

Scaling Stratified Stochastic Gradient Descent for Distributed Matrix Completion

Nabil Abubaker

Nabil Abubaker, M. Ozan Karsavuran and Cevdet Aykanat, “**Scaling Stratified Stochastic Gradient Descent for Distributed Matrix Completion**”, *IEEE Transactions on Knowledge and Data Engineering*, 2023

Code: <https://github.com/nfabubaker/CESSGD>



Bilkent University

Scaling Stratified Stochastic Gradient Descent for
Distributed Matrix Completion

Communication-Efficient Stratified SGD

Nabil Abubaker

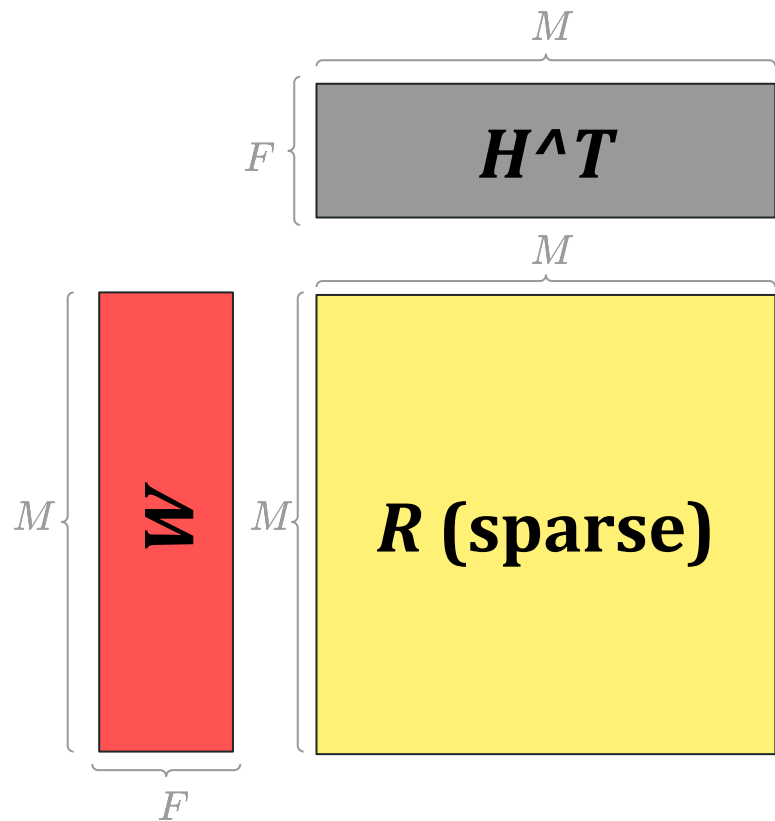
Nabil Abubaker, M. Ozan Karsavuran and Cevdet Aykanat, “**Scaling Stratified Stochastic Gradient Descent for Distributed Matrix Completion**”, *IEEE Transactions on Knowledge and Data Engineering*, 2023

Code: <https://github.com/nfabubaker/CESSGD>



Bilkent University

SGD for Matrix Completion



Goal:

Find low-rank approximation $R \cong WH^T$

A missing entry r_{ij} in R can be approximated (completed) by $r_{ij} = \mathbf{w}_i \mathbf{h}_j^T$

How?

Using Stochastic Gradient Descent to find W and H that minimize:

$$Loss = \sum (r_{ij} - \mathbf{w}_i \mathbf{h}_j^T)^2$$

By:

$$\mathbf{w}_i = \mathbf{w}_i + \epsilon ((r_{ij} - \mathbf{w}_i \mathbf{h}_j^T) \mathbf{h}_j + \gamma \mathbf{w}_i)$$

Step size ϵ and Regularization factor γ are indicated by red arrows.

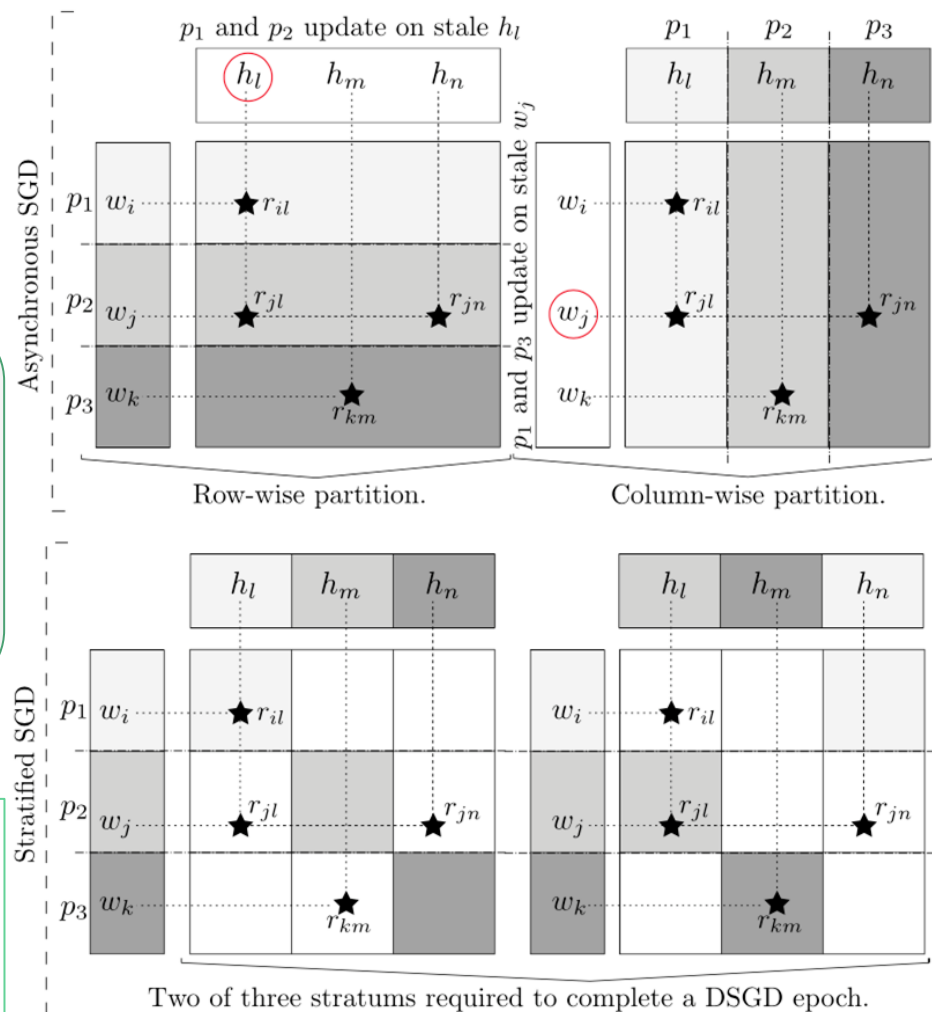
$$\mathbf{h}_j = \mathbf{h}_j + \epsilon ((r_{ij} - \mathbf{w}_i \mathbf{h}_j^T) \mathbf{w}_i + \gamma \mathbf{h}_j)$$

