

A LINEAR-COMPLEXITY TENSOR BUTTERFLY ALGORITHM FOR COMPRESSING HIGH-DIMENSIONAL OSCILLATORY INTEGRAL OPERATORS

P. MICHAEL KIELSTRA*, TIANYI SHI†, HENGRUI LUO ‡, JIANLIANG QIAN §, AND
YANG LIU†

Abstract. This paper presents a multilevel tensor compression algorithm called tensor butterfly algorithm for efficiently representing large-scale and high-dimensional oscillatory integral operators, including Green’s functions for wave equations and integral transforms such as Radon transforms and Fourier transforms. The proposed algorithm leverages a tensor extension of the so-called complementary low-rank property of existing matrix butterfly algorithms. The algorithm partitions the discretized integral operator tensor into subtensors of multiple levels, and factorizes each subtensor at the middle level as a Tucker-like interpolative decomposition, whose factor matrices are formed in a multilevel fashion. For a d -dimensional integral operator discretized into a $2d$ -mode tensor with n^{2d} entries, the overall CPU time and memory requirement scale as $O(n^d)$, in stark contrast to the $O(n^d \log n)$ requirement of existing matrix algorithms such as matrix butterfly algorithm and fast Fourier transforms (FFT), where n is the number of points per direction. When comparing with other tensor algorithms such as quantized tensor train (QTT), the proposed algorithm also shows superior CPU and memory performance for tensor contraction. Remarkably, the tensor butterfly algorithm can efficiently model high-frequency Green’s function interactions between two unit cubes, each spanning 512 wavelengths per direction, which represents over $512\times$ larger problem sizes than existing algorithms. On the other hand, for a problem representing 64 wavelengths per direction, which is the largest size existing algorithms can handle, our tensor butterfly algorithm exhibits $200\times$ speedups and $30\times$ memory reduction comparing with existing ones. Moreover, the tensor butterfly algorithm also permits $O(n^d)$ -complexity FFTs and Radon transforms up to $d = 6$ dimensions.

Key word. butterfly algorithm, tensor algorithm, Tucker decomposition, interpolative decomposition, quantized tensor train (QTT), fast Fourier transforms (FFT), fast algorithm, high-frequency wave equations, integral transforms, Radon transform, low-rank compression, Fourier integral operator, non-uniform FFT (NUFFT)

AMS subject classifications. 15A23, 65F50, 65R10, 65R20

1. Introduction. Oscillatory integral operators (OIOs), such as Fourier transforms and Fourier integral operators [32, 7], are critical computational and theoretical tools for many scientific and engineering applications, such as signal and image processing, inverse problems and imaging, computer vision, quantum mechanics, and analyzing and solving partial differential equations (PDEs). The development of accurate and efficient algorithms for computing OIOs has profound impacts on the evolution of the pertinent research areas including, perhaps mostly remarkably, the invention of the fast Fourier transform (FFT) by Cooley and Tukey in 1965 and the invention of the fast multipole method (FMM) by Greengard and Rokhlin in 1987, both of which were listed among the ten most significant algorithms discovered in the 20th century. Among existing analytical and algebraic methods for OIOs, butterfly algorithms [52, 46, 37, 36, 56] represent an emerging class of multilevel matrix decomposition algorithms that have been proposed for Fourier transforms and Fourier

*Department of Mathematics, University of California, Berkeley, CA, USA.
Email: pmkielstra@berkeley.edu

†Applied Mathematics and Computational Research Division, Lawrence Berkeley National Laboratory, Berkeley, CA, USA.
Email: {tianyishi, liuyangzhuan}@lbl.gov

‡Department of Statistics, Rice University, Houston, TX, USA. Email: hrluo@rice.edu

§Department of Mathematics and Department of CMSE, Michigan State University, East Lansing, MI, USA. Email: jqian@msu.edu

integral operators [8, 68, 67], special function transforms [63, 4, 54], fast iterative [53, 52, 47] and direct [24, 43, 25, 26, 59, 44, 60] solution of surface and volume integral equations for wave equations, high-frequency Green’s function ansatz for inhomogeneous wave equations [45, 41, 48], direct solution of PDE-induced sparse systems [42, 13], and machine learning for inverse problems [33, 35]. The (matrix) butterfly algorithms leverage the so-called complementary low-rank (CLR) property of the matrix representation of OIOs after proper row/column permutation. The CLR states that judiciously selected submatrices exhibit numerical low ranks, known as the butterfly ranks, which stay constant irrespective of the matrix sizes. This permits a multilevel sparse matrix decomposition requiring $O(n \log n)$ factorization time, application time, and storage units with n being the matrix size.

Despite their low asymptotic complexity, the matrix butterfly algorithms often-times exhibit relatively large prefactors, i.e., constant but high butterfly ranks, particularly for higher-dimensional OIOs. Examples include Green’s functions for 3D high-frequency wave equations [59, 45], 3D Radon transforms for linear inverse problems [17], 6D Fourier–Bros–Iagolnitzer transforms for Wigner equations [15, 66], 6D Fourier transforms in diffusion magnetic resonance imaging [11] and plasma physics [18], 4D space-time transforms in quantum field theories [57, 49], and multi-particle Green’s functions in quantum chemistry [21]. For these high-dimensional OIOs, the computational advantage of the matrix butterfly algorithms over other existing algorithms becomes significant only for very large matrices.

More broadly speaking, for large-scale multi-dimensional scientific data and operators, tensor algorithms are typically more efficient than matrix algorithms. Popular low-rank tensor compression algorithms include CANDECOMP/PARAFAC [30], Tucker [16], hierarchical Tucker [28], tensor train (TT) [55], and tensor network [12] decomposition algorithms. See references [34, 23] for a more complete review of available tensor formats and their applications. When applied to the representation of high-dimensional integral operators, tensor algorithms often leverage additional translational- or scaling-invariance property to achieve superior compression performance, including solution of quasi-static wave equations [65, 64, 22, 14], elliptic PDEs [3, 27], many-body Schrödinger equations [31], and quantum Fourier transforms (QFTs) [9]. That being said, most existing tensor decomposition algorithms will break down for OIOs due to their incapability to exploit the oscillatory structure of these operators; therefore, new tensor algorithms are called for.

In this paper, we propose a linear-complexity, low-prefactor tensor decomposition algorithm for large-scale and high-dimensional OIOs. This new tensor algorithm, henceforth dubbed the tensor butterfly algorithm, leverages the intrinsic CLR property of high-dimensional OIOs more effectively than the matrix butterfly algorithm, which is enabled by additional tensor properties such as translational invariance of free-space Green’s functions and dimensional separability of Fourier transforms. The algorithm partitions the OIO tensor into subtensors of multiple levels, and factorizes each subtensor at the middle level as a Tucker-like interpolative decomposition, whose factor matrices are further constructed in a nested fashion. For a d -dimensional OIO (assuming d constant) discretized as a $2d$ -mode tensor with n being the size per mode, the factorization time, application time, and storage cost scale as $O(n^d)$, and the resulting tensor factors have small multi-linear ranks. This is in stark contrast both to the $O(n^d \log n)$ scaling of existing matrix algorithms such as matrix butterfly algorithms and FFTs, and to the super-linear scaling of existing tensor algorithms. We mention that the linear complexity of the factorization time in our proposed algorithm is achieved via a simple random entry evaluation scheme, assuming that any arbitrary

entry can be computed in $O(1)$ time. We remark that, for 3D high-frequency wave equations, the proposed tensor butterfly algorithm can handle discretized Green's function tensors $512\times$ larger than existing algorithms; on the other hand, for the largest sized tensor that can be handled by existing algorithms, our tensor butterfly algorithm is $200\times$ faster than existing ones. Moreover, we claim that the tensor butterfly algorithm instantiates the first linear-complexity implementation of high-dimensional FFTs for arbitrary input data.

1.1. Related Work. *Multi-dimensional butterfly algorithms* represent a version of matrix butterfly algorithms designed for high-dimensional OIOs [38, 10]. Instead of the traditional binary tree partitioning of the matrix rows/columns [52], these algorithms can be viewed as a modern version of [53] that permits quadtree and octree partitioning of the matrix rows/columns, which have been demonstrated on 2D and 3D OIOs. For a general d -dimensional OIO, the d -dimensional tree partitioning leads to a butterfly factorization with a d -fold reduction in the number of levels compared to the binary tree partitioning. However, we note that both the multi-dimensional and binary tree-based butterfly algorithms are still matrix-based algorithms that scale as $O(n^d \log n)$, as opposed to the proposed tensor algorithm that scales as $O(n^d)$.

Quantized tensor train (QTT) algorithms, or simply TT algorithms, are tensor algorithms well-suited for very high-dimensional integral operators. They have been proposed to compress volume integral operators [14] arising from quasi-static wave equations and static PDEs with $O(\log n)$ memory and CPU complexities. However, for high-frequency wave equations, the QTT rank scales proportionally to the wave number [14] leading to deteriorated CPU and memory complexities (see our numerical results in Section 4). Moreover, QTT has been proposed for computing FFT and QFT with $O(\log n)$ memory and CPU complexities [9]. However, after obtaining the QTT-compressed formats of both the volume-integral operator and the Fourier transform, the CPU complexity for contracting such a QTT compressed operator with arbitrary (i.e., non QTT-compressed) input data scales super-linearly. In contrast, our algorithm yields a linear CPU and memory complexity for the contraction operation.

1.2. Contents. In what follows, we first review the matrix low-rank decomposition and butterfly decomposition algorithms in section 2. In subsection 3.1, we introduce the Tucker-like interpolative decomposition algorithm as the building block for the proposed tensor butterfly algorithm detailed in subsection 3.2. The multi-linear butterfly ranks for a few special cases are analyzed in subsection 3.2.1 and the complete complexity analysis is given in subsection 3.2.2. Section 4 shows a variety of numerical examples, including Green's functions for wave equations, Radon transforms, and uniform and non-uniform discrete Fourier transforms, to demonstrate the performance of matrix butterfly, tensor butterfly, Tucker and QTT algorithms.

1.3. Notations. Given a scalar-valued function $f(x)$, its integral transform is defined as

$$(1.1) \quad g(x) = \int_y K(x, y) f(y) dy$$

with an integral kernel $K(x, y)$. The indexing of a matrix \mathbf{K} is denoted by $\mathbf{K}(i, j)$ or $\mathbf{K}(t, s)$, where i, j are indices and t, s are index sets. We use \mathbf{K}^T to denote the transpose of matrix \mathbf{K} . For a sequence of matrices $\mathbf{K}_1, \dots, \mathbf{K}_n$, the matrix product

137 is

$$138 \quad (1.2) \quad \prod_{i=1}^n \mathbf{K}_i = \mathbf{K}_1 \mathbf{K}_2 \dots \mathbf{K}_n,$$

139 the vertical stacking (assuming the same column dimension) is

$$140 \quad (1.3) \quad [\mathbf{K}_i]_i = [\mathbf{K}_1; \mathbf{K}_2; \dots; \mathbf{K}_n],$$

141 and

$$142 \quad (1.4) \quad \text{diag}_i(\mathbf{K}_i) = \text{diag}(\mathbf{K}_1, \mathbf{K}_2, \dots, \mathbf{K}_n)$$

143 is a block diagonal matrix with \mathbf{K}_i being the diagonal blocks. Given an L -level binary-
 144 tree partitioning \mathcal{T}_t of an index set $t = \{1, 2, \dots, n\}$, any node τ at each level is a subset
 145 of t . The parent and children of τ are denoted by p_τ and τ^c ($c = 1, 2$), respectively,
 146 and $\tau = \tau^1 \cup \tau^2$.

147 A multi-index $\mathbf{i} = (i_1, \dots, i_d)$ is a tuple of indices, and similarly a multi-set
 148 $\boldsymbol{\tau} = (\tau_1, \tau_2, \dots, \tau_d)$ is a tuple of index sets. We define

$$149 \quad (1.5) \quad \tau_{k \leftarrow t} = (\tau_1, \tau_2, \dots, \tau_{k-1}, t, \tau_{k+1}, \tau_{k+2}, \dots, \tau_d).$$

150 Given a tuple of nodes (i.e. a multi-set) $\boldsymbol{\tau} = (\tau_1, \tau_2, \dots, \tau_d)$ and a multi-index $\mathbf{c} =$
 151 (c_1, c_2, \dots, c_d) with $c_i \in \{1, 2\}$, the children of $\boldsymbol{\tau}$ are denoted $\boldsymbol{\tau}^c = (\tau_1^{c_1}, \tau_2^{c_2}, \dots, \tau_d^{c_d})$
 152 and the parents of τ_i , $i = 1, 2, \dots, d$ can be simply written as $\mathbf{p}_\tau = (p_{\tau_1}, p_{\tau_2}, \dots, p_{\tau_d})$.
 153 Similar to the above-described notations, we can replace the index i in $[\mathbf{K}_i]_i$ and
 154 $\text{diag}_i(\mathbf{K}_i)$ with an index set τ , a multi-index \mathbf{c} , or a multi-set $\boldsymbol{\tau}$ assuming certain
 155 predefined index ordering.

156 Given complex-valued (or real-valued) functions $f(x)$ of d variables and integral
 157 operators $K(x, y)$, the tensor representations of their discretizations are respec-
 158 tively denoted by $\mathcal{F} \in \mathbb{C}^{n_1 \times n_2 \times \dots \times n_d}$ and $\mathcal{K} \in \mathbb{C}^{m_1 \times m_2 \times \dots \times m_d \times n_1 \times n_2 \times \dots \times n_d}$, where
 159 n_1, \dots, n_d and m_1, \dots, m_d are sizes of discretizations for the corresponding vari-
 160 ables. In this paper, we use *matricization* to denote the reshaping of \mathcal{K} into a
 161 $(\Pi_k m_k) \times (\Pi_k n_k)$ matrix, and the reshaping of \mathcal{F} into a $(\Pi_k n_k) \times 1$ matrix. The entries
 162 of \mathcal{F} and \mathcal{K} are denoted by $\mathcal{F}(\mathbf{i})$ (or equivalently $\mathcal{F}(i_1, i_2, \dots, i_d)$) and $\mathcal{K}(\mathbf{i}, \mathbf{j})$, respec-
 163 tively. Similarly the subtensors are denoted by $\mathcal{F}(\boldsymbol{\tau})$ (or equivalently $\mathcal{F}(\tau_1, \tau_2, \dots, \tau_d)$)
 164 and $\mathcal{K}(\boldsymbol{\tau}, \boldsymbol{\nu})$.

165 Given a d -mode tensor $\mathcal{F} \in \mathbb{C}^{n_1 \times n_2 \times \dots \times n_d}$, the mode- j unfolding is denoted by
 166 $\mathbf{F}^{(j)} \in \mathbb{C}^{(\Pi_{k \neq j} n_k) \times n_j}$, the mode- j tensor-matrix product of \mathcal{F} with a matrix $\mathbf{X} \in$
 167 $\mathbb{C}^{m \times n_j}$ is denoted by $\mathcal{Y} = \mathcal{F} \times_j \mathbf{X}$, or equivalently $\mathbf{Y}^{(j)} = \mathbf{F}^{(j)} \mathbf{X}^T$.

168 **2. Review of Matrix Algorithms.** We consider a d -dimensional OIO kernel
 169 $K(x, y)$ with $x, y \in \mathbb{R}^d$ discretized on point pairs x^i and y^j , $i = 1, 2, \dots, (m_1 m_2 \cdot$
 170 $\dots m_d)$, $j = 1, 2, \dots, (n_1 n_2 \cdot \dots n_d)$, where i (and similarly j) is the flattening of the
 171 corresponding multi-index \mathbf{i} . Such a discretization can be represented as a matrix
 172 $\mathbf{K} \in \mathbb{C}^{(m_1 m_2 \dots m_d) \times (n_1 n_2 \dots n_d)}$. When it is clear in the context, we assume that $m_k =$
 173 $n_k = n$ for $k = 1, \dots, d$. Throughout this paper, we assume that \mathbf{K} (and its tensor
 174 representation) is never fully formed, but instead a function is provided to evaluate any
 175 matrix (or tensor) entry in $O(1)$ time. Next we review matrix compression algorithms
 176 for \mathbf{K} including low-rank and butterfly algorithms.

