

# Broadcast Product: Redefining Shape-aligned Element-wise Multiplication and Beyond

Anonymous authors  
Paper under double-blind review

## Abstract

Broadcast operations are widely used in scientific computing libraries, yet their mathematical formulation is often implicit and inconsistently represented in machine learning literature. This problem frequently leads to invalid equations when element-wise products are written despite mismatched tensor shapes. In this paper, we formalize such operations by introducing the *broadcast product*  $\boxtimes$ , which explicitly extends the Hadamard product through shape-aligned element duplication. We provide a rigorous definition of the broadcast product, analyze its algebraic properties, and show how it can be expressed using standard linear algebra. Building on this framework, we formulate least-squares problems and **sketch a proof-of-concept broadcast decomposition**. **As a preliminary illustration, we show that the formalism enables a new family of decompositions with distinct structural properties from** conventional tensor decompositions. This work establishes a mathematical foundation for broadcast-aware tensor operations, connecting practical implementations with rigorous tensor analysis.

## 1 Introduction

Broadcast operations are a fundamental abstraction in scientific computing libraries for performing mathematical operations on tensors of different shapes. They automatically replicate elements as needed to align tensor shapes before applying the specified operation, enabling concise and expressive tensor computations. For example, in libraries such as `numpy` (Harris et al., 2020), adding a vector  $\mathbf{x}$  to a matrix  $\mathbf{A}$  implicitly replicates  $\mathbf{x}$  to match the shape of  $\mathbf{A}$ . Moreover, languages such as `MATLAB` and `Julia` (Bezanson et al., 2026) support broadcasting at the language level.

Although we cannot directly write broadcast operations into equations, many modern machine learning papers have attempted to do so by mimicking `numpy`'s descriptions. For example, consider selecting a region of a three-channel image  $\mathcal{X} \in \mathbb{R}^{H \times W \times 3}$  using a binary mask  $\mathbf{B} \in \{0, 1\}^{H \times W}$  to obtain  $\mathcal{Y} \in \mathbb{R}^{H \times W \times 3}$ , where  $H$  and  $W$  denote the height and width, respectively. One often describes this operation by mimicking a `numpy`-style expression,  $\mathcal{Y} = \mathcal{X} * \mathbf{B}$ , and reformulating it using the Hadamard product  $\odot$  as follows:

$$\mathcal{Y} = \mathcal{X} \odot \mathbf{B}. \quad (1)$$

However, this expression is incorrect because the Hadamard product is only defined for tensors of the same shape. Therefore, this equation implicitly assumes the mask is broadcasted along the channel dimension. Such incorrect notation not only creates discrepancies between the source code and mathematical expressions but also leads to invalid mathematical manipulations and flawed discussions.

For example, for any two tensors with the same shape  $\mathcal{M}$  and  $\mathcal{N}$ , the  $L_0$  norm of their Hadamard product must be smaller than that of the original tensor:  $\|\mathcal{M} \odot \mathcal{N}\|_0 \leq \|\mathcal{N}\|_0$ . This is because the result of the Hadamard product can only increase the number of zero elements, never decrease it. However, substituting Equation 1 into this inequality gives  $\|\mathcal{Y}\|_0 = \|\mathcal{X} \odot \mathbf{B}\|_0 \leq \|\mathbf{B}\|_0$ , which is clearly incorrect<sup>1</sup>. In this way,

<sup>1</sup>For example, consider a  $1 \times 1$  three-channel image  $\mathcal{X} \in \mathbb{R}^{1 \times 1 \times 3}$  and a mask  $\mathbf{B} \in \{0, 1\}^{1 \times 1}$ , where the all elements of  $\mathcal{X}$  and  $\mathbf{B}$  are non-zero ( $\|\mathbf{B}\|_0 = 1$ ). In this case, if we apply the incorrect Equation 1, all three elements of  $\mathcal{Y}$  also become non-zero, meaning that we have  $\|\mathcal{Y}\|_0 = 3$ . Thus, the inequality  $\|\mathcal{Y}\|_0 \leq \|\mathbf{B}\|_0$ , which should always hold, is no longer satisfied.

improper use of broadcast operations can lead to mathematically flawed reasoning. Many people tend to write this masking formula, and it has even appeared in award-nominated papers at top conferences (e.g., Figure 5 in Li et al. (2023b) and Equation 2 in Lin et al. (2021))<sup>2</sup>.

To address this issue, we redefine the *broadcast product* operator, denoted by  $\boxtimes$ . This operator extends the Hadamard product  $\odot$  by duplicating elements to align tensor shapes before multiplication. **When tensor shapes are given explicitly (i.e., all axis sizes are specified), it is mathematically equivalent to broadcasting in numpy; the difference arises only when tensor orders differ, due to the distinction between the F-convention and C-convention (formally introduced in Definition 1).** The masking operation above can be expressed as:

$$\mathbf{Y} = \mathbf{X} \boxtimes \mathbf{B} = \mathbf{X} \odot \mathbf{B}^\square. \quad (2)$$

Here,  $\mathbf{B}^\square \in \{0, 1\}^{H \times W \times 3}$  denotes the result of duplicating and concatenating the mask  $\mathbf{B}$  along the channel dimension (as detailed later in Definition 3). Now, the inequality discussed above can be correctly written as:  $\|\mathbf{Y}\|_0 = \|\mathbf{X} \boxtimes \mathbf{B}\|_0 = \|\mathbf{X} \odot \mathbf{B}^\square\|_0 \leq \|\mathbf{B}^\square\|_0 = 3\|\mathbf{B}\|_0$ . The proposed operator  $\boxtimes$  can accurately describe complex problems, potentially offering new interpretations of the problem and possibilities for optimization from alternative perspectives.

Furthermore, we analyze the mathematical properties of the broadcast product and position it within the context of tensor decomposition. We first represent the broadcast product in graphical notation widely used in tensor decomposition (Penrose, 1971; Taylor, 2024; Yokota, 2024). Next, we approximate an arbitrary tensor as the broadcast product of two tensors and minimize the reconstruction error, proposing a solution to the least squares problem based on the broadcast product. This **least-squares** formulation allows us to decompose an arbitrary tensor into multiple tensors via broadcast products, which we term *broadcast decomposition*.

Our contributions are summarized as follows:

- **Notation/formalization contribution:** We propose a notation for explicitly representing broadcast-aware element-wise operations. This fills the gap between informal `numpy`-style descriptions and rigorous mathematical notation.
- **Analytical contribution:** Under this notation, we provide a systematic definition of the broadcast product, summarize its basic algebraic properties, and clarify its relationship to existing operations such as the Hadamard product, Khatri–Rao product, and standard tensor decompositions.
- **Proof-of-concept decomposition:** Building on this framework, we sketch a broadcast decomposition (BD) as a preliminary result rather than a central contribution. We do not claim empirical superiority; rather, we show that BD has distinct structural properties from conventional tensor decompositions.

Following Kolda & Bader (2009), scalars, vectors, matrices, and tensors are denoted by  $x \in \mathbb{R}$ ,  $\mathbf{x} \in \mathbb{R}^I$ ,  $\mathbf{X} \in \mathbb{R}^{I \times J}$ , and  $\mathcal{X} \in \mathbb{R}^{I \times J \times K}$ , respectively. The element-wise (Hadamard) product is denoted by  $\mathcal{X} \odot \mathcal{Y}$ , and the element-wise division is denoted by  $\mathcal{X} \oslash \mathcal{Y}$ . An element of a tensor  $\mathcal{X} \in \mathbb{R}^{I \times J \times K}$  is written as  $x_{ijk} \in \mathbb{R}$ . For a third-order tensor  $\mathcal{X} \in \mathbb{R}^{I \times J \times K}$ , the frontal slice ( $k$ -th channel) is denoted by  $\mathbf{X}_k \in \mathbb{R}^{I \times J}$ . For an  $N$ -th order tensor  $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$  and a matrix  $\mathbf{A} \in \mathbb{R}^{J \times I_n}$ , we denote the  $n$ -th mode tensor-matrix product by  $\mathcal{X} \times_n \mathbf{A} \in \mathbb{R}^{I_1 \times \dots \times I_{n-1} \times J \times I_{n+1} \times \dots \times I_N}$  (we retain the term “mode” in tensor-decomposition terminology, such as the  $n$ -th mode product  $\times_n$  and mode- $n$  unfolding, following Kolda & Bader (2009)). We assume the underlying field to be  $\mathbb{R}$ , although the definitions extend straightforwardly to  $\mathbb{C}$ .

## 1.1 Examples

Before going into the detailed definition, let us introduce the broadcast product intuitively. If two input tensors have identical shapes, their broadcast product equals their Hadamard product, e.g., with  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^2$ :

$$\mathbf{x} \boxtimes \mathbf{y} = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \boxtimes \begin{bmatrix} 3 \\ 4 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \odot \begin{bmatrix} 3 \\ 4 \end{bmatrix} = \begin{bmatrix} 3 \\ 8 \end{bmatrix}. \quad (3)$$

<sup>2</sup>We emphasize that this point concerns notation only and does not affect the proposed algorithms.

$$\begin{array}{c}
\begin{array}{ccc}
\begin{array}{c} X_2 \\ \mathbf{X}_1 \end{array} \begin{bmatrix} 13 & 16 & 19 & 22 \\ 14 & 17 & 20 & 23 \\ 15 & 18 & 21 & 24 \\ 16 & 19 & 20 & 21 \\ 17 & 18 & 19 & 20 \end{bmatrix} & \begin{bmatrix} -1 & 2 & 3 & 4 \\ -5 & 6 & 7 & 8 \\ -9 & 10 & 11 & 12 \end{bmatrix} & = \\
\mathbf{X} \quad \square & \mathbf{Y} & \\
\downarrow \text{Marginalization} & \downarrow \text{Marginalization} & \\
\begin{array}{c} \sqrt{170} \\ \sqrt{200} \\ \sqrt{234} \end{array} \begin{array}{c} \sqrt{272} \\ \sqrt{314} \\ \sqrt{360} \end{array} \begin{array}{c} \sqrt{410} \\ \sqrt{464} \\ \sqrt{522} \end{array} \begin{array}{c} \sqrt{584} \\ \sqrt{650} \\ \sqrt{720} \end{array} & \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \end{bmatrix} & \boxed{\|\mathbf{Z}\|_F = \|\mathbf{X} \square \mathbf{Y}\|_F = \|\mathbf{X}^\square \odot \mathbf{Y}^\square\|_F = \|\mathbf{X}_\square \odot \mathbf{Y}_\square\|_F}
\end{array}
\end{array}$$

Figure 1: Example of the broadcast product of a third-order tensor  $\mathcal{X} \in \mathbb{R}^{3 \times 4 \times 2}$  and a matrix  $\mathbf{Y} \in \mathbb{R}^{3 \times 4}$ .  $\mathcal{X}^\square, \mathbf{Y}^\square \in \mathbb{R}^{3 \times 4 \times 2}$  represent their broadcasted forms.  $\mathbf{X}_\square, \mathbf{Y}_\square \in \mathbb{R}^{3 \times 4}$  illustrate their **Frobenius-norm** marginalization.

Next, we show an example of matrices with different shapes,  $\mathbf{X} \in \mathbb{R}^{3 \times 2}$  and  $\mathbf{y} \in \mathbb{R}^{1 \times 2}$ , as follows:

$$\mathbf{X} \square \mathbf{y} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix} \square [7 \ 8] = \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix} \odot \begin{bmatrix} 7 & 8 \\ 7 & 8 \\ 7 & 8 \end{bmatrix} = \begin{bmatrix} 7 & 16 \\ 21 & 32 \\ 35 & 48 \end{bmatrix}. \quad (4)$$

Here, to match the shapes on both sides of  $\square$ ,  $\mathbf{y}$  is duplicated along the first **axis**, producing a matrix with the same shape as  $\mathbf{X}$  (shown in red). The Hadamard product with  $\mathbf{X}$  is then computed. Note that although the result coincides with the Khatri–Rao product (Kolda & Bader, 2009) in this case (the Khatri–Rao product  $\mathbf{B} \odot_{\text{KR}} \mathbf{A}$  concatenates column-wise Kronecker products of two matrices with the same number of columns), the broadcast product and the Khatri–Rao product are generally different. For instance, the Khatri–Rao product is defined for  $\mathbf{X} \in \mathbb{R}^{3 \times 2}$  and  $\mathbf{Y} \in \mathbb{R}^{2 \times 2}$ , whereas the broadcast product is not (see Section 6.1 for details).

