



Bridging the gap between flat and hierarchical low-rank matrix formats: the multilevel BLR format

Patrick Amestoy, Alfredo Buttari, Jean-Yves L'Excellent, Théo Mary

► To cite this version:

Patrick Amestoy, Alfredo Buttari, Jean-Yves L'Excellent, Théo Mary. Bridging the gap between flat and hierarchical low-rank matrix formats: the multilevel BLR format. [Research Report] 1-25, University of Manchester. 2018. hal-01774642v2

HAL Id: hal-01774642

<https://hal.science/hal-01774642v2>

Submitted on 22 Mar 2019 (v2), last revised 10 Oct 2019 (v3)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

BRIDGING THE GAP BETWEEN FLAT AND HIERARCHICAL LOW-RANK MATRIX FORMATS: THE MULTILEVEL BLR FORMAT*

PATRICK R. AMESTOY[†], ALFREDO BUTTARI[‡], JEAN-YVES L'EXCELLENT[§], AND THEO MARY[¶]

Abstract. Matrices possessing a low-rank property arise in numerous scientific applications. This property can be exploited to provide a substantial reduction of the complexity of their LU or LDL^T factorization. Among the possible low-rank formats, the flat Block Low-Rank (BLR) format is easy to use but achieves superlinear complexity. Alternatively, the hierarchical formats achieve linear complexity at the price of a much more complex, hierarchical matrix representation. In this paper, we propose a new format based on multilevel BLR approximations: the matrix is recursively defined as a BLR matrix whose full-rank blocks are themselves represented by BLR matrices. We call this format *multilevel BLR* (MBLR). Contrarily to hierarchical matrices, the number of levels in the block hierarchy is fixed to a given constant; while this format can still be represented within the \mathcal{H} formalism, we show that applying the \mathcal{H} theory to it leads to very pessimistic complexity bounds. We therefore extend the theory to prove better bounds, and show that the MBLR format provides a simple way to finely control the desired complexity of dense factorizations. By striking a balance between the simplicity of the BLR format and the low complexity of the hierarchical ones, the MBLR format bridges the gap between flat and hierarchical low-rank matrix formats. The MBLR format is of particular relevance in the context of sparse direct solvers, for which it is able to trade off the optimal dense complexity of the hierarchical formats to benefit from the simplicity and flexibility of the BLR format while still achieving $O(n)$ sparse complexity. We finally compare our MBLR format with the related BLR- \mathcal{H} (or Lattice- \mathcal{H}) format; our theoretical analysis shows that both formats achieve the same asymptotic complexity for a given top level block size.

Key words. low-rank approximations, matrix factorization, sparse linear algebra, hierarchical matrices

AMS subject classifications. 15a06, 15a23, 65f05, 65f50, 65n30, 65y20

1. Introduction. Efficiently computing the solution of a dense linear system is a fundamental building block of numerous scientific computing applications. Let us refer to such a system as

$$Fu_F = v_F, \quad (1.1)$$

where F is a dense matrix of order m , u_F is the unknown vector of size m , and v_F is the right-hand side vector of size m .

This paper focuses on solving (1.1) with direct approaches based on Gaussian elimination, which consist in factorizing matrix F as $F = LU$ or $F = LDL^T$, depending on whether the matrix is unsymmetric or symmetric, respectively.

In many applications (e.g., Schur complements arising from the discretization of elliptic partial differential equations), the matrix F has been shown to have a low-rank property: many of its off-diagonal blocks can be approximated by low-rank matrices [10].

Several formats have been proposed to exploit this property. The simplest one is the Block Low-Rank (BLR) format [2], which partitions the matrix with a flat, 2D blocking and approximates its off-diagonal blocks by low-rank submatrices, as illustrated in Figure 1.1a. Compared with the cubic $O(m^3)$ complexity of the dense full-rank LU or LDL^T factorizations, the complexity of the dense BLR factorization can be as low as $O(m^2)$ [4].

More advanced formats are based on a hierarchical partitioning of the matrix: the matrix F is partitioned with a 2×2 blocking and the two diagonal blocks are recursively refined, as illustrated in Figures 1.1c and 1.1d. Different hierarchical formats can be defined depending on whether the off-diagonal blocks are directly approximated (so-called *weakly-admissible* formats) or further refined (so-called *strongly-admissible* formats). The most general of the hierarchical formats is the strongly-admissible \mathcal{H} -matrix format [10, 21, 12]; the HODLR format [6] is its weakly-admissible counterpart. These hierarchical formats can factorize a dense matrix in near-linear complexity $O(m \log^q m)$, where q is a small integer that depends on which factorization algorithm is used; in

***Funding:** This work was supported by Engineering and Physical Sciences Research Council grant EP/P020720/1. The opinions and views expressed in this publication are those of the authors, and not necessarily those of the funding bodies.

[†]University of Toulouse, INP-IRIT

[‡]University of Toulouse, CNRS-IRIT

[§]University of Lyon, CNRS, ENS de Lyon, Inria, UCBL, LIP UMR5668, France

[¶]School of Mathematics, The University of Manchester, Manchester M13 9PL, UK (theo.mary@manchester.ac.uk)

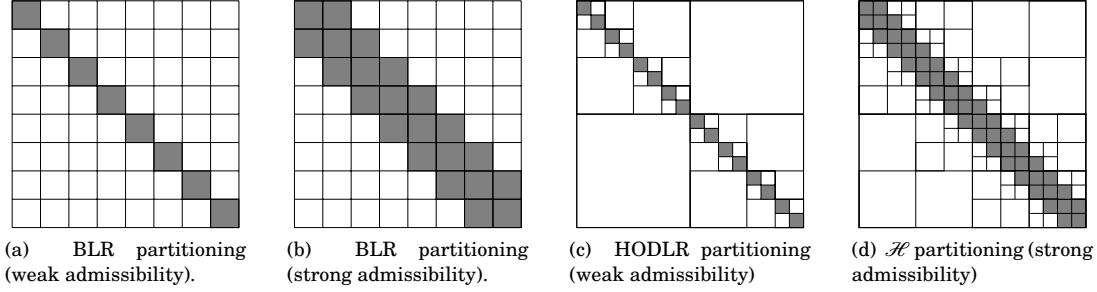


Fig. 1.1: Illustration of different low-rank formats. Gray blocks are stored in full-rank whereas white ones are approximated by low-rank matrices.

the following, we will consider $q = 2$. The log factor can be removed by using a so-called *nested-basis* structure. The strongly-admissible \mathcal{H}^2 -matrix format [12] and the weakly-admissible HSS [31, 13] and HBS [19] formats exploit such nested basis structures to achieve linear complexity $O(m)$.

In this paper, we propose a new format based on multilevel BLR approximations. The matrix F is recursively represented as a BLR matrix whose full-rank blocks are themselves BLR matrices. We call this format *multilevel BLR* (MBLR). In the hierarchical format, the matrix is refined until the diagonal blocks are of constant size; this therefore leads to a number of levels in the block hierarchy which is nonconstant (usually $O(\log_2 m)$). With the MBLR format, we propose to make this number of levels a tunable parameter to be set to a given value ℓ that does not asymptotically depend on the matrix size m . We prove that this parameter provides a simple way to finely control the complexity of the dense MBLR factorization. The complexity varies from $O(m^2)$ for monolevel BLR down to nearly $O(m)$ for an infinite number of levels. By striking a balance between the simplicity of the BLR format and the low complexity of the hierarchical ones, the MBLR format bridges the gap between flat and hierarchical low-rank matrix formats.

We will show that the MBLR format is of particular relevance in the context of sparse direct solvers, which aim to compute the solution of a sparse linear system

$$Au_A = v_A, \quad (1.2)$$

where A is a sparse matrix of order n , u_A is the unknown vector of size n , and v_A is the right-hand side vector of size n . Two widely studied classes of sparse direct methods are the multifrontal [15, 24] and supernodal [9, 14] approaches.

Sparse direct methods rely on a sequence of partial factorizations of dense matrices F , referred to as supernodes or fronts. Therefore, the complexity of sparse direct methods is directly derived from the complexity of the factorization of each dense matrix F . For example, with an adequate reordering, a well-known result is that the dense standard full-rank $O(m^3)$ factorization leads to a $O(n^2)$ sparse complexity for regular 3D problems [16]. The low-rank formats described above can be efficiently exploited within sparse solvers to provide a substantial reduction of their complexity.

The potential of BLR sparse solvers has been first investigated in [2]; the simplicity and flexibility of the BLR format makes it easy to use in the context of a general purpose, algebraic solver, as presented in [5, 3, 28, 26]. [5] focuses on the multicore performance of BLR multifrontal solvers, while [3] and [28] present their use in two real-life industrial applications coming from geosciences. [26] present the use of the BLR format in supernodal solvers. Furthermore, it has been proved in [4] that the theoretical complexity of the BLR multifrontal factorization may be as low as $O(n^{4/3})$ (for 3D problems with constant ranks).

Alternatively, most sparse solvers based on the more complex hierarchical formats have been shown to possess near-linear complexity. To cite a few, [30, 29, 18, 17] are HSS-based, [19] is HBS-based, [7] is HODLR-based, and [27] is \mathcal{H}^2 -based.

However, a critical observation is that achieving $O(n)$ sparse complexity does not actually require a linear dense complexity $O(m)$. For instance, for 3D problems, all that is required is a

dense complexity lower than $O(m^{1.5})$. Therefore, we will prove that the MBLR format is able to trade off the optimal dense complexity of the hierarchical formats to benefit from the simplicity and flexibility of the flat BLR format while still achieving $O(n)$ sparse complexity.

We now describe the organization of the rest of this paper. In Section 2, we provide some background on the BLR factorization and its complexity and we motivate the key idea behind MBLR approximations. We explain in Section 3 how the MBLR format can be described using the cluster tree representation commonly used in the \mathcal{H} literature; this provides a convenient way to explain the key difference between the MBLR and \mathcal{H} formats. We show that, similarly to the BLR case, the \mathcal{H} theoretical formalism leads to MBLR complexity bounds that are very pessimistic. We therefore extend the theory, beginning by the two-level case in Section 4. We prove that two levels can already significantly improve the theoretical complexity of the factorization. In Section 5, we generalize the previous proof to the MBLR format with an arbitrary number of levels; we prove that, for constant ranks, only four levels are already enough to reach $O(n)$ sparse 3D complexity (and three levels already achieve near-linear $O(n \log n)$ complexity). In Section 6, we validate our theoretical results with numerical experiments. We provide our concluding remarks in Section 7. In the main body of this article, we consider for the sake of simplicity the weakly-admissible case, in which only the diagonal blocks are refined. In the appendix, we provide the extension of the MBLR format to the strongly-admissible case, in which off-diagonal full-rank blocks are also recursively refined. We prove that our complexity bounds remain valid in this context.

2. Background and motivation.

2.1. Block Low-Rank approximations. The BLR format is based on a flat, non-hierarchical blocking of the matrix which is defined by conveniently clustering the associated unknowns. A BLR representation \tilde{F} of a dense matrix F is shown in (2.1), where we assume that $p \times p$ blocks have been defined. Off-diagonal blocks F_{ij} ($i \neq j$) of size $m_i \times n_j$ and numerical rank k_{ij}^ε are approximated by a low-rank matrix $\tilde{F}_{ij} = X_{ij}Y_{ij}^T$ at accuracy ε , where X_{ij} is a $m_i \times k_{ij}^\varepsilon$ matrix and Y_{ij} is a $n_j \times k_{ij}^\varepsilon$ matrix. The diagonal blocks F_{ii} are stored as full-rank matrices ($\tilde{F}_{ii} = F_{ii}$).

$$\tilde{F} = \begin{bmatrix} \tilde{F}_{11} & \tilde{F}_{12} & \cdots & \tilde{F}_{1p} \\ \tilde{F}_{21} & \cdots & \cdots & \vdots \\ \vdots & \cdots & \cdots & \vdots \\ \tilde{F}_{p1} & \cdots & \cdots & \tilde{F}_{pp} \end{bmatrix}. \quad (2.1)$$

Throughout this article, we will note F_{ij} the (i, j) -th block of F and $F_{:,k}$ its k -th block-column. We will also assume that all blocks have a size of order b , i.e., $m_i = n_j = O(b)$.

Computing the low-rank approximation \tilde{F}_{ij} to each block, referred to as the *compression* step, can be performed in different ways. We have chosen to use a truncated QR factorization with column pivoting; this corresponds to a QR factorization with pivoting which is truncated as soon as a diagonal coefficient of the R factor falls below the prescribed threshold ε . For a block of size $b \times b$ and rank r , the cost of the compression is $O(b^2 r)$, whereas computing the exact singular value decomposition of the block would require $O(b^3)$ operations. This choice thus allows for a convenient compromise between cost and accuracy of the compression operation.

2.2. Block Low-Rank dense factorization. We describe in Algorithm 2.1 the CUFs variant (the ‘‘CUFs’’ acronym is explained below) of the BLR factorization algorithm for dense matrices, introduced in [4]. Algorithm 2.1 is presented in its LU version, but it can easily be adapted to the symmetric case.

In order to perform the LU or LDL^T factorization of a dense BLR matrix, the standard block LU or LDL^T factorization has to be modified so that the low-rank blocks can be exploited to perform fast operations. Many such algorithms can be defined depending on where the compression step is performed. As described in [4], the CUFs variant achieves the lowest complexity of all BLR variants by performing the compression as early as possible.

