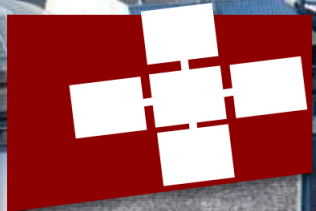


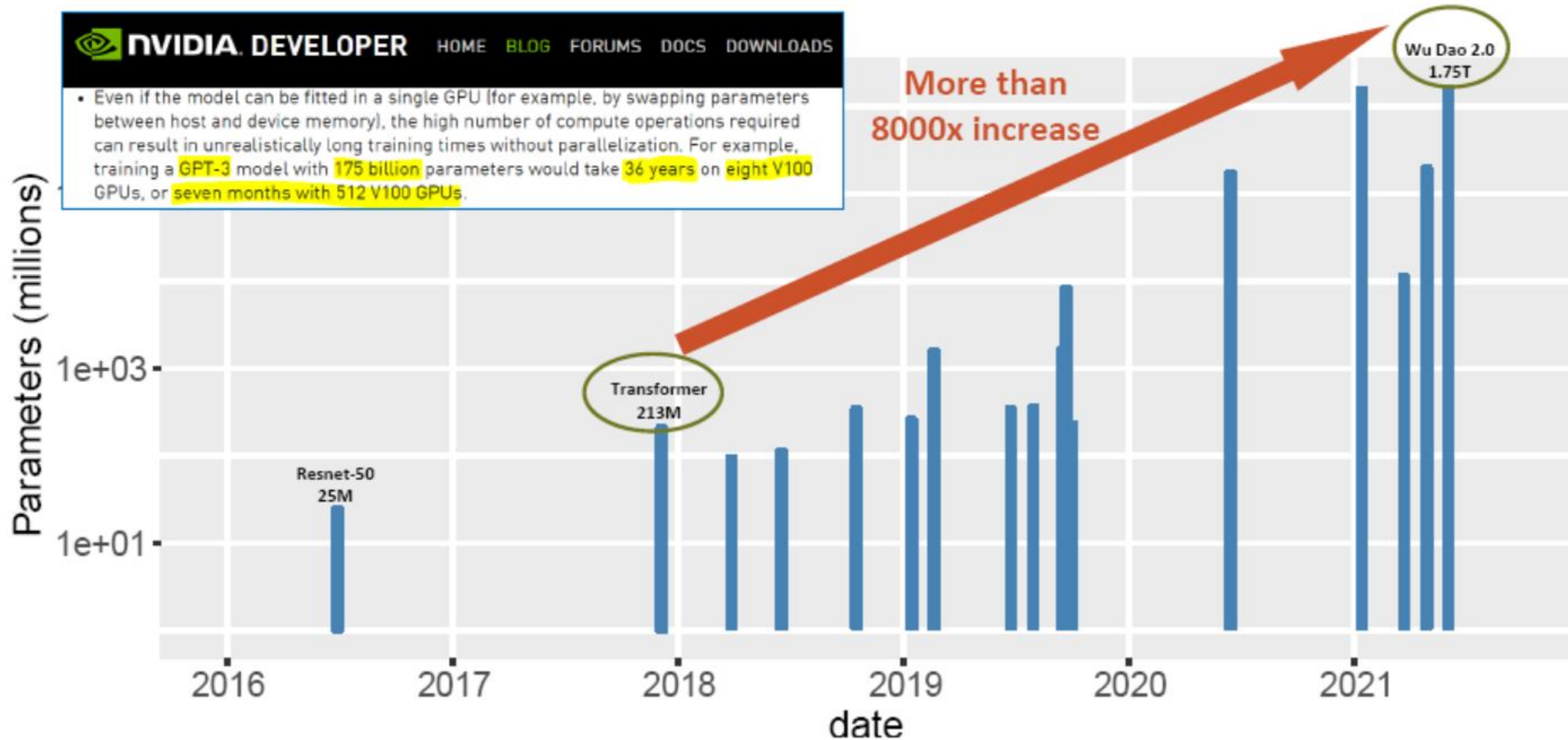
Near-Optimal Sparse Allreduce for Distributed Deep Learning

Shigang Li, Torsten Hoefler
ETH Zurich



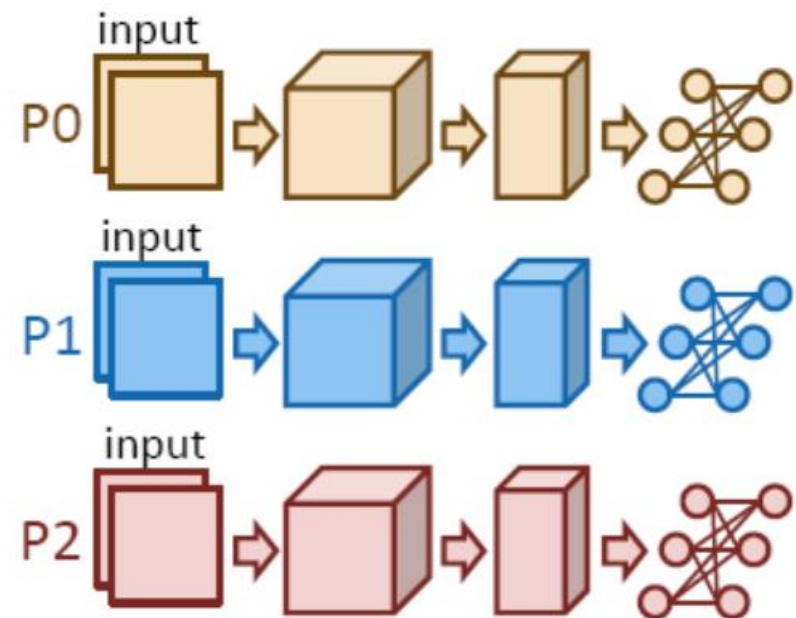
PPoPP'22
April 2-6, 2022
Seoul, Republic of Korea

Model size growing exponentially



Parallel and distributed training

Data parallelism



Pros:

- a. Easy to realize

Cons:

- a. Using DP alone may NOT work for large models, but works with others (**OP, PP**)
- b. **High allreduce overhead**

This work (**Ok-Topk**) aims to solve

Operator parallelism



Pros:

- a. Make large model training feasible

Cons:

- b. Communication for each operator (or each layer)

Pipeline parallelism



Pros:

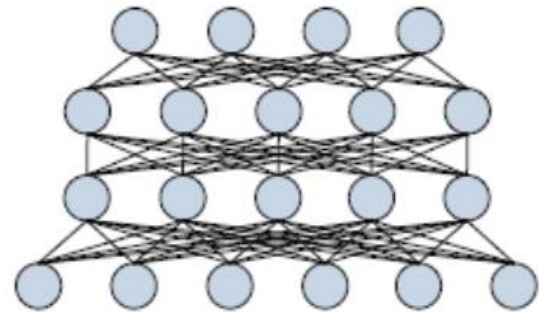
- a. Make large model training feasible
- b. No collective, only P2P

Cons:

- a. Bubbles in pipeline
- b. Removing bubbles leads to stale weights

We have **Chimera** (SC'21) to solve the above issues for pipeline parallelism.

