

Block-Diagonal Coding for Distributed Computing With Straggling Servers

Albin Severinson^{†‡}, Alexandre Graell i Amat[†], and Eirik Rosnes[‡]

[†] Department of Electrical Engineering, Chalmers University of Technology, Gothenburg, Sweden

[‡] University of Bergen/Simula Research Lab, Bergen, Norway

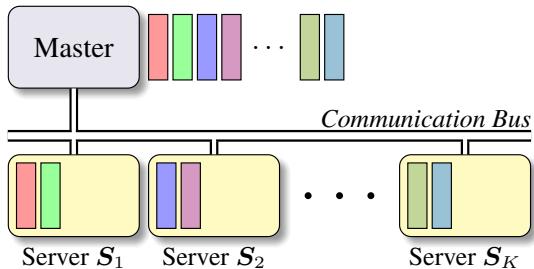
IEEE ITW
Kaohsiung, November, 2017

simula@uib

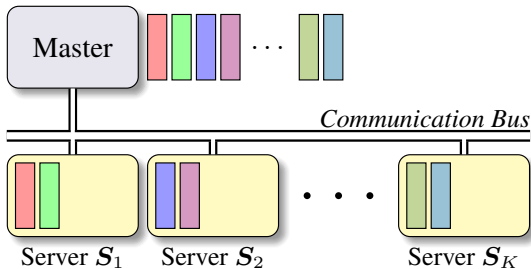


CHALMERS

Motivation



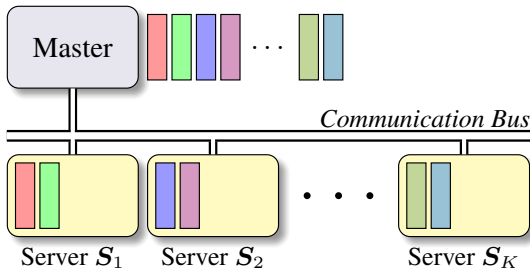
Motivation



Problem addressed

- Given an $m \times n$ matrix A and N vectors x_1, \dots, x_N , we want to compute $y_1 = Ax_1, \dots, y_N = Ax_N$ using K servers.

Motivation



Problem addressed

- Given an $m \times n$ matrix A and N vectors x_1, \dots, x_N , we want to compute $y_1 = Ax_1, \dots, y_N = Ax_N$ using K servers.

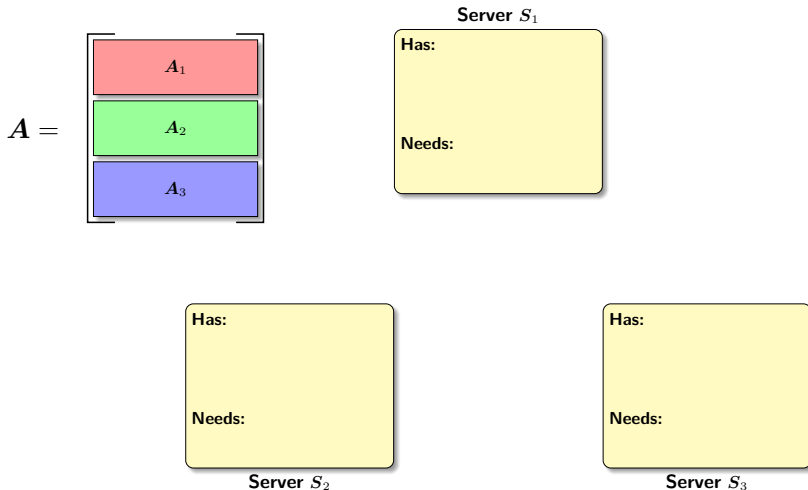
Performance metrics

- Communication load:** Average number of messages sent over the network
- Computational delay:** Average overall runtime of the computation

Bandwidth Scarcity

(Coded MapReduce, Li *et al.*, 2015)

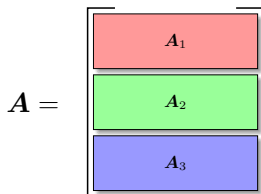
$$y_1 = Ax_1, y_2 = Ax_2, y_3 = Ax_3$$



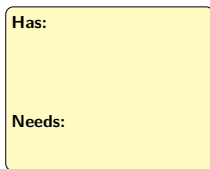
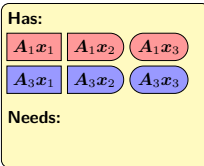
Bandwidth Scarcity

(Coded MapReduce, Li *et al.*, 2015)

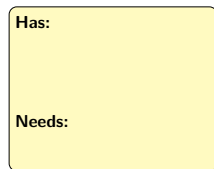
$$y_1 = Ax_1, y_2 = Ax_2, y_3 = Ax_3$$



Server S_1



Server S_2

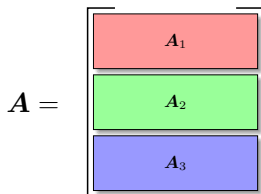


Server S_3

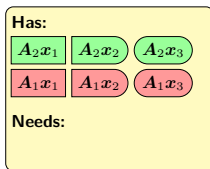
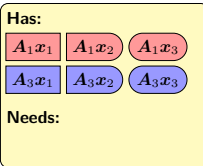
Bandwidth Scarcity

(Coded MapReduce, Li *et al.*, 2015)

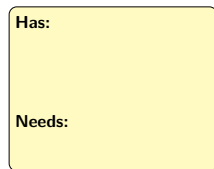
$$y_1 = Ax_1, y_2 = Ax_2, y_3 = Ax_3$$



Server S_1



Server S_2

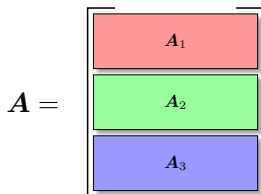


Server S_3

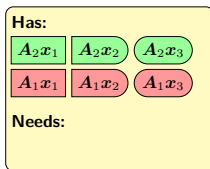
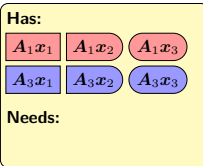
Bandwidth Scarcity