This algorithm is referred to as CUFs (standing for Compress, Update, Factor, and Solve), to indicate the order in which the steps are performed. All low-rank updates of a given block \tilde{F}_{ik} are

Algorithm 2.1 Dense BLR LU factorization: CUFS variant.Input: a $p \times p$ block matrix F of order m .Output: F overwritten by its BLR LU factors \tilde{F} .

```

1: for  $k = 1$  to  $p$  do
2:   for  $i = k + 1$  to  $p$  do
3:     Compress ( $L$ ):  $F_{ik} \leftarrow \tilde{F}_{ik} = X_{ik} Y_{ik}^T$ 
4:     Compress ( $U$ ):  $F_{ki} \leftarrow \tilde{F}_{ki} = Y_{ki} X_{ki}^T$ 
5:   end for
6:   for  $i = k$  to  $p$  do
7:     for  $j = 1$  to  $k - 1$  do
8:       Update ( $L$ ):  $\tilde{F}_{ik} \leftarrow \tilde{F}_{ik} - X_{ij} Y_{ij}^T Y_{jk} X_{jk}^T$ 
9:       Update ( $U$ ):  $\tilde{F}_{ki} \leftarrow \tilde{F}_{ki} - X_{kj} Y_{kj}^T Y_{ji} X_{ji}^T$ 
10:    end for
11:     $\tilde{F}_{ik} \leftarrow \text{Recompress}(\tilde{F}_{ik})$ 
12:     $\tilde{F}_{ki} \leftarrow \text{Recompress}(\tilde{F}_{ki})$ 
13:  end for
14:  Factor:  $F_{kk} = L_{kk} U_{kk}$ 
15:  for  $i = k + 1$  to  $p$  do
16:    Solve ( $L$ ):  $\tilde{F}_{ik} \leftarrow \tilde{F}_{ik} U_{kk}^{-1} = X_{ik} Y_{ik}^T U_{kk}^{-1}$ 
17:    Solve ( $U$ ):  $\tilde{F}_{ki} \leftarrow L_{kk}^{-1} \tilde{F}_{ki} = L_{kk}^{-1} X_{ki} Y_{ki}^T$ 
18:  end for
19: end for

```

accumulated together before being recompressed, in order to achieve the smallest rank possible for \tilde{F}_{ik} .

The CUFS BLR variant is referred to as *fully-structured*, which means the off-diagonal low-rank blocks \tilde{F}_{ik} are never stored in full-rank again after being initially compressed. Furthermore, in the rest of this article, we will assume that the matrix is already available in compressed form, that is, that the cost of the Compress step is negligible with respect to the total complexity of the factorization. This is for example the case when the blocks of the original matrix F are sparse, since their low-rank representation can be computed in only $O(br)$ flops.

One of the main results of [4] is that the storage complexity of the factorization of a dense matrix of order m with off-diagonal blocks of rank at most r is equal to

$$\mathcal{S}_{ds}^1(m, r) = O(m^{1.5} \sqrt{r}) \quad (2.2)$$

and the flop complexity is

$$\mathcal{F}_{ds}^1(m, r) = O(m^2 r). \quad (2.3)$$

The proof of this result will be recalled in Section 2.4. The 1 superscript refers to the monolevel BLR factorization. This notation will be generalized in the next sections.

2.3. BLR sparse factorization. Because sparse direct factorizations such as multifrontal or supernodal approaches rely on dense factorizations, block low-rank approximations can easily be incorporated into the sparse factorization by representing the fronts (or supernodes) with the chosen low-rank format. For example, in the BLR case, the fronts are represented as defined by (2.1), and Algorithm 2.1 is adapted to perform their partial factorization. This is described in detail in [2, 25].

As a consequence, the complexity of the sparse factorization can be directly computed from the complexity of the dense factorization. We consider a matrix of order $n = N^d$, where d denotes the problem dimension, that is reordered using nested dissection [16]: the domain is recursively partitioned by so-called *separators*. The sparse complexity is then computed as follows (assuming cross-shaped separators): at each level ℓ of the separator tree, we need to factorize $(2^d)^\ell$ fronts of order $O((N/2^\ell)^{d-1})$, for ℓ ranging from 0 to $L = \log_2(N)$. Therefore, the flop complexity $\mathcal{F}_{sp}(N)$ is

Table 2.1: Flop and storage complexities of the factorization of a sparse system of $n = N \times N$ (2D case) or $n = N \times N \times N$ (3D case) unknowns, assuming a dense factorization complexity $O(m^\beta)$.

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

Table 2.2: Flop and storage complexities of the FR, BLR, and \mathcal{H} factorizations of a sparse system of $n = N \times N$ (2D case) or $n = N \times N \times N$ (3D case) unknowns, derived from the complexities of the factorization of a dense matrix of order m and with a rank bound r . The BLR variant considered is CUFS.

	$\mathcal{F}_{ds}(m, r)$	$\mathcal{F}_{sp}(n, r)$		$\mathcal{S}_{ds}(m, r)$	$\mathcal{S}_{sp}(n, r)$	
		2D	3D		2D	3D
FR	$O(m^3)$	$O(n^{3/2})$	$O(n^2)$	$O(m^2)$	$O(n \log n)$	$O(n^{4/3})$
BLR	$O(m^2 r)$	$O(n \max(r, \log n))$	$O(n^{4/3} r)$	$O(m^{3/2} r^{1/2})$	$O(n)$	$O(n \max(r^{1/2}, \log n))$
\mathcal{H}	$O(m r^2 \log^2 m)$	$O(\max(n, n^{1/2} r^2))$	$O(\max(n, n^{2/3} r^2))$	$O(m r \log m)$	$O(n)$	$O(\max(n, n^{2/3} r))$

equal to

$$\mathcal{F}_{sp}(N) = \sum_{\ell=0}^L (2^d)^\ell \mathcal{F}_{ds}\left(\left(\frac{N}{2^\ell}\right)^{d-1}\right), \quad (2.4)$$

where $\mathcal{F}_{ds}(m)$ is the dense flop complexity. In BLR, it is given by (2.3). Similarly, the storage complexity $\mathcal{S}_{sp}(N)$ is equal to

$$\mathcal{S}_{sp}(N) = \sum_{\ell=0}^L (2^d)^\ell \mathcal{S}_{ds}\left(\left(\frac{N}{2^\ell}\right)^{d-1}\right), \quad (2.5)$$

where $\mathcal{S}_{ds}(m)$ is the dense storage complexity. In BLR, it is given by (2.2).

Assuming a dense complexity $O(m^\beta)$, it can easily be shown from (2.4) and (2.5) that the sparse complexities only depend on the dense complexity exponent β and the dimension d . This correspondence is reported in Table 2.1. The key observation is that *a linear $O(m)$ dense complexity is not required to achieve a linear $O(n)$ sparse complexity*. In fact, a dense complexity lower than $O(m^2)$ and $O(m^{1.5})$ suffices for 2D and 3D problems, respectively.

Then, the sparse complexities are reported in Table 2.2, for the FR, BLR, and \mathcal{H} factorizations. Thanks to the key observation above, the BLR sparse complexities are not that far from the \mathcal{H} complexities. For example, the BLR 2D storage complexity is already optimal; furthermore, the 2D flop and 3D storage complexities are nearly linear, with only an additional factor depending only on r and $\log n$ compared with \mathcal{H} . More importantly, thanks to the same key observation, we only need a modest improvement of the dense complexity to reach $O(n)$ complexity. Specifically:

- To drop the $\max(r, \log n)$ factor in the 2D flop complexity, the dense complexity $\mathcal{F}_{ds}(m)$ only needs to be strictly inferior to $O(m^2)$;
- Similarly, to drop the $\max(r^{1/2}, \log n)$ factor in the 3D storage complexity, the dense complexity $\mathcal{S}_{ds}(m)$ only needs to be strictly inferior to $O(m^{1.5})$;
- Finally, the superlinear 3D flop complexity can be made linear with a dense complexity $\mathcal{F}_{ds}(m)$ strictly inferior to $O(m^{1.5})$.

The main motivation behind the MBLR format is to find a *simple* modification of the BLR factorization that preserves its simplicity and flexibility, while improving the complexity *just enough* to get the desired exponent. We will prove in Section 5 that this complexity improvement can be controlled by the number of levels.

Table 2.3: Main operations for the BLR factorization of a dense matrix of order m , with blocks of size b , and low-rank blocks of rank at most r . We note $p = m/b$. Type: type of the block(s) on which the operation is performed. $cost_{Step}^1$: cost of performing the operation once. $number_{Step}$: number of times the operation is performed. \mathcal{F}_{Step}^1 : obtained by multiplying the $cost_{Step}^1$ and $number_{Step}$ columns (equation (2.11)). The first expression is given as a function of b , p , and r , while the second is obtained with the assumption that $b = O(m^x)$ (and thus $p = O(m^{1-x})$) and $r = O(m^\alpha)$, for $x, \alpha \in [0, 1]$.

Step	Type	$cost_{Step}^1$	$number_{Step}$	$\mathcal{F}_{Step}^1(b, p, r)$	=	$\mathcal{F}_{Step}^1(m, x, \alpha)$
Factor	FR	$O(b^3)$	$O(p)$	$O(pb^3)$	=	$O(m^{1+2x})$
Solve	FR-LR	$O(b^2r)$	$O(p^2)$	$O(p^2b^2r)$	=	$O(m^{2+\alpha})$
Product	LR-LR	$O(br^2)$	$O(p^3)$	$O(p^3br^2)$	=	$O(m^{3-2x+2\alpha})$
Recompress	LR	$O(bpr^2)$	$O(p^2)$	$O(p^3br^2)$	=	$O(m^{3-2x+2\alpha})$

2.4. Complexity analysis for BLR. To motivate the key idea behind MBLR, let us summarize the dense storage and flop complexity analysis found in [4] that leads to the formulas (2.2) and (2.3). We consider the CUFS variant of the BLR factorization. As indicated in the introduction, we first consider the weakly-admissible case: we assume that all off-diagonal blocks are low-rank. The extension to the strongly-admissible case is provided in the appendix.

We consider a dense BLR matrix of order m . We note b the block size and $p = m/b$ the number of blocks per row and/or column. The amount of storage required to store the factors of such a matrix can be computed as the sum of the storage for the full-rank diagonal blocks and that for the low-rank off-diagonal blocks:

$$\mathcal{S}_{ds}^1(b, p, r) = O(pb^2) + O(p^2br). \quad (2.6)$$

Then, we assume that the block size b is of order $O(m^x)$, where x is a real value in $[0, 1]$, and thus the number of blocks p per row and/or column is of order $O(m^{1-x})$. We also assume that the rank bound is of the form $r = O(m^\alpha)$. By replacing b , p , and r by their expression in (2.6), we obtain an expression of \mathcal{S}_{ds}^1 which depends on (m, x, α) instead of (b, p, r) :

$$\mathcal{S}_{ds}^1(m, x, \alpha) = O(m^{1+x}) + O(m^{2-x+\alpha}). \quad (2.7)$$

We then define x^* as the optimal choice of x which minimizes the asymptotic complexity of (2.7). x^* can be computed as the value which makes each term in (2.7) asymptotically equal. We obtain

$$x^* = (1 + \alpha)/2, \quad (2.8)$$

which means the optimal choice of block size is

$$b^* = O(\sqrt{mr}). \quad (2.9)$$

This leads to the final dense complexity (already given in (2.2))

$$\mathcal{S}_{ds}^1(m, r) = O(m^{1.5}\sqrt{r}). \quad (2.10)$$

Next, to compute the flop complexity, we compute the cost of the Factor, Solve, Product, and Recompress steps and report them in Table 2.3 (third column). This cost depends on the type (full-rank or low-rank) of the block(s) on which the operation is performed (second column). Note that the Product operation can only take the form of a product of two low-rank blocks (LR-LR), because it involves only off-diagonal blocks, which are low-rank in the weakly-admissible case. We also recall that we have assumed that the matrix is already available in compressed form and thus we do not report the Compress step in Table 2.3.

In the weakly-admissible case, there is only one full-rank block on each block-row (the diagonal one); therefore, we can easily count the number of blocks on which each step is performed; we report

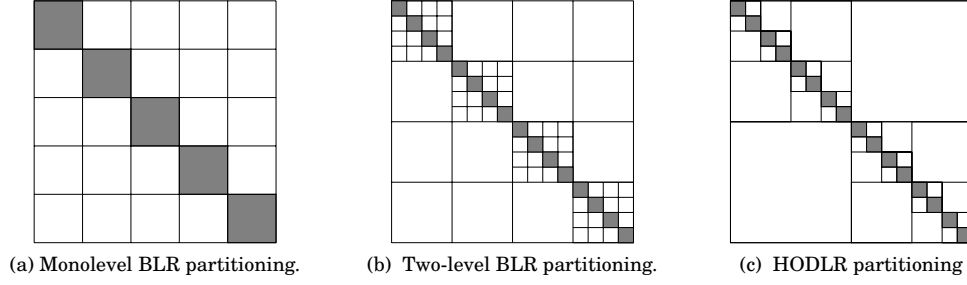


Fig. 2.1: Comparison of BLR, MBLR, and hierarchical formats in the weakly-admissible case.

it on the fourth column of Table 2.3. The BLR factorization cost of each step is then equal to

$$\mathcal{F}_{Step}^1 = cost_{Step}^1 \times number_{Step} \quad (2.11)$$

and is reported in the fifth and sixth columns of Table 2.3. In the fifth column, its expression depends on b , p , and r , while in the sixth column it is given as a function of m , x , and α by substituting b , p , and r by $O(m^x)$, $O(m^{1-x})$, and $O(m^\alpha)$, respectively.