2.1. Interpolative Decomposition. The interpolative decomposition (ID) algorithm [29, 39] is a matrix compression technique that constructs a low-rank decomposition whose factors contain original entries of the matrix. More specifically, consider the matrix $\mathbf{K}(\tau, \nu) \in \mathbb{C}^{m \times n}$ with $m \approx n$, $\tau = \{1, 2, \dots, m\}$, $\nu = \{1, 2, \dots, n\}$, the column ID of \mathbf{K} (the index sets τ and ν are omitted for clarity in context) is

$$(2.1) \quad \mathbf{K} \approx \mathbf{K}(:, \bar{\nu})\mathbf{V},$$

where the skeleton matrix $\mathbf{K}(:, \bar{\nu})$ contains r skeleton columns indexed by $\bar{\nu} \subseteq \nu$ and the interpolation matrix \mathbf{V} has bounded entries. Here the numerical rank r is chosen such that

$$(2.2) \quad \|\mathbf{K} - \mathbf{K}(:, \bar{\nu})\mathbf{V}\|_F^2 \leq O(\epsilon^2)\|\mathbf{K}\|_F^2$$

for a prescribed relative tolerance ϵ . In practice, the column ID can be computed via rank-revealing QR decomposition with a relative tolerance ϵ [39]. Similarly, the row ID of the matrix \mathbf{K} reads

$$(2.3) \quad \mathbf{K} \approx \mathbf{U}\mathbf{K}(\bar{\tau}, :),$$

where the skeleton matrix $\mathbf{K}(\bar{\tau}, :)$ contains r skeleton rows indexed by $\bar{\tau} \subseteq \tau$ and the interpolation matrix \mathbf{U} has bounded entries. The row ID can be simply computed by the column ID of \mathbf{K}^T . Combining the column and row ID in (2.1) and (2.3) gives

$$(2.4) \quad \mathbf{K} \approx \mathbf{U}\mathbf{K}(\bar{\tau}, \bar{\nu})\mathbf{V}.$$

It is straightforward to note that the memory and CPU complexities of ID scale as $O(nr)$ and $O(n^2r)$, respectively. The CPU complexity can be reduced to $O(nr^2)$ when properly selected proxy rows in (2.1) and columns in (2.3) are used in the rank-revealing QR. Common strategies of choosing proxy rows/columns (henceforth called *proxy index* strategies) for integral operators include evenly spaced or uniform random samples, and more generally the use of Chebyshev nodes and proxy surfaces (where new rows $K(\bar{x}, y^j)$ other than original rows of \mathbf{K} are used with \bar{x} denoting the proxies). However, for large OIOs, the rank r depends on the size n of the matrix; consequently, ID is not an efficient compression algorithm. Next, we review the matrix butterfly algorithm capable of achieving quasi-linear memory and CPU complexities for OIOs.

2.2. Matrix Butterfly Algorithm. Letting $t^0 = \{1, 2, \dots, m\}$, $s^0 = \{1, 2, \dots, n\}$, and $m = n$, the L -level butterfly representation of the discretized OIO $\mathbf{K}(t^0, s^0)$ is based on two binary trees, \mathcal{T}_{t^0} and \mathcal{T}_{s^0} , and the CLR property of the OIO takes the following form: at any level $0 \leq l \leq L$, for any node τ at level l of \mathcal{T}_{t^0} and any node ν at level $L - l$ of \mathcal{T}_{s^0} , the subblock $\mathbf{K}(\tau, \nu)$ is numerically low-rank with rank $r_{\tau, \nu}$ bounded by a small number r called the butterfly rank [46, 36, 37, 56].

For any subblock $\mathbf{K}(\tau, \nu)$, CLR permits a low-rank representation using ID in (2.4) as

$$(2.5) \quad \mathbf{K}(\tau, \nu) \approx \mathbf{U}_{\tau, \nu}\mathbf{K}(\bar{\tau}, \bar{\nu})\mathbf{V}_{\tau, \nu},$$

where the skeleton rows and columns are indexed by $\bar{\tau}$ and $\bar{\nu}$, respectively. It is worth noting that given a node ν , the selection of skeleton columns $\bar{\nu}$ depends on the node τ . However, the notation $\bar{\cdot}$ does not reflect the dependency when it is clear in the context.

Without loss of generality, we assume that L is an even number so that $L^c = L/2$ denotes the middle level. At levels $l = 0, \dots, L^c$, the interpolation matrices $\mathbf{V}_{\tau, \nu}$ are computed as follows:

At level $l = 0$, $\mathbf{V}_{\tau, \nu}$ are explicitly formed. While at level $0 < l \leq L^c$, they are represented in a nested fashion. To see this, consider a node pair (τ, ν) at level $l > 0$ and let ν^1, ν^2 and p_τ be the children and parent of ν and τ , respectively. Let s be the ancestor of ν at level L^c of \mathcal{T}_{s^0} and let \mathcal{T}_s denote the subtree rooted at s .

By CLR, we have

$$\begin{aligned} \mathbf{K}(\tau, \nu) &= [\mathbf{K}(\tau, \nu^1) \quad \mathbf{K}(\tau, \nu^2)] \\ &\approx [\mathbf{K}(\tau, \bar{\nu}^1) \quad \mathbf{K}(\tau, \bar{\nu}^2)] \begin{bmatrix} \mathbf{V}_{p_\tau, \nu^1}^s & \\ & \mathbf{V}_{p_\tau, \nu^2}^s \end{bmatrix} \\ &\approx \mathbf{K}(\tau, \bar{\nu}) \mathbf{W}_{\tau, \nu}^s \begin{bmatrix} \mathbf{V}_{p_\tau, \nu^1}^s & \\ & \mathbf{V}_{p_\tau, \nu^2}^s \end{bmatrix}. \end{aligned}$$

Here $\mathbf{W}_{\tau, \nu}^s$ and $\bar{\nu}$ are the interpolation matrix and skeleton columns from the ID of $\mathbf{K}(\tau, \bar{\nu}^1 \cup \bar{\nu}^2)$, respectively. $\mathbf{W}_{\tau, \nu}$ is henceforth referred to as the transfer matrix for ν in the rest of this paper. Note that we have added an additional superscript s to $\mathbf{V}_{p_\tau, \nu^c}$ and $\mathbf{W}_{\tau, \nu}$, for notation convenience in the later context. From (2.6), it is clear that the interpolation matrix $\mathbf{V}_{\tau, \nu}^s$ can be expressed in terms of its parent p_τ 's and children ν^1, ν^2 's interpolation matrices as

$$\mathbf{V}_{\tau, \nu}^s = \mathbf{W}_{\tau, \nu}^s \begin{bmatrix} \mathbf{V}_{p_\tau, \nu^1}^s & \\ & \mathbf{V}_{p_\tau, \nu^2}^s \end{bmatrix}.$$

Note that the interpolation matrices $\mathbf{V}_{\tau, \nu}^s$ at level $l = 0$ and transfer matrices $\mathbf{W}_{\tau, \nu}^s$ at level $0 < l \leq L^c$ do not require the column ID on the full subblocks $\mathbf{K}(\tau, \nu)$ and $\mathbf{K}(\tau, \bar{\nu}^1 \cup \bar{\nu}^2)$, which would lead to at least an $O(mn)$ computational complexity.

In practice, one can select $O(r_{\tau, \nu})$ proxy rows $\hat{\tau} \subset \tau$ to compute $\mathbf{V}_{\tau, \nu}^s$ and $\mathbf{W}_{\tau, \nu}^s$ via ID as:

$$\mathbf{K}(\hat{\tau}, \nu) \approx \mathbf{K}(\hat{\tau}, \bar{\nu}) \mathbf{V}_{\tau, \nu}^s, \quad l = 0,$$

$$\mathbf{K}(\hat{\tau}, \bar{\nu}^1 \cup \bar{\nu}^2) \approx \mathbf{K}(\hat{\tau}, \bar{\nu}) \mathbf{W}_{\tau, \nu}^s, \quad 0 < l \leq L^c.$$

The viable choices for proxy rows have been discussed in several existing papers [45, 56, 59, 8].

At levels $l = L^c, \dots, L$, the interpolation matrices $\mathbf{U}_{\tau, \nu}$ are computed by performing similar operations on \mathbf{K}^T . We only provide their expressions here and omit the redundant explanation. Let t be the ancestor of ν at level L^c of \mathcal{T}_{t^0} and let \mathcal{T}_t be the subtree rooted at t . At level $l = L$, $\mathbf{U}_{\tau, \nu}^t$ are explicitly formed. At level $L^c \leq l < L$, only the transfer matrices $\mathbf{P}_{\tau, \nu}^t$ are computed from the column ID of $\mathbf{K}^T(\nu, \bar{\tau}^1 \cup \bar{\tau}^2)$ satisfying

$$\mathbf{U}_{\tau, \nu}^t = \begin{bmatrix} \mathbf{U}_{\tau^1, p_\nu}^t & \\ & \mathbf{U}_{\tau^2, p_\nu}^t \end{bmatrix} \mathbf{P}_{\tau, \nu}^t.$$

Combining (2.5), (2.8) and (2.11), the matrix butterfly decomposition can be expressed for each node pair (t, s) at level L^c of \mathcal{T}_{t^0} and \mathcal{T}_{s^0} as

$$\mathbf{K}(t, s) \approx \bar{\mathbf{U}}^t \left(\prod_{l=1}^{L^c} \bar{\mathbf{P}}_l^{t, s} \right) \mathbf{K}(\bar{t}, \bar{s}) \left(\prod_{l=L^c}^1 \bar{\mathbf{W}}_l^{t, s} \right) \bar{\mathbf{V}}^s.$$

Here, \bar{t} and \bar{s} represent the skeleton rows and columns of the ID of $\mathbf{K}(t, s)$. The interpolation factors $\bar{\mathbf{U}}^t$ and $\bar{\mathbf{V}}^s$ in (2.12) are

$$(2.13) \quad \bar{\mathbf{U}}^t = \text{diag}_\tau(\mathbf{U}_{\tau, s^0}^t), \quad \tau \text{ at level } L^c \text{ of } \mathcal{T}_t,$$

$$(2.14) \quad \bar{\mathbf{V}}^s = \text{diag}_\nu(\mathbf{V}_{t^0, \nu}^s), \quad \nu \text{ at level } L^c \text{ of } \mathcal{T}_s,$$

and the transfer factors $\bar{\mathbf{P}}_l^{t,s}$ and $\bar{\mathbf{W}}_l^{t,s}$ for $l = 1, \dots, L^c$ consist of transfer matrices $\mathbf{W}_{\tau, \nu}^s$ and $\mathbf{P}_{\tau, \nu}^s$:

$$(2.15) \quad \bar{\mathbf{W}}_l^{t,s} = \text{diag}_\tau \left(\left[\text{diag}_\nu(\mathbf{W}_{\tau^c, \nu^c}^s) \right]_c \right), \quad \begin{array}{l} \tau \text{ at level } l-1 \text{ of } \mathcal{T}_{t^0}, \text{ and } t \subseteq \tau, \\ \nu \text{ at level } L^c - l \text{ of } \mathcal{T}_s; \end{array}$$

$$(2.16) \quad (\bar{\mathbf{P}}_l^{t,s})^T = \text{diag}_\nu \left(\left[\text{diag}_\tau((\mathbf{P}_{\tau, \nu^c}^t)^T) \right]_c \right), \quad \begin{array}{l} \tau \text{ at level } L^c - l \text{ of } \mathcal{T}_t, \\ \nu \text{ at level } l-1 \text{ of } \mathcal{T}_{s^0}, \text{ and } s \subseteq \nu. \end{array}$$

Here τ^c and ν^c with $c = 1, 2$ are children of τ and ν , respectively.

The CPU and memory requirement for computing the matrix butterfly decomposition can be briefly analyzed as follows. Note that we only need to analyze the costs for $\mathbf{V}_{\tau, \nu}^s$, $\mathbf{W}_{\tau, \nu}^s$ and $\mathbf{K}(\bar{t}, \bar{s})$ as those for $\mathbf{U}_{\tau, \nu}^t$ and $\mathbf{P}_{\tau, \nu}^t$ are similar. By CLR assumption, we assume that $r_{\tau, \nu} \leq r, \forall \tau, \nu$ for some constant r . Thanks to the use of the proxy rows and columns, the computation of one individual $\mathbf{V}_{\tau, \nu}^s$ and $\mathbf{W}_{\tau, \nu}^s$ by ID only operates on $O(r) \times O(r)$ matrices, hence its memory and CPU requirements are $O(r^2)$ and $O(r^3)$, respectively. In total, there are $O(2^{L^c})$ middle-level nodes s each having $O(2^{L^c})$ numbers of $\mathbf{V}_{\tau, \nu}^s$ and $O(L^c 2^{L^c})$ numbers of $\mathbf{W}_{\tau, \nu}^s$. Similarly, each $\mathbf{K}(\bar{t}, \bar{s})$ requires $O(r^2)$ CPU and memory costs, and there are in total $O(2^L)$ middle-level node pairs (t, s) . These numbers sum up to the overall $O(nr^2 \log n)$ memory and $O(nr^3 \log n)$ CPU complexities for matrix butterfly algorithms.

For d -dimensional discretized OIOs $\mathbf{K} \in \mathbb{C}^{(m_1 m_2 \dots m_d) \times (n_1 n_2 \dots n_d)}$ with $m_k = n_k = n$, we can assume that $n = C_b 2^L$ with some constant C_b . For the above-described binary-tree-based butterfly algorithm, the leaf nodes of the trees are of size C_b^d and this leads to a dL -level butterfly factorization. The memory and CPU complexities for this algorithm become $O(dn^d r^2 \log n)$ and $O(dn^d r^3 \log n)$, respectively. On the other hand, the multi-dimensional tree-based butterfly algorithm [38, 10] leads to a L -level factorization with $O(2^d n^d r^2 \log n)$ memory and $O(2^d n^d r^3 \log n)$ CPU complexities. Despite their quasi-linear complexity for high-dimensional OIOs, the butterfly rank r is constant but high, leading to very large prefactors of these binary and multi-dimensional tree-based algorithms. In the following, we turn to tensor decomposition algorithms to reduce both the prefactor and asymptotic scaling of matrix butterfly algorithms.

3. Proposed Tensor Algorithms. In this section, we assume that the d -dimensional discretized OIO in section 2 is directly represented as a $2d$ -mode tensor $\mathcal{K} \in \mathbb{C}^{m_1 \times m_2 \times \dots \times m_d \times n_1 \times n_2 \times \dots \times n_d}$. We first extend the matrix ID algorithm in subsection 2.1 to its tensor variant, which serves as the building block for the proposed tensor butterfly algorithm.