Parallel SGD

- ❖ **Asynchronous**: Allows staleness
- ❖ **Stratified (SSGD)**: Doesn't allow staleness
 - Serializable
 - Better convergence
 - Well-studied behaviour

Stratified SGD proposed in:

R. Gemulla, E. Nijkamp, P. J. Haas, and Y. Sismanis, "Large-scale matrix factorization with distributed stochastic gradient descent," in Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining, pp. 69–77, 2011.



s1

p_1 / s_1	p_1 / s_2	p_1 / s_3	p_1 / s_4
p_2 / s_4	p_2 / s_1	p_2 / s_2	p_2 / s_3
p_3 / s_3	p_3 / s_4	p_3 / s_1	p_3 / s_2
p_4 / s_2	p_4 / s_3	p_4 / s_4	p_4 / s_1

p_1 / s_1	p_1 / s_2	p_1 / s_3	p_1 / s_4
p_2 / s_4	p_2 / s_1	p_2 / s_2	p_2 / s_3
p_3 / s_3	p_3 / s_4	p_3 / s_1	p_3 / s_2
p_4 / s_2	p_4 / s_3	p_4 / s_4	p_4 / s_1

s2

p_1 / s_1	p_1 / s_2	p_1 / s_3	p_1 / s_4
p_2 / s_4	p_2 / s_1	p_2 / s_2	p_2 / s_3
p_3 / s_3	p_3 / s_4	p_3 / s_1	p_3 / s_2
p_4 / s_2	p_4 / s_3	p_4 / s_4	p_4 / s_1

p_1 / s_1	p_1 / s_2	p_1 / s_3	p_1 / s_4
p_2 / s_4	p_2 / s_1	p_2 / s_2	p_2 / s_3
p_3 / s_3	p_3 / s_4	p_3 / s_1	p_3 / s_2
p_4 / s_2	p_4 / s_3	p_4 / s_4	p_4 / s_1

p_1 / s_1	p_1 / s_2	p_1 / s_3	p_1 / s_4
p_2 / s_4	p_2 / s_1	p_2 / s_2	p_2 / s_3
p_3 / s_3	p_3 / s_4	p_3 / s_1	p_3 / s_2
p_4 / s_2	p_4 / s_3	p_4 / s_4	p_4 / s_1

s3

p_1 / s_1	p_1 / s_2	p_1 / s_3	p_1 / s_4
p_2 / s_4	p_2 / s_1	p_2 / s_2	p_2 / s_3
p_3 / s_3	p_3 / s_4	p_3 / s_1	p_3 / s_2
p_4 / s_2	p_4 / s_3	p_4 / s_4	p_4 / s_1

p_1 / s_1	p_1 / s_2	p_1 / s_3	p_1 / s_4
p_2 / s_4	p_2 / s_1	p_2 / s_2	p_2 / s_3
p_3 / s_3	p_3 / s_4	p_3 / s_1	p_3 / s_2
p_4 / s_2	p_4 / s_3	p_4 / s_4	p_4 / s_1

p_1 / s_1	p_1 / s_2	p_1 / s_3	p_1 / s_4
p_2 / s_4	p_2 / s_1	p_2 / s_2	p_2 / s_3
p_3 / s_3	p_3 / s_4	p_3 / s_1	p_3 / s_2
p_4 / s_2	p_4 / s_3	p_4 / s_4	p_4 / s_1

s4

p_1 / s_1	p_1 / s_2	p_1 / s_3	p_1 / s_4
p_2 / s_4	p_2 / s_1	p_2 / s_2	p_2 / s_3
p_3 / s_3	p_3 / s_4	p_3 / s_1	p_3 / s_2
p_4 / s_2	p_4 / s_3	p_4 / s_4	p_4 / s_1

p_1 / s_1	p_1 / s_2	p_1 / s_3	p_1 / s_4
p_2 / s_4	p_2 / s_1	p_2 / s_2	p_2 / s_3
p_3 / s_3	p_3 / s_4	p_3 / s_1	p_3 / s_2
p_4 / s_2	p_4 / s_3	p_4 / s_4	p_4 / s_1

s1

p_1 / s_1	p_1 / s_2	p_1 / s_3	p_1 / s_4
p_2 / s_4	p_2 / s_1	p_2 / s_2	p_2 / s_3
p_3 / s_3	p_3 / s_4	p_3 / s_1	p_3 / s_2
p_4 / s_2	p_4 / s_3	p_4 / s_4	p_4 / s_1

s2

p_1 / s_1	p_1 / s_2	p_1 / s_3	p_1 / s_4
p_2 / s_4	p_2 / s_1	p_2 / s_2	p_2 / s_3
p_3 / s_3	p_3 / s_4	p_3 / s_1	p_3 / s_2
p_4 / s_2	p_4 / s_3	p_4 / s_4	p_4 / s_1

Ring
Schedule

Or

Ring
Strata

s3

p_1 / s_1	p_1 / s_2	p_1 / s_3	p_1 / s_4
p_2 / s_4	p_2 / s_1	p_2 / s_2	p_2 / s_3
p_3 / s_3	p_3 / s_4	p_3 / s_1	p_3 / s_2
p_4 / s_2	p_4 / s_3	p_4 / s_4	p_4 / s_1

s4

p_1 / s_1	p_1 / s_2	p_1 / s_3	p_1 / s_4
p_2 / s_4	p_2 / s_1	p_2 / s_2	p_2 / s_3
p_3 / s_3	p_3 / s_4	p_3 / s_1	p_3 / s_2
p_4 / s_2	p_4 / s_3	p_4 / s_4	p_4 / s_1



H_1 H_2 H_3 H_4

s_1

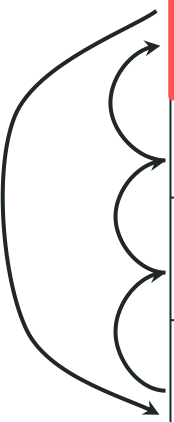
p_1 / s_1	p_1 / s_2	p_1 / s_3	p_1 / s_4
p_2 / s_4	p_2 / s_1	p_2 / s_2	p_2 / s_3
p_3 / s_3	p_3 / s_4	p_3 / s_1	p_3 / s_2
p_4 / s_2	p_4 / s_3	p_4 / s_4	p_4 / s_1



H_1 H_2 H_3 H_4

s_2

p_1 / s_1	p_1 / s_2	p_1 / s_3	p_1 / s_4
p_2 / s_4	p_2 / s_1	p_2 / s_2	p_2 / s_3
p_3 / s_3	p_3 / s_4	p_3 / s_1	p_3 / s_2
p_4 / s_2	p_4 / s_3	p_4 / s_4	p_4 / s_1



Send up-to-date H_1 to p_4

Send up-to-date H_4 to p_3

Problem with distributed SSGD:

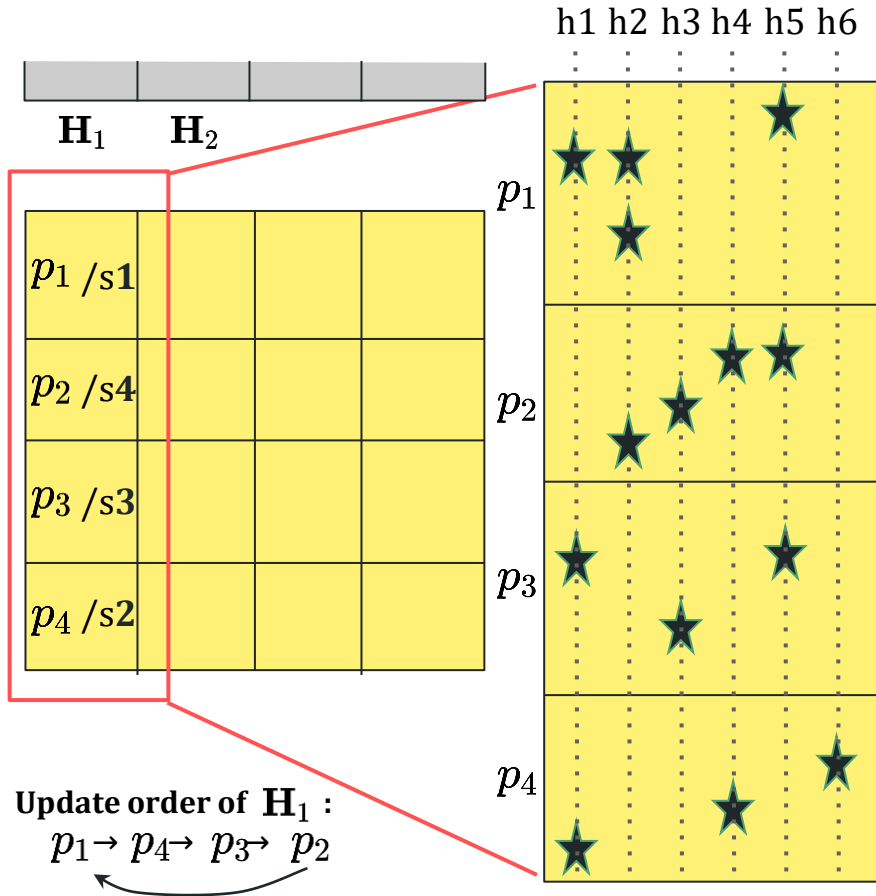
Existing implementations communicate **blocks** of factor matrix rows:

- p_1 updates block H_j
- p_1 sends **all rows in H_j** to the processor that updates it next
- Data is **sparse** → extra **unnecessary** data movement

Proposal:

Communicate only factor matrix rows that are **essential** for the correctness of the SGD algorithm **using P2P messages**

Finding essential communication

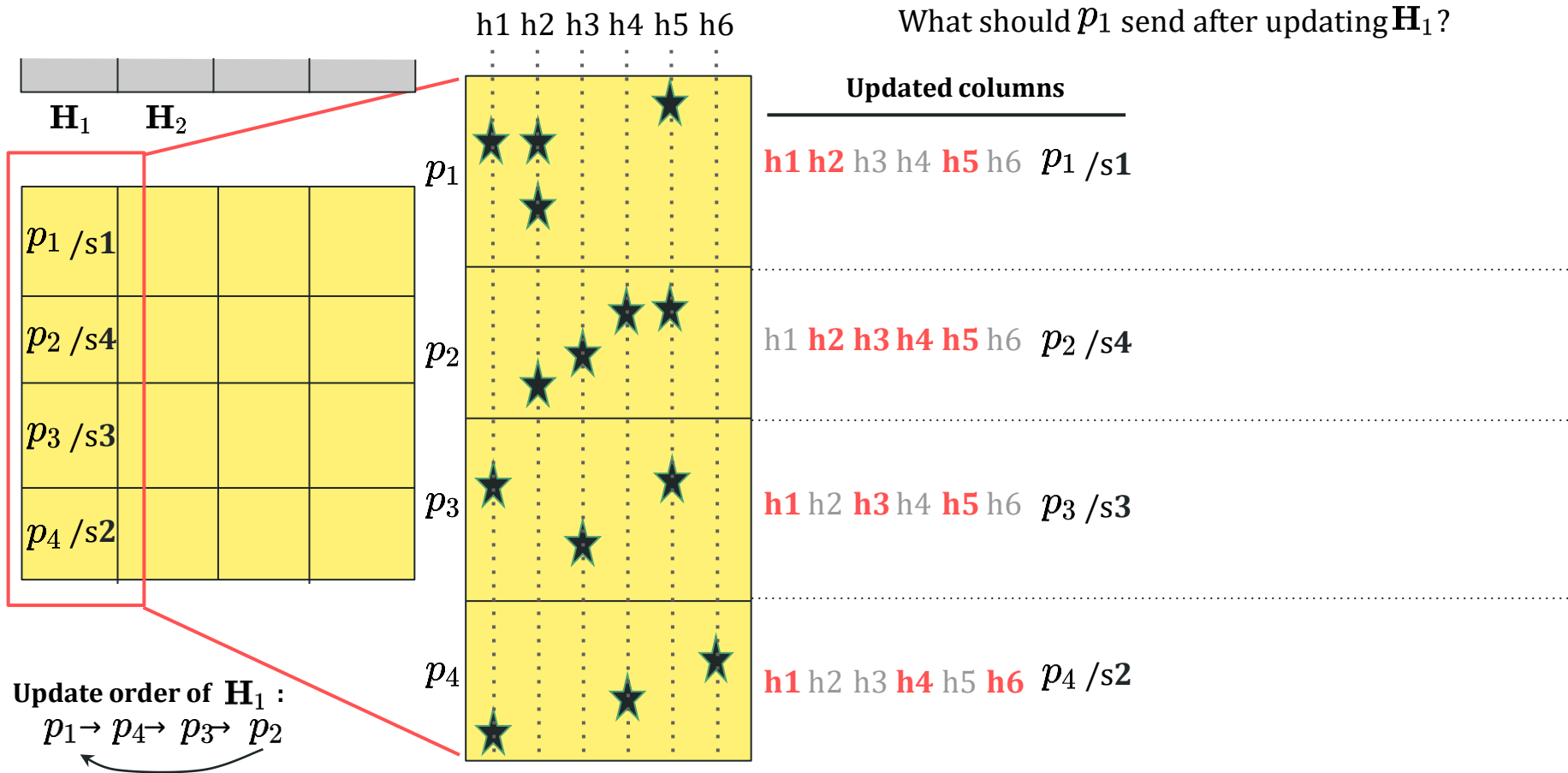


h_i is sent from p_x to p_y if :