The total flop complexity of the dense BLR factorization is equal to the sum of the cost of all steps

$$\mathcal{F}_{ds}^1(m, x, \alpha) = O(m^{1+2x} + m^{2+\alpha} + m^{3-2x+2\alpha}). \quad (2.12)$$

We compute x^* , the optimal choice of x which minimizes the complexity, and find again $x^* = (1 + \alpha)/2$, which means that the same x^* value minimizes both the storage and flop complexities, a valuable property. We finally obtain

$$\mathcal{F}_{ds}^1(m, r) = O(m^2 r). \quad (2.13)$$

2.5. Key idea of the MBLR format. Let us now consider each step of Table 2.3 with the objective of reducing the total cost of the factorization. The Product and Recompress steps involve exclusively low-rank blocks and their cost is already optimal as it is linear with respect to the block size b . Therefore, we focus on the Factor and Solve steps. These steps have superlinear cost with respect to b because they involve the diagonal full-rank blocks.

Thus, the key idea of the MBLR format is to further refine these full-rank diagonal blocks by replacing them by BLR matrices. This is illustrated in Figure 2.1b. Compared with the BLR format (Figure 2.1a), we will show in Section 4 that this more compressed representation decreases the cost of performing the Factor and Solve steps, which allows to increase the block size to achieve a lower complexity. However, it also remains very different from a hierarchical format (Figure 2.1c). First, the diagonal blocks are represented by the simple BLR format, rather than a more complex hierarchical format. Second, while larger than in the BLR case, the off-diagonal blocks are in general still much smaller than in the \mathcal{H} case. This has several advantages:

- No relative order is needed between blocks; this allows the clustering to easily be computed and delivers a great flexibility to distribute the data in a parallel environment.
- The size of the blocks can be small enough to fit on a single shared-memory node; therefore, in a parallel environment, each processor can efficiently and concurrently work on different blocks.
- To perform numerical pivoting, the quality of the pivot candidates in the off-diagonal blocks F_{ik} can be estimated with the entries of its low-rank basis Y_{ik} , provided that X_{ik} is an orthonormal matrix. However, as explained in [25], the quality of this estimation depends on the size of the block; smaller blocks lead to a tighter estimate. This makes the BLR format particularly suitable to handle numerical pivoting, a critical feature lacking in most hierarchical solvers presented in the literature.

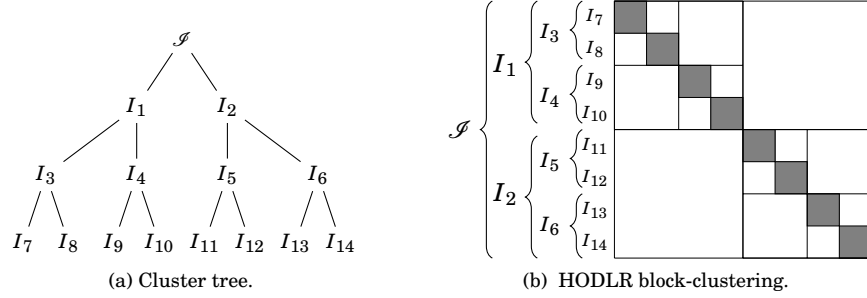


Fig. 3.1: An example of cluster tree and its associated HODLR block-clustering.

While many advantages of the BLR and MBLR formats lie in their efficiency and flexibility in the context of a parallel execution, the parallel implementation of the MBLR format is out of the scope of this paper. Similarly, we omit a detailed description of the algorithms designed to handle numerical pivoting; see [25] for a thorough discussion. In this paper, we focus on the theoretical complexity analysis of the MBLR factorization.

3. Difference between MBLR and \mathcal{H} matrices. In this section, we first explain how MBLR matrices can be represented using the cluster tree modelization typically used in the \mathcal{H} literature; this provides a convenient tool to formalize the key difference between the MBLR and hierarchical formats that was informally presented in the previous section: the number of levels in the cluster tree is a constant in the MBLR format, while it is logarithmically dependent on the problem size in the \mathcal{H} format.

Using this formalism, \mathcal{H} theory is thus applicable to the MBLR format; however, we show in Section 3.3 that it leads to very pessimistic complexity bounds. We must therefore develop a new theory to compute satisfying bounds, which is the object of Sections 4 and 5. Since the proofs and computations in these sections are *not* based on the \mathcal{H} theory, this section may be skipped by the reader, at least on a first read.

3.1. The hierarchical case. Here, we briefly recall the definition of cluster trees, and how they are used to represent hierarchical partitionings. We refer to [10, 22] for a formal and detailed presentation.

Let us note \mathcal{I} the set of unknowns. We assume that the sets of row and column indices of the matrix are the same, for the sake of simplicity, and because we do not need to distinguish them for the purpose of this section.

Computing a recursive partition $\mathfrak{S}(\mathcal{I})$ of \mathcal{I} can be modeled with a so-called *cluster tree*.

DEFINITION 3.1 (Cluster tree). *Let \mathcal{I} be a set of unknowns and $T_{\mathcal{I}}$ a tree whose nodes v are associated with subsets σ_v of \mathcal{I} . $T_{\mathcal{I}}$ is said to be a cluster tree iff:*

- *The root of $T_{\mathcal{I}}$, noted r , is associated with $\sigma_r = \mathcal{I}$;*
- *For each non-leaf node $v \in T_{\mathcal{I}}$, with children noted $C_{T_{\mathcal{I}}}(v)$, the subsets associated with the children form a partition of σ_v , that is, the σ_c subsets are disjoint and satisfy*

$$\bigcup_{c \in C_{T_{\mathcal{I}}}(v)} \sigma_c = \sigma_v.$$

An example of cluster tree is provided in Figure 3.1a. As illustrated, cluster trees establish a *hierarchy* between clusters. A given cluster tree uniquely defines a weakly-admissible, hierarchical block-clustering (so-called HODLR matrix), as illustrated in Figure 3.1b.

Note that the cluster tree is not enough to define general, strongly-admissible hierarchical block-clusterings such as an \mathcal{H} block-clustering. In that case, a new tree structure, so-called block-cluster tree, must be introduced. For the sake of simplicity, we do not discuss them here, since cluster trees (and weakly-admissible hierarchical block-clusterings) are enough to explain the difference between MBLR and hierarchical matrices.

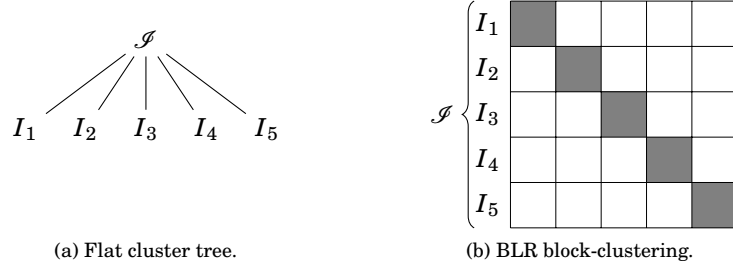


Fig. 3.2: An example of flat cluster tree and its associated BLR block-clustering.

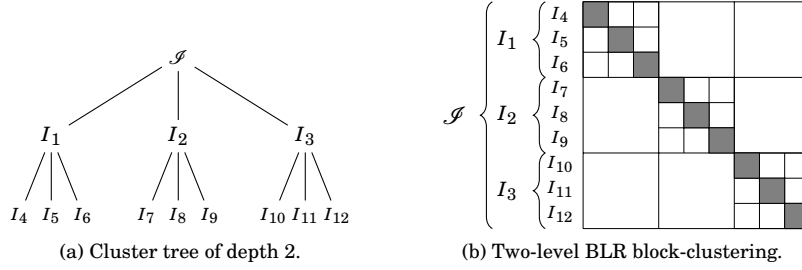


Fig. 3.3: An example of cluster tree of depth 2 and its associated two-level BLR block-clustering.

3.2. The BLR and MBLR formats, explained with cluster trees. It is interesting to first mention how the BLR format can be viewed as a very particular kind of \mathcal{H} -matrix. In [4], we made a first attempt to model BLR partitionings with cluster trees, where the BLR partitioning is defined using only the leaves of the cluster tree. However, this model is inadequate because the other nodes of the tree are unnecessary (and thus there are several possible cluster trees for a unique BLR partitioning).

Therefore, we newly propose to define a given BLR partitioning with a unique, flat cluster tree of depth 1, as illustrated in Figure 3.2.

With this definition, it is straightforward to extend the model to the MBLR format. An ℓ -level BLR block-clustering can be represented by a cluster tree of depth ℓ . This is illustrated in Figure 3.3 in the two-level case.

It is thus clear that, just like the BLR format, the MBLR one can also be viewed as a very particular kind of \mathcal{H} -matrix, since any MBLR matrix can be represented with a cluster tree that satisfies Definition 3.1. However, while \mathcal{H} -matrices are indeed also represented by cluster trees, in practice, they are virtually always built with an implicit assumption: the number of children of any node in the cluster tree is constant with respect to the size of the matrix. Most commonly, cluster trees are binary trees, and thus this number is 2. Since the diagonal blocks of an \mathcal{H} -matrix are of constant size, this leads to $O(\log m)$ levels in the cluster tree.

The MBLR format is based on the converse approach: the number of levels is fixed to some constant $\ell = O(1)$, which in turn leads to a number of children per node that is asymptotically dependent on m . For example, assuming a constant rank $r = O(1)$, we will prove in Section 4 that a two-level BLR matrix is represented by a cluster tree with $O(m^{1/3})$ nodes on the first level (children of the root node), and each of these nodes has $O((m^{2/3})^{1/2}) = O(m^{1/3})$ children on the second level.

While this is technically not prohibited by the general \mathcal{H} -format definition, it has, to the best of our knowledge, never been considered in the \mathcal{H} literature. As a matter of fact, in his recent book, Hackbusch writes ([22], p.83): “The partition should contain as few blocks as possible since the storage cost increases with the number of blocks”. While that is true, we believe that the smaller size and greater number of blocks can provide a gain in flexibility that can be useful in a

parallel solver.

One may in fact find it surprising that this simple idea of having a constant number of levels has never been proposed before. We believe that this is mainly explained by the fact that the \mathcal{H} theoretical formalism does not provide a satisfying result if applied to the MBLR format, as explained in the following.

3.3. \mathcal{H} theory leads to pessimistic MBLR complexity bounds. In [4, Section 3.2], we explained that applying the \mathcal{H} theoretical complexity formulas to the BLR format does not provide a satisfying complexity bound. Here, we show this remains the case for the MBLR format.

The storage and flop complexities of the factorization of a dense \mathcal{H} -matrix of order m have been shown [20] to be

$$\mathcal{S}_{ds}^{\mathcal{H}}(m) = O(mLc_{sp} \max(r_{max}, b_{diag})), \quad (3.1)$$

$$\mathcal{F}_{ds}^{\mathcal{H}}(m) = O(mL^2 c_{sp}^2 \max(r_{max}, b_{diag})^2). \quad (3.2)$$

L is the depth of the cluster tree. c_{sp} is the so-called sparsity constant, defined as the maximum number of blocks of a given level in the cluster tree that are in the same row or column of the matrix; it is a measure of how much the original matrix has been refined. r_{max} is the maximal rank of all the blocks of the matrix, while b_{diag} is the size of the diagonal blocks. Note that the power to which c_{sp} is elevated to in the flop complexity depends on the reference considered. The term c_{sp}^3 is found in [20]; a reduced complexity with c_{sp}^2 instead is given in [11]. Since our aim is to show that the \mathcal{H} complexity formulas do not provide a satisfying result when applied to MBLR matrices, in the following, we use the best case formula with c_{sp}^2 ; the formula with c_{sp}^3 would lead to an even more pessimistic complexity bound.

Under some assumptions on how the partition $\mathfrak{S}(\mathcal{T})$ is built [20, Lemma 4.5], the sparsity constant can be bounded by $O(1)$ in the \mathcal{H} case. By recursively refining the diagonal blocks until they become of constant size, we have $b_{diag} = O(1)$ and $L = O(\log m)$. Therefore, (3.1) and (3.2) lead to $\mathcal{S}_{ds}^{\mathcal{H}}(m) = O(rm \log m)$ and $\mathcal{F}_{ds}^{\mathcal{H}}(m) = O(r^2 m \log^2 m)$.

In the BLR case, we have $L = 1$ and $b_{diag} = b$; the sparsity constant is equal to the number of blocks per row m/b , a high value that translates the fact that BLR matrices are much more refined than hierarchical ones. This leads to $\mathcal{S}_{ds}^{\mathcal{H},1}(m) = O(m^2)$ and $\mathcal{F}_{ds}^{\mathcal{H},1}(m) = O(m^3)$, where the notation $\mathcal{S}_{ds}^{\mathcal{H},1}$ and $\mathcal{F}_{ds}^{\mathcal{H},1}$ signifies “ \mathcal{H} complexity formula applied to the monolevel BLR case”. We thus obtain a very pessimistic complexity bound, which is due to the fact that the diagonal blocks are of the same size as all the other blocks, i.e., $b_{diag} = b$.

