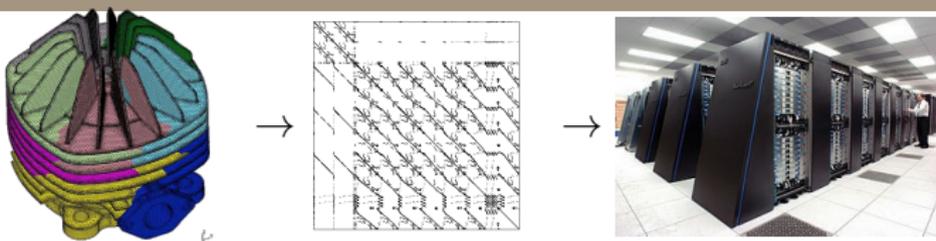


# Bridging the gap between flat and hierarchical low-rank matrix formats

P. Amestoy<sup>1</sup>   A. Buttari<sup>2</sup>   J.-Y. L'Excellent<sup>3</sup>   T. Mary<sup>4</sup>

<sup>1</sup>INP-IRIT   <sup>2</sup>CNRS-IRIT   <sup>3</sup>INRIA-LIP   <sup>4</sup>University of Manchester

Structured Matrix Days, Lyon, 14-15 May 2018



Linear system  $Ax = b$

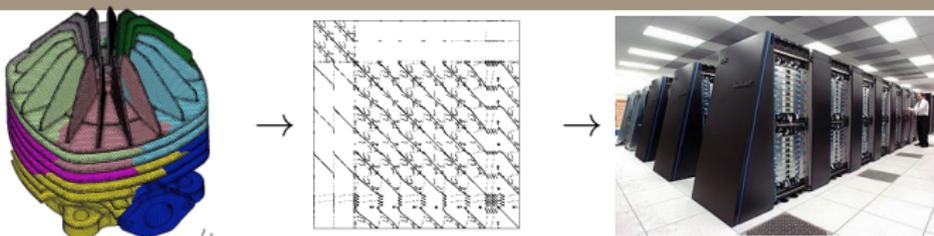
Often a keystone in **scientific computing applications**  
(discretization of PDEs, step of an optimization method, ...)

Direct methods

Factorize  $A = LU$  and solve  $LUx = b$

😊 Numerically reliable

☹ Computational cost



Linear system  $Ax = b$

Often a keystone in **scientific computing applications**  
(discretization of PDEs, step of an optimization method, ...)

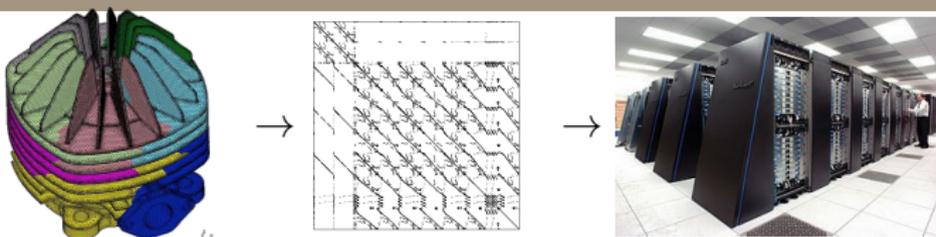
Direct methods

Factorize  $A = LU$  and solve  $LUx = b$

😊 Numerically reliable

☹ Computational cost

**Objective:**  
**reduce the cost of direct methods ...**  
**...while maintaining their numerical reliability**

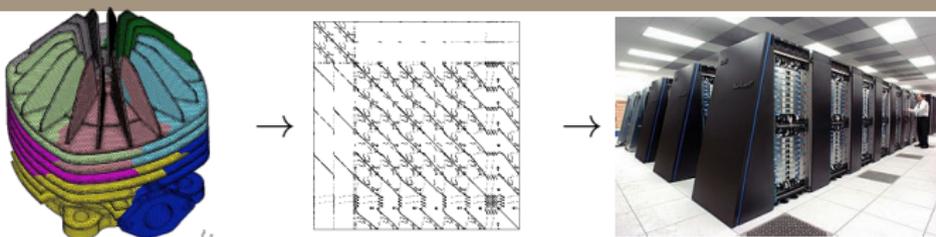


## Large scale applications

- Target size is  $n \sim 10^9$  for sparse  $\Rightarrow m \sim 10^6$  for dense
  - $O(m^2)$  storage complexity and  $O(m^3)$  flop complexity  
 $m \sim 10^6 \Rightarrow$  **TeraBytes** of storage and **ExaFlops** of computation!
- $\Rightarrow$  Need to reduce the asymptotic complexity

## Large scale systems

- Increasingly **large numbers of cores** available
- $\Rightarrow$  Need to design **parallel algorithms**



