# Coding for Distributed Computing

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

† Department of Electrical Engineering, Chalmers University of Technology, Gothenburg, Sweden
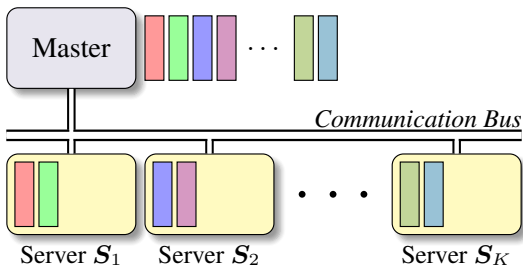‡ University of Bergen/Simula Research Lab, Bergen, Norway

Finse, May 09, 2018

simula@uib

CHALMERS

Introduction | Block-diagonal coding | LT code-based scheme | Numerical results | Conclusion | One More Thing...

simula@uib

## Motivation



## Challenges

- Straggler problem: May induce a large computational delay.
- Bandwidth scarcity: Need to reduce the communication load.

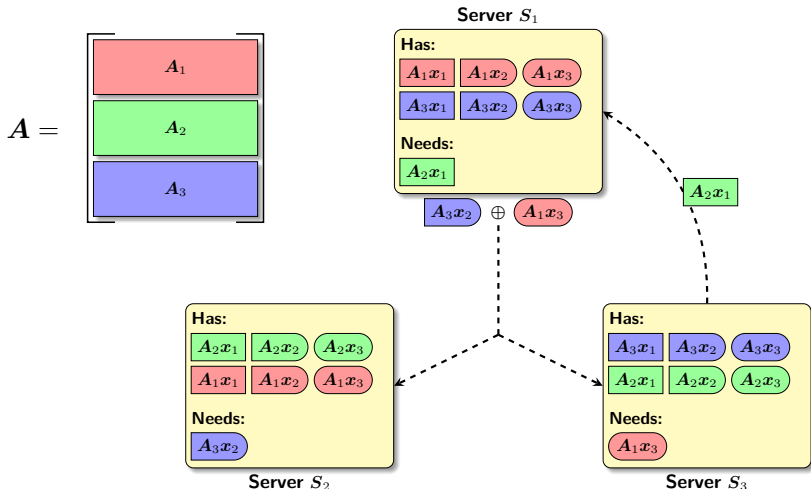## Problem addressed: Matrix multiplication

- Given an $m \times n$ matrix $\boldsymbol{A}$ and $N$ vectors $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_N$, we want to compute $\boldsymbol{y}_1 = \boldsymbol{A}\boldsymbol{x}_1, \ldots, \boldsymbol{y}_N = \boldsymbol{A}\boldsymbol{x}_N$ using $K$ servers.

## Bandwidth Scarcity
(Coded MapReduce, Li *et al.*, 2015)

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

Introduction | Block-diagonal coding | LT code-based scheme | Numerical results | Conclusion | One More Thing...
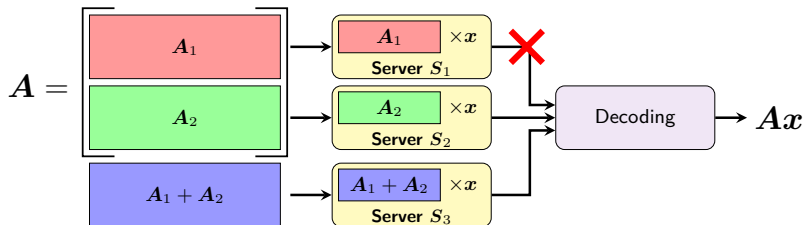
simula@uib