Applying the formulas to the MBLR case leads to the same problem. Assuming that $r = O(1)$ for the sake of simplicity, let us consider the two-level case and denote by b_1 and b_2 the outer and inner blocks sizes, respectively. We have $L = 2$, $b_{diag} = b_2$, and $c_{sp} = \max(m/b_1, b_1/b_2)$. The storage complexity

$$\mathcal{S}_{ds}^{\mathcal{H},2}(m) = O(\max(m^2 b_2/b_1, m b_1))$$

is minimized for $b_1 = O(\sqrt{m})$ and $b_2 = O(1)$, and equal to $O(m^{3/2})$, which is not a satisfying result since we have proven in Section 2.4 that this is the BLR complexity. A similarly unsatisfying result is obtained for the flop complexity, again minimized for $b_1 = O(\sqrt{m})$ and $b_2 = O(1)$ and equal to

$$\mathcal{F}_{ds}^{\mathcal{H},2}(m) = O(\max(m^3 b_2^2/b_1^2, m b_1^2)) = O(m^2).$$

Does the problem persist with a higher number of levels ℓ ? The answer is unfortunately positive. Indeed, we have

$$\mathcal{S}_{ds}^{\mathcal{H},\ell}(m) = O\left(\max_{i=0,\ell-1} \left(\frac{m b_i b_\ell}{b_{i+1}}\right)\right),$$

with the notation $b_0 = m$. To minimize the expression, we equilibrate each term in the maximum, which leads to the recursive relation $b_i^2 = b_{i-1} b_{i+1}$. Noting $b_1 = b$ the outer block size, it is clear that the closed-form expression of b_i is $b_i = b^i / m^{i-1}$. Since the smallest block size b_ℓ must be at least $O(1)$, this leads to

$$\frac{b^\ell}{m^{\ell-1}} \geq O(1)$$

Table 3.1: Applying the \mathcal{H} theoretical complexity formulas to the BLR or MBLR cases does not provide a satisfying result. We have assumed that $r = O(1)$ for the sake of clarity, but the bounds would remain pessimistic for general ranks.

	\mathcal{H}	BLR	MBLR
L	$O(\log m)$	1	ℓ
c_{sp}	$O(1)$	m/b	$\max_{i=1,\ell} \frac{b_{i-1}}{b_i}^*$
b_{diag}	$O(1)$	b	b_ℓ
$\mathcal{S}_{ds}^{\mathcal{H}}$	$O(m \log m)$	$O(m^2)$	$O(m^{(\ell+1)/\ell})$
$\mathcal{F}_{ds}^{\mathcal{H}}$	$O(m \log^2 m)$	$O(m^3)$	$O(m^{(\ell+2)/\ell})$

*with the notation $b_0 = m$

and thus $b \geq O(m^{(\ell-1)/\ell})$, which in turn leads to the complexity bound

$$\mathcal{S}_{ds}^{\mathcal{H},\ell}(m) \geq O\left(\frac{m^2 b_\ell}{b_1}\right) \geq O(m^{(\ell+1)/\ell}).$$

In Section 5, we will prove that for $r = O(1)$ we have $\mathcal{S}_{ds}^\ell(m) = O(m^{(\ell+2)/(\ell+1)})$, which is a better bound, especially when considering a small number of levels ℓ . The analysis for the flop complexity leads to the same expression of b_i and to the complexity bound $\mathcal{F}_{ds}^{\mathcal{H},\ell}(m) = O(m^{(\ell+2)/\ell})$, which is again overly pessimistic compared with the bound $\mathcal{F}_{ds}^\ell(m) = O(m^{(\ell+3)/(\ell+1)})$ that we will prove in Section 5.

We summarize the result of applying the \mathcal{H} complexity formulas in Table 3.1. It is interesting to note that these bounds overestimate the actual complexity by *exactly one level*. This is actually not surprising: in the \mathcal{H} theory, all the blocks at the bottom level are assumed to be non-admissible, which accounts for the “loss” of one level in the complexity bounds.

It is therefore clear that we must develop a new theory extending the \mathcal{H} formalism to be able to prove better MBLR complexity bounds, just as we did in [4] for the BLR case. We begin by the two-level case in the next section.

4. Two-level BLR matrices. In this section, we compute the theoretical complexity of the two-level BLR factorization. The proofs and computations on this particular two-level case are meant to be illustrative of those of the general multilevel case with an arbitrary number of levels, which is discussed in Section 5.

We remind the reader that in this section, we are considering the weakly-admissible case only.

4.1. Two-level kernels description. In order to adapt Algorithm 2.1 to two-level BLR matrices, two modifications must be performed.

First, the Factor kernel is not a full-rank factorization of the diagonal blocks F_{kk} anymore, but a BLR factorization. Thus, line 14 must be replaced by

$$\tilde{F}_{kk} = \tilde{L}_{kk} \tilde{U}_{kk} = \text{BLR-Factor}(F_{kk}), \quad (4.1)$$

where BLR-Factor refers to the BLR factorization described in Algorithm 2.1.

Second, the FR-LR Solve kernel at lines 16 and 17 must be changed to a BLR-LR Solve:

$$\tilde{F}_{ik} \leftarrow \text{BLR-LR-Solve}(\tilde{U}_{kk}, \tilde{F}_{ik}) \quad (4.2)$$

$$\tilde{F}_{ki} \leftarrow \text{BLR-LR-Solve}(\tilde{L}_{kk}, \tilde{F}_{ki}), \quad (4.3)$$

where \tilde{F}_{ik} and \tilde{F}_{ki} are LR matrices and \tilde{L}_{kk} and \tilde{U}_{kk} are lower and upper triangular BLR matrices.

We describe in Algorithm 4.1 the BLR-LR-Solve kernel, in the upper triangular case, omitting the lower triangular case which is very similar. The kernel consists in applying a triangular solve with a upper triangular BLR matrix \tilde{U} to a low-rank matrix $\tilde{B} = \Phi \Psi^T$. We recall that \tilde{U}_{ij} denotes

Algorithm 4.1 BLR-LR-Solve kernel (upper triangular case)

Input: a $p \times p$ block upper triangular BLR matrix \tilde{U} ; $\tilde{U} = [\tilde{U}_{ij}]_{i=1:j,j=1:p}$ such that $\tilde{U}_{ij} = Y_{ij}X_{ij}^T$ for $i \neq j$ and $\tilde{U}_{jj} = U_{jj}$ is a full-rank matrix; a LR matrix $\tilde{B} = \Phi\Psi^T$.

Output: overwritten Ψ (modified in-place) corresponding to the operation $\tilde{B} \leftarrow \tilde{B}\tilde{U}^{-1}$.

```

1: for  $j = 1$  to  $p$  do
2:    $\Psi_{j,:}^T \leftarrow \Psi_{j,:}^T - \sum_{i=1}^{j-1} \Psi_{i,:}^T Y_{ij} X_{ij}^T$ 
3:    $\Psi_{j,:}^T \leftarrow \Psi_{j,:}^T U_{jj}^{-1}$ 
4: end for

```

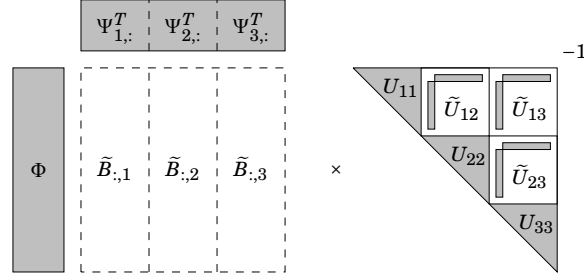


Fig. 4.1: The BLR-LR-Solve kernel

the (i, j) -th low-rank sub-block of \tilde{U} (with the notation $\tilde{U}_{jj} = U_{jj}$), and that $\Psi_{j,:}$ denotes the j -th block-row of Ψ .

Two main operations must be performed: a triangular solve using the full-rank diagonal blocks U_{jj} of \tilde{U} , and an update using the low-rank off-diagonal blocks $\tilde{U}_{ij} = Y_{ij}X_{ij}^T$ of \tilde{U} . Both are applied on the low-rank block-columns $\tilde{B}_{:,j} = \Phi\Psi_{j,:}^T$ of \tilde{B} . These two operations take place at lines 3 and 2 of Algorithm 4.1, respectively.

The FR-LR triangular solve can be written as

$$\tilde{B}_{:,j} \leftarrow \tilde{B}_{:,j} U_{jj}^{-1} = \Phi \left(\Psi_{j,:}^T U_{jj}^{-1} \right) \quad (4.4)$$

and thus only $\Psi_{j,:}^T$ needs to be updated, as shown at line 3.

The LR-LR update takes the following form:

$$\begin{aligned} \tilde{B}_{:,j} \leftarrow \tilde{B}_{:,j} - \sum_{i=1}^{j-1} \tilde{B}_{:,i} \tilde{U}_{ij} &= \Phi \Psi_{j,:}^T - \sum_{i=1}^{j-1} \Phi \left(\Psi_{i,:}^T Y_{ij} X_{ij}^T \right) \\ &= \Phi \left(\Psi_{j,:}^T - \sum_{i=1}^{j-1} \Psi_{i,:}^T Y_{ij} X_{ij}^T \right) \end{aligned} \quad (4.5)$$

and thus, again, only $\Psi_{j,:}^T$ needs to be updated, as shown at line 2.

The BLR-LR-Solve kernel is illustrated in Figure 4.1.

It can easily be computed that the cost of applying the BLR-LR-Solve kernel once is equal to the storage complexity times the rank bound r :

$$\text{cost}_{\text{Solve}}^2 = O(r) \times \mathcal{S}_{ds}^1(b, r). \quad (4.6)$$

Injecting (2.10) into (4.6) leads to

$$\text{cost}_{\text{Solve}}^2 = O(b^{3/2} r^{3/2}). \quad (4.7)$$

4.2. Two-level complexity analysis. We now compute the complexity of the two-level BLR factorization of a dense matrix F of order m . We reuse the monolevel notations for the top level. We denote by b the size of the first-level blocks and $p = m/b$ the number of block-rows and block-columns. We assume that b is of the form $b = m^x$, for $x \in [0, 1]$. We also assume that the ranks of the off-diagonal blocks are bounded by $r = O(m^a)$.

Table 4.1: Two-level equivalent of Table 2.3. The legend of the table is the same. The differences between the two tables are highlighted in gray.

Step	Type	$cost_{Step}^2$	$number_{Step}$	$\mathcal{F}_{Step}^2(b, p, r)$	$=$	$\mathcal{F}_{Step}^2(m, x, \alpha)$
Factor	BLR	$O(b^2 r)$	$O(p)$	$O(p b^2 r)$	$=$	$O(m^{1+x+\alpha})$
Solve	BLR-LR	$O(b^{3/2} r^{3/2})$	$O(p^2)$	$O(p^2 b^{3/2} r^{3/2})$	$=$	$O(m^{2-x/2+3\alpha/2})$
Product	LR-LR	$O(b r^2)$	$O(p^3)$	$O(p^3 b r^2)$	$=$	$O(m^{3-2x+2\alpha})$
Recompress	LR	$O(b p r^2)$	$O(p^2)$	$O(p^3 b r^2)$	$=$	$O(m^{3-2x+2\alpha})$

The size required to store the factors is again the sum of the storage for the diagonal and off-diagonal blocks, as in (2.6). The difference is that the diagonal blocks are not full-rank but BLR matrices. They are further refined into smaller blocks whose size should be chosen of order $O(\sqrt{br})$, as determined by (2.9) in the BLR complexity analysis. Therefore, in the two-level case, (2.6) becomes

$$\mathcal{S}_{ds}^2(b, p, r) = p \times \mathcal{S}_{ds}^1(b, r) + O(p^2 b r). \quad (4.8)$$

By replacing $\mathcal{S}_{ds}^1(b, r)$ by its second expression computed in (2.10), we obtain

$$\mathcal{S}_{ds}^2(b, p, r) = O(p b^{3/2} r^{1/2}) + O(p^2 b r) \quad (4.9)$$

Finally, we replace b , p , and r by their expression $O(m^x)$, $O(m^{1-x})$, and $O(m^\alpha)$, respectively, to obtain

$$\mathcal{S}_{ds}^2(m, x, \alpha) = O(m^{1+(x+\alpha)/2} + m^{2-x+\alpha}). \quad (4.10)$$

This leads to

$$x^* = (2 + \alpha)/3, \quad (4.11)$$

where x^* defines the optimal choice of the first-level block size b . We thus obtain a final two-level storage complexity of

$$\mathcal{S}_{ds}^2(m, r) = O(m^{4/3} r^{2/3}). \quad (4.12)$$

We now compute the flop complexity $\mathcal{F}_{ds}^2(m, r)$ of the two-level BLR dense factorization. We compute the cost of each step and report it in Table 4.1, which is the two-level equivalent of Table 2.3; the differences between the two tables are highlighted in gray. The Product and Recompress steps have not changed and thus have the same cost. The Factor step consists in factorizing the p diagonal blocks which are now represented by BLR matrices; thus its cost is directly derived from the BLR complexity computed in (2.13):

$$\mathcal{F}_{Factor}^2(b, p, r) = p \times \mathcal{F}_{ds}^1(b, r) = O(p b^2 r). \quad (4.13)$$

It remains to compute the cost of the Solve step, which now takes the form of $O(p^2)$ calls to the BLR-LR-Solve kernel whose cost is given by (4.7). Thus, the overall cost of the Solve step is

$$\mathcal{F}_{Solve}^2(b, p, r) = O(p^2 b^{3/2} r^{3/2}). \quad (4.14)$$

This concludes the computations for the third, fourth, and fifth columns of Table 4.1. Just as in the monolevel case, the sixth column is obtained by replacing b , p , and r by their expression $O(m^x)$, $O(m^{1-x})$, and $O(m^\alpha)$, respectively.

