

Tensor Factorization and Tensor Networks for Machine Learning

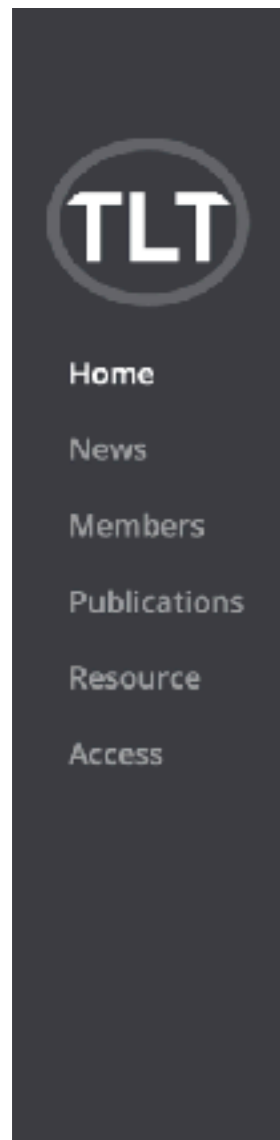
Qibin Zhao

Tensor Learning Team
RIKEN AIP

<https://qibinzhao.github.io>



<https://qibinzhao.github.io>



Tensor Learning Team

RIKEN Center for Advanced Intelligence Project (AIP).

Email: [qibin.zhao \[at\] riken.jp](mailto:qibin.zhao@riken.jp)

About Us

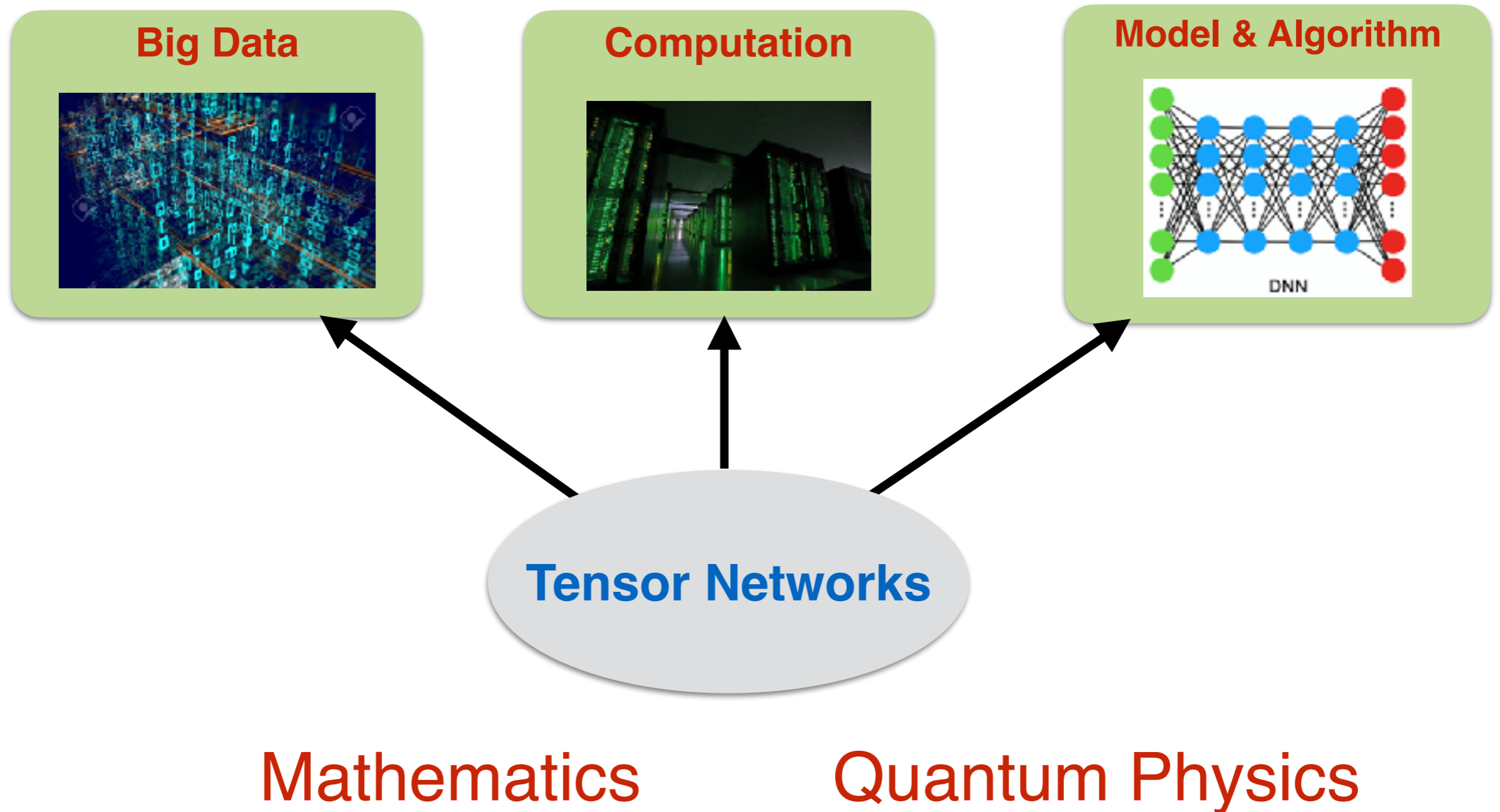
We study various tensor-based machine learning technologies, e.g., tensor decomposition, multilinear latent variable model, tensor regression and classification, tensor networks, deep tensor learning, and Bayesian tensor learning, with aim to facilitate the learning from high-order structured data or large-scale latent space. Our goal is to develop innovative, scalable and efficient tensor learning algorithms supported by theoretical principles. The novel applications in computer vision and brain data analysis will also be explored to provide new insights into tensor learning methods. ([RIKEN AIP TLT Web](#))