Data parallelism



$$f(\mathbf{w}) = \mathbb{E}_{\xi \sim D} F(\mathbf{w}; \xi)$$

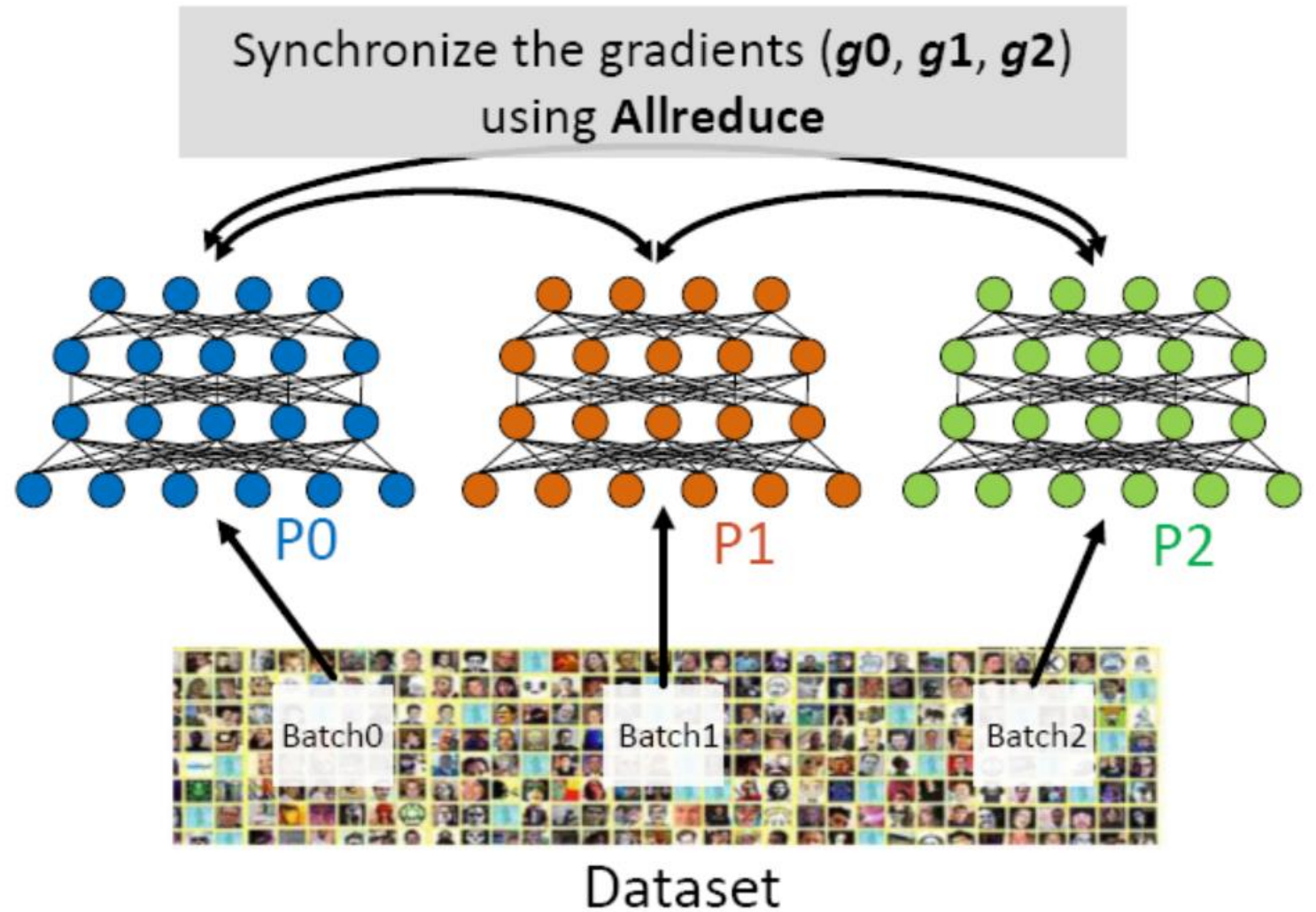
\mathbf{w} denotes the model parameters.

F is the loss function.

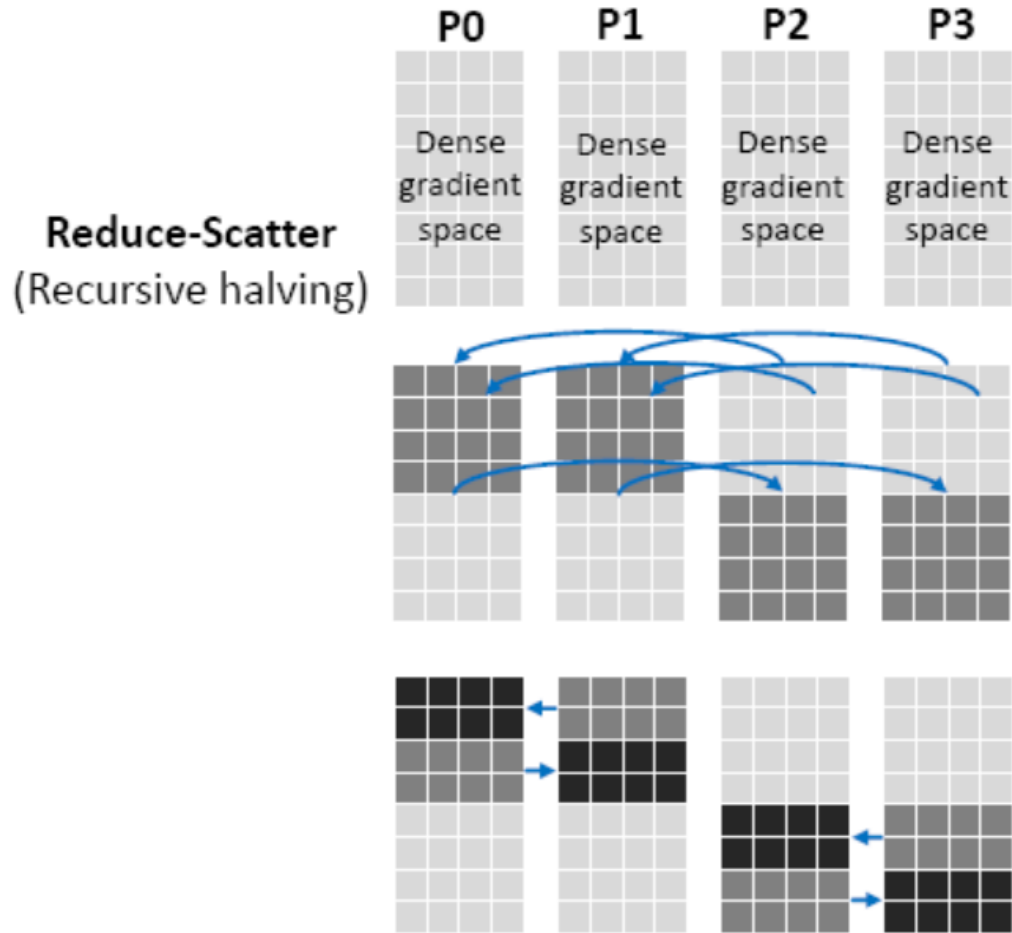
ξ is a data batch sampled from a distribution D .

Training: update \mathbf{w} to minimize f (e.g., SGD).

$$g_t = \frac{1}{b} \sum_{i=0}^b \nabla F(\mathbf{w}_t, \xi_i) \quad \mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t g_t$$

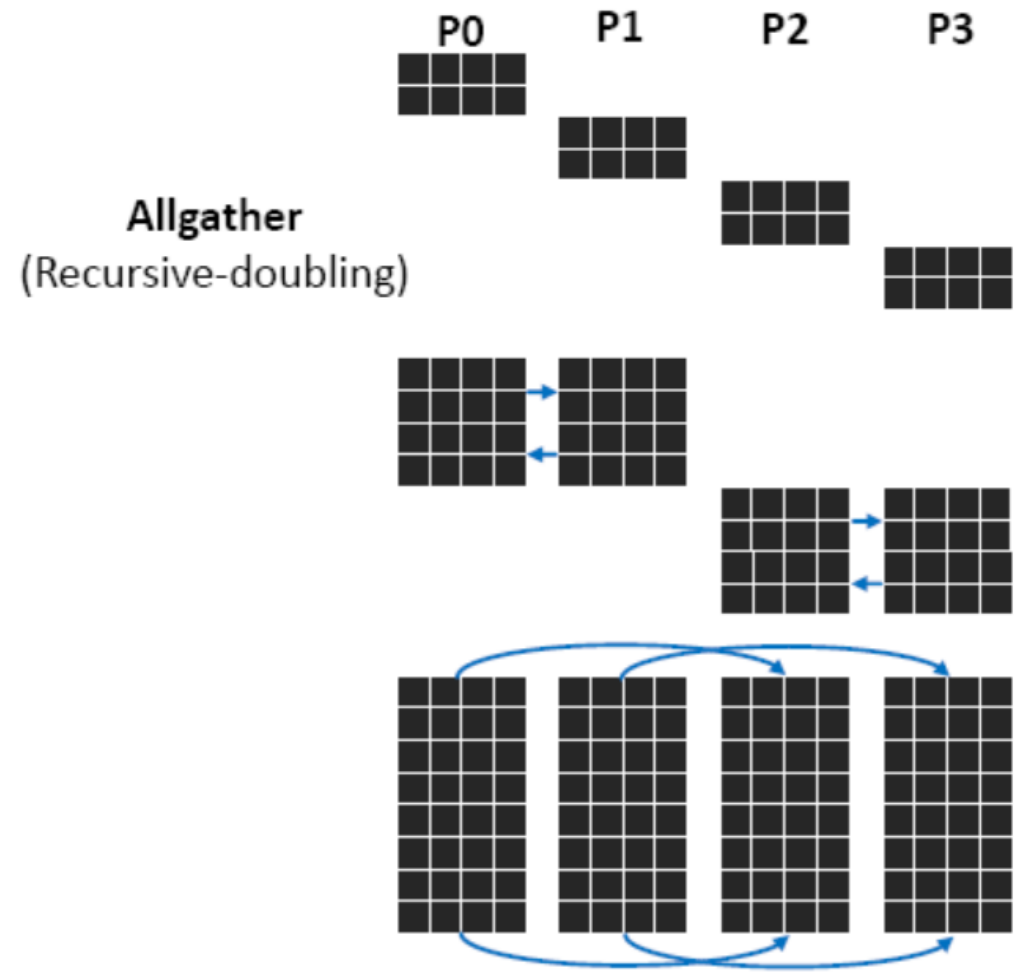


Revisit Dense Allreduce (Rabenseifner's algorithm)