(Coded MapReduce, Li *et al.*, 2015)

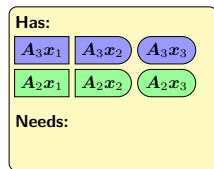
$$y_1 = Ax_1, y_2 = Ax_2, y_3 = Ax_3$$



Server S_1



Server S_2

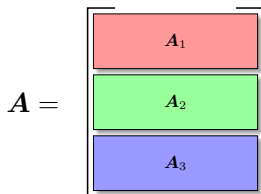


Server S_3

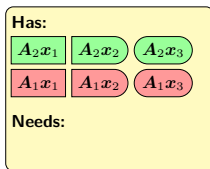
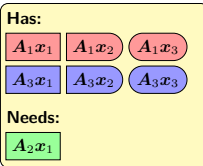
Bandwidth Scarcity

(Coded MapReduce, Li *et al.*, 2015)

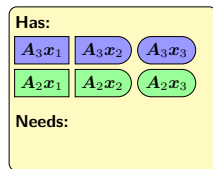
$$y_1 = Ax_1, y_2 = Ax_2, y_3 = Ax_3$$



Server S_1



Server S_2

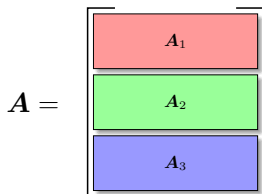


Server S_3

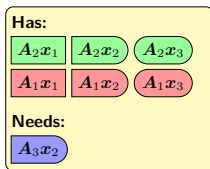
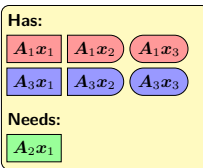
Bandwidth Scarcity

(Coded MapReduce, Li *et al.*, 2015)

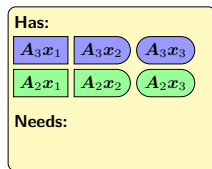
$$y_1 = Ax_1, y_2 = Ax_2, y_3 = Ax_3$$



Server S_1



Server S_2

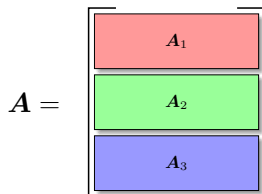


Server S_3

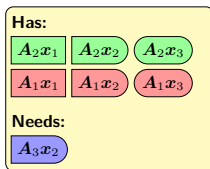
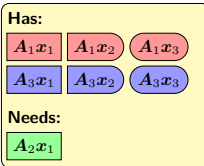
Bandwidth Scarcity

(Coded MapReduce, Li *et al.*, 2015)

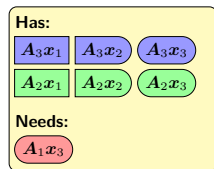
$$y_1 = Ax_1, y_2 = Ax_2, y_3 = Ax_3$$



Server S_1



Server S_2

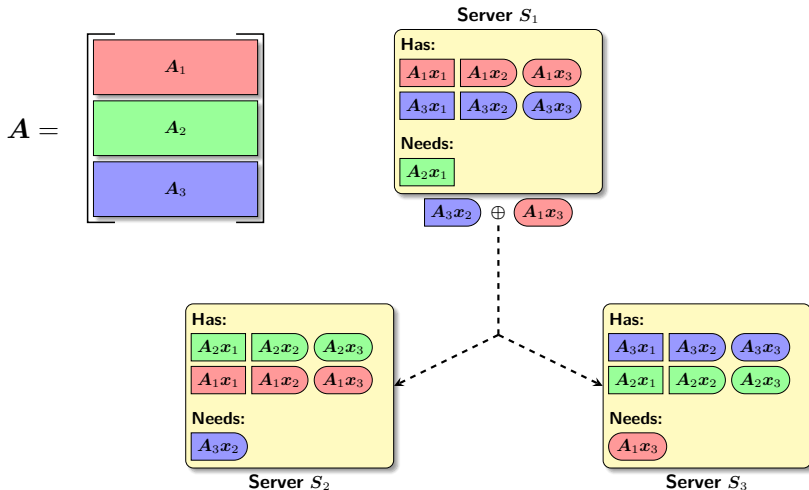


Server S_3

Bandwidth Scarcity

(Coded MapReduce, Li *et al.*, 2015)

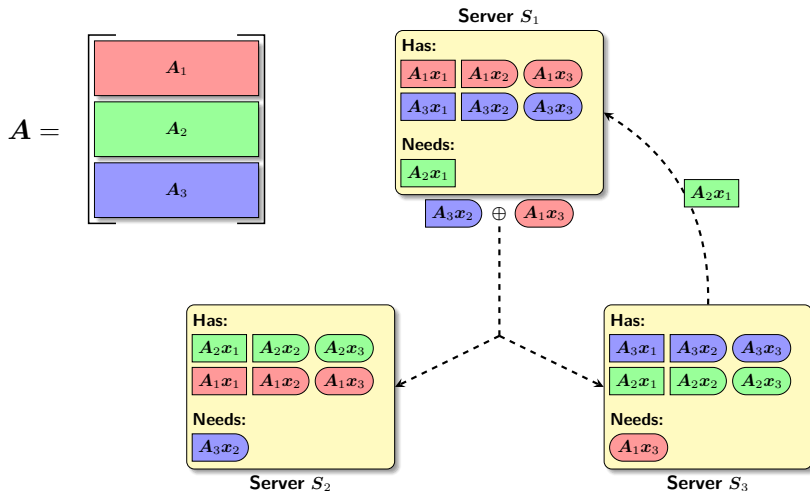
$$y_1 = Ax_1, y_2 = Ax_2, y_3 = Ax_3$$



Bandwidth Scarcity

(Coded MapReduce, Li *et al.*, 2015)

$$y_1 = Ax_1, y_2 = Ax_2, y_3 = Ax_3$$

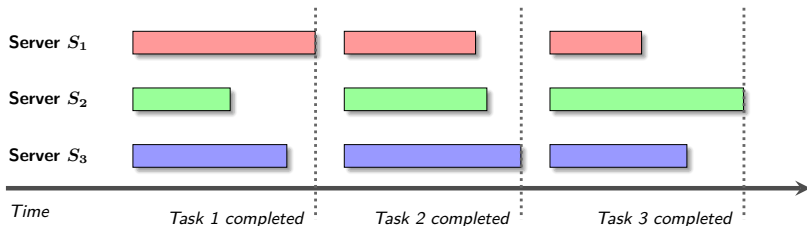


The straggler problem

(Speeding up Distributed Machine Learning Using Codes, Lee *et al.*, 2016)