Next, let’s look at an example of the product of a third-order tensor  $\mathcal{X} \in \mathbb{R}^{3 \times 4 \times 2}$  and a matrix  $\mathbf{Y} \in \mathbb{R}^{3 \times 4}$ :

$$\mathbf{X}_1 = \begin{bmatrix} 1 & 4 & 7 & 10 \\ 2 & 5 & 8 & 11 \\ 3 & 6 & 9 & 12 \end{bmatrix}, \quad \mathbf{X}_2 = \begin{bmatrix} 13 & 16 & 19 & 22 \\ 14 & 17 & 20 & 23 \\ 15 & 18 & 21 & 24 \end{bmatrix}, \quad \mathbf{Y} = \begin{bmatrix} -1 & 2 & 3 & 4 \\ -5 & 6 & 7 & 8 \\ -9 & 10 & 11 & 12 \end{bmatrix}. \quad (5)$$

$\mathcal{X} \square \mathbf{Y}$  is equivalent to duplicating  $\mathbf{Y}$  along the third **axis** to match shapes and computing the Hadamard product. This example closely resembles the masking operation discussed in Equation (2), which is illustrated in the upper part of Figure 1, where we will define the notation of  $\mathcal{X}^\square$  and  $\mathbf{Y}^\square$  in Section 2.

## 2 Definition

Let us define the broadcast product in detail. We first introduce two conventions for identifying tensors of different orders.

**Definition 1** (F-convention and C-convention). *When tensors of different orders are compared, a shape-normalization convention is needed to equalize their orders by inserting singleton axes.*

**F-convention** (Fortran-style): *trailing singleton axes are appended. For example,  $\mathbb{R}^{3 \times 4}$  is identified with  $\mathbb{R}^{3 \times 4 \times 1}$  (and recursively with  $\mathbb{R}^{3 \times 4 \times 1 \times 1}$ , etc.). In particular, writing  $\mathbf{x} \in \mathbb{R}^3$  denotes a column vector, i.e.,  $\mathbf{x} \in \mathbb{R}^{3 \times 1}$ . Mathematical notation commonly used in machine learning and computer vision, as well as *Julia*, follow this convention.*

**C-convention** (C-style): leading singleton axes are prepended. For example,  $\mathbb{R}^{3 \times 4}$  is identified with  $\mathbb{R}^{1 \times 3 \times 4}$ . `numpy` and `pytorch` follow this convention.

Throughout this paper we adopt the F-convention. See Appendix F for a detailed comparison and concrete examples of the discrepancy between the two formats.

Following Definition 1, we define the *broadcast condition*, which determines whether the broadcast operation can be applied to two tensors.

**Definition 2** (Broadcast Condition). *Under the F-convention (Definition 1), two tensors  $\mathcal{X}$  and  $\mathcal{Y}$  are first brought to the same order  $N$  by appending trailing singleton axes as needed (e.g.,  $\mathcal{Y} \in \mathbb{R}^{I_1 \times \dots \times I_M}$  with  $M < N$  becomes  $\mathcal{Y} \in \mathbb{R}^{I_1 \times \dots \times I_M \times 1 \times \dots \times 1}$ ). Writing the resulting shapes as  $\mathcal{X} \in \mathbb{R}^{I_1 \times \dots \times I_N}$  and  $\mathcal{Y} \in \mathbb{R}^{J_1 \times \dots \times J_N}$ , the pair  $(\mathcal{X}, \mathcal{Y})$  satisfies the broadcast condition if for every  $n \in \{1, 2, \dots, N\}$ , at least one of the following holds: (1)  $I_n = J_n$ , (2)  $I_n = 1$ , or (3)  $J_n = 1$ .*

For example, the following meet the broadcast condition.

- Same shape:  $\mathcal{X} \in \mathbb{R}^{3 \times 2}, \mathcal{Y} \in \mathbb{R}^{3 \times 2}$ .
- The length of an axis ( $J_2$ ) is one:  $\mathcal{X} \in \mathbb{R}^{3 \times 2}, \mathcal{Y} \in \mathbb{R}^{3 \times 1}$ .
- The lengths of some axes are one ( $I_1 = J_2 = 1$ ):  $\mathcal{X} \in \mathbb{R}^{1 \times 2 \times 5}, \mathcal{Y} \in \mathbb{R}^{3 \times 1 \times 5}$ .
- Multiple ones:  $\mathcal{X} \in \mathbb{R}^{1 \times 1 \times 5}, \mathcal{Y} \in \mathbb{R}^{3 \times 1 \times 5}$ .
- Different orders (F-convention):  $\mathcal{X} \in \mathbb{R}^{5 \times 4 \times 3}, \mathcal{Y} \in \mathbb{R}^{5 \times 4}$  (identified with  $\mathbb{R}^{5 \times 4 \times 1}$ ).

The following do not satisfy the broadcast condition.

- The length of an axis is different and not one ( $I_2 = 2, J_2 = 3$ ):  $\mathcal{X} \in \mathbb{R}^{3 \times 2}, \mathcal{Y} \in \mathbb{R}^{3 \times 3}$ .
- Different orders with first-axis mismatch:  $\mathcal{X} \in \mathbb{R}^{3 \times 2}, \mathcal{Y} \in \mathbb{R}^{4 \times 2 \times 5}$  (under F-convention,  $\mathcal{X}$  is identified with  $\mathbb{R}^{3 \times 2 \times 1}$ ; since  $I_1 = 3 \neq 4 = J_1$  and neither equals 1, the condition fails at axis 1).

Next, we define the *broadcast product* as follows.

**Definition 3** (Broadcast Product). *Let  $\mathcal{X}$  and  $\mathcal{Y}$  be two tensors satisfying the broadcast condition (Definition 2), with their shapes brought to the same order  $N$  under the F-convention (Definition 1). Writing  $\mathcal{X} \in \mathbb{R}^{I_1 \times \dots \times I_N}$  and  $\mathcal{Y} \in \mathbb{R}^{J_1 \times \dots \times J_N}$ , the broadcast product  $\mathcal{Z} = \mathcal{X} \square \mathcal{Y}$  is defined as follows, where  $\mathcal{Z}, \mathcal{X}^\square, \mathcal{Y}^\square \in \mathbb{R}^{\max(I_1, J_1) \times \dots \times \max(I_N, J_N)}$ .*

$$\mathcal{Z} = \mathcal{X} \square \mathcal{Y} := \text{bc}(\mathcal{X}, \text{size}(\mathcal{Y})) \odot \text{bc}(\mathcal{Y}, \text{size}(\mathcal{X})) = \mathcal{X}^\square \odot \mathcal{Y}^\square.$$

Here, “size” returns the input tensor’s shape as a tuple, e.g.,  $\text{size}(\mathcal{Y}) = (J_1, J_2, \dots, J_N)$ , and “bc” is a function to duplicate a tensor. For  $\mathcal{X}$  and  $\mathcal{Y}$  satisfying the broadcast condition, bc inputs  $\mathcal{X}$  itself and the shape of  $\mathcal{Y}$ , and outputs  $\mathcal{X}^\square$ , i.e.,  $\mathcal{X}^\square = \text{bc}(\mathcal{X}, \text{size}(\mathcal{Y}))$ . Here,  $\mathcal{X}^\square$  refers to the broadcasted version of  $\mathcal{X}$ , which means the following. For all  $n$ , if  $I_n = 1$ , replicate all elements of  $\mathcal{X}$  along the  $n$ -th axis  $J_n$  times. This operation is explicitly defined using index notation as follows. With  $k_n \in \{1, 2, \dots, \max(I_n, J_n)\}$  for all  $n$ , we write the  $(k_1, k_2, \dots, k_N)$ -th element of  $\mathcal{X}^\square$  as  $x_{k_1 k_2 \dots k_N}^\square = x_{i'_1 i'_2 \dots i'_N}$ , where

$$i'_n = \begin{cases} 1 & (I_n = 1), \\ k_n & (I_n \neq 1). \end{cases} \quad (6)$$

The same applies to  $\mathcal{Y}^\square$ . In the end, we obtain

$$z_{k_1, k_2, \dots, k_N} = x_{k_1 k_2 \dots k_N}^\square y_{k_1 k_2 \dots k_N}^\square = x_{i'_1 i'_2 \dots i'_N} y_{j'_1 j'_2 \dots j'_N}, \quad \text{where } j'_n = \begin{cases} 1 & (J_n = 1), \\ k_n & (J_n \neq 1). \end{cases} \quad (7)$$

In other words, placing  $\square$  on the superscript means duplicating the elements appropriately so that the shapes of  $\mathcal{X}^\square$  and  $\mathcal{Y}^\square$  match, i.e.,  $\text{size}(\mathcal{X}^\square) = \text{size}(\mathcal{Y}^\square)$ . Note that  $\mathcal{X}^\square$  and  $\mathcal{Y}^\square$  are considered as a shorthand notation, and should only be used when the interpretation is unique and obvious from the context.

The result of the broadcast product is uniquely determined. In the example of Equation (4), we can write:

$$\mathbf{X}^\square = \mathbf{X} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix}, \quad \mathbf{Y}^\square = \begin{bmatrix} \mathbf{y} \\ \mathbf{y} \\ \mathbf{y} \end{bmatrix} = \begin{bmatrix} 7 & 8 \\ 7 & 8 \\ 7 & 8 \end{bmatrix}. \quad (8)$$

In the example of Equation (5) and Figure 1, we obtain

$$\mathcal{X}^\square = \mathcal{X}, \quad \mathbf{Y}_1^\square = \mathbf{Y}_2^\square = \mathbf{Y} = \begin{bmatrix} -1 & 2 & 3 & 4 \\ -5 & 6 & 7 & 8 \\ -9 & 10 & 11 & 12 \end{bmatrix}, \quad (9)$$

where stacking  $\mathbf{Y}_1^\square$  and  $\mathbf{Y}_2^\square$  along the third axis leads  $\mathcal{Y}^\square$ . The following is an example of duplication for both  $\mathcal{X}$  and  $\mathcal{Y}$ . Let us define  $\mathcal{X} \in \mathbb{R}^{1 \times 2 \times 3}$  and  $\mathcal{Y} \in \mathbb{R}^{4 \times 2 \times 1}$ :

$$\mathbf{X}_1 = [1, 2], \quad \mathbf{X}_2 = [3, 4], \quad \mathbf{X}_3 = [5, 6], \quad \mathbf{Y}_1 = \begin{bmatrix} 7 & 8 \\ 9 & 10 \\ 11 & 12 \\ 13 & 14 \end{bmatrix}. \quad (10)$$

In this case,  $\mathcal{X}^\square, \mathcal{Y}^\square \in \mathbb{R}^{4 \times 2 \times 3}$  are written as

$$\mathbf{X}_1^\square = \begin{bmatrix} 1 & 2 \\ 1 & 2 \\ 1 & 2 \\ 1 & 2 \end{bmatrix}, \quad \mathbf{X}_2^\square = \begin{bmatrix} 3 & 4 \\ 3 & 4 \\ 3 & 4 \\ 3 & 4 \end{bmatrix}, \quad \mathbf{X}_3^\square = \begin{bmatrix} 5 & 6 \\ 5 & 6 \\ 5 & 6 \\ 5 & 6 \end{bmatrix}, \quad \mathbf{Y}_1^\square = \mathbf{Y}_2^\square = \mathbf{Y}_3^\square = \begin{bmatrix} 7 & 8 \\ 9 & 10 \\ 11 & 12 \\ 13 & 14 \end{bmatrix}. \quad (11)$$

One alternative to  $\square$  is to write every expression explicitly using the sup-box:  $\mathcal{X}^\square \odot \mathcal{Y}^\square$ . This is always mathematically correct and may be clearer when the broadcast structure of a particular expression is the focus. The  $\square$  operator is preferable when the same broadcast pattern recurs across multiple formulas, since it attaches the broadcasting semantics to the operation itself rather than the operands. For example,  $\mathcal{X} \square \mathcal{Y}$  remains valid regardless of whether  $\mathcal{X}$  has one or many channels, without rewriting the formula each time.

### 3 Properties

We present various mathematical properties of the broadcast product. These properties help convert mathematical expressions using the broadcast product into the standard forms used in linear algebra.

#### 3.1 Basic properties

We first show the basic properties of the broadcast product  $\square$ .

**Compatibility:** When  $\mathcal{X}$  and  $\mathcal{Y}$  have the same shapes,  $\square$  is equivalent to  $\odot$ .

$$\mathcal{X} \square \mathcal{Y} = \mathcal{X} \odot \mathcal{Y}. \quad (12)$$

**Commutativity:** For  $\mathcal{X}$ ,  $\mathcal{Y}$ , and  $\mathbf{0}$  satisfying the broadcast conditions, the following holds, where  $k \in \mathbb{R}$ .

$$\mathcal{X} \square \mathcal{Y} = \mathcal{Y} \square \mathcal{X}. \quad (13)$$

$$(k\mathcal{X}) \square \mathcal{Y} = \mathcal{X} \square (k\mathcal{Y}) = k(\mathcal{X} \square \mathcal{Y}). \quad (14)$$

$$\mathcal{X} \square \mathbf{0} = \mathbf{0} \square \mathcal{X} = \mathbf{0}. \quad (15)$$

**Associativity:** When  $\mathcal{X}$ ,  $\mathcal{Y}$ , and  $\mathcal{Z}$  pairwise satisfy the broadcast conditions, we obtain

$$(\mathcal{X} \boxtimes \mathcal{Y}) \boxtimes \mathcal{Z} = \mathcal{X} \boxtimes (\mathcal{Y} \boxtimes \mathcal{Z}) = \mathcal{X} \boxtimes \mathcal{Y} \boxtimes \mathcal{Z}. \quad (16)$$

Be careful that, even if  $(\mathcal{X}, \mathcal{Y})$  and  $(\mathcal{X}, \mathcal{Z})$  satisfy the broadcast conditions,  $(\mathcal{Y}, \mathcal{Z})$  do not necessarily satisfy the broadcast conditions, e.g.,  $\mathcal{X} \in \mathbb{R}^{3 \times 1}$ ,  $\mathcal{Y} \in \mathbb{R}^{3 \times 2}$ , and  $\mathcal{Z} \in \mathbb{R}^{3 \times 5}$ .