Latency-bandwidth model:

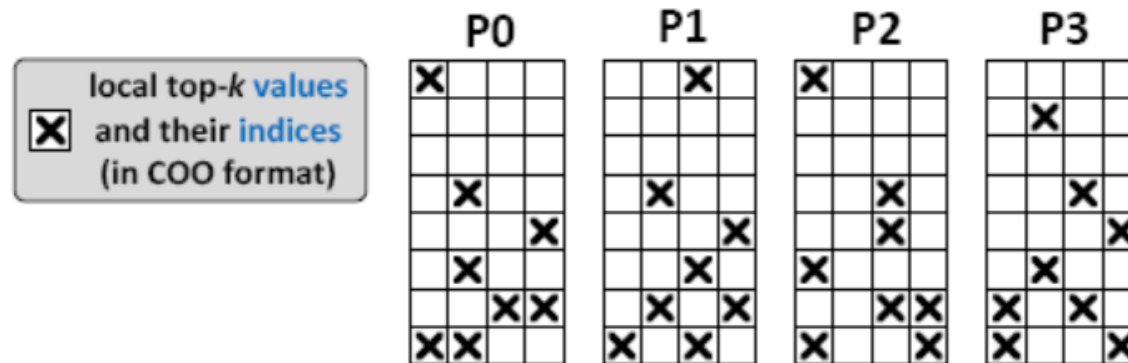
α is the latency; β is the transfer time per word;
 p is the number of processes; n is the message size.



Comm. Cost = $2(\lg p)\alpha + 2((p-1)/p)n\beta$
Bandwidth optimal and scalable, but w.r.t. n

Gradient sparsification (Topk SGD)

Topk SGD: each process only selects the largest (absolute value) k of n components from the gradients, and usually the density k/n is around 1% or less.

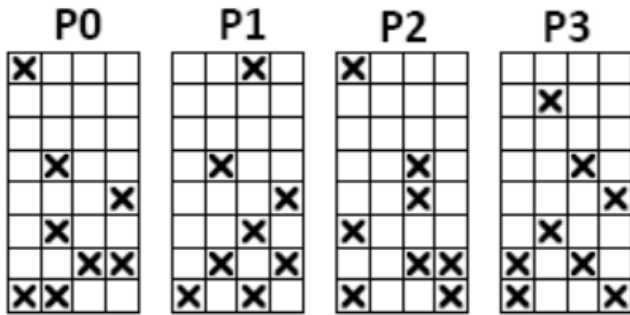


How to **Allreduce** these **sparse** gradients?

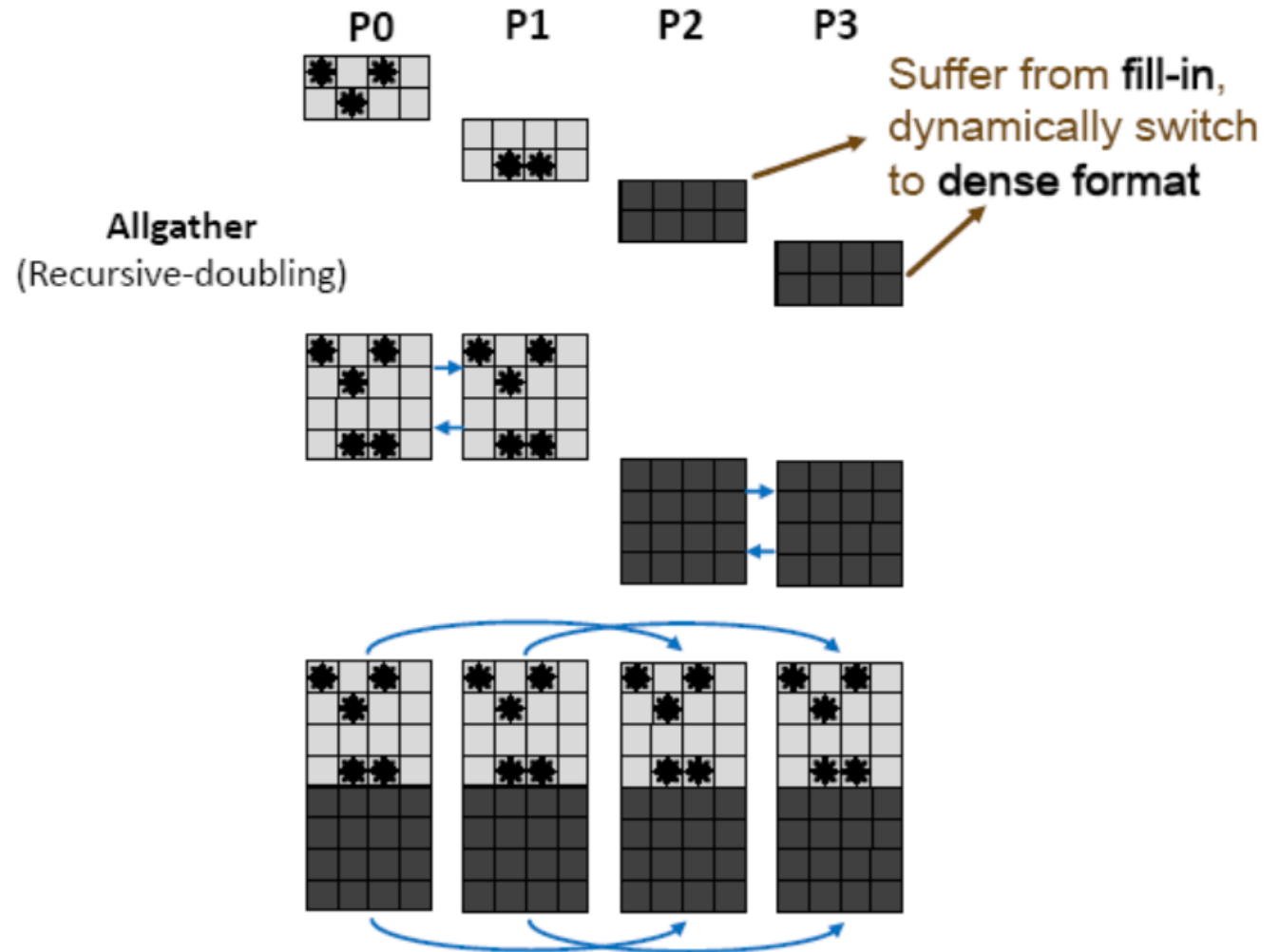
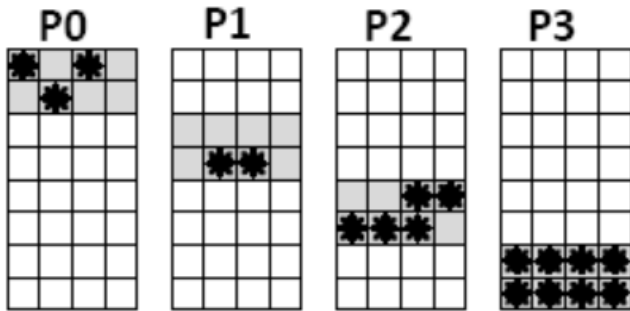
Algorithm 2: Dynamic Sparse Allreduce (TopkDSA, SC'19)

Inspired by dense Rabenseifner's algorithm

local top- k values and their indices (in COO format)