## Large scale applications

- Target size is  $n \sim 10^9$  for sparse  $\Rightarrow m \sim 10^6$  for dense
  - $O(m^2)$  storage complexity and  $O(m^3)$  flop complexity  
 $m \sim 10^6 \Rightarrow$  **TeraBytes** of storage and **ExaFlops** of computation!
- $\Rightarrow$  Need to reduce the asymptotic complexity

## Large scale systems

- Increasingly **large numbers of cores** available
- $\Rightarrow$  Need to design **parallel algorithms**

**These two objectives are not always compatible**

## 1. Introduction

- The  $\mathcal{H}$  format: very good complexity
- The BLR format: very good parallelism

## 2. Motivation

- Why we need a new format to bridge the gap

## 3. The MBLR format

- Complexity analysis
- Numerical results

## 4. Conclusion

### Preprint



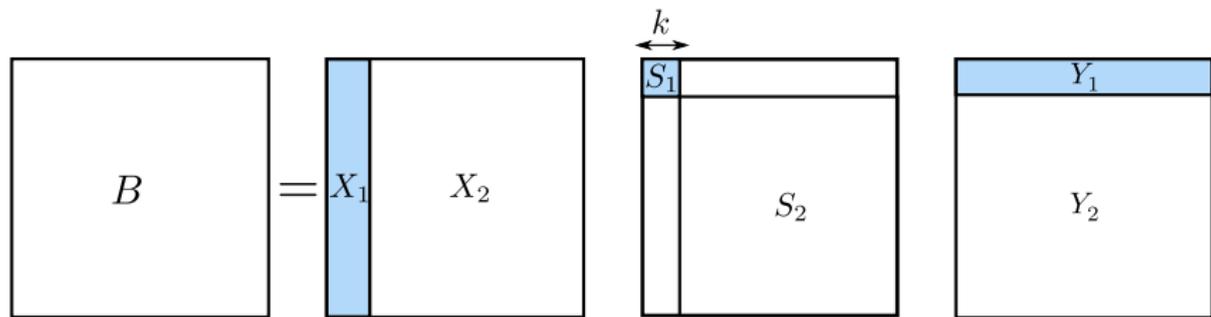
P. Amestoy, A. Buttari, J.-Y. L'Excellent, and T. Mary, *Bridging the gap between flat and hierarchical low-rank matrix formats: the multilevel BLR format*, submitted (2018).

# Introduction

Take a dense matrix  $B$  of size  $b \times b$  and compute its SVD  $B = XSY$ :

$$\boxed{B} = \boxed{X} \boxed{S} \boxed{Y}$$

Take a dense matrix  $B$  of size  $b \times b$  and compute its SVD  $B = XSY$ :



$k = \min \{k \leq n; \sigma_{k+1} \leq \varepsilon\}$  is the **numerical rank at accuracy  $\varepsilon$**

Take a dense matrix  $B$  of size  $b \times b$  and compute its SVD  $B = XSY$ :



$k = \min \{k \leq n; \sigma_{k+1} \leq \varepsilon\}$  is the **numerical rank at accuracy  $\varepsilon$**

$\tilde{B} = X_1 S_1 Y_1$  is a **low-rank approximation** to  $B$ :  $\|B - \tilde{B}\|_2 \leq \varepsilon$

Take a dense matrix  $B$  of size  $b \times b$  and compute its SVD  $B = XSY$ :



$k = \min \{k \leq n; \sigma_{k+1} \leq \varepsilon\}$  is the **numerical rank at accuracy  $\varepsilon$**

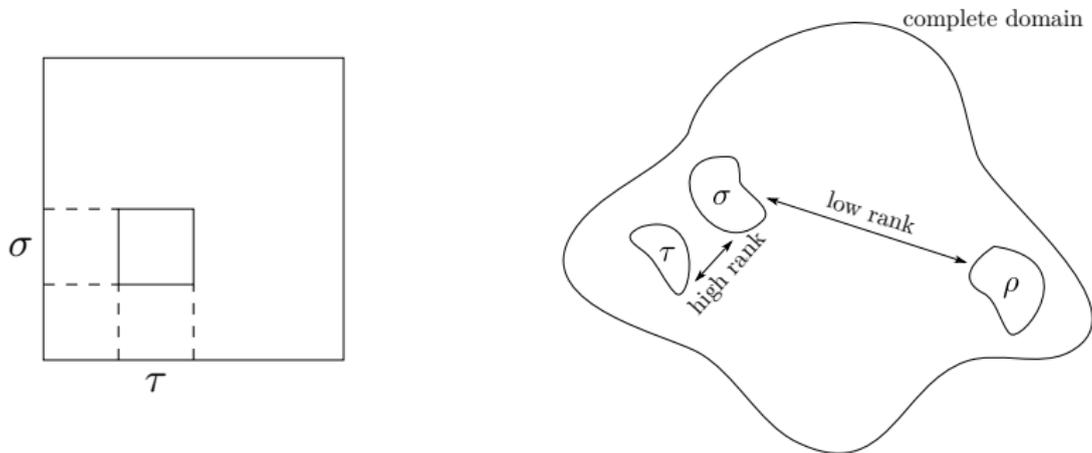
$\tilde{B} = X_1 S_1 Y_1$  is a **low-rank approximation** to  $B$ :  $\|B - \tilde{B}\|_2 \leq \varepsilon$

