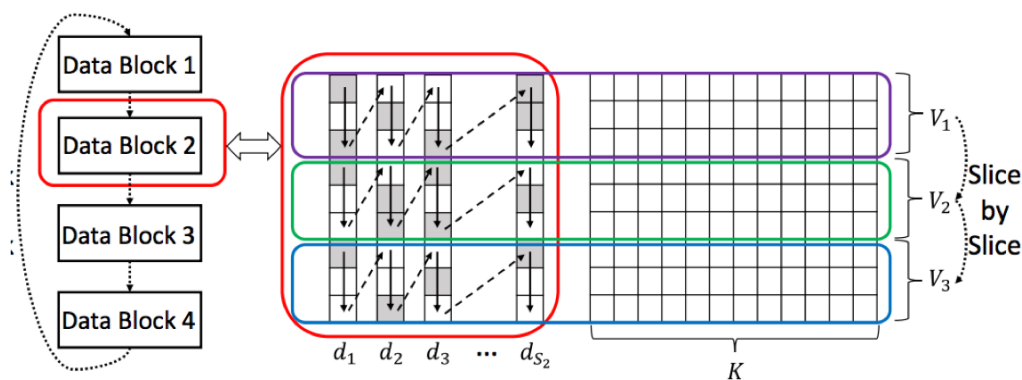
- Model parameters are represented as a matrix



- Other examples: Topic Model, Multiclass Logistic Regression, Distance Metric Learning, Sparse Coding, Group Lasso, etc.

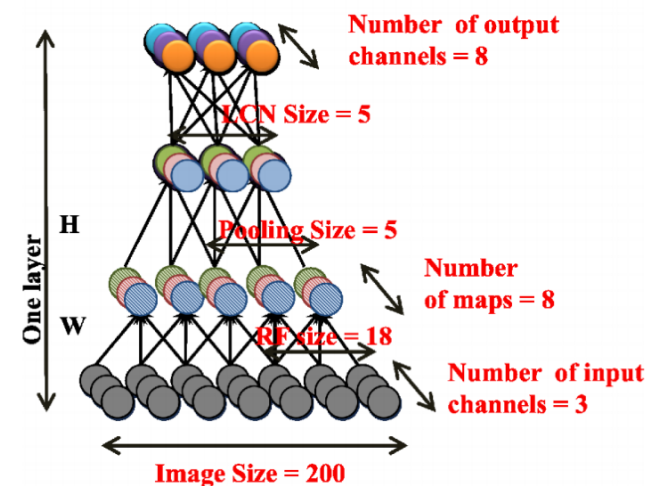
Parameter Matrices Could Be Very Large

LightLDA Topic Model
(Yuan et al. 2015)



The topic matrix has **50 billion** entries.

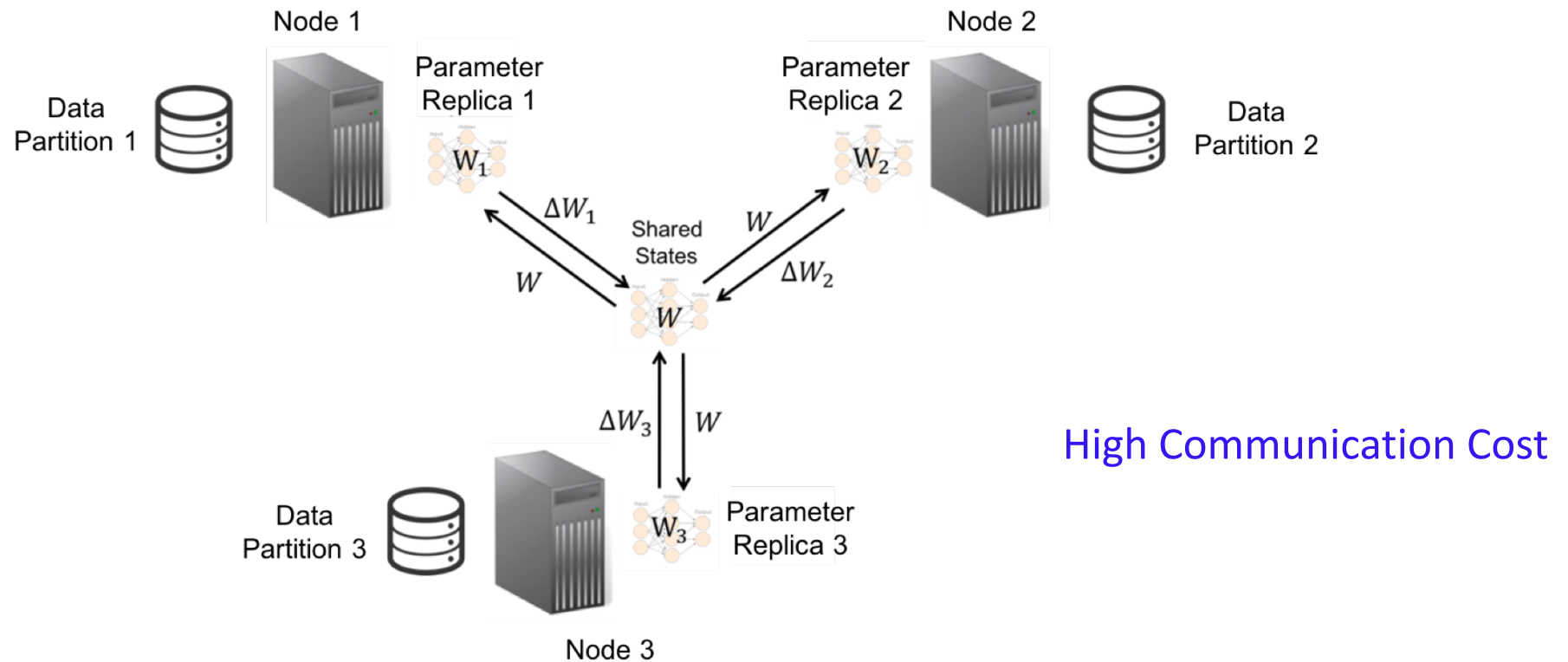
Google Brain Neural Network
(Le et al. 2012)