We can finally compute the total flop complexity as the sum of the costs of all steps

$$\mathcal{F}_{ds}^2(m, x, \alpha) = O(m^{1+x+\alpha} + m^{2-x/2+3\alpha/2} + m^{3-2x+2\alpha}) \quad (4.15)$$

We then compute x^* which is again equal to the same value as the one that minimizes the storage complexity, $x^* = (2 + \alpha)/3$. Therefore, the final two-level dense flop complexity is

$$\mathcal{F}_{ds}^2(m, r) = O(m^{5/3} r^{4/3}). \quad (4.16)$$

Table 4.2: Flop and storage complexities of the monolevel and two-level BLR factorizations of a sparse system of $n = N \times N$ (2D case) or $n = N \times N \times N$ (3D case) unknowns, derived from the complexities of the factorization of a dense matrix of order m . We consider a constant rank bound $r = O(1)$. The BLR variant considered is CUFS.

	$\mathcal{F}_{ds}(m)$	$\mathcal{F}_{sp}(n)$		$\mathcal{S}_{ds}(m)$	$\mathcal{S}_{sp}(n)$	
		2D	3D		2D	3D
Monolevel BLR	$O(m^2)$	$O(n \log n)$	$O(n^{4/3})$	$O(m^{1.5})$	$O(n)$	$O(n \log n)$
Two-level BLR	$O(m^{5/3})$	$O(n)$	$O(n^{10/9})$	$O(m^{4/3})$	$O(n)$	$O(n)$

The two-level BLR format therefore significantly improves the asymptotic storage and flop complexity compared with the monolevel format. Our analysis shows that the top level block size should be set to $b^* = O(m^{x^*}) = O(m^{2/3}r^{1/3})$. This is an asymptotically larger value than the monolevel optimal block size computed in (2.9), which translates the fact that by refining the diagonal blocks, we can afford to take a larger block size to improve the overall asymptotic complexity. However, contrarily to hierarchical matrices, b^* remains asymptotically much lower than $O(m)$; this makes the format much more flexible for the reasons described in Section 2.5. In short, *we have traded off some of the flexibility of the monolevel format to improve the asymptotic complexity.*

This improvement of the dense flop and storage complexities is translated into an improvement of the sparse complexities. Assuming a rank bound in $O(1)$, we quantify this improvement in Table 4.2. Compared with the monolevel BLR format, the two-level BLR format drops the $O(\log n)$ factor in the 2D flop and 3D storage complexities, which become linear and thus optimal. The two-level BLR format can thus achieve, in these two cases, the same $O(n)$ complexity as the hierarchical formats while being almost as simple and flexible as flat formats.

Finally, the 3D flop complexity remains superlinear but is significantly reduced, from $O(n^{4/3})$ to $O(n^{10/9})$. As the problem size gets larger and larger, even small asymptotic improvements can make a big difference. Therefore, we now generalize the two-level analysis to the multilevel case with an arbitrary number of levels.

5. Generalization to ℓ -level BLR matrices. In this section, we generalize two-level proof and computations of the previous section to an arbitrary number of levels ℓ by computing recursive complexity formulas. We remind the reader that in this section, we are still considering the weakly-admissible case only.

5.1. Recursive complexity analysis. Just as for the two-level case, one can compute the three-level asymptotic complexities, and so on until the general formula becomes clear. We state the result for an arbitrary number of levels ℓ in the following theorem. Note that it is important to assume that the number of levels ℓ is constant ($\ell = O(1)$), since the constants hidden in the big O depend on ℓ .

THEOREM 5.1 (Storage and flop complexity of the ℓ -level BLR factorization). *Let us consider a dense ℓ -level BLR matrix of order m . We note b the size of the top level blocks, and $p = m/b$. Let $r = O(m^\alpha)$ be the bound on the maximal rank of any block on any level. Then, the optimal choice of the top level block size is $b = O(m^{x^*})$, with $x^* = (\ell + \alpha)/(\ell + 1)$, which leads to the following storage and flop complexities:*

$$\mathcal{S}_{ds}^\ell(m, r) = O(m^{(\ell+2)/(\ell+1)} r^{\ell/(\ell+1)}); \quad (5.1)$$

$$\mathcal{F}_{ds}^\ell(m, r) = O(m^{(\ell+3)/(\ell+1)} r^{2\ell/(\ell+1)}). \quad (5.2)$$

Proof. The proof is inductive. We carefully track the constants hidden in the big O so as to

Table 5.1: ℓ -level equivalent of Tables 2.3 and 4.1. The legend of the table is the same. The differences between this table and the previous two are highlighted in gray.

Step	Type	$cost_{Step}^\ell$	$number_{Step}$	$\mathcal{F}_{Step}^\ell(b, p, r)$	$= \mathcal{F}_{Step}^\ell(m, x, \alpha)$
Factor	BLR $_{(\ell-1)}$	$O(b^{(\ell+2)/\ell} r^{2(\ell-1)/\ell})$	$O(p)$	$O(p b^{(\ell+2)/\ell} r^{2(\ell-1)/\ell})$	$= O(m^{1+2x/\ell+2\alpha(\ell-1)/\ell})$
Solve	BLR $_{(\ell-1)}$ -LR	$O(b^{(\ell+1)/\ell} r^{2(\ell-1)/\ell})$	$O(p^2)$	$O(p^2 b^{(\ell+1)/\ell} r^{2(\ell-1)/\ell})$	$= O(m^{2+(1-x)/\ell+\alpha(2\ell-1)/\ell})$
Product	LR-LR	$O(b r^2)$	$O(p^3)$	$O(p^3 b r^2)$	$= O(m^{3-2x+2\alpha})$
Recompress	LR	$O(b p r^2)$	$O(p^2)$	$O(p^3 b r^2)$	$= O(m^{3-2x+2\alpha})$

check their asymptotic dependence on ℓ . We therefore seek to prove the recursive bounds

$$\begin{aligned}\mathcal{S}_{ds}^\ell(m, r) &\leq C_\ell m^{(\ell+2)/(\ell+1)} r^{\ell/(\ell+1)}, \\ \mathcal{F}_{ds}^\ell(m, r) &\leq C'_\ell m^{(\ell+3)/(\ell+1)} r^{2\ell/(\ell+1)},\end{aligned}$$

where C_ℓ and C'_ℓ are constants independent of m . The formulas hold for $\ell = 1$ with $C_1 = 2$ and $C'_1 = 4$. Let us assume that they are true for the $(\ell - 1)$ -level BLR factorization and prove that they still hold for the ℓ -level one.

$\mathcal{S}_{ds}^\ell(m, r)$ can be computed as the storage cost for the off-diagonal low-rank blocks (whose number is less than p^2) plus that of the p diagonal blocks, which are represented as $(\ell - 1)$ -level BLR matrices. Therefore, by induction,

$$\mathcal{S}_{ds}^\ell(b, p, r) \leq p \times \mathcal{S}_{ds}^{\ell-1}(b, r) + p^2 b r \leq p C_{\ell-1} b^{(\ell+1)/\ell} r^{(\ell-1)/\ell} + p^2 b r.$$

We replace b , p , and r by their expression m^x , m^{1-x} , and m^α , respectively, to obtain

$$\mathcal{S}_{ds}^\ell(m, x, \alpha) \leq C_{\ell-1} m^{1-x+(\ell+1)x/\ell+\alpha(\ell-1)/\ell} + m^{2-x+\alpha}.$$

For $x^* = (\ell + \alpha)/(\ell + 1)$, we obtain

$$\mathcal{S}_{ds}^\ell(m, \alpha) \leq C_{\ell-1} m^{(\ell+2+\alpha)/(\ell+1)} + m^{(\ell+2+\alpha)/(\ell+1)}$$

and therefore

$$\mathcal{S}_{ds}^\ell(m, r) \leq C_\ell m^{(\ell+2)/(\ell+1)} r^{\ell/(\ell+1)},$$

with $C_\ell = C_{\ell-1} + 1$ (and thus $C_\ell = \ell + 1$).

We now consider the flop complexity. $\mathcal{F}_{ds}^\ell(m, r)$ can be computed as the sum of the costs of the Factor, Solve, Product and Recompress steps. We provide the ℓ -level equivalent of Tables 2.3 and 4.1 in Table 5.1. The Product and Recompress steps do not depend on ℓ and their cost is less than $p^3 b r^2 \leq m^{3-2x+2\alpha}$. The Solve step consists in applying less than p^2 times the BLR $_{(\ell-1)}$ -LR-Solve kernel, described in Algorithm 5.1, where BLR $_{(\ell-1)}$ denotes a $(\ell - 1)$ -level BLR matrix (with the convention that BLR $_0$ denotes a FR matrix). It can easily be proven by induction that the cost of applying the BLR $_{(\ell-1)}$ -LR-Solve kernel is less than $r \times \mathcal{S}_{ds}^{\ell-1}(b, r)$. Therefore, it holds

$$\mathcal{F}_{Solve}^\ell(p, b, r) \leq p^2 r \times \mathcal{S}_{ds}^{\ell-1}(b, r) \leq C_{\ell-1} p^2 b^{(\ell+1)/\ell} r^{2(\ell-1)/\ell},$$

and thus

$$\mathcal{F}_{Solve}^\ell(m, \alpha, x) \leq C_{\ell-1} m^{2+(1-\ell)x/\ell+\alpha(2\ell-1)/\ell}.$$

The Factor step consists in factorizing p diagonal $(\ell - 1)$ -level BLR matrices and therefore, by induction,

$$\mathcal{F}_{Factor}^\ell(p, b, r) \leq p C'_{\ell-1} b^{(\ell+2)/\ell} r^{2(\ell-1)/\ell} \leq C'_{\ell-1} m^{1+2x/\ell+2\alpha(\ell-1)/\ell}.$$

Summing the cost of all steps yields

$$\mathcal{F}_{ds}^\ell(m, \alpha, x) \leq C'_{\ell-1} m^{1+2x/\ell+2\alpha(\ell-1)/\ell} + C_{\ell-1} m^{2+(1-\ell)x/\ell+\alpha(2\ell-1)/\ell} + 2m^{3-2x+2\alpha}.$$

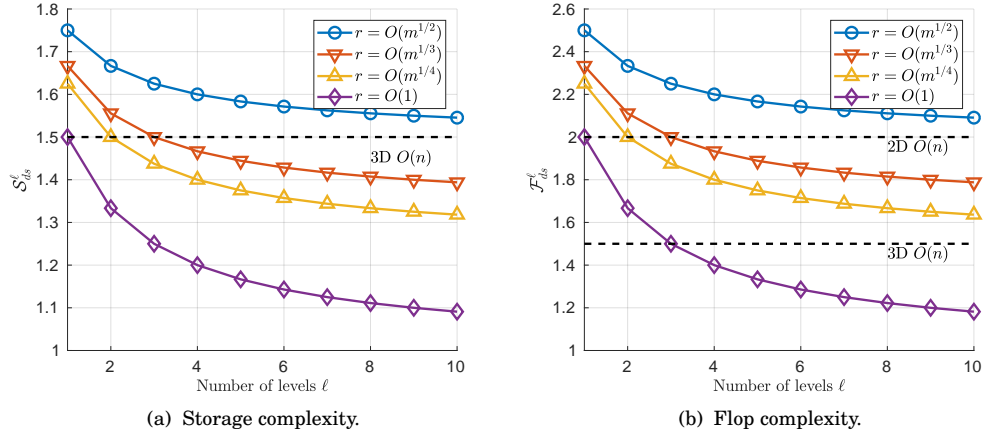


Fig. 5.1: Theoretical asymptotic exponent of the storage and flop complexity of the dense MBLR factorization

Taking $x^* = (\ell + \alpha)/(\ell + 1)$, we obtain

$$\mathcal{F}_{ds}^{\ell}(m, \alpha) \leq (C'_{\ell-1} + C_{\ell-1} + 2)m^{(\ell+3+2\alpha)/(\ell+1)}$$

and thus

$$\mathcal{F}_{ds}^{\ell}(m, r) \leq C'_{\ell} m^{(\ell+3)/(\ell+1)} r^{2\ell/(\ell+1)},$$

with $C'_{\ell} = C'_{\ell-1} + C_{\ell-1} + 2$ (and thus $C'_{\ell} = \ell(\ell + 1)/2 + \ell + 2$). Since the number of levels ℓ is assumed to be constant, we have $C_{\ell} = O(C'_{\ell}) = O(1)$ which concludes the proof. \square

Algorithm 5.1 BLR $_{\ell}$ -LR-Solve kernel, $\ell > 1$ (upper triangular case)

Input: a $p \times p$ block upper triangular BLR $_{\ell}$ matrix \tilde{U} ; $\tilde{U} = [\tilde{U}_{ij}]_{i=1:j, j=1:p}$ such that $\tilde{U}_{ij} = Y_{ij}X_{ij}^T$ for $i \neq j$ and \tilde{U}_{jj} is a BLR $_{(\ell-1)}$ matrix; a LR matrix $\tilde{B} = \Phi\Psi^T$.