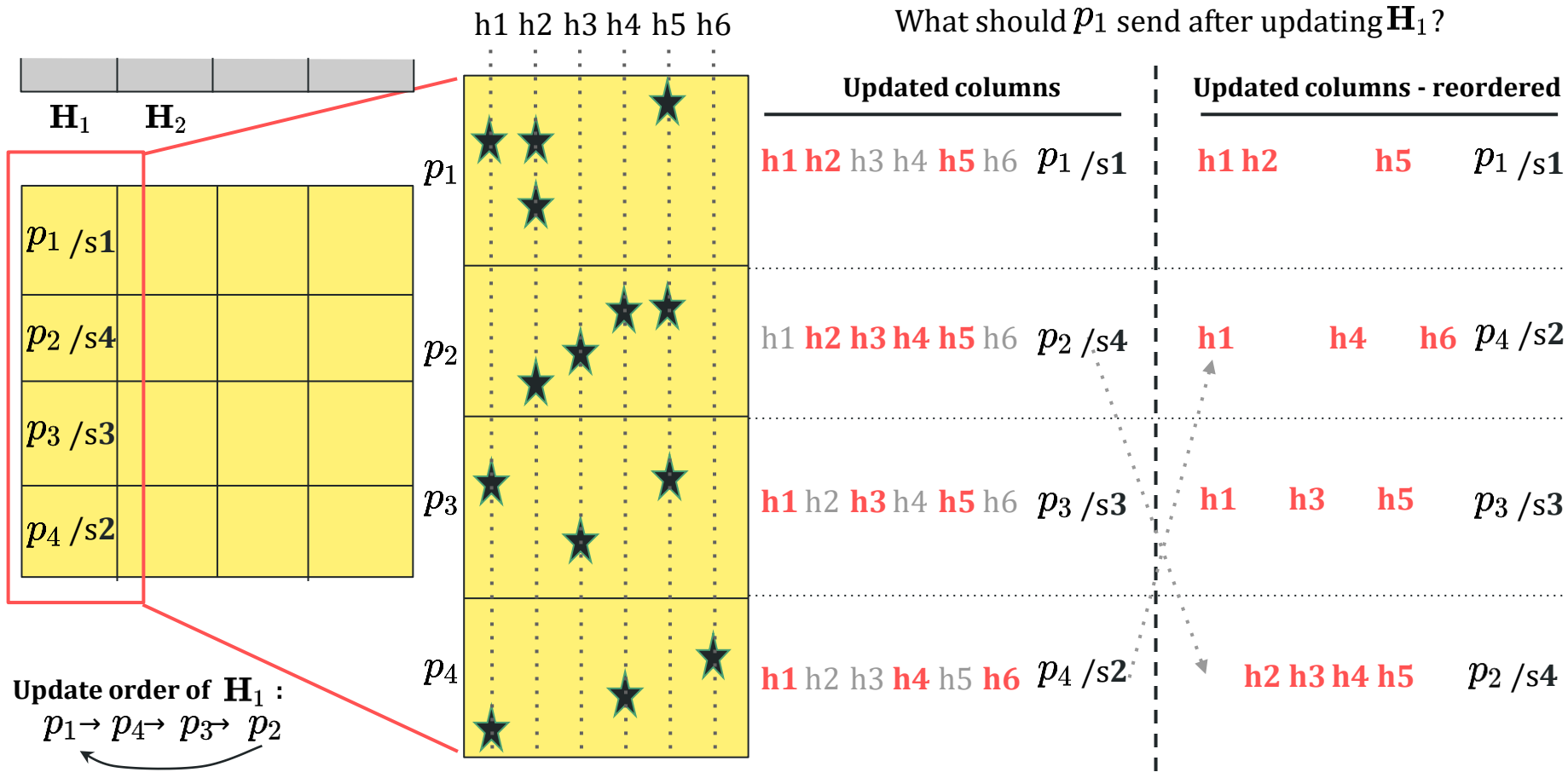
- both p_x and p_y update h_i and
- **no** processor **in between** does.