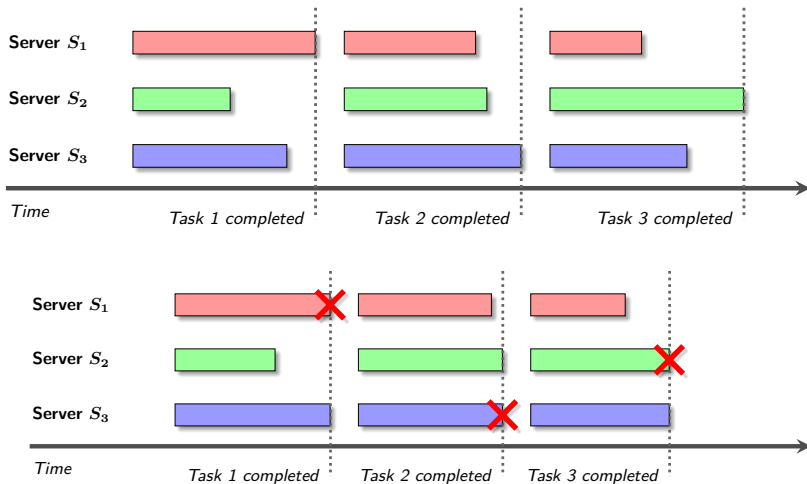
The straggler problem

(Speeding up Distributed Machine Learning Using Codes, Lee *et al.*, 2016)



The straggler problem

(Speeding up Distributed Machine Learning Using Codes, Lee *et al.*, 2016)