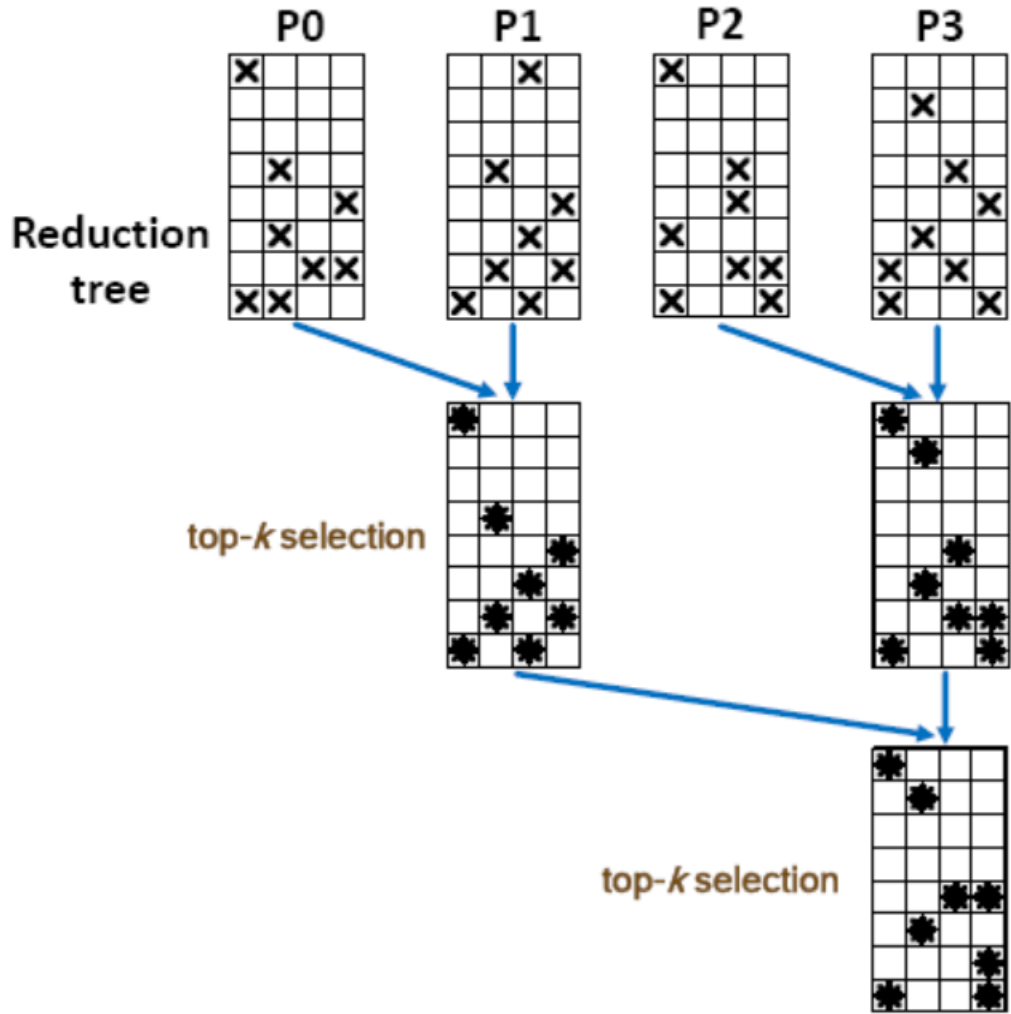


reduced top- k values and their indices (in COO format)

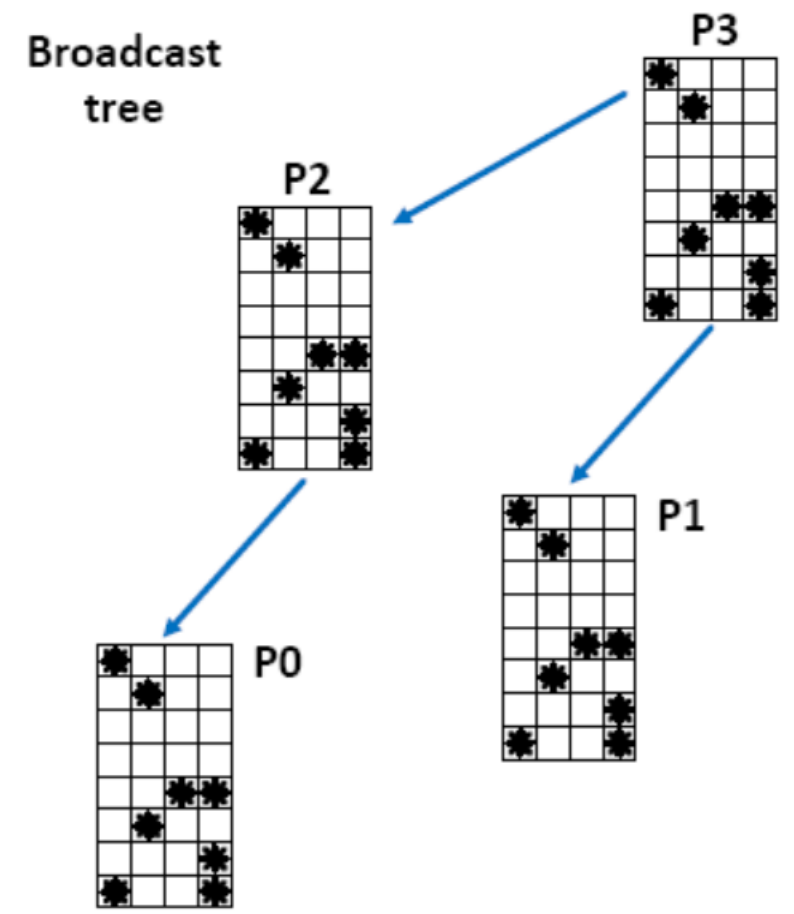


Comm. Cost = $(p + 2 \lg p) \alpha + [(4k(p-1)/p) \beta, ((2k+n)(p-1)/p) \beta]$ Suffer from fill-in issue when scaling up!

Algorithm 3: Global Topk (gTopk, ICDCS'19)



$$\text{Comm. Cost} = (2 \lg p) \alpha + 4k (\log p) \beta$$

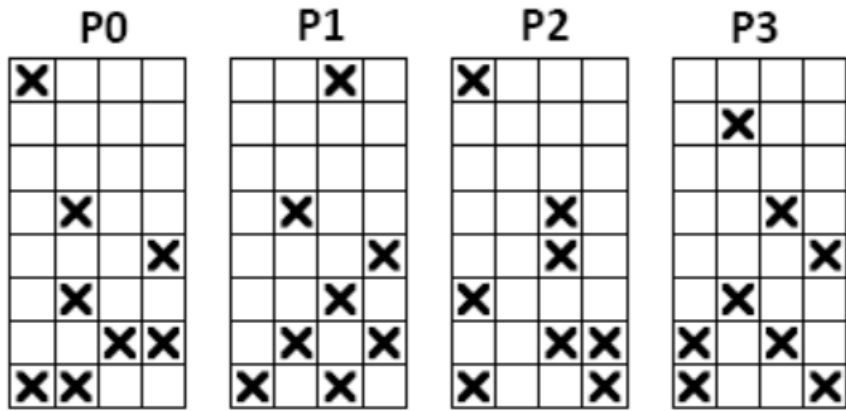


Solve the fill-in issue, but
 (1) high top- k selection cost, and
 (2) suboptimal bandwidth cost.

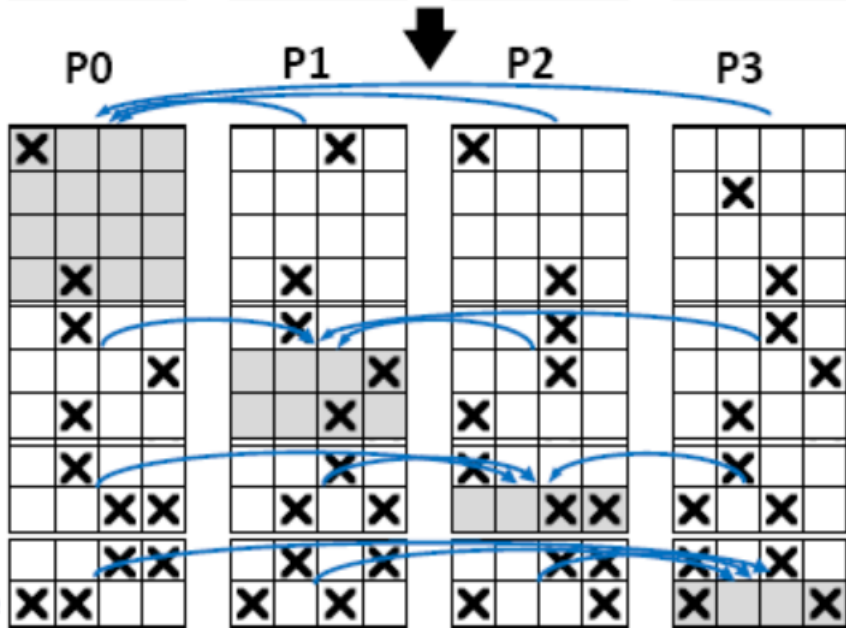
O(k) sparse Allreduce

local top-k values and their indices (in COO format)

