



лекция

# Tensor and matrix methods

спикер



Konstantin  
Sozykin



Skoltech

# Content

Introduction →

# Content

Introduction → Matrix Decompositions →

# Content

Introduction → Matrix Decompositions →  
**Tensor Decompositions →**

# Content

Introduction → Matrix Decompositions →  
Tensor Decompositions → **Advanced  
Ideas** →

# Content

Introduction → Matrix Decompositions →  
Tensor Decompositions → Advanced  
Ideas → **Low-rank transformers** →

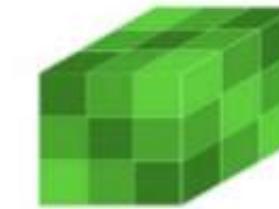
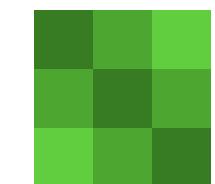
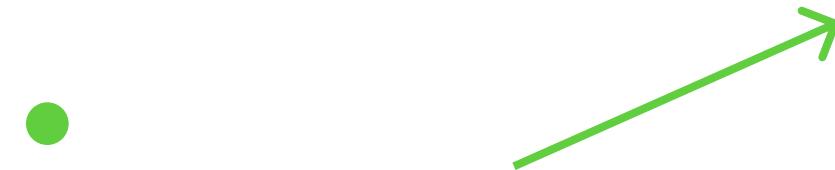
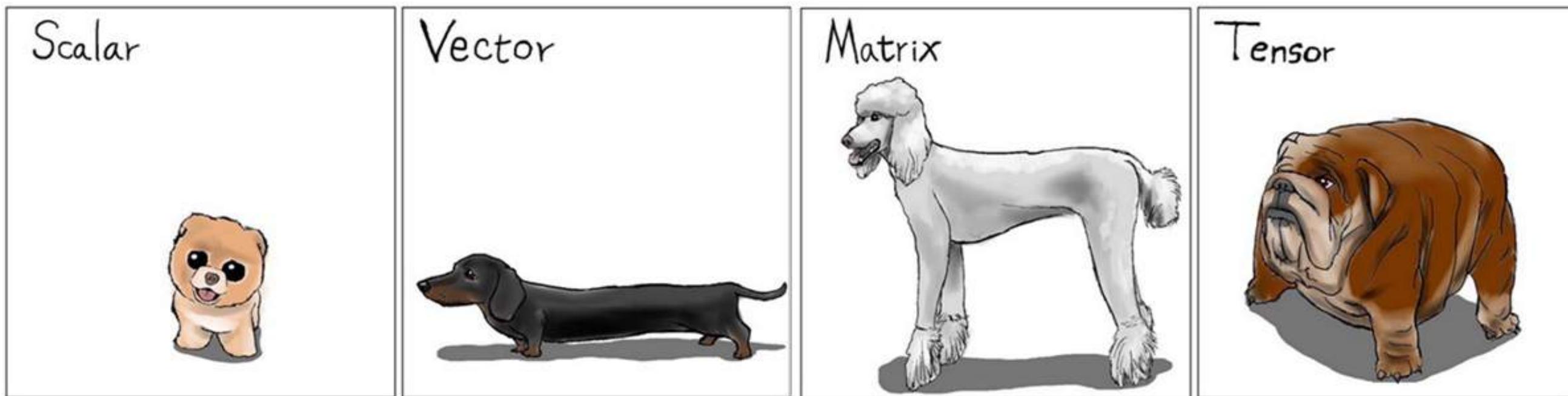
# Content

Introduction → Matrix Decompositions →  
Tensor Decompositions → Advanced  
Ideas → Low-rank transformers →  
**Links and References →**

# Introduction

# Data and tensors

- A multidimensional array
- A multilinear map

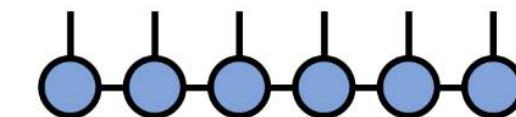


# Why tensor networks

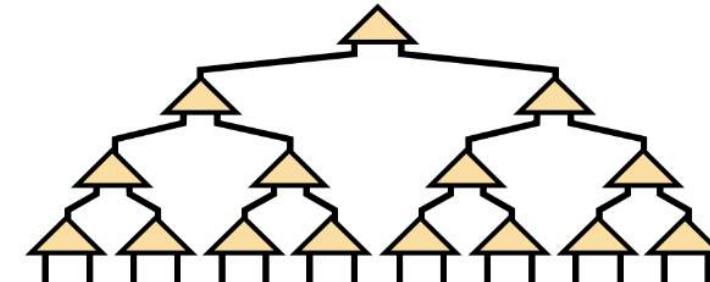
Decompositions, factorizations

- To handle big data and noisy data
- To describe and represent some complex physical and biological systems
- Initially it is more quantum physics and applied mathematics tool

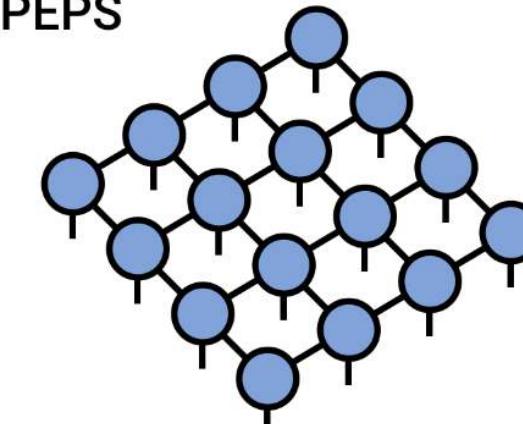
Matrix Product State /  
Tensor Train



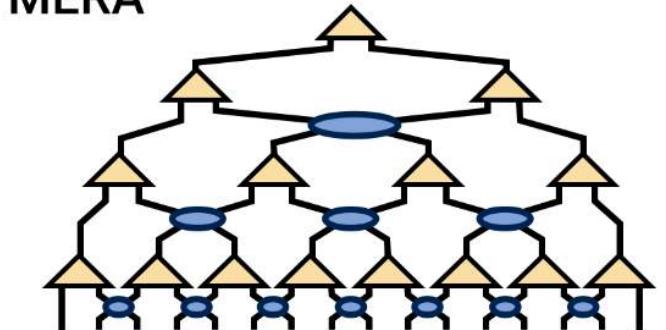
Tree Tensor Network /  
Hierarchical Tucker



PEPS



MERA



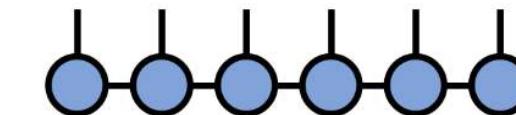
# Why tensor networks

Decompositions, factorizations

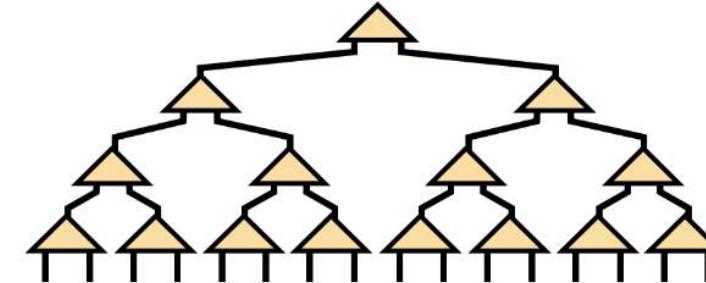
- To compress and restore (complete) data
- To reduce computational complexity (almost today topic)
- To able to use an algebraic approach

[Tensor Network ↗](#)

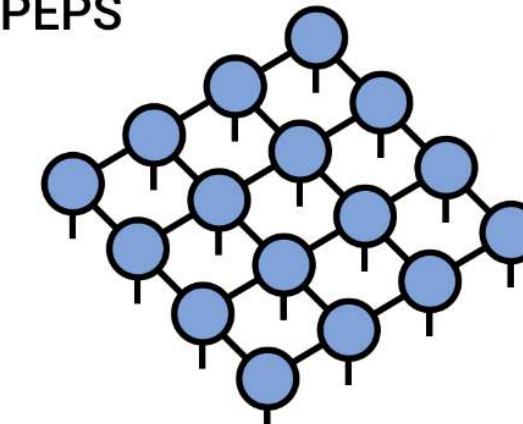
Matrix Product State /  
Tensor Train



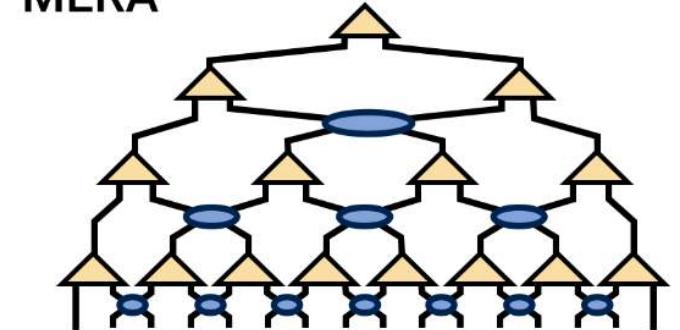
Tree Tensor Network /  
Hierarchical Tucker