Output: overwritten Ψ (modified in-place) corresponding to the operation $\tilde{B} \leftarrow \tilde{B}\tilde{U}^{-1}$.

- 1: **for** $j = 1$ **to** p **do**
 - 2: $\Psi_{j,:}^T \leftarrow \Psi_{j,:}^T - \sum_{i=1}^{j-1} \Psi_{i,:}^T Y_{ij} X_{ij}^T$
 - 3: BLR $_{(\ell-1)}$ -LR-Solve $(\tilde{U}_{jj}, \Phi\Psi_{j,:}^T)$
 - 4: **end for**
-

5.2. Influence of the number of levels ℓ . With the formulas from Theorem 5.1 proven, we now analyze their practical implications. It is clear that both the storage and flop asymptotic complexities decrease monotonically with the number of levels, while the top level block size increases.

In Figure 5.1, we plot the value of the exponent of the asymptotic complexities as a function of the number of levels, for different rank bounds $r = O(m^{\alpha})$.

Let us first consider the case of $r = O(1)$ (i.e., $\alpha = 0$). We have already shown that, in this case, the two-level BLR factorization has $O(m^{5/3})$ flop complexity, which leads to $O(n^{10/9})$ 3D sparse complexity. With a third level, the dense complexity decreases to $O(m^{3/2})$, which is precisely the 3D sparse breaking point (see Table 2.1), and thus leads to $O(n \log n)$ complexity. The log factor can be dropped by adding a fourth level. The dense complexity tends towards $O(m)$ as the number of levels increases, but the sparse complexity cannot be further improved after the optimal $O(n)$ has been reached. This illustrates that only a small number of levels is necessary to reach low sparse complexity. In particular, with four levels, the number of blocks on the top level is $p = O(m^{1/5})$,

which is still a quite large number which, as previously described, provides more flexibility to address issues such as data distribution, parallel implementation, numerical pivoting, etc.

The picture is different with higher rank bounds. Indeed, the higher the rank, the more difficult it is to reach a low complexity and thus more levels are required. For example, with $r = O(\sqrt{m})$ ($\alpha = 1/2$), it is actually not possible to reach a $O(n)$ 3D flop complexity since the dense complexity is at best $O(mr^2) = O(m^2)$. This is the dense complexity achieved by the hierarchical formats, as well as the MBLR format with an infinite number of levels, and leads to $O(n^{4/3})$ sparse complexity. It is therefore not possible to achieve this complexity with a constant number of levels. However, a crucial observation is that the rate of improvement of the exponent, which follows $(\ell + 3 + 2\alpha\ell)/(\ell + 1)$, is rapidly decreasing as ℓ increases. For example, with $\alpha = 1/2$, one level decreases the full-rank $O(m^3)$ complexity to the BLR $O(m^{2.5})$ complexity; it would require an infinite number of levels to achieve another $O(m^{0.5})$ factor of gain. Similarly, adding two more levels leads to $O(m^{2.25})$, achieving a $O(m^{0.25})$ gain which can only be achieved again with infinitely more levels! This illustrates the critical observation that *the first few levels achieve most of the asymptotic gain*. We therefore believe that the MBLR factorization with only a small number of levels can be of practical interest, even for problems with larger ranks.

5.3. Comparison with the BLR- \mathcal{H} format. We conclude this section by comparing our MBLR format to the related BLR- \mathcal{H} format, sometimes also referred to as “Lattice- \mathcal{H} ”. It consists in representing the matrix using the BLR format, and then approximating its diagonal blocks with \mathcal{H} -matrices. The BLR- \mathcal{H} format has been considered as a simple way to use hierarchical matrices in a distributed-memory setting [23, 1] but has been little studied from a theoretical standpoint. The question is whether refining the diagonal blocks with additional levels (as \mathcal{H} rather than MBLR matrices) improves the asymptotic complexity of the format. In the following, we prove that this is not the case and therefore recommend the use of the MBLR format over that of the BLR- \mathcal{H} one.

The complexity of the BLR- \mathcal{H} format is entirely determined by the block size b used for the BLR partitioning. Indeed, the storage complexity can be computed as the sum of the storage for the off-diagonal low-rank blocks and that of the diagonal \mathcal{H} blocks:

$$\mathcal{S}_{ds}^{BLR-\mathcal{H}}(p, b, r) = O(p^2 br) + O(pbr \log b) = O(p^2 br),$$

where $p = m/b$. Thus, the term corresponding to the off-diagonal low-rank blocks is dominant and we obtain

$$\mathcal{S}_{ds}^{BLR-\mathcal{H}}(m, b, r) = O\left(\frac{m^2 r}{b}\right). \quad (5.3)$$

Similarly, the flop complexity of the BLR- \mathcal{H} factorization is dominated by the LR-LR-Product and Recompress steps which cost

$$\mathcal{F}_{ds}^{BLR-\mathcal{H}}(m, b, r) = O(p^3 br^2) = O\left(\frac{m^3 r^2}{b^2}\right). \quad (5.4)$$

Applying (5.3) and (5.4) with the optimal choice of block size for the ℓ -level BLR format, $b = O(m^{\ell/(\ell+1)} r^{1/(\ell+1)})$, we obtain the same asymptotic complexity as that proved in Section 5 (Theorem 5.1):

$$\begin{aligned} \mathcal{S}_{ds}^{BLR-\mathcal{H}}(m, r) &= \mathcal{S}_{ds}^{\ell}(m, r) = O(m^{(\ell+2)/(\ell+1)} r^{\ell/(\ell+1)}); \\ \mathcal{F}_{ds}^{BLR-\mathcal{H}}(m, r) &= \mathcal{F}_{ds}^{\ell}(m, r) = O(m^{(\ell+3)/(\ell+1)} r^{2\ell/(\ell+1)}). \end{aligned}$$

This result can be interpreted as follows: for any given block size b , there exists a constant number of levels ℓ_b such that representing the diagonal blocks of the matrix as ℓ_b -level BLR matrices suffices to achieve the lowest possible complexity. As far as asymptotic complexity is concerned, it is thus not necessary to represent these diagonal blocks with the \mathcal{H} format. Note that the value of ℓ_b can easily be computed as

$$\ell_b = \min \left\{ \ell : m^{\ell/(\ell+1)} r^{1/(\ell+1)} \geq b \right\} - 1.$$

6. Numerical experiments. In this section we compare the experimental complexities of the full-rank, BLR, and MBLR formats (with different numbers of levels) for the factorization of dense matrices arising from Schur complements of sparse problems.

We have developed a MATLAB code to perform the BLR and MBLR factorization of a dense matrix, which we use to run all experiments. The objective of this experimental section is purely to validate the theoretical complexity bounds that we have computed in Theorem 5.1. Assessing the practical performance of the MBLR factorization is complex, not our focus in this paper and will be the objective of future work; numerical experiments with sparse matrix factorizations are also out of the scope of this article.

6.1. Experimental setting. All the experiments were performed on the brunch computer from the LIP laboratory (ENS Lyon), a shared-memory machine equipped with 24 Haswell processors and 1.5 TB of memory.

To validate our theoretical complexity results, we use a Poisson problem, which generates the symmetric positive definite matrix A from a 7-point finite-difference discretization of equation

$$-\Delta u = f$$

on a 3D domain of size $n = N \times N \times N$ with Dirichlet boundary conditions. We compute the dense MBLR factorization of the matrices F corresponding to the root separator of the nested dissection partitioning, which are of order $m = N^2$.

To compute the low-rank form of the blocks, we perform a truncated QR factorization with column pivoting (i.e., a truncated version of LAPACK's [8] `_geqp3` routine). We use a mix of fixed-accuracy and fixed-rank truncation: we stop the factorization after either an accuracy of ε has been achieved or at most r_{\max} columns have been computed. In the following experiments, we have set $\varepsilon = 10^{-14}$ and we compare two choices of r_{\max} , 10 and 40.

Note that in the weakly admissible context, fixing the rank to some constant yields a limited accuracy, measured by either $\|F - LU\|$ or $\|I - U^{-1}L^{-1}F\|$: in our case, both quantities are of the order of 10^{-1} , both for $r_{\max} = 10$ and $r_{\max} = 40$. Indeed, increasing r_{\max} barely improves the accuracy because some off-diagonal blocks have almost full numerical rank. While an error of the order of 10^{-1} could be enough to build a preconditioner, achieving a higher accuracy would require a strongly admissible implementation. The generalization of the MBLR format to the strongly admissible case is therefore of great importance and, while its practical implementation is outside the scope of this article, we provide its algorithmic description in the appendix. We prove that its theoretical complexity is identical to that of the weakly admissible case; we also include some very preliminary experiments showing that its complexity is in agreement with the theory while leading to a much higher accuracy.

6.2. Experimental estimates of the complexity. We first experimentally estimate the exponents in the asymptotic complexities of the dense MBLR factorization, using a number of levels ℓ varying from 1 to 4. To do so, we use the least-squares estimation of the coefficients β_1, β_2 of the regression function

$$X_{fit} = e^{\beta_1^*} N^{\beta_2^*} \text{ with } \beta_1^*, \beta_2^* = \underset{\beta_1, \beta_2}{\operatorname{argmin}} \|\log X_{obs} - \beta_1 - \beta_2 \log N\|^2.$$

such that X_{fit} fits the observed data X_{obs} .

We report in Figure 6.1 the measured storage and flops costs, along with the estimate of the exponent (β_1) obtained by the fitting. Figure 6.1 corresponds to $r_{\max} = 40$; the results for $r_{\max} = 10$ (not shown) are similar. In Table 6.1 we summarize these estimated exponents and compare them to the theoretical ones. The estimated exponents slightly differ from their theoretical counterpart, and may depend on various factors such as the choice of the block sizes, or the actual value of the numerical ranks, but overall these results are in relatively good agreement with the theory. The crucial observation is that, for both choices of r_{\max} , an asymptotic gain is achieved by each addition of a new level, at least up to $\ell = 4$.

We now analyze the sharpness of the constant prefactors in the theoretical bounds (5.1) and (5.2) computed in Theorem 5.1. To do so, we report in Figure 6.2 the ratio between the measured

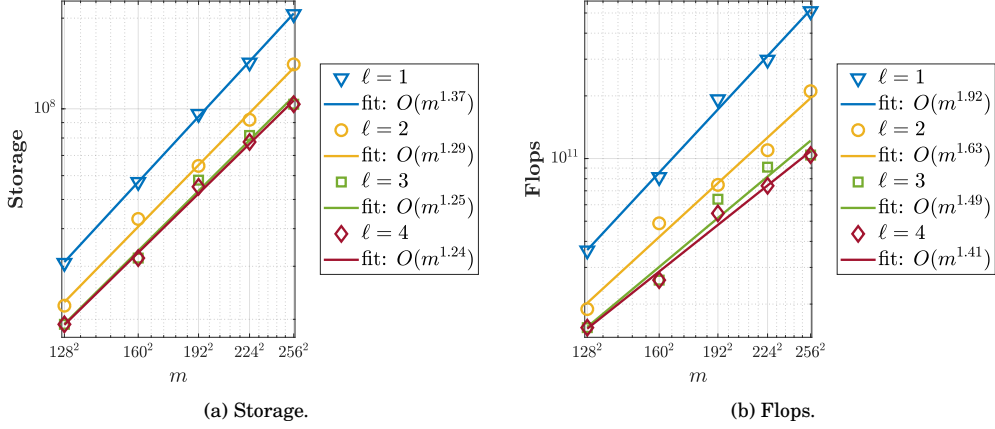


Fig. 6.1: Experimentally estimated exponents in the asymptotic complexities of the dense MBLR factorization of the root separator (of order $m = N^2$) of a 3D Poisson problem (of order $n = N^3$).

Table 6.1: Comparison between theoretical and experimental complexities.