**Distributivity:** When  $\mathcal{Y}$  and  $\mathcal{Z}$  have identical shapes, and  $\mathcal{X}$  and  $\mathcal{Y}$  satisfy the broadcast condition,

$$\mathcal{X} \boxtimes (\mathcal{Y} + \mathcal{Z}) = \mathcal{X} \boxtimes \mathcal{Y} + \mathcal{X} \boxtimes \mathcal{Z}. \quad (17)$$

Even if  $(\mathcal{X}, \mathcal{Y})$  and  $(\mathcal{X}, \mathcal{Z})$  satisfy the broadcast conditions, differing shapes of  $\mathcal{Y}$  and  $\mathcal{Z}$  make the left-hand side uncomputable, though the right-hand side may be; e.g.,  $\mathcal{X} \in \mathbb{R}^{2 \times 3}$ ,  $\mathcal{Y} \in \mathbb{R}^{2 \times 1}$ , and  $\mathcal{Z} \in \mathbb{R}^{1 \times 3}$ .

**Generalization of the Hadamard Product:** Whenever the Hadamard product is well-defined (i.e., the operands have identical shapes), it can be replaced with the broadcast product. That is, the equation  $\mathcal{A} \odot \mathcal{B}$  can always be rewritten as  $\mathcal{A} \boxtimes \mathcal{B}$ . This is because the use of the Hadamard product implies that  $\mathcal{A}$  and  $\mathcal{B}$  have the same shape, allowing it to be rewritten as the broadcast product via Equation (12). Note that, however, the Hadamard product is often misused as discussed in Equation (1).

### 3.2 Frobenius-norm Marginalization and the Frobenius norm

We introduce **Frobenius-norm** marginalization for the broadcast product, enabling the efficient computation of the Frobenius norm and the derivation of the Cauchy–Schwarz inequality. We first define **Frobenius-norm marginalization** formally. Here, we write  $:$  to denote the full index range of an axis, following MATLAB/Julia convention (Golub & Van Loan, 2012) (e.g.,  $\mathbf{X}_{:,j}$  is the  $j$ -th column of  $\mathbf{X}$ ).

**Definition 4** (Frobenius-norm Marginalization). *When two  $N$ -th order tensors  $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$  and  $\mathcal{Y} \in \mathbb{R}^{J_1 \times J_2 \times \dots \times J_N}$  satisfy the broadcast condition, the marginalized tensors are written as*

$$\mathcal{X}_{\square}, \mathcal{Y}_{\square} \in \mathbb{R}^{\min(I_1, J_1) \times \min(I_2, J_2) \times \dots \times \min(I_N, J_N)}.$$

*These are defined using index notation. With  $k_n \in \{1, 2, \dots, \min(I_n, J_n)\}$  for all  $n$ ,*

$$x_{\square k_1 k_2 \dots k_N} = \|\mathcal{X}_{\bar{i}_1 \bar{i}_2 \dots \bar{i}_N}\|_F, \quad \bar{i}_n = \begin{cases} : & (J_n = 1), \\ k_n & (J_n \neq 1). \end{cases}$$

That is, the marginalized tensor  $\mathcal{X}_{\square}$  is obtained by taking the Frobenius norm along each axis of  $\mathcal{X}$  if the axis’s length is longer than that of  $\mathcal{Y}$ . Here,  $\mathcal{Y}_{\square}$  is defined similarly. This **Frobenius-norm marginalization** (hereafter simply *marginalization*) can also be regarded as a process that shrinks the shape of a tensor while preserving the Frobenius norm. Analogously to marginalizing a probability distribution (which preserves the  $L_1$  norm), our operation reduces a tensor by one or more axes while preserving its Frobenius norm. As is clear from the construction, we have  $\|\mathcal{X}\|_F = \|\mathcal{X}_{\square}\|_F \leq \|\mathcal{X}^{\square}\|_F$  in general. Also, all elements of the marginalized tensor are non-negative:  $\mathcal{X}_{\square} \geq \mathbf{0}$ . In the example of Equation (4), we obtain

$$\mathbf{x}_{\square} = [\sqrt{35} \quad \sqrt{56}], \quad \mathbf{y}_{\square} = [7 \quad 8]. \quad (18)$$

Figure 1 shows the example of Equation (5). As this example of  $\mathbf{Y}$  shows, negative elements must become positive even if the shape does not change. For Equation (10), we obtain

$$\mathbf{x}_{\square} = [\sqrt{35} \quad \sqrt{56}], \quad \mathbf{y}_{\square} = [\sqrt{420} \quad \sqrt{504}]. \quad (19)$$

Using the marginalized tensors, we can write the Frobenius norm of the broadcast product as follows:

$$\|\mathcal{X} \boxtimes \mathcal{Y}\|_F = \|\mathcal{X}^{\square} \odot \mathcal{Y}^{\square}\|_F = \|\mathcal{X}_{\square} \odot \mathcal{Y}_{\square}\|_F. \quad (20)$$

With this, one can compute the Frobenius norm efficiently. Computing the norm requires “enlarging” the tensors via  $\mathcal{X}^{\square}$ , but one can compute it using smaller tensors by first “shrinking” them via  $\mathcal{X}_{\square}$  (see Figure 1).

From this, we can derive that the Cauchy–Schwarz inequality also holds for the broadcast product:

**Theorem 3.1** (Cauchy–Schwarz Inequality for the Broadcast Product). *Let  $\mathcal{X}$  and  $\mathcal{Y}$  be tensors satisfying the broadcast conditions. Then:*

$$\|\mathcal{X} \boxtimes \mathcal{Y}\|_F \leq \|\mathcal{X}\|_F \cdot \|\mathcal{Y}\|_F.$$

*Proof.* This follows directly from Equation (20) and the standard Cauchy–Schwarz inequality. See Appendix A for the complete proof.  $\square$

### 3.3 Properties of lower-order tensors

We describe the properties for lower-order tensors. By employing these transformations, one can express the description of the broadcast product using standard linear algebra notation.

**Scalar:** For any tensor  $\mathcal{X}$  and a scalar  $y \in \mathbb{R}$ , we obtain

$$\mathcal{X} \boxtimes y = y\mathcal{X}. \quad (21)$$

**Vector and vector:** For vectors  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^I$  of equal length, the following holds:

$$\mathbf{x} \boxtimes \mathbf{y} = \mathbf{x} \odot \mathbf{y} = \text{diag}(\mathbf{x})\mathbf{y} \in \mathbb{R}^I. \quad (22)$$

$$\mathbf{x} \boxtimes \mathbf{y}^\top = \mathbf{x}\mathbf{y}^\top \in \mathbb{R}^{I \times I}. \quad (23)$$

For vectors  $\mathbf{x} \in \mathbb{R}^I$  and  $\mathbf{y} \in \mathbb{R}^J$  of different lengths,  $\mathbf{x} \boxtimes \mathbf{y}$  cannot be defined, but  $\mathbf{x} \boxtimes \mathbf{y}^\top$  can be defined:

$$\mathbf{x} \boxtimes \mathbf{y}^\top = \mathbf{x}\mathbf{y}^\top \in \mathbb{R}^{I \times J}. \quad (24)$$

**Matrix and vector:** For a matrix  $\mathbf{X} \in \mathbb{R}^{I \times J}$ , vectors  $\mathbf{y} \in \mathbb{R}^I$  and  $\mathbf{z} \in \mathbb{R}^J$ , the following holds:

$$\mathbf{X} \boxtimes \mathbf{y} = \mathbf{X} \odot [\mathbf{y}, \mathbf{y}, \dots, \mathbf{y}] = \mathbf{X} \odot (\mathbf{y}\mathbf{1}_J^\top) = \text{diag}(\mathbf{y})\mathbf{X}. \quad (25)$$

$$\mathbf{X} \boxtimes \mathbf{z}^\top = \mathbf{X} \odot \begin{bmatrix} \mathbf{z}^\top \\ \vdots \\ \mathbf{z}^\top \end{bmatrix} = \mathbf{X} \odot (\mathbf{1}_I \mathbf{z}^\top) = \mathbf{X} \text{diag}(\mathbf{z}). \quad (26)$$

Additionally, the following holds for the Frobenius norm:

$$\|\mathbf{X} \boxtimes \mathbf{y}\|_F^2 = \|\text{diag}(\mathbf{y})\mathbf{X}\|_F^2 = \text{tr}(\mathbf{X}^\top \text{diag}^2(\mathbf{y})\mathbf{X}). \quad (27)$$

**Third-order tensor and matrix:** For a third-order tensor  $\mathcal{X} \in \mathbb{R}^{I \times J \times K}$  and a matrix  $\mathbf{Y} \in \mathbb{R}^{I \times J}$ , the following holds. **Here,  $\text{fold}_n$  denotes the inverse of mode- $n$  unfolding; see Appendix B for the definition and a visual illustration.** Figure 1 is a visual example.

$$\mathcal{X} \boxtimes \mathbf{Y} = \mathcal{X} \odot \text{fold}_1(\mathbf{1}_K^\top \otimes \mathbf{Y}) = \text{fold}_3(\mathbf{X}_{(3)} \text{diag}(\text{vec}(\mathbf{Y}))). \quad (28)$$

This operation is typical in computer vision, e.g.,  $\mathcal{X}$  represents an image and  $\mathbf{Y}$  serves as a grayscale mask.

### 3.4 Graphical notations

Here, we show how to draw the broadcast product in graphical notations of tensors (Penrose, 1971; Taylor, 2024; Yokota, 2024). Here, we draw tensors as nodes with multiple edges, represented by black circles and lines. Additionally, we draw super-diagonal tensors using slashed circles,  $\oslash$ , instead of black circles for arbitrary tensors. **An  $N$ -th order tensor  $\mathcal{I} \in \mathbb{R}^{I \times \dots \times I}$  is called *super-diagonal* if its  $(k_1, k_2, \dots, k_N)$ -th element equals 1 when  $k_1 = k_2 = \dots = k_N$  and 0 otherwise; this generalises the identity matrix to higher orders.** For example, the third-order super-diagonal tensor  $\mathcal{I} \in \mathbb{R}^{3 \times 3 \times 3}$  is given by

$$\mathbf{I}_1 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}, \quad \mathbf{I}_2 = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}, \quad \mathbf{I}_3 = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}. \quad (29)$$

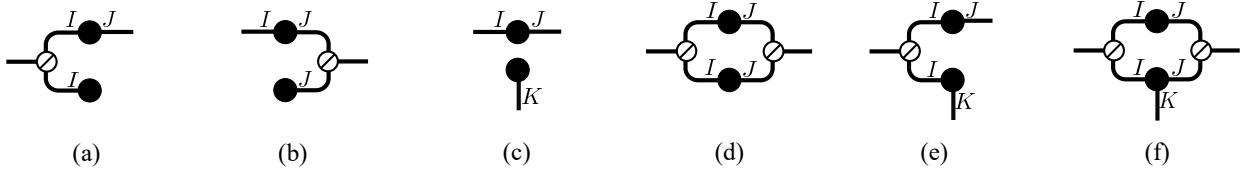


Figure 2: Graphical notations of various patterns of broadcast product: Diagrams (a)-(f) represent broadcast products between  $(I, J)$ -matrix with (a)  $(I, 1)$ -vector, (b)  $(1, J)$ -vector, (c)  $(1, 1, K)$ -vector, (d)  $(I, J)$ -matrix, (e)  $(I, 1, K)$ -matrix, and (f)  $(I, J, K)$ -tensor, respectively. Especially, case (c) is equivalent to outer product, and case (d) equivalent to Hadamard product.

Figure 2 shows the graphical notation of various patterns of the broadcast product. In this way, the broadcast product can be represented as a product (network) of tensors via super-diagonal tensors.