Efficient local top-k selection

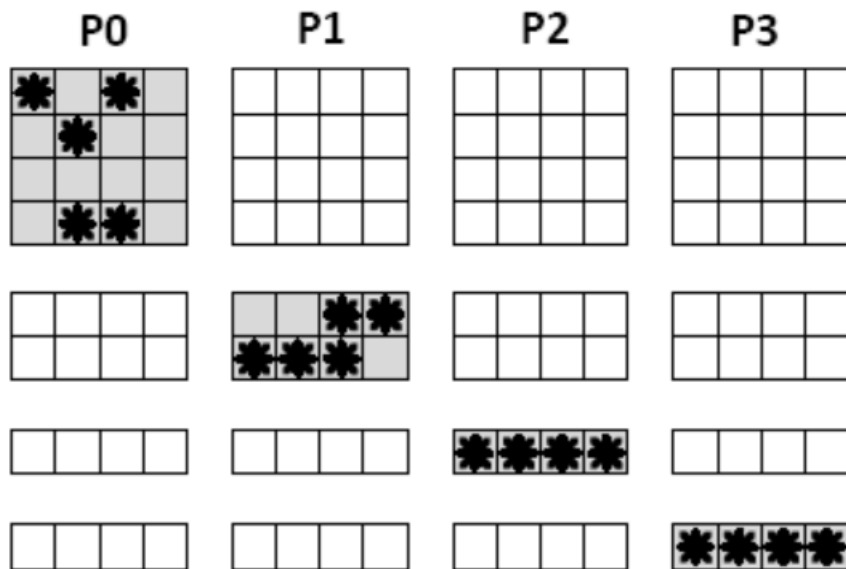


Balanced communication volume



(1) Balanced Split-and-Reduction

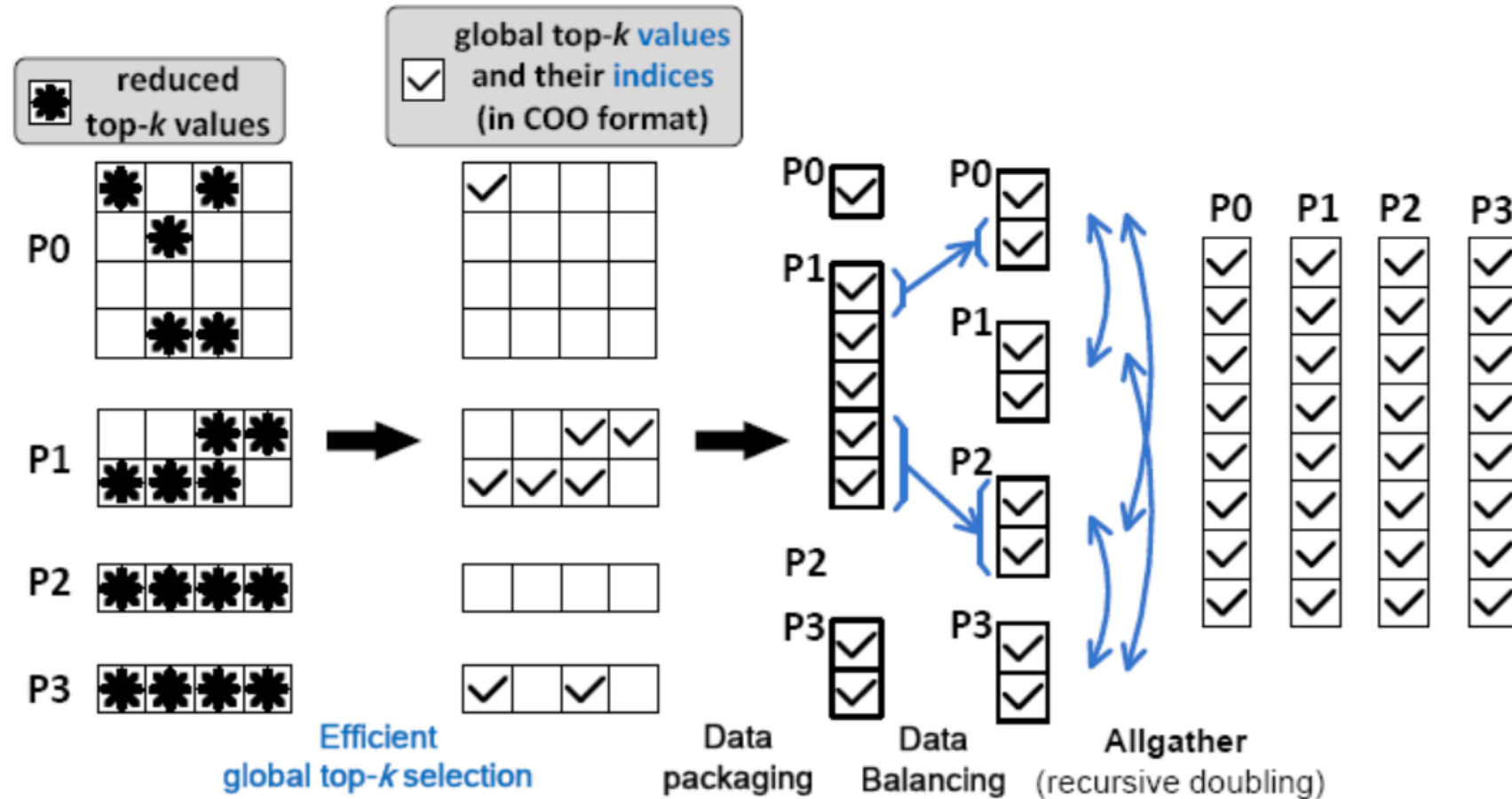
reduced top-k values



The cost of balanced *split_and_reduce* = $(P - 1)\alpha + 2k ((P - 1)/P) \beta$

O(k) sparse Allreduce

(2) Global top-k selection & balanced Allgather

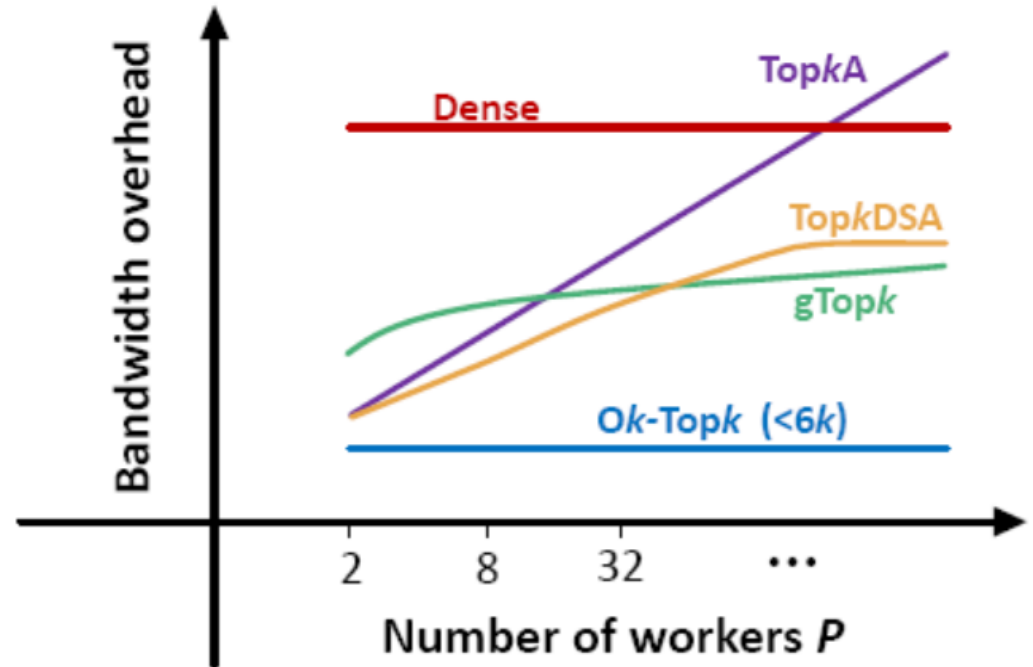


$$\text{Total Cost} \leq (2p + 2\lg p)\alpha + 6k((p-1)/p)\beta$$

Less than $6k$, **asymptotically optimal**

Scalability analysis for dense/sparse Allreduce algorithms

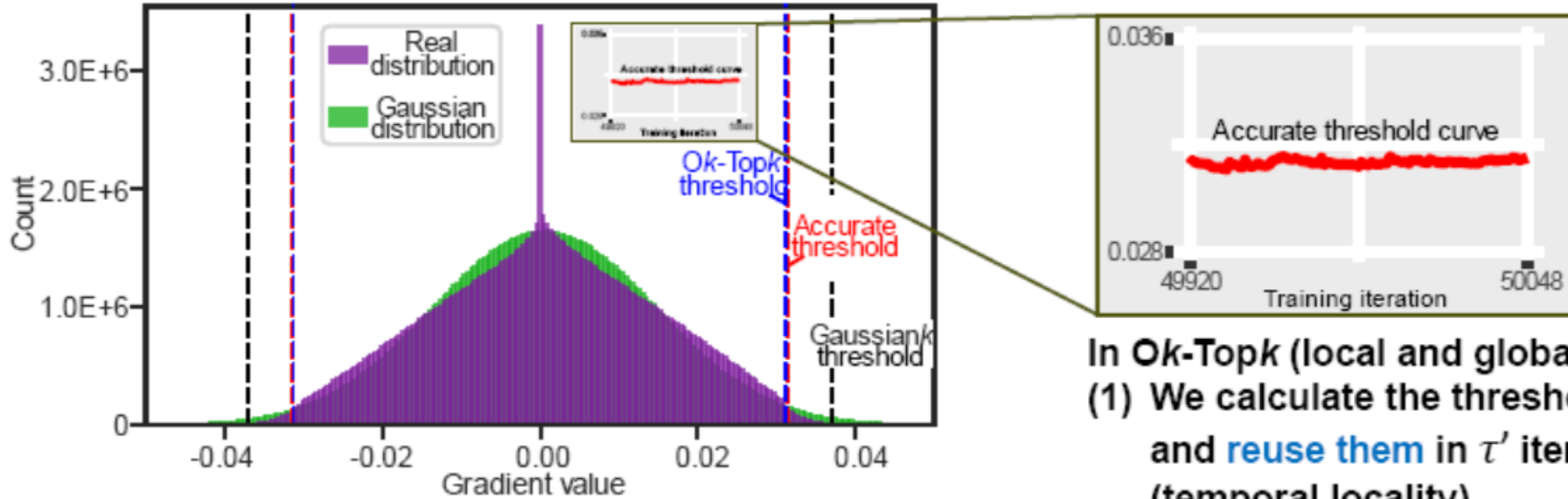
Algorithms	Bandwidth	Scalability
Dense [12]	$2n \frac{P-1}{P} \beta$	
TopkA [36, 47]	$2k(P-1)\beta$	👎👎
TopkDSA [36]	$[4k \frac{P-1}{P} \beta, (2k + n) \frac{P-1}{P} \beta]^1$	👎
gTopk [42]	$4k(\log P)\beta$	👎
Gaussiank [41]	$2k(P-1)\beta$	👎👎
Ok-Topk (ours)	$[2k \frac{P-1}{P} \beta, 6k \frac{P-1}{P} \beta]^1$	👍



Existing sparse Allreduce algorithms suffer from **scalability issue**.
Ok-Topk solves the issue!

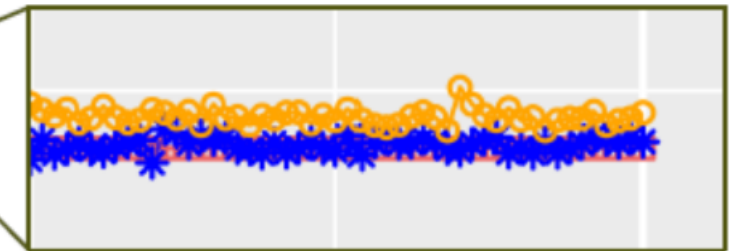