3.1. Tucker-like Interpolative Decomposition. Given the $2d$ -mode tensor $\mathcal{K}(\tau, \nu)$ with $\tau_k = \{1, 2, \dots, m_k\}$ and $\nu_k = \{1, 2, \dots, n_k\}$ for $k = 1, \dots, d$, the proposed tensor ID decomposition compresses each dimension independently via the column ID of the unfolding of \mathcal{K} along the k -th dimension,

$$(3.1) \quad \mathbf{K}^{(k)} = \mathbf{K}^{(k)}(:, \bar{\tau}_k) \mathbf{U}^k, \quad \mathbf{K}^{(d+k)} = \mathbf{K}^{(d+k)}(:, \bar{\nu}_k) \mathbf{V}^k, \quad k = 1, \dots, d,$$

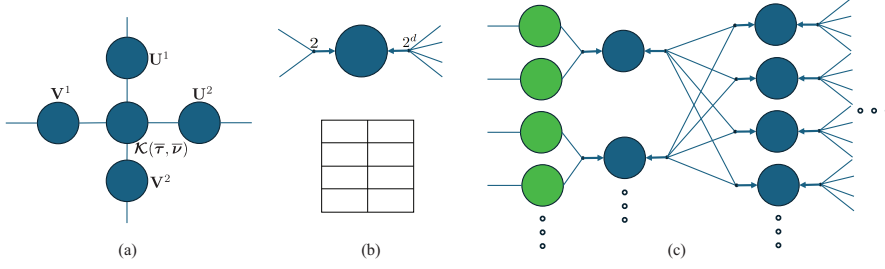


Fig. 3.1: Tensor diagrams for (a) the Tucker-ID decomposition of a 4-mode tensor, and (b) the matrix partitioner corresponding to a $2^d \times 2$ partitioning with $d = 2$ used in the tensor butterfly decomposition of a $2d$ -mode tensor, such as $\left[\mathbf{W}_{\tau^c, \nu}^{s, k}\right]_{\mathbf{c}}$ in (3.12) for fixed \mathbf{s}, τ, k and ν , or $\left[\mathbf{P}_{\tau, \nu^c}^{t, k}\right]_{\mathbf{c}}$ in (3.11) for fixed \mathbf{t}, ν, k and τ . (c) The tensor diagram involving blocks $\mathbf{V}_{t^0, \nu}^{s, k}$ (in green) and blocks $\left[\mathbf{W}_{\tau^c, \nu}^{s, k}\right]_{\mathbf{c}}$ (in blue) for fixed \mathbf{s} and k for the tensor butterfly decomposition of a $2d$ -mode tensor.

where $\mathbf{K}^{(k)} \in \mathbb{C}^{(\prod_{j \neq k} n_j) \times n_k}$ is the mode- k unfolding, or equivalently

$$(3.2) \quad \mathbf{K} = \mathbf{K}(\tau_{k \leftarrow \bar{\tau}_k}, \nu) \times_k \mathbf{U}^k, \quad \mathbf{K} = \mathbf{K}(\tau, \nu_{k \leftarrow \bar{\nu}_k}) \times_{d+k} \mathbf{V}^k, \quad k = 1, \dots, d.$$

Here, $\bar{\tau}_k$ and $\bar{\nu}_k$ denote the skeleton indices along modes k and $d+k$ of \mathbf{K} , respectively, while $\tau_{k \leftarrow \bar{\tau}_k}$ and $\nu_{k \leftarrow \bar{\nu}_k}$ denote multi-sets that replace τ_k and ν_k , respectively, with $\bar{\tau}_k$ and $\bar{\nu}_k$. Combining (3.2) for all dimensions yields the following proposed tensor interpolation decomposition,

$$(3.3) \quad \mathbf{K} = \mathbf{K}(\bar{\tau}, \bar{\nu}) \left(\prod_{k=1}^d \times_k \mathbf{U}^k \right) \left(\prod_{k=1}^d \times_{d+k} \mathbf{V}^k \right),$$

where, $\bar{\tau} = (\bar{\tau}_1, \bar{\tau}_2, \dots, \bar{\tau}_d)$, $\bar{\nu} = (\bar{\nu}_1, \bar{\nu}_2, \dots, \bar{\nu}_d)$, the core tensor $\mathbf{K}(\bar{\tau}, \bar{\nu})$ is a subtensor of \mathbf{K} , and \mathbf{U}^k and \mathbf{V}^k are the factor matrices for modes k and $d+k$, respectively.

Note that the tensor diagram of (3.3) is exactly the same as Tucker decompositions or high-order singular value decompositions (HOSVD) [16]. Both decompositions provide a canonical form of “core and factor product” tensor approximation. See Figure 3.1(a) for the tensor diagram of (3.3) for a 4-mode tensor. Unlike the Tucker decomposition that leads to orthonormal factor matrices, the proposed decomposition leads to factor matrices with bounded entries and the core tensor with the original tensor entries. Therefore, the proposed decomposition is named Tucker-like interpolative decomposition (Tucker ID). It is worth noting that there exist several interpolative tensor decomposition algorithms [6, 50, 51, 58]. However they either use original tensor entries in the factor matrices (instead of the core tensor) [50, 58, 6] or rely on a different tensor diagram [51]. As will be seen in subsection 3.2, the Tucker ID algorithm is a unique and essential building block of the tensor butterfly algorithm.

The memory and CPU complexities of Tucker-ID can be briefly analyzed as follows. Assuming that $m_k = n_k = n$ and $\max_k |\bar{\tau}_k| = \max_k |\bar{\nu}_k| = r$ is a constant (we will discuss the case of non-constant r in subsection 3.2.3), the memory requirement is simply $O(drn + r^{2d})$, where the first and second term represents the storage units for the factor matrices and the core tensor, respectively. The CPU cost for naive computation of Tucker-ID is $O(drn^{2d} + r^{2d})$, where the first term represents the cost

of rank-revealing QR of the unfolding matrices in (3.1), and the second term represents the cost forming the core tensor $\mathcal{K}(\bar{\tau}, \bar{\nu})$. In practice, however, the unfolding matrices do not need to be fully formed and one can leverage the idea of proxy rows in subsection 2.2 to reduce the cost for computing the factor matrices to $O(dnr^{2d})$. We will explain this in more detail in the context of the proposed tensor butterfly decomposition algorithm.

Just like the matrix ID algorithm, Tucker-ID is also not suitable for representing large-sized OIOs as the rank r depends on the size n . That said, the Tucker-ID rank is typically significantly smaller than the matrix ID rank, as it exploits more compressibility properties across dimensions by leveraging e.g. translational-invariance or dimensional-separability properties of OIOs; see subsection 3.2.1 for a few of such examples. In what follows, we use Tucker-ID as the building block for constructing a linear-complexity tensor butterfly decomposition algorithm for large-sized OIOs.

3.2. Tensor Butterfly Algorithm. Consider a 2d-mode OIO tensor $\mathcal{K}(\mathbf{t}^0, \mathbf{s}^0)$ with $\mathbf{t}^0 = (t_1^0, t_1^0, \dots, t_d^0)$, $\mathbf{s}^0 = (s_1^0, s_1^0, \dots, s_d^0)$, $t_k^0 = \{1, 2, \dots, m_k\}$, $s_k^0 = \{1, 2, \dots, n_k\}$, $k = 1, 2, \dots, d$. Without loss of generality, we assume that $m_k = n_k = n$. We further assume that each t_k^0 (and s_k^0) is binary partitioned with a tree $\mathcal{T}_{t_k^0}$ (and $\mathcal{T}_{s_k^0}$) of L levels for $k = 1, 2, \dots, d$.

To start with, we first define the tensor CLR property as follows:

- For any level $0 \leq l \leq L^c$, any multi-set $\tau = (\tau_1, \tau_2, \dots, \tau_d)$ with $\tau_i, i \leq d$ at level l of $\mathcal{T}_{t_i^0}$, any multi-set $\mathbf{s} = (s_1, s_2, \dots, s_d)$ with $s_i, i \leq d$ at level L^c of $\mathcal{T}_{s_i^0}$, any mode $1 \leq k \leq d$, and any node ν at level $L^c - l$ of \mathcal{T}_{s_k} , the mode- $(d+k)$ unfolding of the subtensor $\mathcal{K}(\tau, \mathbf{s}_{k \leftarrow \nu})$ is numerically low-rank (with rank bounded by r), permitting an ID via (3.2):

$$(3.4) \quad \mathcal{K}(\tau, \mathbf{s}_{k \leftarrow \nu}) \approx \mathcal{K}(\tau, \mathbf{s}_{k \leftarrow \bar{\nu}}) \times_{d+k} \mathbf{V}_{\tau, \nu}^{\mathbf{s}, k}.$$

- For any level $0 \leq l \leq L^c$, any multi-set $\nu = (\nu_1, \nu_2, \dots, \nu_d)$ with $\nu_i, i \leq d$ at level l of $\mathcal{T}_{s_i^0}$, any multi-set $\mathbf{t} = (t_1, t_2, \dots, t_d)$ with $t_i, i \leq d$ at level L^c of $\mathcal{T}_{t_i^0}$, any mode $1 \leq k \leq d$, and any node τ at level $L^c - l$ of \mathcal{T}_{t_k} , the mode- k unfolding of the subtensor $\mathcal{K}(\mathbf{t}_{k \leftarrow \tau}, \nu)$ is numerically low-rank (with rank bounded by r), permitting an ID via (3.2):

$$(3.5) \quad \mathcal{K}(\mathbf{t}_{k \leftarrow \tau}, \nu) \approx \mathcal{K}(\mathbf{t}_{k \leftarrow \bar{\tau}}, \nu) \times_k \mathbf{U}_{\tau, \nu}^{\mathbf{t}, k}.$$

In essence, the tensor CLR in (3.4) and (3.5) investigates the unfolding of judiciously selected subtensors rather than the matricization used in the matrix CLR. Moreover, the tensor CLR requires fixing $d - 1$ modes of the $2d$ -mode subtensors to be of size $O(\sqrt{n})$ while changing the remaining $d + 1$ modes with respect to l . Therefore each ID computation can operate on larger subtensors compared to the matrix CLR. In subsection 3.2.1 we provide two examples, namely a free-space Green's function tensor and a high-dimensional Fourier transform, to explain why the tensor CLR is valid, and in subsection 3.2.2 we will see that the tensor CLR essentially reduces the quasilinear complexity of the matrix butterfly algorithm to linear complexity. Here, assuming that the tensor CLR holds true, we describe the tensor butterfly algorithm. We note that there may be alternative ways to define the tensor CLR different from (3.4) and (3.5), and we leave that as a future work.

In what follows, we focus on the computation of $\mathbf{V}_{\tau, \nu}^{\mathbf{s}, k}$ (corresponding to the mid-level multi-set \mathbf{s}), as $\mathbf{U}_{\tau, \nu}^{\mathbf{t}, k}$ (corresponding to the mid-level multi-set \mathbf{t}) can be computed in a similar fashion. At level $l = 0$, $\mathbf{V}_{\tau, \nu}^{\mathbf{s}, k}$ are explicitly formed. At level

0 < l ≤ L^c, they are represented in a nested fashion. Let $\mathbf{p}_\tau = (p_{\tau_1}, p_{\tau_2}, \dots, p_{\tau_d})$ consist of parents of $\tau = (\tau_1, \tau_2, \dots, \tau_d)$ in (3.4).

By the tensor CLR property, we have

$$\begin{aligned} \mathcal{K}(\tau, \mathbf{s}_{k \leftarrow \nu}) &\approx \mathcal{K}(\tau, \mathbf{s}_{k \leftarrow \overline{\nu^1 \cup \nu^2}}) \times_{d+k} \begin{bmatrix} \mathbf{V}_{\mathbf{p}_\tau, \nu^1}^{s,k} & \\ & \mathbf{V}_{\mathbf{p}_\tau, \nu^2}^{s,k} \end{bmatrix} \\ &\approx \mathcal{K}(\tau, \mathbf{s}_{k \leftarrow \overline{\nu}}) \times_{d+k} \left(\mathbf{W}_{\tau, \nu}^{s,k} \begin{bmatrix} \mathbf{V}_{\mathbf{p}_\tau, \nu^1}^{s,k} & \\ & \mathbf{V}_{\mathbf{p}_\tau, \nu^2}^{s,k} \end{bmatrix} \right). \end{aligned}$$

Comparing (3.6) and (3.4), one realizes that the interpolation matrix $\mathbf{V}_{\tau, \nu}^{s,k}$ is represented as the product of the transfer matrix $\mathbf{W}_{\tau, \nu}^{s,k}$ and $\text{diag}_c(\mathbf{V}_{\mathbf{p}_\tau, \nu^c}^{s,k})$. Here, the transfer matrix $\mathbf{W}_{\tau, \nu}^{s,k}$ is computed as the interpolation matrix of the column ID of the mode-(d+k) unfolding of $\mathcal{K}(\tau, \mathbf{s}_{k \leftarrow \overline{\nu^1 \cup \nu^2}})$. As mentioned in section 3, in practice one never forms the unfolding matrix in full, but instead considers the unfolding of $\mathcal{K}(\hat{\tau}, \hat{\mathbf{s}}_{k \leftarrow \overline{\nu^1 \cup \nu^2}})$, where $\hat{\tau} = (\hat{\tau}_1, \hat{\tau}_2, \dots, \hat{\tau}_d)$ and $\hat{\mathbf{s}} = (\hat{s}_1, \hat{s}_2, \dots, \hat{s}_d)$; here $\hat{\tau}_i$ and \hat{s}_i consist of $O(r)$ judiciously selected indices along modes i and $d+i$, respectively. Note that \hat{s}_k is never used as it is replaced by $\overline{\nu^1 \cup \nu^2}$ in (3.6). The same proxy index strategy can be used to obtain $\mathbf{V}_{\tau, \nu}^{s,k}$ at the level $l = 0$. For each $\mathbf{W}_{\tau, \nu}^{s,k}$ or $\mathbf{V}_{\tau, \nu}^{s,k}$, its computation requires $O(r^{2d+1})$ CPU time.

Similarly in (3.5), $\mathbf{U}_{\tau, \nu}^{t,k}$ is explicitly formed at $l = 0$ and constructed via the transfer matrix $\mathbf{P}_{\tau, \nu}^{t,k}$ at level $0 < l \leq L^c$:

$$\begin{aligned} \mathcal{K}(t_{k \leftarrow \tau}, \nu) &\approx \mathcal{K}(t_{k \leftarrow \overline{\tau^1 \cup \tau^2}}, \nu) \times_k \begin{bmatrix} \mathbf{U}_{\tau^1, \mathbf{p}_\nu}^{t,k} & \\ & \mathbf{U}_{\tau^2, \mathbf{p}_\nu}^{t,k} \end{bmatrix} \\ &\approx \mathcal{K}(t_{k \leftarrow \overline{\tau}}, \nu) \times_k \left(\mathbf{P}_{\tau, \nu}^{t,k} \begin{bmatrix} \mathbf{U}_{\tau^1, \mathbf{p}_\nu}^{t,k} & \\ & \mathbf{U}_{\tau^2, \mathbf{p}_\nu}^{t,k} \end{bmatrix} \right). \end{aligned}$$

Putting together (3.4), (3.5), (3.6) and (3.7), the proposed tensor butterfly decomposition can be expressed, for any multi-set $\mathbf{t} = (t_1, t_2, \dots, t_d)$ with t_i at level L^c of \mathcal{T}_i^0 and any multi-set $\mathbf{s} = (s_1, s_2, \dots, s_d)$ with s_i at level L^c of $\mathcal{T}_{s_i}^0$, by forming a Tucker-ID for the (\mathbf{t}, \mathbf{s}) pair:

$$\mathcal{K}(\mathbf{t}, \mathbf{s}) \approx \mathcal{K}(\bar{\mathbf{t}}, \bar{\mathbf{s}}) \left(\prod_{k=1}^d \times_k \left(\prod_{l=L^c}^1 \overline{\mathbf{P}}_l^{t, s, k} \overline{\mathbf{U}}^{t, k} \right) \right) \left(\prod_{k=1}^d \times_{d+k} \left(\prod_{l=L^c}^1 \overline{\mathbf{W}}_l^{t, s, k} \overline{\mathbf{V}}^{s, k} \right) \right).$$

Here, $\bar{\mathbf{t}}$ and $\bar{\mathbf{s}}$ represent the skeleton indices of the Tucker-ID of $\mathcal{K}(\mathbf{t}, \mathbf{s})$. The interpolation factors $\overline{\mathbf{U}}^{t, k}$ and $\overline{\mathbf{V}}^{s, k}$ in (3.8) are:

$$\overline{\mathbf{U}}^{t, k} = \text{diag}_\tau(\mathbf{U}_{\tau, \mathbf{s}^0}^{t, k}), \quad \tau \text{ at level } L^c \text{ of } \mathcal{T}_{t_k},$$

$$\overline{\mathbf{V}}^{s, k} = \text{diag}_\nu(\mathbf{V}_{\mathbf{t}^0, \nu}^{s, k}), \quad \nu \text{ at level } L^c \text{ of } \mathcal{T}_{s_k},$$

and the transfer factors $\overline{\mathbf{P}}_l^{t, s, k}$ and $\overline{\mathbf{W}}_l^{t, s, k}$ for $l = 1, \dots, L^c$ are:

$$\overline{\mathbf{P}}_l^{t, s, k} = \text{diag}_\nu \left(\left[\text{diag}_\tau(\mathbf{P}_{\tau, \nu^c}^{t, k}) \right]_c \right), \quad \begin{array}{l} \tau \text{ at level } L^c - l \text{ of } \mathcal{T}_{t_k}, \\ \nu_i \text{ at level } l - 1 \text{ of } \mathcal{T}_{s_i^0}, s_i \subseteq \nu_i, i \leq d; \end{array}$$

