

Sparse factorization of the square all-ones matrix of arbitrary order

Xin Jiang* Edward Duc Hien Nguyen† César A. Uribe‡ Bicheng Ying§

January 29, 2024

Abstract

In this paper, we study sparse factorization of the (scaled) square all-ones matrix J of arbitrary order. We introduce the concept of hierarchically banded matrices and propose two types of hierarchically banded factorization of J : the reduced hierarchically banded (RHB) factorization and the doubly stochastic hierarchically banded (DSHB) factorization. Based on the DSHB factorization, we propose the sequential doubly stochastic (SDS) factorization, in which J is decomposed as a product of sparse, doubly stochastic matrices. Finally, we discuss the application of the proposed sparse factorizations to the decentralized average consensus problem and decentralized optimization.

1 Introduction

We study sparse factorization of the real $n \times n$ matrix $J := \frac{1}{n} \mathbf{1} \mathbf{1}^\top \in \mathbb{R}^{n \times n}$; that is, we seek to find a (finite) sequence of matrices $\{W^{(k)}\}_{k=1}^q \subset \mathbb{R}^{n \times n}$ such that

$$W^{(q)} W^{(q-1)} \dots W^{(1)} = \frac{1}{n} \begin{bmatrix} 1 & 1 & \dots & 1 \\ 1 & 1 & \dots & 1 \\ \vdots & \vdots & & \vdots \\ 1 & 1 & \dots & 1 \end{bmatrix}. \quad (1)$$

This problem finds applications in graph theory, systems and control, decentralized optimization, and other fields [4, 7, 8]. In this paper, we consider the general case where n is an arbitrary integer and propose several types of sparse factorization.

Previous work on the sparse factorization of J , or the all-ones matrix $\tilde{J} = nJ = \mathbf{1} \mathbf{1}^\top$, can be roughly divided into two categories. The first class considers the case in which all the factors are identical, *i.e.*, $W^{(k)} = W$ for all $k \in [q]$. For example, the binary square root of \tilde{J} (when $n = p^2$ for some $p \in \mathbb{N}_{\geq 2}$) is studied in [2]. The De Bruijn matrix, first proposed in [3], serves as the q -th root of \tilde{J} when $n = p^q$, and has been extensively studied in the literature [4, 10]. For general n , the g -circulant binary solutions to $W^q = \tilde{J}$ have also been investigated [5, 6, 11, 12].

*Department of Industrial and Systems Engineering, Lehigh University. Email: xjiang@lehigh.edu.

†Corresponding author. Department of Electrical and Computer Engineering, Rice University. Email: en18@rice.edu.

‡Department of Electrical and Computer Engineering, Rice University. Email: cauribe@rice.edu.

§Google Inc. Email: ybc@google.com.

The second class of solutions allows for differing factors of J . Among these solutions include one-peer exponential graphs (when $n = 2^q$) [13], one-peer hyper-cubes (when $n = 2^q$) [8], and p -peer hyper-cuboids [7, 9]. The p -peer hyper-cuboids serve as a factorization of J for arbitrary n , but the sparsity of the factors depends on the prime factorization of n . In particular, p -peer hyper-cuboids are no longer sparse when the matrix order equals a large prime factor. Allowing for different factors of J , in general, gives greater control over the sparsity of the factors compared to the case in which all the factors are identical [7].

In this paper, we consider the general case where $n \in \mathbb{N}_{\geq 2}$ is an arbitrary integer and study sparse factorization of J in the form

$$J = J_0 A J_0, \quad (2)$$

where $J_0 = J_1 \oplus \dots \oplus J_\tau$ with $J_k := \frac{1}{n_k} \mathbf{1}\mathbf{1}^\top \in \mathbb{R}^{n_k \times n_k}$, $k \in [\tau]$. (Here, \oplus denotes the direct sum of two matrices.) Throughout the paper, it is assumed that the partition $n = \sum_{k=1}^{\tau} n_k$ is given, with conditions that will be specified later (see (3)). Factorization (2) holds for arbitrary matrix order n and is inspired by the applications of J in decentralized averaging (and optimization). In decentralized averaging, for example, there is a group of agents where each holds a piece of information and cooperates with other agents to compute a global quantity. The communication between agents is modeled by a graph (or a sequence of graphs) $\mathcal{G}^{(k)} = (\mathcal{V}, W^{(k)}, \mathcal{E}^{(k)})$. If the weight matrices $\{W^{(k)}\}$ satisfy (1), then the *exact* global average is computed in q communication rounds. In modern application scenarios, agents can be abstracted as high-performance computing (HPC) resources and can be naturally formed into clusters [14]. Such clustering structure is captured by the proposed form of factorization (2). The block diagonal matrix J_0 models the intra-cluster communication, and each sub-block J_k , $k \in [\tau]$, can be further decomposed as (1) into, *e.g.*, g -circulant matrices or p -peer hyper-cuboids. In contrast, the A -factor models the more expensive inter-cluster communication, and the main focus of this paper is to design *sparse* A -factors to reduce the communication overhead across clusters. Sparsity in A is desirable in decentralized averaging (and optimization) as the communication overhead is related to the total number of nonzeros $\text{nnz}(A)$ as well as the largest node degree $d_{\max}(A) = \max_i \{A_{i,:}\}$, where $A_{i,:}$ is the i th row of A .

Contributions In this paper, we study the form of factorization in (2) for arbitrary matrix order n and propose three types of A -factors. In the first two types, the sparse factor A has the so-called *hierarchically banded (HB)* structure, and additional properties of A distinguish these two types of HB factorization: (density) reduced HB and doubly stochastic HB. The third one is called the *sequential doubly stochastic (SDS) factorization* and admits an asymmetric, doubly stochastic factor A , which can be further decomposed as a product of several symmetric, doubly stochastic matrices. When applied to decentralized optimization, the proposed sparse factorizations provide more flexibility to balance communication costs and the total number of communication rounds in a decentralized optimization algorithm.

Notation Let \mathbb{R} denote the set of real numbers (*i.e.*, scalars). Let \mathbb{R}^n denote the set of n -dimensional (column) vectors. (In this paper, all vectors are column vectors.) Let $\mathbb{R}^{m \times n}$ denote the set of m -by- n real matrices, and let \mathbb{D}^n denote the set of $n \times n$ diagonal matrices. The set of natural numbers is denoted as $\mathbb{N} := \{0, 1, 2, \dots\}$, and let $\mathbb{N}_{\geq r}$ denote the set of natural numbers greater than or equal to $r \in \mathbb{N}$. For any $n \in \mathbb{N}$, let $[n] := \{1, 2, \dots, n\}$. Let $\mathbf{1}$ denote the all-ones (column) vector of compatible size. The direct sum of two matrices $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{p \times q}$ forms the block diagonal matrix $A \oplus B := \text{blkdiag}(A, B) \in \mathbb{R}^{(m+p) \times (n+q)}$.

Outline In [Section 2](#), we propose the notion of hierarchically banded (HB) matrices. [Sections 3](#) and [4](#) study two types of HB factorization. The sequential doubly stochastic (SDS) factorization is discussed in [Section 5](#). In [Section 6](#), we present the potential usefulness of these sparse factorizations in decentralized averaging and optimization, and concluding remarks are offered in [Section 7](#).

2 Hierarchically banded matrices

Factorization of the form (2) relies on a partition of $n \in \mathbb{N}_{\geq 2}$:

$$n = \sum_{k=1}^{\tau} n_k, \text{ where } \{n_k\}_{k=1}^{\tau} \subset \mathbb{N}_{\geq 1}, \text{ and } n_k \geq \sum_{j=k+1}^{\tau} n_j =: m_k \text{ for all } k \in [\tau - 1]. \quad (3)$$

Such a partition can be constructed systematically, *e.g.*, via the base- p representation of n (with $p \in \mathbb{N}_{\geq 2}$). Overloading the binary representation, we denote the base- p representation of n as $(i_{\tau-1}i_{\tau-2}\cdots i_1i_0)_p$, where $\tau = \lceil \log_p(n) \rceil + 1$. Then, any integer $n \in \mathbb{N}_{\geq 2}$ can be written as $n = \sum_{k=1}^{\tau} n_k$, where $n_k = i_{\tau-k}p^{\tau-k}$. In this case, the condition $n_k \geq m_k$, for all $k \in [\tau - 1]$, directly follows from the property of the base- p representation. A simple example is $(n, p) = (15, 2)$, and $(n_1, n_2, n_3, n_4) = (8, 4, 2, 1)$, which follows from the binary representation $15 = (1111)_2$.

Given such a decomposition (3), we study the factorization in the form of (2), where the matrix $A \in \mathbb{R}^{n \times n}$ has the so-called *hierarchically banded* structure.

Definition 1 (Hierarchically banded matrices). *Given $n \in \mathbb{N}_{\geq 2}$ and a partition (3), a real symmetric $n \times n$ matrix A is called hierarchically banded (HB) if there exists a sequence of symmetric matrices $A^{(k)} \in \mathbb{R}^{m_k \times m_k}$, $k \in [\tau]$, such that the following three conditions hold.*

- $A^{(1)} = A$.
- $A^{(\tau)} \in \mathbb{D}^{n_{\tau}}$ is diagonal.
- For all $k \in [\tau - 1]$, the matrix $A^{(k)}$ can be partitioned as

$$A^{(k)} = \begin{bmatrix} A_{11}^{(k)} & A_{12}^{(k)} \\ (A_{12}^{(k)})^{\top} & A_{22}^{(k)} \end{bmatrix}, \quad (4)$$

where $A_{11}^{(k)} \in \mathbb{D}^{n_k}$, $A_{12}^{(k)} \in \mathbb{R}^{n_k \times m_k}$ have nonzero entries only on the diagonals, and the last submatrix satisfies $A_{22}^{(k)} = A^{(k+1)}$.