PEPS



MERA



# Matrix Decompositions

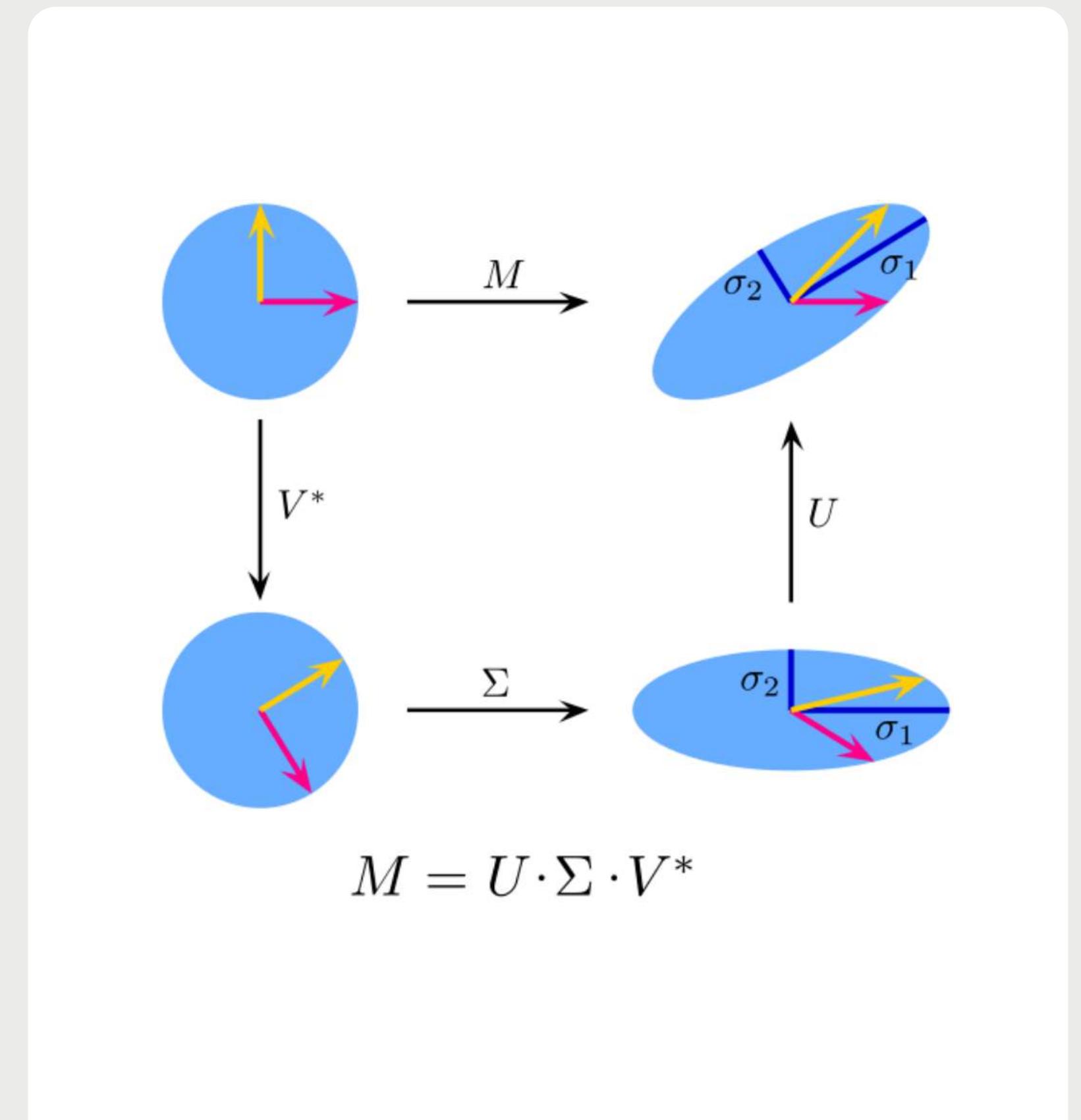
# Why tensor networks

Decompositions, factorizations

Illustration of the singular value decomposition  $\mathbf{U}\Sigma\mathbf{V}^*$  of a real  $2 \times 2$  matrix  $\mathbf{M}$ .

- **Top:** The action of  $\mathbf{M}$ , indicated by its effect on the unit disc  $D$  and the two canonical unit vectors  $e_1$  and  $e_2$
- **Left:** The action of  $\mathbf{V}^*$ , a rotation, on  $D$ ,  $e_1$ , and  $e_2$ .
- **Bottom:** The action of  $\Sigma$ , a scaling by the singular values  $\sigma_1$  horizontally and  $\sigma_2$  vertically.
- **Right:** The action of  $\mathbf{U}$ , another rotation.

Singular value decomposition ↗

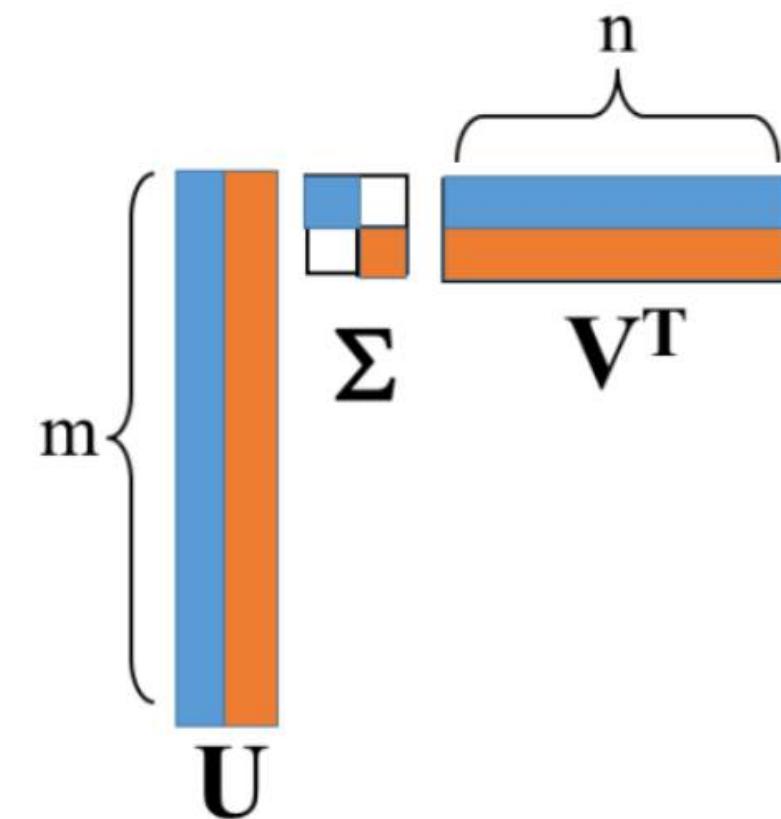


# SVD and Linear Layer

- Neural Networks overparameterized
- Recover that Linear Layer is  $W @ x + b$
- then it can be expressed such as ~
- $U @ [\Sigma @ V^T] @ x + b$
- instead  $m * n$  parameters we have only  $r * (1+m+n)$  or  $r * (m+n)$ , where  $r$  singular values here
- Then we **just replace and finetune the layer or whole network**
- Be aware about transpose of matrix, when doing practical implementation

$$A = U \Sigma V^T$$

Singular Value Decomposition (SVD)



$$A_{[m \times n]} = U_{[m \times r]} \Sigma_{[r \times r]} (V_{[n \times r]})^T$$