The weight matrices have **1.3 billion** entries.

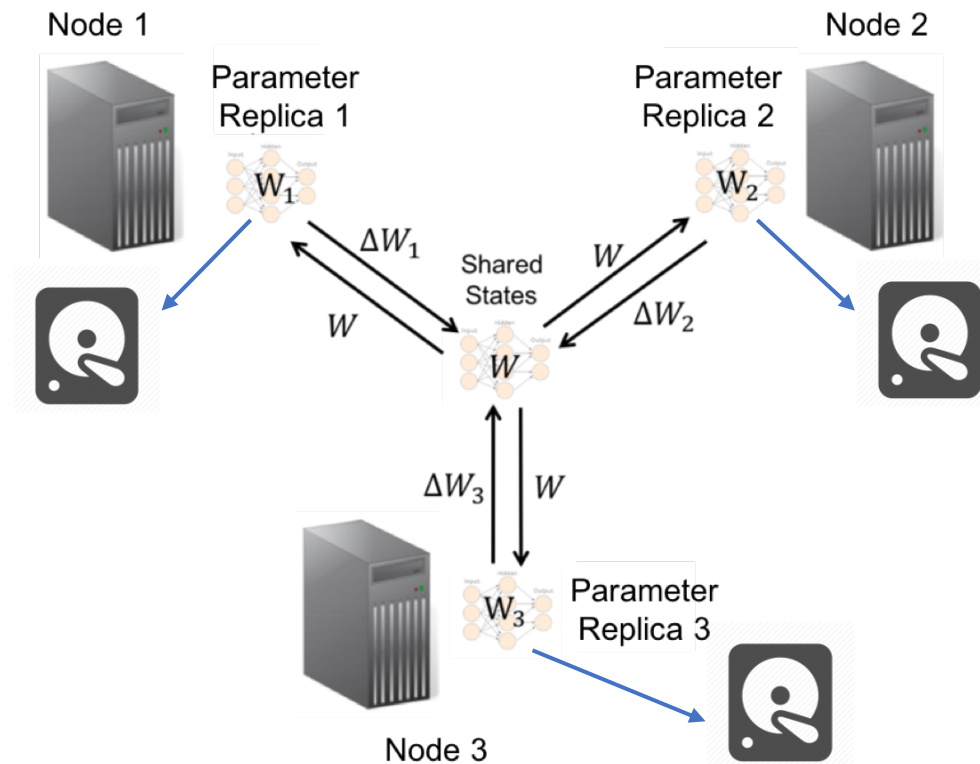
Existing Approaches

- Parameter server frameworks communicate matrices for parameter synchronization.



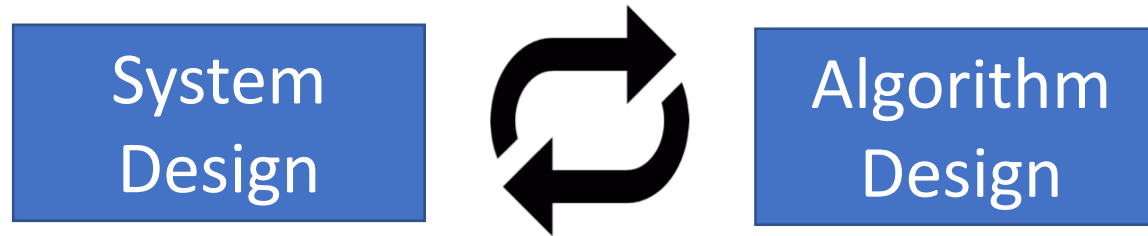
Existing Approaches (Cont'd)

- Parameter matrices are checkpointed to stable storage for fault tolerance.