Such a sequence $\{A^{(k)}\}_{k=1}^{\tau}$ is called the hierarchically banded (HB) sequence of A , and the set of $n \times n$ hierarchically banded matrices is denoted by $\mathbb{H}\mathbb{B}^n$.

The hierarchically banded structure is illustrated in [Figure 1](#). The word ‘‘hierarchically’’ means that the matrix can be hierarchically partitioned, and this term is inspired by the notion of hierarchical matrices (or \mathcal{H} -matrices) (see, *e.g.*, [1]). In addition, recall that a symmetric $n \times n$ banded matrix A satisfies

$$A_{ij} = 0 \quad \text{if } j < i - r \text{ or } j > i + r,$$

where $r \in [n]$ is called the *bandwidth* of A .

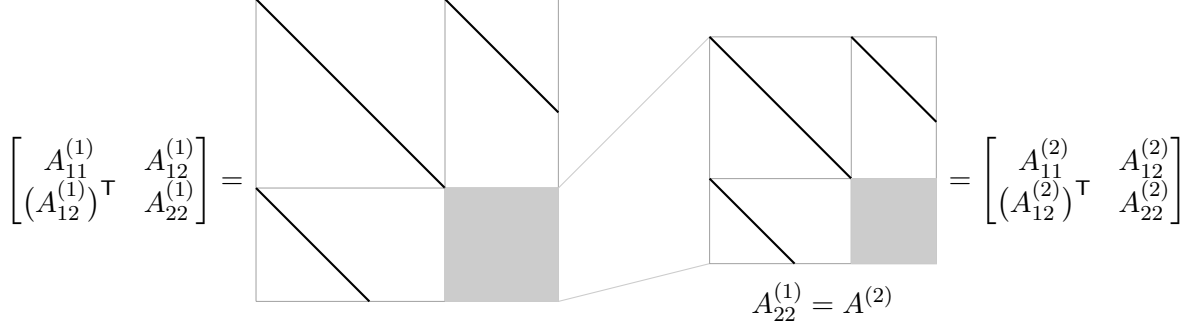


Figure 1: Illustration of a hierarchically banded matrix.

Factorization (2) with a hierarchically banded factor A is called the *hierarchically banded (HB) factorization* of J , and the matrix $A \in \mathbb{H}\mathbb{B}^n$ is called the *hierarchically banded (HB) factor* of J . It turns out that the hierarchically banded factor A is not unique. In this paper, we study the following two types of hierarchically banded factorization, characterized by additional properties of A or the HB sequence $A^{(k)}$ defined in (4).

- *Reduced hierarchically banded (RHB) factorization.* To further promote the sparsity of the A -factor, we impose an additional condition that only a few elements in the two bands of each $A^{(k)}$ are nonzero:

$$A_{12}^{(k)}[j, j] \neq 0, \quad \text{if } j = 1, 1 + n_{k+1}, 1 + n_{k+1} + n_{k+2}, \dots, 1 + \sum_{\ell=1}^{\tau-1} n_{k+\ell},$$

for all $k \in [\tau - 1]$. In $A^{(1)}$, for example, this condition means that the largest cluster (the one of size n_1) communicates with exactly *one* agent from each of the other clusters. Such a condition would further reduce the communication overhead, and the HB factor A designed for this purpose is called the (density) reduced HB factor, which is studied in [Section 3](#).

- *Doubly stochastic hierarchically banded (DSHB) factorization.* In this case, the factor A is both hierarchically banded and *doubly stochastic*, *i.e.*, all the entries in A are nonnegative, $A\mathbf{1} = \mathbf{1}$, and $A^\top\mathbf{1} = \mathbf{1}$. This additional property of A would be useful in decentralized optimization. The details are discussed in [Section 4](#).

Moreover, the DSHB factorization inspires another sparse factorization (2) of J , which is called the *sequential doubly stochastic (SDS) factorization*. In this factorization, the SDS factor A is doubly stochastic and additionally can be written as the product of a sequence of symmetric, doubly stochastic matrices. Although the SDS factor is not hierarchically banded (nor symmetric), it is closely related to the DSHB factorization and finds its application in decentralized optimization. The details of the SDS factorization are presented in [Section 5](#).

3 Reduced hierarchically banded factorization

As motivated in [Section 2](#), the (density) reduced hierarchically banded (RHB) factorization further promotes sparsity in the HB factor A by requiring

$$A_{12}^{(k)}[j, j] = \beta_k \neq 0, \quad \text{if } j = 1, 1 + n_{k+1}, 1 + n_{k+1} + n_{k+2}, \dots, 1 + \sum_{\ell=1}^{\tau-1} n_{k+\ell}, \quad (5)$$

i.e., only a few nonzeros exist in the diagonal entries of $A_{12}^{(k)}$, and these nonzeros are all equal to some $\beta \in \mathbb{R}$. In addition, it is also assumed that only one diagonal entry in each $A_{11}^{(k)}$ is not one:

$$A_{11}^{(k)} = \text{diag}(\alpha_k, 1, \dots, 1), \quad (6)$$

for some $\alpha_k \in \mathbb{R}$. (Other requirements on the diagonal submatrices $\{A_{11}^{(k)}\}$ can be applied, and they do not affect the idea of density reduction in the RHB factorization. So (6) is chosen for simplicity.) The RHB factorization is illustrated in [Section 3.1](#) via the simple example where $\tau = 2$, and [Section 3.2](#) presents an algorithm for the RHB factorization in the general case.

3.1 A two-block example

To illustrate the idea of the RHB factorization, we consider the simple case: $\tau = 2$. In this case, suppose that $n = n_1 + n_2$ with $(n_1, n_2) \in \mathbb{N}_{\geq n_2} \times \mathbb{N}_{\geq 1}$. Then, the HB factorization (2) reduces to

$$J = \begin{bmatrix} J_1 & \\ & J_2 \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} \\ A_{12}^\top & A_{22} \end{bmatrix} \begin{bmatrix} J_1 & \\ & J_2 \end{bmatrix} = \begin{bmatrix} J_1 A_{11} J_1 & J_1 A_{12} J_2 \\ (J_1 A_{12} J_2)^\top & J_2 A_{22} J_2 \end{bmatrix}, \quad (7)$$

where $J_1 = \frac{1}{n_1} \mathbf{1} \mathbf{1}^\top \in \mathbb{R}^{n_1 \times n_1}$, $J_2 = \frac{1}{n_2} \mathbf{1} \mathbf{1}^\top \in \mathbb{R}^{n_2 \times n_2}$, $A_{11} \in \mathbb{D}^{n_1}$, $A_{12} \in \mathbb{R}^{n_1 \times n_2}$, and $A_{22} \in \mathbb{D}^{n_2}$. Expanding (7) yields

$$J_1 A_{11} J_1 = \frac{1}{n} \mathbf{1}_{n_1} \mathbf{1}_{n_1}^\top, \quad J_1 A_{12} J_2 = \frac{1}{n} \mathbf{1}_{n_1} \mathbf{1}_{n_2}^\top, \quad J_2 A_{22} J_2 = \frac{1}{n} \mathbf{1}_{n_2} \mathbf{1}_{n_2}^\top. \quad (8)$$

Recall that the condition (6) requires A_{11} to take the form $A_{11} = \text{diag}(\alpha_1, 1, \dots, 1)$ for some $\alpha_1 \in \mathbb{R}$. Substituting it into $J_1 A_{11} J_1$ gives

$$J_1 A_{11} J_1 = \frac{1}{n_1^2} \mathbf{1}_{n_1} \left(\mathbf{1}_{n_1}^\top \text{diag}(\alpha_1, 1, \dots, 1) \mathbf{1}_{n_1} \right) \mathbf{1}_{n_1}^\top = \frac{\alpha_1 + n_1 - 1}{n_1^2} \mathbf{1}_{n_1} \mathbf{1}_{n_1}^\top.$$

Then, combining it with the first condition in (8) yields

$$\alpha_1 = \frac{n_1^2}{n} - n_1 + 1.$$

Similarly, one obtains that $A_{22} = \text{diag}(\alpha_2, 1, \dots, 1)$ with $\alpha_2 = \frac{n_2^2}{n} - n_2 + 1$. Finally, the condition (5) implies that A_{12} is nonzero only at the first element: $A_{12}[1, 1] := \beta$, and then expanding the second block $J_1 A_{12} J_2$ with the condition (6) gives

$$J_1 A_{12} J_2 = \frac{1}{n_1 n_2} \mathbf{1}_{n_1} \left(\mathbf{1}_{n_1}^\top A_{12} \mathbf{1}_{n_2} \right) \mathbf{1}_{n_2}^\top = \frac{\beta}{n_1 n_2} \mathbf{1}_{n_1} \mathbf{1}_{n_2}^\top,$$

Combining it with the second condition in (8) gives $A_{12}[1, 1] = \beta = \frac{n_1 n_2}{n}$.