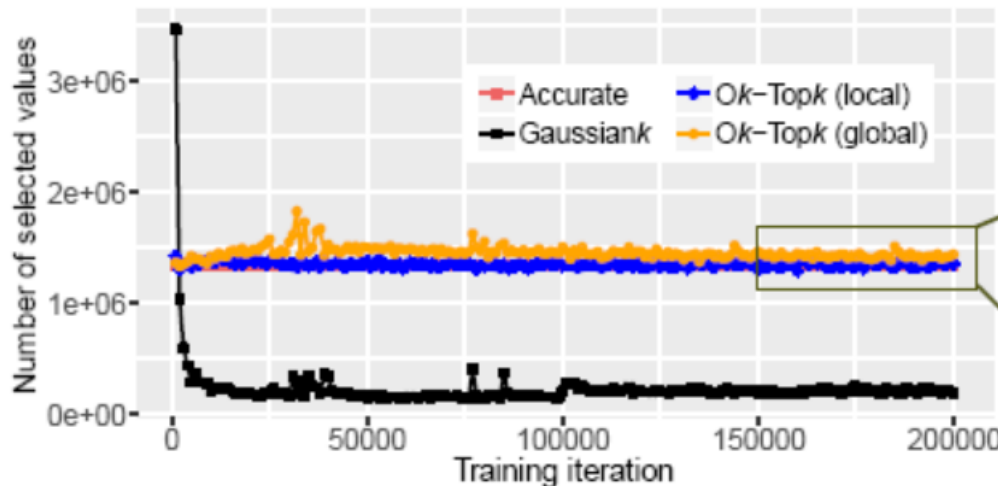
Efficient local and global top- k selection

BERT on Wikipedia
with density = 1.0%



- In **Ok-Topk** (local and global top- k):
- (1) We calculate the thresholds once and **reuse them** in τ' iterations (temporal locality)
 - (2) $O(n)$ overhead rather than $O(n \log n)$ in sorting based top- k selection
 - (3) More **friendly to GPU** than sorting

BERT on Wikipedia on 32 GPUs
(density=1.0%, $\tau' = 128$)



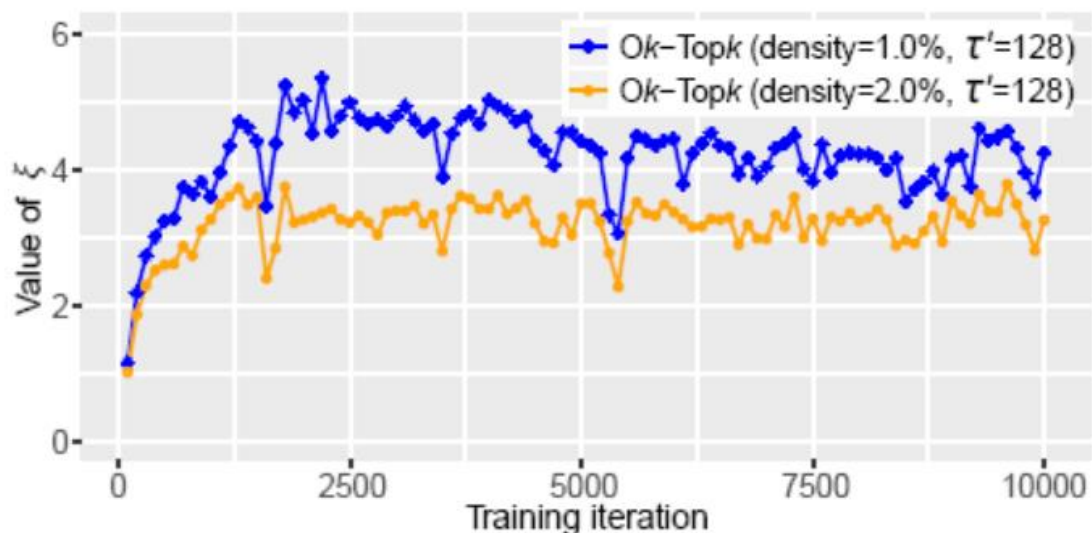
Ok-Topk parallel SGD algorithm

-
- 1: **Inputs:** stochastic gradient $G^i(\cdot)$ at worker i , value k , learning rate α .
 - 2: Initialize $\epsilon_0^i = 0, G_0^i = 0$
 - 3: **for** $t = 1$ **to** T **do**
 - 4: $acc_t^i = \epsilon_{t-1}^i + \alpha G_{t-1}^i(w_{t-1})$ \triangleright Accumulate residuals
 - 5: $u_t, indexes = Ok_sparse_allreduce(acc_t^i, t, k)$
 - 6: $\epsilon_t^i = acc_t^i - acc_t^i(indexes)$ \triangleright Update residuals
 - 7: $w_t = w_{t-1} - \frac{1}{p} u_t$ \triangleright Apply the model update
 - 8: **end for**
-