High Disk IO

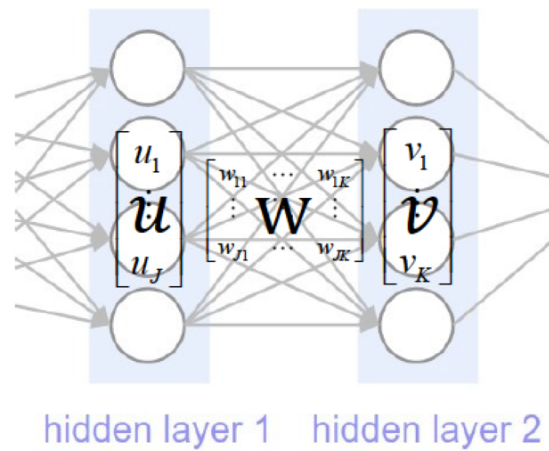
System and Algorithm Co-design



- System design should be tailored to the unique mathematical properties of ML algorithms
- Algorithms can be re-designed to better exploit the system architecture

Sufficient Vectors (SVs)

- Parameter-update matrix can be computed from a few vectors (referred to as sufficient vectors)

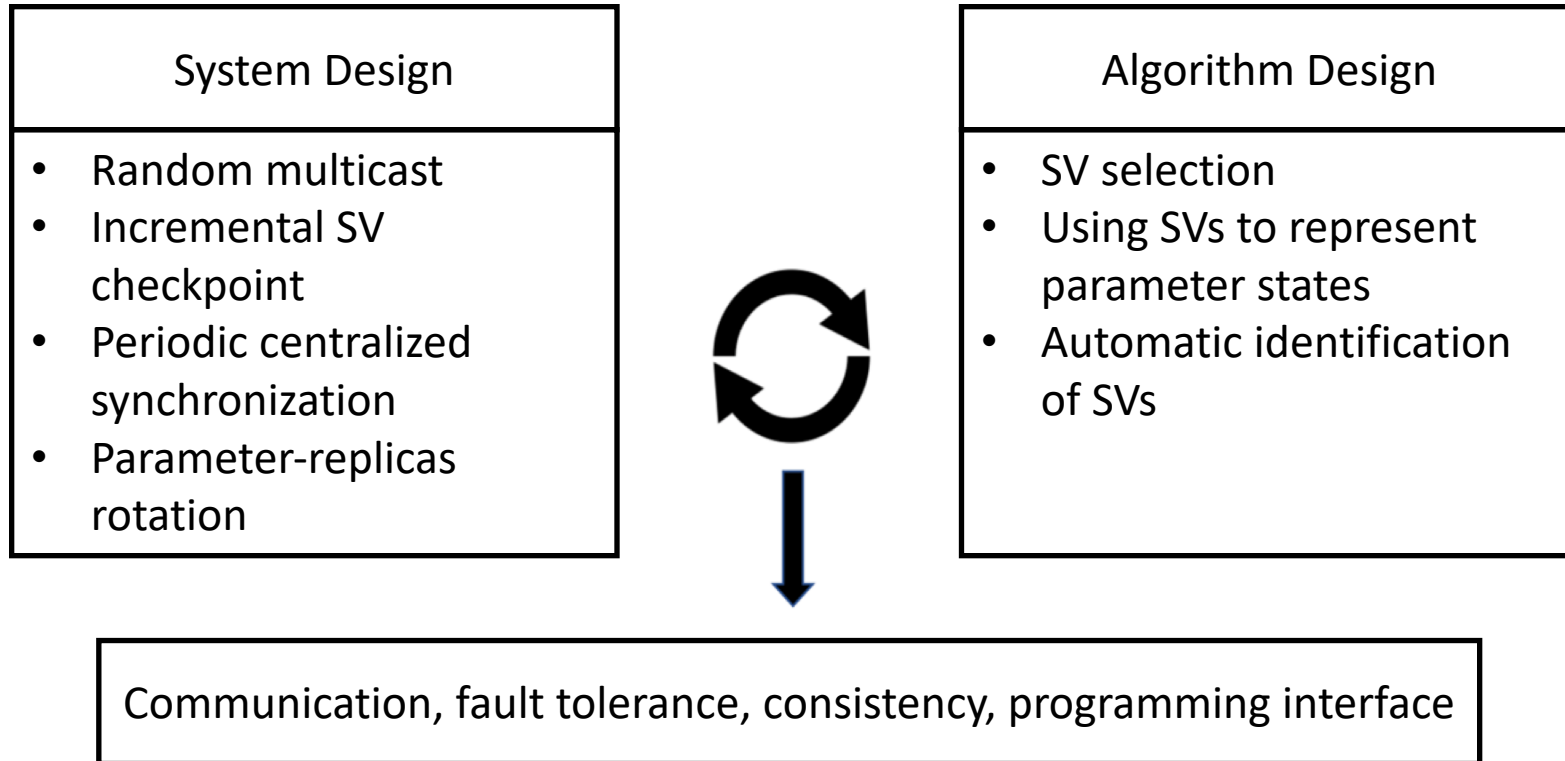


$J \times K$ Entries $J + K$ Entries Sufficient Vectors

$$\Delta W = \mathbf{u} \otimes \mathbf{v}$$
$$\begin{bmatrix} \Delta w_{11} & \cdots & \Delta w_{1K} \\ \vdots & \ddots & \vdots \\ \Delta w_{J1} & \cdots & \Delta w_{JK} \end{bmatrix} \begin{bmatrix} u_1 \\ \vdots \\ u_J \end{bmatrix} \begin{bmatrix} v_1 & \cdots & v_K \end{bmatrix}$$

(Xie et al. 2016)