According to the update order

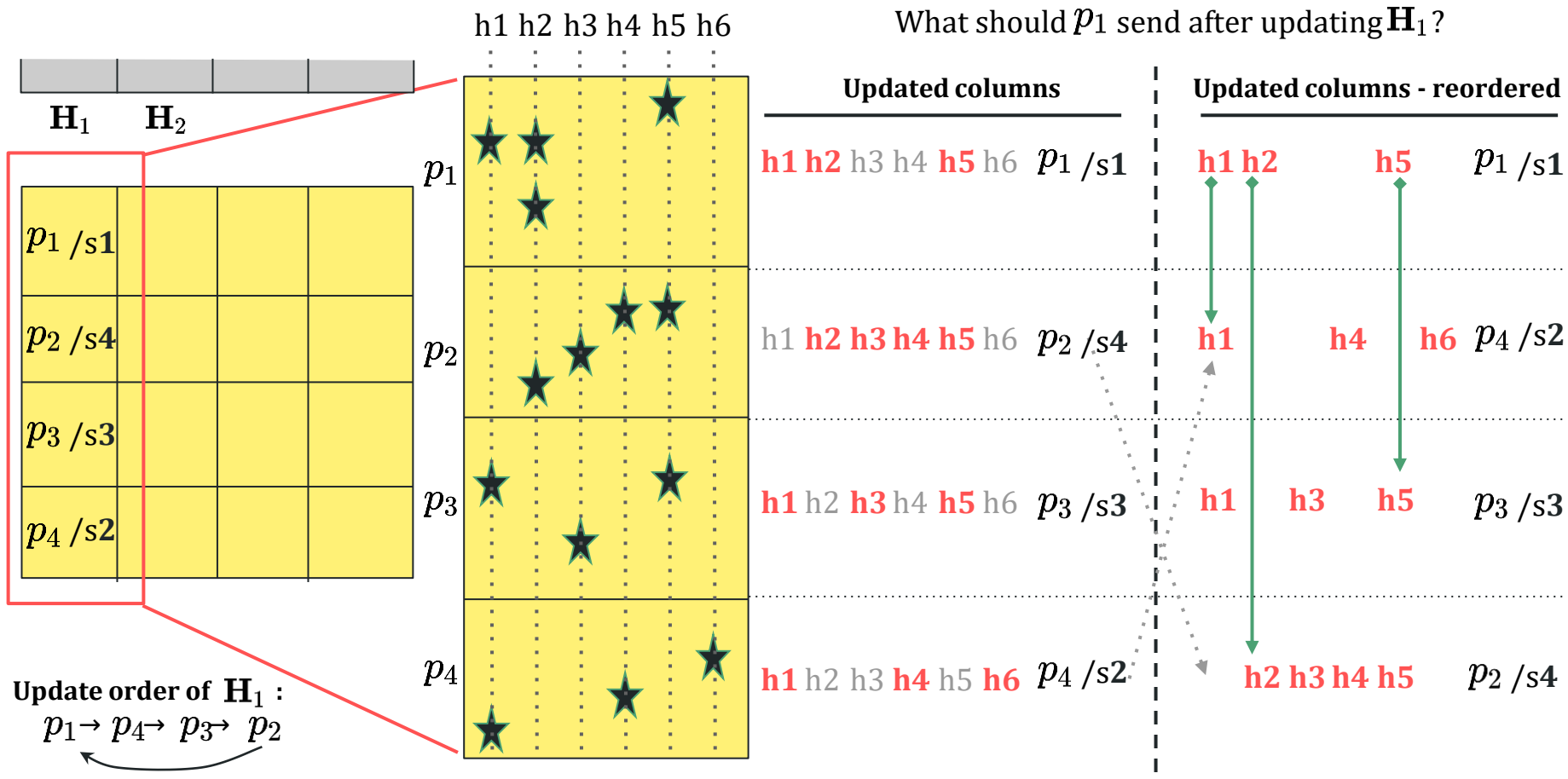
Finding essential communication



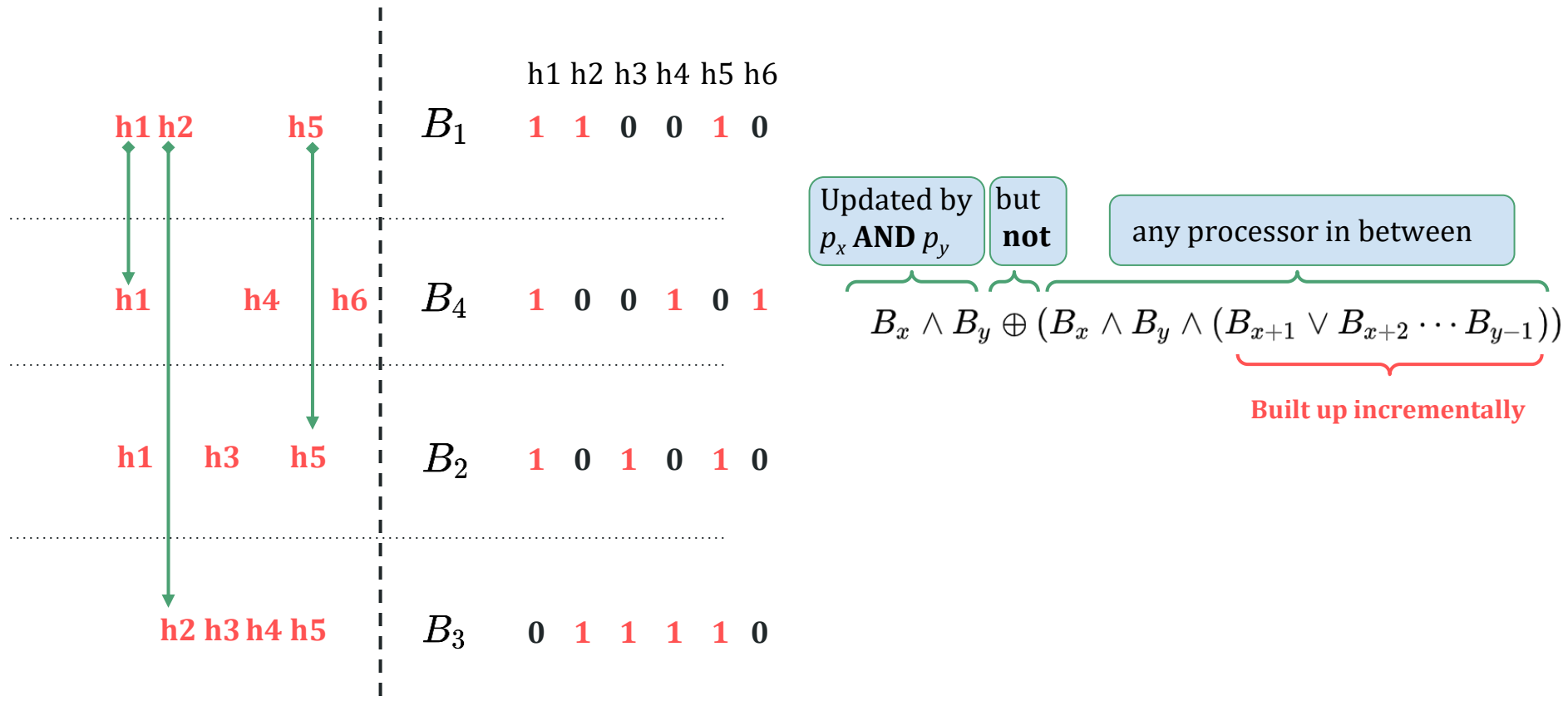
Finding essential communication



Finding essential communication



Efficiently finding essential communication



So far:

- ❖ Finding essential communication can be used to send P2P messages.
- ❖ **Invaluable** for reducing the volume of communication
- ❖ **Problem:** Number of messages significantly increase compared to block-wise communication:
 - **With block-wise:** each processor sends **1** msg per sub-epoch; **total of K** messages per processor.
 - **With P2P,** each processor sends up to **$K-1$** msgs per sub-epoch; **total of $O(K^2)$** messages per processor.

So far:

- ❖ Finding essential communication can be used to send P2P messages.
- ❖ **Invaluable** for reducing the volume of communication
- ❖ **Problem:** Number of messages significantly increase compared to block-wise communication:
 - **With block-wise:** each processor sends **1** msg per sub-epoch; **total of K** messages per processor.
 - **With P2P,** each processor sends up to **$K-1$** msgs per sub-epoch; **total of $O(K^2)$** messages per processor.

Research Question

Is it possible to exchange the same essential communication with message count asymptotically less than **$O(K^2)$**

So far:

- ❖ Finding essential communication can be used to send P2P messages.
- ❖ **Invaluable** for reducing the volume of communication
- ❖ **Problem:** Number of messages significantly increase compared to block-wise communication:
 - **With block-wise:** each processor sends **1** msg per sub-epoch; **total of K** messages per processor.
 - **With P2P,** each processor sends up to **$K-1$** msgs per sub-epoch; **total of $O(K^2)$** messages per processor.

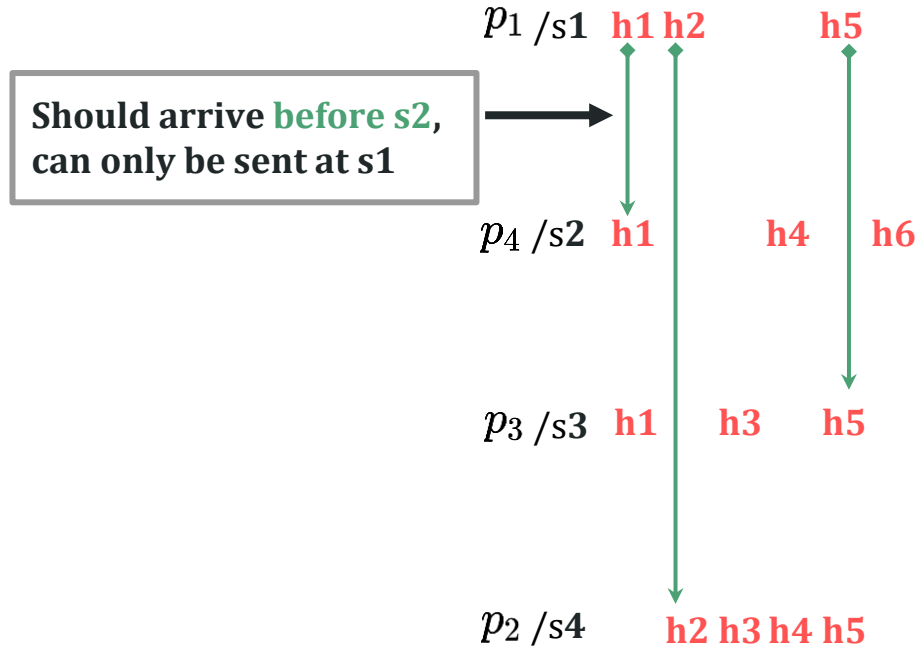
Research Question

Is it possible to exchange the same essential communication with message count asymptotically less than **$O(K^2)$**