Storage savings:  $b^2/2bk = b/2k$

Similar flops savings when used in most linear algebra kernels

# Low-rank blocks

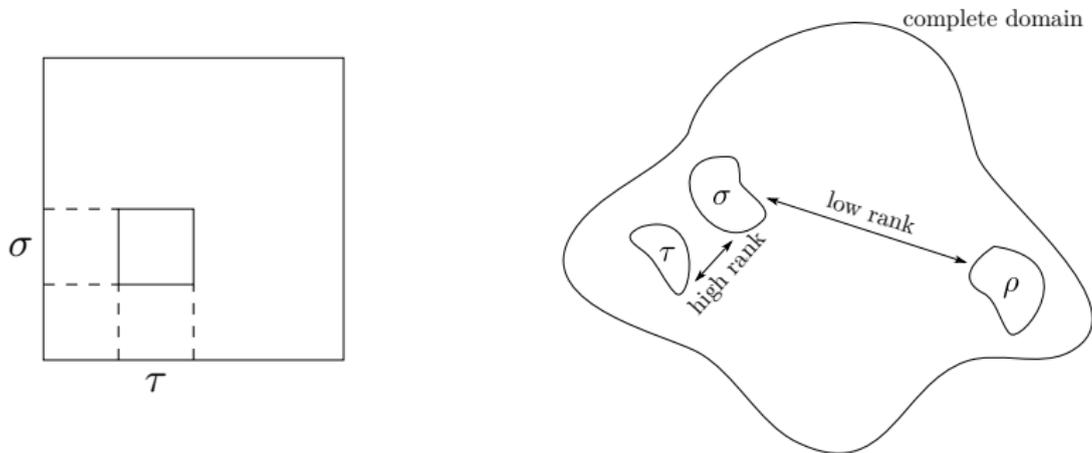
Most matrices are not low-rank in general but in some applications they exhibit **low-rank blocks**



A block  $B$  represents the interaction between two subdomains  $\sigma$  and  $\tau$ .

**Small diameter** and **far away**  $\Rightarrow$  low numerical rank.

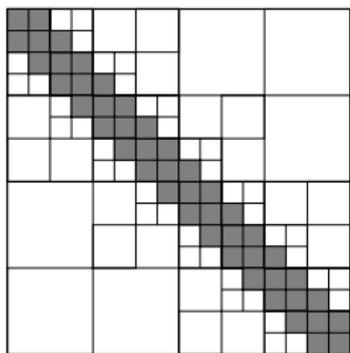
Most matrices are not low-rank in general but in some applications they exhibit **low-rank blocks**



A block  $B$  represents the interaction between two subdomains  $\sigma$  and  $\tau$ .

**Small diameter** and **far away**  $\Rightarrow$  low numerical rank.

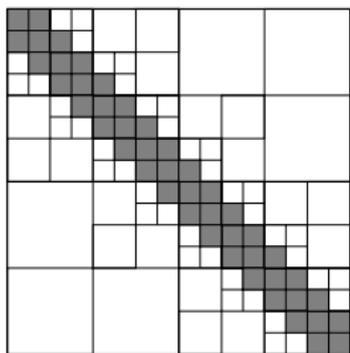
**How to choose a good block partitioning of the matrix?**



$\mathcal{H}$ -matrix

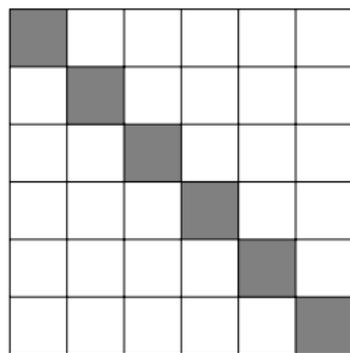
- Nearly linear complexity
- Complex, hierarchical structure