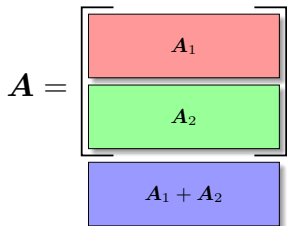
The Straggler Problem

$$y = Ax$$

$$A = \begin{bmatrix} A_1 \\ A_2 \end{bmatrix}$$

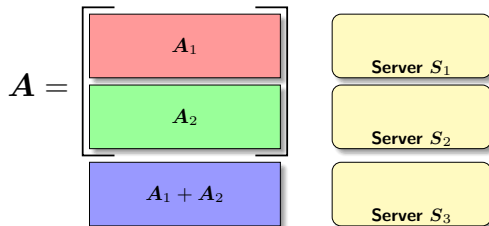
The Straggler Problem

$$y = Ax$$



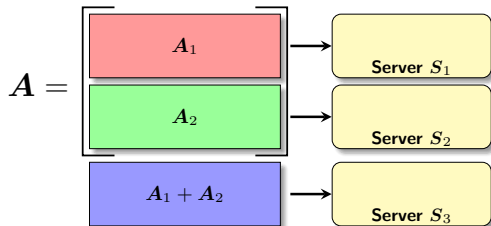
The Straggler Problem

$$y = Ax$$



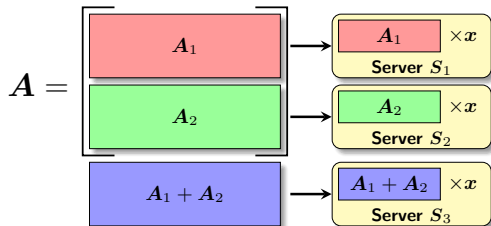
The Straggler Problem

$$y = Ax$$



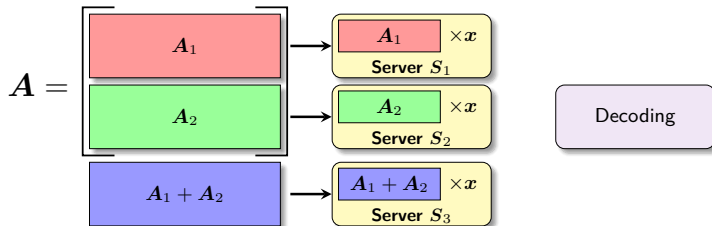
The Straggler Problem

$$y = Ax$$



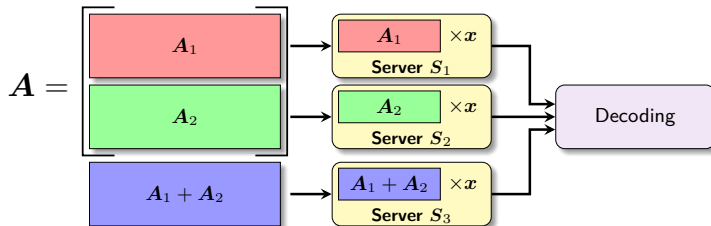
The Straggler Problem

$$y = Ax$$



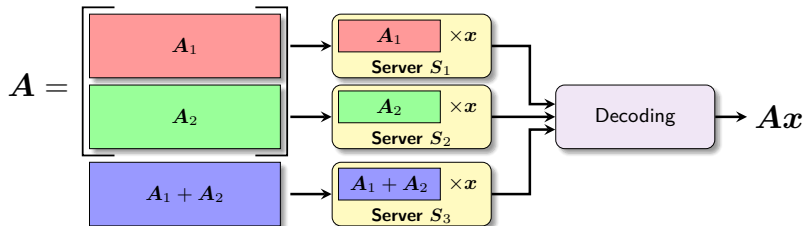
The Straggler Problem

$$y = Ax$$



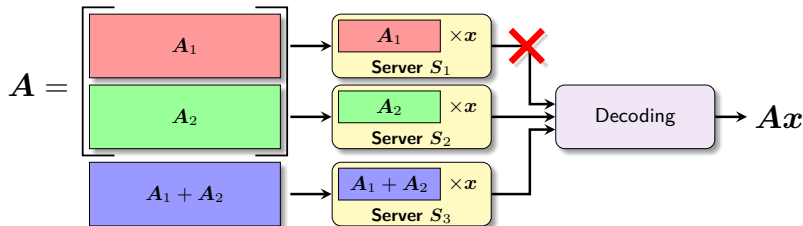
The Straggler Problem

$$y = Ax$$



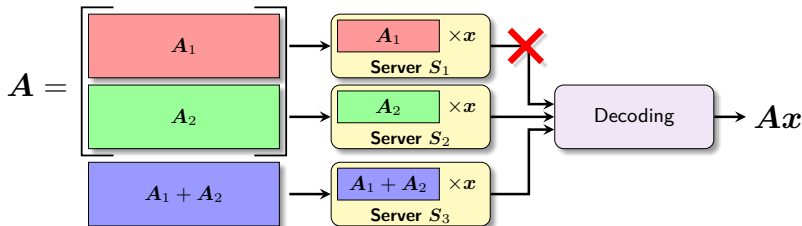
The Straggler Problem

$$y = Ax$$



The Straggler Problem

$$y = Ax$$

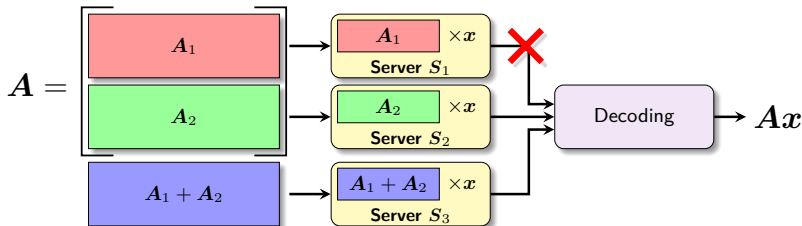


In general

- Introduce redundancy by encoding the input matrix A .

The Straggler Problem

$$y = Ax$$



In general

- Introduce redundancy by encoding the input matrix A .
- Each server is given **more work**. However, this may still **lower the computational delay!**

The Unified Coded Framework

(A Unified Coding Framework for Distributed Computing with Straggling Servers, Li *et al.*, 2016)

Coded approach by *Li et al.*

The Unified Coded Framework

(A Unified Coding Framework for Distributed Computing with Straggling Servers, Li *et al.*, 2016)

Coded approach by Li *et al.*

- Encode the columns of $\mathbf{A} \in \mathbb{F}^{m \times n}$ using an (r, m) MDS code by multiplying \mathbf{A} with an $r \times n$ encoding matrix Ψ_{MDS} , i.e., $\mathbf{C} = \Psi_{\text{MDS}} \mathbf{A}$.

The Unified Coded Framework

(A Unified Coding Framework for Distributed Computing with Straggling Servers, Li *et al.*, 2016)

Coded approach by Li *et al.*

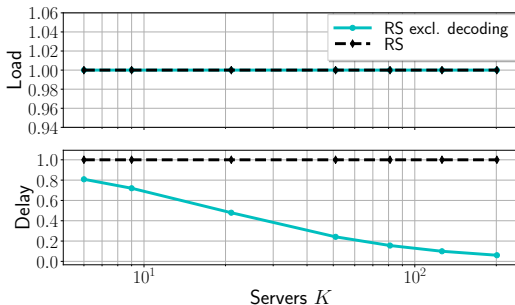
- Encode the columns of $\mathbf{A} \in \mathbb{F}^{m \times n}$ using an (r, m) MDS code by multiplying \mathbf{A} with an $r \times n$ encoding matrix Ψ_{MDS} , i.e., $\mathbf{C} = \Psi_{\text{MDS}} \mathbf{A}$.
- But long MDS codes may be complex to decode \rightarrow decoding delay!

The Unified Coded Framework

(A Unified Coding Framework for Distributed Computing with Straggling Servers, Li *et al.*, 2016)

Coded approach by **Li *et al.***

- Encode the columns of $A \in \mathbb{F}^{m \times n}$ using an (r, m) MDS code by multiplying A with an $r \times n$ encoding matrix Ψ_{MDS} , i.e., $C = \Psi_{\text{MDS}} A$.
- But long MDS codes may be complex to decode \rightarrow **decoding delay!**



- 2000 rows assigned to each server, $n = 10000$ columns, and code rate $m/r = 2/3$.

In this talk

Block-Diagonal Coding

- We propose a **block-diagonal** encoding scheme with lower decoding complexity.

In this talk

Block-Diagonal Coding

- We propose a **block-diagonal** encoding scheme with lower decoding complexity.

Idea

- **Partition** \mathbf{A} into T disjoint submatrices and apply **smaller MDS codes** to each submatrix,

$$\mathbf{C} = \Psi_{\text{BDC}} \mathbf{A}, \quad \Psi_{\text{BDC}} = \begin{bmatrix} \psi_1 & & \\ & \ddots & \\ & & \psi_T \end{bmatrix}.$$

In this talk

Block-Diagonal Coding

- We propose a **block-diagonal** encoding scheme with lower decoding complexity.

Idea

- **Partition** A into T disjoint submatrices and apply **smaller MDS codes** to each submatrix,

$$C = \Psi_{\text{BDC}} A, \quad \Psi_{\text{BDC}} = \begin{bmatrix} \psi_1 & & \\ & \ddots & \\ & & \psi_T \end{bmatrix}.$$

Outcome

- No impact on the **communication load** or the **computational delay of the matrix multiplication** up to a given T .

In this talk

Block-Diagonal Coding

- We propose a **block-diagonal** encoding scheme with lower decoding complexity.

Idea

- Partition** A into T disjoint submatrices and apply **smaller MDS codes** to each submatrix,

$$C = \Psi_{\text{BDC}} A, \quad \Psi_{\text{BDC}} = \begin{bmatrix} \psi_1 & & \\ & \ddots & \\ & & \psi_T \end{bmatrix}.$$

Outcome

- No impact on the **communication load** or the **computational delay of the matrix multiplication** up to a given T .
- Overall **computational delay** is lower than that of the scheme by Li *et al.*

In this talk

Block-Diagonal Coding

- We propose a **block-diagonal** encoding scheme with lower decoding complexity.