Singular Value Decomposition (SVD) ↗

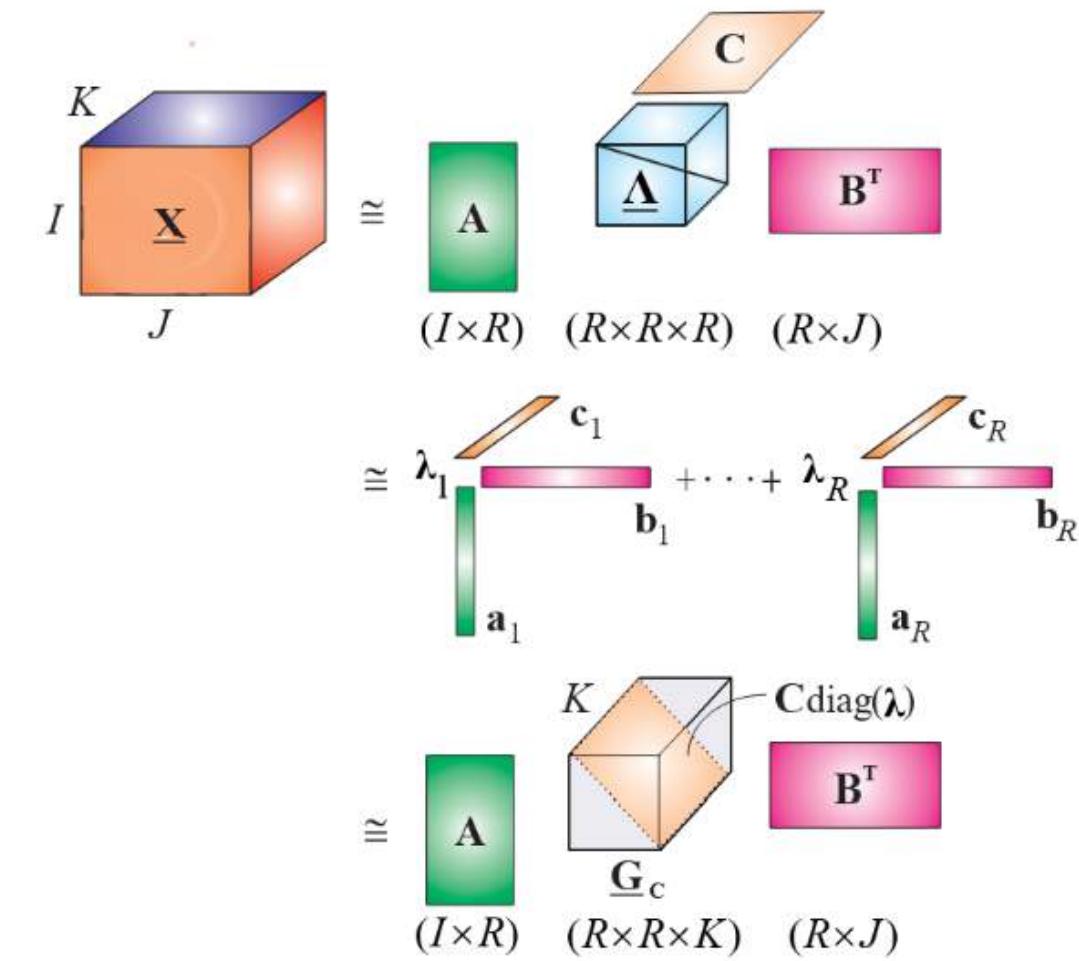
# Tensor Decompositions

# Canonical Polyadic Decomposition

- CPD approximates given tensor  $\mathbf{X}$  as a sum of  $R$  rank-one tensors (vectors)
- Optimized with alternating least squares technique
- Rank selection is NP hard!

Now publishers — Tensor Networks for Dimensionality Reduction and Large-scale Optimization: Part 1 Low-Rank Tensor Decompositions ↗

(a) Standard block diagram for CP decomposition of a 3rd-order tensor



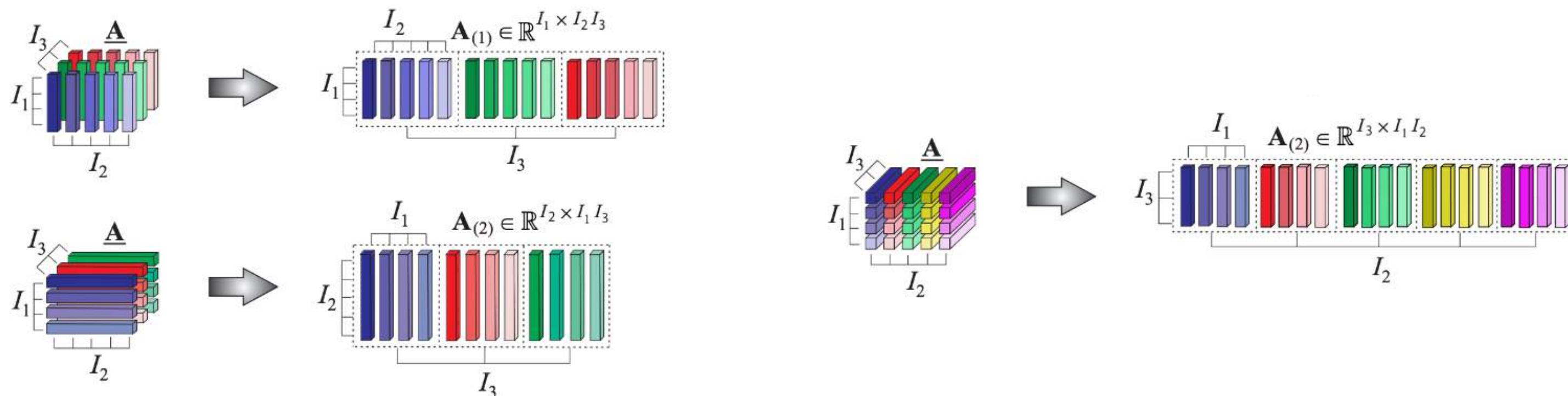
# Alternated Least Squares for CPD

- Loss functions looks easy, but its non-convex
- The ALS approach fixes B and C to solve for A, then fixes A and C to solve for B then etc.

$$\min_{\hat{\mathbf{x}}} \|\mathbf{x} - \hat{\mathbf{x}}\| \quad \text{with} \quad \hat{\mathbf{x}} = \sum_{r=1}^R \lambda_r \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r = [\![\boldsymbol{\lambda}; \mathbf{A}, \mathbf{B}, \mathbf{C}]\!].$$

# Unfolding

→ Matricization (flattening, unfolding)  
used in tensor reshaping



# Additional operations

- Khatri-Rao — blockwise
- Hadamard — elementwise
- Pseudo-inverse matrix

- ▶ Khatri-Rao product  $\mathbf{A} \odot \mathbf{B}$  is the columns-wise Kronecker product.

Setting  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_{J_1}]$  and  $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_{J_1}]$ ,

$$\mathbf{A} \odot \mathbf{B} = [\mathbf{a}_1 \boxtimes \mathbf{b}_1, \dots, \mathbf{a}_{J_1} \boxtimes \mathbf{b}_{J_1}].$$

- ▶ Example  $\mathbf{A} = \begin{bmatrix} a & b & c \\ d & e & f \end{bmatrix}, \mathbf{B} = \begin{bmatrix} g & h & i \\ j & k & l \end{bmatrix}$ .

$$\begin{aligned} \mathbf{A} \odot \mathbf{B} &= \begin{bmatrix} \mathbf{a}_1 \boxtimes \mathbf{b}_1 & \mathbf{a}_2 \boxtimes \mathbf{b}_2 & \mathbf{a}_3 \boxtimes \mathbf{b}_3 \end{bmatrix} \\ &= \begin{bmatrix} \begin{bmatrix} a \\ d \end{bmatrix} \boxtimes \begin{bmatrix} g \\ j \end{bmatrix} & \begin{bmatrix} b \\ e \end{bmatrix} \boxtimes \begin{bmatrix} h \\ k \end{bmatrix} & \begin{bmatrix} c \\ f \end{bmatrix} \boxtimes \begin{bmatrix} i \\ l \end{bmatrix} \end{bmatrix} \\ &= \begin{bmatrix} ag & bh & ci \\ aj & bk & cl \\ dg & eh & fi \\ dj & ek & fl \end{bmatrix}. \end{aligned}$$

“ $\odot$ ” denotes the Khatri-Rao product.

“ $\boxtimes$ ” represents the element-wise Hadamard product.

# Alternated Least Squares for CPD

- Pseudocode of basic ALS
- $X : I \times J \times K$
- $X_{(1)} : I \times JK$
- $B \text{ khatrio } C : JK \times R$
- $A^T A : R \times R$
- A update :  $(I \times JK) \times (JK \times R) \times (R \times R)$

---

**Algorithm 1: Basic ALS for the CP decomposition of a 3rd-order tensor**

---

**Input:** Data tensor  $\underline{X} \in \mathbb{R}^{I \times J \times K}$  and rank  $R$   
**Output:** Factor matrices  $\mathbf{A} \in \mathbb{R}^{I \times R}$ ,  $\mathbf{B} \in \mathbb{R}^{J \times R}$ ,  $\mathbf{C} \in \mathbb{R}^{K \times R}$ , and scaling vector  $\lambda \in \mathbb{R}^R$

- 1: Initialize  $\mathbf{A}, \mathbf{B}, \mathbf{C}$
- 2: **while** not converged or iteration limit is not reached **do**
- 3:    $\mathbf{A} \leftarrow \underline{X}_{(1)}(\mathbf{C} \odot \mathbf{B})(\mathbf{C}^T \mathbf{C} \circledast \mathbf{B}^T \mathbf{B})^\dagger$
- 4:   Normalize column vectors of  $\mathbf{A}$  to unit length (by computing the norm of each column vector and dividing each element of a vector by its norm)
- 5:    $\mathbf{B} \leftarrow \underline{X}_{(2)}(\mathbf{C} \odot \mathbf{A})(\mathbf{C}^T \mathbf{C} \circledast \mathbf{A}^T \mathbf{A})^\dagger$
- 6:   Normalize column vectors of  $\mathbf{B}$  to unit length
- 7:    $\mathbf{C} \leftarrow \underline{X}_{(3)}(\mathbf{B} \odot \mathbf{A})(\mathbf{B}^T \mathbf{B} \circledast \mathbf{C}^T \mathbf{C})^\dagger$
- 8:   Normalize column vectors of  $\mathbf{C}$  to unit length, store the norms in vector  $\lambda$
- 9: **end while**
- 10: **return**  $\mathbf{A}, \mathbf{B}, \mathbf{C}$  and  $\lambda$ .