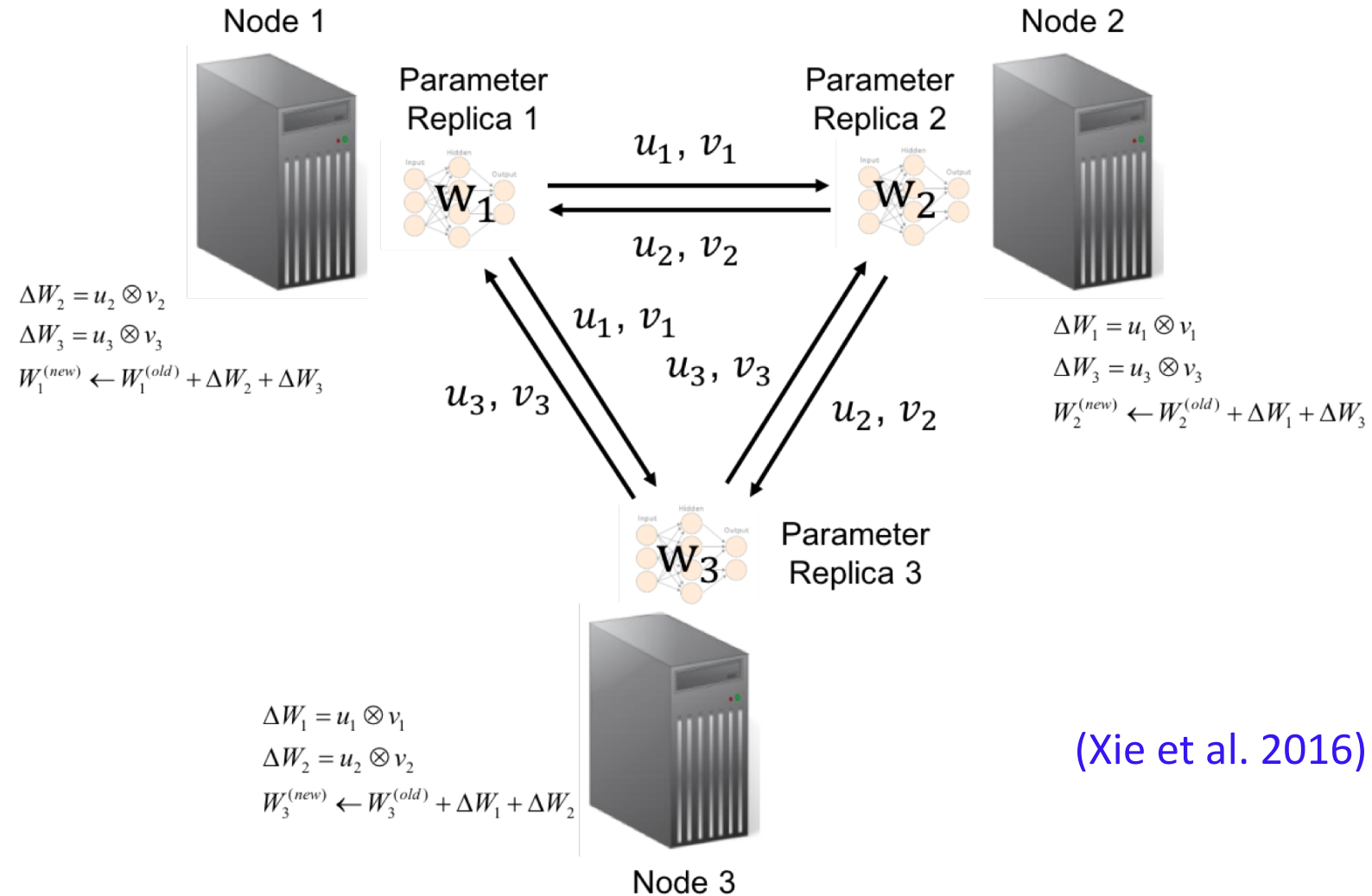
System and Algorithm Co-design



Outline

- Introduction
- **Communication**
- Fault tolerance
- Evaluation
- Conclusions

Peer-to-Peer Transfer of SVs



Cost Comparison

	Size of one message	Number of messages	Network Traffic
P2P SV-Transfer	$O(J + K)$	$O(P^2)$	$O((J + K)P^2)$
Parameter Server	$O(JK)$	$O(P)$	$O(JKP)$

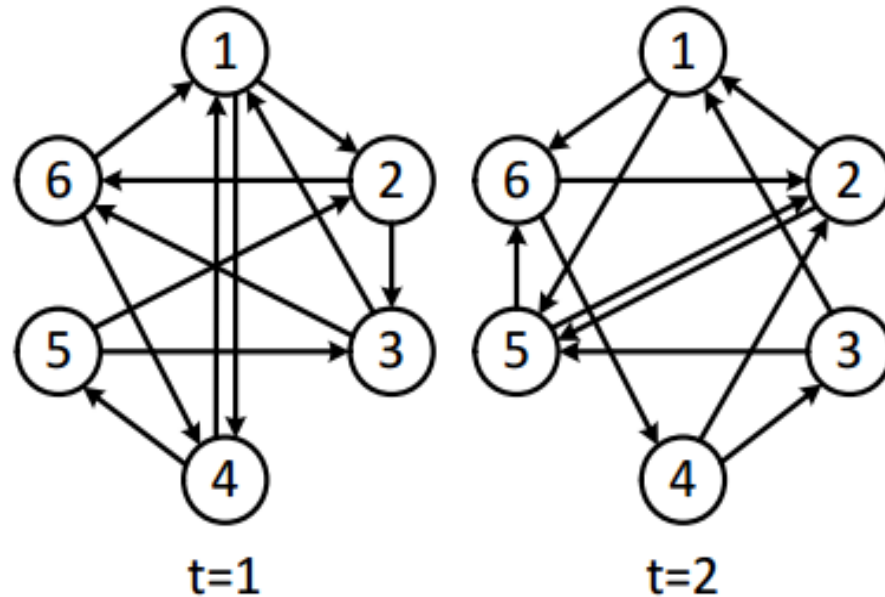
J, K: dimensions of the parameter matrix