# $\mathcal{H}$ and BLR matrices



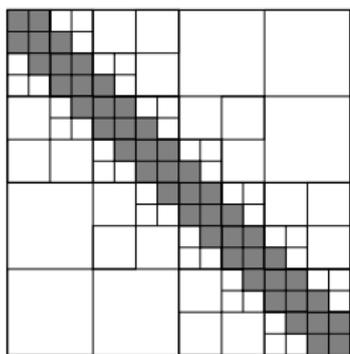
$\mathcal{H}$ -matrix

- Nearly linear complexity
- Complex, hierarchical structure

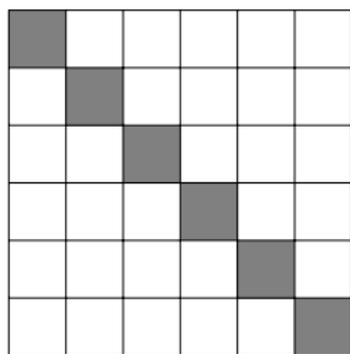


BLR matrix

- Superlinear complexity
- Simple, flat structure



$\mathcal{H}$ -matrix

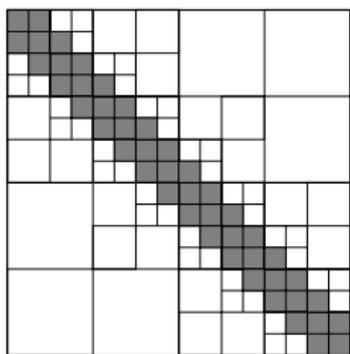


BLR matrix

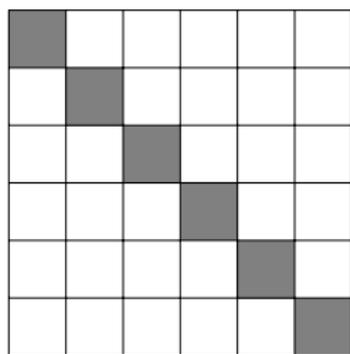
- Nearly linear complexity
- Complex, hierarchical structure
- Superlinear complexity
- Simple, flat structure

**BLR is a compromise between complexity and performance:**

- Small blocks  $\Rightarrow$  can fit on **single shared-memory** node
- No global order between blocks  $\Rightarrow$  **flexible data distribution**
- Easy to handle **numerical pivoting**



$\mathcal{H}$ -matrix



BLR matrix

- Nearly linear complexity
- Complex, hierarchical structure
- Superlinear complexity
- Simple, flat structure

## **BLR is a compromise between complexity and performance:**

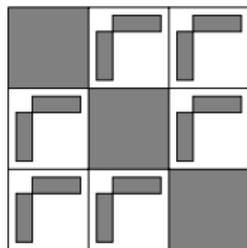
- Small blocks  $\Rightarrow$  can fit on **single shared-memory** node
- No global order between blocks  $\Rightarrow$  **flexible data distribution**
- Easy to handle **numerical pivoting**

## **Can we find an even better compromise?**

Motivation

# Computing the BLR complexity

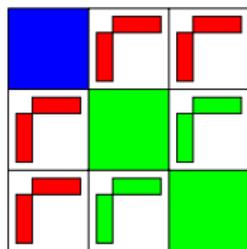
Assume all off-diagonal blocks are low-rank. Then:



$$\begin{aligned} \text{Storage} &= \text{cost}_{LR} * nb_{LR} + \text{cost}_{FR} * nb_{FR} \\ &= O(br) * O\left(\left(\frac{m}{b}\right)^2\right) + O(b^2) * O\left(\frac{m}{b}\right) \\ &= O(m^2r/b + mb) \\ &= \mathbf{O(m^{3/2}r^{1/2})} \text{ for } b = (mr)^{1/2} \end{aligned}$$

# Computing the BLR complexity

Assume all off-diagonal blocks are low-rank. Then:



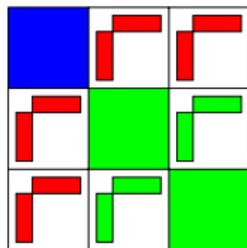
getrf  
trsm  
gemm

$$\begin{aligned} \text{Storage} &= \text{cost}_{LR} * nb_{LR} + \text{cost}_{FR} * nb_{FR} \\ &= O(br) * O\left(\left(\frac{m}{b}\right)^2\right) + O(b^2) * O\left(\frac{m}{b}\right) \\ &= O(m^2r/b + mb) \\ &= \mathbf{O(m^{3/2}r^{1/2})} \text{ for } b = (mr)^{1/2} \end{aligned}$$

$$\begin{aligned} \text{FlopLU} &= \text{cost}_{\text{getrf}} * nb_{\text{getrf}} + \text{cost}_{\text{trsm}} * nb_{\text{trsm}} + \text{cost}_{\text{gemm}} * nb_{\text{gemm}} \\ &= O(b^3) * O\left(\frac{m}{b}\right) + O(b^2r) * O\left(\left(\frac{m}{b}\right)^2\right) + O(br^2) * O\left(\left(\frac{m}{b}\right)^3\right) \\ &= O(mb^2 + m^2r + m^3r^2/b^2) \\ &= \mathbf{O(m^2r)} \text{ for } b = (mr)^{1/2} \end{aligned}$$