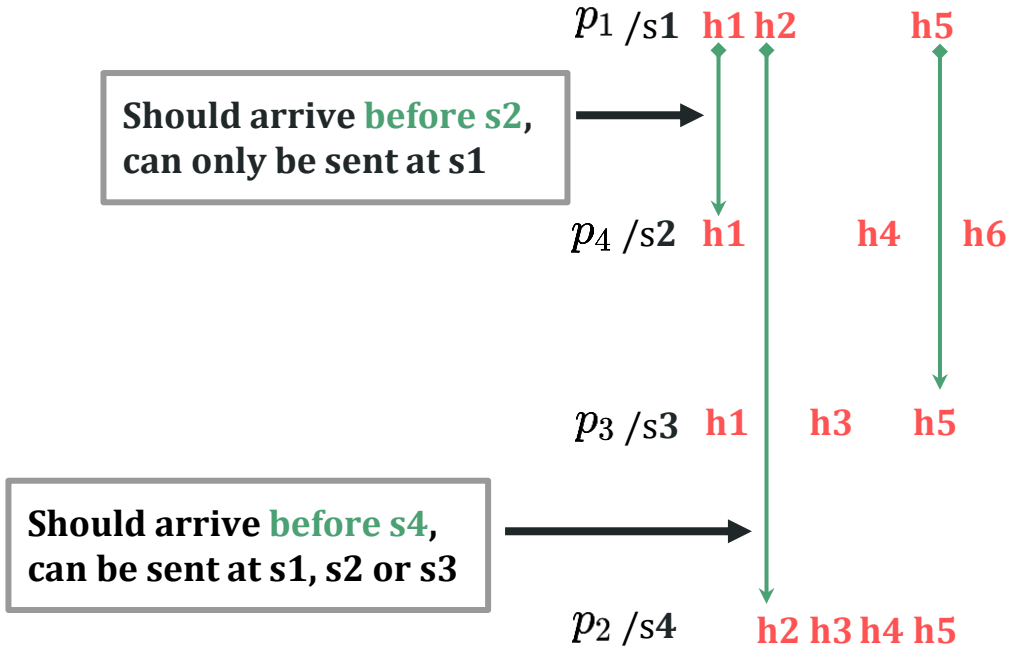
Key Observation

H-matrix rows sent via P2P messages are not always immediately needed in the next sub-epoch.

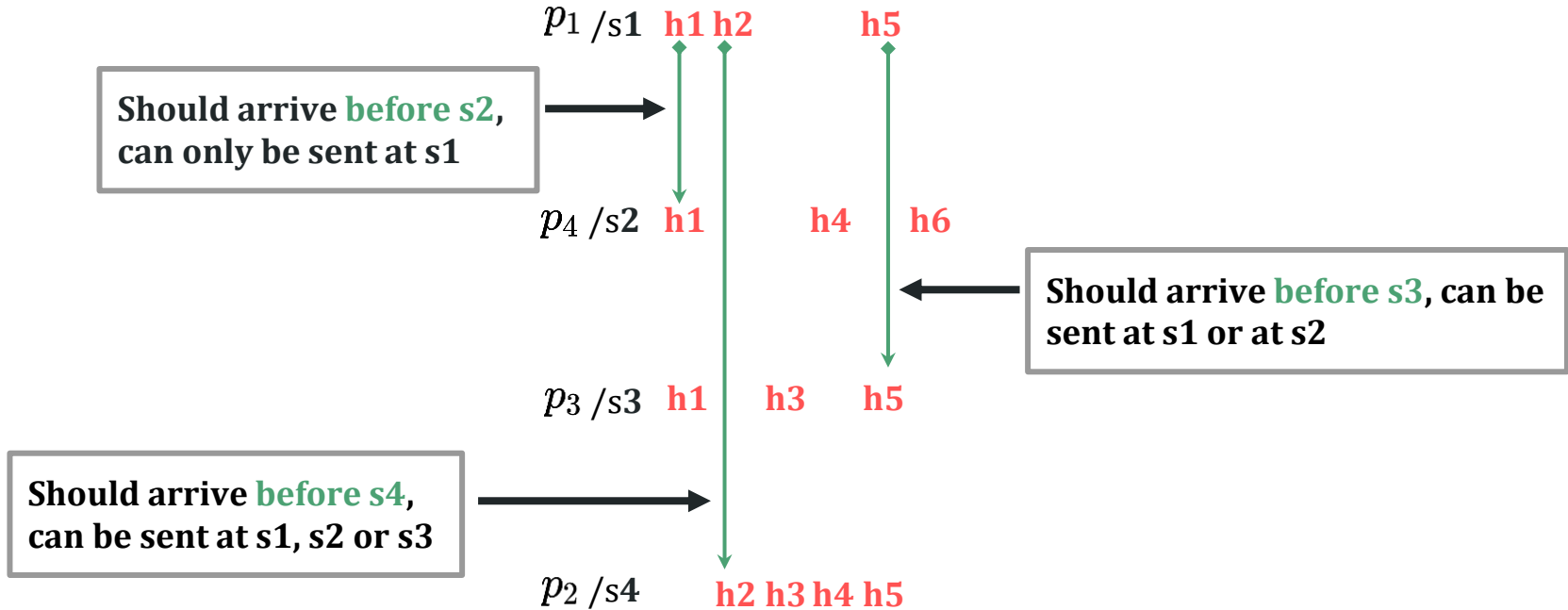
Key Observation: A closer look



Key Observation: A closer look



Key Observation: A closer look



The Hold & Combine Algorithm

Key idea: Hold H-matrix rows to be sent to the same processor at different sub-epochs and combine/send them in one message

Schedule of p_1 :

p_4				
p_3				
p_2				

The Hold & Combine Algorithm

Key idea: Hold H-matrix rows to be sent to the same processor at different sub-epochs and combine/send them in one message

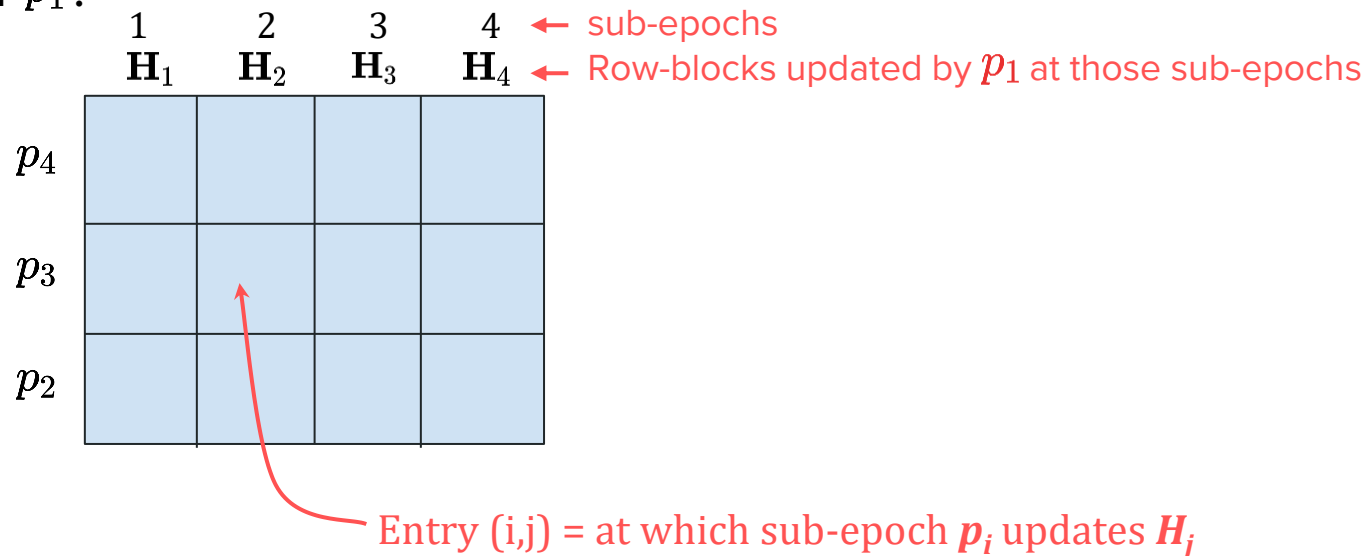
Schedule of p_1 :