P: number of machines

How to reduce the number of messages in P2P?

Random Multicast

- Send SVs to a random subset of Q ($Q \ll P$) machines
- Reduce number of messages from $O(P^2)$ to $O(PQ)$



Random Multicast (Cont'd)

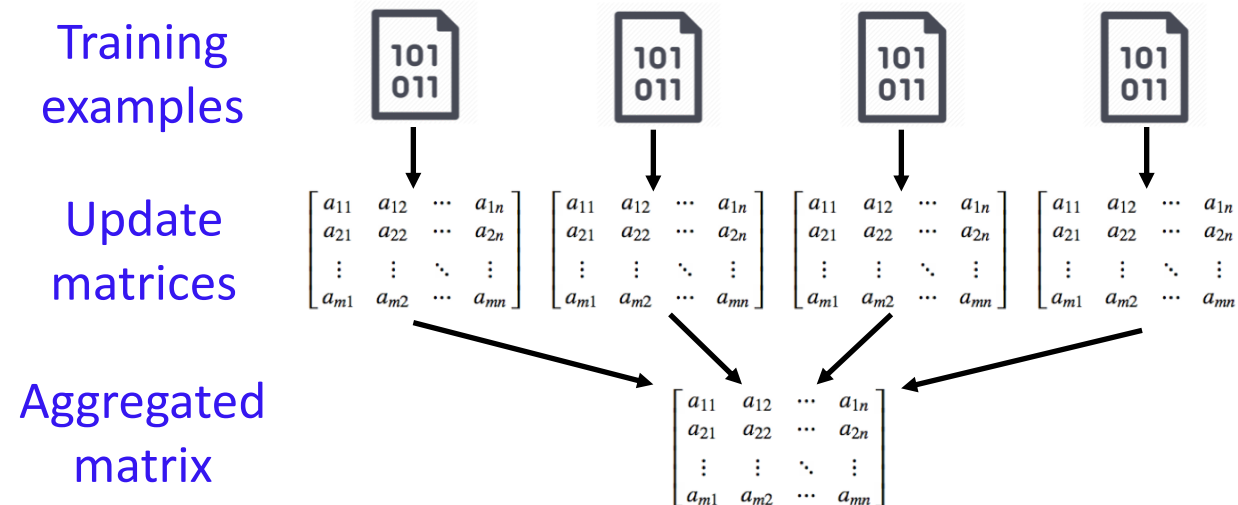
- Correctness is guaranteed due to the error-tolerant nature of ML.

Corollary 1. Denote $\{\mathbf{W}_p^c\}$, $p = 1, \dots, P$, and $\{\mathbf{W}^c\}$ as the local sequences and the auxiliary sequence generated by multicast SFB, respectively, and assume that each machine aggregates the stochastic gradients from $Q < P$ other machines at every iteration. Let Assumption 1 hold. Then, with learning rate $\eta_c \equiv O(\frac{1}{L(P-Q)C})$, we have

$$\min_{c=1, \dots, C} \mathbb{E} \left\| \sum_{p=1}^P \nabla f_p(\mathbf{W}_p^c) \right\|^2 \leq O \left(L(P-Q) + \frac{(Ls + L_f)\sigma^2}{KPL(P-Q)C} \right).$$

Mini-Batch

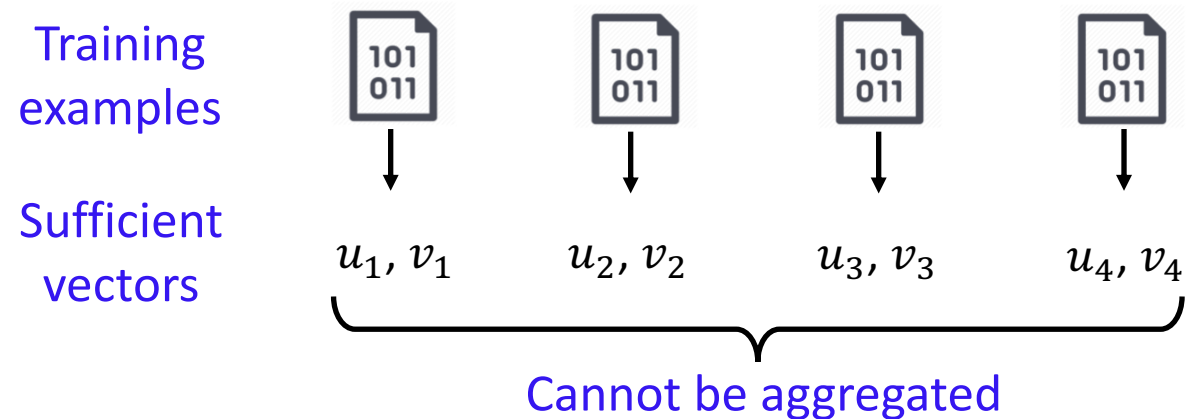
- It is common to use a mini-batch of training examples (instead of one) to compute updates
- If represented as matrices, the updates computed w.r.t different samples can be aggregated into a single update matrix to communicate