# Computing the BLR complexity

Assume all off-diagonal blocks are low-rank. Then:



getrf  
trsm  
gemm

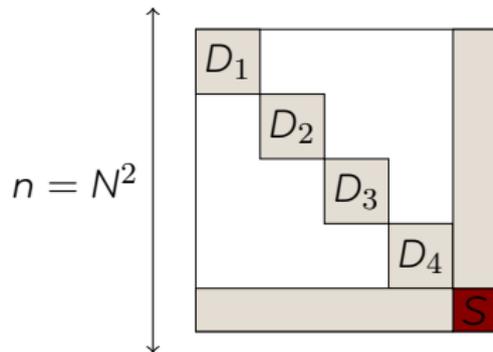
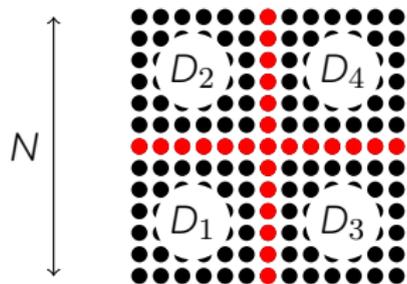
$$\begin{aligned} \text{Storage} &= \text{cost}_{LR} * nb_{LR} + \text{cost}_{FR} * nb_{FR} \\ &= O(br) * O\left(\left(\frac{m}{b}\right)^2\right) + O(b^2) * O\left(\frac{m}{b}\right) \\ &= O(m^2r/b + mb) \\ &= \mathbf{O(m^{3/2}r^{1/2})} \text{ for } b = (mr)^{1/2} \end{aligned}$$

$$\begin{aligned} \text{FlopLU} &= \text{cost}_{\text{getrf}} * nb_{\text{getrf}} + \text{cost}_{\text{trsm}} * nb_{\text{trsm}} + \text{cost}_{\text{gemm}} * nb_{\text{gemm}} \\ &= O(b^3) * O\left(\frac{m}{b}\right) + O(b^2r) * O\left(\left(\frac{m}{b}\right)^2\right) + O(br^2) * O\left(\left(\frac{m}{b}\right)^3\right) \\ &= O(mb^2 + m^2r + m^3r^2/b^2) \\ &= \mathbf{O(m^2r)} \text{ for } b = (mr)^{1/2} \end{aligned}$$

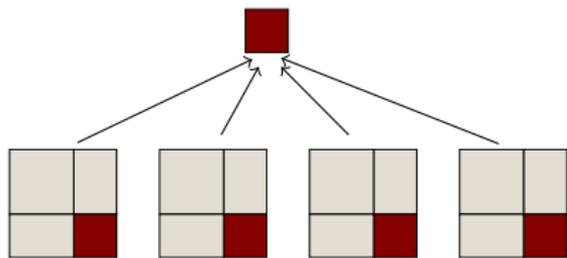
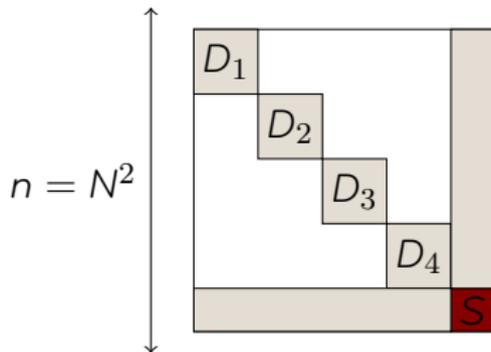
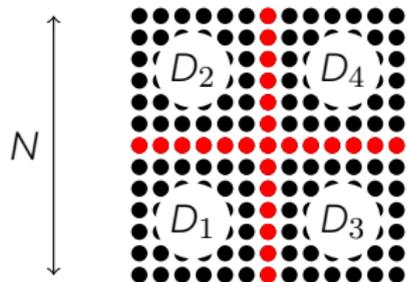
Result holds if a **constant** number of off-diag. blocks is full-rank.

 P. Amestoy, A. Buttari, J.-Y. L'Excellent, and T. Mary, *On the Complexity of the Block Low-Rank Multifrontal Factorization*, SIAM J. Sci. Comput. (2017).

# From dense to sparse: nested dissection



# From dense to sparse: nested dissection



Proceed recursively to  
compute **separator tree**

Factorizing a sparse matrix  
amounts to factorizing a  
sequence of dense matrices

⇒

**sparse complexity is directly  
derived from dense one**

$$\mathbf{2D:} \quad C_{\text{sparse}} = \sum_{\ell=0}^{\log N} 4^{\ell} C_{\text{dense}}\left(\frac{N}{2^{\ell}}\right)$$

# Nested dissection complexity formulas

$$\mathbf{2D:} \quad C_{\text{sparse}} = \sum_{\ell=0}^{\log N} 4^{\ell} C_{\text{dense}}\left(\frac{N}{2^{\ell}}\right)$$

$$\mathbf{3D:} \quad C_{\text{sparse}} = \sum_{\ell=0}^{\log N} 8^{\ell} C_{\text{dense}}\left(\frac{N^2}{4^{\ell}}\right)$$