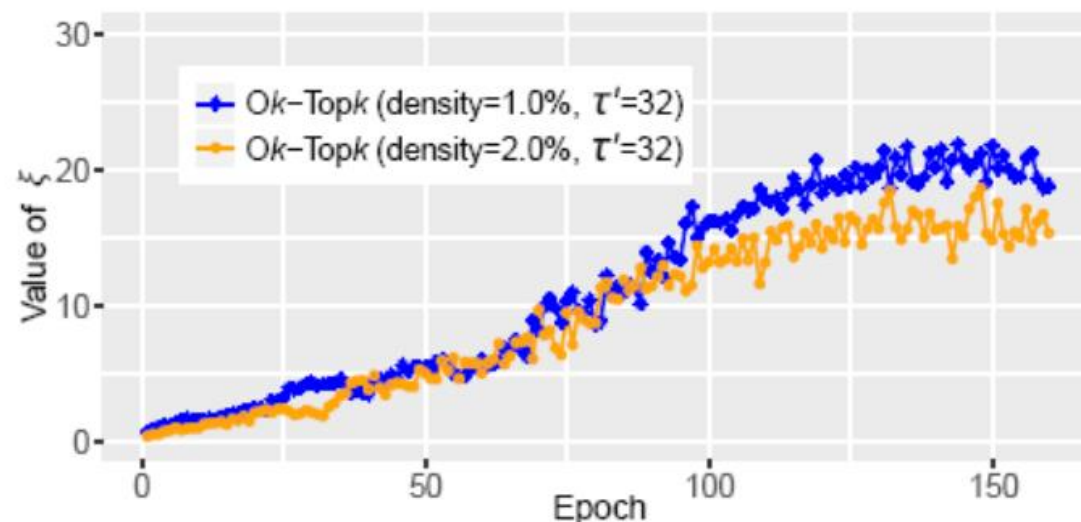
Convergence analysis for Ok-Topk SGD

$$\left\| \underbrace{\text{Topk}\left(\frac{1}{P} \sum_{i=0}^{P-1} \left(\alpha G_t^i(w_t) + \epsilon_t^i\right)\right)}_{\text{True global top-}k \text{ gradient}} - \underbrace{\text{Topk}\left(\frac{1}{P} \sum_{i=0}^{P-1} \text{Topk}\left(\alpha G_t^i(w_t) + \epsilon_t^i\right)\right)}_{\text{Ok-topk gradient}} \right\| \leq \underbrace{\xi}_{\text{Dense gradient}} \|\alpha G_t(w_t)\|$$

For full proof refer to:
 Dan Alistarh, et al., The convergence of sparsified gradient methods, NeurIPS'18



BERT on Wikipedia on 32 GPU nodes



LSTM on AN4 running on 32 GPU nodes

The effect of ξ is dampened by both small learning rates and P .

Evaluation

- **CSCS Piz Daint** supercomputer
 - Each node contains an Intel Xeon E5-2690 CPU, and one NVIDIA Tesla **P100 GPU**
 - **Cray Aries** interconnected network
 - **mpi4py** as the communication library, built against Cray-MPICH 7.7.16



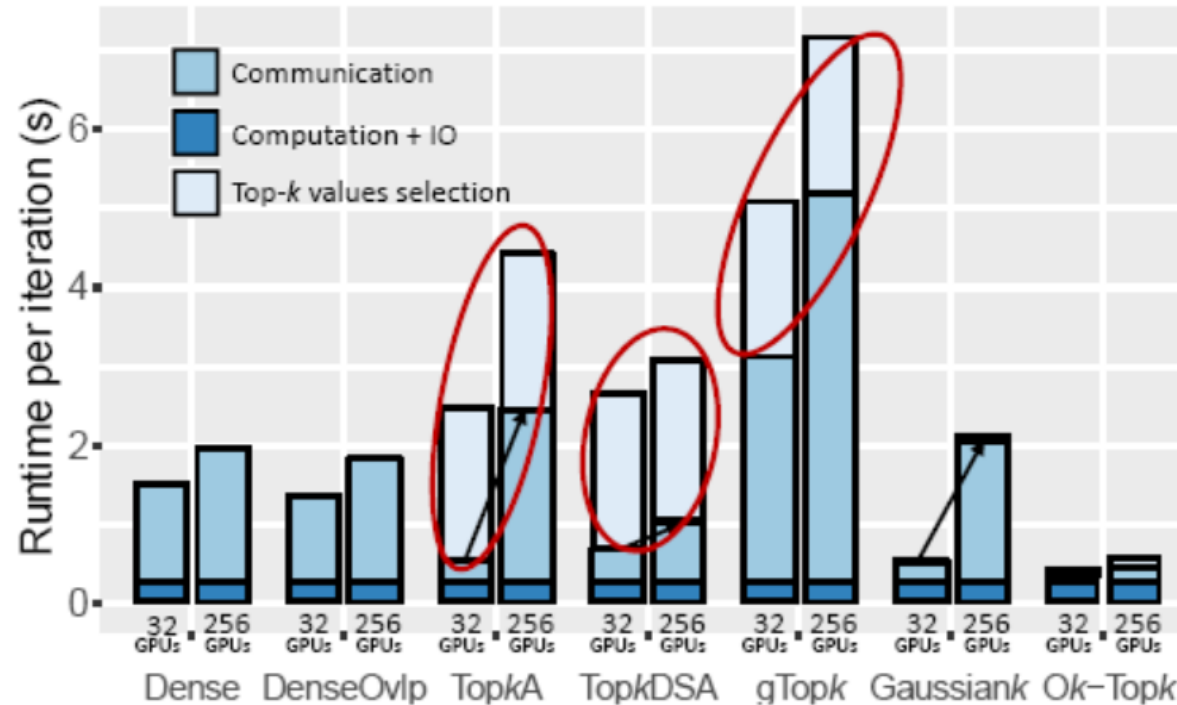
Dense/Sparse algorithms used in evaluation

Algorithms	Bandwidth
Dense [12]	$2n \frac{P-1}{P} \beta$
TopkA [36, 47]	$2k(P-1)\beta$
TopkDSA [36]	$[4k \frac{P-1}{P} \beta, (2k + n) \frac{P-1}{P} \beta]^1$
gTopk [42]	$4k(\log P)\beta$
Gaussiank [41]	$2k(P-1)\beta$
Ok-Topk (ours)	$[2k \frac{P-1}{P} \beta, 6k \frac{P-1}{P} \beta]^1$

Neural networks used for evaluation

Tasks	Models	Parameters	Dataset
Image classification	VGG-16 [44]	14,728,266	Cifar-10
Speech recognition	LSTM [21]	27,569,568	AN4 [1]
Language processing	BERT [13]	133,547,324	Wikipedia [13]