	$\ell = 1$	$\ell = 2$	$\ell = 3$	$\ell = 4$
Storage complexity				
Theoretical	$O(m^{1.50})$	$O(m^{1.33})$	$O(m^{1.25})$	$O(m^{1.20})$
Experimental ($r_{max} = 10$)	$O(m^{1.47})$	$O(m^{1.36})$	$O(m^{1.32})$	$O(m^{1.27})$
Experimental ($r_{max} = 40$)	$O(m^{1.37})$	$O(m^{1.29})$	$O(m^{1.25})$	$O(m^{1.24})$
Flop complexity				
Theoretical	$O(m^{2.00})$	$O(m^{1.67})$	$O(m^{1.50})$	$O(m^{1.40})$
Experimental ($r_{max} = 10$)	$O(m^{1.97})$	$O(m^{1.68})$	$O(m^{1.62})$	$O(m^{1.51})$
Experimental ($r_{max} = 40$)	$O(m^{1.92})$	$O(m^{1.63})$	$O(m^{1.49})$	$O(m^{1.41})$

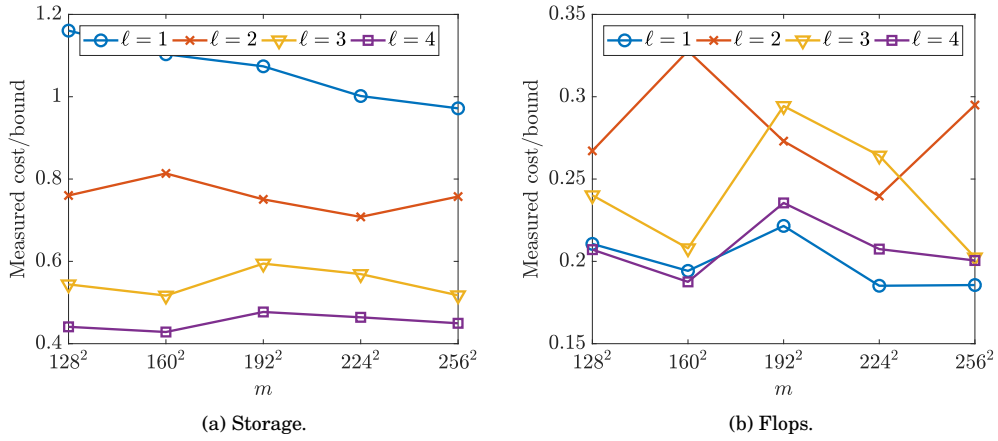


Fig. 6.2: Ratio between the measured costs and their theoretical bound for the MBLR factorization of the root separator (of order $m = N^2$) of a 3D Poisson problem (of order $n = N^3$), with $r_{max} = 40$.

cost and these bounds in the case $r_{\max} = 40$; the results for $r_{\max} = 10$ are again similar. The figure shows that the bounds are almost sharp: the actual constants are in all cases within an order of magnitude of their bound. While their sharpness seems to decrease as the number of levels ℓ increases (this is especially the case for the storage cost), the important observation is that the ratio remains more or less constant as n increases (flat curves), which means that the asymptotic cost of the MBLR factorization behaves as expected. Note that for $\ell = 1$, the storage cost is actually slightly over the bound: this may be explained by a suboptimal choice of block size.

Overall, these experimental results therefore support the capacity of the MBLR format to significantly reduce the asymptotic complexity of the factorization, even when only a small number of levels is used.

7. Conclusion. We have proposed a new multilevel BLR (MBLR) format to bridge the gap between flat and hierarchical low-rank matrix formats. Contrarily to hierarchical formats for which the number of levels in the block hierarchy is logarithmically dependent on the size of the problem, the MBLR format only uses a constant number of levels ℓ .

We had previously explained why the \mathcal{H} -matrix theory, while applicable to the BLR format, leads to very pessimistic complexity bounds and is therefore not suitable. Here, we have shown that this remains true for the MBLR format and we therefore extended the theory to compute better bounds. We proved that both the storage and flop complexities of the factorization can be finely controlled by ℓ . We theoretically showed that the first few levels achieve most of the asymptotic gain that can be expected. In particular, for a sparse 3D problem with constant ranks, two levels suffice to achieve $O(n)$ storage complexity and three levels achieve $O(n \log n)$ flop complexity, suggesting that a small number of levels may be enough in practice. Our numerical experiments confirm this trend.

Having a small number of levels leads to a greater freedom to distribute data in parallel; in particular blocks are small enough for several of them to fit in shared-memory, allowing an efficient parallelization. Finally, a small number of levels greatly simplifies the implementation of the format, making it easy to handle important features such as dynamic data structures and numerical pivoting. The related BLR- \mathcal{H} (or Lattice- \mathcal{H}) format targets a similar objective; however, our theoretical analysis shows that using more levels to refine the diagonal blocks actually does not improve the asymptotic complexity with respect to the MBLR format.

In short, the MBLR format aims to strike a balance between asymptotic complexity and actual performance on parallel computers by trading off the optimal hierarchical complexity to retrieve some of the simplicity and flexibility of flat monolevel formats. We believe that this increased simplicity and flexibility will prove to be useful in a parallel, algebraic, fully-featured, general purpose sparse direct solver. The implementation of the MBLR format in such a solver will be the object of future work.

Acknowledgements. We thank Cleve Ashcraft for useful discussions.

REFERENCES

- [1] *Private communication with R. Kriemann.*
- [2] P. R. AMESTOY, C. ASHCRAFT, O. BOITEAU, A. BUTTARI, J.-Y. L'EXCELLENT, AND C. WEISBECKER, *Improving multifrontal methods by means of block low-rank representations*, SIAM Journal on Scientific Computing, 37 (2015), pp. A1451–A1474.
- [3] P. R. AMESTOY, R. BROSSIER, A. BUTTARI, J.-Y. L'EXCELLENT, T. MARY, L. MÉTIVIER, A. MINIUSSI, AND S. OPERTO, *Fast 3D frequency-domain full waveform inversion with a parallel Block Low-Rank multifrontal direct solver: application to OBC data from the North Sea*, Geophysics, 81 (2016), pp. R363 – R383.
- [4] P. R. AMESTOY, A. BUTTARI, J.-Y. L'EXCELLENT, AND T. MARY, *On the complexity of the Block Low-Rank multifrontal factorization*, SIAM Journal on Scientific Computing, 39 (2017), pp. A1710–A1740, <https://doi.org/10.1137/16M1077192>.
- [5] P. R. AMESTOY, A. BUTTARI, J.-Y. L'EXCELLENT, AND T. MARY, *Performance and Scalability of the Block Low-Rank Multifrontal Factorization on Multicore Architectures*, ACM Transactions on Mathematical Software, (2018). Accepted for publication.
- [6] A. AMINFAR, S. AMBIKASARAN, AND E. DARVE, *A fast block low-rank dense solver with applications to finite-element matrices*, Journal of Computational Physics, 304 (2016), pp. 170–188.
- [7] A. AMINFAR AND E. DARVE, *A fast, memory efficient and robust sparse preconditioner based on a multifrontal approach with applications to finite-element matrices*, International Journal for Numerical Methods in Engineering,

- 107 (2016), pp. 520–540.
- [8] E. ANDERSON, Z. BAI, C. BISCHOF, S. BLACKFORD, J. DEMMEL, J. DONGARRA, J. DUCROZ, A. GREENBAUM, S. HAMMARLING, A. MCKENNEY, AND D. SORENSSEN, *LAPACK Users' Guide*, SIAM Press, Philadelphia, PA, third ed., 1995.
 - [9] C. ASHCRAFT, R. G. GRIMES, J. G. LEWIS, B. W. PEYTON, AND H. D. SIMON, *Progress in sparse matrix methods for large linear systems on vector computers*, Int. Journal of Supercomputer Applications, 1(4) (1987), pp. 10–30.
 - [10] M. BEBENDORF, *Hierarchical Matrices: A Means to Efficiently Solve Elliptic Boundary Value Problems*, vol. 63 of Lecture Notes in Computational Science and Engineering (LNCSE), Springer-Verlag, 2008.
 - [11] S. BÖRM, *Hierarchical matrix arithmetic with accumulated updates*, arXiv preprint arXiv:1703.09085, (2017).
 - [12] S. BÖRM, L. GRASEDYCK, AND W. HACKBUSCH, *Introduction to hierarchical matrices with applications*, Engineering analysis with boundary elements, 27 (2003), pp. 405–422, [https://doi.org/10.1016/S0955-7997\(02\)00152-2](https://doi.org/10.1016/S0955-7997(02)00152-2).
 - [13] S. CHANDRASEKARAN, M. GU, AND T. PALS, *A fast ULV decomposition solver for hierarchically semiseparable representations*, SIAM Journal on Matrix Analysis and Applications, 28 (2006), pp. 603–622.
 - [14] J. W. DEMMEL, S. C. EISENSTAT, J. R. GILBERT, X. S. LI, AND J. W. H. LIU, *A supernodal approach to sparse partial pivoting*, SIAM Journal on Matrix Analysis and Applications, 20 (1999), pp. 720–755.
 - [15] I. S. DUFF AND J. K. REID, *The multifrontal solution of indefinite sparse symmetric linear systems*, ACM Transactions on Mathematical Software, 9 (1983), pp. 302–325.
 - [16] J. A. GEORGE, *Nested dissection of a regular finite-element mesh*, SIAM Journal on Numerical Analysis, 10 (1973), pp. 345–363.
 - [17] P. GHYSELS, X. S. LI, C. GORMAN, AND F. H. ROUET, *A robust parallel preconditioner for indefinite systems using hierarchical matrices and randomized sampling*, in 2017 IEEE International Parallel and Distributed Processing Symposium (IPDPS), May 2017, pp. 897–906, <https://doi.org/10.1109/IPDPS.2017.21>.
 - [18] P. GHYSELS, X. S. LI, F.-H. ROUET, S. WILLIAMS, AND A. NAPOV, *An efficient multicore implementation of a novel hss-structured multifrontal solver using randomized sampling*, SIAM Journal on Scientific Computing, 38 (2016), pp. S358–S384, <https://doi.org/10.1137/15M1010117>.
 - [19] A. GILLMAN, P. YOUNG, AND P.-G. MARTINSSON, *A direct solver with $\mathcal{O}(N)$ complexity for integral equations on one-dimensional domains*, Frontiers of Mathematics in China, 7 (2012), pp. 217–247.
 - [20] L. GRASEDYCK AND W. HACKBUSCH, *Construction and Arithmetics of \mathcal{H} -Matrices*, Computing, 70 (2003), pp. 295–334, <https://doi.org/10.1007/s00607-003-0019-1>.
 - [21] W. HACKBUSCH, *A sparse matrix arithmetic based on \mathcal{H} -matrices. Part I: introduction to \mathcal{H} -matrices*, Computing, 62 (1999), pp. 89–108, <https://doi.org/http://dx.doi.org/10.1007/s006070050015>.
 - [22] W. HACKBUSCH, *Hierarchical matrices : algorithms and analysis*, vol. 49 of Springer series in computational mathematics, Springer, Berlin, 2015, <https://doi.org/10.1007/978-3-662-47324-5>.
 - [23] A. IDA, *Efficient Low-rank Solver for Integral Equations on Distributed Memory Systems*, in SIAM Conference on Parallel Processing (SIAM PP18), Tokyo, Japan, March 2018.
 - [24] J. W. H. LIU, *The multifrontal method for sparse matrix solution: Theory and Practice*, SIAM Review, 34 (1992), pp. 82–109.
 - [25] T. MARY, *Block Low-Rank multifrontal solvers: complexity, performance, and scalability*, PhD thesis, Université de Toulouse, November 2017, <http://thesesups.ups-tlse.fr/3854/>.
 - [26] G. PICHON, E. DARVE, M. FAVERGE, P. RAMET, AND J. ROMAN, *Sparse Supernodal Solver Using Block Low-Rank Compression*, in 18th IEEE International Workshop on Parallel and Distributed Scientific and Engineering Computing (PDSEC 2017), Orlando, United States, jun 2017.
 - [27] H. POURANSARI, P. COULIER, AND E. DARVE, *Fast hierarchical solvers for sparse matrices using low-rank approximation*, arXiv preprint arXiv:1510.07363, (2015).
 - [28] D. SHANTSEV, P. JAYSAVAL, S. DE LA KETHULLE DE RYHOVE, P. R. AMESTOY, A. BUTTARI, J.-Y. L'EXCELLENT, AND T. MARY, *Large-scale 3D EM modeling with a Block Low-Rank multifrontal direct solver*, Geophysical Journal International, 209 (2017), pp. 1558–1571.
 - [29] J. XIA, *Efficient structured multifrontal factorization for general large sparse matrices*, SIAM Journal on Scientific Computing, 35 (2013), pp. A832–A860, <https://doi.org/10.1137/120867032>.
 - [30] J. XIA, S. CHANDRASEKARAN, M. GU, AND X. S. LI, *Superfast multifrontal method for large structured linear systems of equations*, SIAM Journal on Matrix Analysis and Applications, 31 (2009), pp. 1382–1411, <https://doi.org/10.1137/09074543X>.
 - [31] J. XIA, S. CHANDRASEKARAN, M. GU, AND X. S. LI, *Fast algorithms for hierarchically semiseparable matrices*, Numerical Linear Algebra with Applications, 17 (2010), pp. 953–976.

Appendix A. Extension to the strongly-admissible case. In Sections 4 and 5, we were considering the weakly-admissible case only. In this section, we extend the proofs and computations to the strongly-admissible case. This requires the introduction of new kernels involving computations on off-diagonal non-admissible blocks that are represented as MBLR matrices. We prove that the theoretical complexity results established in the weakly-admissible case still hold in the strongly-admissible one. The proof involves two main ingredients. First, the fact that the number of non-admissible blocks per row or column on any level of the block hierarchy can be bounded by a constant. Second, the computation of the cost of the new kernels, which we show do not asymptotically increase the overall cost of the factorization.

A.1. New kernels description. In the strongly admissible case, some off-diagonal blocks may be non-admissible and therefore represented as MBLR matrices. This introduces five new kernels: the BLR_ℓ -LR and BLR_ℓ - BLR_ℓ additions, the BLR_ℓ -LR and BLR_ℓ - BLR_ℓ products, and the BLR_ℓ - BLR_ℓ triangular solve. The practical implementation of these kernels is outside our scope, but we describe them theoretically in Algorithms A.1 through A.5.

Note that the BLR_ℓ -LR (addition or product) algorithms can easily be adapted to define similar LR- BLR_ℓ algorithms. The BLR_ℓ - BLR_ℓ triangular solve (Algorithm A.5) is described in its upper triangular version; the lower triangular version is similar and is omitted for the sake of conciseness.