# Nested dissection complexity formulas

$$\mathbf{2D:} \quad C_{\text{sparse}} = \sum_{\ell=0}^{\log N} 4^{\ell} C_{\text{dense}}\left(\frac{N}{2^{\ell}}\right) \quad \rightarrow \text{common ratio } 2^{2-\beta}$$

$$\mathbf{3D:} \quad C_{\text{sparse}} = \sum_{\ell=0}^{\log N} 8^{\ell} C_{\text{dense}}\left(\frac{N^2}{4^{\ell}}\right) \quad \rightarrow \text{common ratio } 2^{3-2\beta}$$

Assume  $C_{\text{dense}} = O(m^{\beta})$ . Then:

2D		3D	
	$C_{\text{sparse}}(n)$		$C_{\text{sparse}}(n)$
$\beta > 2$	$O(n^{\beta/2})$	$\beta > 1.5$	$O(n^{2\beta/3})$
$\beta = 2$	$O(n \log n)$	$\beta = 1.5$	$O(n \log n)$
$\beta < 2$	$O(n)$	$\beta < 1.5$	$O(n)$

$\Rightarrow$  **Key motivation:**  $C_{\text{dense}} < O(m^2)$  (2D) or  $O(m^{3/2})$  (3D)  
is enough to get optimal sparse complexity!

# Bridging the gap between flat and hierarchical formats

	$\mathcal{C}_{dense}$	Storage		$\mathcal{C}_{dense}$	Flop LU	
		2D	$\mathcal{C}_{sparse}$ 3D		2D	$\mathcal{C}_{sparse}$ 3D
FR	$O(m^2)$	$O(n \log n)$	$O(n^{4/3})$	$O(m^3)$	$O(n^{3/2})$	$O(n^2)$
BLR	$O(m^{3/2})$	$O(n)$	$O(n \log n)$	$O(m^2)$	$O(n \log n)$	$O(n^{4/3})$
$\mathcal{H}$	$O(m \log m)$	$O(n)$	$O(n)$	$O(m \log^2 m)$	$O(n)$	$O(n)$

Motivation:

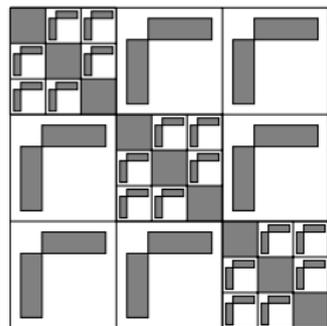
- 2D flop and 3D storage complexity: can we find a simple way to improve **just a little**  $\mathcal{C}_{dense}$ ?
- 3D flop complexity: still a large **gap** between BLR and  $\mathcal{H}$

**We propose a multilevel BLR format to bridge the gap**

# The MBLR format

# Complexity of the two-level BLR format

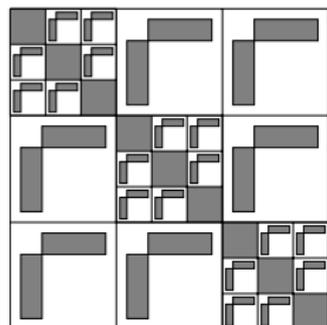
Assume all off-diagonal blocks are low-rank. Then:



$$\begin{aligned}\text{Storage} &= \text{cost}_{LR} * nb_{LR} + \text{cost}_{BLR} * nb_{BLR} \\ &= O(br) * O\left(\left(\frac{m}{b}\right)^2\right) + O(b^{3/2}r^{1/2}) * O\left(\frac{m}{b}\right) \\ &= O(m^2r/b + m(br)^{1/2}) \\ &= \mathbf{O(m^{4/3}r^{2/3})} \text{ for } b = (m^2r)^{1/3}\end{aligned}$$

# Complexity of the two-level BLR format

Assume all off-diagonal blocks are low-rank. Then:



$$\begin{aligned}\text{Storage} &= \text{cost}_{LR} * nb_{LR} + \text{cost}_{BLR} * nb_{BLR} \\ &= O(br) * O\left(\left(\frac{m}{b}\right)^2\right) + O(b^{3/2}r^{1/2}) * O\left(\frac{m}{b}\right) \\ &= O(m^2r/b + m(br)^{1/2}) \\ &= \mathbf{O(m^{4/3}r^{2/3})} \text{ for } b = (m^2r)^{1/3}\end{aligned}$$

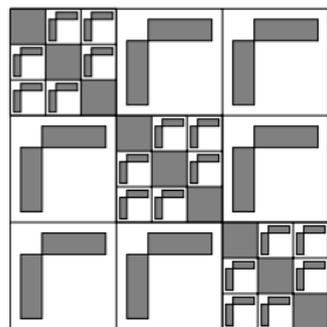
Similarly, we can prove:

$$\text{FlopLU} = \mathbf{O(m^{5/3}r^{4/3})} \text{ for } b = (m^2r)^{1/3}$$

Result holds if a **constant** number of off-diag. blocks is BLR.