Figure 2a represents a case of Equation (25). Since a tensor product of a super diagonal tensor with a vector is a diagonal matrix

$$\mathcal{I} \times_3 \mathbf{y}^\top = \text{diag}(\mathbf{y}), \quad (30)$$

the property Equation (25) can be immediately verified. The property Equation (26) can also be verified in a similar way. Figure 2f represents a case of Equation (28). For the proof, see Appendix B. **Although Penrose graphical notation (Penrose, 1971) has not yet been widely adopted in mainstream machine learning, it is the standard language in tensor-network communities spanning computational physics (Ran et al., 2020) and machine learning (Stoudenmire & Schwab, 2016; Cichocki et al., 2016; 2017; Li et al., 2023a; Núñez Fernández et al., 2025). We adopt it here because it makes the axis structure and contractions visually explicit, enabling an intuitive understanding of the broadcast product and broadcast decomposition in Section 5.2.**

### 3.5 Broadcast sum, difference, and division

The broadcast sum ( $\boxplus$ ), difference ( $\boxminus$ ), and division ( $\boxdiv$ ) are similarly defined for  $\mathcal{X}$  and  $\mathcal{Y}$  that meet the broadcast condition:

$$\mathcal{X} \boxplus \mathcal{Y} := \mathcal{X}^\square + \mathcal{Y}^\square, \quad \mathcal{X} \boxminus \mathcal{Y} := \mathcal{X}^\square - \mathcal{Y}^\square, \quad \mathcal{X} \boxdiv \mathcal{Y} := \mathcal{X}^\square \oslash \mathcal{Y}^\square, \quad (31)$$

where  $\mathcal{Y}$  must not have zero elements for  $\boxdiv$ . See Appendix C for the basic properties of these operators.

## 4 Real-world Examples

This section demonstrates that the broadcast product naturally formalizes a wide range of operations commonly described using informal or ambiguous notation in machine learning. Through the broadcast operations, various expressions can be represented in a unified manner, offering new perspectives on them.

### 4.1 Blending

Similar to the masking operation in Equation (2), multi-channel blending can be expressed using the broadcast product. Consider two color images  $\mathcal{X}, \mathcal{Y} \in \mathbb{R}^{H \times W \times 3}$  with height  $H$  and width  $W$  and a mask  $\mathbf{A} \in [0, 1]^{H \times W}$  representing an alpha channel. The blending operation is written as

$$\mathcal{X} \boxdiv \mathbf{A} + \mathcal{Y} \boxdiv (\mathbf{1} - \mathbf{A}), \quad (32)$$

where  $\mathbf{1}$  denotes an all-ones matrix of the same shape as  $\mathbf{A}$ . This formulation naturally extends to batched inputs  $\mathcal{X}, \mathcal{Y} \in \mathbb{R}^{N \times H \times W \times C}$ , where  $N$  denotes the batch size, by defining the mask as  $\mathcal{A} \in [0, 1]^{1 \times H \times W}$  (or equivalently  $[0, 1]^{1 \times H \times W \times 1}$ ), and applying Equation (32) directly.

## 4.2 Batch Normalization

Batch Normalization (Ioffe & Szegedy, 2015), Layer Normalization (Ba et al., 2016), and Instance Normalization (Ulyanov et al., 2016) can all be described in a unified manner using broadcast operations. We first demonstrate that the following commonly used notation is not mathematically rigorous:

$$\text{BN}(x) = \left( \frac{x - \mu}{\sigma} \right) \gamma + \beta. \quad (33)$$

Here,  $x \in \mathbb{R}^{N \times C \times H \times W}$  denotes the input tensor, where  $N$ ,  $C$ ,  $H$ , and  $W$  are the batch size, number of channels, height, and width, respectively, following the channel-second convention of Wu & He (2018). Here,  $\mu, \sigma \in \mathbb{R}^C$  denote the batch-wise mean and standard deviation, and  $\gamma, \beta \in \mathbb{R}^C$  are learnable parameters.

This expression is mathematically incorrect, as tensor operations with different shapes rely on implicit broadcasting, e.g., we cannot divide a tensor by a vector ( $\sigma$ ).

Let us now rewrite Batch Normalization explicitly **using index notation**. Let the input and output tensors be  $\mathcal{X}, \mathcal{Y} \in \mathbb{R}^{N \times C \times H \times W}$ . Batch Normalization is defined as follows. For each channel  $c \in \{1, 2, \dots, C\}$ , let us define the mean  $\mu_c \in \mathbb{R}$  and standard deviation  $\sigma_c \in \mathbb{R}$  by:

$$\mu_c = \frac{1}{NHW} \sum_{n=1}^N \sum_{h=1}^H \sum_{w=1}^W X_{n,c,h,w}, \quad \sigma_c = \sqrt{\frac{1}{NHW} \sum_{n=1}^N \sum_{h=1}^H \sum_{w=1}^W (X_{n,c,h,w} - \mu_c)^2 + \varepsilon}. \quad (34)$$

Here,  $\varepsilon$  is a small constant introduced to prevent division by zero. Using the channel-wise learned parameters  $\gamma_1, \dots, \gamma_C \in \mathbb{R}$  and  $\beta_1, \dots, \beta_C \in \mathbb{R}$ , Batch Normalization can be written element-wise as:

$$Y_{n,c,h,w} = \left( \frac{X_{n,c,h,w} - \mu_c}{\sigma_c} \right) \gamma_c + \beta_c. \quad (35)$$

This representation is rigorous. However, because it focuses on individual elements, the notation becomes cluttered with indices. The resulting tensor  $\mathcal{Y}$  cannot be easily substituted into other equations.

By contrast, using the broadcast operations, Equation (35) can be concisely expressed at the tensor level as:

$$\mathcal{Y} = ((\mathcal{X} \boxminus \boldsymbol{\mu}) \boxdiv \boldsymbol{\sigma}) \boxtimes \boldsymbol{\gamma} \boxplus \boldsymbol{\beta}. \quad (36)$$

Here,  $\boldsymbol{\mu} \in \mathbb{R}^{1 \times C \times 1 \times 1}$  is the tensor containing  $\mu_1, \dots, \mu_C$ , and  $\boldsymbol{\sigma}, \boldsymbol{\gamma}, \boldsymbol{\beta} \in \mathbb{R}^{1 \times C \times 1 \times 1}$  are defined similarly. Now, Batch Normalization is written in a mathematically rigorous manner using broadcast operations.

**Under the F-convention (Definition 1), trailing singleton axes can be omitted; for example,  $\boldsymbol{\mu} \in \mathbb{R}^{1 \times C \times 1 \times 1}$  can be written as  $\boldsymbol{\mu} \in \mathbb{R}^{1 \times C}$ .** For clarity and consistency, we explicitly specify all shapes in the upcoming normalization examples.

## 4.3 Layer Normalization (for CNN)

Next, let us consider Layer Normalization (Ba et al., 2016). Here, we follow the notation convention used in Wu & He (2018). Let the input tensor be  $\mathcal{X} \in \mathbb{R}^{N \times C \times H \times W}$ . In Layer Normalization for CNNs, the mean over all elements of each instance  $n \in \{1, \dots, N\}$  is defined as  $\mu_n = \frac{1}{CHW} \sum_{c=1}^C \sum_{h=1}^H \sum_{w=1}^W X_{n,c,h,w}$ . The standard deviation is defined as  $\sigma_n = \sqrt{\frac{1}{CHW} \sum_{c=1}^C \sum_{h=1}^H \sum_{w=1}^W (X_{n,c,h,w} - \mu_n)^2 + \varepsilon}$ . Stacking these values gives  $\boldsymbol{\mu}, \boldsymbol{\sigma} \in \mathbb{R}^{N \times 1 \times 1 \times 1}$ . As in Batch Normalization, the per-channel learnable parameters are  $\boldsymbol{\gamma}, \boldsymbol{\beta} \in \mathbb{R}^{1 \times C \times 1 \times 1}$ . The output tensor is computed using the same expression as in Equation (36).

This formulation involves the input tensor  $\mathcal{X} \in \mathbb{R}^{N \times C \times H \times W}$ , the means  $\boldsymbol{\mu} \in \mathbb{R}^{N \times 1 \times 1 \times 1}$ , and the affine parameters  $\boldsymbol{\gamma} \in \mathbb{R}^{1 \times C \times 1 \times 1}$ . Despite these differing shapes, the entire computation can be written in a single equation using broadcast operations. This directly and rigorously reflects the actual implementation used in libraries such as `numpy` and `pytorch`.

#### 4.4 Layer Normalization (for Transformer)

Layer Normalization is also widely used in Transformer models (Phuong & Hutter, 2022), but its usage differs slightly from the CNN case. Let the input tensor be  $\mathcal{X} \in \mathbb{R}^{N \times L \times D}$ , where  $N$  is the batch size,  $L$  is the number of tokens, and  $D$  is the embedding dimension. Unlike the CNN case, where normalization is performed over all elements of an instance, Transformer-style Layer Normalization computes statistics only over the embedding dimension. For each  $n \in \{1, \dots, N\}$  and  $l \in \{1, \dots, L\}$ , the mean is defined as  $\mu_{n,l} = \frac{1}{D} \sum_{d=1}^D X_{n,l,d}$ . The standard deviation is defined as  $\sigma_{n,l} = \sqrt{\frac{1}{D} \sum_{d=1}^D (X_{n,l,d} - \mu_{n,l})^2} + \varepsilon$ . Stacking these values gives  $\boldsymbol{\mu}, \boldsymbol{\sigma} \in \mathbb{R}^{N \times L \times 1}$ . In this setting, the learnable parameters are defined per embedding dimension as  $\boldsymbol{\gamma}, \boldsymbol{\beta} \in \mathbb{R}^{1 \times 1 \times D}$ . The output tensor is again computed using the same expression in Equation (36).

#### 4.5 Instance Normalization

Finally, let us consider Instance Normalization (Ulyanov et al., 2016). Let the input tensor be  $\mathcal{X} \in \mathbb{R}^{N \times C \times H \times W}$ . The mean over the spatial dimensions is defined as  $\mu_{n,c} = \frac{1}{HW} \sum_{h=1}^H \sum_{w=1}^W X_{n,c,h,w}$ . The standard deviation is defined as  $\sigma_{n,c} = \sqrt{\frac{1}{HW} \sum_{h=1}^H \sum_{w=1}^W (X_{n,c,h,w} - \mu_{n,c})^2} + \varepsilon$ . Stacking these values gives  $\boldsymbol{\mu}, \boldsymbol{\sigma} \in \mathbb{R}^{N \times C \times 1 \times 1}$ . With the learnable parameters  $\boldsymbol{\gamma}, \boldsymbol{\beta} \in \mathbb{R}^{1 \times C \times 1 \times 1}$ , the output of Instance Normalization is computed using the same expression Equation (36) as for Batch Normalization and Layer Normalization.

In this way, Batch, Layer, and Instance Normalization can all be written in a unified and mathematically rigorous form by appropriately choosing the shapes of the variables. In contrast, the commonly used abbreviated notation in Equation (33) cannot precisely express these relationships and introduces ambiguity. For example, the fact that  $\boldsymbol{\sigma}$  is a vector rather than a scalar is not apparent from that notation alone.

#### 4.6 AdaIN

Adaptive Instance Normalization (AdaIN) (Huang & Belongie, 2017) is an extension of Instance Normalization that is widely used in style transfer. Using broadcast operations, AdaIN can be written in the same form as the normalization expression in Equation (36).

Let the input tensor be  $\mathcal{X} \in \mathbb{R}^{N \times C \times H \times W}$ . The transferred style information is represented by channel-wise affine parameters, which are given as  $\boldsymbol{\gamma}, \boldsymbol{\beta} \in \mathbb{R}^{1 \times C \times 1 \times 1}$ . As in standard Instance Normalization, the mean and standard deviation of the input tensor are defined per instance and per channel, yielding  $\boldsymbol{\mu}, \boldsymbol{\sigma} \in \mathbb{R}^{N \times C \times 1 \times 1}$ . With these definitions, Adaptive Instance Normalization is computed using the same broadcast-based expression as in Equation (36).

#### 4.7 FiLM

The final example is Feature-wise Linear Modulation (FiLM) (Perez et al., 2018), widely used in image generation:  $\text{FiLM}(\mathbf{F}_{i,c} \mid \gamma_{i,c}, \beta_{i,c}) = \gamma_{i,c} \mathbf{F}_{i,c} + \beta_{i,c}$ . Here,  $\mathbf{F}_{i,c}$  denotes the  $c$ -th channel response of the  $i$ -th instance, and  $\gamma_{i,c}$  and  $\beta_{i,c}$  are scaling and shifting weights. The operation is difficult to express in matrix notation, making the expression cumbersome.

With the broadcast product, we can intuitively express FiLM. Let  $\mathcal{F} \in \mathbb{R}^{N \times C \times H \times W}$  be a feature volume with the number of instance  $N$ ,  $C$  channels, height  $H$ , width  $W$ . The weights are represented as  $\boldsymbol{\gamma}, \boldsymbol{\beta} \in \mathbb{R}^{N \times C \times 1 \times 1}$ . FiLM can be expressed as:

$$\text{FiLM}(\mathcal{F} \mid \boldsymbol{\gamma}, \boldsymbol{\beta}) = \mathcal{F} \boxtimes \boldsymbol{\gamma} \boxplus \boldsymbol{\beta}. \quad (37)$$

This formulation removes explicit indexing and places FiLM in the same broadcast-based framework as normalization layers, revealing FiLM as an instance transformation without normalization. We can see the clear connection to Equation (36), where  $\boldsymbol{\mu}$  and  $\boldsymbol{\sigma}$  are removed. FiLM thus belongs to the same class of broadcast-based transformations as normalization layers, differing in the absence of normalization statistics.