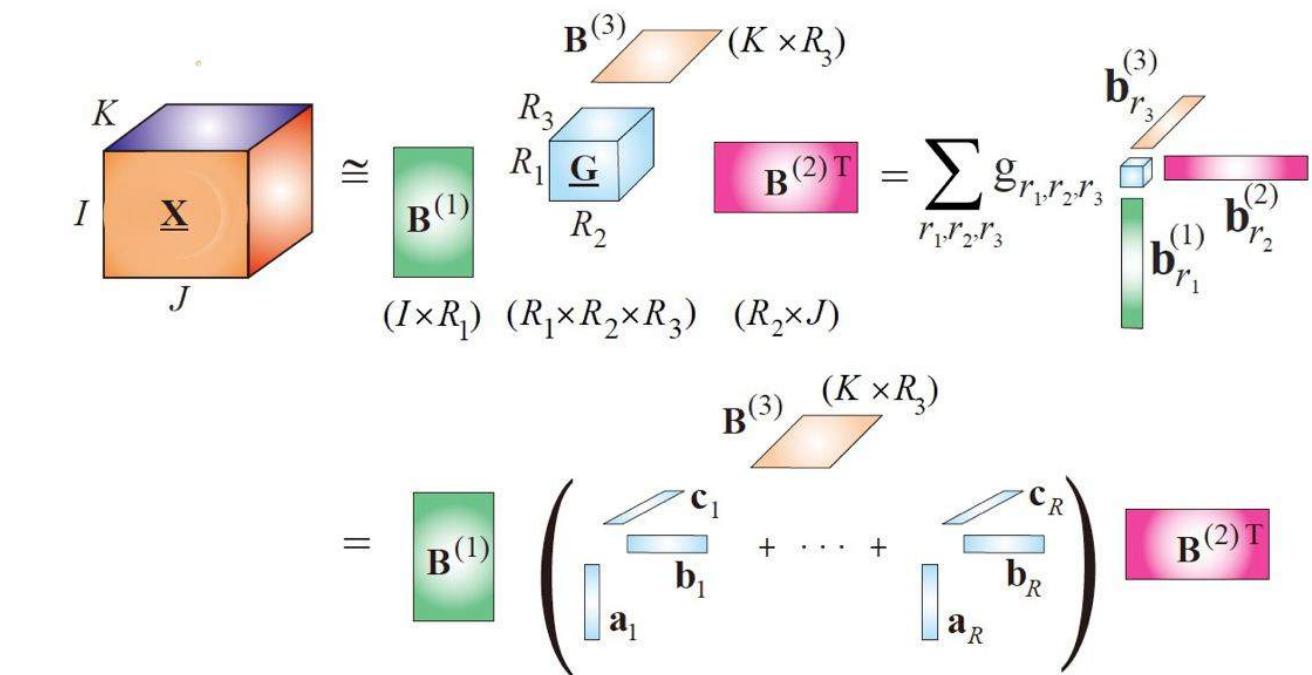
# Tucker decomposition

- $\mathbf{X}$  decomposed as a set of matrices (2D tensors) and one small core tensor (small 3D tensor)

A New Truncation Strategy for the Higher-Order Singular Value Decomposition | SIAM Journal on Scientific Computing ↗

Now publishers — Tensor Networks for Dimensionality Reduction and Large-scale Optimization: Part 1 Low-Rank Tensor Decompositions ↗

(a) Standard block diagrams of Tucker (top) and Tucker-CP (bottom) decompositions for a 3rd-order tensor

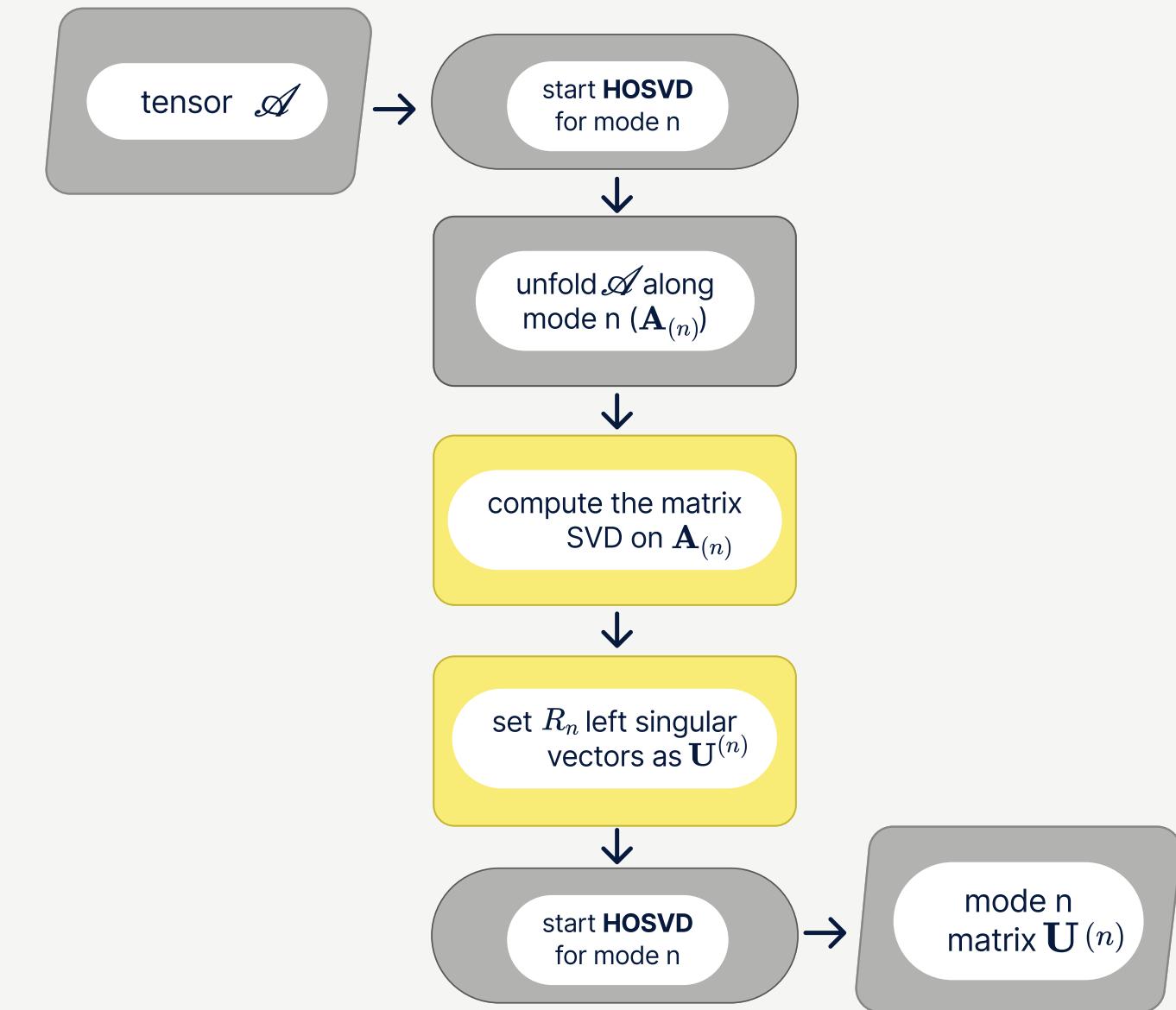


# Tucker decomposition

- Theorem by De Lathauwer, De Moor and Vandewalle proved in following ↗ paper

- SVD on every mode's tensor unfolding  $\mathbf{A}_{(n)}$
- Set basis factor matrices  $\mathbf{U}^{(n)}$  as  $R_n$  leading left singular vectors of  $\mathbf{A}_{(n)}$
- Derive core  $\mathcal{B}$  from original data and inverse factor matrices
- Defines a Tucker model with  $\mathcal{B} \mathbf{U}^{(n)}$

$$\mathcal{B} = \mathcal{A} \times_1 \mathbf{U}^{(1)-1} \times_2 \mathbf{U}^{(2)-1} \times_3 \mathbf{U}^{(3)-1}$$



# Tucker decomposition

- HOSVD is not the only solution,  
for instance see STHOSVD method ↗
- Sequential reduce original tensor
- No inverse matrices

Algorithm	4	Sequentially Truncated	HOSVD
<b>(STHOSVD) Algorithm</b>			
<b>Input :</b> A data tensor $\underline{\mathbf{X}} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ and a tolerance $\epsilon$ ;			
<b>Output:</b> Approximative HOSVD of the tensor $\underline{\mathbf{X}}$ $\underline{\mathbf{X}} \cong [\underline{\mathbf{S}}; \mathbf{Q}^{(1)}, \mathbf{Q}^{(2)}, \dots, \mathbf{Q}^{(N)}]$ and multilinear rank $(R_1, R_2, \dots, R_N)$			
<hr/>			
1	Set $\underline{\mathbf{S}} = \underline{\mathbf{X}}$		
2	<b>for</b> $n = 1, 2, \dots, N$ <b>do</b>		
3	$[\mathbf{Q}^{(n)}, \sim, \sim] = \text{truncated-svd}(\mathbf{S}_{(n)}, \frac{\epsilon}{N})$		
4	$\underline{\mathbf{S}} = \underline{\mathbf{S}} \times_n \mathbf{Q}^{(n)T}$		
5	<b>end</b>		