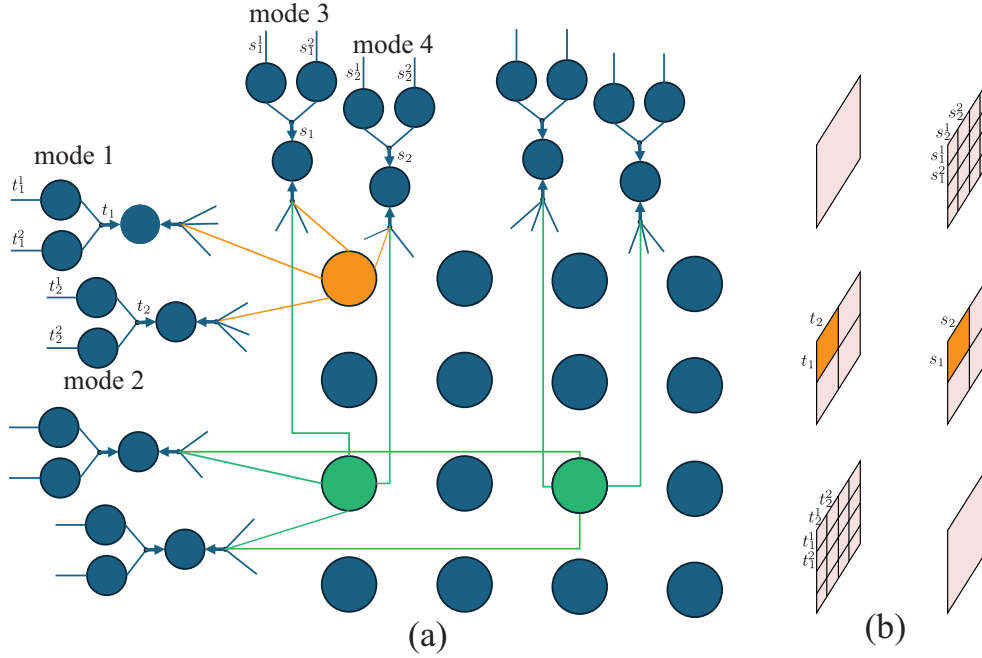


Fig. 3.2: (a) Tensor diagram for the tensor butterfly decomposition of $L = 2$ levels of a 4-mode OIO tensor representing (b) high-frequency Green's function interactions between parallel facing 2D unit squares. Only the full connectivity regarding three middle-level node pairs is shown (the two green circles and one orange circle in (a)). The orange circle in (a) represents the core tensor $\mathcal{K}(\bar{\mathbf{t}}, \bar{\mathbf{s}})$ for a mid-level pair (\mathbf{t}, \mathbf{s}) with $\mathbf{t} = (t_1, t_2)$, $\mathbf{s} = (s_1, s_2)$ highlighted in orange in (b).

$$(3.12) \quad \bar{\mathbf{W}}_l^{\mathbf{t}, \mathbf{s}, k} = \text{diag}_{\boldsymbol{\tau}} \left(\left[\text{diag}_{\nu} (\mathbf{W}_{\boldsymbol{\tau}^c, \nu}^{\mathbf{s}, k}) \right]_c \right), \quad \begin{array}{l} \tau_i \text{ at level } l-1 \text{ of } \mathcal{T}_{t_i^0}, t_i \subseteq \tau_i, i \leq d, \\ \nu \text{ at level } L^c - l \text{ of } \mathcal{T}_{s_k}. \end{array}$$

One can verify that when $d = 1$, the tensor butterfly algorithm (3.8) reduces to the matrix butterfly algorithm (2.12). But when $d > 1$, the tensor butterfly algorithm has a distinct algorithmic structure and the computational complexity can be significantly reduced compared with the matrix butterfly algorithm. Detailed computational complexity analysis is provided in subsection 3.2.2.

To better understand the structure of the tensor butterfly in (3.8), (3.9), (3.10), (3.11), and (3.12), we describe its tensor diagram here. We first create the tensor diagram for a *matrix partitioner* as shown in Figure 3.1(b), which represents a $2^d \times 2$ block partitioning of a matrix such as $\left[\mathbf{W}_{\boldsymbol{\tau}^c, \nu}^{\mathbf{s}, k} \right]_c$ in (3.12) for fixed $\mathbf{s}, \boldsymbol{\tau}, k$ and ν , or $\left[\mathbf{P}_{\boldsymbol{\tau}, \nu^c}^{\mathbf{t}, k} \right]_c$ in (3.11) for fixed $\mathbf{t}, \boldsymbol{\nu}, k$ and $\boldsymbol{\tau}$. In other words, there are 2 legs on the column dimension and 2^d legs on the row dimension. The diagram in Figure 3.1(c) shows the connectivity for all $\mathbf{V}_{\mathbf{t}^0, \nu}^{\mathbf{s}, k}$ (the green circles) and $\left[\mathbf{W}_{\boldsymbol{\tau}^c, \nu}^{\mathbf{s}, k} \right]_c$ (the blue circles) for fixed \mathbf{s} and k . The multiplication or contraction of all matrices in Figure 3.1(c) results in $\mathbf{V}_{\mathbf{t}, s_k}^{\mathbf{s}, k}$ for all mid-level multi-sets \mathbf{t} , which are of course not explicitly formed.

As an example, consider an OIO representing the free-space Green's function in-

416 teraction between two parallel facing unit square plates in Figure 3.2. The tensor is
 417 $\mathcal{K}(i, j) = K(x^i, y^j) = \frac{\exp(-i\omega\rho)}{\rho}$ where $x^i = (\frac{i_1}{n}, \frac{i_2}{n}, 0)$, $y^j = (\frac{j_1}{n}, \frac{j_2}{n}, 1)$, $\rho = |x^i - y^j|$
 418 and ω is the wavenumber. Here 1 represents the distance between the two plates.
 419 Consider an $L=2$ -level tensor butterfly decomposition, with a total of 16 middle-level
 420 multi-set pairs. Let (\mathbf{t}, \mathbf{s}) denote one middle-level multi-set pair with $\mathbf{t} = (t_1, t_2)$ and
 421 $\mathbf{s} = (s_1, s_2)$ as highlighted in orange in Figure 3.2(b). Their children are $t_1^1, t_2^1, t_1^2, t_2^2$
 422 and $s_1^1, s_2^1, s_1^2, s_2^2$. Leveraging the representations in Figure 3.1(b)-(c), the full di-
 423 agram for $\mathcal{K}(\mathbf{t}, \mathbf{s})$ consists of one 4-mode tensor $\mathcal{K}(\bar{\mathbf{t}}, \bar{\mathbf{s}})$ (highlighted in orange in
 424 Figure 3.2(a)), one transfer matrix per mode, and two factor matrices per mode. In
 425 addition, we plot the full connectivity for two other multi-set pairs (highlighted in
 426 Green in Figure 3.2(a)). It is important to note that the factor matrices and transfer
 427 matrices are shared among the multi-set pairs.

428 The proposed tensor butterfly algorithm is fully described in Algorithm 3.1 for a
 429 $2d$ -mode tensor $\mathcal{K} \in \mathbb{C}^{m_1 \times m_2 \times \dots \times m_d \times n_1 \times n_2 \times \dots \times n_d}$, which consists of three steps: (1)
 430 computation of $\mathbf{V}_{\tau, \nu}^{s, k}$ and $\mathbf{W}_{\tau, \nu}^{s, k}$ starting at Line 1, (2) computation of $\mathbf{U}_{\tau, \nu}^{t, k}$ and $\mathbf{P}_{\tau, \nu}^{t, k}$
 431 starting at Line 17, and (3) computation of $\mathcal{K}(\bar{\mathbf{t}}, \bar{\mathbf{s}})$ starting at Line 33. We note that,
 432 after each $\mathcal{K}(\bar{\mathbf{t}}, \bar{\mathbf{s}})$ is formed, we leverage floating-point compression tools such as the
 433 ZFP software [40] to further compress it.

434 Once \mathcal{K} is compressed, any input tensor $\mathcal{F} \in \mathbb{C}^{n_1 \times n_2 \times \dots \times n_d \times n_v}$ can contract with
 435 it to compute $\mathcal{G} = \mathcal{K} \times_{d+1, d+2, \dots, 2d} \mathcal{F}$. It is clear to see that the contraction is
 436 equivalent to matrix-matrix multiplication $\mathbf{G} = \mathbf{K}\mathbf{F}$, where $\mathbf{G} \in \mathbb{C}^{\prod_k m_k \times n_v}$, $\mathbf{K} \in$
 437 $\mathbb{C}^{\prod_k m_k \times \prod_k n_k}$, and $\mathbf{F} \in \mathbb{C}^{\prod_k n_k \times n_v}$ are matricizations of \mathcal{G} , \mathcal{K} and \mathcal{F} , respectively,
 438 and n_v is the number of columns of \mathbf{F} . The contraction algorithm is described in
 439 Algorithm 3.2 which consists of three steps:

- 440 (1) Contraction with $\mathbf{V}_{\tau, \nu}^{s, k}$ and $\mathbf{W}_{\tau, \nu}^{s, k}$. For each level $l = 0, 1, \dots, L^c$, one notices that,
 441 since the contraction operation for each multi-set τ with τ_i at level l of $\mathcal{T}_{t_i^0}$ and
 442 the middle-level multi-set \mathbf{s} is independent of each other, one needs a separate
 443 tensor $\mathcal{F}_{\tau, \mathbf{s}}$ to store the contraction result for each multi-set pair (τ, \mathbf{s}) . $\mathcal{F}_{\tau, \mathbf{s}}$
 444 can be computed by mode-by-mode contraction with the factor matrices $\bar{\mathbf{V}}^{s, k}$ for
 445 $l = 0$ (Line 6) and the transfer matrices $\text{diag}_{\nu}(\mathbf{W}_{\tau, \nu}^{s, k})$ for $l > 0$ (Line 8).
- 446 (2) Contraction with $\mathcal{K}(\bar{\mathbf{t}}, \bar{\mathbf{s}})$ at the middle level. Tensors at the middle level $\mathcal{F}_{\mathbf{t}, \mathbf{s}}$
 447 are contracted with each subtensor $\mathcal{K}(\bar{\mathbf{t}}, \bar{\mathbf{s}})$ separately, resulting in tensors $\mathcal{G}_{\mathbf{t}, \mathbf{s}} =$
 448 $\mathcal{K}(\bar{\mathbf{t}}, \bar{\mathbf{s}}) \times_{d+1, d+2, \dots, 2d} \mathcal{F}_{\mathbf{t}, \mathbf{s}}$.
- 449 (3) Contraction with $\mathbf{U}_{\tau, \nu}^{t, k}$ and $\mathbf{P}_{\tau, \nu}^{t, k}$. As Step (1), for each level $l = L^c, L^{c-1}, \dots, 0$,
 450 the contraction operation for each multi-set ν with ν_i at level l of $\mathcal{T}_{s_i^0}$ and middle-
 451 level multi-set \mathbf{t} is independent. At level $l > 0$, the contribution of tensors $\mathcal{G}_{\mathbf{t}, \nu}$ is
 452 accumulated into $\mathcal{G}_{\mathbf{t}, \mathbf{p}_{\nu}}$ (Line 26); at level $l = 0$, the contraction results are stored
 453 in the final output tensor $\mathcal{G}(\mathbf{t}, 1 : n_v)$ (Line 24).

454 **3.2.1. Rank Estimate.** In this subsection, we use two specific high-dimensional
 455 examples, namely high-frequency free-space Green's functions for wave equations and
 456 uniform discrete Fourier transforms (DFTs) to investigate the matrix and tensor CLR
 457 properties, and compare the matrix and tensor butterfly ranks r_m and r_t , respectively.
 458 For the Green's function example, the tensor CLR property is a result of matrix CLR
 459 and translational invariance, and r_t is much smaller than r_m ; for the DFT example, the
 460 tensor CLR property is a result of matrix CLR and dimensionality separability, and
 461 r_t is exactly the same as r_m of 1D DFTs. For more-general OIOs, such as analytical
 462 and numerical Green's functions for inhomogeneous media, Radon transforms, non-
 463 uniform DFTs, and general Fourier integral operators, rigorous rank analysis is non-
 464 trivial and we rely on numerical experiments in section 4 to demonstrate the efficacy

Algorithm 3.1 Construction algorithm for the tensor butterfly decomposition of a $2d$ -mode tensor $\mathcal{K} \in \mathbb{C}^{m_1 \times m_2 \times \dots \times m_d \times n_1 \times n_2 \times \dots \times n_d}$

Input: A function to evaluate a $2d$ -mode tensor $\mathcal{K}(\mathbf{i}, \mathbf{j})$ for arbitrary multi-indices (\mathbf{i}, \mathbf{j}) , binary partitioning trees of L levels $\mathcal{T}_{t_k^0}$ and $\mathcal{T}_{s_k^0}$ with roots $t_k^0 = \{1, 2, \dots, m_k\}$ and $s_k^0 = \{1, 2, \dots, n_k\}$, a relative compression tolerance ϵ .

Output: Tensor butterfly decomposition of \mathcal{K} : (1) $\mathbf{V}_{\tau, \nu}^{s, k}$ at $l = 0$ and $\mathbf{W}_{\tau, \nu}^{s, k}$ at $1 \leq l \leq L^c$ of $k \leq d$ for multi-set τ with node τ_i at level l of $\mathcal{T}_{t_i^0}$, multi-set s with node s_i at level L^c of $\mathcal{T}_{s_i^0}$, and node ν at level $L^c - l$ of subtree \mathcal{T}_{s_k} , (2) $\mathbf{U}_{\tau, \nu}^{t, k}$ at $l = 0$ and $\mathbf{P}_{\tau, \nu}^{t, k}$ at $1 \leq l \leq L^c$ of $k \leq d$ for multi-set ν with node ν_i at level l of $\mathcal{T}_{s_i^0}$, multi-set t with node t_i at level L^c of $\mathcal{T}_{t_i^0}$, and node τ at level $L^c - l$ of subtree \mathcal{T}_{t_k} , and (3) subtensors $\mathcal{K}(\bar{t}, \bar{s})$ at $l = L^c$.

```

1: (1) Compute  $\mathbf{V}_{\tau, \nu}^{s, k}$  and  $\mathbf{W}_{\tau, \nu}^{s, k}$ :
2: for level  $l = 0, \dots, L^c$  do
3:   for multi-set  $s = (s_1, \dots, s_d)$  with  $s_i$  at level  $L^c$  of  $\mathcal{T}_{s_i^0}$  do
4:     for multi-set  $\tau = (\tau_1, \tau_2, \dots, \tau_d)$  with  $\tau_i$  at level  $l$  of  $\mathcal{T}_{t_i^0}$  do
5:       for mode index  $k = 1, \dots, d$  do
6:         for node  $\nu$  at level  $L^c - l$  of  $\mathcal{T}_{s_k}$  do
7:           if  $l = 0$  then ▷ Use (3.4) with proxies  $\hat{\tau}, \hat{s}$  and tolerance  $\epsilon$ 
8:             Compute  $\mathbf{V}_{\tau, \nu}^{s, k}$  and  $\bar{\nu}$  via mode- $(d + k)$  unfolding of  $\mathcal{K}(\hat{\tau}, \hat{s}_{k \leftarrow \nu})$ 
9:           else ▷ Use (3.6) with proxies  $\hat{\tau}, \hat{s}$  and tolerance  $\epsilon$ 
10:            Compute  $\mathbf{W}_{\tau, \nu}^{s, k}$  and  $\bar{\nu}$  via mode- $(d + k)$  unfolding of
11:               $\mathcal{K}(\hat{\tau}, \hat{s}_{k \leftarrow \nu^1 \cup \nu^2})$ 
12:            end if
13:          end for
14:        end for
15:      end for
16:    end for
17: (2) Compute  $\mathbf{U}_{\tau, \nu}^{t, k}$  and  $\mathbf{P}_{\tau, \nu}^{t, k}$ :
18: for level  $l = 0, \dots, L^c$  do
19:   for multi-set  $t = (t_1, \dots, t_d)$  with  $t_i$  at level  $L^c$  of  $\mathcal{T}_{t_i^0}$  do
20:     for multi-set  $\nu = (\nu_1, \nu_2, \dots, \nu_d)$  with  $\nu_i$  at level  $l$  of  $\mathcal{T}_{s_i^0}$  do
21:       for mode index  $k = 1, \dots, d$  do
22:         for node  $\tau$  at level  $L^c - l$  of  $\mathcal{T}_{t_k}$  do
23:           if  $l = 0$  then ▷ Use (3.5) with proxies  $\hat{t}, \hat{\nu}$  and tolerance  $\epsilon$ 
24:             Compute  $\mathbf{U}_{\tau, \nu}^{t, k}$  and  $\bar{\tau}$  via mode- $k$  unfolding of  $\mathcal{K}(\hat{t}_{k \leftarrow \tau}, \hat{\nu})$ 
25:           else ▷ Use (3.7) with proxies  $\hat{t}, \hat{\nu}$  and tolerance  $\epsilon$ 
26:             Compute  $\mathbf{P}_{\tau, \nu}^{t, k}$  and  $\bar{\tau}$  via mode- $k$  unfolding of  $\mathcal{K}(\hat{t}_{k \leftarrow \tau^1 \cup \tau^2}, \hat{\nu})$ 
27:           end if
28:         end for
29:       end for
30:     end for
31:   end for
32: end for
33: (3) Compute  $\mathcal{K}(\bar{t}, \bar{s})$ :
34: for multi-set  $s = (s_1, \dots, s_d)$  with  $s_i$  at level  $L^c$  of  $\mathcal{T}_{s_i^0}$  do
35:   for multi-set  $t = (t_1, \dots, t_d)$  with  $t_i$  at level  $L^c$  of  $\mathcal{T}_{t_i^0}$  do
36:     Compute  $\mathcal{K}(\bar{t}, \bar{s})$  and ZFP compress it
37:   end for
38: end for

```

Algorithm 3.2 Contraction algorithm for a tensor butterfly decomposition with an input tensor

Input: The tensor butterfly decomposition of a $2d$ -mode tensor

$\mathcal{K} \in \mathbb{C}^{m_1 \times m_2 \times \dots \times m_d \times n_1 \times n_2 \times \dots \times n_d}$, and a (full) $d+1$ -mode input tensor

$\mathcal{F} \in \mathbb{C}^{n_1 \times n_2 \times \dots \times n_d \times n_v}$ where n_v denotes the number of columns of $\mathbf{F}^{(d+1)}$.