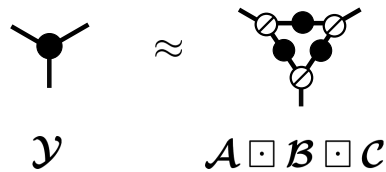


Figure 3: Broadcast decomposition

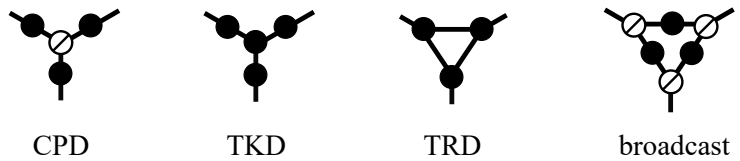


Figure 4: Graphical notation of canonical polyadic decomposition (CPD), Tucker decomposition (TKD), tensor ring decomposition (TRD), and broadcast decomposition (BD).

## 5 Optimizations

Up to this point, we have defined the broadcast product and its basic properties. In this section, we delve into the broadcast product in greater detail, particularly in the context of tensor decomposition, and address the new optimization challenges that arise from it, including matrix factorization techniques and their applications in dimensionality reduction.

### 5.1 Least squares

Let us consider three tensors  $\mathcal{Y} \in \mathbb{R}^{I \times J \times K}$ ,  $\mathcal{A} \in \mathbb{R}^{I \times J \times 1}$ ,  $\mathcal{Z} \in \mathbb{R}^{1 \times J \times K}$  and the following least squares (LS) problem:

$$\underset{\mathcal{A}}{\text{minimize}} \|\mathcal{Y} - \mathcal{A} \boxtimes \mathcal{Z}\|_F^2, \quad (38)$$

then the solution can be given by

$$\hat{\mathcal{A}} = \mathcal{P}_3(\mathcal{Y} \boxtimes \mathcal{Z}) \boxtimes \mathcal{P}_3(\mathcal{Z} \boxtimes \mathcal{Z}), \quad (39)$$

where  $\mathcal{P}_k(\mathcal{X}) := \mathcal{X} \times_k \mathbf{1}^\top$  performs a sum of the entries of an input tensor along the  $k$ -th mode. The proof is in Appendix D. This least-squares solution can be easily generalized to  $N$ -th order tensors by simply changing the **axes** of  $\mathcal{P}$  to match the shape of  $\mathcal{A}$  (i.e., summing along the **axes** with length 1 of  $\mathcal{A}$ ).

### 5.2 Tensor decomposition

We emphasize that broadcast decomposition is a proof-of-concept intended to illustrate the mathematical productivity of the broadcast product operator, not a mature decomposition method. We propose a new tensor decomposition, called broadcast decomposition (BD), based on broadcast products:

$$\mathcal{Y} \approx \mathcal{A} \boxtimes \mathcal{B} \boxtimes \mathcal{C}, \quad (40)$$

where  $\mathcal{A}$ ,  $\mathcal{B}$ , and  $\mathcal{C}$  mutually satisfy the broadcast conditions. For minimizing squared errors  $\|\mathcal{Y} - \mathcal{A} \boxtimes \mathcal{B} \boxtimes \mathcal{C}\|_F^2$ , the alternating least squares (ALS) algorithm can be easily derived using Equation (39). For example, when updating  $\mathcal{A}$ , set  $\mathcal{Z} = \mathcal{B} \boxtimes \mathcal{C}$  and make  $\mathcal{P}$  correspond to the shape of  $\mathcal{A}$ . Let us consider the sizes of tensors as  $\mathcal{A} \in \mathbb{R}^{I \times J \times 1}$ ,  $\mathcal{B} \in \mathbb{R}^{I \times 1 \times K}$ ,  $\mathcal{C} \in \mathbb{R}^{1 \times J \times K}$ , then these update rules<sup>3</sup> can be given by

$$\mathcal{A} \leftarrow \mathcal{P}_3(\mathcal{Y} \boxtimes \mathcal{B} \boxtimes \mathcal{C}) \oslash \mathcal{P}_3(\mathcal{B} \boxtimes \mathcal{B} \boxtimes \mathcal{C} \boxtimes \mathcal{C}); \quad (41)$$

$$\mathcal{B} \leftarrow \mathcal{P}_2(\mathcal{Y} \boxtimes \mathcal{A} \boxtimes \mathcal{C}) \oslash \mathcal{P}_2(\mathcal{A} \boxtimes \mathcal{A} \boxtimes \mathcal{C} \boxtimes \mathcal{C}); \quad (42)$$

$$\mathcal{C} \leftarrow \mathcal{P}_1(\mathcal{Y} \boxtimes \mathcal{A} \boxtimes \mathcal{B}) \oslash \mathcal{P}_1(\mathcal{A} \boxtimes \mathcal{A} \boxtimes \mathcal{B} \boxtimes \mathcal{B}); \quad (43)$$

Here, the **axis for the sum operation**  $\mathcal{P}$  is determined according to the shape of the update tensor. It corresponds to the **axis** with length 1 of the update tensor. Figures 3 and 4 show tensor networks of BD and other tensor decompositions in graphical notation.

<sup>3</sup>More generally  $\oslash$  can be replaced with  $\boxtimes$  if the shapes are different.

**Computational complexity.** The computational cost of BD is  $O(IJK)$  per ALS step for a tensor of size  $(I, J, K)$ , both for the least-squares solution (Equation (39)) and for each factor matrix update. This is comparable to a rank-1 CP decomposition, since one ALS update step for a rank- $R$  CP decomposition costs  $O(IJKR) + O(R^3)$ , which reduces to  $O(IJK)$  at  $R = 1$ . Other properties such as sensitivity to initialization remain important directions for future work.

Furthermore, the expressive power of the model can be improved by considering the sum of BDs:

$$\mathcal{Y} \approx \sum_{r=1}^R \mathcal{A}^{(r)} \boxtimes \mathcal{B}^{(r)} \boxtimes \mathcal{C}^{(r)}, \quad (44)$$

where  $\mathcal{A}^{(r)} \in \mathbb{R}^{I \times J \times 1}$ ,  $\mathcal{B}^{(r)} \in \mathbb{R}^{I \times 1 \times K}$ , and  $\mathcal{C}^{(r)} \in \mathbb{R}^{1 \times J \times K}$  mutually satisfy the broadcast conditions for each  $r \in \{1, 2, \dots, R\}$ . The sum of BDs shown in Equation (44) is an extension of the outer product in CP decomposition (Hitchcock, 1927; Carroll & Chang, 1970; Harshman, 1970) to the broadcast product. For minimizing squared errors  $\|\mathcal{Y} - \sum_{r=1}^R \mathcal{A}^{(r)} \boxtimes \mathcal{B}^{(r)} \boxtimes \mathcal{C}^{(r)}\|_F^2$ , the hierarchical ALS (HALS) (Cichocki et al., 2007) can be adapted. Objective function of the sub-problem for  $k$ -th components is given by  $\|\mathcal{Y}_k - \mathcal{A}^{(k)} \boxtimes \mathcal{B}^{(k)} \boxtimes \mathcal{C}^{(k)}\|_F^2$ , where  $\mathcal{Y}_k := \mathcal{Y} - \sum_{r \neq k} \mathcal{A}^{(r)} \boxtimes \mathcal{B}^{(r)} \boxtimes \mathcal{C}^{(r)}$ , then the update rules can be derived in the same way of Equation (41), Equation (42), and Equation (43).

### 5.3 Difference from conventional TDs

In this section, we provide a preliminary illustration of how BD differs from conventional tensor decompositions (TDs). Note that the tensors we want to approximate may represent some real data or they may represent parameters of a neural network in real applications, but the concrete case study is beyond the scope of this paper. Our goal here is to show, at least in part, that BD has distinct structural properties from other decomposition models. Note that synthetic results below are generated from the BD model and are therefore expected to favor BD by construction; the real-data results should be interpreted as preliminary.

First, we constructed a synthetic tensor  $\mathcal{W} \in \mathbb{R}^{32 \times 32 \times 32}$  as follows:

$$\mathcal{W} = \mathcal{A} \boxtimes \mathcal{B} \boxtimes \mathcal{C} + \sigma \mathcal{E}, \quad (45)$$

where  $\mathcal{A} \in \mathbb{R}^{32 \times 32 \times 1}$ ,  $\mathcal{B} \in \mathbb{R}^{32 \times 1 \times 32}$ ,  $\mathcal{C} \in \mathbb{R}^{1 \times 32 \times 32}$  and noise  $\mathcal{E} \in \mathbb{R}^{32 \times 32 \times 32}$  are randomly generated, and  $\sigma > 0$ . We also used a real traffic speed data tensor<sup>4</sup> of size  $(32, 21, 24)$  measured every 24 hours at 32 locations over 21 days. The data tensor is treated as the ground truth, and a small amount of Gaussian noise is added to it, which is then treated as the observed data.

Next, we applied CP decomposition (CPD), Tucker decomposition (TKD) (Tucker, 1966; De Lathauwer et al., 2000), tensor-ring decomposition (TRD) (Zhao et al., 2016), and the proposed sum of BDs to a tensor  $\mathcal{W}$ . The optimization algorithm was applied with various values of the rank parameters, and the signal-to-noise ratio (SNR) of the reconstructed tensors was evaluated by

$$\text{SNR} = 10 \log \frac{\|\mathcal{W}_0\|_F^2}{\|\mathcal{W}_0 - \widehat{\mathcal{W}}\|_F^2}, \quad (46)$$

where  $\mathcal{W}_0$  is a ground-truth synthetic or real tensor without noise (i.e.,  $\mathcal{W}_0 = \mathcal{A} \boxtimes \mathcal{B} \boxtimes \mathcal{C}$  for synthetic) and  $\widehat{\mathcal{W}}$  is a reconstructed tensor obtained by TD algorithms. The number of model parameters and SNRs are shown in Figure 5. It can be seen that while the proposed BD succeeds in achieving accurate approximation, other low-rank TD models have difficulty achieving efficient approximation in synthetic data. In contrast, the results from real data

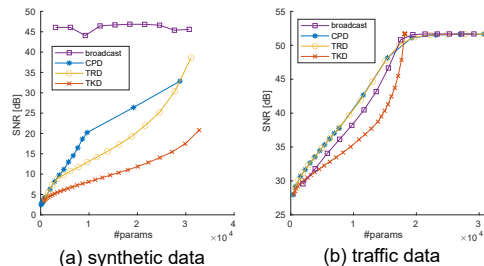


Figure 5: Dimensionality reduction of synthetic and traffic tensors: This suggests the existence of a tensor structure that favors broadcast decomposition.

while the proposed BD succeeds in achieving accurate approximation, other low-rank TD models have difficulty achieving efficient approximation in synthetic data. In contrast, the results from real data

<sup>4</sup>The data is publicly available at <http://www.openits.cn/openData2/792.jhtml>

suggest that the BD is intermediate between CPD/TRD and TKD. Although BD and TDs are similar in the sense that they are compact representations with few parameters, the properties of tensors are significantly different.

## 6 Related Work

### 6.1 Relationship with existing tensor products

The prototype of broadcast product has been proposed as the penetrating face (PF) product (Slyusar, 1999). The PF product is only defined for a matrix  $\mathbf{X} \in \mathbb{R}^{I \times J}$  and a tensor  $\mathcal{Y} \in \mathbb{R}^{I \times J \times K}$  as

$$\mathbf{X} \odot_{\text{PF}} \mathcal{Y} = \mathbf{X} \square \mathcal{Y}. \quad (47)$$

This can be regarded as a generalization of the Hadamard product between two matrices and a special case of the broadcast product at the same time. The broadcast product gives an integrated general definition of both the Hadamard product and the PF product.

From the perspective of graphical notation, broadcast product can be thought of as an operation that connects the corresponding modes of two tensors via third-order super-diagonal tensors (see Figure 2). When the third-order super-diagonal tensors are  $(1, 1, 1)$  in size, it is equivalent to a simple outer product. The broadcast product is notable for not including an inner-product-like operation (operation involving the  $\sum$  symbol) that directly connects the edges of tensors. Its mathematical structure is completely different from that of ordinary tensor products, Einstein products Brazell et al. (2013), and t-products (Zhang et al., 2014; Zhang & Aeron, 2016).