# Complexity of the two-level BLR format

Assume all off-diagonal blocks are low-rank. Then:



$$\begin{aligned}
 \text{Storage} &= \text{cost}_{LR} * nb_{LR} + \text{cost}_{BLR} * nb_{BLR} \\
 &= O(br) * O\left(\left(\frac{m}{b}\right)^2\right) + O(b^{3/2}r^{1/2}) * O\left(\frac{m}{b}\right) \\
 &= O(m^2r/b + m(br)^{1/2}) \\
 &= \mathbf{O(m^{4/3}r^{2/3})} \text{ for } b = (m^2r)^{1/3}
 \end{aligned}$$

Similarly, we can prove:

$$\text{FlopLU} = \mathbf{O(m^{5/3}r^{4/3})} \text{ for } b = (m^2r)^{1/3}$$

Result holds if a **constant** number of off-diag. blocks is BLR.

		FR	BLR	2-BLR	...	$\mathcal{H}$
storage	dense	$O(m^2)$	$O(m^{1.5})$	$O(m^{1.33})$	...	$O(m \log m)$
	sparse	$O(n^{1.33})$	$O(n \log n)$	$O(n)$	...	$O(n)$
flop LU	dense	$O(m^3)$	$O(m^2)$	$O(m^{1.66})$	...	$O(m \log^3 m)$
	sparse	$O(n^2)$	$O(n^{1.33})$	$O(n^{1.11})$	...	$O(n)$

## Main result

For  $b = m^{\ell/(\ell+1)}r^{1/(\ell+1)}$ , the  $\ell$ -level complexities are:

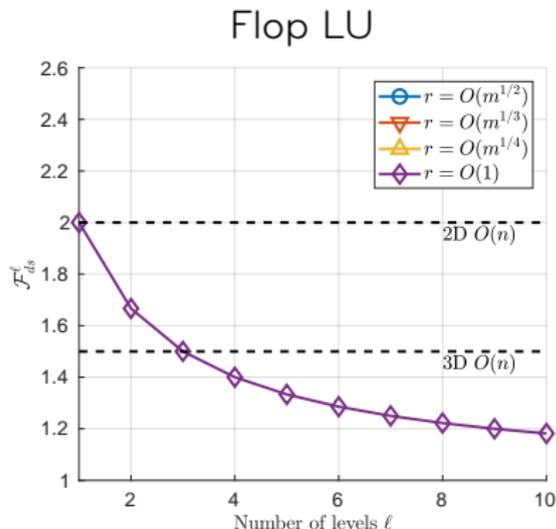
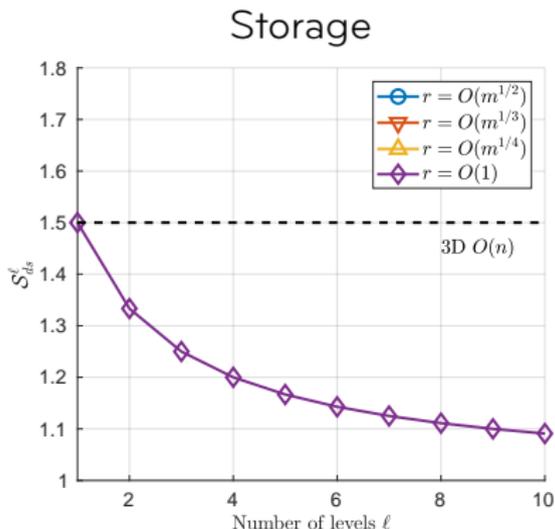
$$\text{Storage} = \mathbf{O}(m^{(\ell+2)/(\ell+1)}r^{\ell/(\ell+1)})$$