Output: The $d+1$ -mode output tensor $\mathcal{G} = \mathcal{K} \times_{d+1, d+2, \dots, 2d} \mathcal{F}$ where

$\mathcal{G} \in \mathbb{C}^{m_1 \times m_2 \times \dots \times m_d \times n_v}$.

```

1: (1) Multiply with  $\mathbf{V}_{\tau, \nu}^{s, k}$  and  $\mathbf{W}_{\tau, \nu}^{s, k}$ :
2: for level  $l = 0, \dots, L^c$  do
3:   for multi-set  $\mathbf{s} = (s_1, s_2, \dots, s_d)$  with  $s_i$  at level  $L^c$  of  $\mathcal{T}_{s_i^0}$  do
4:     for multi-set  $\boldsymbol{\tau} = (\tau_1, \tau_2, \dots, \tau_d)$  with  $\tau_i$  at level  $l$  of  $\mathcal{T}_{t_i^0}$  do
5:       if  $l = 0$  then
6:          $\mathcal{F}_{\boldsymbol{\tau}, \mathbf{s}} = \mathcal{F}(\mathbf{s}, 1 : n_v) \prod_{k=1}^d \times_k \overline{\mathbf{V}}^{s, k}$ 
7:       else
8:          $\mathcal{F}_{\boldsymbol{\tau}, \mathbf{s}} = \mathcal{F}_{\boldsymbol{p}_{\boldsymbol{\tau}}, \mathbf{s}} \prod_{k=1}^d \times_k \text{diag}_{\nu}(\mathbf{W}_{\tau, \nu}^{s, k})$   $\triangleright \nu$  at level  $L^c - l$  of  $\mathcal{T}_{s_k}$ 
9:       end if
10:    end for
11:  end for
12: end for
13: (2) Contract with  $\mathcal{K}(\bar{\mathbf{t}}, \bar{\mathbf{s}})$ :
14: for multi-set  $\mathbf{t} = (t_1, t_2, \dots, t_d)$  with  $t_i$  at level  $L^c$  of  $\mathcal{T}_{t_i^0}$  do
15:   for multi-set  $\mathbf{s} = (s_1, s_2, \dots, s_d)$  with  $s_i$  at level  $L^c$  of  $\mathcal{T}_{s_i^0}$  do
16:     ZFP decompress  $\mathcal{K}(\bar{\mathbf{t}}, \bar{\mathbf{s}})$  and compute  $\mathcal{G}_{\mathbf{t}, \mathbf{s}} = \mathcal{K}(\bar{\mathbf{t}}, \bar{\mathbf{s}}) \times_{d+1, d+2, \dots, 2d} \mathcal{F}_{\mathbf{t}, \mathbf{s}}$ 
17:   end for
18: end for
19: (3) Multiply with  $\mathbf{U}_{\tau, \nu}^{t, k}$  and  $\mathbf{P}_{\tau, \nu}^{t, k}$ :
20: for level  $l = L^c, \dots, 0$  do
21:   for multi-set  $\mathbf{t} = (t_1, t_2, \dots, t_d)$  with  $t_i$  at level  $L^c$  of  $\mathcal{T}_{t_i^0}$  do
22:     for multi-set  $\boldsymbol{\nu} = (\nu_1, \nu_2, \dots, \nu_d)$  with  $\nu_i$  at level  $l$  of  $\mathcal{T}_{s_i^0}$  do
23:       if  $l = 0$  then  $\triangleright$  Compute and return  $\mathcal{G}$ 
24:          $\mathcal{G}(\mathbf{t}, 1 : n_v) = \mathcal{G}_{\mathbf{t}, \boldsymbol{\nu}} \prod_{k=1}^d \times_k \overline{\mathbf{U}}^{t, k}$ 
25:       else
26:          $\mathcal{G}_{\mathbf{t}, \boldsymbol{\nu}} += \mathcal{G}_{\mathbf{t}, \boldsymbol{\nu}} \prod_{k=1}^d \times_k \text{diag}_{\tau}(\mathbf{P}_{\tau, \nu}^{t, k})$   $\triangleright \tau$  at level  $L^c - l$  of  $\mathcal{T}_{t_k}$ 
27:       end if
28:     end for
29:   end for
30: end for
```

of the tensor butterfly algorithm.

High-frequency Green's functions. We use an example similar to the one used in subsection 3.2. Consider an OIO representing the free-space Green's function interaction between two parallel-facing unit-square plates. The $n \times n \times n \times n$ tensor is

$$(3.13) \quad \mathcal{K}(\mathbf{i}, \mathbf{j}) = K(x^{\mathbf{i}}, y^{\mathbf{j}}) = \frac{\exp(-i\omega\rho)}{\rho},$$

where $x^{\mathbf{i}} = (\frac{i_1}{n}, \frac{i_2}{n}, 0)$, $y^{\mathbf{j}} = (\frac{j_1}{n}, \frac{j_2}{n}, \rho_{\min})$, ω is the wavenumber, and $\rho = |x^{\mathbf{i}} - y^{\mathbf{j}}|$. Here ρ_{\min} represents the distance between the two plates assumed to be sufficiently large. In the high-frequency setting, $n = C_p \omega$ with a constant C_p independent of n and ω , and the grid size is $\delta_x = \delta_y = \frac{1}{n}$ per dimension. It has been well studied

475 [52, 53, 20, 5] that for any multi-set pair (τ, ν) leading to a subtensor $\mathcal{K}(\tau, \nu)$ of
 476 sizes $m_1 \times m_2 \times n_1 \times n_2$ with $m_i, n_i \leq n$, the numerical rank of its matricization
 477 $\mathbf{K} \in \mathbb{C}^{m_1 m_2 \times n_1 n_2}$ can be estimated as

$$478 \quad (3.14) \quad r_m \approx \omega^2 a^2 \theta \phi \approx \frac{\omega^2 a^2 n_1 n_2}{n^2 \rho_{\min}^2}.$$

479 Here a is the radius of the sphere enclosing the target domain of physical sizes $m_1 \delta_x \times$
 480 $m_2 \delta_y$. $\theta \approx \frac{n_1}{n \rho_{\min}}$, $\phi \approx \frac{n_2}{n \rho_{\min}}$, and the product $\theta \phi$ represents the solid angle covered
 481 by the source domain as seen from the center of the target domain. Note that $\frac{\omega a}{\rho_{\min}}$
 482 approximately represents the Nyquist sampling rate per direction needed in the source
 483 domain. The matrix and tensor butterfly ranks can be estimated as follows:

484 • *Matrix butterfly rank:* Consider a matrix butterfly factorization of matricization of
 485 \mathcal{K} . By design, for any node pair at each level, $m_1 n_1 = m_2 n_2 = C_b n$, where C_b^2
 486 represents the size of the leaf nodes. Therefore, the matrix butterfly rank can be
 487 estimated from (3.14) as

$$488 \quad (3.15) \quad r_m \approx \frac{C_b^2}{2 C_p^2 \rho_{\min}^2}$$

489 Here we have assumed $a = \frac{m_1}{\sqrt{2n}}$. Note that r_m is a constant independent of n , and
 490 therefore the matrix CLR property holds true.

491 • *Tensor butterfly rank:* Consider an L -level tensor butterfly factorization of \mathcal{K} . We
 492 just need to check the tensor rank, e.g., the rank of the mode-4 unfolding of the
 493 corresponding subtensors at Step (1) of Algorithm 3.1, as the unfolding for the
 494 other modes can be investigated in a similar fashion. Figure 3.3(a) shows an exam-
 495 ple of $L = 2$, where the target and source domains are partitioned at $l = 0$ (top)
 496 and $l = L^c = 1$ (bottom) at Step (1) of Algorithm 3.1. Consider a multi-set pair
 497 $(\tau, \mathbf{s}_{k \leftarrow \nu})$ with $k = 4$ required by the tensor CLR property in (3.4). Figure 3.3(a)
 498 highlights in orange one multi-set pair at $l = 0$ (top) and one multi-set pair at
 499 $l = L^c$ (bottom). Mode 4 is highlighted in red, which needs to be skeletonized by
 500 ID. By (3.14), the rank of the matricization of $\mathcal{K}(\tau, \mathbf{s}_{k \leftarrow \nu})$ is no longer a constant as
 501 the tensor butterfly algorithm needs to keep $n_1 = |s_1| = n/2^{L^c}$ (see Figure 3.3(b)).
 502 However, due to translational invariance of the free-space Green's function, i.e.,
 503 $K(x^i, y^j) = K(\tilde{x}, \tilde{y})$, where $\tilde{x} = (0, \frac{j_2}{n}, 0)$, $\tilde{y} = (\frac{j_1 - i_1}{n}, \frac{j_2}{n}, \rho_{\min})$, the mode-4 un-
 504 folding of $\mathcal{K}(\tau, \mathbf{s}_{k \leftarrow \nu})$ is the matrix representing the Green's function interaction
 505 between an enlarged target domain of sizes $(m_1 + n_1) \delta_x \times m_2 \delta_y$ and a source line
 506 segment of length $n_2 \delta_y$. Therefore its rank (hence the tensor rank) can be estimated
 507 as

$$508 \quad (3.16) \quad r_t \approx \omega a' \phi \approx \frac{\omega a' n_2}{n \rho_{\min}} \leq \frac{\sqrt{2} C_b}{C_p \rho_{\min}},$$

509 where a' is the radius of the sphere enclosing the enlarged target domain and $\frac{\omega a'}{\rho_{\min}}$
 510 approximately represents the Nyquist sampling rate on the source line segment.
 511 The last inequality is a result of $a' \approx \frac{m_1 + n_1}{\sqrt{2n}} \leq \frac{\sqrt{2} m_1}{n}$ and $m_2 n_2 = C_b n$. Here, the
 512 critical condition $n_1 \leq m_1$ is a direct result of the setup of the tensor CLR in (3.4):
 513 $l \leq L^c$ and $n_1 = |s_1| = n/2^{L^c}$ (i.e., s_1 is fixed as the middle level set as l changes).
 514 One can clearly see from (3.16) that r_t is independent of n , and thus the tensor
 515 CLR property holds true.

We remark that the tensor butterfly rank r_t in (3.16) is significantly smaller than the matrix butterfly rank r_m in (3.15) with $r_t \approx 2\sqrt{r_m}$. One can perform similar analysis of r_m and r_t for different geometrical settings, such as a pair of well-separated 3D unit cubes, or a pair of co-planar 2D unit-square plates. We leave these exercises to the readers.

Discrete Fourier Transform. Our second example is the high-dimensional discrete Fourier transform (DFT) defined by

$$(3.17) \quad \mathcal{K}(\mathbf{i}, \mathbf{j}) = \exp(2\pi i \mathbf{x}^{\mathbf{i}} \cdot \mathbf{y}^{\mathbf{j}})$$

with $\mathbf{x}^{\mathbf{i}} = (i_1 - 1, i_2 - 1, \dots, i_d - 1)$ and $\mathbf{y}^{\mathbf{j}} = (\frac{j_1-1}{n}, \frac{j_2-1}{n}, \dots, \frac{j_d-1}{n})$. We first notice that, since

$$(3.18) \quad \exp(2\pi i \mathbf{x}^{\mathbf{i}} \cdot \mathbf{y}^{\mathbf{j}}) = \prod_{k=1}^d \exp\left(\frac{2\pi i (i_k - 1)(j_k - 1)}{n}\right),$$

to carry out arbitrary high-dimensional DFTs one can simply perform 1D DFTs one dimension at a time (while fixing the indices of the other dimensions) by either 1D FFT or 1D matrix butterfly algorithms. We choose the 1D butterfly approach as our reference algorithm. For each node pair at dimension k discretized into a $m_k \times n_k$ matrix, we assume that $m_k n_k = C_b n$. It has been proved in [8, 68] that this leads to the matrix CLR property and each 1D DFT (fixing indices in other dimensions) can be computed by the matrix butterfly algorithm in $O(n \log n)$ time with a constant butterfly rank r_m . Overall this approach requires $O(dn^d \log n)$ operations.

In contrast, the tensor butterfly algorithm relies on direct compression of e.g., mode- k unfolding of subtensors $\mathcal{K}(\boldsymbol{\tau}, \mathbf{s}_{k \leftarrow \nu})$. Consider any submatrix $\mathbf{K}_{sub} \in \mathbb{C}^{m_k \times n_k}$ of this unfolding matrix $\mathbf{K}^{(k)}$; by fixing i_p and j_p with $p \neq k$, its entry is simply

$$\exp\left(\frac{2\pi i (i_k - 1)(j_k - 1)}{n}\right)$$

scaled by a constant factor

$$\prod_{p \neq k} \exp\left(\frac{2\pi i (i_p - 1)(j_p - 1)}{n}\right)$$

of modulus 1. Therefore the tensor butterfly rank is

$$(3.19) \quad r_t = \text{rank}(\mathbf{K}^{(k)}) = \text{rank}(\mathbf{K}_{sub}) = r_m.$$

The tensor CLR property thus holds true, and the tensor rank is exactly the same as the 1D butterfly algorithm per dimension. However, as we will see in subsection 3.2.2, our tensor butterfly algorithm yields a linear instead of quasi-linear CPU complexity for high-dimensional DFTs.

3.2.2. Complexity Analysis. Here we provide an analysis of computational complexity and memory requirement of the proposed construction algorithm (Algorithm 3.1) and contraction algorithm (Algorithm 3.2), assuming that the tensor butterfly rank r_t is a small constant and $d > 1$.