The broadcast product is closely related to the Khatri-Rao (KR) product. The KR product is defined between two matrices with the same number of columns  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_R] \in \mathbb{R}^{I \times R}$  and  $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_R] \in \mathbb{R}^{J \times R}$ :

$$\mathbf{B} \odot_{\text{KR}} \mathbf{A} = [\mathbf{b}_1 \otimes \mathbf{a}_1, \dots, \mathbf{b}_R \otimes \mathbf{a}_R] \in \mathbb{R}^{IJ \times R}, \quad (48)$$

where  $\otimes$  is the Kronecker product. The following relationship holds:

$$\mathbf{B} \odot_{\text{KR}} \mathbf{A} = \text{fold}_3(\tilde{\mathcal{A}} \square \tilde{\mathcal{B}})^\top, \quad (49)$$

where  $\tilde{\mathcal{A}} \in \mathbb{R}^{I \times 1 \times R}$  and  $\tilde{\mathcal{B}} \in \mathbb{R}^{1 \times J \times R}$  are reshaped tensors of  $\mathbf{A}$  and  $\mathbf{B}$ . In other words, the KR product and the broadcast product have the same mathematical structure because  $\text{fold}_3(\cdot)$  and  $\cdot^\top$  are just tensor reshaping operations. Noting that the KR product is defined only for matrices, the broadcast product can therefore be regarded as a generalization of the KR product in a broader sense.

### 6.2 Notation for describing complex mathematical expressions

Some papers that aim for accurate descriptions have already introduced the concept of the broadcast product, such as Wang et al. (2022) and You et al. (2024). Unlike us, they have not discussed the mathematical properties. Also, Wang et al. (2022) used  $\otimes$  as the symbol for the broadcast product, but  $\otimes$  is generally used for the Kronecker product. This confusion can be avoided by using our  $\square$ .

The Einops notation (Rogozhnikov, 2021) and the detailed Transformer description (Phuong & Hutter, 2022) are valuable references for writing clear mathematical descriptions. Looking at practical implementations, `TensorLy` is a modern library for tensor processing (Kossaifi et al., 2019).

### 6.3 Named tensor notation

Named Tensor Notation (Chiang et al., 2023) shares the most similar motivation with ours. By explicitly naming each **axis**, Named Tensor Notation allows complex tensor operations to be expressed concisely and code-like. For example, the masking operation discussed in Section 1 is described as follows:

$$Y = X \odot B, \quad \text{where } Y, X \in \mathbb{R}^{\text{height} \times \text{width} \times \text{chans}}, B \in \mathbb{R}^{\text{height} \times \text{width}} \quad (50)$$

Broadcasting is automatically performed when the **axes** do not match (Def. 6 in Chiang et al. (2023)), so broadcast products can be represented with the above simple expression. Named Tensor Notation has several valuable properties, such as omitting **axis** ordering and hiding **axis** lengths.

For comparison, the same masking operation in our notation is:

$$\mathcal{Y} = \mathcal{X} \boxtimes \mathcal{B}, \quad \mathcal{X} \in \mathbb{R}^{H \times W \times 3}, \quad \mathcal{B} \in \mathbb{R}^{H \times W}. \quad (51)$$

Shapes are specified numerically, and the  $\boxtimes$  operator handles the broadcasting. Because the resulting expression is standard tensor notation, it is directly substitutable into mathematical inequalities and identities (as illustrated in Sections 1 and 3).

**Trade-offs.** The two approaches have complementary strengths. Our broadcast product is fully compatible with standard mathematical notation: results can be substituted directly into equations, no axis-naming scheme is required, and it connects naturally to existing tensor algebra such as Frobenius-norm properties, least-squares solutions, and decompositions (see Sections 3 and 5). The main cost is that shapes must be specified explicitly for each operand. Named Tensor Notation, on the other hand, offers richer expressiveness (axis naming, reductions, and advanced indexing in a concise, code-like style) and leaves axis ordering implicit, reducing notational overhead for complex multi-axis operations. Its limitations are that a consistent axis-naming scheme must be maintained throughout the paper, and that results cannot be directly substituted into standard mathematical expressions.

**When to prefer each.** When a broadcast operation must be connected to standard tensor algebra (for instance, when deriving norm inequalities or least-squares solutions), our notation provides a cleaner formulation. It is also preferable when the same broadcast pattern recurs across multiple equations or when broadcasting is a structural feature of the model, as in broadcast decomposition. When the goal is a concise description of complex multi-axis operations, or when axis ordering should be left implicit, Named Tensor Notation may be preferable. When a single broadcast alignment must be specified in isolation, explicit shape annotation (e.g.,  $\mathcal{B}^{\square} \in \mathbb{R}^{H \times W \times 3}$ ) can also be more readable than either. Named Tensor Notation and our broadcast product represent different approaches with the same motivation; **developing** a unifying theory would be a direction for future work.

## 7 Conclusion

We redefined the broadcast operation commonly used in scientific computing libraries. The broadcast product enhances the accuracy of mathematical expressions in machine learning and other fields. Additionally, we contextualized this product within tensor decomposition, highlighting its potential for developing new types of tensor decompositions.

## References

- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- Dennis S Bernstein. *Matrix Mathematics: Theory, Facts, and Formulas*. Princeton University press, second edition, 2009.
- Jeff Bezanson, Stefan Karpinski, Viral B. Shah, and Other Contributors. Julia documentation, broadcasting. <https://docs.julialang.org/en/v1/manual/arrays/#Broadcasting>, 2026. Accessed: 2026-02-20.
- Michael Brazell, Na Li, Carmeliza Navasca, and Christino Tamon. Solving multilinear systems via tensor inversion. *SIAM Journal on Matrix Analysis and Applications*, 34(2):542–570, 2013.
- J Douglas Carroll and Jih-Jie Chang. Analysis of individual differences in multidimensional scaling via an N-way generalization of "Eckart-Young" decomposition. *Psychometrika*, 35(3):283–319, 1970.

- David Chiang, Alexander M Rush, and Boaz Barak. Named tensor notation. *Transactions on Machine Learning Research*, 2023.
- Andrzej Cichocki, Rafal Zdunek, and Shun-ichi Amari. Hierarchical ALS algorithms for nonnegative matrix and 3d tensor factorization. In *Proceedings of the International Conference on Independent Component Analysis and Signal Separation*, pp. 169–176. Springer, 2007.
- Andrzej Cichocki, Rafal Zdunek, Anh Huy Phan, and Shun-Ichi Amari. *Nonnegative Matrix and Tensor Factorizations: Applications to Exploratory Multi-Way Data Analysis and Blind Source Separation*. John Wiley & Sons, Ltd, 2009.
- Andrzej Cichocki, Namgil Lee, Ivan Oseledets, Anh-Huy Phan, Qibin Zhao, and Danilo P. Mandic. 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.
- Andrzej Cichocki, Anh-Huy Phan, Qibin Zhao, Namgil Lee, Ivan Oseledets, Masashi Sugiyama, and Danilo P. Mandic. Tensor networks for dimensionality reduction and large-scale optimization: Part 2 applications and future perspectives. *Foundations and Trends in Machine Learning*, 9(6):431–673, 2017.
- Lieven De Lathauwer, Bart De Moor, and Joos Vandewalle. On the best rank-1 and rank- $(r_1, r_2, \dots, r_n)$  approximation of higher-order tensors. *SIAM Journal on Matrix Analysis and Applications*, 21(4):1324–1342, 2000.
- Gene H Golub and Charles F Van Loan. *Matrix Computations*. JHU press, fourth edition, 2012.
- Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT press, 2016.
- Charles R Harris, K Jarrod Millman, Stéfan J Van Der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J Smith, et al. Array programming with numpy. *Nature*, 585(7825):357–362, 2020.
- Richard A Harshman. Foundations of the PARAFAC procedure: Models and conditions for an "explanatory" multimodal factor analysis. *UCLA Working Paper in Phonetics*, 16:1–84, 1970.
- Frank L. Hitchcock. The expression of a tensor or a polyadic as a sum of products. *Journal of Mathematics and Physics*, 6(1-4):164–189, 1927.
- Roger A Horn and Charles R Johnson. *Topics in matrix analysis*. Cambridge University Press, 1994.
- Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*, 2017.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015.
- Anil K. Jain. *Fundamentals of Digital Image Processing*. Prentice-Hall, 1989.
- Eric Kernfeld, Misha Kilmer, and Shuchin Aeron. Tensor-tensor products with invertible linear transforms. *Linear Algebra and its Applications*, 485:545–570, 2015.
- MJ Kochenderfer. *Algorithms for Optimization*. MIT Press, 2019.
- Tamara G. Kolda and Brett W. Bader. Tensor decompositions and applications. *SIAM Review*, 51(3):455–500, 2009.
- Jean Kossaifi, Yannis Panagakis, Anima Anandkumar, and Maja Pantic. Tensorly: Tensor learning in python. *Journal of Machine Learning Research*, 20(26):1–6, 2019.
- Chao Li, Junhua Zeng, Chunmei Li, Cesar F. Caiafa, and Qibin Zhao. Alternating local enumeration (TnALE): Solving tensor network structure search with fewer evaluations. In *International Conference on Machine Learning (ICML)*, 2023a.

- Zhengqi Li, Qianqian Wang, Forrester Cole, Richard Tucker, and Noah Snavely. Dynibar: Neural dynamic image-based rendering. In *CVPR*, 2023b.
- Shanchuan Lin, Andrey Ryabtsev, Soumyadip Sengupta, Brian L Curless, Steven M Seitz, and Ira Kemelmacher-Shlizerman. Real-time high-resolution background matting. In *CVPR*, 2021.
- Ivan Markovsky. *Low Rank Approximation: Algorithms, Implementation, Applications*. Springer, 2012.
- Yuriel Núñez Fernández, Marc K. Ritter, Matthieu Jeannin, Jheng-Wei Li, Thomas Kloss, Thibaud Louvet, Satoshi Terasaki, Olivier Parcollet, Jan von Delft, Hiroshi Shinaoka, and Xavier Waintal. Learning tensor networks with tensor cross interpolation: New algorithms and libraries. *SciPost Physics*, 18(3):104, 2025.
- Roger Penrose. Applications of negative dimensional tensors. *Combinatorial Mathematics and its Applications*, 1:221–244, 1971.
- Ethan Perez, Florian Strub, Harm de Vries, Vincent Dumoulin, and Aaron C. Courville. Film: Visual reasoning with a general conditioning layer. In *AAAI*, 2018.
- Mary Phuong and Marcus Hutter. Formal algorithms for transformers. *arXiv:2207.09238*, 2022.
- Shi-Ju Ran, Emanuele Tirrito, Cheng Peng, Xi Chen, Luca Tagliacozzo, Gang Su, and Maciej Lewenstein. *Tensor Network Contractions: Methods and Applications to Quantum Many-Body Systems*. Springer, 2020.
- Alex Rogozhnikov. Einops: Clear and reliable tensor manipulations with einstein-like notation. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.
- VI Slyusar. A family of face products of matrices and its properties. *Cybernetics and Systems Analysis*, 35(3):379–384, 1999.
- Edwin Stoudenmire and David J. Schwab. Supervised learning with tensor networks. In *Advances in Neural Information Processing Systems*, 2016.
- Gilbert Strang. *Linear Algebra and Learning from Data*. SIAM, 2019.
- Jordan K Taylor. An introduction to graphical tensor notation for mechanistic interpretability. *arXiv preprint arXiv:2402.01790*, 2024.
- Sergios Theodoridis. *Machine learning: a Bayesian and Optimization Perspective*. Academic press, second edition, 2020.
- Ledyard R Tucker. Some mathematical notes on three-mode factor analysis. *Psychometrika*, 31(3):279–311, 1966.
- Dmitry Ulyanov, Andrea Vedaldi, and Victor S. Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2016.
- Fanyi Wang, Haotian Hu, Cheng Shen, Tianpeng Feng, and Yandong Guo. BAM: a balanced attention mechanism to optimize single image super-resolution. *Journal of Real-Time Image Processing*, 19(5):941–955, 2022.
- Yuxin Wu and Kaiming He. Group normalization. In *ECCV*, 2018.
- Tatsuya Yokota. Very basics of tensors with graphical notations: Unfolding, calculations, and decompositions. *arXiv preprint arXiv:2411.16094*, 2024.
- Kaichao You, Guo Qin, Anchang Bao, Meng Cao, Ping Huang, Jiulong Shan, and Mingsheng Long. Efficient convbn blocks for transfer learning and beyond. In *ICLR*, 2024.
- Zemin Zhang and Shuchin Aeron. Exact tensor completion using t-svd. *IEEE Transactions on Signal Processing*, 65(6):1511–1526, 2016.

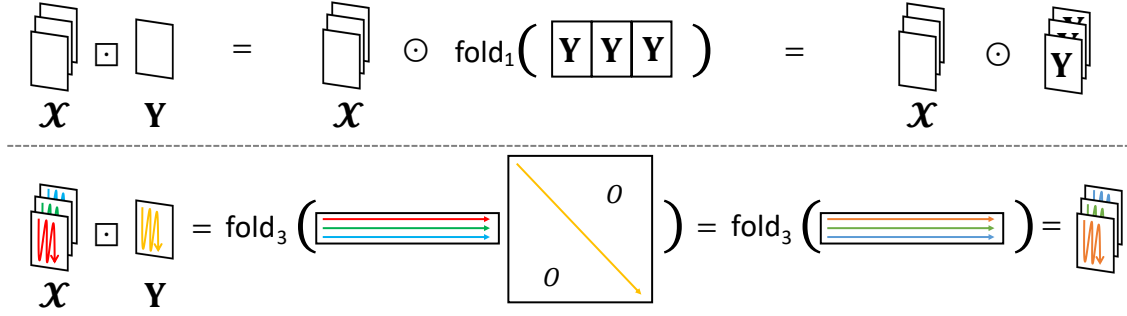


Figure 6: The broadcast product of a third-order tensor and a matrix Equation (28):  $\mathcal{X} \square \mathbf{Y} = \mathcal{X} \odot \text{fold}_1(\mathbf{1}_K^\top \otimes \mathbf{Y}) = \text{fold}_3(\mathbf{X}_{(3)} \text{diag}(\text{vec}(\mathbf{Y})))$

Zemin Zhang, Gregory Ely, Shuchin Aeron, Ning Hao, and Misha Kilmer. Novel methods for multilinear data completion and de-noising based on tensor-SVD. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3842–3849, 2014.