Idea

- **Partition** A into T disjoint submatrices and apply **smaller MDS codes** to each submatrix,

$$C = \Psi_{\text{BDC}} A, \quad \Psi_{\text{BDC}} = \begin{bmatrix} \psi_1 & & \\ & \ddots & \\ & & \psi_T \end{bmatrix}.$$

Outcome

- No impact on the **communication load** or the **computational delay of the matrix multiplication** up to a given T .
- Overall **computational delay** is lower than that of the scheme by Li *et al.*
- Larger T may reduce **computational delay** further at the expense of higher **communication load**.

Block-Diagonal Coding

$$\Psi_{\text{BDC}} = \begin{bmatrix} \psi_1 & & \\ & \psi_2 & \\ & & \psi_3 \end{bmatrix}$$

$r \times m$

- Block-diagonal encoding with $T = 3$ partitions.

Block-Diagonal Coding

$$\Psi_{\text{BDC}} \mathbf{A} = \begin{bmatrix} \boxed{\psi_1} & & \\ & \boxed{\psi_2} & \\ & & \boxed{\psi_3} \end{bmatrix} \begin{bmatrix} \boxed{\mathbf{A}_1} \\ \boxed{\mathbf{A}_2} \\ \boxed{\mathbf{A}_3} \end{bmatrix}$$

$r \times m$ $m \times n$

- Block-diagonal encoding with $T = 3$ partitions.

Block-Diagonal Coding

$$\Psi_{\text{BDC}} \mathbf{A} = \begin{bmatrix} \boxed{\psi_1} & & \\ & \boxed{\psi_2} & \\ & & \boxed{\psi_3} \end{bmatrix} \begin{bmatrix} \boxed{\mathbf{A}_1} \\ \boxed{\mathbf{A}_2} \\ \boxed{\mathbf{A}_3} \end{bmatrix} = \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix}$$

$r \times m$
 $m \times n$
 $r \times n$

- Block-diagonal encoding with $T = 3$ partitions.

Block-Diagonal Coding

$$\Psi_{\text{BDC}} \mathbf{A} = \begin{bmatrix} \boxed{\psi_1} & & \\ & \boxed{\psi_2} & \\ & & \boxed{\psi_3} \end{bmatrix} \begin{bmatrix} \boxed{\mathbf{A}_1} \\ \boxed{\mathbf{A}_2} \\ \boxed{\mathbf{A}_3} \end{bmatrix} = \begin{bmatrix} \boxed{\psi_1 \mathbf{A}_1} \\ & & \\ & & \end{bmatrix}$$

$r \times m$
 $m \times n$
 $r \times n$

- Block-diagonal encoding with $T = 3$ partitions.

Block-Diagonal Coding

$$\Psi_{\text{BDC}} \mathbf{A} = \begin{bmatrix} \boxed{\psi_1} & & \\ & \boxed{\psi_2} & \\ & & \boxed{\psi_3} \end{bmatrix} \begin{bmatrix} \boxed{A_1} \\ \boxed{A_2} \\ \boxed{A_3} \end{bmatrix} = \begin{bmatrix} \boxed{\psi_1 A_1} \\ \boxed{\psi_2 A_2} \end{bmatrix}$$

$r \times m$
 $m \times n$
 $r \times n$

- Block-diagonal encoding with $T = 3$ partitions.

Block-Diagonal Coding

$$\Psi_{\text{BDC}} \mathbf{A} = \begin{bmatrix} \boxed{\psi_1} & & \\ & \boxed{\psi_2} & \\ & & \boxed{\psi_3} \end{bmatrix} \begin{bmatrix} \boxed{\mathbf{A}_1} \\ \boxed{\mathbf{A}_2} \\ \boxed{\mathbf{A}_3} \end{bmatrix} = \begin{bmatrix} \boxed{\psi_1 \mathbf{A}_1} \\ \boxed{\psi_2 \mathbf{A}_2} \\ \boxed{\psi_3 \mathbf{A}_3} \end{bmatrix}$$

$r \times m$ $m \times n$ $r \times n$

- Block-diagonal encoding with $T = 3$ partitions.

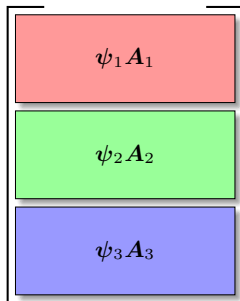
Block-Diagonal Coding

$$\Psi_{\text{BDC}} \mathbf{A} = \begin{bmatrix} \boxed{\psi_1} & & \\ & \boxed{\psi_2} & \\ & & \boxed{\psi_3} \end{bmatrix} \begin{bmatrix} \boxed{A_1} \\ \boxed{A_2} \\ \boxed{A_3} \end{bmatrix} = \begin{bmatrix} \boxed{\psi_1 A_1} \\ \boxed{\psi_2 A_2} \\ \boxed{\psi_3 A_3} \end{bmatrix}$$

$r \times m$
 $m \times n$
 $r \times n$

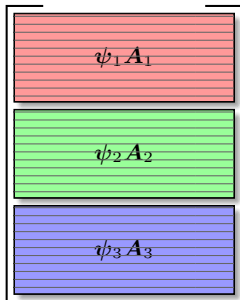
- Block-diagonal encoding with $T = 3$ partitions.
- Need any m/T out of r/T rows from **each partition** to decode.

Block-Diagonal Coding



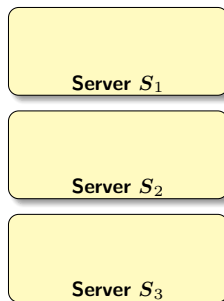
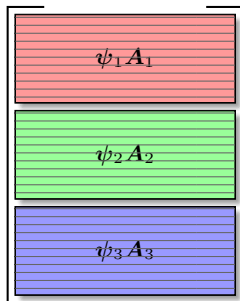
- Block-diagonal encoding with $T = 3$ partitions.
- Need any m/T out of r/T rows from **each partition** to decode.