At Step (1) of Algorithm 3.1, each level $1 \leq l \leq L^c$ has $\#\mathbf{s} = O(\sqrt{n}^d)$, $\#\boldsymbol{\tau} = 2^{dl}$, $\#\nu = O(\sqrt{n}/2^l)$ for each mode $k \leq d$. Each $\mathbf{W}_{\boldsymbol{\tau}, \nu}^{s, k}$ requires $O(r_t^2)$ storage, and

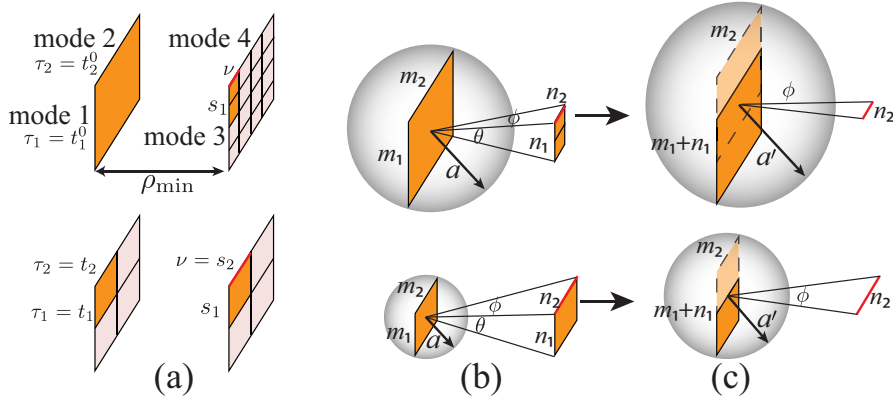


Fig. 3.3: Illustration of the tensor CLR property with $L = 2$ for a 4-mode tensor representing free-space Green's function interactions between parallel facing unit square plates. (a) The target and source domains are partitioned at $l = 0$ (top) and $l = L^c = 1$ (bottom) with a multi-set pair $(\tau, s_{k \leftarrow \nu})$ highlighted in orange for the skeletonization along mode 4. The sizes of the nodes are $|\tau_1| = m_1, |\tau_2| = m_1, |s_1| = n_1$ and $|\nu| = n_2$. (b) Illustration of the rank of the matricization of $\mathcal{K}(\tau, s_{k \leftarrow \nu})$. (c) Illustration of the rank of the mode-4 unfolding of $\mathcal{K}(\tau, s_{k \leftarrow \nu})$.

553 $O(r_t^{2d+1})$ computational time when proxy indices $\hat{\tau}, \hat{s}$ are being used. The storage
 554 requirement and computational cost for $\mathbf{W}_{\tau, \nu}^{s, k}$ are:

555 (3.20)
$$\text{mem}_W = \sum_{l=1}^{L^c} dO(\sqrt{n}^d) 2^{dl} O(\sqrt{n}/2^l) O(r_t^2) = O(dn^d r_t^2),$$

556 (3.21)
$$\text{time}_W = \sum_{l=1}^{L^c} dO(\sqrt{n}^d) 2^{dl} O(\sqrt{n}/2^l) O(r_t^{2d+1}) = O(dn^d r_t^{2d+1}).$$

557 One can easily verify that the computation and storage of $\mathbf{V}_{\tau, \nu}^{s, k}$ at $l = 0$ is less
 558 dominant than $\mathbf{W}_{\tau, \nu}^{s, k}$ at $l > 0$ and we skip its analysis.

559 At Step (2) of Algorithm 3.1, we have $\#s = O(\sqrt{n}^d)$ and $\#t = O(\sqrt{n}^d)$, and
 560 each $\mathcal{K}(\bar{t}, \bar{s})$ requires $O(r_t^{2d})$ computation time and storage units (even if it is further
 561 ZFP compressed to reduce storage requirement), which adds up to

562 (3.22)
$$\text{mem}_K = O(\sqrt{n}^d) O(\sqrt{n}^d) O(r_t^{2d}) = O(n^d r_t^{2d}),$$

563 (3.23)
$$\text{time}_K = O(\sqrt{n}^d) O(\sqrt{n}^d) O(r_t^{2d}) = O(n^d r_t^{2d}).$$

564 Step (3) of Algorithm 3.1 has similar computational cost and memory requirement
 565 to Step (1) when contracting with the intermediate matrices $\mathbf{P}_{\tau, \nu}^{t, k}$, with $\text{mem}_P \sim$
 566 mem_W and $\text{time}_P \sim \text{time}_W$.

567 Overall, Algorithm 3.1 requires

568 (3.24)
$$\text{mem} = \text{mem}_W + \text{mem}_K + \text{mem}_P = O(n^d r_t^{2d}),$$

569 (3.25)
$$\text{time} = \text{time}_W + \text{time}_K + \text{time}_P = O(dn^d r_t^{2d+1}).$$

570 Following a similar analysis, one can estimate the computational cost of Algo-
 571 rithm 3.2 as $O(n^d r_t^{2d} n_v)$, which is essentially of the similar order as mem of Algo-

Algorithm	Factor time		Apply time		r	
	$d = 2$	$d = 3$	$d = 2$	$d = 3$	$d = 2$	$d = 3$
Tensor butterfly	n^2	n^3	n^2	n^3	1	1
Matrix butterfly	$n^2 \log n$	$n^3 \log n$	$n^2 \log n$	$n^3 \log n$	1	1
Tucker ID	n^4	$n^4 - n^{6*}$	n^4	$n^4 - n^{6*}$	n	n
QTT (Green&Radon)	$n^3 \log n$	$n^3 \log n$	$n^4 \log n$	$n^5 \log n$	n	n
QTT (DFT)	$\log n$	$\log n$	$n^2 \log n$	$n^3 \log n$	1	1

Table 3.1: CPU complexity of the tensor butterfly algorithm, matrix butterfly algorithm, Tucker ID and QTT when applied to high-frequency Green’s functions ($d = 2$ represents two parallel facing unit square plates and $d = 3$ represents two separated unit cubes), DFT and Radon transforms. Here we assume that tensor butterfly, matrix butterfly and Tucker ID algorithms use proxy indices, and the QTT algorithm uses TT-cross. The big O notation is assumed. *: for $d = 3$, the complexity of Tucker ID is n^6 for Radon transform and DFT, and n^4 for Green’s function.

algorithm 3.1, except an extra factor n_v representing the size of the last dimension of the input tensor.

One critical observation is that the time and storage complexity of the tensor butterfly algorithm is *linear* in n^d with smaller ranks r_t , while that of the matrix butterfly algorithm is *quasi-linear* in n^d with much larger ranks r_m . This leads to a significantly superior algorithm, as will be demonstrated with the numerical results in section 4. That being said, one can verify that there is no difference between the two algorithms when $d = 1$.

3.2.3. Comparison with Tucker-ID and QTT. Here we make a comparison of the computational complexities of the matrix butterfly algorithm, tensor butterfly algorithm, Tucker-ID and QTT for several frequently encountered OIOs with $d = 2, 3$, namely Green’s functions for high-frequency wave equations (where $d = 2$ represents two parallel facing unit square plates and $d = 3$ represents two separated unit cubes), Radon transforms (a type of Fourier integral operators), and DFT. We first summarize the computational complexities of the factorization and application of matrix and tensor butterfly algorithms in Table 3.1. Here we use r to denote the maximum rank of the submatrices or (unfolding and matricization of) subtensors associated with each algorithm. In other words, we drop the subscript of r_m and r_t in this subsection. We note that $r = O(1)$ for butterfly algorithms, and the computational complexity for matrix and tensor butterfly algorithms is, respectively, $O(dn^d \log n)$ and $O(dn^d)$, for all OIOs considered here.

The Tucker ID algorithm in subsection 3.1 (even with the use of proxy indices to accelerate the factorization), always leads to $r = O(n)$ for OIOs and hence almost always $O(n^{2d})$ factorization and application complexities (see Table 3.1). One exception is perhaps the Green’s function for $d = 3$, where one can easily show that 4 out of the 6 unfolding matrices have a rank of $O(n)$ and the remaining 2 have a rank of $O(1)$, leading to the $O(n^4)$ computational complexity. Overall, we remark that Tucker-like decomposition algorithms are typically the least efficient tensor algorithms for OIOs.

The QTT algorithm, on the other hand, is a more subtle algorithm to compare with. Assuming that the maximum rank among all steps in QTT is r , we first summarize the computational complexities of the factorization and application of QTT. For factorization, we only consider the TT-cross type of algorithms, which yields the best

known computational complexity among all TT-based algorithms. The computational complexity of TT-cross is $O(dr^3 \log n)$ [14, 55]. Once factorized, the application cost of the QTT factorization with a full input tensor is $O(dr^2 n^d \log n)$ [14]. This complexity can be reduced to $O(dr^2 r_i^2 \log n)$ when the input tensor is also in the QTT format with TT rank r_i . However, an arbitrary input tensor can have a TT rank up to $r_i = O(n^{d/2})$ (which leads to the same application cost as contraction with a full input tensor). Therefore in our comparative study, we stick with the $O(dr^2 n^d \log n)$ application complexity.

For high-frequency Green's functions and general-form Fourier integral operators (e.g. Radon transforms), the TT rank in general behaves as $r = O(n)$ [14], leading to a factorization cost of $O(dn^3 \log n)$ and an application cost of $O(dn^{2+d} \log n)$, as detailed in Table 3.1. It is worth mentioning that, treating DFTs as a special type of Fourier integral operators, QTT can achieve $r = O(1)$ when a proper bit-reversal ordering is used [9], leading to a factorization cost of $O(d \log n)$ and an application cost of $O(dn^d \log n)$, as shown in Table 3.1. In contrast, the proposed tensor butterfly algorithm can always yield $O(dn^d)$ factorization and $O(n^d)$ application costs.

4. Numerical Results. This section provides several numerical examples to demonstrate the accuracy and efficiency of the proposed tensor butterfly algorithm when applied to large-scale and high-dimensional OIOs including Green's function tensors for high-frequency Helmholtz equations (subsection 4.1), Radon transform tensors (subsection 4.2), and high-dimensional DFTs (subsection 4.3). We compare our algorithm with a few existing matrix and tensor algorithms including the matrix butterfly algorithm in subsection 2.2, the Tucker ID algorithm in subsection 3.1, the QTT algorithm [55], the FFT algorithm implemented in the heFFTe package [1], and the non-uniform FFT (NUFFT) algorithm implemented in the FINUFFT package [2]. All of these algorithms except for Tucker ID and FINUFFT are tested in distributed-memory parallelism. It is worth noting that currently there is no single package that can both compute and apply the QTT decomposition in distributed-memory parallelism. In our tests, we perform the factorization using a distributed-memory TT code [61] that parallelizes a cross interpolation algorithm [19], and then we implement the distributed-memory QTT contraction via the CTF package [62]. All experiments are performed using 4 CPU nodes of the Perlmutter machine at NERSC in Berkeley, where each node has two 64-core AMD EPYC 7763 processors and 128GB of 2133MHz DDR4 memory.

4.1. Green's functions for high-frequency Helmholtz equations. In this subsection, we consider the tensor discretized from 3D free-space Green's functions for high-frequency Helmholtz equations. Specifically, the tensor entry is

$$(4.1) \quad \mathcal{K}(\mathbf{i}, \mathbf{j}) = \frac{\exp(-i\omega\rho)}{\rho}, \quad \rho = |\mathbf{x}^{\mathbf{i}} - \mathbf{y}^{\mathbf{j}}|,$$

where ω represents the wave number. Two tests are performed: (1) A 4-way tensor representing the Green's function interaction between two parallel facing unit plates with distance 1, i.e., $\mathbf{x}^{\mathbf{i}} = (\frac{i_1}{n}, \frac{i_2}{n}, 0)$, $\mathbf{y}^{\mathbf{j}} = (\frac{j_1}{n}, \frac{j_2}{n}, 1)$, and $d = 2$. (2) A 6-way tensor representing the Green's function interaction between two unit cubes with the distance between their centers set to 2, i.e., $\mathbf{x}^{\mathbf{i}} = (\frac{i_1}{n}, \frac{i_2}{n}, \frac{i_3}{n})$, $\mathbf{y}^{\mathbf{j}} = (\frac{j_1}{n}, \frac{j_2}{n}, \frac{j_3}{n} + 2)$, and $d = 3$. For both tests, the wave number is chosen such that the number of points per wave length is 4, i.e., $2\pi n/\omega = 4$ or $C_p = 2/\pi$. We first perform compression using the tensor butterfly, Tucker ID and QTT algorithms, and then perform application/contraction

using a random input tensor \mathcal{F} . We also add results for the matrix butterfly algorithm using the corresponding matricization of \mathcal{K} and \mathcal{F} .

Figure 4.1 (left) shows the factorization time, application time and memory usage of each algorithm using a compression tolerance $\epsilon = 10^{-6}$ for the parallel plate case. For QTT, we show the memory of the factorization (labeled as “QTT(Factor)”) and application (labeled as “QTT(Apply)”) separately. Note that although QTT factorization requires sub-linear memory usage, QTT contraction becomes super-linear due to the full QTT rank of the input tensor. Overall, we achieve the expected complexities listed in Table 3.1 for the butterfly and Tucker ID algorithms. For QTT, however, instead of an $O(n)$ rank scaling, we observe an $O(n^{3/4})$ rank scaling, leading to slightly better complexities compared with Table 3.1. We leave this as a future investigation. That said, the tensor butterfly algorithm achieves the linear CPU and memory complexities for both factorization and application with a much smaller prefactor compared to all the other algorithms. Remarkably, the tensor butterfly algorithm achieves a 30x memory reduction and 15x speedup, capable of handling 64x larger-sized tensors compared with the matrix butterfly algorithm.

Figure 4.1 (right) shows the factorization time, application time and memory usage of each algorithm using a compression tolerance $\epsilon = 10^{-2}$ for the cube case. Overall, we achieve the expected complexities listed in Table 3.1 for all four algorithms. The tensor butterfly algorithm achieves the linear CPU and memory complexities for both factorization and application with a much smaller prefactor compared to all the other algorithms. Remarkably, the tensor butterfly algorithm achieves a $30\times$ memory reduction and 200x speedup, capable of handling $512\times$ larger-sized tensors compared with the matrix butterfly algorithm. The largest data point $n = 2048$ corresponds to 512 wavelengths per physical dimension. The results in Figure 4.1 suggest the superiority of the tensor butterfly algorithm in solving high-frequency wave equations in 3D volumes and on 3D surfaces.

Next, we demonstrate the effect of changing compression tolerance ϵ for both test cases in Table 4.1. Here the error is measured by

$$(4.2) \quad \text{error} = \frac{\|\mathcal{K} \times_{d+1,d+2,\dots,2d} \mathcal{F}_e - \mathcal{K}_{BF} \times_{d+1,d+2,\dots,2d} \mathcal{F}_e\|_F}{\|\mathcal{K} \times_{d+1,d+2,\dots,2d} \mathcal{F}_e\|_F}$$

where \mathcal{K}_{BF} is the tensor butterfly representation of \mathcal{K} , $\mathcal{F}_e(j) = 1$ for a small set of random entries j and 0 elsewhere. This way, \mathcal{K} does not need to be fully formed to compute the error. Table 4.1 shows the minimum rank (r_{\min}) and maximum rank (r), error, factorization time, application time and memory usage of varying ϵ , for $n = 16384, d = 2$ and $n = 512, d = 3$. Overall, the errors are close to the prescribed tolerances and the costs increase for smaller ϵ , as expected. We also note that keeping r as low as possible is critical in maintaining small prefactors of the tensor butterfly algorithm, particularly for higher dimensions.