- Communication cost does not grow with mini-batch size

Mini-Batch (Cont'd)

- If represented as SVs, the updates computed w.r.t different samples cannot be aggregated into a single SV



- The SVs must be transmitted individually
- Communication cost grows linearly with mini-batch size

SV Selection

- Select a subset of “representative” SVs to communicate
- Reduce communication cost
- Does not hurt the correctness of updates
 - The aggregated update computed from the selected SVs are close to that from the entire mini-batch
 - The selected SVs can well represent the others

SV Selection (Cont'd)

- Algorithm: joint matrix column subset selection

$$\min_I \sum_{r=1}^R \left\| X^{(r)} - S_I^{(r)} (S_I^{(r)})^\dagger X^{(r)} \right\|_2$$

Algorithm 2 Joint Matrix Column Subset Selection

Input: $\{X^{(r)}\}_{r=1}^R$

Initialize: $\forall r, X_0^{(r)} = X^{(r)}, S_0^{(r)} = []$

for $t \in \{1, \dots, C\}$ **do**

 Compute the squared L2 norm of column vectors in

$\{X_{t-1}^{(r)}\}_{r=1}^R$

 Sample a column index i_t

$\forall r, X_t^{(r)} \leftarrow X_{t-1}^{(r)} / \{x_{i_t}^{(r)}\}, S_t^{(r)} \leftarrow S_{t-1}^{(r)} \cup \{x_{i_t}^{(r)}\}$

$\forall r, X_t^{(r)} \leftarrow X_t^{(r)} - S_t^{(r)} (S_t^{(r)})^\dagger X_t^{(r)}$

end for

Output: $\{S^{(r)}\}_{r=1}^R$

Outline

- Introduction
- Communication
- **Fault tolerance**
- Evaluation
- Conclusions

SV-based Representation

- SV-based representation of parameters
 - At iteration T , the state W_T of the parameter matrix is

$$\begin{array}{l} \text{Initialization} \qquad \qquad \text{Update matrices} \\ \downarrow \qquad \qquad \qquad \downarrow \\ W_T = W_0 + \overline{\Delta W_1 + \dots + \Delta W_T} \\ \downarrow \qquad \qquad \qquad \downarrow \qquad \qquad \downarrow \\ \text{SV Representation (SVR)} \quad W_T = u_0 v_0^T + u_1 v_1^T \quad \dots \quad u_T v_T^T \end{array}$$

Fault Tolerance

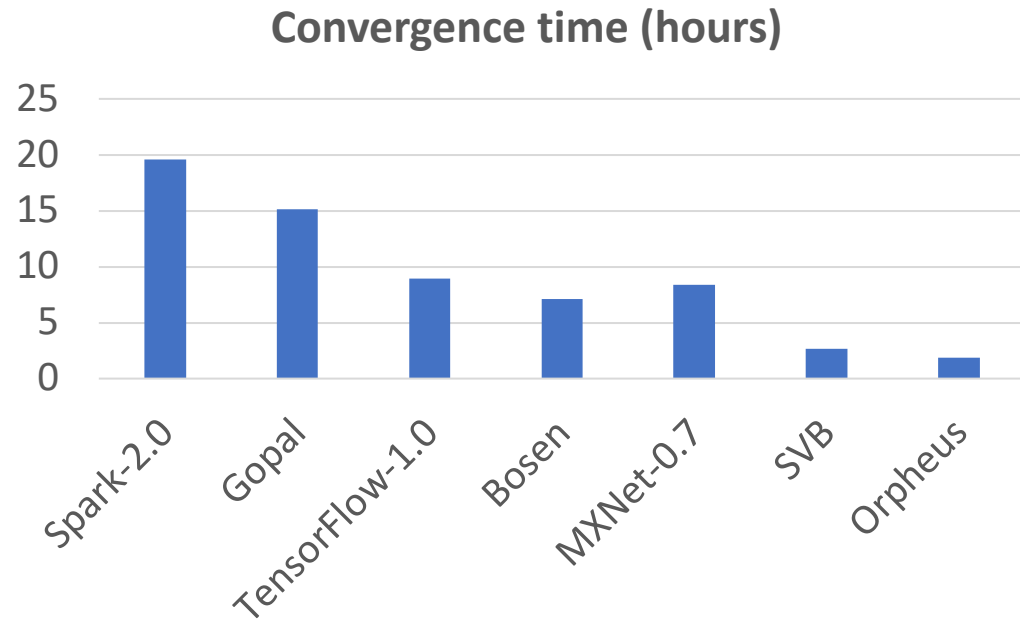
- **SV-based checkpoint:** save SVs computed in each clock on disk
 - Consume little disk bandwidth
 - Do not halt computation
- **Recovery:** transform saved SVs into parameter matrix
 - Can rollback to the state of every clock