Block-Diagonal Coding



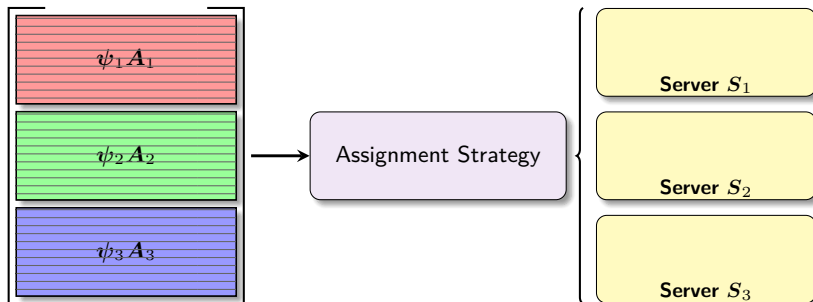
- Block-diagonal encoding with $T = 3$ partitions.
- Need any m/T out of r/T rows from **each partition** to decode.

Block-Diagonal Coding



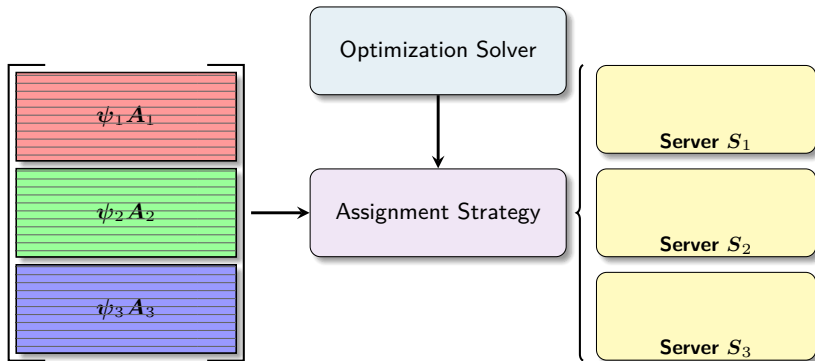
- Block-diagonal encoding with $T = 3$ partitions.
- Need any m/T out of r/T rows from **each partition** to decode.

Block-Diagonal Coding



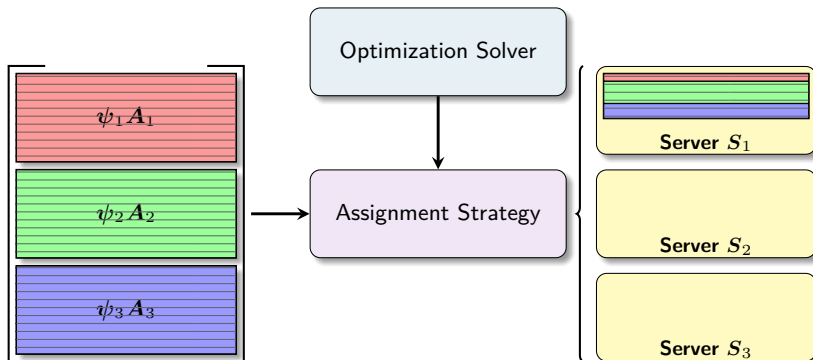
- Block-diagonal encoding with $T = 3$ partitions.
- Need any m/T out of r/T rows from **each partition** to decode.
- Need to **assign** coded rows to servers very carefully in some instances (such as when the number of servers is small).

Block-Diagonal Coding



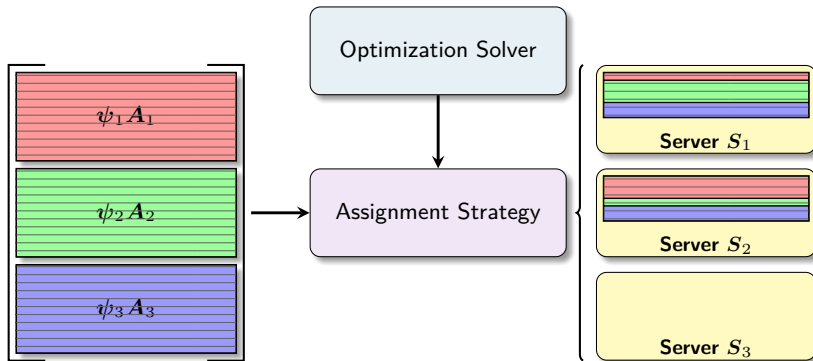
- Block-diagonal encoding with $T = 3$ partitions.
- Need any m/T out of r/T rows from **each partition** to decode.
- Need to **assign** coded rows to servers very carefully in some instances (such as when the number of servers is small).
- This assignment can be formulated as an **optimization** problem.

Block-Diagonal Coding



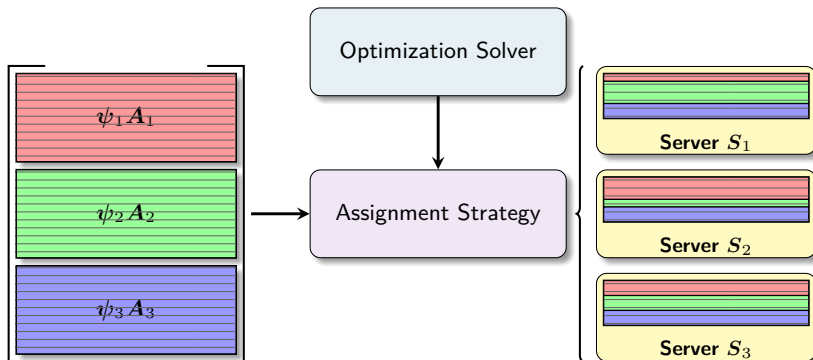
- Block-diagonal encoding with $T = 3$ partitions.
- Need any m/T out of r/T rows from **each partition** to decode.
- Need to **assign** coded rows to servers very carefully in some instances (such as when the number of servers is small).
- This assignment can be formulated as an **optimization** problem.