# Connections

→ There is analogy/connection between matrix and tensor decompositions

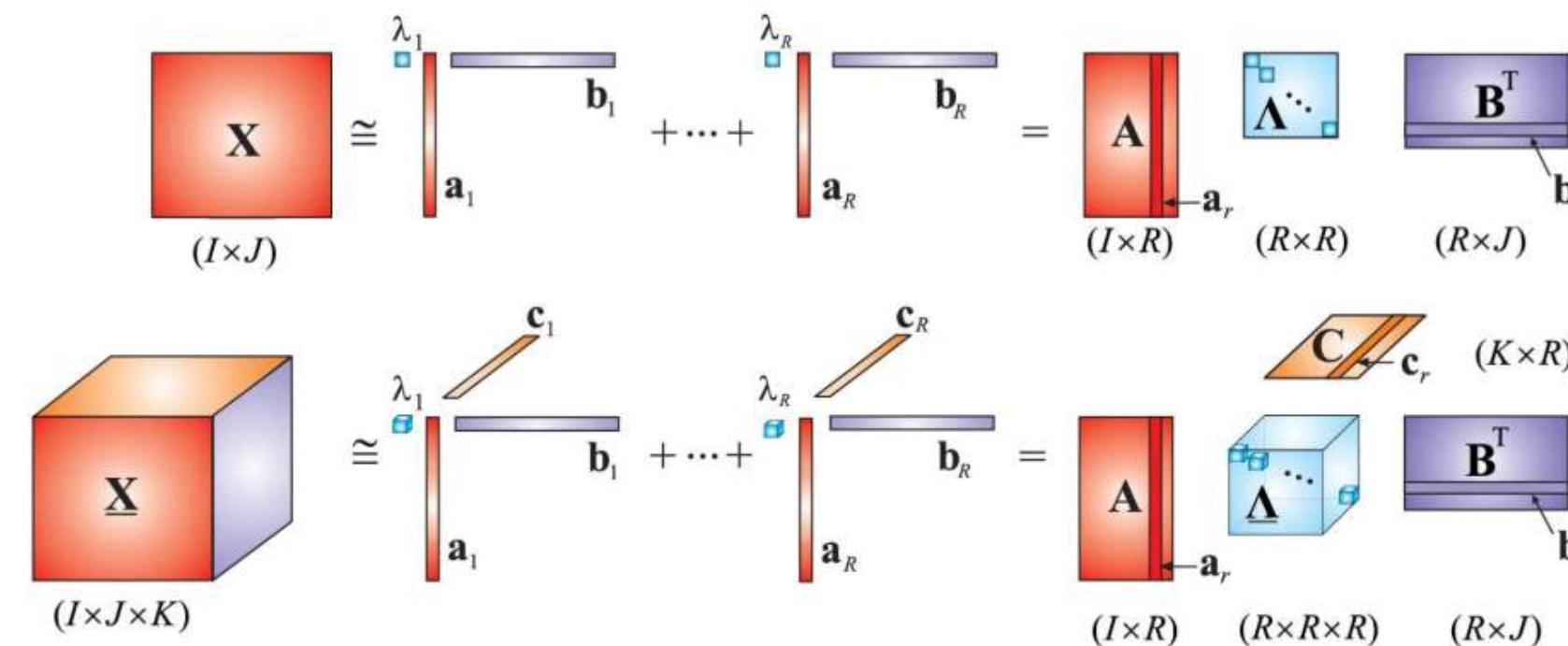


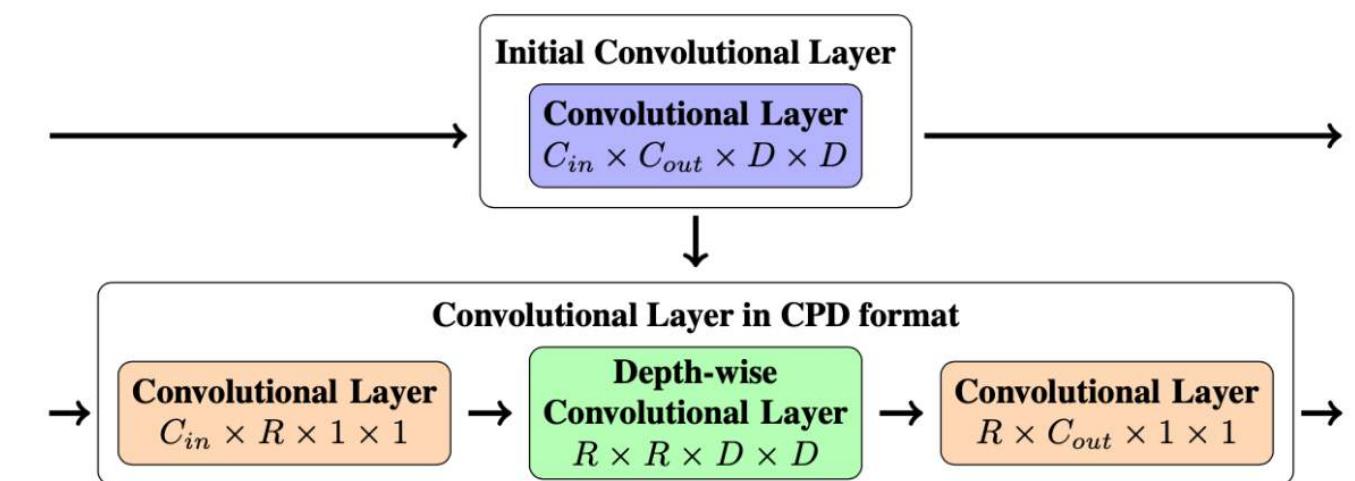
Figure 3.2: Analogy between a low-rank matrix factorization,  $\mathbf{X} \cong \mathbf{A}\Lambda\mathbf{B}^T = \sum_{r=1}^R \lambda_r \mathbf{a}_r \circ \mathbf{b}_r$  (top), and a simple low-rank tensor factorization (CP decomposition),  $\underline{\mathbf{X}} \cong \underline{\Lambda} \times_1 \mathbf{A} \times_2 \mathbf{B} \times_3 \mathbf{C} = \sum_{r=1}^R \lambda_r \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r$  (bottom).

# CPD as a depthwise convolution

- CPD allows to reduce number of parameters from  $S \times T \times D^2 \times D$  to  $(S + T + D) \times R$ .
- Applying ordinary CPD to convolutional kernel leads to instability (will cover later) for  $\hat{\mathcal{K}} = [\![\mathbf{A}, \mathbf{B}, \mathbf{C}]\!]$  the neural network

$$\mathcal{K} \approx \begin{matrix} [D^2 \times R] \\ C \\ A \\ B \end{matrix} \quad [S \times T \times D^2] \quad [S \times R] \quad [R \times T]$$