In conclusion, when $n = n_1 + n_2$, the RHB factor A of J is given by

$$A = \left[\begin{array}{ccc|ccc} \alpha_1 & & & \beta & & \\ & 1 & & & 0 & \\ & & 1 & & & \ddots \\ & & & \ddots & & \\ & & & & & 0 \\ & & & & & 1 \\ \hline \beta & & & \alpha_2 & & \\ & 0 & & & 1 & \\ & & \ddots & & & \ddots \\ & & & 0 & & 1 \end{array} \right], \quad (9)$$

where

$$\alpha_1 = \frac{n_1^2}{n} - n_1 + 1, \quad \alpha_2 = \frac{n_2^2}{n} - n_2 + 1, \quad \beta = \frac{n_1 n_2}{n}.$$

3.2 The RHB factorization algorithm

In this section, we extend the key idea in [Section 3.1](#) to handle the general case $n = \sum_{k=1}^{\tau} n_k$, where we denote $m_k := \sum_{i=k+1}^{\tau} n_i$ for $k \in [\tau - 1]$ and assume that $n_k \geq m_k$ for all $k \in [\tau - 1]$. Then, the construction of the RHB factorization of J is summarized in [Algorithm 1](#), which outputs the RHB factor A that satisfies (2), (6), and (5).

To verify the correctness of [Algorithm 1](#), we start with the case $k = 1$ and write out the equality $J = J_0 A J_0$ for the partitioned matrices:

$$\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top = \begin{bmatrix} J_1 & \\ & \bar{J}_1 \end{bmatrix} \begin{bmatrix} A_{11}^{(1)} & A_{12}^{(1)} \\ (A_{12}^{(1)})^\top & A_{22}^{(1)} \end{bmatrix} \begin{bmatrix} J_1 & \\ & \bar{J}_1 \end{bmatrix}, \quad (12)$$

where $\bar{J}_1 := J_2 \oplus J_3 \oplus \cdots \oplus J_\tau \in \mathbb{R}^{m_1 \times m_1}$. Expanding the above equation gives three conditions similar to (8):

$$J_1 A_{11}^{(1)} J_1 = \frac{1}{n} \mathbf{1}_{n_1} \mathbf{1}_{n_1}^\top, \quad J_1 A_{12}^{(1)} \bar{J}_1 = \frac{1}{n} \mathbf{1}_{n_1} \mathbf{1}_{m_1}^\top, \quad \bar{J}_1 A_{22}^{(1)} \bar{J}_1 = \frac{1}{n} \mathbf{1}_{m_1} \mathbf{1}_{m_1}^\top. \quad (13)$$

It then follows from the condition (6) that

$$J_1 A^{(1)} J_1 = \frac{1}{n_1^2} \mathbf{1}_{n_1} \left(\mathbf{1}_{n_1}^\top \text{diag}(\alpha_1, 1, \dots, 1) \mathbf{1}_{n_1} \right) \mathbf{1}_{n_1}^\top = \frac{\alpha_1 + n_1 - 1}{n_1^2} \mathbf{1}_{n_1} \mathbf{1}_{n_1}^\top.$$

Combining it with the first condition in (13) yields $\alpha_1 = \frac{n_1^2}{n} - n_1 + 1$.

We now consider the (1,2)-block $J_1 A_{12}^{(1)} \bar{J}_1$. Recall from the condition (5) that the matrix

Algorithm 1 Reduced hierarchically banded (RHB) factorization algorithm

1: **Input:** $n \in \mathbb{N}_{\geq 2}$, and the factors $\{n_k\}_{k=1}^\tau$ satisfying $n = \sum_{k=1}^\tau n_k$ and $n_k \geq m_k = \sum_{i=k+1}^\tau n_i$ for all $k \in [\tau - 1]$.

2: **Output:** The RHB factor A of J , and the associated HB sequence $\{A^{(k)}\}_{k=1}^\tau$.

3: Set $m_0 \leftarrow n$.

4: **for** $k = 1, 2, \dots, \tau - 2$ **do**

5: Compute the (1, 1)-block $A_{11}^{(k)} \in \mathbb{D}^{n_k}$ of $A^{(k)}$:

$$A_{11}^{(k)} \leftarrow \text{diag} \left(\frac{n_k^2}{m_{k-1}} - n_k + 1, 1, \dots, 1 \right).$$

6: Compute the (1, 2)-block $A_{12}^{(k)} \in \mathbb{R}^{n_k \times m_k}$:

$$A_{12}^{(k)}[i, j] \leftarrow \begin{cases} \frac{n_k n_{k+1}}{n} & \text{if } i = j = 1 \\ \frac{n_k n_{k+\ell}}{n} & \text{if } i = j = 1 + \sum_{r=1}^{\ell} n_{k+r} \text{ for } \ell = 1, 2, \dots, \tau - k - 1 \\ 0 & \text{otherwise.} \end{cases}$$

7: Compute the (2, 2)-block $A_{22}^{(k)} = A^{(k+1)}$ as the RHB factorization:

$$\frac{1}{m_k} \mathbf{1}_{m_k} \mathbf{1}_{m_k}^\top = \bar{J}_k A^{(k+1)} \bar{J}_k, \quad (10)$$

where $\bar{J}_k := J_{k+1} \oplus \dots \oplus J_\tau$, and the OHHG factor $A^{(k+1)} = A_{22}^{(k)}$ is partitioned as

$$A^{(k+1)} = \begin{bmatrix} A_{11}^{(k+1)} & A_{12}^{(k+1)} \\ (A_{12}^{(k+1)})^\top & A_{22}^{(k+1)} \end{bmatrix}. \quad (11)$$

8: **end for**

9: Set the RHB factor A : $A \leftarrow A^{(1)}$.

$A_{12}^{(1)} \in \mathbb{R}^{n_1 \times m_1}$ can be partitioned as

$$A_{12}^{(1)} = \begin{bmatrix} B_2^{(1)} & & & & \\ & B_3^{(1)} & & & \\ & & \ddots & & \\ & & & B_{\tau-1}^{(1)} & \\ \mathbf{0}_2 & \mathbf{0}_3 & \cdots & \mathbf{0}_{\tau-1} & \mathbf{0}_\tau \end{bmatrix},$$

where $B_j^{(1)} = \text{diag}(\beta_j^{(1)}, 0, \dots, 0) \in \mathbb{D}^{n_j}$, and $\mathbf{0}_j$ is the all-zeros matrix of size $(n_1 - m_1) \times n_j$, for $j = 2, \dots, \tau$. In addition, we denote the diagonal entries of $B_j^{(1)}$ by the n_j -vector $b_j^{(1)} = (\beta_j^{(1)}, 0, \dots, 0)$. Then, it holds that

$$\mathbf{1}_{n_1}^\top A_{12}^{(1)} = \left[(b_2^{(1)})^\top \quad (b_3^{(1)})^\top \quad \cdots \quad (b_\tau^{(1)})^\top \right] \in \mathbb{R}^{1 \times m_1}.$$

Then, it holds that

$$\begin{aligned}
J_1 A_{12}^{(1)} \bar{J} &= \frac{1}{n_1} \mathbf{1}_{n_1} (\mathbf{1}_{n_1}^\top A_{12}^{(1)}) \bar{J}_1 \\
&= \frac{1}{n_1} \mathbf{1}_{n_1} \left[(b_2^{(1)})^\top \quad (b_3^{(1)})^\top \quad \cdots \quad (b_\tau^{(1)})^\top \right] (J_2 \oplus J_3 \oplus \cdots \oplus J_\tau) \\
&= \frac{1}{n_1} \mathbf{1}_{n_1} \left[\frac{1}{n_2} (b_2^{(1)})^\top \mathbf{1}_{n_2} \mathbf{1}_{n_2}^\top \quad \frac{1}{n_3} (b_3^{(1)})^\top \mathbf{1}_{n_3} \mathbf{1}_{n_3}^\top \quad \cdots \quad \frac{1}{n_\tau} (b_\tau^{(1)})^\top \mathbf{1}_{n_\tau} \mathbf{1}_{n_\tau}^\top \right] \\
&= \frac{1}{n_1} \mathbf{1}_{n_1} \left[\frac{\beta_2^{(1)}}{n_2} \mathbf{1}_{n_2} \mathbf{1}_{n_2}^\top \quad \frac{\beta_3^{(1)}}{n_2} \mathbf{1}_{n_3} \mathbf{1}_{n_3}^\top \quad \cdots \quad \frac{\beta_\tau^{(1)}}{n_\tau} \mathbf{1}_{n_\tau} \mathbf{1}_{n_\tau}^\top \right] \\
&= \left[\frac{\beta_2^{(1)}}{n_1 n_2} \mathbf{1}_{n_1} \mathbf{1}_{n_2}^\top \quad \frac{\beta_3^{(1)}}{n_1 n_3} \mathbf{1}_{n_1} \mathbf{1}_{n_3}^\top \quad \cdots \quad \frac{\beta_\tau^{(1)}}{n_1 n_\tau} \mathbf{1}_{n_1} \mathbf{1}_{n_\tau}^\top \right] \in \mathbb{R}^{n_1 \times m_1}.
\end{aligned}$$

Hence, to satisfy the second condition in (13), we must have for all $j = 2, \dots, \tau$ that

$$\frac{\beta_j^{(1)}}{n_1 n_j} = \frac{1}{n} \quad \iff \quad \beta_j^{(1)} = \frac{n_1 n_j}{n}.$$