# The straggler problem

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

Introduction
○○○●○○○

Block-diagonal coding
○○○

LT code-based scheme
○

Numerical results
○○○

Conclusion
○

One More Thing…
○
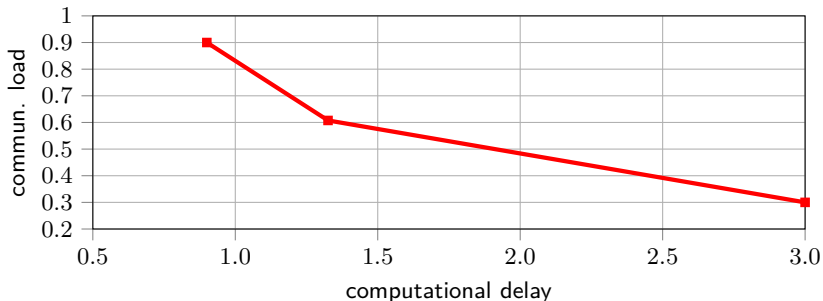
simula@uib

# 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!

Introduction    Block-diagonal coding    LT code-based scheme    Numerical results    Conclusion    One More Thing…
○○○○●○○    ○○○    ○    ○○○    ○    ○

simula@uib

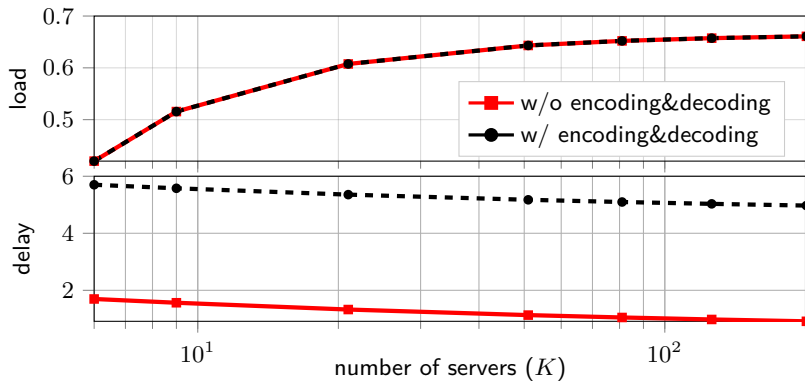## Coding for distributed computing

- [Lee *et al.* '17]: Introduce redundant computations using MDS codes to alleviate the straggler problem.
- [Li, Maddah-Ali, Avestimehr '17]: A fundamental tradeoff between computational delay and communication load. A unified coding framework trading higher computational delay for lower communication load.

## Unified coding framework [Li, Maddah-Ali, Avestimehr '17]

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

- Code length $r$ proportional to number of rows of $A \rightarrow$ high overall delay!



- $A$ with $n = 10000$ columns and $m = 2000K/3$ rows, $N = 2000K/3$ vectors, and code rate $2/3$ (2000 rows assigned to each server).

simula@uib

# In this talk

Two coding schemes to reduce the overall computational delay

- Block-diagonal coding scheme, based on a block-diagonal encoding matrix and shorter MDS codes.
- LT code-based scheme under inactivation decoding.

Outcome
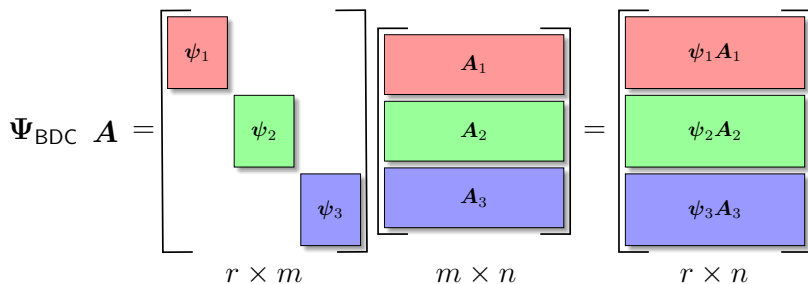
- Block-diagonal coding scheme: Significantly lower overall computational delay than the scheme by [Li, Maddah-Ali, Avestimehr '17] with no or little impact on communication load.
- LT code-based scheme: Very good performance when requiring to meet a deadline with high probability, at the expense of an increased communication load.

Introduction   Block-diagonal coding   LT code-based scheme   Numerical results   Conclusion   One More Thing...
○○○○○○○       ●○○                     ○                      ○○○                 ○            ○

simula@uib

## Block-diagonal coding scheme

**Idea**

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

$$C = \Psi_{\mathrm{BDC}} A, \quad \Psi_{\mathrm{BDC}} = \begin{bmatrix} \psi_1 & & \\ & \ddots & \\ & & \psi_T \end{bmatrix}, \quad \psi_i : \left(\frac{r}{T}, \frac{m}{T}\right) \text{ MDS code.}$$



$$\Psi_{\mathrm{BDC}} \ A = \qquad \qquad \qquad \qquad =$$

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

- Need any $m/T$ out of $r/T$ rows from each partition to decode.