CPD scheme:  $\hat{\mathcal{K}} = [\![\mathbf{A}, \mathbf{B}, \mathbf{C}]\!]$



CP-decomposed Convolutional Layer

# CPD/Tucker and Convolutions

→ Tensor Networks can be combined to achieve better performance

## Stable Tensor Deomposition for Compression of Cnn

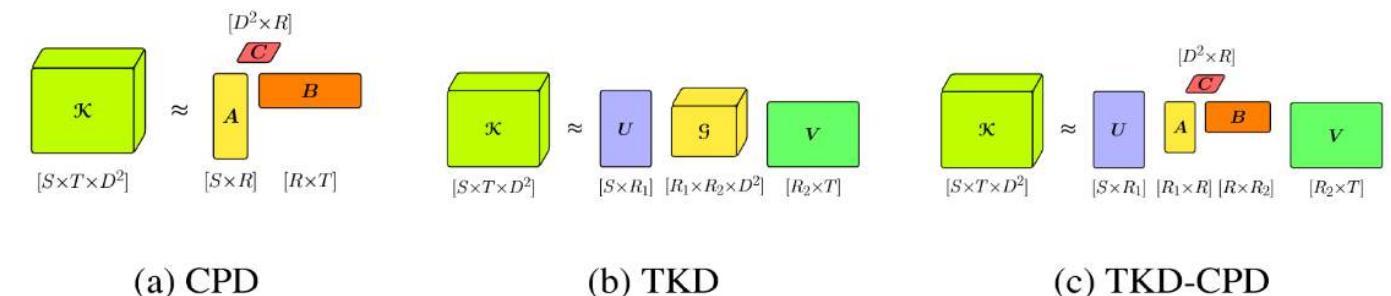
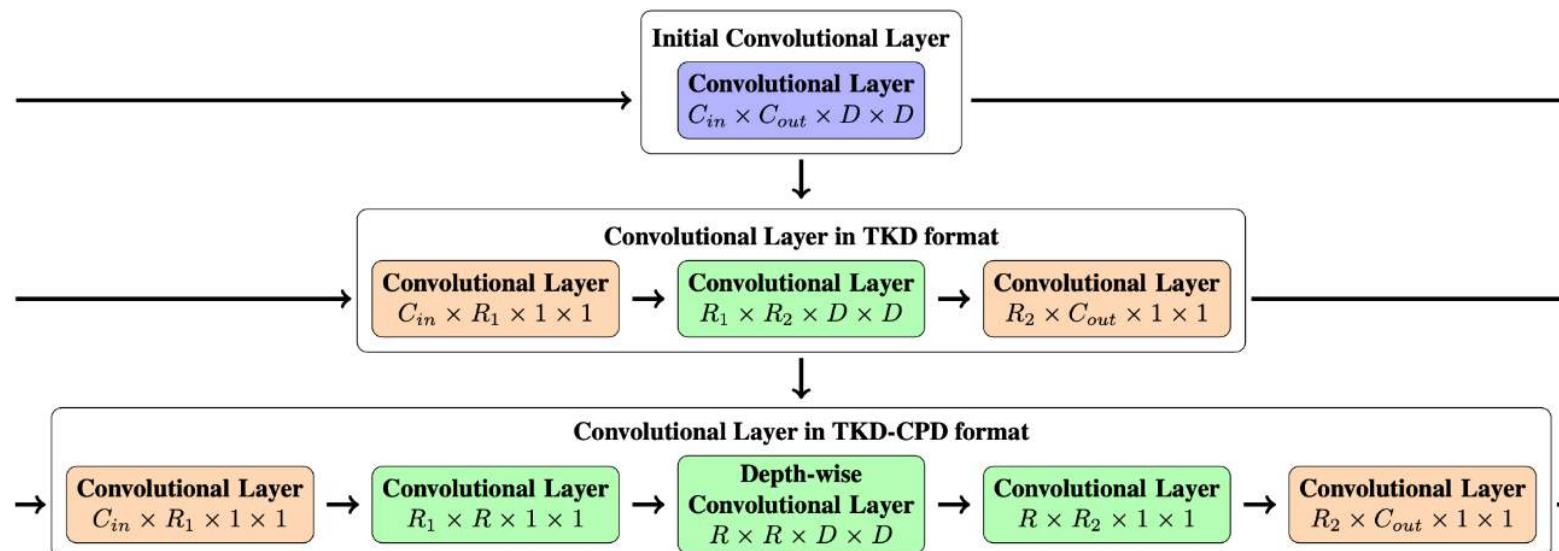
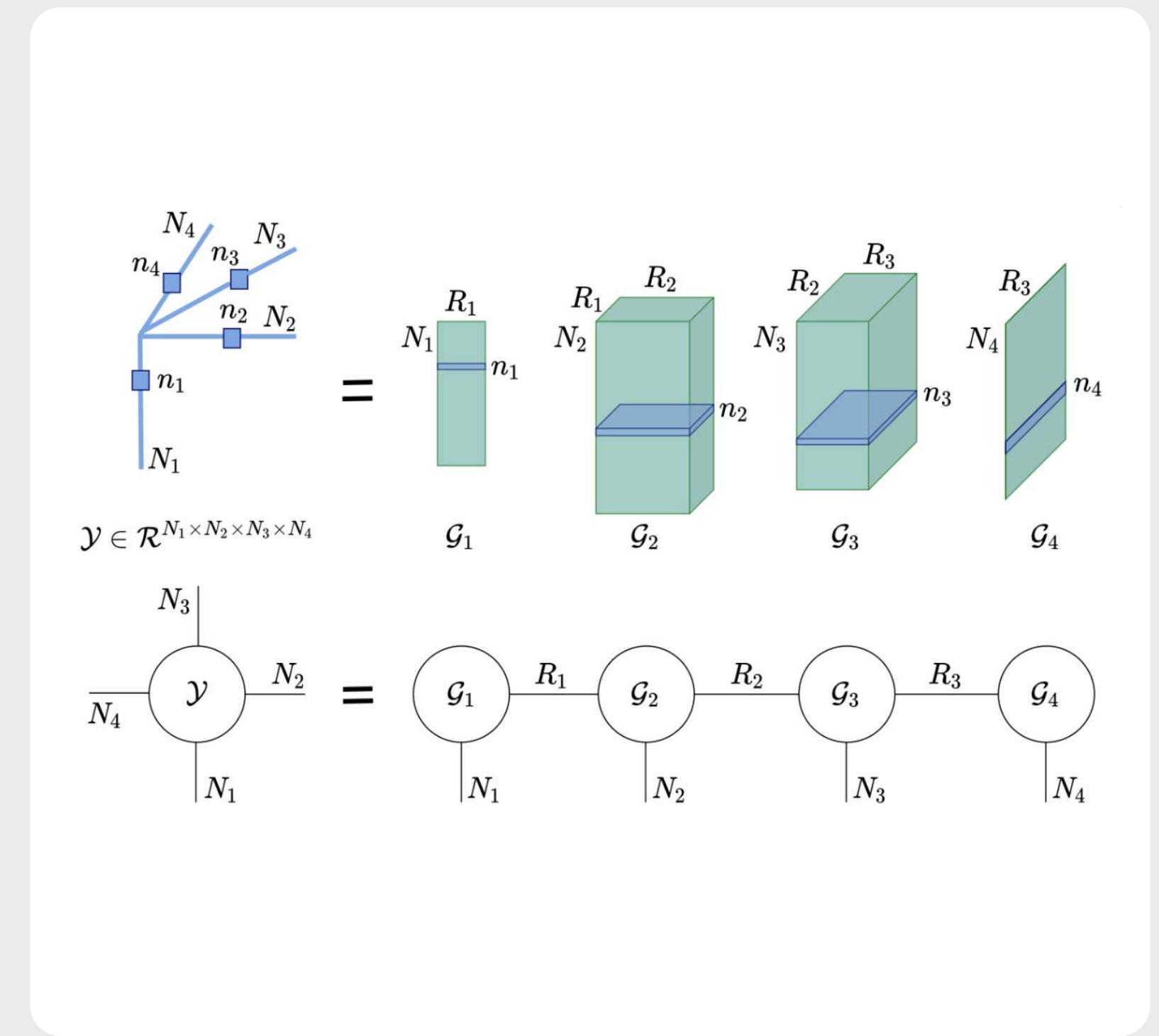


Fig. 1: Approximation of a third-order tensor using Canonical Polyadic tensor decomposition (CPD), Tucker-2 tensor decomposition (TKD), and their combination (TKD-CPD). CPD and TKD are common methods applied for CNN compression.

# Tensor Train

- **Tensor Train (TT)** is a contraction (linear chain) of 3d cores/factors (except first and last which are 2D)

$$\tilde{\mathbf{W}}_{c,y,x,t}^{\text{TT}} = \sum_{r_1=1}^{R_1} \sum_{r_2=1}^{R_2} \sum_{r_3=1}^{R_3} C_{c,r_1} \mathbf{Y}_{r_1,y,r_2} \mathbf{X}_{r_2,x,r_3} \mathbf{T}_{r_3,t}.$$



# Comparison of TT vs CPD and Tucker

→ Notice common index  
in every factor

$$\tilde{\mathbf{W}}_{c,y,x,t}^{\text{CP}} = \sum_{r=1}^R \mathbf{C}_{c,r} \mathbf{Y}_{y,r} \mathbf{X}_{x,r} \mathbf{T}_{t,r}.$$

$$\tilde{\mathbf{W}}_{c,y,x,t}^{\text{Tucker}} = \sum_{r_1=1}^{R_1} \sum_{r_2=1}^{R_2} \sum_{r_3=1}^{R_3} \sum_{r_4=1}^{R_4} \mathbf{G}_{r_1, r_2, r_3, r_4} \mathbf{C}_{c, r_1} \mathbf{Y}_{y, r_2} \mathbf{X}_{x, r_3} \mathbf{T}_{t, r_4}$$

$$\tilde{\mathbf{W}}_{c,y,x,t}^{\text{TT}} = \sum_{r_1=1}^{R_1} \sum_{r_2=1}^{R_2} \sum_{r_3=1}^{R_3} \mathbf{C}_{c, r_1} \mathbf{Y}_{r_1, y, r_2} \mathbf{X}_{r_2, x, r_3} \mathbf{T}_{r_3, t}.$$



# Tensor Network Libraries

[Teneva documentation ↗](#)

---

1

More than 80 packages on tensor networks available

---

2

See [THE LANDSCAPE OF SOFTWARE FOR TENSOR COMPUTATIONS ↗](#)

---

3

For ML and optimization notable packages are TensorLy TTPY, Teneva etc

---

# Practical and Theoretical Questions

[Teneva documentation ↗](#)

Which decomposition to choose



How to optimize/decompose given full tensor



How to select ranks



Greedy approach — expensive



AutoML/Hyperparameter Optimization  
approach — Optuna, RAY  
etc. less expensive, but harder to implement



How approximation error related to final accuracy



Layer-by-layer or a single shot?

# Advanced Ideas

# MUSCO

Multi-Stage Compression of neural networks

→ Iteratively we do:

1. Compress with **gradually decreased ranks**, using past approximations

2. finetune

→ automatic rank selection with **bayesian optimization**

→ opensource — you can try to play

[Click to read full PhD Thesis by E. Ponomarev ↗](#)

[MUSCO: Multi-Stage COmpression of neural networks ↗](#)

Чтобы сжать слой с весовым тензором  $\theta$  при заданном ранге  $R$  первый раз, решается следующая задача минимизации нормы Фробениуса между исходным тензором  $\theta$  и его факторизованным приближением  $\hat{\theta}^R$ :

$$\min_{\theta_1^R, \dots, \theta_N^R} \|\theta - \hat{\theta}^R\|, \quad \text{такие, что} \\ \mathcal{F}_{\text{fact}}(\hat{\theta}^R) = (\theta_1^R, \dots, \theta_N^R), \quad (1.6)$$

где  $\theta_1^R, \dots, \theta_N^R$  обозначают компоненты тензора в факторизованной форме с рангом  $R$  и числом факторов  $N$ .

Так как предложенный алгоритм сжатия MUSCO заведомо итеративный, то для второго и следующих сжатий слоя производится операция пересчета факторов  $\{\theta_n^R\}_{n=1}^N$ , полученных на предыдущем шаге метода. А именно, при заданном новом ранге  $R' : R' < R$  решается следующая задача минимизации:

$$\min_{\theta_1^{R'}, \dots, \theta_N^{R'}} \|\mathcal{F}_{\text{full}}(\theta_{\text{fact}}) - \hat{\theta}^{R'}\|, \quad \text{такие, что} \\ \theta_{\text{fact}} = (\theta_1^R, \dots, \theta_N^R), \\ \mathcal{F}_{\text{fact}}(\hat{\theta}^{R'}) = (\theta_1^{R'}, \dots, \theta_N^{R'}), \quad (1.7)$$

где  $\theta_1^{R'}, \dots, \theta_N^{R'}$  обозначают обновленные факторы исходного тензора весов выделенного слоя нейросетевой модели.