Finally, we consider the last condition in (13). Denote $A^{(2)} := A_{22}^{(1)}$ and consider the partition (11). Also notice that $\bar{J}_1 = J_2 \oplus (J_3 \oplus \cdots \oplus J_\tau) := J_2 \oplus \bar{J}_2$. Then, we write out the last condition (13) in the partitioned form:

$$\frac{1}{n} \mathbf{1}_{m_1} \mathbf{1}_{m_1}^\top = \begin{bmatrix} J_2 & \\ & \bar{J}_2 \end{bmatrix} \begin{bmatrix} A_{11}^{(2)} & A_{12}^{(2)} \\ (A_{12}^{(2)})^\top & A_{22}^{(2)} \end{bmatrix} \begin{bmatrix} J_2 & \\ & \bar{J}_2 \end{bmatrix},$$

which takes the same form as (12). We can repeat the above process for $k = 1, 2, \dots, \tau - 2$. When Algorithm 1 reaches iteration $k = \tau - 2$, Line Algorithm 1 computes the RHB factorization of the matrix

$$A^{(\tau-1)} = \begin{bmatrix} A_{11}^{(\tau-1)} & A_{12}^{(\tau-1)} \\ (A_{12}^{(\tau-1)})^\top & A_{22}^{(\tau-1)} \end{bmatrix},$$

which is the two-block case studied in Section 3.1. Thus, the RHB factor of $A^{(\tau-1)}$ is in the form of (9) with

$$\alpha_1 = \frac{n_{\tau-1}^2}{n_{\tau-1} + n_\tau} - n_{\tau-1} + 1, \quad \alpha_2 = \frac{n_\tau^2}{n_{\tau-1} + n_\tau} - n_\tau + 1, \quad \beta = \frac{n_{\tau-1} n_\tau}{n_{\tau-1} + n_\tau}.$$

From the above discussion, we obtain the following result.

Theorem 1. *The $n \times n$ matrix A_{RHB} generated by Algorithm 1 is hierarchically banded, satisfies the conditions (5)–(6), and satisfies $J = J_0 A_{\text{RHB}} J_0$. In addition, the total number of nonzeros is $\text{nnz}(A_{\text{RHB}}) = n + \tau(\tau - 1)$, and the largest node degree is $d_{\max}(A_{\text{RHB}}) = \tau$.*

Proof. (The subscript in A_{RHB} (and $\tilde{A}_{\text{RHB}}^{(k)}$) is omitted in the proof for readability.) The hierarchically banded structure of A follows from the recursive nature of Algorithm 1, and in particular, the recursive partition of $A^{(k)}$ in Algorithm 1. Similarly, A satisfies the conditions (5)–(6) due to the assignment of values in $A_{11}^{(k)}$ and $A_{12}^{(k)}$ in Algorithm 1. Next, the factorization $J = J_0 A J_0 = J_0 A^{(1)} J_0$

holds by recursively applying (10) for $k = \tau - 1, \tau - 2, \dots, 1$. Finally, one has for all $k \in [\tau - 1]$ that $\text{nnz}(A_{11}^{(k)}) = n_k$, $\text{nnz}(A_{12}^{(k)}) = \tau - k$, and $\text{nnz}(A^{(\tau)}) = n_\tau$. So, the total number of nonzeros is

$$\text{nnz}(A) = \sum_{k=1}^{\tau-1} (n_k + 2(\tau - k)) + n_\tau = n + \tau(\tau - 1).$$

The row with the most nonzeros is row $n - n_\tau - n_{\tau-1} + 1$, where $A_{n-n_\tau-n_{\tau-1}+1,j} \neq 0$ if $j = 1, 1 + n_1, 1 + n_1 + n_2, \dots, 1 + \sum_{k=1}^{\tau-1} n_k$, and thus $d_{\max}(A) = \tau$. \square

4 Doubly stochastic hierarchically banded factorization

Another useful type of hierarchically banded factorization, especially in decentralized optimization (see Section 6.2 for details), requires the HB factor A to be *doubly stochastic*, i.e., all the entries are nonnegative and $A\mathbf{1} = \mathbf{1}$ (and $A^\top \mathbf{1} = \mathbf{1}$, which is guaranteed by the symmetry of A). Yet in this case, the HB sequence $\{A^{(k)}\}$ is *not* doubly stochastic. Instead, we show that each matrix in the *scaled* HB sequence $\{\tilde{A}^{(k)}\}_{k=1}^\tau$ remains doubly stochastic, where

$$\tilde{A}^{(k)} := \frac{n}{m_{k-1}} A^{(k)} \in \text{HIB}^{m_k}, \quad k \in [\tau]. \quad (14)$$

Again, the doubly stochastic hierarchically banded (DSHB) factorization is illustrated in Section 4.1 via the simple example where $\tau = 2$, and Section 4.2 presents an algorithm for DSHB factorization in the general case.

4.1 A two-block example

Similar to Section 3.1, we start with the simple case where $\tau = 2$, and assume $n = n_1 + n_2$ with $(n_1, n_2) \in \mathbb{N}_{\geq n_2} \times \mathbb{N}_{\geq 1}$. Then, the HB factorization takes the form of (7), which can be partitioned as in (8). Recall that in Section 3.1 we require both submatrices A_{11} and A_{12} to have only one nonzero entry. In the context of decentralized optimization, a larger cluster will communicate with exactly one agent from each of the smaller clusters. This section considers a different setting where all agents in subgroup 2 (recall $n_2 \leq n_1$) can communicate across subgroups. In particular, the submatrix $A_{12} \in \mathbb{R}^{n_1 \times n_2}$ takes the following form:

$$A_{12} = \begin{bmatrix} \text{diag}(\beta \mathbf{1}_{n_2}) \\ 0 \end{bmatrix}.$$

Substituting into $J_1 A_{12} J_2$ gives

$$J_1 A_{12} J_2 = \frac{1}{n_1 n_2} \mathbf{1}_{n_1} \left(\mathbf{1}_{n_1}^\top A_{12} \mathbf{1}_{n_2} \right) \mathbf{1}_{n_2}^\top = \frac{\beta}{n_1}.$$

Then, the second condition in (8) implies that

$$\beta = \frac{n_1}{n} = \frac{n_1}{n_1 + n_2}.$$

With the subblock A_{12} settled, the doubly stochastic property of A implies that the subblocks $A_{11} \in \mathbb{D}^{n_1}$ and $A_{22} \in \mathbb{D}^{n_2}$ are diagonal matrices satisfying

$$A_{11} = \text{diag}(\underbrace{1 - \beta, \dots, 1 - \beta}_{n_2}, \underbrace{1, \dots, 1}_{n_1 - n_2}), \quad A_{22} = \text{diag}((1 - \beta) \mathbf{1}_{n_2}).$$

Finally, we confirm that this choice of A_{11} and A_{22} also satisfies the first and third conditions in (8):

$$\begin{aligned}\frac{1}{n_1} \mathbb{1}_{n_1} \mathbb{1}_{n_1}^\top A_{11} \mathbb{1}_{n_1} \mathbb{1}_{n_1}^\top &= \frac{\mathbb{1}_{n_1}^\top A_{11} \mathbb{1}_{n_1}}{n_1^2} \mathbb{1}_{n_1} \mathbb{1}_{n_1}^\top = \frac{(1-\beta)n_2 + (n_1 - n_2)}{n_1^2} \mathbb{1}_{n_1} \mathbb{1}_{n_1}^\top = \frac{1}{n} \mathbb{1}_{n_1} \mathbb{1}_{n_1}^\top, \\ \frac{1}{n_2} \mathbb{1}_{n_2} \mathbb{1}_{n_2}^\top A_{22} \mathbb{1}_{n_2} \mathbb{1}_{n_2}^\top &= \frac{\mathbb{1}_{n_2}^\top A_{22} \mathbb{1}_{n_2}}{n_2^2} \mathbb{1}_{n_2} \mathbb{1}_{n_2}^\top = \frac{(1-\beta)n_2}{n_2^2} \mathbb{1}_{n_2} \mathbb{1}_{n_2}^\top = \frac{1}{n} \mathbb{1}_{n_2} \mathbb{1}_{n_2}^\top.\end{aligned}$$

In conclusion, when $n = n_1 + n_2$, the doubly stochastic HB factor A of J is

$$A = \left[\begin{array}{cc|c} \frac{n_2}{n} I_{n_2} & 0 & \frac{n_1}{n} I_{n_2} \\ 0 & I_{n_1 - n_2} & 0 \\ \hline \frac{n_1}{n} I_{n_2} & 0 & \frac{n_2}{n} I_{n_1} \end{array} \right]. \quad (15)$$

4.2 The DSHB factorization algorithm

We extend the key idea in Section 3.1 to handle the general case where $n = \sum_{k=1}^{\tau} n_k$. The construction of the DSHB factor A , as well as the associated (scaled) HB sequence $\{A^{(k)}\}$ ($\{\tilde{A}^{(k)}\}$ in (14)), is summarized in Algorithm 2.