Introduction  Block-diagonal coding  LT code-based scheme  Numerical results  Conclusion  One More Thing...
○○○○○○○  ○●○  ○  ○○○  ○  ○

simula@uib

## Assignment of coded rows to servers



- 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.

Introduction | Block-diagonal coding | LT code-based scheme | Numerical results | Conclusion | One More Thing...
0000000      000●                      0                     000               0           0

simula@uib

# Lossless partitioning

### Theorem

For $T \leq r / \binom{K}{\mu q}$, there exists an assignment matrix such that the communication load and the computational delay (not taking encoding/decoding delay into account) are equal to those of the unpartitioned scheme by [Li, Maddah-Ali, Avestimehr '17].

### However...

The overall computational delay of the block-diagonal coding scheme is much lower than that of the scheme by Li *et al.* due to its lower encoding and decoding complexity.
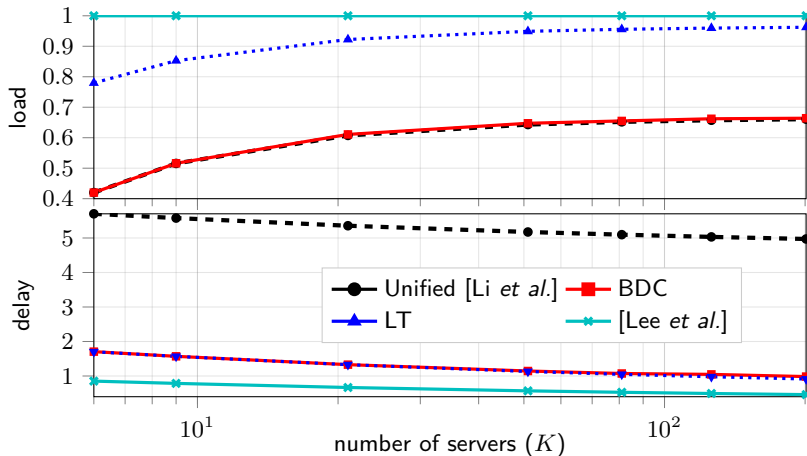
Introduction | Block-diagonal coding | LT code-based scheme | Numerical results | Conclusion | One More Thing...
0000000    000                    ●                         000                 0                  0

simula@uib

## Luby-transform code-based scheme

**LT code-based scheme**

- Encode $A$ as $C = \Psi_{\text{LT}} A$; $\Psi_{\text{LT}}$ corresponds to an LT code of fixed rate.
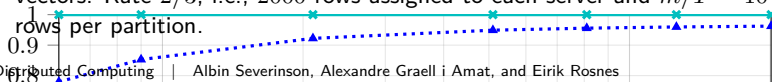- Decode the LT code using inactivation decoding.

**Code design**

- Design the LT code for a minimum overhead $\epsilon_{\text{min}}$ and a target failure probability $P_{\text{f,target}}$, such that $P_{\text{f}}(\epsilon_{\text{min}}) \leq P_{\text{f,target}}$.
- Increasing $\epsilon_{\text{min}}$ leads to lower encoding/decoding complexity but increased communication load and may require waiting for more servers $\longrightarrow$ optimal $\epsilon_{\text{min}}$ depends on the scenario.
- For a given $\epsilon_{\text{min}}$ and $P_{\text{f,target}}$, optimize the LT code so that the decoding complexity is minimized: for a fixed computational delay of $C\boldsymbol{x}_1, \ldots, C\boldsymbol{x}_N$, minimize the computational delay of the decoding phase.