4.2. Radon transforms. In this subsection, we consider 2D and 3D discretized Radon transforms similar to those presented in [8]. Specifically, the tensor entry is

$$(4.3) \quad \mathcal{K}(i, j) = \exp(2\pi i \phi(x^i, y^j))$$

with $x^i = (\frac{i_1}{n}, \frac{i_2}{n}, \dots, \frac{i_d}{n})$ and $y^j = (j_1 - \frac{n}{2}, j_2 - \frac{n}{2}, \dots, j_d - \frac{n}{2})$. For $d = 2$, we consider

$$(4.4) \quad \begin{aligned} \phi(x, y) &= x \cdot y + \sqrt{c_1^2 y_1^2 + c_2^2 y_2^2}, \\ c_1 &= (2 + \sin(2\pi x_1) \sin(2\pi x_2))/16, \end{aligned}$$

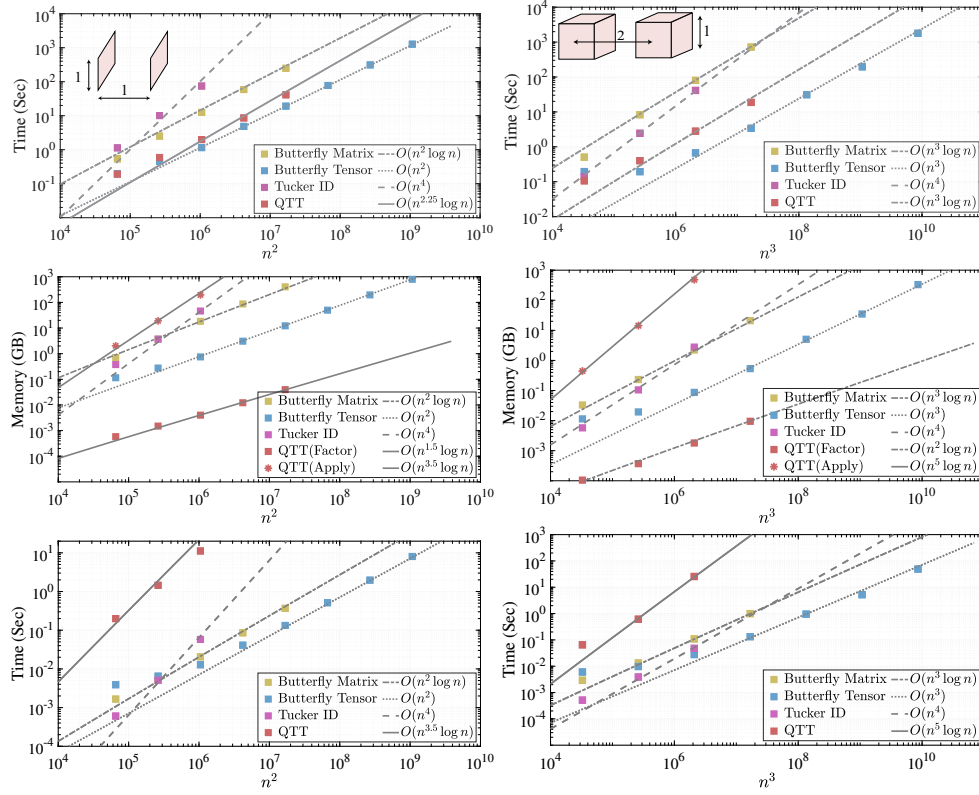


Fig. 4.1: Helmholtz equation: Computational complexity comparison among butterfly matrix, butterfly tensor, Tucker ID and QTT for compressing (left) a 4-way Green's function tensor for interactions between two parallel 2D plates and (right) a 6-way Green's function tensor for interactions between two 3D cubes. The geometries are discretized with 4 points per wavelength. (Top): Factor time. (Middle): Factor and apply memory. (Bottom): Apply time. The largest data points correspond to 8192 wavelengths per direction for the 2D tests (left) and 512 wavelengths per direction for the 3D tests (right).

$$c_2 = (2 + \cos(2\pi x_1) \cos(2\pi x_2))/16.$$

For $d = 3$, we consider

$$(4.5) \quad \begin{aligned} \phi(x, y) &= x \cdot y + c|y|, \\ c &= (3 + \sin(2\pi x_1) \sin(2\pi x_2) \sin(2\pi x_3))/100. \end{aligned}$$

We first perform compression using the matrix butterfly, tensor butterfly, and QTT algorithms, and then perform application/contraction using a random input tensor \mathcal{F} .

Figure 4.2 shows the factorization time, application time and memory usage of each algorithm using a compression tolerance $\epsilon = 10^{-3}$ for the 2D transform (left) and 3D transform (right). Overall, we achieve the expected complexities listed in Table 3.1 for all three algorithms. The QTT algorithm can only obtain the first 2 or 3 data points due to its high memory usage and large QTT ranks. In comparison,

n^d	ϵ	r_{\min}	r	error	$T_f(\text{sec})$	$T_a(\text{sec})$	$Mem(\text{MB})$
16384^2	1E-02	5	8	1.49E-02	6.83E+01	1.16E+00	2.40E+04
16384^2	1E-03	6	10	2.19E-03	1.17E+02	1.89E+00	4.69E+04
16384^2	1E-04	7	11	1.84E-04	1.57E+02	2.80E+00	7.49E+04
16384^2	1E-05	8	12	3.46E-05	2.29E+02	4.03E+00	1.21E+05
16384^2	1E-06	9	13	9.26E-06	3.18E+02	5.92E+00	1.96E+05
512^3	1E-02	2	5	2.01E-02	1.18E+02	1.42E+00	1.19E+04
512^3	1E-03	2	6	1.18E-03	3.46E+02	4.08E+00	4.87E+04
512^3	1E-04	2	7	8.39E-05	6.26E+02	9.85E+00	1.49E+05
512^3	1E-05	3	8	9.21E-06	1.25E+03	2.40E+01	4.07E+05

Table 4.1: The technical data for a 4-way Green’s function tensor of $n = 16384$ and a 6-way Green’s function tensor of $n = 512$ for the Helmholtz equation using the proposed tensor butterfly algorithm of varying compression tolerance ϵ . The table shows the maximum rank r and minimum rank r_{\min} across all ID operations, relative error in (4.2), factor time T_f , apply time T_a , and memory usage Mem .

n^d	ϵ	r_{\min}	r	error	$T_f(\text{sec})$	$T_a(\text{sec})$	$Mem(\text{MB})$
2048^2	1E-02	4	18	2.04E-02	9.32E+01	7.20E-01	1.25E+04
2048^2	1E-03	4	20	1.51E-03	1.61E+02	1.28E+00	2.40E+04
2048^2	1E-04	4	22	1.49E-04	2.55E+02	2.05E+00	4.26E+04
2048^2	1E-05	4	23	2.45E-05	3.73E+02	3.12E+00	6.95E+04
128^3	1E-02	2	6	4.31E-02	3.89E+01	8.57E-01	1.59E+04
128^3	1E-03	2	8	1.00E-02	1.31E+02	3.74E+00	9.44E+04
128^3	1E-04	2	9	1.68E-03	2.42E+02	8.28E+00	2.38E+05
128^3	1E-05	2	11	1.48E-04	4.30E+02	2.05E+01	6.06E+05

Table 4.2: The technical data for a 4-way Radon transform tensor of $n = 2048$ in (4.4) and a 6-way Radon transform tensor of $n = 128$ in (4.5) using the proposed tensor butterfly algorithm of varying compression tolerance ϵ . The table shows the maximum rank r and minimum rank r_{\min} across all ID operations, relative error in (4.2), factor time T_f , apply time T_a , and memory usage Mem .

the tensor butterfly algorithm achieves the linear CPU and memory complexities for both factorization and application with a much smaller prefactor compared to all the other algorithms. Note that the Radon transform kernels in (4.4) and (4.5) are not translational invariant, but the tensor butterfly algorithm can still attain small ranks. As a result, the tensor butterfly algorithm can handle 64x larger-sized Radon transforms compared with the matrix butterfly algorithm, showing their superiority for solving linear inverse problems in tomography and seismic imaging.

Next, we demonstrate the effect of changing compression tolerance ϵ for both test cases in Table 4.2 with the error defined by (4.2). Table 4.2 shows the minimum and maximum ranks, error, factorization time, application time and memory usage of varying ϵ , for $n = 2048$ with $d = 2$ and $n = 128$ with $d = 3$, respectively. Overall, the errors are close to the prescribed tolerances and the costs increase for smaller ϵ , as expected. Just like the Green’s function example, it is critical to keep r a low

719 constant, particularly for higher dimensions.

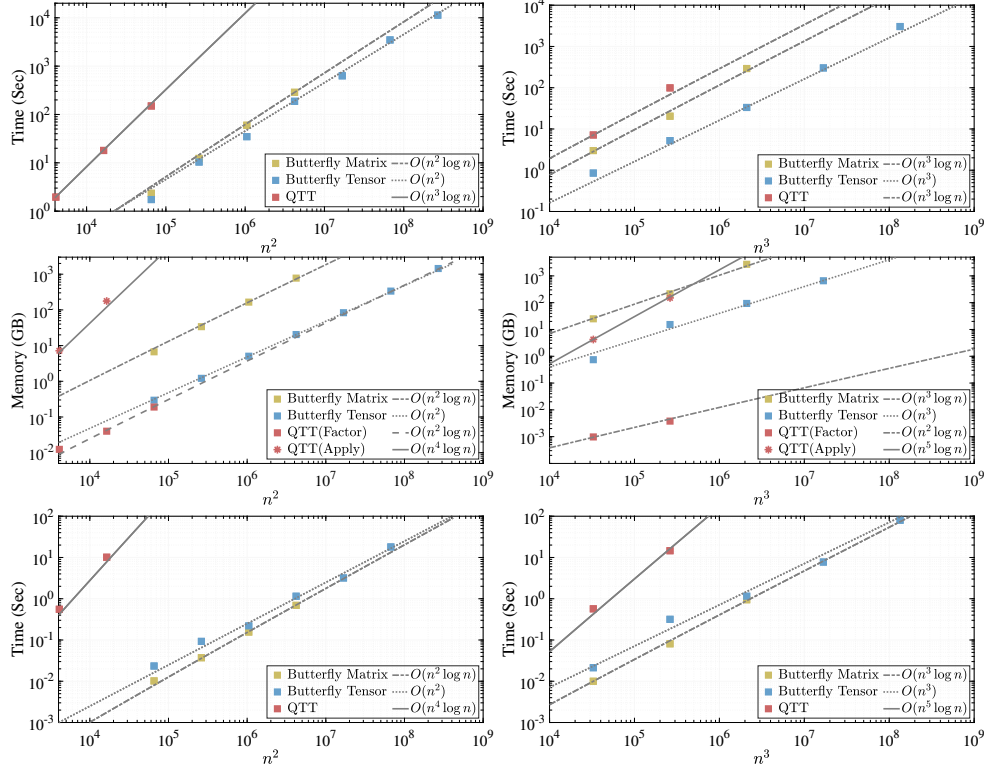


Fig. 4.2: Radon transforms: Computational complexity comparison among butterfly matrix, butterfly tensor and QTT for compressing (left) a 2D Radon transform tensor and (right) a 3D Radon transform tensor. (Top): Factor time. (Middle): Factor and apply memory. (Bottom): Apply time.

720 **4.3. High-dimensional discrete Fourier transform.** Finally, we consider
721 high-dimensional DFTs defined as

722 (4.6)
$$\mathcal{K}(\mathbf{i}, \mathbf{j}) = \exp(2\pi i \mathbf{x}^{\mathbf{i}} \cdot \mathbf{y}^{\mathbf{j}}),$$

723 where we choose $\mathbf{x}^{\mathbf{i}} = (i_1 - 1, i_2 - 1, \dots, i_d - 1)$ and $\mathbf{y}^{\mathbf{j}} = (\frac{j_1-1}{n}, \frac{j_2-1}{n}, \dots, \frac{j_d-1}{n})$ for
724 uniform DFTs, and we choose $\mathbf{x}^{\mathbf{i}}$ to be random (in the sense that $x_k^{\mathbf{i}} \in [0, n-1]$ for
725 $k \leq d$ is a random number) and $\mathbf{y}^{\mathbf{j}} = (\frac{j_1-1}{n}, \frac{j_2-1}{n}, \dots, \frac{j_d-1}{n})$ for type-2 non-uniform
726 DFTs. For high-dimensional DFTs with $d = 3, 4, 5, 6$, we perform compression using
727 the tensor butterfly algorithms (with the bit-reversal ordering for each dimension),
728 and perform application/contraction using a random input tensor \mathcal{F} . In comparison,
729 for $d = 3$ we perform FFT via the heFFTe package for the uniform DFT example and
730 NUFFT via the FINUFFT package for the type-2 non-uniform DFT example.

731 Figure 4.3 shows the factorization time for the butterfly algorithm (or equiva-
732 lently the plan creation time for heFFTe/FINUFFT), application time and memory
733 usage of each algorithm using a compression tolerance $\epsilon = 10^{-3}$ (for butterfly and
734 FINUFFT) for the uniform (left) and nonuniform (right) transforms. Overall, the ten-
735 sor butterfly algorithm can obtain $O(n^d)$ CPU and memory complexities compared

with the $O(n^d \log n)$ complexities of FFT and NUFFT. It is also worth mentioning that QTT can attain logarithmic-complexity uniform DFTs [9] when the input tensor \mathcal{F} is also in the QTT form with low TT ranks. However, for a general input tensor, the complexity of QTT falls back to $O(n^d \log n)$. Although the proposed tensor butterfly algorithm can obtain the best computational complexity among all existing algorithms, we observe that for the $d = 3$ case, FFT or NUFFT shows a memory usage similar to the tensor butterfly algorithm but much smaller prefactors for plan creation and application time. That said, the tensor butterfly algorithm provides a unique capability to perform higher dimensional DFTs (i.e., $d \geq 4$) with optimal asymptotic complexities.

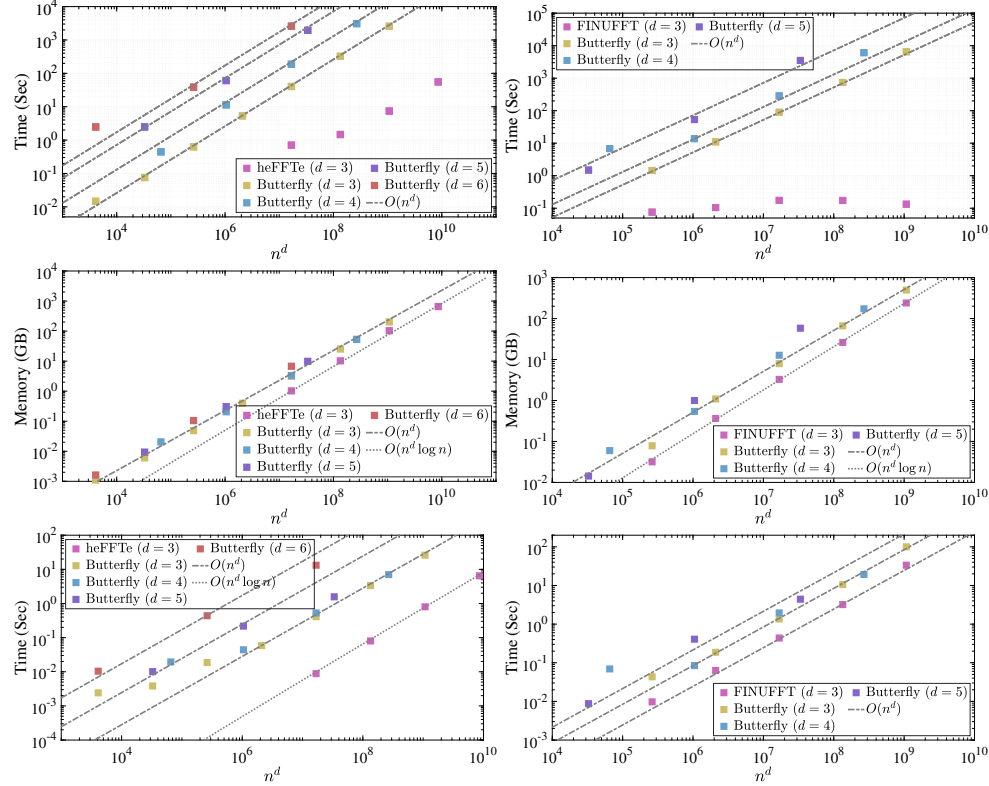


Fig. 4.3: Fourier transforms: Computational complexity of (left) butterfly tensor and heFFTe for compressing the high-dimensional DFT tensor and (right) butterfly tensor and FINUFFT for compressing the high-dimensional NUFFT tensor. (Top): Factor time of butterfly tensor and plan creation time for heFFTe/FINUFFT. (Middle): Factor memory. (Bottom): Apply time.

5. Conclusion. We present a new tensor butterfly algorithm efficiently compressing and applying large-scale and high-dimensional OIOs, such as Green's functions for wave equations and integral transforms, including Radon transforms and Fourier transforms. The tensor butterfly algorithm leverages an essential tensor CLR property to achieve both improved asymptotic computational complexities and lower leading constants. For the contraction of high-dimensional OIOs with arbitrary input tensors, the tensor butterfly algorithm achieves the optimal linear CPU and memory

complexities; this is in huge contrast with both existing matrix algorithms and fast transform algorithms. The former includes the matrix butterfly algorithm, and the latter contains FFT, NUFFT, and other tensor algorithms such as Tucker-like decompositions and QTT. Nevertheless, all these algorithms exhibit higher asymptotic complexities and larger leading constants. As a result, the tensor butterfly algorithm can efficiently model high-frequency 3D Green’s function interactions with over $512 \times$ larger problem sizes than existing algorithms; for the largest sized tensor that can be handled by existing algorithms, the tensor butterfly algorithm requires $200 \times$ less CPU time and $30 \times$ less memory than existing algorithms. Moreover, it can perform linear-complexity Radon transforms and DFTs with up to $d = 6$ dimensions. These OIOs are frequently encountered in the solution of high-frequency wave equations, X-ray and MRI-based inverse problems, seismic imaging and signal processing; therefore, we expect the tensor butterfly algorithm developed here to be both theoretically attractive and practically useful for many applications.