To verify the correctness of Algorithm 2, we start with the iteration $k = 1$ and write out the equality $J = J_0 A J_0$ for the partitioned matrices:

$$\frac{1}{n} \mathbb{1}_n \mathbb{1}_n^\top = \begin{bmatrix} J_1 & \\ & \bar{J}_1 \end{bmatrix} \begin{bmatrix} A_{11}^{(1)} & A_{12}^{(1)} \\ (A_{12}^{(1)})^\top & A_{22}^{(1)} \end{bmatrix} \begin{bmatrix} J_1 & \\ & \bar{J}_1 \end{bmatrix},$$

where recall $\bar{J}_1 := J_2 \oplus J_3 \oplus \dots \oplus J_\tau \in \mathbb{R}^{m_1 \times m_1}$ and $A^{(1)} = \tilde{A}^{(1)}$. Expanding the above equation gives three conditions similar to (8):

$$J_1 A_{11}^{(1)} J_1 = \frac{1}{n} \mathbb{1}_{n_1} \mathbb{1}_{n_1}^\top, \quad J_1 A_{12}^{(1)} \bar{J}_1 = \frac{1}{n} \mathbb{1}_{n_1} \mathbb{1}_{m_1}^\top, \quad \bar{J}_1 A_{22}^{(1)} \bar{J}_1 = \frac{1}{n} \mathbb{1}_{m_1} \mathbb{1}_{m_1}^\top. \quad (19)$$

For the (1, 2)-block $A_{12}^{(1)}$, we follow the convention in Section 4.1 and assume that it has the structure

$$A_{12}^{(1)} = \begin{bmatrix} \text{diag}(\beta^{(1)} \mathbb{1}_{m_1}) \\ 0 \end{bmatrix}.$$

Substituting into $J_1 A_{12}^{(1)} \bar{J}_1$ gives

$$J_1 A_{12}^{(1)} \bar{J}_1 = \frac{1}{n_1} \mathbb{1}_{n_1} (\mathbb{1}_{n_1}^\top A_{12}^{(1)}) \bar{J}_1 = \frac{\beta^{(1)}}{n_1} \mathbb{1}_{n_1} \mathbb{1}_{m_1}^\top \bar{J}_1 = \frac{\beta^{(1)}}{n_1} \mathbb{1}_{n_1} \mathbb{1}_{m_1}^\top.$$

Combining it with the second condition in (19) yields $\beta^{(1)} = \frac{n_1}{n}$. Then, the doubly stochastic property of A implies that

$$A_{11}^{(1)} = \text{diag} \left(\underbrace{\frac{m_1}{n}, \dots, \frac{m_1}{n}}_{m_1}, \underbrace{1, \dots, 1}_{n_1 - m_1} \right), \quad A_{22}^{(1)} \mathbb{1}_{m_1} = (1 - \beta^{(1)}) \mathbb{1}_{m_1} = \frac{m_1}{n} \mathbb{1}_{m_1}.$$

The second equation above is equivalent to the doubly stochastic property of the *scaled* matrix

$$\tilde{A}^{(2)} \mathbb{1}_{m_1} = \mathbb{1}_{m_1}, \quad \text{where } \tilde{A}^{(2)} := \frac{n}{m_1} A_{22}^{(1)} \in \mathbb{R}^{m_1 \times m_1}.$$

Algorithm 2 Doubly stochastic hierarchically banded (DSHB) factorization algorithm

- 1: **Input:** $n \in \mathbb{N}_{\geq 2}$, and the factors $\{n_k\}_{k=1}^{\tau}$ satisfying $n = \sum_{k=1}^{\tau} n_k$ and $n_k \geq m_k = \sum_{i=k+1}^{\tau} n_i$ for all $k \in [\tau - 1]$.
- 2: **Output:** The doubly stochastic HB factor A of J , and the associated HB sequence $\{A^{(k)}\}_{k=1}^{\tau}$.
- 3: Set $m_{-1} \leftarrow n$ and $m_0 \leftarrow n$.
- 4: **for** $k = 1, 2, \dots, \tau - 2$ **do**
- 5: Compute the (1, 1)-block $\tilde{A}_{11}^{(k)} \in \mathbb{D}^{n_k}$ of $\tilde{A}^{(k)}$:

$$\tilde{A}_{11}^{(k)} \leftarrow \text{diag} \left(\underbrace{\frac{m_k}{m_{k-1}}, \dots, \frac{m_k}{m_{k-1}}}_{m_k}, \underbrace{1, \dots, 1}_{n_k - m_k} \right). \quad (16)$$

- 6: Compute the (1, 2)-block $\tilde{A}_{12}^{(k)} \in \mathbb{R}^{n_k \times m_k}$:

$$\tilde{A}_{12}^{(k)}[i, j] \leftarrow \begin{cases} \frac{n_k}{m_{k-1}} & \text{if } i = j = 1, 2, \dots, m_k \\ 0 & \text{otherwise.} \end{cases} \quad (17)$$

- 7: Compute the (2, 2)-block $\tilde{A}_{22}^{(k)}$ from the DSHB factorization:

$$\frac{1}{m_k} \mathbf{1}_{m_k} \mathbf{1}_{m_k}^{\top} = \bar{J}_k \tilde{A}^{(k+1)} \bar{J}_k, \quad (18)$$

where $\bar{J}_k := J_{k+1} \oplus \dots \oplus J_{\tau}$, and the DSHB factor $\tilde{A}^{(k+1)} \leftarrow \frac{m_{k-1}}{m_k} \tilde{A}_{22}^{(k)}$ is partitioned as

$$\tilde{A}^{(k+1)} = \begin{bmatrix} \tilde{A}_{11}^{(k+1)} & \tilde{A}_{12}^{(k+1)} \\ (\tilde{A}_{12}^{(k+1)})^{\top} & \tilde{A}_{22}^{(k+1)} \end{bmatrix}.$$

- 8: **end for**

- 9: Set the DSHB factor $A \leftarrow A^{(1)} \equiv \tilde{A}^{(1)}$ and the associated HB sequence $A^{(k)} \leftarrow \frac{m_{k-1}}{n} \tilde{A}^{(k)}$, for all $k \in [\tau]$.
-

Similarly, the third condition in (19) can be written in terms of $\tilde{A}^{(2)}$ as

$$\bar{J}_1 \tilde{A}^{(2)} \bar{J}_1 = \frac{1}{m_1} \mathbf{1}_{m_1} \mathbf{1}_{m_1}^{\top}. \quad (20)$$

Therefore, to find a doubly stochastic, hierarchically banded matrix $\tilde{A}^{(2)}$ that satisfies (20), we need to construct the DSHB factorization of $\frac{1}{m_1} \mathbf{1}_{m_1} \mathbf{1}_{m_1}^{\top}$, which requires recursive execution of the above process for $k = 1, 2, \dots, \tau - 2$.

When [Algorithm 2](#) reaches iteration $k = \tau - 2$, [Line 7](#) computes the DSHB factor

$$\tilde{A}^{(\tau-1)} = \begin{bmatrix} \tilde{A}_{11}^{(\tau-1)} & \tilde{A}_{12}^{(\tau-1)} \\ (\tilde{A}_{12}^{(\tau-1)})^{\top} & \tilde{A}_{22}^{(\tau-1)} \end{bmatrix},$$

which is the two-block case studied in [Section 4.1](#). Thus, the DSHB factor of $\frac{1}{m_{\tau-2}} \mathbf{1}_{m_{\tau-2}} \mathbf{1}_{m_{\tau-2}}^{\top}$ is

in the form of (15):

$$\tilde{A}^{(\tau-1)} = \left[\begin{array}{cc|c} \alpha I_{n_\tau} & 0 & \beta I_{n_\tau} \\ 0 & I_{n_{\tau-1}-n_\tau} & 0 \\ \hline \beta I_{n_\tau} & 0 & \alpha I_{n_\tau} \end{array} \right], \text{ where } \alpha = \frac{n_\tau}{n_{\tau-1} + n_\tau} \text{ and } \beta = \frac{n_{\tau-1}}{n_{\tau-1} + n_\tau}.$$

From the above discussion, we obtain the following result.

Theorem 2. *The $n \times n$ matrix A_{DSHB} generated by Algorithm 2 is doubly stochastic, hierarchically banded, and satisfies $J = J_0 A_{\text{DSHB}} J_0$. Each matrix in the scaled HB sequence $\{\tilde{A}_{\text{DSHB}}^{(k)}\}_{k=1}^\tau$ generated by Algorithm 2 is doubly stochastic. In addition, $\text{nnz}(A_{\text{DSHB}}) = \sum_{k=1}^\tau kn_k$, and $d_{\max}(A_{\text{DSHB}}) = \tau$.*

Proof. (The subscript in A_{DSHB} (and $\tilde{A}_{\text{DSHB}}^{(k)}$) is omitted in the proof for readability.) The doubly stochastic property of A and $\{\tilde{A}^{(k)}\}$ follows from the assignments (16)–(17) and condition (18). The hierarchically banded structure of A follows from the recursive nature of Algorithm 2, and in particular, the recursive partition of $\tilde{A}^{(k)}$ in Algorithm 2. Next, the factorization $J = J_0 A J_0$ holds by recursively applying (18) for $k = \tau - 1, \tau - 2, \dots, 1$. Finally, the number of nonzeros and the largest node degree can be calculated using the same approach as for the RHB factor in Theorem 1. \square

5 Sequential doubly stochastic factorization

The DSHB factorization inspires another type of factorization for J , in which the factor $A \in \mathbb{R}^{n \times n}$ in $J = J_0 A J_0$ is no longer symmetric (nor hierarchically banded) but remains doubly stochastic. Since the asymmetric, doubly stochastic matrix A can be written as the product of a sequence of doubly stochastic matrices, such a factorization is called the *sequential doubly stochastic (SDS) factorization* of J .