Block-Diagonal Coding



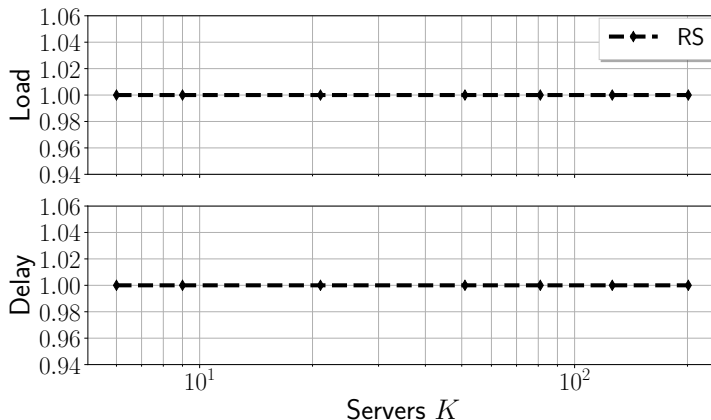
- Block-diagonal encoding with $T = 3$ partitions.
- Need any m/T out of r/T rows from **each partition** to decode.
- Need to **assign** coded rows to servers very carefully in some instances (such as when the number of servers is small).
- This assignment can be formulated as an **optimization** problem.

Block-Diagonal Coding



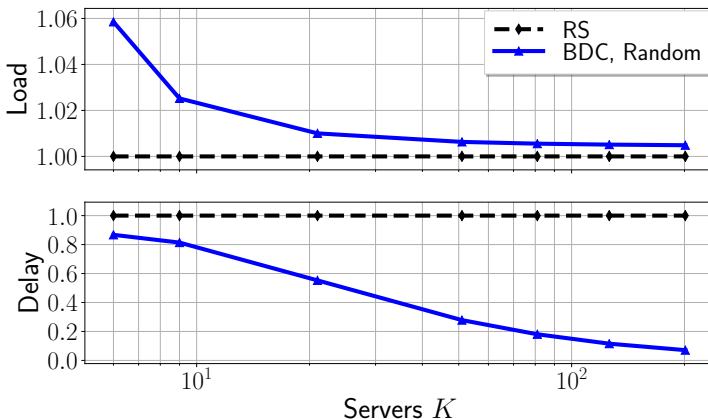
- Block-diagonal encoding with $T = 3$ partitions.
- Need any m/T out of r/T rows from **each partition** to decode.
- Need to **assign** coded rows to servers very carefully in some instances (such as when the number of servers is small).
- This assignment can be formulated as an **optimization** problem.

Numerical results



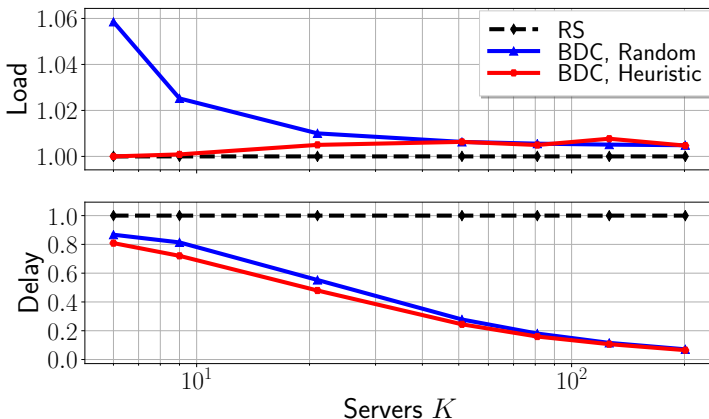
- 2000 rows assigned to each server, $n = 10000$ columns, $m/T = 10$ rows per partition, and code rate $m/r = 2/3$.

Numerical results



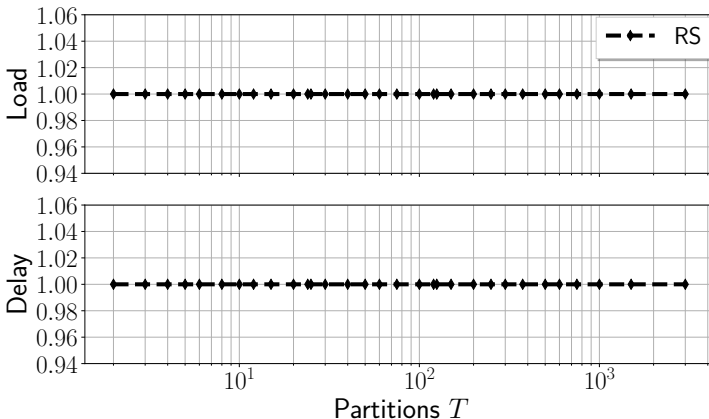
- 2000 rows assigned to each server, $n = 10000$ columns, $m/T = 10$ rows per partition, and code rate $m/r = 2/3$.

Numerical results



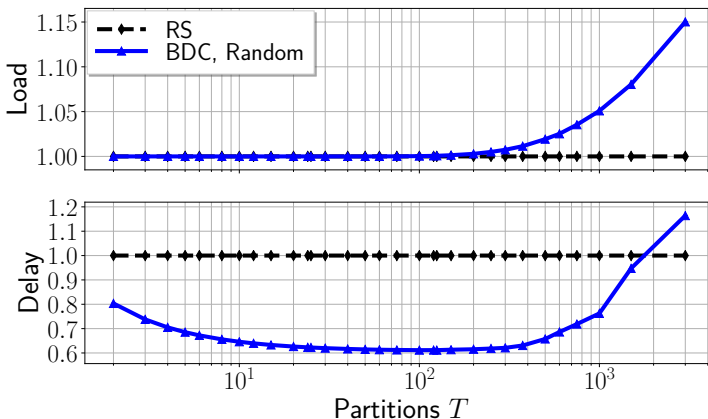
- 2000 rows assigned to each server, $n = 10000$ columns, $m/T = 10$ rows per partition, and code rate $m/r = 2/3$.