Outline

- Introduction
- Communication
- Fault tolerance
- **Evaluation**
- Conclusions

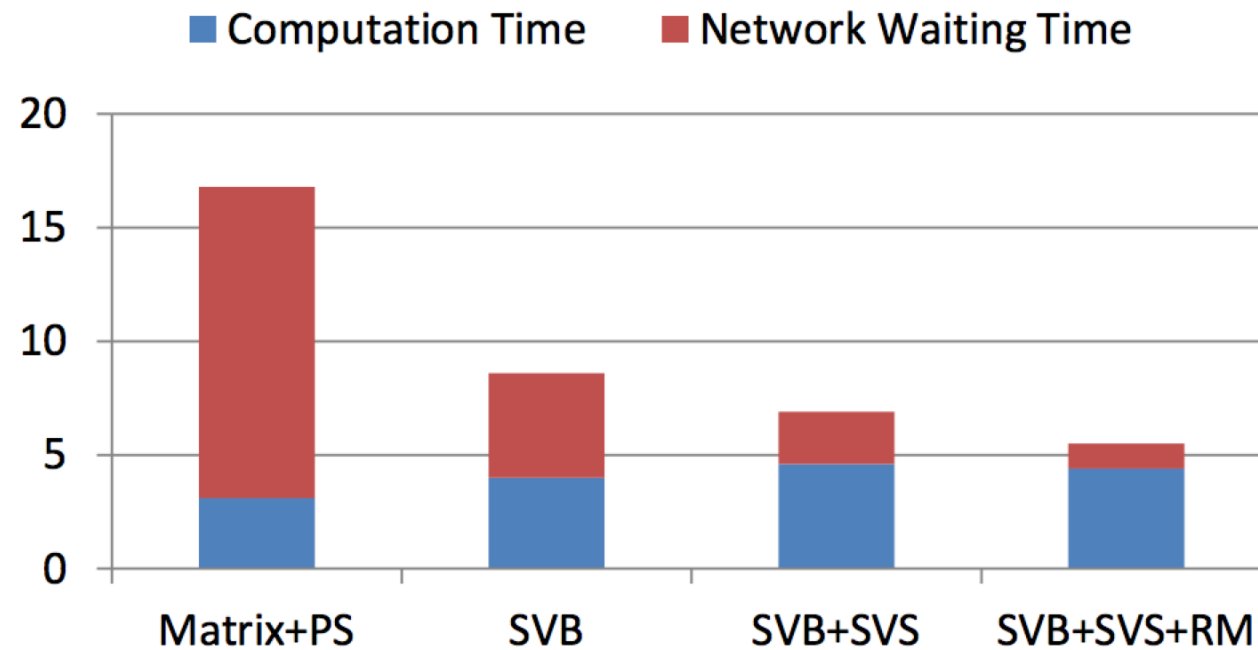
Convergence Speed

Multi-class Logistic Regression (MLR)
Weight matrix: 325K-by-20K

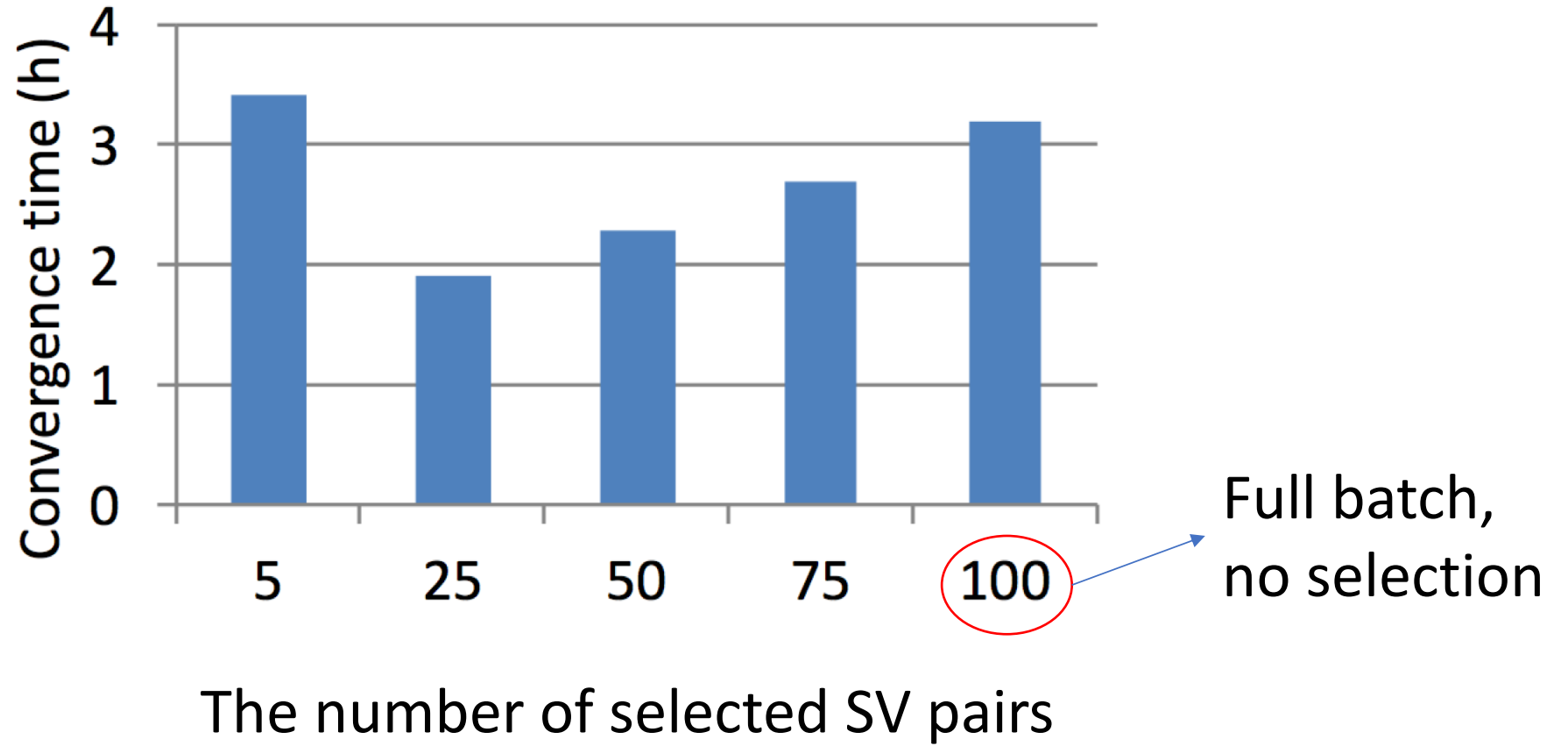


	Convergence time (h)
Matrix+PS	7.1
SVB	2.7
SVB+SVS	2.3
SVB+SVS+RM	2.1

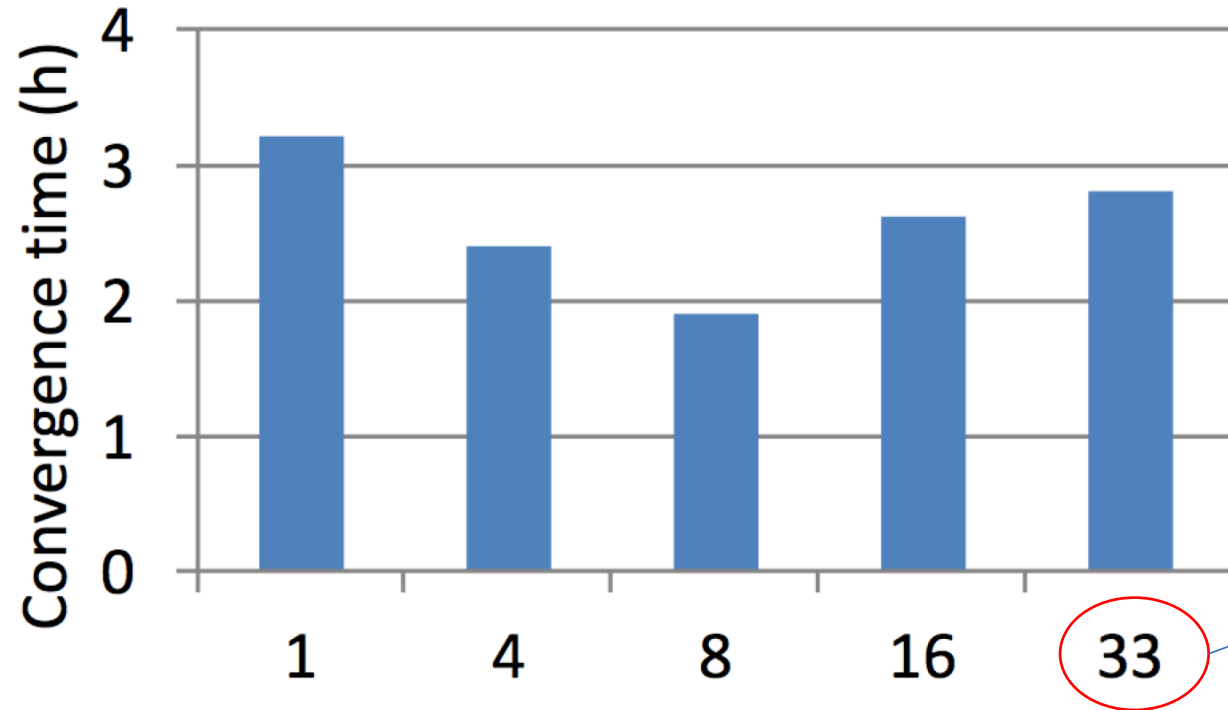
Breakdown of Network Waiting Time and Computation Time



SV Selection



Random Multicast



Full broadcast,
no selection

The number of destinations each
machine sends messages to

Fault Tolerance

	MLR
No checkpoint	1.9
Matrix-based checkpoint	2.4
SV-based checkpoint	1.9
Matrix-based recovery	2.9
SV-based recovery	2.3

Conclusions