Theorem 3 (Sequential doubly stochastic (SDS) factorizations of J). *Let $A \in \mathbb{H}\mathbb{B}^n$ be the DSHB factor of J and $\{\tilde{A}^{(k)}\}_{k=1}^\tau$ the associated scaled HB sequence, constructed via Algorithm 2. For all $k \in [\tau - 1]$, define*

$$T^{(k)} := \begin{bmatrix} \tilde{A}_{11}^{(k)} & \tilde{A}_{12}^{(k)} \\ (\tilde{A}_{12}^{(k)})^\top & \frac{m_k}{m_{k-1}} I_{m_k} \end{bmatrix} \in \mathbb{R}^{m_{k-1} \times m_{k-1}}, \quad (21)$$

with the convention $m_0 := n$, and $T^{(\tau)} := \tilde{A}^{(\tau)} \equiv I_{n_\tau}$. The augmented matrices $\{\hat{T}^{(k)}\}_{k=1}^\tau \subset \mathbb{R}^{n \times n}$ are defined as

$$\hat{T}^{(k)} = I_{n_1} \oplus I_{n_2} \oplus \dots \oplus I_{n_{k-1}} \oplus T^{(k)}, \quad k \in [\tau].$$

Then, the matrices $\{T^{(k)}\}_{k=1}^\tau$ (and $\{\hat{T}^{(k)}\}_{k=1}^\tau$) are all symmetric, doubly stochastic, and the matrix $J = \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top$ can be factored as

$$J = J_0 A_L J_0 = J_0 A_R J_0, \quad (22)$$

where

$$A_L := \hat{T}^{(1)} \hat{T}^{(2)} \dots \hat{T}^{(\tau)} \quad (23a)$$

$$= T^{(1)} \cdot (I_{n_1} \oplus (T^{(2)} \cdot (I_{n_2} \oplus \dots (T^{(\tau-1)} \cdot (I_{n_\tau} \oplus T^{(\tau)}))))),$$

$$A_R := \hat{T}^{(\tau)} \hat{T}^{(\tau-1)} \dots \hat{T}^{(1)} \quad (23b)$$

$$= (I_{n_1} \oplus \cdots \oplus (I_{n_{\tau-1}} \oplus (I_{n_\tau} \oplus T^{(\tau)}) \cdot T^{(\tau-1)}) \cdot T^{(\tau-2)}) \cdots T^{(1)}. \quad (23c)$$

In addition, both factors A_L and A_R are doubly stochastic.

By definition, the matrices $\{T^{(k)}\}$ have nonzero entries only in three subdiagonals. The first factorization in (22) is called the *left SDS factorization* of J , and the second is called the *right SDS factorization*.

Proof. Define $\bar{J}_0 := J_0$, $m_0 := n$, and

$$\begin{aligned} V^{(k)} &:= T^{(k)} \cdot (I_{n_k} \oplus (T^{(k+1)} \cdot (I_{n_{k+1}} \oplus \cdots (T^{(\tau-1)} \cdot (I_{n_\tau} \oplus T^{(\tau)})))))) \\ &= T^{(k)} \cdot (I_{n_k} \oplus V^{(k+1)}) \in \mathbb{R}^{m_{k-1} \times m_{k-1}}, \end{aligned} \quad (24)$$

for all $k \in [\tau - 1]$, and $V^{(\tau)} := T^{(\tau)} \equiv I_{n_\tau}$. By definition, each matrix $V^{(k)}$ is doubly stochastic, because $T^{(k)}$ is doubly stochastic.

First, we apply mathematical induction to prove that for $k = \tau - 1, \dots, 1$,

$$\bar{J}_{k-1} V^{(k)} \bar{J}_{k-1} = \frac{1}{m_{k-1}} \mathbf{1}_{m_{k-1}} \mathbf{1}_{m_{k-1}}^\top, \quad (25)$$

The base case $k \leftarrow \tau - 1$ holds because

$$\bar{J}_{\tau-2} V^{(\tau-1)} \bar{J}_{\tau-2} = \begin{bmatrix} J_{\tau-1} & \\ & J_\tau \end{bmatrix} T^{(\tau-1)} (I_{n_\tau} \oplus T^{(\tau)}) \begin{bmatrix} J_{\tau-1} & \\ & J_\tau \end{bmatrix} \quad (26a)$$

$$= \begin{bmatrix} J_{\tau-1} & \\ & J_\tau \end{bmatrix} T^{(\tau-1)} \begin{bmatrix} J_{\tau-1} & \\ & J_\tau \end{bmatrix} \quad (26b)$$

$$= \begin{bmatrix} J_{\tau-1} & \\ & J_\tau \end{bmatrix} \begin{bmatrix} \tilde{A}_{11}^{(\tau-1)} & \tilde{A}_{12}^{(\tau-1)} \\ (\tilde{A}_{12}^{(\tau-1)})^\top & \frac{m_{\tau-1}}{m_{\tau-2}} I \end{bmatrix} \begin{bmatrix} J_{\tau-1} & \\ & J_\tau \end{bmatrix} \quad (26c)$$

$$\begin{aligned} &= \begin{bmatrix} J_{\tau-1} \tilde{A}_{11}^{(\tau-1)} J_{\tau-1} & J_{\tau-1} \tilde{A}_{12}^{(\tau-1)} J_\tau \\ (J_{\tau-1} \tilde{A}_{12}^{(\tau-1)} J_\tau)^\top & \frac{m_{\tau-1}}{m_{\tau-2}} J_\tau^2 \end{bmatrix} \\ &= \frac{1}{m_{\tau-2}} \mathbf{1}_{m_{\tau-2}} \mathbf{1}_{m_{\tau-2}}^\top. \end{aligned} \quad (26d)$$

The first equation (26a) uses $\bar{J}_{\tau-2} = J_{\tau-1} \oplus \bar{J}_{\tau-1} = J_{\tau-1} \oplus J_\tau$. Then, (26b) and (26c) use the definition $T^{(\tau)} = I_{n_\tau}$ and (21). Finally, (26d) follows from the updates (16) and (17) in Algorithm 2, as well as the fact $J_\tau^2 = J_\tau$.

Next, suppose the identity (25) holds for $k \in [\tau - 1]$, and we establish the same identity with $k \leftarrow k - 1$:

$$\begin{aligned} &\bar{J}_{k-2} T^{(k-1)} (I_{n_k} \oplus V^{(k)}) \bar{J}_{k-2} \\ &= \begin{bmatrix} J_{k-1} & \\ & \bar{J}_{k-1} \end{bmatrix} \begin{bmatrix} \tilde{A}_{11}^{(k-1)} & \tilde{A}_{12}^{(k-1)} \\ (\tilde{A}_{12}^{(k-1)})^\top & \frac{m_{k-1}}{m_{k-2}} I \end{bmatrix} \begin{bmatrix} I_{n_{k-1}} & \\ & V^{(k)} \end{bmatrix} \begin{bmatrix} J_{k-1} & \\ & \bar{J}_{k-1} \end{bmatrix} \end{aligned} \quad (27a)$$

$$= \begin{bmatrix} J_{k-1} A_{11}^{(k-1)} J_{k-1} & J_{k-1} A_{12}^{(k-1)} V^{(k)} \bar{J}_{k-1} \\ (J_{k-1} A_{12}^{(k-1)} \bar{J}_{k-1})^\top & \frac{m_{k-1}}{m_{k-2}} \bar{J}_{k-1} V^{(k)} \bar{J}_{k-1} \end{bmatrix} \quad (27b)$$

$$\begin{aligned}
&= \begin{bmatrix} \frac{1}{m_{k-2}} \mathbb{1}_{n_{k-1}} \mathbb{1}_{n_{k-1}}^\top & \frac{1}{m_{k-2}} \mathbb{1}_{n_{k-1}} \mathbb{1}_{m_{k-1}}^\top \\ \frac{1}{m_{k-2}} \mathbb{1}_{m_{k-1}} \mathbb{1}_{n_{k-1}}^\top & \frac{1}{m_{k-2}} \mathbb{1}_{m_{k-1}} \mathbb{1}_{m_{k-1}}^\top \end{bmatrix} \\
&= \frac{1}{m_{k-2}} \mathbb{1}_{m_{k-2}} \mathbb{1}_{m_{k-2}}^\top.
\end{aligned} \tag{27c}$$

In (27a), we use $\bar{J}_{k-2} = J_{k-1} \oplus \bar{J}_{k-1}$, the definition of $T^{(k)}$ in (21), and the definition of direct sum. After multiplying out all the matrices in (27b), the third equation (27c) follows from the definition of $A_{11}^{(k-1)}$ and $A_{12}^{(k-1)}$ (see (16)–(17)), the definition of $V^{(k)}$ in (24), and the assumption that identity (25) holds for $k \in [\tau-1]$. In particular, the (1,2)-subblock of (27b) is further simplified as follows:

$$\begin{aligned}
&J_{k-1} A_{12}^{(k-1)} V^{(k)} \bar{J}_{k-1} \\
&= \frac{1}{n_{k-1}} \mathbb{1}_{n_{k-1}} \left(\mathbb{1}_{n_{k-1}}^\top \left[\text{diag} \left(\frac{n_{k-1}}{m_{k-2}} \mathbb{1}_{m_{k-1}} \right) \right] \right) V^{(k)} \bar{J}_{k-1}
\end{aligned} \tag{28a}$$

$$= \frac{1}{m_{k-2}} \mathbb{1}_{n_{k-1}} (\mathbb{1}_{m_{k-1}}^\top V^{(k)}) \bar{J}_{k-1} \tag{28b}$$

$$= \frac{1}{m_{k-2}} \mathbb{1}_{n_{k-1}} \mathbb{1}_{m_{k-1}}^\top \bar{J}_{k-1} \tag{28c}$$

$$= \frac{1}{m_{k-2}} \mathbb{1}_{n_{k-1}} \mathbb{1}_{m_{k-1}}^\top. \tag{28d}$$

In (28a), we use the definition of $\tilde{A}_{12}^{(k-1)}$ in (17), and (28b) writes out $\mathbb{1}_{n_{k-1}}^\top \tilde{A}_{12}^{(k-1)} = \frac{n_{k-1}}{m_{k-2}} \mathbb{1}_{m_{k-1}}^\top$. Then, (28c) and (28d) use the doubly stochastic property of $V^{(k)}$ and \bar{J}_{k-1} , respectively.