# MUSCO

## Multi-Stage Compression of neural networks

Model	FLOPs	mAP
<b>FASTER R-CNN (VGG-16) @ VOC2007</b>		
[7] baseline	1×	68.7
<i>Channel Pruning</i> [7]	4×	66.9(-1.8)
<i>Accelerating VD</i> [23]	4×	67.8(-0.9)
<i>AutoML Compression</i> [6]	4×	68.8(+0.1)
Used baseline	1.0×	71.1
Tucker2-iter (nx, 3.16)	3.16×	70.7(-0.4)
<b>MUSCO(nx, 1.77, 2)</b>	<b>3.72×</b>	<b>70.4(-0.7)</b>
<b>MUSCO(nx, 2, 2)</b>	<b>4.69×</b>	<b>70.1(-1.0)</b>
Tucker2-iter (nx, 10)	9.67×	68.6(-2.5)
<b>MUSCO(nx, 3.16, 2)</b>	<b>10.49×</b>	<b>69.2(-1.9)</b>
<b>MUSCO(nx, 1.77, 4)</b>	<b>13.95×</b>	<b>68.3(-2.8)</b>

Table 1: Comparison of Faster R-CNN (with VGG-16 backbone) compressed models on VOC2007 evaluation dataset. MUSCO(nx, 3.16, 2) is a compressed model obtained after two compression steps using  $3.16\times$  parameter reduction at each step

**Алгоритм 1** Алгоритм итеративной аппроксимации малого ранга для автоматического сжатия сети

**Input:** Предварительно обученная исходная модель,  $M$

**Output:** Тонко настроенная сжатая модель,  $M^*$ .

- 1:  $M^* \leftarrow M$
- 2: **while** желаемая степень сжатия не достигнута или автоматически выбранные ранги не стабилизировались **do**
- 3:      $R \leftarrow$  автоматически выбранные ранги для низкоранговых тензорных аппроксимаций сверточных и полно связных весовых тензоров.
- 4:      $M \leftarrow$  (дополнительно) сжатая модель, полученная из  $M$  заменой весов слоев их тензорными аппроксимациями ранга  $R$ .
- 5:      $M^* \leftarrow$  тонко настроенная модель  $M$ .
- 6: **end while**

# Solution Sensitivity

Sensitivity of the tensor  $\mathcal{T} = [\![\mathbf{A}, \mathbf{B}, \mathbf{C}]\!]$  of size  $I \times J \times K$  is a measure of factorized tensor norm change with respect to perturbations in individual factor matrices

$$\text{ss}([\![\mathbf{A}, \mathbf{B}, \mathbf{C}]\!]) = \lim_{\sigma^2 \rightarrow 0} \frac{1}{R\sigma^2} E\{\|\mathcal{T} - [\![\mathbf{A} + \delta\mathbf{A}, \mathbf{B} + \delta\mathbf{B}, \mathbf{C} + \delta\mathbf{C}]\!]\|_F^2\}$$

where  $\delta\mathbf{A}, \delta\mathbf{B}, \delta\mathbf{C}$  have random i.i.d. elements from  $N(0, \sigma^2)$ .

Sensitivity can be computed as

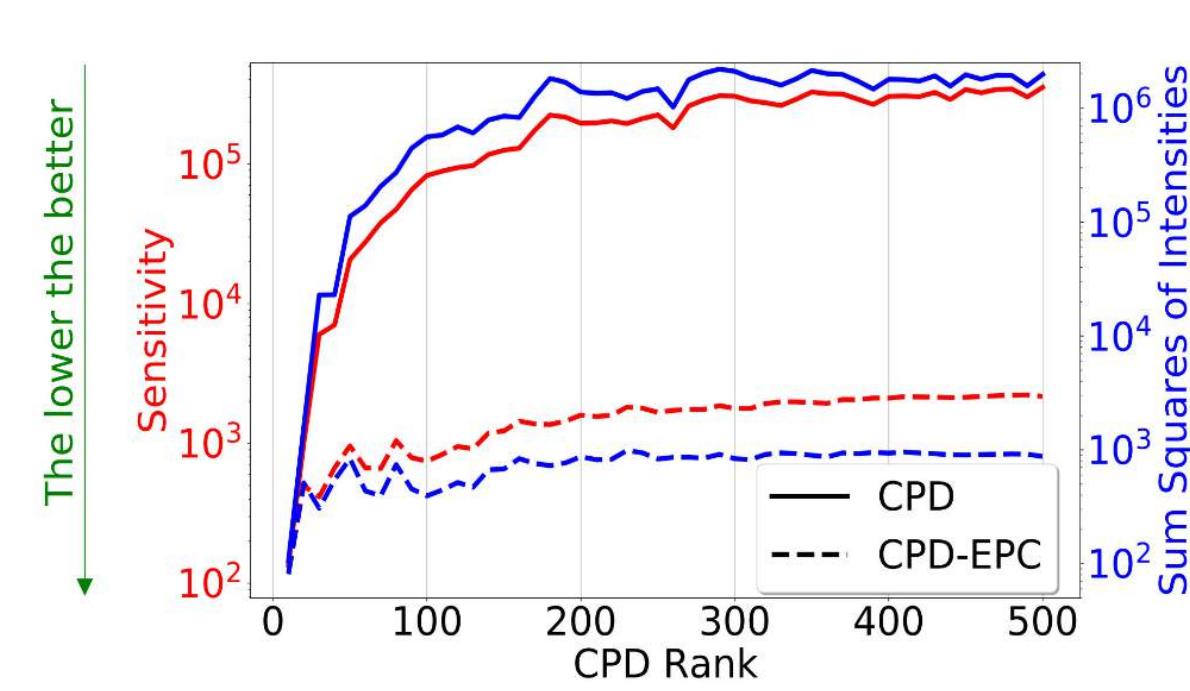
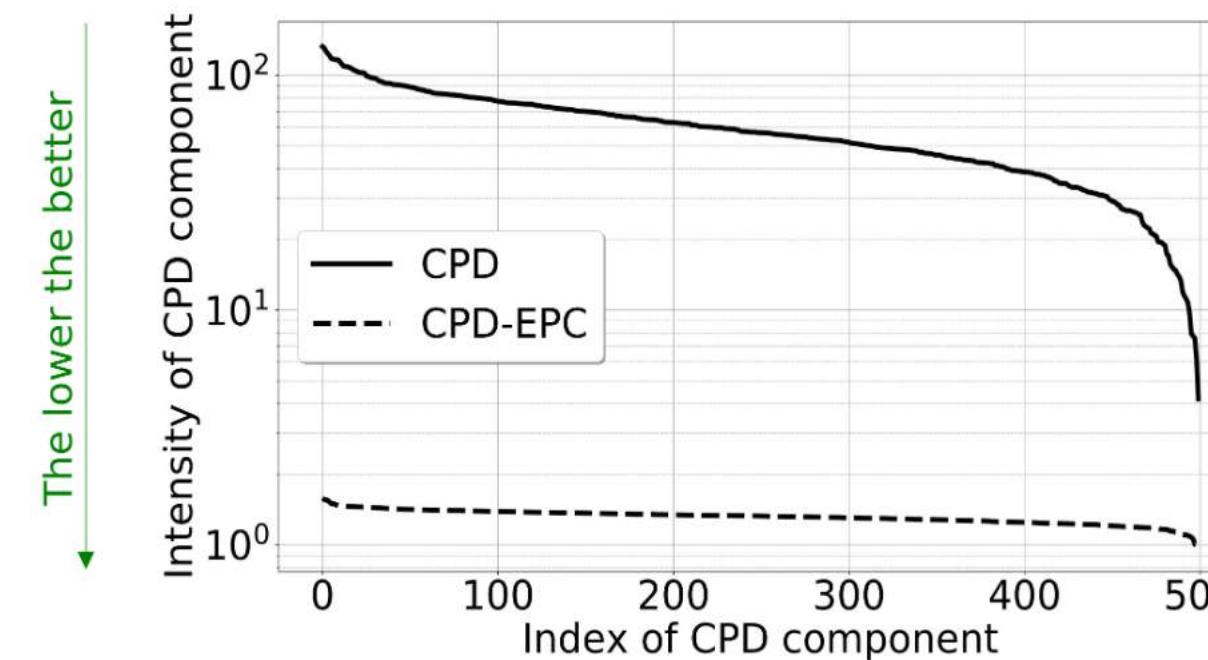
$$\text{ss}([\![\mathbf{A}, \mathbf{B}, \mathbf{C}]\!]) = K \text{tr}\{(\mathbf{A}^T \mathbf{A}) \circledast (\mathbf{B}^T \mathbf{B})\} + I \text{tr}\{(\mathbf{B}^T \mathbf{B}) \circledast (\mathbf{C}^T \mathbf{C})\} + J \text{tr}\{(\mathbf{A}^T \mathbf{A}) \circledast (\mathbf{C}^T \mathbf{C})\}$$

# Solution CPD-EPC

CPD  
with minimal  
sensitivity:

$$\begin{aligned} & \min_{\{\mathbf{A}, \mathbf{B}, \mathbf{C}\}} && \text{ss}([\![\mathbf{A}, \mathbf{B}, \mathbf{C}]\!]) \\ & \text{s.t.} && \|\mathcal{K} - [\![\mathbf{A}, \mathbf{B}, \mathbf{C}]\!]\|_F^2 \leq \delta^2 \end{aligned}$$

Bound  $\delta^2$  can be the approximation error of the CPD with diverging components.