$$\text{FlopLU} = \mathbf{O}(m^{(\ell+3)/(\ell+1)}r^{2\ell/(\ell+1)})$$

Proof: by induction.  $\square$

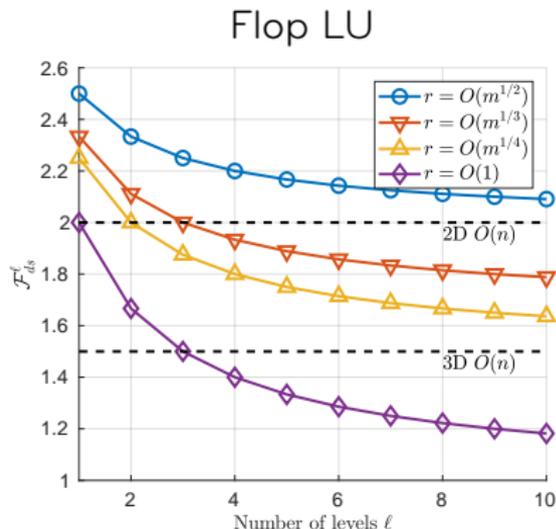
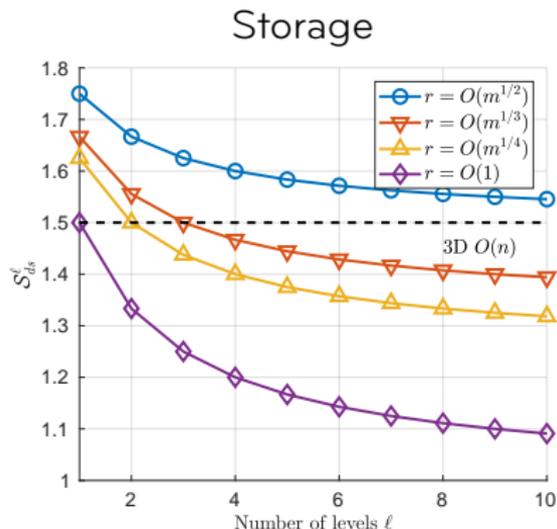
- Simple way to **finely control** the desired complexity
- Block size  $b \propto O(m^{\ell/(\ell+1)}) \ll O(m)$   
 $\Rightarrow$  may be efficiently processed in shared-memory
- Number of blocks per row/column  $\propto O(m^{1/(\ell+1)}) \gg O(1)$   
 $\Rightarrow$  flexibility to distribute data in parallel

# Influence of the number of levels $\ell$



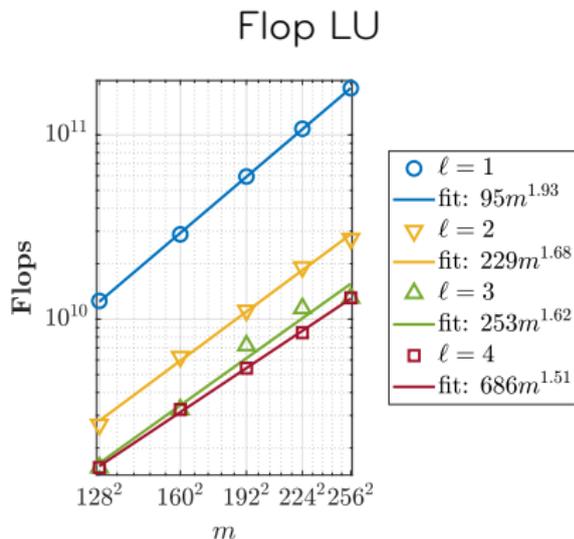
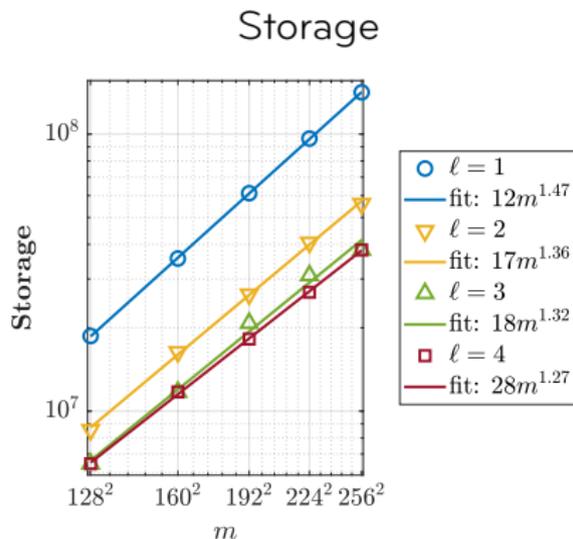
- If  $r = O(1)$ , can achieve  $O(n)$  storage complexity with only two levels and  $O(n \log n)$  flop complexity with three levels

# Influence of the number of levels $\ell$



- If  $r = O(1)$ , can achieve  $O(n)$  storage complexity with only two levels and  $O(n \log n)$  flop complexity with three levels
- For higher ranks, optimal sparse complexity is not attainable with constant  $\ell$  but improvement rate is rapidly decreasing: the first few levels achieve most of the asymptotic gain

# Numerical experiments (Poisson)



- Experimental complexity in relatively good agreement with theoretical one
- Asymptotic gain decreases with levels

Conclusion

## A new multilevel format to...

- **Finely control** desired complexity between BLR's and  $\mathcal{H}$ 's
- **Strike a balance** between BLR's simplicity and  $\mathcal{H}$ 's complexity
- Trade off  $\mathcal{H}$ 's nearly linear dense complexity and still achieve  $\mathcal{C}_{sparse} = O(n)$

## Future work: high-performance implementation

- Implementation of the MBLR format in a parallel, algebraic, general purpose sparse solver (e.g. **MUMPS**)



Thank you for  
your attention

Slides and paper available here:  
[personalpages.manchester.ac.uk/staff/theo.mary/](https://personalpages.manchester.ac.uk/staff/theo.mary/)