Therefore, the induction hypothesis is proved, and (25) holds for all $k \in [\tau-1]$. In particular, (25) with $k \leftarrow 1$ gives the left SDS factorization:

$$J_0 A_L J_0 = \bar{J}_0 V^{(1)} \bar{J}_0 = \frac{1}{m_0} \mathbb{1}_{m_0} \mathbb{1}_{m_0}^\top = J.$$

The first equation uses the convention $\bar{J}_0 = J_0$ and the relation $A_L = T^{(1)}(I_{n_1} \oplus V^{(2)}) = V^{(1)}$. The second equation applies (25) with $k = 1$, and the last one follows from the convention $m_0 = n$. Then, the right SDS factorization follows directly from the fact that $A_R = A_L^\top$ and thus

$$J = (J_0 A_L J_0)^\top = J_0 A_L^\top J_0 = J_0 A_R J_0.$$

Finally, the doubly stochastic property of A_L (and A_R) follows from that of $\{T^{(k)}\}$ and the fact that the product of doubly stochastic matrices is still doubly stochastic. \square

In the context of decentralized optimization, if communication is modeled by the T -factors, then at each round of communication, each agent only needs to communicate with at most one neighbor (as $d_{\max}(T^{(k)}) = 2$ for all $k \in [\tau-1]$). Such a property is called ‘‘one-peer’’ in decentralized optimization and holds for one-peer hyper-cubes [8] and one-peer exponential graphs [13].

We also note that the matrices $\{T^{(k)}\}$ represent the base- $(p+1)$ graphs introduced in [9]. Yet, the original work [9] fails to provide an explicit matrix representation for the base- $(p+1)$ graphs and does not prove that the weight matrices of their proposed base- $(p+1)$ graphs can be used to factorize the J matrix. Moreover, as explained in Section 2, the construction of all the matrices (A_L , A_R ,

$\{T^{(k)}\}$, and $\{\tilde{A}^{(k)}\}$) does not necessarily rely on the base- p representation of the integer $n \in \mathbb{N}_{\geq 2}$, and only needs a decomposition $n = \sum_{k=1}^{\tau} n_k$ with $n_k \geq m_k = \sum_{i=k+1}^{\tau} n_i$ for all $k \in [\tau - 1]$. So, the original name “base- $(p + 1)$ ” does not fully reveal the flexibility of the sequential doubly stochastic factorization proposed in this paper.

The following corollary presents the basic properties of the two SDS factors and the T -factors.

Corollary 4. *The total number of nonzeros in the matrix $T^{(k)}$ is $\text{nnz}(T^{(k)}) = n_k + 2 \sum_{i=k+1}^{\tau} n_i$, for $k \in [\tau]$, and the largest node degree is $d_{\max}(T^{(k)}) = 2$. In addition,*

$$\text{nnz}(A_L) = \text{nnz}(A_R) = \sum_{k=1}^{\tau} (2^k - 1)n_k, \quad d_{\max}(A_L) = \tau, \quad d_{\max}(A_R) = 2^{\tau-1}.$$

Proof. The total number of nonzeros in the matrix $T^{(k)}$ and the largest node degree $d_{\max}(T^{(k)})$ hold from the definition (21).

It follows from the definition of $V^{(k)}$ (24) that

$$\text{nnz}(V^{(k)}) = n_k + m_k + 2\text{nnz}(V^{(k+1)}), \quad \text{for all } k \in [\tau - 1],$$

and $\text{nnz}(V^{(\tau)}) = n_{\tau}$. Then, recursion over k yields

$$\begin{aligned} \text{nnz}(A_L) = \text{nnz}(A_R) &= \text{nnz}(V^{(1)}) = n_1 + m_1 + 2\text{nnz}(V^{(2)}) \\ &= n_1 + m_1 + 2(n_2 + m_2) + 4\text{nnz}(V^{(3)}) \\ &\quad \vdots \\ &= \sum_{k=1}^{\tau-1} 2^{k-1}(n_k + m_k) + 2^{\tau-1}\text{nnz}(V^{(\tau)}) \\ &= \sum_{k=1}^{\tau-1} 2^{k-1}(n_k + m_k) + 2^{\tau-1}n_{\tau} \\ &= \sum_{k=1}^{\tau} 2^{k-1}n_k + \sum_{k=1}^{\tau-1} 2^{k-1}m_{\tau} \\ &= \sum_{k=1}^{\tau} 2^{k-1}n_k + \sum_{k=1}^{\tau-1} 2^{k-1} \sum_{i=k+1}^{\tau} n_i \\ &= \sum_{k=1}^{\tau} 2^{k-1}n_k + \sum_{k=2}^{\tau} (2^{k-1} - 1)n_k \\ &= \sum_{k=1}^{\tau} (2^k - 1)n_k. \end{aligned}$$

Similarly, the largest node degree of A_L (and A_R) can be calculated as

$$\begin{aligned} d_{\max}(A_L) &= d_{\max}(V^{(1)}) = d_{\max}(V^{(2)}) + 1 = \dots = d_{\max}(V^{(\tau)}) + \tau - 1 = \tau, \\ d_{\max}(A_R) &= d_{\max}((V^{(1)})^{\top}) = 2d_{\max}(V^{(2)}) = \dots = 2^{\tau-1}d_{\max}(V^{(\tau)}) = 2^{\tau-1}. \end{aligned}$$

□

6 Application in decentralized average consensus and optimization

In this section, we show how the presented factorizations of the form (2) can be used in decentralized averaging (in Section 6.1) and then describe extensions to decentralized optimization (in Section 6.2).

6.1 Decentralized average consensus

We first formulate the decentralized average consensus problem as follows. In a group of n agents, each one holds a piece of information, denoted by $x_i^{(0)} \in \mathbb{R}^d$, and the entire group aims to compute the average $\bar{x} := \frac{1}{n} \sum_{i=1}^n x_i^{(0)}$ via communication. The communication (or connection) between agents is modeled by a sequence of (undirected) graphs (or topologies) $\mathcal{G}^{(k)} = (\mathcal{V}, W^{(k)}, \mathcal{E}^{(k)})$, where $\mathcal{V} = \{1, \dots, n\}$ is the vertex set representing agents and each $\mathcal{E}^{(k)} \subseteq \mathcal{V} \times \mathcal{V}$ is the set of edges (or connections). It is assumed that the set of agents remains static while the set of edges can be time-varying. The entry $w_{ij}^{(k)} \in \mathbb{R}_{\geq 0}$ in the *mixing matrix* $W^{(k)}$ applies a weighting factor to the information exchanged between agent j and agent i . If $w_{ij}^{(k)} = 0$, it means agent i is not a neighbor of agent j in $\mathcal{G}^{(k)}$; *i.e.*, $(i, j) \notin \mathcal{E}^{(k)}$. The state of agent i (or the information held by agent i) at iteration k is designated as $x_i^{(k)}$ and evolves according to the following recursion: for $k \in \mathbb{N}$,

$$x_i^{(k+1)} = \sum_{j: (i,j) \in \mathcal{E}^{(k)}} w_{ij}^{(k)} x_j^{(k)}, \quad \text{for all } i \in [n].$$

the above recursion can be written more compactly as

$$X^{(k+1)} = W^{(k)} X^{(k)}, \quad \text{where } X^{(k)} = \begin{bmatrix} x_1^{(k)} & x_2^{(k)} & \dots & x_n^{(k)} \end{bmatrix}^\top \in \mathbb{R}^{n \times d}. \quad (29)$$

We say average consensus is achieved if either of the following conditions is satisfied.

1. The limit of each $x_i^{(k)}$ is \bar{x} :

$$\lim_{k \rightarrow \infty} x_i^{(k)} = \bar{x}, \quad \text{for all } i \in [n].$$

2. There exists $\bar{k} \in \mathbb{N}$ such that $X^{(\bar{k})} = \bar{x} \mathbf{1}^\top$ and $X^{(k)} = \bar{x} \mathbf{1}^\top$ for all $k \in \mathbb{N}_{\geq \bar{k}}$.

In modern application scenarios involving GPUs and high-performance computing (HPC) resources, we can design and alter the communication topology in the decentralized averaging process. In this case, sparse factorization of J helps design topologies with cheap communication cost per iteration and achieve consensus within a finite number of iterations (29). To see this, consider a set of sparse matrices $\{W^{(i)}\}_{i=1}^q$ that satisfies $J = W^{(q)} \dots W^{(2)} W^{(1)}$. When the associated graph sequence $\{\mathcal{G}^{(i)}\}_{i=1}^q$ is used as the (time-varying) topologies for decentralized averaging, the iteration (29) yields

$$X^{(q)} = W^{(q)} W^{(q-1)} \dots W^{(2)} W^{(1)} X^{(0)} = \frac{1}{n} \mathbf{1} \mathbf{1}^\top X^{(0)} = \bar{x} \mathbf{1}^\top.$$

Therefore, unlike classical results where consensus is achieved only asymptotically, *finite-time consensus* (*i.e.*, consensus in *exactly* q iterations) in decentralized averaging is achieved by exploiting sparse factorization of J .