Qibin Zhao, Guoxu Zhou, Shengli Xie, Liqing Zhang, and Andrzej Cichocki. Tensor ring decomposition. *arXiv preprint arXiv:1606.05535*, 2016.

## A Proof of Theorem 3.1

We prove Theorem 3.1. For two tensors  $\mathcal{X}, \mathcal{Y}$  that satisfy the broadcast condition,

$$\|\mathcal{X} \square \mathcal{Y}\|_F = \|\mathcal{X} \square \odot \mathcal{Y}\|_F \quad (\text{Equation (20)}) \quad (52)$$

$$\leq \|\mathcal{X}\|_F \cdot \|\mathcal{Y}\|_F \quad (\text{Cauchy-Schwarz for the Hadamard product}) \quad (53)$$

$$= \|\mathcal{X}\|_F \cdot \|\mathcal{Y}\|_F. \quad (\text{Definition of marginalization}) \quad (54)$$

Thus, the claim follows.

Here, Cauchy-Schwarz for the Hadamard product (Horn & Johnson, 1994) means that for any tensors  $\mathcal{A}, \mathcal{B}$  of the same shape,  $\|\mathcal{A} \odot \mathcal{B}\|_F \leq \|\mathcal{A}\|_F \cdot \|\mathcal{B}\|_F$ . This is obtained by considering  $\mathbf{A} = \text{diag}(\mathcal{A})$  and  $\mathbf{B} = \text{diag}(\mathcal{B})$ , and substituting them into  $\|\mathbf{A}\mathbf{B}\|_F \leq \|\mathbf{A}\|_F \cdot \|\mathbf{B}\|_F$  (the sub-multiplicativity of the Frobenius norm).

## B Visualizing Third-order Tensors and Matrices

Figure 6 visualizes Equation (28), the broadcast product of a third-order tensor and a matrix. Here, the Kronecker product is denoted by  $\mathbf{X} \otimes \mathbf{Y}$ . The **mode- $n$  unfolding**  $\mathbf{X}_{(n)} \in \mathbb{R}^{I_n \times \prod_{k \neq n} I_k}$  rearranges all entries of  $\mathcal{X}$  into a matrix whose rows are indexed by the  $n$ -th axis (Kolda & Bader, 2009; Yokota, 2024); for example,  $\mathbf{X}_{(1)} \in \mathbb{R}^{I \times JK}$ . Its inverse is denoted  $\text{fold}_n(\cdot)$ , so that  $\text{fold}_n(\mathbf{X}_{(n)}) = \mathcal{X}$ .

Here,  $\mathbf{1}_K^\top \otimes \mathbf{Y} \in \mathbb{R}^{I \times JK}$  represents  $\mathbf{Y}$  repeated  $K$  times horizontally. By  $\text{fold}_1$ , the shape is aligned with that of  $\mathcal{X}$ . Also, one can express this computation by multiplying the matrix  $\mathbf{X}_{(3)} \in \mathbb{R}^{K \times IJ}$  by the diagonal matrix  $\text{diag}(\text{vec}(\mathbf{Y})) \in \mathbb{R}^{IJ \times IJ}$ . Here, we first unfold each frontal slice of  $\mathcal{X}$  and arrange them to form a matrix. Then, for each row, we take the product with the unfolded  $\mathbf{Y}$ . Finally, the result is folded back into its original shape using  $\text{fold}_3$ .

In addition, the broadcast product of a third-order tensor and a matrix shown in Equation (28) can be analyzed using graphical notation in Figure 7. Here, unfolding and folding operators are denoted by half circles with three lines. The analysis is based on the topological equivalence of graphs and properties of operators of unfolding, folding with super-diagonal tensors. For examples, operations of unfolding-and-folding and folding-and-unfolding are identity mappings (see Figure 7a). The outer product of two super-diagonal tensors and their unfolding results in another super-diagonal tensor (see Figure 7b).

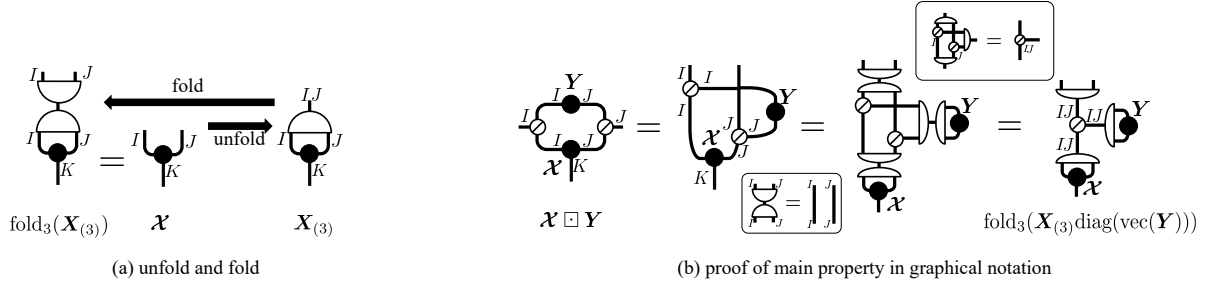


Figure 7: Proof of property Equation (28) using graphical notation. (a) Unfolding-and-folding is an identity mapping, and it can be applied to proof in (b). In addition, the outer product of two super-diagonal tensors and their unfolding results in another super-diagonal tensor. Then, we can obtain  $\mathcal{X} \square \mathcal{Y} = \text{fold}_3(\mathcal{X}_{(3)} \text{diag}(\text{vec}(\mathcal{Y})))$  in graphical notation.

## C Properties of broadcast sum, difference, and division

The broadcast difference can be expressed using the broadcast sum through the following trivial transformation. For tensors  $\mathcal{X}$  and  $\mathcal{Y}$  satisfying the broadcast condition,

$$\mathcal{X} \boxminus \mathcal{Y} = \mathcal{X} \boxplus (-\mathcal{Y}). \quad (55)$$

Similarly, the broadcast division can be expressed using the broadcast product through the following transformation. For tensors  $\mathcal{X}$  and  $\mathcal{Y}$  satisfying the broadcast condition (with  $\mathcal{Y}$  containing no zero elements),

$$\mathcal{X} \boxdiv \mathcal{Y} = \mathcal{X} \boxtimes (\mathbf{1} \odot \mathcal{Y}), \quad (56)$$

where  $\mathbf{1}$  is a tensor of the same shape as  $\mathcal{Y}$ , with all elements equal to 1. Therefore, focusing on the properties of the broadcast sum and product is sufficient.

Here, we introduce the basic properties of broadcast sum. First, Equations (12), (13) and (16) also hold if  $\boxtimes$  and  $\odot$  are replaced with  $\boxplus$  and  $+$ , respectively.

**Scalar:** For any tensor  $\mathcal{X}$  and a scalar  $y$ , we obtain

$$\mathcal{X} \boxplus y = \mathcal{X} + y\mathbf{1}, \quad (57)$$

where  $\mathbf{1}$  is an all-one tensor having the same shape as  $\mathcal{X}$ . This equation may seem extremely trivial initially, but it is important. Adding a constant to all tensor elements is the most basic form of broadcasting, and in `numpy` it is written as a simple addition like `X + y`. Unfortunately, many papers directly write this in mathematical notation as  $X + y$ . Such a description is incorrect and must be written as  $\mathcal{X} + y\mathbf{1}$  as shown above. In addition, there are frequent mistakes where this operation is misunderstood, and the identity matrix  $I$  is incorrectly used, resulting in expressions like  $X + yI$ .

**Vector and vector:** For vectors  $\mathbf{x} \in \mathbb{R}^I$  and  $\mathbf{y} \in \mathbb{R}^J$  of different lengths, we obtain

$$\mathbf{x} \boxplus \mathbf{y}^\top = \mathbf{x}\mathbf{1}_J^\top + \mathbf{1}_I\mathbf{y}^\top \in \mathbb{R}^{I \times J}. \quad (58)$$

**Matrix and vector:** For a matrix  $\mathbf{X} \in \mathbb{R}^{I \times J}$ , vectors  $\mathbf{y} \in \mathbb{R}^I$  and  $\mathbf{z} \in \mathbb{R}^J$ , the following holds:

$$\mathbf{X} \boxplus \mathbf{y} = \mathbf{X} + [\mathbf{y}|\mathbf{y}|\dots|\mathbf{y}] = \mathbf{X} + \mathbf{y}\mathbf{1}_J^\top. \quad (59)$$

$$\mathbf{X} \boxplus \mathbf{z}^\top = \mathbf{X} + \begin{bmatrix} \mathbf{z}^\top \\ \vdots \\ \mathbf{z}^\top \end{bmatrix} = \mathbf{X} + \mathbf{1}_I\mathbf{z}^\top. \quad (60)$$

## D Proof of Least Squares Solution

### D.1 Case of the third-order tensors

Here, we show the derivation of the LS solution Equation (39):

$$\widehat{\mathcal{A}} = \underset{\mathcal{A}}{\operatorname{argmin}} \|\mathcal{Y} - \mathcal{A} \boxtimes \mathcal{Z}\|_F^2 \quad (61)$$

for  $\mathcal{Y} \in \mathbb{R}^{I \times J \times K}$ ,  $\mathcal{A} \in \mathbb{R}^{I \times J \times 1}$ , and  $\mathcal{Z} \in \mathbb{R}^{1 \times J \times K}$ . Let us put  $\mathbf{Y}_j := \mathcal{Y}_{:j} \in \mathbb{R}^{I \times K}$ ,  $\mathbf{a}_j := \mathcal{A}_{:j1} \in \mathbb{R}^I$  and  $\mathbf{z}_j := \mathcal{Z}_{1j} \in \mathbb{R}^K$ , the squared errors can be transformed as

$$\|\mathcal{Y} - \mathcal{A} \boxtimes \mathcal{Z}\|_F^2 = \sum_{j=1}^J \|\mathbf{Y}_j - \mathbf{a}_j \mathbf{z}_j^\top\|_F^2. \quad (62)$$

Then the solution of  $\mathbf{a}_j$  can be independently obtained by

$$\widehat{\mathbf{a}}_j = \underset{\mathbf{a}_j}{\operatorname{argmin}} \|\mathbf{Y}_j - \mathbf{a}_j \mathbf{z}_j^\top\|_F^2 = \mathbf{Y}_j \mathbf{z}_j (\mathbf{z}_j^\top \mathbf{z}_j)^{-1} \quad (63)$$

for each  $j \in \{1, 2, \dots, J\}$ . Since  $(i, j, 1)$ -th entry of  $\widehat{\mathcal{A}}$  corresponds to  $i$ -th entry of  $\widehat{\mathbf{a}}_j$ , we have

$$\widehat{a}_{ij1} = \widehat{\mathbf{a}}_j(i) = \left( \sum_{k=1}^K y_{ijk} z_{jk} \right) \left( \sum_{k=1}^K z_{jk}^2 \right)^{-1} \iff \widehat{\mathcal{A}} = \mathcal{P}_3(\mathcal{Y} \boxtimes \mathcal{Z}) \boxtimes \mathcal{P}_3(\mathcal{Z} \boxtimes \mathcal{Z}). \quad (64)$$

### D.2 Case of the $N$ -th order tensors

Let us consider  $N$ -th order tensors  $\mathcal{W} \in \mathbb{R}^{D_1 \times D_2 \times \dots \times D_N}$ ,  $\mathcal{H} \in \mathbb{R}^{F_1 \times F_2 \times \dots \times F_N}$  and  $\mathcal{X} \in \mathbb{R}^{\max(D_1, F_1) \times \max(D_2, F_2) \times \dots \times \max(D_N, F_N)}$ , then the problem can be written by

$$\widehat{\mathcal{W}} = \underset{\mathcal{W}}{\operatorname{argmin}} \|\mathcal{X} - \mathcal{W} \boxtimes \mathcal{H}\|_F^2. \quad (65)$$