To avoid a profusion of different test cases, we use the compact notation $\tilde{A} \otimes \tilde{B}$ and $\tilde{A} \oplus \tilde{B}$ to denote the product and addition of two matrices \tilde{A} and \tilde{B} that may be either LR or MBLR. If at least one of the two is MBLR, then \oplus corresponds to a recursive call to Algorithms A.1 or A.2, while \otimes corresponds to a recursive call to Algorithms A.3 or A.4. If both \tilde{A} and \tilde{B} are LR, then \otimes simply corresponds to a regular LR-LR product, while \oplus corresponds to accumulating the two LR matrices together and recompressing the result. In particular, note that Algorithm A.1 requires some recompression operations corresponding to line 4.

Algorithm A.1 BLR_ℓ -LR-Addition kernel

Input: a $p \times p$ BLR_ℓ matrix \tilde{A} and a LR matrix $\tilde{B} = \Phi\Psi^T$.

Output: a BLR_ℓ matrix $\tilde{C} = \tilde{A} + \tilde{B}$

```

1: for  $i = 1$  to  $p$  do
2:   for  $j = 1$  to  $p$  do
3:     if  $\tilde{A}_{ij}$  is LR then
4:        $\tilde{C}_{ij} \leftarrow \text{Recompress}(\tilde{A}_{ij} + \Phi_i\Psi_j^T)$ 
5:     else
6:        $\tilde{C}_{ij} \leftarrow \text{BLR}_{(\ell-1)}\text{-LR-Addition}(\tilde{A}_{ij}, \Phi_i\Psi_j^T)$ 
7:     end if
8:   end for
9: end for

```

Algorithm A.2 BLR_ℓ - BLR_ℓ -Addition kernel

Input: two $p \times p$ BLR_ℓ matrices \tilde{A} and \tilde{B} .

Output: a $p \times p$ BLR_ℓ matrix $\tilde{C} = \tilde{A} + \tilde{B}$.

```

1: for  $i = 1$  to  $p$  do
2:   for  $j = 1$  to  $p$  do
3:      $\tilde{C}_{ij} \leftarrow \tilde{A}_{ij} \oplus \tilde{B}_{ij}$ 
4:   end for
5: end for

```

Algorithm A.3 BLR_ℓ -LR-Product kernel

Input: a $p \times p$ BLR_ℓ matrix \tilde{A} and a LR matrix $\tilde{C} = \Phi\Psi^T$.

Output: a LR matrix $\tilde{C}^{out} = \Phi^{out}\Psi^T = \tilde{A}\tilde{C}$

```

1: for  $i = 1$  to  $p$  do
2:    $\Phi_{i,:}^{out}\Psi^T \leftarrow \text{BLR}_{(\ell-1)}\text{-LR-Product}(\tilde{A}_{ii}, \Phi_{i,:}\Psi^T)$ 
3:    $\Phi_{i,:}^{out} \leftarrow \Phi_{i,:}^{out} + \sum_{j=1:j \neq i}^p X_{ij}Y_{ij}^T\Phi_{j,:}$ 
4: end for

```

To prove that the strong admissibility case does not increase the asymptotic complexity of the ℓ -level BLR factorization, we thus need to compute the cost of these new kernels. To do so, two ingredients are necessary. First, we prove in Section A.2 that the number of non-admissible blocks

Algorithm A.4 BLR $_{\ell}$ -BLR $_{\ell}$ -Product kernelInput: two $p \times p$ BLR $_{\ell}$ matrices \tilde{A} and \tilde{B} .Output: a $p \times p$ BLR $_{\ell}$ matrix $\tilde{C} = \tilde{A}\tilde{B}$.

```

1: for  $i = 1$  to  $p$  do
2:   for  $j = 1$  to  $p$  do
3:      $\tilde{C}_{ij} \leftarrow []$  (rank-0 matrix)
4:     for  $k = 1$  to  $p$  do
5:        $\tilde{C}_{ij} \leftarrow \tilde{C}_{ij} \oplus \tilde{A}_{ik} \otimes \tilde{B}_{kj}$ 
6:     end for
7:   end for
8: end for

```

Algorithm A.5 BLR $_{\ell}$ -BLR $_{\ell}$ -Solve kernel (upper triangular case)Input: two $p \times p$ BLR $_{\ell}$ matrices \tilde{U} and \tilde{A} ; \tilde{U} is upper triangular.Output: overwritten \tilde{A} (modified in-place) corresponding to the operation $\tilde{A} \leftarrow \tilde{A}\tilde{U}^{-1}$.

```

1: for  $k = 1$  to  $p$  do
2:   for  $i = 1$  to  $p$  do
3:     if  $\tilde{A}_{ik}$  is LR then
4:        $\tilde{A}_{ik} \leftarrow \text{BLR}_{(\ell-1)\text{-LR-Solve}}(\tilde{U}_{kk}, \tilde{A}_{ik})$ 
5:     else
6:        $\tilde{A}_{ik} \leftarrow \text{BLR}_{(\ell-1)\text{-BLR}_{(\ell-1)\text{-Solve}}}(\tilde{U}_{kk}, \tilde{A}_{ik})$ 
7:     end if
8:     for  $j = k + 1$  to  $p$  do
9:        $\tilde{A}_{ij} \leftarrow \tilde{A}_{ij} \oplus (-\tilde{A}_{ik} \otimes \tilde{U}_{kj})$ 
10:    end for
11:   end for
12: end for

```

at any level of the block hierarchy is bounded by a constant (Lemma A.1). Then, in Section A.3, we compute their cost by induction (Lemma A.2).

A.2. Boundedness of the number of non-admissible blocks. In this section, we seek to compute a bound on N_{na} , the maximal number of non-admissible blocks per row and column at any level of the block-hierarchy.

The result is directly derived from Lemma 2 in [4], whose result on BLR matrices we reproduce here.

LEMMA A.1. *For any BLR matrix, it is possible to build a partition such that the number of non-admissible blocks (as defined in [20]) per row and column is bounded by $N_{na} = O(1)$.*

Proof. See [4], Lemma 2. □

Since this result holds for any BLR matrix, it can be recursively applied to the non-admissible blocks of a ℓ -level matrix, and therefore the result trivially holds at any level of the block hierarchy.

As a consequence, the ℓ -level BLR factorization of a dense matrix clustered into $p = m/b$ blocks requires $O(p^2)$ BLR $_{(\ell-1)$ -LR products and additions and $O(p)$ BLR $_{(\ell-1)$ -BLR $_{(\ell-1)}$ products, additions, and triangular solves to be performed.

Now that we have bounded the number of times these kernels are performed, let us bound the cost of performing them once.

A.3. Recursive complexity analysis of the new kernels.

LEMMA A.2. *The costs of Algorithms A.1 through A.5 are*

$$\mathcal{F}_{A.1}^{\ell}(m, r) = \mathcal{F}_{A.2}^{\ell}(m, r) = \mathcal{F}_{A.3}^{\ell}(m, r) = O(r) \times \mathcal{S}_{ds}^{\ell}(m, r), \quad (\text{A.1})$$

$$\mathcal{F}_{A.4}^{\ell}(m, r) = \mathcal{F}_{A.5}^{\ell}(m, r) = \mathcal{F}_{ds}^{\ell}(m, r), \quad (\text{A.2})$$

where $\mathcal{F}_{ds}^\ell(m, r)$ and $\mathcal{S}_{ds}^\ell(m, r)$ are given by Theorem 5.1.

Proof. We proceed by induction. Let us note $b = O(m^x)$ the top level block size of the matrices and $p = m/b = O(m^{1-x})$. Let us also assume that $r = O(m^\alpha)$.

We begin with the initial case $\ell = 1$. Algorithm A.1 requires $O(p^2)$ LR-LR additions of cost $O(br^2)$ and $O(p)$ LR-FR additions of cost $O(b^2r)$. Algorithm A.2 requires the same operations, plus $O(p)$ FR-FR additions of cost $O(b^2)$. Algorithm A.3 requires $O(p^2)$ LR-LR products of cost $O(br^2)$ and $O(p)$ FR-LR products of cost $O(b^2r)$. Therefore, all three algorithms have a total cost of $O(p^2br^2 + pb^2r) = r \times \mathcal{S}_{ds}^1(b, p, r)$. Algorithm A.4 requires $O(p^3)$ LR-LR additions and products of cost $O(br^2)$, $O(p^2)$ LR-FR additions and products of cost $O(b^2r)$, and $O(p)$ FR-FR products of cost $O(b^3)$. Overall, the total cost of Algorithm A.4 is therefore $O(p^3br^2 + p^2b^2r + pb^3) = \mathcal{F}_{ds}^1(b, p, r)$. Algorithm A.5 requires the same computations and therefore has the same asymptotic cost.

We now assume that the formulas hold for $(\ell - 1)$ -level BLR matrices and prove them for ℓ -level ones. The previous analysis in the case $\ell = 1$ still holds by replacing every “FR” occurrence by “BLR $_{(\ell-1)}$ ”. Thus, Algorithms A.1, A.2, and A.3 require $O(p^2)$ LR-LR additions and products, $O(p)$ BLR $_{(\ell-1)}$ -LR additions and products, and $O(p)$ BLR $_{(\ell-1)}$ -BLR $_{(\ell-1)}$ additions. By induction, their cost is therefore $O(p^2br^2) + O(pr) \times \mathcal{S}_{ds}^{\ell-1}(b, r)$, which is equal to $r \times \mathcal{S}_{ds}^\ell(m, r)$ for $x = x^* = (\ell + \alpha)/(\ell + 1)$. Similarly, Algorithm A.4 requires $O(p^3)$ LR-LR additions and products, $O(p^2)$ BLR $_{(\ell-1)}$ -LR additions and products and $O(p)$ BLR $_{(\ell-1)}$ -BLR $_{(\ell-1)}$ products. By induction, its cost is therefore $O(p^3br^2) + O(p^2r) \times \mathcal{S}_{ds}^{\ell-1}(b, r) + O(p) \times \mathcal{F}_{ds}^{\ell-1}(b, r)$, which is equal to $\mathcal{F}_{ds}^\ell(m, r)$ for $x = x^*$. Algorithm A.5 requires the same computations and therefore has the same asymptotic cost. \square

THEOREM A.3. *The complexity formulas of Theorem 5.1 established in the weakly-admissible case also hold in the strongly-admissible case.*

Proof. Lemma A.1 states that there are only $O(1)$ non-admissible blocks per row and column on any level of the block hierarchy. From this we can already deduce that the storage complexity \mathcal{S}_{ds}^ℓ is asymptotically the same in the weakly- and strongly-admissible cases. Moreover, using the costs computed in Lemma A.2, we can compute the overall cost associated with the new kernels required to be performed in the strongly-admissible case as

$$\mathcal{F}_{New\ Kernels}^\ell(b, p, r) = O(p) \times \mathcal{F}_{ds}^{\ell-1}(b, r) + O(p^2r) \times \mathcal{S}_{ds}^{\ell-1}(b, r) \quad (\text{A.3})$$

$$= \mathcal{F}_{Factor}^\ell(b, p, r) + \mathcal{F}_{Solve}^\ell(b, p, r) \quad (\text{A.4})$$

which clearly does not asymptotically increase the total, thus proving that \mathcal{F}_{ds}^ℓ is also unchanged by the strong admissibility condition. \square

A.4. Preliminary storage complexity experiments. The practical implementation of the kernels presented in this appendix are outside our scope and therefore, so is the experimental study of the flop complexity of the strongly admissible MBLR factorization. However, we can easily estimate its storage complexity by first factorizing the matrix with a weakly admissible partitioning, and then compressing it with a strongly admissible partitioning. Preliminary results are reported in Figure A.1, where the blocks are compressed at accuracy $\varepsilon = 10^{-10}$ and with no limit on the rank r . This leads to a backward error $\|F - LU\|$ of the order of 10^{-9} , regardless of the number of levels ℓ . Because this is a well conditioned problem, the quantity $\|I - U^{-1}L^{-1}F\|$ is quite close to the backward error, of the order of 10^{-8} . As shown in the figure, the asymptotic storage complexity results lead to the same conclusions as those with the weak admissibility.

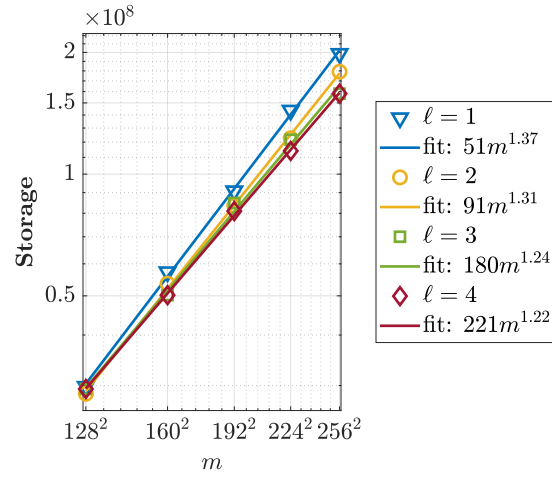


Fig. A.1: Storage complexity of the strongly admissible MBLR format, on the root separator of size $m = N^2$ of a Poisson problem of size N^3 , and with $\varepsilon = 10^{-10}$.