To this end, the proposed HB and SDS factorizations of J can be used to construct sparse graph sequences with the desirable finite-time consensus property for *arbitrary* number of agents $n \in \mathbb{N}_{\geq 2}$.

Matrices in phase 2	A_{RHB}	A_{DSHB}	A_{L}	A_{R}	T -factors
Num. of nonzeros	$n + \tau(\tau - 1)$	$\sum_{k=1}^{\tau} kn_k$	$\sum_{k=1}^{\tau} (2^k - 1)n_k$	$\sum_{k=1}^{\tau} (2^k - 1)n_k$	$n_k + 2 \sum_{i=k+1}^{\tau} n_i$
Largest node degree d_{max}	τ	τ	τ	$2^{\tau-1}$	2
Num. of iter. in Phase 2	1	1	1	1	$\tau - 1$

Table 1: Trade-offs between the communication cost (modeled by the largest node degree d_{max}) and the number of iterations in phase 2.

This is in contrast to most previous work, which has requirements on the matrix order n (e.g., $n = p^\tau$ for some $(p, \tau) \in \mathbb{N}_{\geq 2} \times \mathbb{N}_{\geq 1}$). Below, we describe in detail how to exploit the factorization $J = J_0 A J_0$ to construct graph sequences $\{W^{(i)}\}$ with finite-time consensus, and then discuss two additional advantages of the proposed HB and SDS factorizations.

- **Phase 1.** The communication network is constructed via a sparse factorization of $J_0 = J_1 \oplus \dots \oplus J_\tau$. For example, each smaller matrix $J_j \in \mathbb{R}^{n_j \times n_j}$ can be decomposed as product of p -peer hyper-cuboids [7]. Then, each mixing matrix in Phase 1 is a direct sum of several p -peer hyper-cuboids (and identity matrices).
- **Phase 2.** This phase corresponds to the A matrix in (2), which can be the RHB factor, the DSHB factor, the (left or right) SDS factor, or even a sequence of T -factors. A detailed comparison between these choices is discussed in the next paragraph and presented in Table 1.
- **Phase 3.** It corresponds to a sparse factorization of J_0 , and can be the same as Phase 1.

In addition to the ability to handle an arbitrary number of agents, the proposed factorizations (RHB, DSHB, SDS) provide more flexibility to balance the communication costs and the number of communication rounds toward consensus. Recall that the communication cost involved in the decentralized averaging iteration (29) is related to the total number of nonzeros and the largest node degree in the communication graph. Thus, using sparser mixing matrices would reduce communication costs but likely increase the total number of averaging iterations toward consensus. For example, using the left (or right) SDS factor A_{L} (or A_{R}) completes Phase 2 in one iteration, while using the “one-peer” T -factors in (21) results in $\tau - 1$ iterations in Phase 2. Such a trade-off in the choice of Phase-2 matrices is summarized in Table 1.

Moreover, the proposed form of factorization (2) can resolve a practical concern omitted by classical settings of decentralized averaging (and optimization). It is typically assumed that the distance between agents is equidistant and that each agent is indistinguishable from another. In practice, however, it may not be the case. Consider the scenario where agents are HPC resources. Agents A and B may be allocated on the same physical machine while agent C is on another physical machine. Consequently, the communication cost between agents A and B is cheaper than that between agents A and C. It is natural to model the network as several sub-networks where each sub-network is a cluster of relatively “close” agents. Such a structure is easily exploited by factorization $J = J_0 A J_0$. Communication in phases 1 and 3 is all intra-cluster and can be modeled by different sparse factorizations of J_k , $k \in [\tau]$. The more expensive inter-cluster communication only happens in Phase 2 and is modeled by the sparse matrix A , which can be the RHB factor, the DSHB factor, the SDS factors (A_{L} or A_{R}), or a sequence of T -factors in (21). The proposed factorization form (2) promotes cheap, intra-cluster communications and limits the more expensive, inter-cluster ones.

6.2 Decentralized optimization

Besides decentralized averaging, sparse factorization of J is also useful in decentralized optimization. In decentralized optimization, agents collaborate to solve the following optimization problem

$$\text{minimize } f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x), \quad (30)$$

where the optimization variable is $x \in \mathbb{R}^d$, and each component function $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$ is continuously differentiable and potentially nonconvex. Each agent $i \in \mathcal{V}$ only has access to one component function f_i , and agents communicate with each other via (time-varying) topologies $\{\mathcal{G}^{(k)}\}$. It can be shown that the decentralized average consensus problem is a special case of the optimization Problem (30) with $f_i(x) = \frac{1}{2} \|x - x_i^{(0)}\|_2^2$.

In the context of decentralized optimization, a sparse factorization of J offers sequences of graphs that satisfy the finite-time consensus property, and incorporating such graph sequences in decentralized optimization algorithms could significantly reduce the per-iteration communication cost in the algorithm while achieving a comparable convergence rate (compared with decentralized algorithms using traditional communication protocols) [7, 13]. Also note that most existing analyses for decentralized optimization algorithms need the assumption that the weight (mixing) matrices $\{W^{(k)}\}$ are doubly stochastic, which is satisfied by the DSHB factor and the T -factor (21) in the SDS factorization. So both the DSHB factorization (in Section 4) and the SDS factorization (in Section 5) are helpful in decentralized optimization, while the RHB factorization (in Section 3) is not suitable in this scenario.

7 Conclusion

In this paper, we study the sparse factorization $J = J_0 A J_0$, where $J = \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T$ is the (scaled) all-ones matrix and $J_0 = J_1 \oplus \dots \oplus J_\tau$ is the direct sum of several smaller all-ones matrices. We introduce the hierarchically banded structure of a symmetric matrix, based on which we present two types of hierarchically banded factorization of J : the reduced hierarchically banded (RHB) factorization and the doubly stochastic hierarchically banded (DSHB) factorization. Moreover, inspired by the DSHB factorization, we propose the sequential doubly stochastic (SDS) factorization which further factorizes the matrix A as the product of a sequence of symmetric, doubly stochastic matrices. We then discuss the usefulness of the proposed factorizations in the decentralized average consensus problem and decentralized optimization. The presented three types of sparse factorization offer much flexibility in handling the trade-off between the per-iteration communication cost and the total number of communication rounds in decentralized averaging (and optimization).

Finally, recall that the partition $n = \sum_{k=1}^{\tau} n_k$ is assumed to be given and fixed throughout the paper. Further investigation is needed in the design of this partition to fully leverage the power of the proposed sparse factorizations in decentralized optimization.

References

- [1] Steffen Börm, Lars Grasedyck, and Wolfgang Hackbusch. Introduction to hierarchical matrices with applications. *Engineering Analysis with Boundary Elements*, 27(5):405–422, 2003.

- [2] Frank E. Curtis, John Drew, Chi-Kwong Li, and Daniel Prigel. Central groupoids, central digraphs, and zero-one matrices A satisfying $A^2 = J$. *Journal of Combinatorial Theory, Series A*, 105(1):35–50, 2004.
- [3] Nicolaas G. de Bruijn. A combinatorial problem. *Proceedings of the Section of Sciences of the Koninklijke Nederlandse Akademie van Wetenschappen te Amsterdam*, 49(7):758–764, 1946.
- [4] Jean-Charles Delvenne, Ruggero Carli, and Sandro Zampieri. Optimal strategies in the average consensus problem. *Systems & Control Letters*, 58(10-11):759–765, 2009.
- [5] Fenn King and Kai Wang. On the g -circulant solutions to the matrix equation $A^m = \lambda J$. *Journal of Combinatorial Theory, Series A*, 38(2):182–186, 1985.
- [6] S. L. Ma and William C. Waterhouse. The g -circulant solutions of $A^m = \lambda J$. *Linear Algebra and its Applications*, 85:211–220, 1987.
- [7] Edward Duc Hien Nguyen, Xin Jiang, Bicheng Ying, and César A. Uribe. On graphs with finite-time consensus and their use in gradient tracking. *arXiv preprint*, arXiv:2311.01317, 2023.
- [8] Guodong Shi, Bo Li, Mikael Johansson, and Karl Henrik Johansson. Finite-time convergent gossiping. *IEEE/ACM Transactions on Networking*, 24(5):2782–2794, 2016.
- [9] Yuki Takezawa, Ryoma Sato, Han Bao, Kenta Niwa, and Makoto Yamada. Beyond exponential graph: Communication-efficient topologies for decentralized learning via finite-time convergence. *arXiv preprint*, arXiv:2305.11420, 2023.
- [10] Maguy Trefois, Paul Van Dooren, and Jean-Charles Delvenne. Binary factorizations of the matrix of all ones. *Linear Algebra and its Applications*, 468:63–79, 2015.
- [11] Kai Wang. On the g -circulant solutions to the matrix equation $A^m = \lambda J$. *Journal of Combinatorial Theory, Series A*, 33(3):287–296, 1982.
- [12] Yao-Kun Wu, Rui-Zhong Jia, and Qiao Li. g -Circulant solutions to the $(0, 1)$ matrix equation $A^m = J_n$. *Linear Algebra and its Applications*, 345(1):195–224, 2002.
- [13] Bicheng Ying, Kun Yuan, Yiming Chen, Hanbin Hu, Pan Pan, and Wotao Yin. Exponential graph is provably efficient for decentralized deep training. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P. S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 13975–13987. Curran Associates, Inc., 2021.
- [14] Bicheng Ying, Kun Yuan, Hanbin Hu, Yiming Chen, and Wotao Yin. Bluefog: Make decentralized algorithms practical for optimization and deep learning, 2021.