# Stable CPD

- No python source code available unfortunately due to industrial copyright
- There is matlab prototype for Stable CPD

GitHub — phananhuy/TensorBox: Matlab package for Tensor decompositions ↗

Table 1: Comparison of different model compression methods on ILSVRC-12 validation dataset. The baseline models are taken from Torchvision.

Model	Method	↓ FLOPs	Δ top-1	Δ top-5
VGG-16	Asym. [57]	≈ 5.00	-	-1.00
	TKD+VBMF [26]	4.93	-	-0.50
	<b>Our</b> (EPS <sup>1</sup> =0.005)	<b>5.26</b>	<b>-0.92</b>	<b>-0.34</b>
ResNet-18	Channel Gating NN [24]	1.61	-1.62	-1.03
	Discrimination-aware Channel Pruning [58]	1.89	-2.29	-1.38
	FBS [13]	1.98	-2.54	-1.46
	MUSCO [14]	2.42	-0.47	-0.30
ResNet-50	<b>Our</b> (EPS <sup>1</sup> =0.00325)	<b>3.09</b>	<b>-0.69</b>	<b>-0.15</b>
	<b>Our</b> (EPS <sup>1</sup> =0.0028)	<b>2.64</b>	<b>-1.47</b>	<b>-0.71</b>

<sup>1</sup> EPS: accuracy drop threshold. Rank of the decomposition is chosen to maintain the drop in accuracy lower than EPS.

# Data and tensors

→ Inference time on different hardware

Table 2: Inference time and acceleration for ResNet-50 on different platforms.

Platform	Model inference time	
	Original	Compressed
Intel® Xeon®Silver 4114 CPU 2.20 GHz	$3.92 \pm 0.02$ s	$2.84 \pm 0.02$ s
NVIDIA®Tesla®V100	$102.3 \pm 0.5$ ms	$89.5 \pm 0.2$ ms
Qualcomm®Snapdragon™845	$221 \pm 4$ ms	$171 \pm 4$ ms

# Low-rank transformers

# Stable CPD

- Nice Paper from Huawei
- Many other papers but it might hard-to-follow for beginners
- kinda MUSCO + KD + some tricks

In progressive LRD, we use a rank selection approach in which, we start from a lower compression ratio and apply low rank decomposition to the entire model altogether (just like the single-shot method). The model is then fine-tuned for some epochs in order to recover the accuracy. The rank  $R$  is decreased incrementally in the next steps so that we have a higher compression ratio. This process is repeated until the desired compression ratio is achieved for the model. Starting from the second compression step, we only decompose the second decomposed layer of previously decomposed layers and multiply them so that it results in two decomposed layers again. This way, the number of layers will not increase again in the next stages of decomposition as we use matrix multiplication for merging the newly generated weight matrices in order to keep the number of layers the same as the first stage of LRD. This process is shown in Fig. 1.

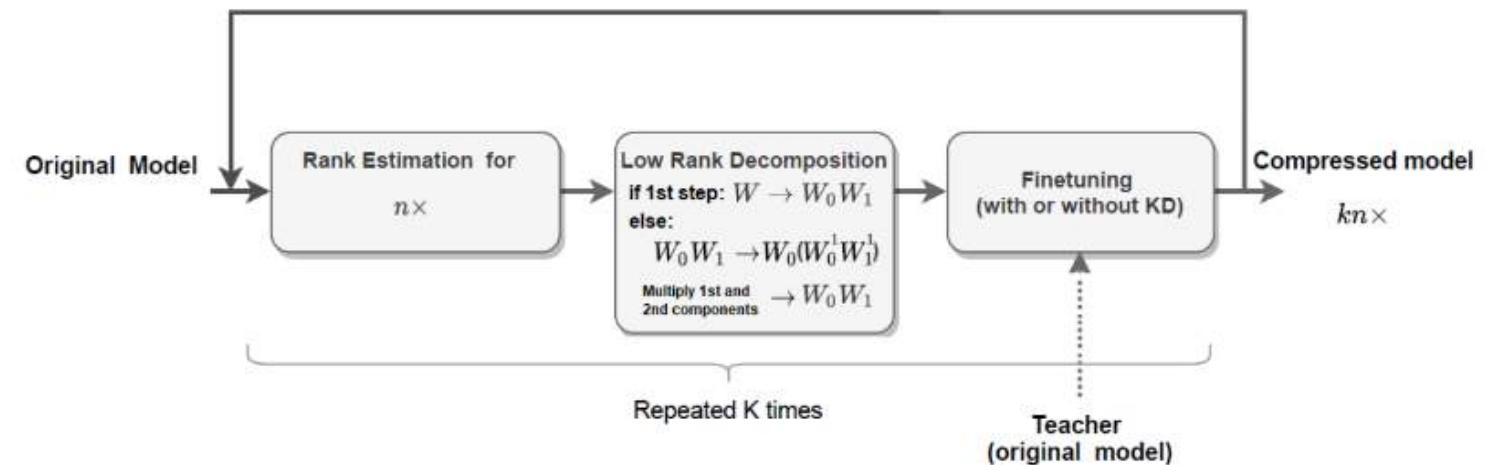


Figure 1: Progressive decomposition scheme

Strategies for Applying Low Rank Decomposition to Transformer-Based Models ↗

# Progressive low-rank compression for ViTå

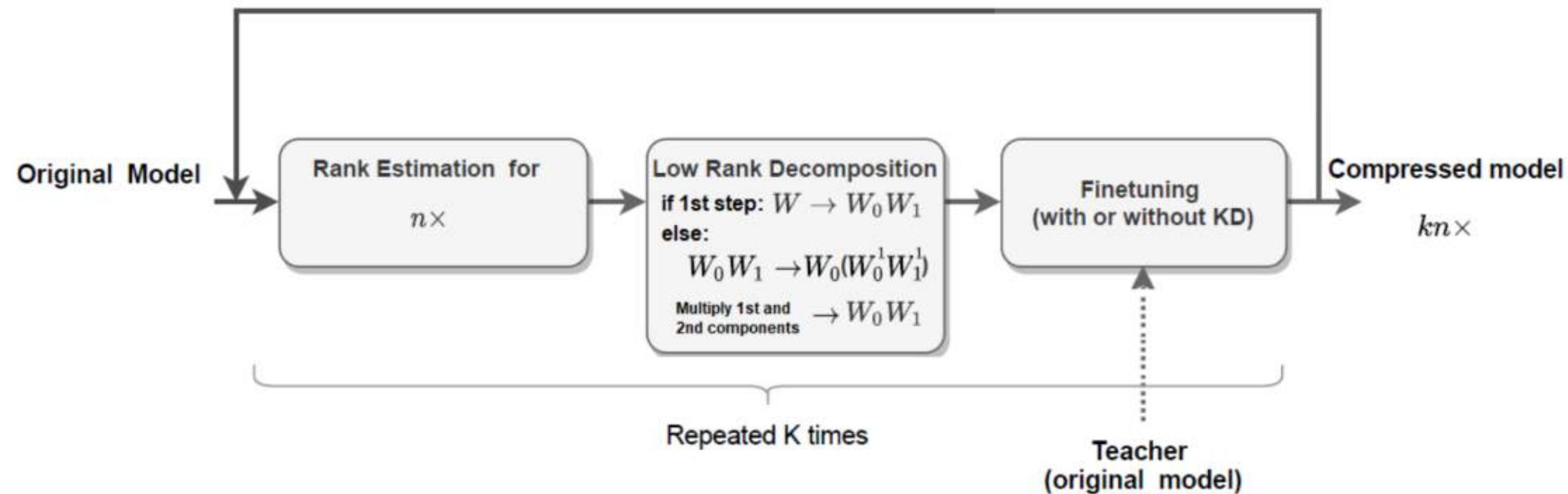


Figure 1: Progressive decomposition scheme

# Progressive low-rank compression for ViT<sup>å</sup>

Table 1: Experimental results for ViT model decomposed LRD using different strategies including single-shot, layer-by-layer and progressive LRD. The model is pretrained on ImageNet-21k and finetuned on CIFAR-10 and CIFAR-100 datasets.

<b>Method</b>	<b># LRD Steps</b>	<b>Comp Ratio</b>	<b>CIFAR-10 Top-1 (%)</b>	<b>CIFAR-100 Top-1 (%)</b>
ViT (Org)	0	1×	98.76	92.95
Single-Shot	1	2.63×	94.53	87.09
Layer-by-Layer	4	2.63×	95.36	87.51
Progressive	4	2.63×	96.01	88.62

# Summary

---

Tensor methods provides:

An elegant tool to reduce overparameterization of neural networks

---

# Homework

Using provided template and previous homework:

Implement SVD layer

1

Replace linear layers in decoder head

2

Compare it within different settings

3

# Links and References

# Further readings

- 
1. [Tensor Networks for Dimensionality Reduction and Large-scale Optimization: Part 1](#) ↗
  2. [Tensor Networks for Dimensionality Reduction and Large-scale Optimization: Part 2](#) ↗
  3. [Stable Low-rank Tensor Decomposition for Compression of Convolutional Neural Network](#) ↗
  4. [MUSCO: Multi-Stage Compression of neural networks](#) ↗
  5. [Quantization Aware Factorization for Deep Neural Network Compression](#) ↗
  6. [How Informative is the Approximation Error from Tensor Decomposition for Neural Network Compression? | OpenReview](#) ↗
-