Numerical results



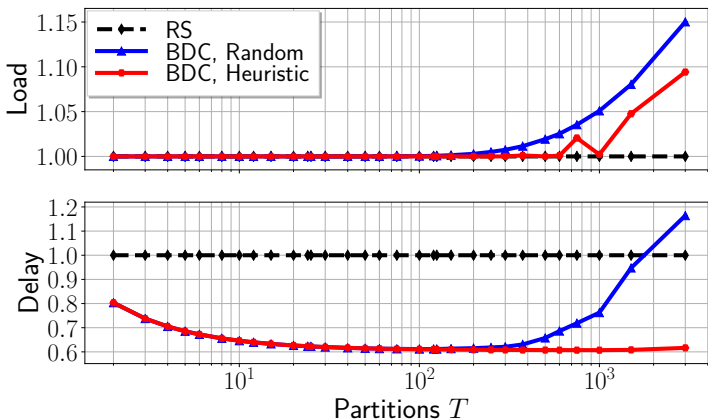
- $m = 6000$ source matrix rows, $n = 6000$ columns, $N = 6$ vectors, $K = 9$ servers, and code rate $m/r = 2/3$.

Numerical results



- $m = 6000$ source matrix rows, $n = 6000$ columns, $N = 6$ vectors, $K = 9$ servers, and code rate $m/r = 2/3$.

Numerical results



- $m = 6000$ source matrix rows, $n = 6000$ columns, $N = 6$ vectors, $K = 9$ servers, and code rate $m/r = 2/3$.

Optimal Assignment

Theorem

- Up to a given number of partitions T , the
 - communication load, and the
 - computational delay of the matrix multiplication, i.e., not taking the decoding time into account,
- of the BDC scheme are equal to those of the scheme by Li *et al.*.

Optimal Assignment

Theorem

- Up to a given number of partitions T , the
 - communication load, and the
 - computational delay of the matrix multiplication, i.e., not taking the decoding time into account,
- of the BDC scheme are equal to those of the scheme by Li *et al.*

The overall computational delay of our scheme is much lower than that of the scheme by Li *et al.* due to its lower decoding complexity.

Conclusion

Take-home message...

- The **delay** of the scheme by Li *et al.* is dominated by the decoding time. Our scheme addresses this issue with little to no impact on the **load**.

Conclusion

Take-home message...

- The **delay** of the scheme by Li *et al.* is dominated by the decoding time. Our scheme addresses this issue with little to no impact on the **load**.
 - $\approx 40\%$ reduced **delay** over the scheme by Li *et al.* for a matrix with 6000 rows and columns with **no** impact on **load**.

Conclusion

Take-home message...

- The **delay** of the scheme by Li *et al.* is dominated by the decoding time. Our scheme addresses this issue with little to no impact on the **load**.
 - $\approx 40\%$ reduced **delay** over the scheme by Li *et al.* for a matrix with 6000 rows and columns with **no** impact on **load**.
 - $> 90\%$ reduced delay over the scheme by Li *et al.* for a matrix with 134000 rows and 10000 columns with $< 1\%$ increased **load**.

Conclusion

Take-home message...

- The **delay** of the scheme by Li *et al.* is dominated by the decoding time. Our scheme addresses this issue with little to no impact on the **load**.
 - $\approx 40\%$ reduced **delay** over the scheme by Li *et al.* for a matrix with 6000 rows and columns with **no** impact on **load**.
 - $> 90\%$ reduced delay over the scheme by Li *et al.* for a matrix with 134000 rows and 10000 columns with $< 1\%$ increased **load**.
- Full version of the paper will soon be available on arXiv.

Conclusion

Take-home message...

- The **delay** of the scheme by Li *et al.* is dominated by the decoding time. Our scheme addresses this issue with little to no impact on the **load**.
 - $\approx 40\%$ reduced **delay** over the scheme by Li *et al.* for a matrix with 6000 rows and columns with **no** impact on **load**.
 - $> 90\%$ reduced delay over the scheme by Li *et al.* for a matrix with 134000 rows and 10000 columns with $< 1\%$ increased **load**.
- Full version of the paper will soon be available on arXiv.
- Encoding delay, LT codes under inactivation decoding...

Conclusion

Take-home message...

- The **delay** of the scheme by Li *et al.* is dominated by the decoding time. Our scheme addresses this issue with little to no impact on the **load**.
 - $\approx 40\%$ reduced **delay** over the scheme by Li *et al.* for a matrix with 6000 rows and columns with **no** impact on **load**.
 - $> 90\%$ reduced delay over the scheme by Li *et al.* for a matrix with 134000 rows and 10000 columns with $< 1\%$ increased **load**.
- Full version of the paper will soon be available on arXiv.
- Encoding delay, LT codes under inactivation decoding...
- LT codes may reduce delay further at the expense of increased load.

Conclusion

Take-home message...

- The **delay** of the scheme by Li *et al.* is dominated by the decoding time. Our scheme addresses this issue with little to no impact on the **load**.
 - $\approx 40\%$ reduced **delay** over the scheme by Li *et al.* for a matrix with 6000 rows and columns with **no** impact on **load**.
 - $> 90\%$ reduced delay over the scheme by Li *et al.* for a matrix with 134000 rows and 10000 columns with $< 1\%$ increased **load**.
- Full version of the paper will soon be available on arXiv.
- Encoding delay, LT codes under inactivation decoding...
- LT codes may reduce delay further at the expense of increased load.
- Slides and code on Github: github.com/severinson/coded-computing-tools

References

- [1] K. Lee et al. “Speeding up distributed machine learning using codes”. In: *Proc. IEEE Int. Symp. Inf. Theory*. Barcelona, Spain, July 2016, pp. 1143–1147. DOI: 10.1109/ISIT.2016.7541478.
- [2] S. Li, M. A. Maddah-Ali, and A. S. Avestimehr. “Coded MapReduce”. In: *Proc. Annual Allerton Conf. Commun., Control, and Computing*. Monticello, IL, Sept. 2015, pp. 964–971. DOI: 10.1109/ALLERTON.2015.7447112.
- [3] Songze Li, Mohammad Ali Maddah-Ali, and Amir Salman Avestimehr. “A Unified Coding Framework for Distributed Computing with Straggling Servers”. In: *Proc. Workshop Network Coding and Appl.* Washington, DC, Dec. 2016. DOI: 10.1109/GLOCOMW.2016.7848828.