Without loss of generality, we can reduce the problem with  $N$ -th order tensors  $(\mathcal{X}, \mathcal{W}, \mathcal{H})$  in Equation (65) to the problem with third-order tensors  $(\mathcal{Y}, \mathcal{A}, \mathcal{Z})$  in Equation (61). Since  $\mathcal{W}$  and  $\mathcal{H}$  satisfy the broadcast condition, the  $N$  modes can be divided into three categories:

$$\mathcal{L} = \{ n \mid D_n > 1, F_n = 1 \}, \quad (66)$$

$$\mathcal{S} = \{ n \mid D_n = F_n \}, \quad (67)$$

$$\mathcal{R} = \{ n \mid D_n = 1, F_n > 1 \}. \quad (68)$$

$\mathcal{L}$  is the set of broadcasting modes for  $\mathcal{H}$ , corresponding to the first mode of  $\mathcal{Z}$  in Equation (61).  $\mathcal{S}$  is the set of non-broadcasting modes, corresponding to the second mode in Equation (61).  $\mathcal{R}$  is the set of broadcasting modes for  $\mathcal{W}$ , corresponding to the third mode of  $\mathcal{A}$  in Equation (61). Then, we convert  $N$ -th order tensors to third-order tensors based on  $(\mathcal{L}, \mathcal{S}, \mathcal{R})$  using mode permutation and tensor unfolding as follow:

$$\mathcal{Y} = \operatorname{unfold}_{(I, J, K)} \operatorname{permute}_{(\mathcal{L}, \mathcal{S}, \mathcal{R})}(\mathcal{X}) \in \mathbb{R}^{I \times J \times K}, \quad (69)$$

$$\mathcal{A} = \operatorname{unfold}_{(I, J, 1)} \operatorname{permute}_{(\mathcal{L}, \mathcal{S}, \mathcal{R})}(\mathcal{W}) \in \mathbb{R}^{I \times J \times 1}, \quad (70)$$

$$\mathcal{Z} = \operatorname{unfold}_{(1, J, K)} \operatorname{permute}_{(\mathcal{L}, \mathcal{S}, \mathcal{R})}(\mathcal{H}) \in \mathbb{R}^{1 \times J \times K}, \quad (71)$$

where  $I = \prod_{n \in \mathcal{L}} D_n$ ,  $J = \prod_{n \in \mathcal{S}} D_n$ , and  $K = \prod_{n \in \mathcal{R}} F_n$ . The LS solution of third-order tensors  $\widehat{\mathcal{A}}$  can be obtained by Equation (64). By converting  $\widehat{\mathcal{A}}$  back to an  $N$ -th order tensor, the solution can be obtained as follows:

$$\widehat{\mathcal{W}} = \operatorname{permute}_{(\mathcal{L}, \mathcal{S}, \mathcal{R})}^{-1} \operatorname{unfold}_{(I, J, 1)}^{-1}(\widehat{\mathcal{A}}) = \mathcal{P}_{\mathcal{R}}(\mathcal{X} \boxtimes \mathcal{H}) \boxtimes \mathcal{P}_{\mathcal{R}}(\mathcal{H} \boxtimes \mathcal{H}), \quad (72)$$

where  $\text{permute}_{(\mathcal{L}, \mathcal{S}, \mathcal{R})}^{-1}$  and  $\text{unfold}_{(I, J, 1)}^{-1}$  are inverse of  $\text{permute}_{(\mathcal{L}, \mathcal{S}, \mathcal{R})}$  and  $\text{unfold}_{(I, J, 1)}$ , and  $\mathcal{P}_{\mathcal{R}}(\cdot)$  is a sum operation along the modes in  $\mathcal{R}$ . For example, let be

$$\begin{aligned}\mathcal{X} &\in \mathbb{R}^{10 \times 20 \times 30 \times 40 \times 50 \times 60}, \\ \mathcal{W} &\in \mathbb{R}^{10 \times 20 \times 1 \times 40 \times 50 \times 1}, \\ \mathcal{H} &\in \mathbb{R}^{10 \times 1 \times 30 \times 1 \times 50 \times 60}, \\ \mathcal{L} &= \{2, 4\}, \mathcal{S} = \{1, 5\}, \mathcal{R} = \{3, 6\},\end{aligned}$$

then the permutation operation outputs

$$\begin{aligned}\tilde{\mathcal{X}} &= \text{permute}_{(\{2,4\}, \{1,5\}, \{3,6\})}(\mathcal{X}) \in \mathbb{R}^{20 \times 40 \times 10 \times 50 \times 30 \times 60}, \\ \tilde{\mathcal{W}} &= \text{permute}_{(\{2,4\}, \{1,5\}, \{3,6\})}(\mathcal{W}) \in \mathbb{R}^{20 \times 40 \times 10 \times 50 \times 1 \times 1}, \\ \tilde{\mathcal{H}} &= \text{permute}_{(\{2,4\}, \{1,5\}, \{3,6\})}(\mathcal{H}) \in \mathbb{R}^{1 \times 1 \times 10 \times 50 \times 30 \times 60},\end{aligned}$$

the unfolding operation outputs

$$\begin{aligned}\mathcal{Y} &= \text{unfold}_{(800, 500, 1800)}(\tilde{\mathcal{X}}) \in \mathbb{R}^{800 \times 500 \times 1800}, \\ \mathcal{A} &= \text{unfold}_{(800, 500, 1)}(\tilde{\mathcal{W}}) \in \mathbb{R}^{800 \times 500 \times 1}, \\ \mathcal{Z} &= \text{unfold}_{(1, 500, 1800)}(\tilde{\mathcal{H}}) \in \mathbb{R}^{1 \times 500 \times 1800},\end{aligned}$$

and the sum operation outputs

$$\begin{aligned}\mathcal{P}_{\{3,6\}}(\mathcal{X} \boxtimes \mathcal{H}) &= (\mathcal{X} \boxtimes \mathcal{H}) \times_3 \mathbf{1}^\top \times_6 \mathbf{1}^\top \in \mathbb{R}^{10 \times 20 \times 1 \times 40 \times 50 \times 1}, \\ \mathcal{P}_{\{3,6\}}(\mathcal{H} \boxtimes \mathcal{H}) &= (\mathcal{H} \boxtimes \mathcal{H}) \times_3 \mathbf{1}^\top \times_6 \mathbf{1}^\top \in \mathbb{R}^{10 \times 1 \times 1 \times 1 \times 50 \times 1}.\end{aligned}$$

## E Conflicts of Mathematical Symbols

Here we discuss the issue of symbols for element-wise multiplication. As shown in Table 1, the symbols used to represent element-wise multiplication (Hadamard product) are very diverse and most of them conflict with other mathematical operations.

Since the symbol  $\boxtimes$  for the broadcast product proposed in this paper does not conflict and can also be used for the Hadamard product, it may solve such problems. In addition, we can represent both the Hadamard product and the broadcast product using only one symbol  $\boxtimes$ . A promising set of conflict-free notations would be given as follows.

- outer product  $\circ$
- Kronecker product  $\otimes$
- Khatri-Rao product  $\odot$
- Hadamard product (with broadcast option)  $\boxtimes$
- element-wise division (with broadcast option)  $\boxdiv$
- convolution  $\circledast$
- t-product  $*$

Note that this paper does not follow the above notation set, but it was necessary because of the definition of the broadcast product using the Hadamard product and discussing their relationship.

## F Translation from `numpy`/Julia

The proposed broadcast product is almost identical to `numpy`'s broadcast operation, allowing  $\mathbf{A} * \mathbf{B}$  in `numpy` to be represented as  $\mathcal{A} \boxtimes \mathcal{B}$  in equations. However, there is a difference in how cases are handled when the tensor orders differ, but the broadcasting condition is satisfied.

The key difference lies in the shape-normalization convention (Definition 1): our definition follows the **F-convention** (trailing singletons are appended), whereas `numpy` follows the **C-convention** (leading singletons

Table 1: List of symbols used for element-wise multiplication

symbol	usage / conflict
◦	Used in (Slyusar, 1999; Bernstein, 2009; Theodoridis, 2020). Conflicts with the outer product in (Cichocki et al., 2009; Kolda & Bader, 2009; Strang, 2019).
⊙	Used in (Markovsky, 2012; Goodfellow et al., 2016; Kochenderfer, 2019). Conflicts with the Khatri–Rao product in (Cichocki et al., 2009; Kolda & Bader, 2009; Strang, 2019).
⊗	Used in (Cichocki et al., 2009). Conflicts with convolution in (Jain, 1989; Strang, 2019).
*	Used in (Kolda & Bader, 2009). Conflicts with convolution (Strang, 2019) and the t-product (Kernfeld et al., 2015).
.*	Used in (Golub & Van Loan, 2012; Cichocki et al., 2009; Strang, 2019). No conflicts reported.
⊠ (ours)	No conflicts reported (including Hadamard product).

are prepended)<sup>5</sup>. For instance,  $\mathbb{R}^{2 \times 3}$  is **identified with  $\mathbb{R}^{2 \times 3 \times 1}$  under the F-convention, but with  $\mathbb{R}^{1 \times 2 \times 3}$  under the C-convention**. Below, we show some examples of Python codes.

As in the following example, if we explicitly specify the tensors’ shapes including modes of length one, all computations work as expected. We can directly translate  $A * B$  in `numpy` to  $\mathcal{A} \boxtimes \mathcal{B}$  in equations.

```
A = np.random.rand(2, 3, 4) # A.shape == (2, 3, 4), i.e.,  $\mathcal{A} \in \mathbb{R}^{2 \times 3 \times 4}$ 
B = np.random.rand(2, 3, 1) # B.shape == (2, 3, 1), i.e.,  $\mathcal{B} \in \mathbb{R}^{2 \times 3 \times 1}$ 

# ( $\mathbb{R}^{2 \times 3 \times 4}, \mathbb{R}^{2 \times 3 \times 1}$ ) satisfies the broadcast condition. Thus we can compute  $\mathcal{A} \boxtimes \mathcal{B}$ 
A * B # OK
```

However, if we don’t explicitly write down the **axis of length one**, `numpy`’s behavior is in the exact opposite way to our definition as follows.

```
C = np.random.rand(2, 3) # C.shape == (2, 3), i.e.,  $\mathcal{C} \in \mathbb{R}^{2 \times 3}$ 
D = np.random.rand(3, 4) # D.shape == (3, 4), i.e.,  $\mathcal{D} \in \mathbb{R}^{3 \times 4}$ 

# F-convention identifies  $\mathbb{R}^{2 \times 3}$  with  $\mathbb{R}^{2 \times 3 \times 1}$ , so  $\mathcal{A} \boxtimes \mathcal{C}$  can be computed.
# However, numpy's C-convention identifies  $\mathbb{R}^{2 \times 3}$  with  $\mathbb{R}^{1 \times 2 \times 3}$ , so  $A * C$  fails.
A * C # ValueError: operands could not be broadcast together with shapes (2,3,4) (2,3)

# F-convention identifies  $\mathbb{R}^{3 \times 4}$  with  $\mathbb{R}^{3 \times 4 \times 1}$ , so  $\mathcal{A} \boxtimes \mathcal{D}$  cannot be computed.
# However, numpy's C-convention identifies  $\mathbb{R}^{3 \times 4}$  with  $\mathbb{R}^{1 \times 3 \times 4}$ , so  $A * D$  succeeds.
A * D # OK
```

If confusion happens when writing the broadcast product, we recommend explicitly defining the shape even for the **axis with length one**. For example, an image with a single channel can be written as  $\mathbb{R}^{H \times W \times 1}$ .

In fact, since **Julia follows the F-convention (Definition 1)**, its behavior is exactly the same as our broadcast product. In **Julia**, the broadcast product is expressed by `.*` as follows.

```
A = rand(2, 3, 4) # size(A) == (2, 3, 4), i.e.,  $\mathcal{A} \in \mathbb{R}^{2 \times 3 \times 4}$ 
B = rand(2, 3, 1) # size(B) == (2, 3, 1), i.e.,  $\mathcal{B} \in \mathbb{R}^{2 \times 3 \times 1}$ 
# ( $\mathbb{R}^{2 \times 3 \times 4}, \mathbb{R}^{2 \times 3 \times 1}$ ) satisfies the broadcast condition. Thus we can compute  $\mathcal{A} \boxtimes \mathcal{B}$ 
A .* B # OK
```

<sup>5</sup><https://numpy.org/doc/stable/user/basics.broadcasting.html>

```
C = rand(2, 3) # size(C) == (2, 3), i.e.,  $\mathcal{C} \in \mathbb{R}^{2 \times 3}$ 
D = rand(3, 4) # size(D) == (3, 4), i.e.,  $\mathcal{D} \in \mathbb{R}^{3 \times 4}$ 
# Under the F-convention,  $\mathbb{R}^{2 \times 3}$  is identified with  $\mathbb{R}^{2 \times 3 \times 1}$ , so  $\mathcal{A} \boxtimes \mathcal{C}$  can be computed.
A .* C # OK
# Under the F-convention,  $\mathbb{R}^{3 \times 4}$  is identified with  $\mathbb{R}^{3 \times 4 \times 1}$ , so  $\mathcal{A} \boxtimes \mathcal{D}$  cannot be computed.
A .* D # ERROR: DimensionMismatch("arrays could not be broadcast to a common size; ..
```