The limitation of the tensor butterfly algorithm is the requirement for a tensor grid, and hence its extension for unstructured meshes will be a future work. Also, the mid-level subtensors represent a memory bottleneck and need to be compressed with more efficient algorithms.

Acknowledgements. This research has been supported by the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research, Mathematical Multifaceted Integrated Capability Centers (MMICCs) program, under Contract No. DE-AC02-05CH11231 at Lawrence Berkeley National Laboratory. This work was also supported in part by the NSF Mathematical Sciences Graduate Internship (NSF MSGI) program. This research used resources of the National Energy Research Scientific Computing Center, a DOE Office of Science User Facility supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC02-05CH11231 using NERSC award ASCR-ERCAP0024170. Qian’s research was partially supported by NSF grants 2152011 and 2309534 and an MSU SPG grant.

REFERENCES

- [1] Alan Ayala, Stanimire Tomov, Piotr Luszczek, Sebastien Cayrols, Gerald Ragghianti, and Jack Dongarra. Analysis of the communication and computation cost of fft libraries towards exascale. Technical Report ICL-UT-22-07, 2022-07 2022.
- [2] Alexander H Barnett, Jeremy Magland, and Ludvig af Klinteberg. A parallel nonuniform fast fourier transform library based on an “exponential of semicircle” kernel. *SIAM Journal on Scientific Computing*, 41(5):C479–C504, 2019.
- [3] David Joseph Biagioni. *Numerical construction of Green’s functions in high dimensional elliptic problems with variable coefficients and analysis of renewable energy data via sparse and separable approximations*. PhD thesis, University of Colorado at Boulder, 2012.
- [4] James Bremer, Ze Chen, and Haizhao Yang. Rapid Application of the Spherical Harmonic Transform via Interpolative Decomposition Butterfly Factorization. *arXiv preprint arXiv:2004.11346*, 2020.
- [5] Ovidio M Bucci and Giorgio Franceschetti. On the degrees of freedom of scattered fields. *IEEE transactions on Antennas and Propagation*, 37(7):918–926, 1989.
- [6] HanQin Cai, Keaton Hamm, Longxiu Huang, and Deanna Needell. Mode-wise tensor decompositions: Multi-dimensional generalizations of cur decompositions. *Journal of machine learning research*, 22(185):1–36, 2021.
- [7] Emmanuel Candes, Laurent Demanet, and Lexing Ying. Fast computation of fourier integral operators. *SIAM Journal on Scientific Computing*, 29(6):2464–2493, 2007.
- [8] Emmanuel Candès, Laurent Demanet, and Lexing Ying. A fast butterfly algorithm for the computation of Fourier integral operators. *Multiscale Model. Sim.*, 7(4):1727–1750, 2009.
- [9] Jieliun Chen and Michael Lindsey. Direct interpolative construction of the discrete fourier transform as a matrix product operator. *arXiv preprint arXiv:2404.03182*, 2024.

- [10] Ze Chen, Juan Zhang, Kenneth L Ho, and Haizhao Yang. Multidimensional phase recovery and interpolative decomposition butterfly factorization. *Journal of Computational Physics*, 412:109427, 2020.
- [11] Jian Cheng, Dinggang Shen, Peter J Basser, and Pew-Thian Yap. Joint 6D kq space compressed sensing for accelerated high angular resolution diffusion MRI. In *Information Processing in Medical Imaging: 24th International Conference, IPMI 2015, Sabhal Mor Ostaig, Isle of Skye, UK, June 28-July 3, 2015, Proceedings*, pages 782–793. Springer, 2015.
- [12] Andrzej Cichocki, Namgil Lee, Ivan Oseledets, Anh-Huy Phan, Qibin Zhao, Danilo P Mandic, et al. Tensor networks for dimensionality reduction and large-scale optimization: Part 1 low-rank tensor decompositions. *Foundations and Trends® in Machine Learning*, 9(4-5):249–429, 2016.
- [13] Lisa Claus, Pieter Ghysels, Yang Liu, Thái Anh Nhan, Ramakrishnan Thirumalaisamy, Amneet Pal Singh Bhalla, and Sherry Li. Sparse approximate multifrontal factorization with composite compression methods. *ACM Transactions on Mathematical Software*, 49(3):1–28, 2023.
- [14] Eduardo Corona, Abtin Rahimian, and Denis Zorin. A tensor-train accelerated solver for integral equations in complex geometries. *Journal of Computational Physics*, 334:145–169, 2017.
- [15] Maurice A De Gosson. *The Wigner Transform*. World Scientific Publishing Company, 2017.
- [16] Lieven De Lathauwer, Bart De Moor, and Joos Vandewalle. A multilinear singular value decomposition. *SIAM journal on Matrix Analysis and Applications*, 21(4):1253–1278, 2000.
- [17] Stanley R Deans. *The Radon transform and some of its applications*. Courier Corporation, 2007.
- [18] Gian Luca Delzanno. Multi-dimensional, fully-implicit, spectral method for the vlasov–maxwell equations with exact conservation laws in discrete form. *Journal of Computational Physics*, 301:338–356, 2015.
- [19] Sergey Dolgov and Dmitry Savostyanov. Parallel cross interpolation for high-precision calculation of high-dimensional integrals. *Computer Physics Communications*, 246:106869, 2020.
- [20] Björn Engquist and Lexing Ying. Fast directional multilevel algorithms for oscillatory kernels. *SIAM Journal on Scientific Computing*, 29(4):1710–1737, 2007.
- [21] Alexander L Fetter and John Dirk Walecka. *Quantum theory of many-particle systems*. Courier Corporation, 2012.
- [22] Ilias I Giannakopoulos, Mikhail S Litsarev, and Athanasios G Polimeridis. Memory footprint reduction for the fft-based volume integral equation method via tensor decompositions. *IEEE Transactions on Antennas and Propagation*, 67(12):7476–7486, 2019.
- [23] Lars Grasedyck, Daniel Kressner, and Christine Tobler. A literature survey of low-rank tensor approximation techniques. *GAMM-Mitteilungen*, 36(1):53–78, 2013.
- [24] Han Guo, Jun Hu, and Eric Michielssen. On MLMDA/butterfly compressibility of inverse integral operators. *IEEE Antennas Wirel. Propag. Lett.*, 12:31–34, 2013.
- [25] Han Guo, Yang Liu, Jun Hu, and Eric Michielssen. A butterfly-based direct integral-equation solver using hierarchical LU factorization for analyzing scattering from electrically large conducting objects. *IEEE Trans. Antennas Propag.*, 65(9):4742–4750, 2017.
- [26] Han Guo, Yang Liu, Jun Hu, and Eric Michielssen. A butterfly-based direct solver using hierarchical LU factorization for Poggio-Miller-Chang-Harrington-Wu-Tsai equations. *Microw Opt Technol Lett.*, 60:1381–1387, 2018.
- [27] Wolfgang Hackbusch and Boris N Khoromskij. Tensor-product approximation to multidimensional integral operators and green’s functions. *SIAM journal on matrix analysis and applications*, 30(3):1233–1253, 2008.
- [28] Wolfgang Hackbusch and Stefan Kühn. A new scheme for the tensor representation. *Journal of Fourier analysis and applications*, 15(5):706–722, 2009.
- [29] Nathan Halko, Per-Gunnar Martinsson, and Joel A. Tropp. Finding structure with randomness: probabilistic algorithms for constructing approximate matrix decompositions. *SIAM Review*, 53(2):217–288, January 2011.
- [30] Richard A Harshman et al. Foundations of the parafac procedure: Models and conditions for an “explanatory” multi-modal factor analysis. *UCLA working papers in phonetics*, 16(1):84, 1970.
- [31] Rui Hong, Ya-Xuan Xiao, Jie Hu, An-Chun Ji, and Shi-Ju Ran. Functional tensor network solving many-body Schrödinger equation. *Phys. Rev. B*, 105:165116, Apr 2022.
- [32] L Hörmander. Fourier integral operators. i. In *Mathematics Past and Present Fourier Integral Operators*, pages 23–127. Springer, 1994.
- [33] Yuehaw Khoo and Lexing Ying. Switchnet: a neural network model for forward and inverse

- scattering problems. *SIAM Journal on Scientific Computing*, 41(5):A3182–A3201, 2019.
- [34] Tamara G Kolda and Brett W Bader. Tensor decompositions and applications. *SIAM review*, 51(3):455–500, 2009.
- [35] Matthew Li, Laurent Demanet, and Leonardo Zepeda-Núñez. Wide-band butterfly network: stable and efficient inversion via multi-frequency neural networks. *Multiscale Modeling & Simulation*, 20(4):1191–1227, 2022.
- [36] Yingzhou Li and Haizhao Yang. Interpolative butterfly factorization. *SIAM J. Sci. Comput.*, 39(2):A503–A531, 2017.
- [37] Yingzhou Li, Haizhao Yang, Eileen R Martin, Kenneth L Ho, and Lexing Ying. Butterfly factorization. *Multiscale Model. Sim.*, 13(2):714–732, 2015.
- [38] Yingzhou Li, Haizhao Yang, and Lexing Ying. Multidimensional butterfly factorization. *Applied and Computational Harmonic Analysis*, 44(3):737–758, 2018.
- [39] E. Liberty, F. Woolfe, P.-G. Martinsson, V. Rokhlin, and M. Tygert. Randomized algorithms for the low-rank approximation of matrices. *Proc. Natl. Acad. Sci. USA*, 104:20167–20172, 2007.
- [40] Peter Lindstrom. Fixed-rate compressed floating-point arrays. *IEEE Transactions on Visualization and Computer Graphics*, 20(12):2674–2683, 2014.
- [41] Yang Liu. A comparative study of butterfly-enhanced direct integral and differential equation solvers for high-frequency electromagnetic analysis involving inhomogeneous dielectrics. In *2022 3rd URSI Atlantic and Asia Pacific Radio Science Meeting (AT-AP-RASC)*, pages 1–4. IEEE, 2022.
- [42] Yang Liu, Pieter Ghysels, Lisa Claus, and Xiaoye Sherry Li. Sparse approximate multifrontal factorization with butterfly compression for high-frequency wave equations. *SIAM Journal on Scientific Computing*, 0(0):S367–S391, 2021.
- [43] Yang Liu, Han Guo, and Eric Michielssen. An HSS matrix-inspired butterfly-based direct solver for analyzing scattering from two-dimensional objects. *IEEE Antennas Wirel. Propag. Lett.*, 16:1179–1183, 2017.
- [44] Yang Liu, Tianhuan Luo, Aman Rani, Hengrui Luo, and Xiaoye Sherry Li. Detecting resonance of radio-frequency cavities using fast direct integral equation solvers and augmented bayesian optimization. *IEEE Journal on Multiscale and Multiphysics Computational Techniques*, 2023.
- [45] Yang Liu, Jian Song, Robert Burrige, and Jianliang Qian. A fast butterfly-compressed Hadamard-Babich integrator for high-frequency Helmholtz equations in inhomogeneous media with arbitrary sources. *Multiscale Modeling & Simulation*, 21(1):269–308, 2023.
- [46] Yang Liu, Xin Xing, Han Guo, Eric Michielssen, Pieter Ghysels, and Xiaoye Sherry Li. Butterfly factorization via randomized matrix-vector multiplications. *SIAM Journal on Scientific Computing*, 43(2):A883–A907, 2021.
- [47] Yang Liu and Haizhao Yang. A hierarchical butterfly LU preconditioner for two-dimensional electromagnetic scattering problems involving open surfaces. *J. Comput. Phys.*, 401:109014, 2020.
- [48] W. Lu, J. Qian, and R. Burrige. Babich’s expansion and the fast Huygens sweeping method for the Helmholtz wave equation at high frequencies. *J. Comput. Phys.*, 313:478–510, 2016.
- [49] Axel Maas. Two and three-point Green’s functions in two-dimensional Landau-gauge Yang-Mills theory. *Phys. Rev. D*, 75:116004, 2007.
- [50] Michael W Mahoney, Mauro Maggioni, and Petros Drineas. Tensor-cur decompositions for tensor-based data. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 327–336, 2006.
- [51] Osman Asif Malik and Stephen Becker. Fast randomized matrix and tensor interpolative decomposition using counts sketch. *Advances in Computational Mathematics*, 46(6):76, 2020.
- [52] Eric Michielssen and Amir Boag. Multilevel evaluation of electromagnetic fields for the rapid solution of scattering problems. *Microw Opt Technol Lett.*, 7(17):790–795, 1994.
- [53] Eric Michielssen and Amir Boag. A multilevel matrix decomposition algorithm for analyzing scattering from large structures. *IEEE Trans. Antennas Propag.*, 44(8):1086–1093, 1996.
- [54] Michael O’Neil, Franco Woolfe, and Vladimir Rokhlin. An algorithm for the rapid evaluation of special function transforms. *Appl. Comput. Harmon. A.*, 28(2):203 – 226, 2010. Special Issue on Continuous Wavelet Transform in Memory of Jean Morlet, Part I.
- [55] Ivan V Oseledets. Tensor-train decomposition. *SIAM Journal on Scientific Computing*, 33(5):2295–2317, 2011.
- [56] Qiyuan Pang, Kenneth L. Ho, and Haizhao Yang. Interpolative decomposition butterfly factorization. *SIAM J. Sci. Comput.*, 42(2):A1097–A1115, 2020.
- [57] Michael E Peskin. *An introduction to quantum field theory*. CRC press, 2018.
- [58] Arvind K Saibaba. Hoid: higher order interpolatory decomposition for tensors based on tucker

- representation. *SIAM Journal on Matrix Analysis and Applications*, 37(3):1223–1249, 2016.
- [59] Sadeed Bin Sayed, Yang Liu, Luis J. Gomez, and Abdulkadir C. Yucel. A butterfly-accelerated volume integral equation solver for broad permittivity and large-scale electromagnetic analysis. *IEEE Transactions on Antennas and Propagation*, 70(5):3549–3559, 2022.
- [60] Weitian Sheng, Abdulkadir C Yucel, Yang Liu, Han Guo, and Eric Michielssen. A domain decomposition based surface integral equation simulator for characterizing EM wave propagation in mine environments. *IEEE Transactions on Antennas and Propagation*, 2023.
- [61] Tianyi Shi, Daniel Hayes, and Jing-Mei Qiu. Distributed memory parallel adaptive tensor-train cross approximation. *arXiv preprint arXiv:2407.11290*, 2024.
- [62] Edgar Solomonik, Devin Matthews, Jeff Hammond, and James Demmel. Cyclops tensor framework: Reducing communication and eliminating load imbalance in massively parallel contractions. In *2013 IEEE 27th International Symposium on Parallel and Distributed Processing*, pages 813–824. IEEE, 2013.
- [63] Mark Tygert. Fast algorithms for spherical harmonic expansions, III. *J. Comput. Phys.*, 229(18):6181 – 6192, 2010.
- [64] Mingyu Wang, Cheng Qian, Enrico Di Lorenzo, Luis J Gomez, Vladimir Okhmatovski, and Abdulkadir C Yucel. Supervoxhenry: Tucker-enhanced and fft-accelerated inductance extraction for voxelized superconducting structures. *IEEE Transactions on Applied Superconductivity*, 31(7):1–11, 2021.
- [65] Mingyu Wang, Cheng Qian, Jacob K White, and Abdulkadir C Yucel. Voxcap: Fft-accelerated and tucker-enhanced capacitance extraction simulator for voxelized structures. *IEEE Transactions on Microwave Theory and Techniques*, 68(12):5154–5168, 2020.
- [66] E. Wigner. On the quantum correction for thermodynamic equilibrium. *Phys. Rev.*, 40:749–759, Jun 1932.
- [67] Haizhao Yang. A unified framework for oscillatory integral transforms: When to use NUFFT or butterfly factorization? *J. Comput. Phys.*, 388:103 – 122, 2019.
- [68] Lexing Ying. Sparse Fourier Transform via Butterfly Algorithm. *SIAM J. Sci. Comput.*, 31(3):1678–1694, 2009.