Weak scaling evaluation

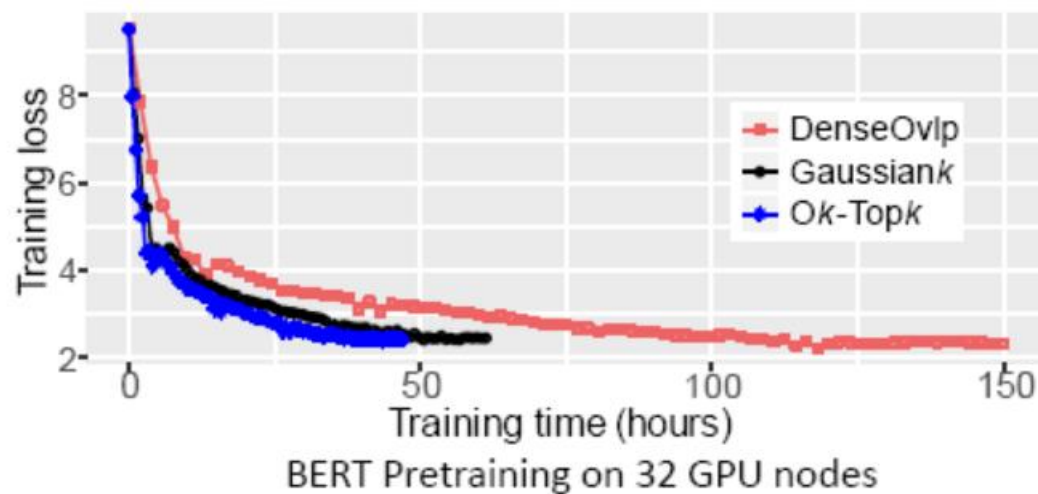
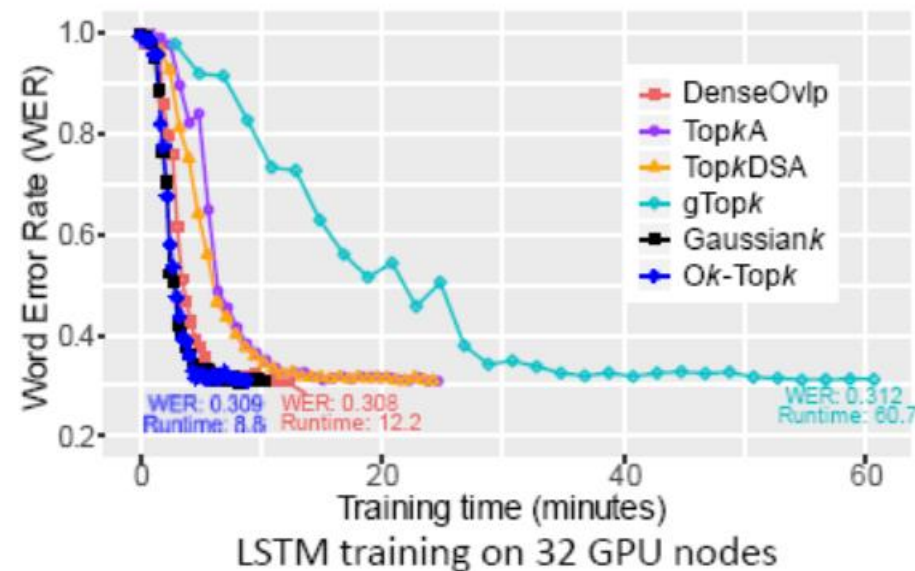
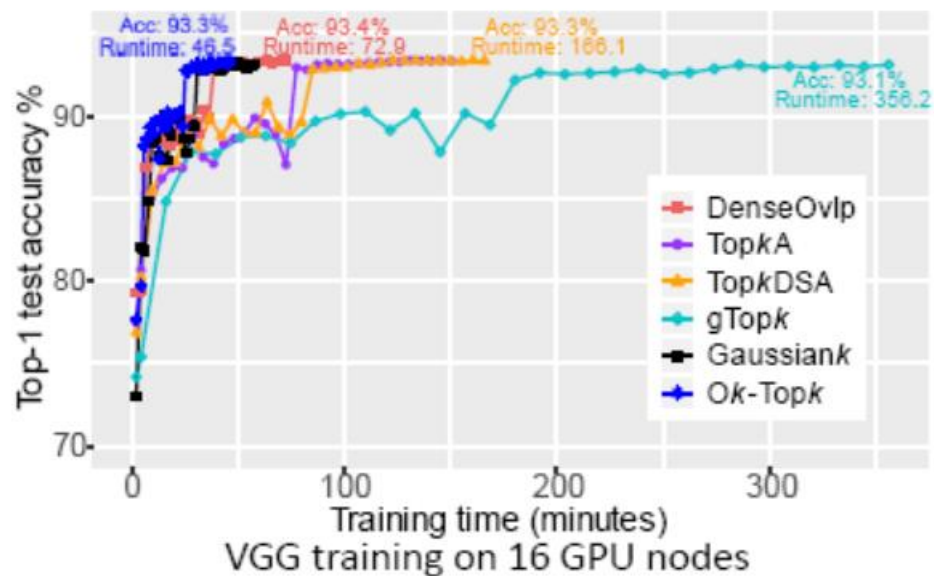


BERT training time,
 scaling from 32 GPUs to 256 GPUs

For Ok-Topk :

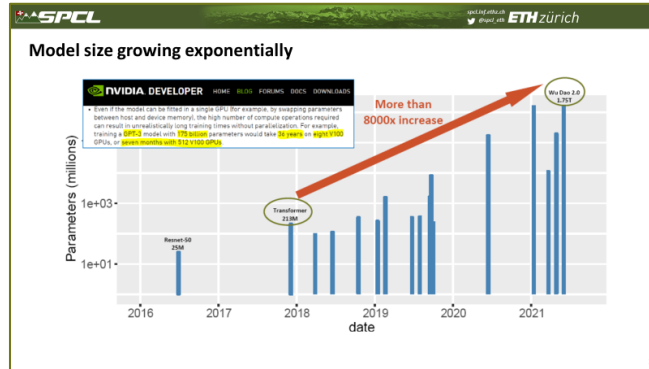
- (1) Much better scalability for the communication overhead than the others.
- (2) Threshold reuse strategy for top-k selection is very effective.

Model accuracy evaluation

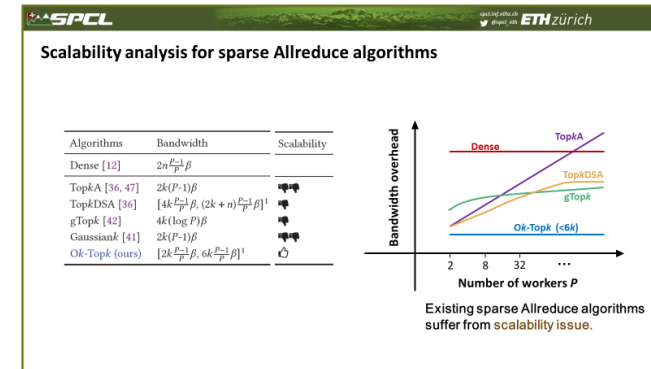


Conclusion

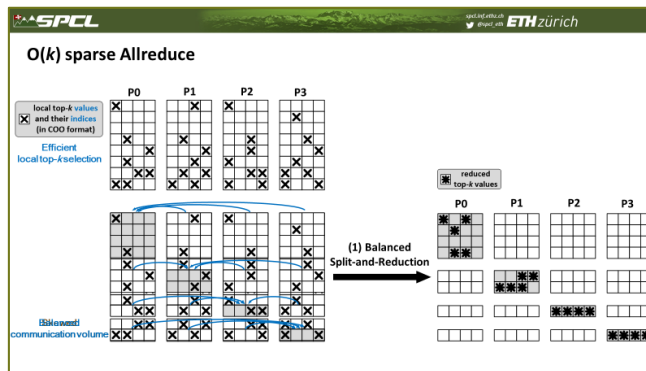
1. Model size rapidly grows



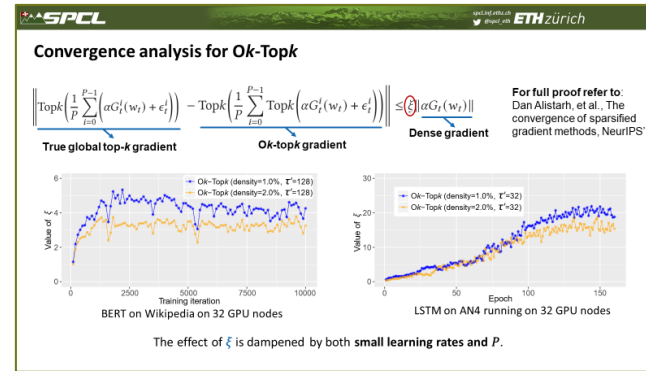
2. Sparse algorithms suffer from scalability issue



3. O(k) sparse Allreduce



4. Ok-Topk SGD and convergence analysis



5. Evaluation on supercomputer

Evaluation

- CS2S Piz Daint supercomputer
 - Each node contains an Intel Xeon E5-2690 CPU, and one NVIDIA Tesla P100 GPU
- Cray Aries interconnected network
- mpi4py as the communication library, built against Cray-MPICH 7.7.16

Dense/Sparse algorithms used in evaluation

Algorithms	Bandwidth
Dense [12]	$2n P^2 \beta$
TopkA [36, 47]	$2k(P-1)\beta$
TopkDSA [36]	$[4k \frac{P-1}{P} \beta, (2k+n) \frac{P-1}{P} \beta]^1$
gTopk [42]	$4k(\log P)\beta$
Gaussiank [41]	$2k(P-1)\beta$
Ok-Topk (ours)	$[2k \frac{P-1}{P} \beta, 6k \frac{P-1}{P} \beta]^1$

Neural networks used for evaluation

Tasks	Models	Parameters	Dataset
Image classification	VGG-16 [44]	14,728,266	Cifar-10
Speech recognition	LSTM [21]	27,569,568	AN4 [11]
Language processing	BERT [13]	133,547,324	Wikipedia [13]

6. For the future work, we will study how to use Ok-Topk with a hybrid data and pipeline parallelism.

For more questions: shigangli.cs@gmail.com