	1	2	3	4	
	\mathbf{H}_1	\mathbf{H}_2	\mathbf{H}_3	\mathbf{H}_4	← sub-epochs
					← Row-blocks updated by p_1 at those sub-epochs
p_4					
p_3					
p_2					

The Hold & Combine Algorithm

Key idea: Hold H-matrix rows to be sent to the same processor at different sub-epochs and combine/send them in one message

Schedule of p_1 :



The Hold & Combine Algorithm

Key idea: Hold H-matrix rows to be sent to the same processor at different sub-epochs and combine/send them in one message

Schedule of p_1 :

	1 H_1	2 H_2	3 H_3	4 H_4
p_4	2	3	4	1 ⁺
p_3	3	4	1 ⁺	2 ⁺
p_2	4	1 ⁺	2 ⁺	3 ⁺

The Hold & Combine Algorithm

Key idea: Hold H-matrix rows to be sent to the same processor at different sub-epochs and combine/send them in one message

Schedule of p_1 :

	1 H_1	2 H_2	3 H_3	4 H_4
p_4	2 ★	3 ★	4 ★	1 ⁺ ★
p_3	3 ★	4 ★	1 ⁺ ★	2 ⁺ ★
p_2	4 ★	1 ⁺ ★	2 ⁺ ★	3 ₊ ★

The Hold & Combine Algorithm

Key idea: Hold H-matrix rows to be sent to the same processor at different sub-epochs and combine/send them in one message

Schedule of p_1 :

	1 H_1	2 H_2	3 H_3	4 H_4
p_4	2	3	4	1 ⁺
p_3	3	4	1 ⁺	2 ⁺
p_2	4	1 ⁺	2 ⁺	3 ₊

$$\sum_{i=1}^{\frac{K}{2}} 2 + \sum_{i=1}^{\frac{K-1}{2}} \frac{K}{i}$$

The Hold & Combine Algorithm

Key idea: Hold H-matrix rows to be sent to the same processor at different sub-epochs and combine/send them in one message

Schedule of p_1 :

	1 H_1	2 H_2	3 H_3	4 H_4
p_4	2	3	4	1 ⁺
p_3	3	4	1 ⁺	2 ⁺
p_2	4	1 ⁺	2 ⁺	3 ⁺

$$\sum_{i=1}^{\frac{K}{2}} 2 + \sum_{i=1}^{\frac{K-1}{2}} \frac{K}{i}$$

No H&C: up to 12 messages ($N * N - 1$)
With H&C: up to 8 messages ($N * \lg N$)

The Hold & Combine Algorithm

Key idea: Hold H-matrix rows to be sent to the same processor at different sub-epochs and combine/send them in one message

Schedule of p_1 :

	1 H_1	2 H_2	3 H_3	4 H_4
p_4	2	3	4	1 ⁺
p_3	3	4	1 ⁺	2 ⁺
p_2	4	1 ⁺	2 ⁺	3 ₊

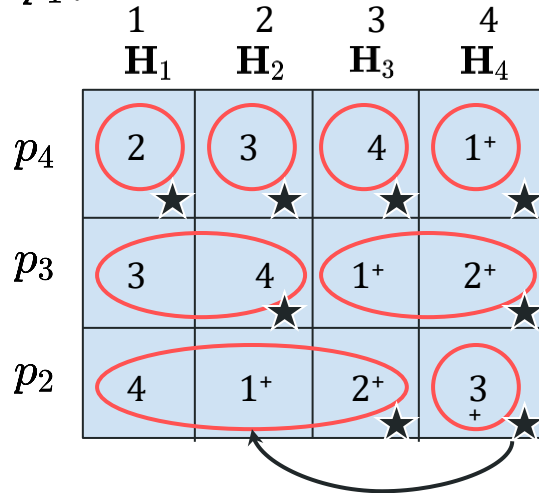
$$\sum_{i=1}^{\frac{K}{2}} 2 + \sum_{i=1}^{\frac{K-1}{2}} \frac{K}{i}$$

No H&C: up to 12 messages ($N * N - 1$)
With H&C: up to 8 messages ($N * \lg N$)

The Hold & Combine Algorithm

Key idea: Hold H-matrix rows to be sent to the same processor at different sub-epochs and combine/send them in one message

Schedule of p_1 :



$$\sum_{i=1}^{\frac{K}{2}} 2 + \sum_{i=1}^{\frac{K-1}{2}} \frac{K}{i}$$

No H&C: up to 12 messages ($N * N - 1$)
With H&C: up to 8 messages ($N * \lg N$)

Greedy choice: Send each message in the iteration just before it is required → leads to balanced distribution of messages among sub-epochs

Experiments and Key Results

6 real-world sparse **rating** matrices. $5M < \text{nnz} < 475M$ nonzeros

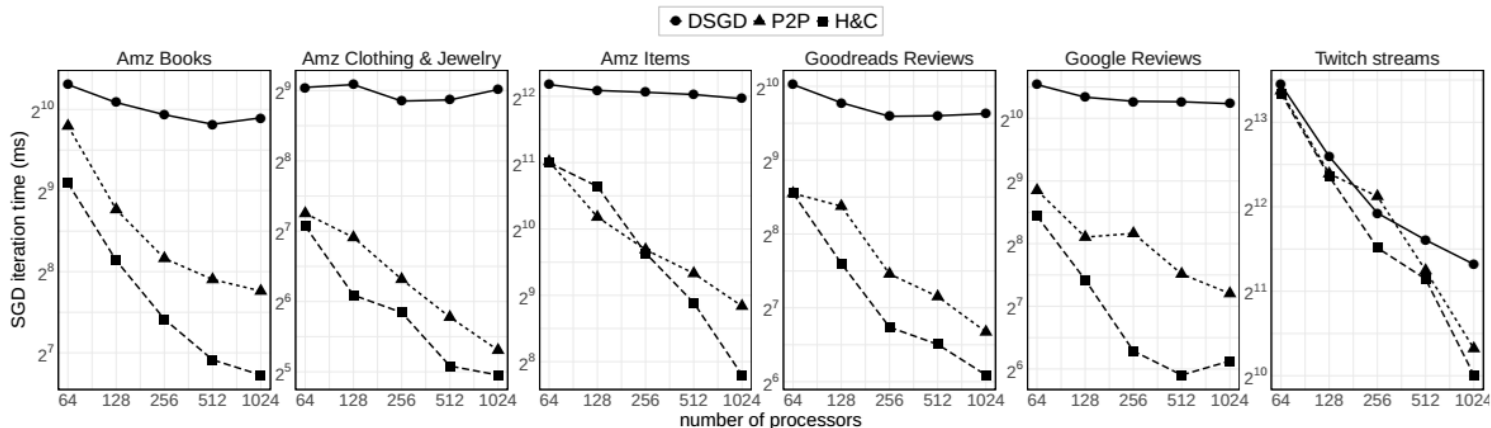
In terms of bandwidth:

- ❖ Using P2P **reduces** total volume by **10x - 120x** (P2P has the same volume with or w/o H&C)

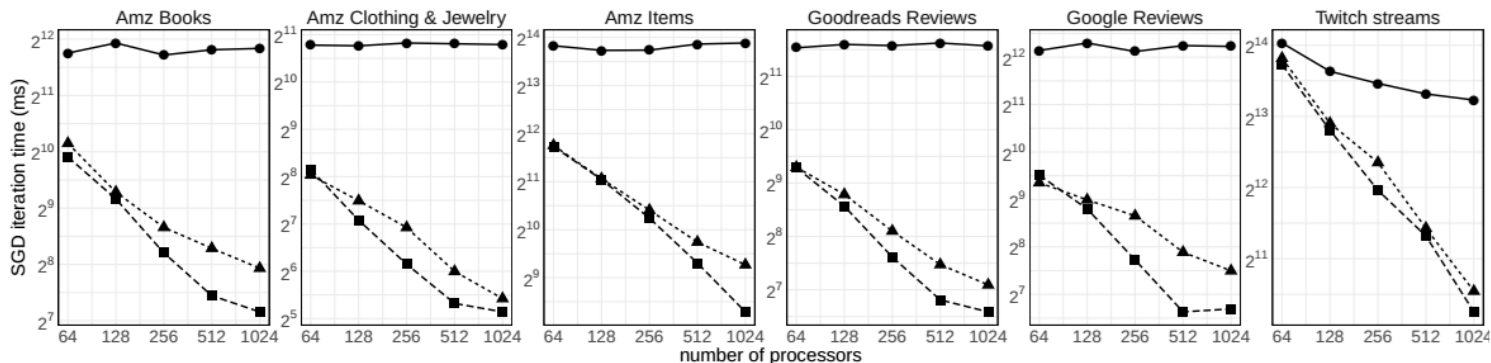
In terms of latency:

- ❖ Using random-based P2P **increases** total messages by **3.5x - 57x**
- ❖ Using **H&C increases** total messages by **3x - 8x**

Experiments and Key Results - Cont'd

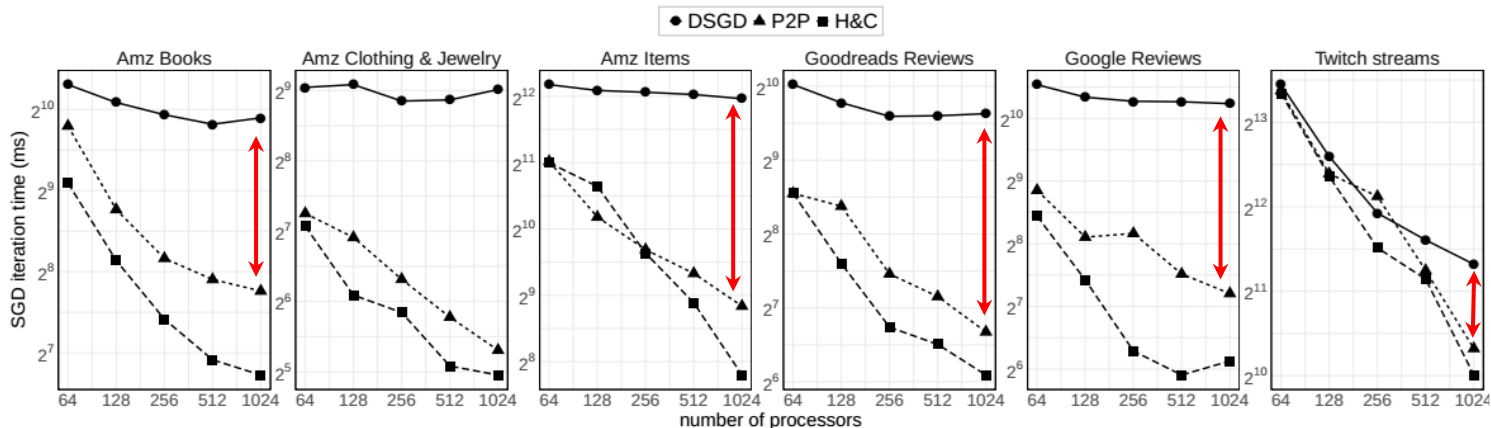


(a) $F = 16$



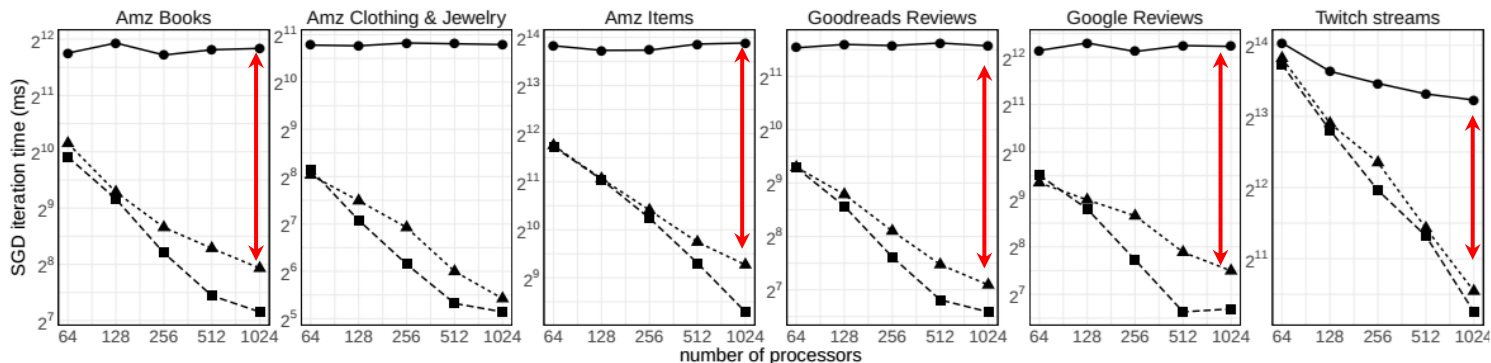
(b) $F = 64$

Experiments and Key Results - Cont'd



(a) $F = 16$

Gap
increase



(b) $F = 64$

Thank you!

Questions ?

More of SPCL's research:

 youtube.com/@spcl **180+ Talks**

 twitter.com/spcl_eth **1.4K+ Followers**

 github.com/spcl **3.8K+ Stars**

... or spcl.ethz.ch