Introduction | Block-diagonal coding | LT code-based scheme | Numerical results | Conclusion | One More Thing…
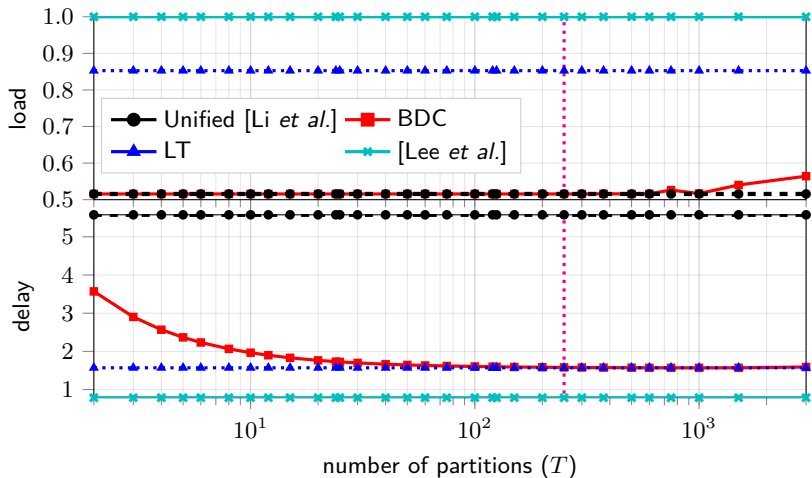0000000 | 000 | 0 | ●00 | 0 | 0

simula@uib

## Computational delay and communication load



- $\boldsymbol{A}$ with $n = 10000$ columns and $m = 2000K/3$ rows. $N = 2000K/3$ vectors. Rate $2/3$, i.e., $2000$ rows assigned to each server and $m/T = 10$ rows per partition.

Introduction | Block-diagonal coding | LT code-based scheme | Numerical results | Conclusion | One More Thing…
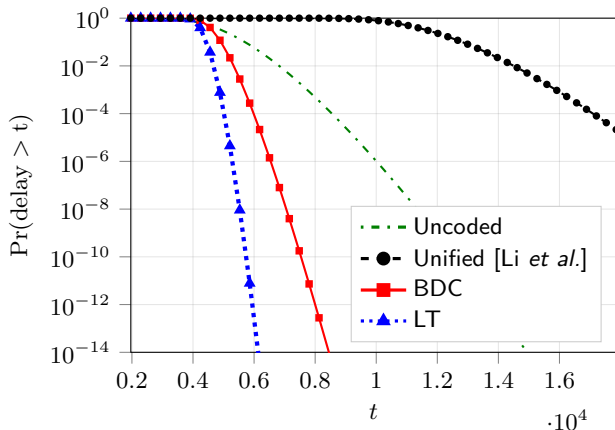0000000 | 000 | 0 | 0●0 | 0 | 0

simula@uib

## Performance as a function of the number of partitions



- $\boldsymbol{A}$ with $m = 6000$ rows and $n = 6000$ columns, $N = 6$ vectors, $K = 9$ servers, and code rate $2/3$.

Introduction   Block-diagonal coding   LT code-based scheme   Numerical results   Conclusion   One More Thing…
○○○○○○○       ○○○                    ○                      ○○●                 ○              ○

simula@uib

## Distributed computing under a deadline



- $\boldsymbol{A}$ with $m = 134000$ rows and $n = 10000$ columns, $N = 134000$ vectors, $K = 201$ servers, $T = 13400$ partitions, and code rate $2/3$.

Introduction   Block-diagonal coding   LT code-based scheme   Numerical results   Conclusion   One More Thing…
○○○○○○○   ○○○   ○   ○○○   ●   ○

simula@uib

## Conclusion

**Take-home message…**

- The encoding and decoding delay may contribute significantly to the overall computational delay.

- The BDC scheme yields significantly lower computational delay (up to 70%–80%) with no or little impact on the communication load.

- The LT code-based scheme achieves very good performance when needing to meet a deadline with high probability.

- Paper available in arxiv:
A. Severinson, A. Graell i Amat, and E. Rosnes, "Block-Diagonal and LT Codes for Distributed Computing With Straggling Servers".

- Code on Github: github.com/severinson/coded-computing-tools

Introduction
0000000
Block-diagonal coding
000
LT code-based scheme
0
Numerical results
000
Conclusion
0
One More Thing...
•

simula@uib

# One More Thing...