Job Opportunities

Seeking Openings for Research Scientists or Postdoctoral Researchers or Technical Scientist or Internship Student. ([Job](#))



Success of Machine Learning



Role of Tensor and Tensor Networks

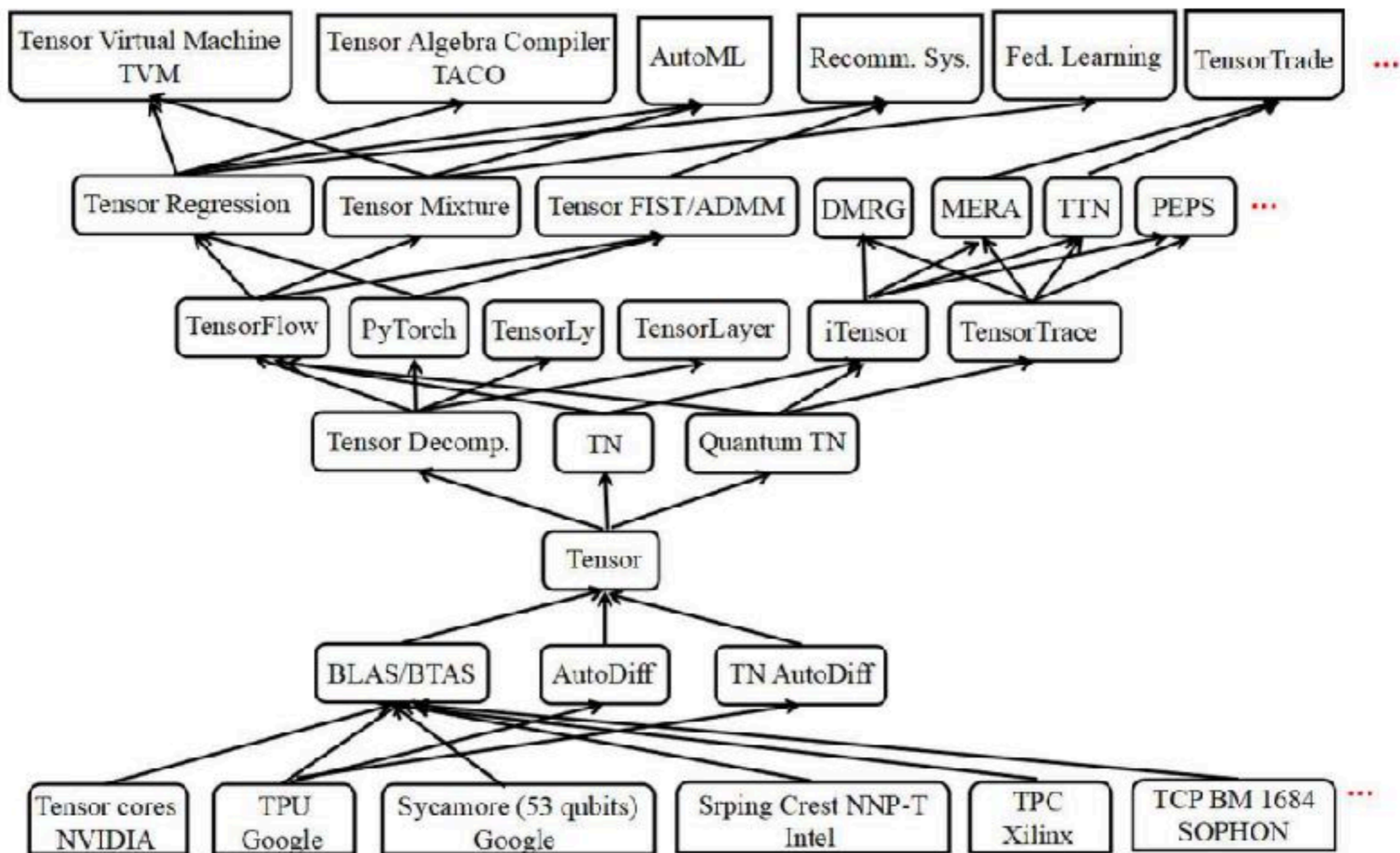


Figure 1: The hourglass architecture of tensor, tensor networks and quantum tensor networks in machine learning.

(Liu et al. IJCAI Workshop 2020)

Agenda

To introduce tensor networks fundamental, and to overview recent progress of tensor networks applied to machine learning.

Part I

**Tensor Methods
for Data
Representation**

Part II

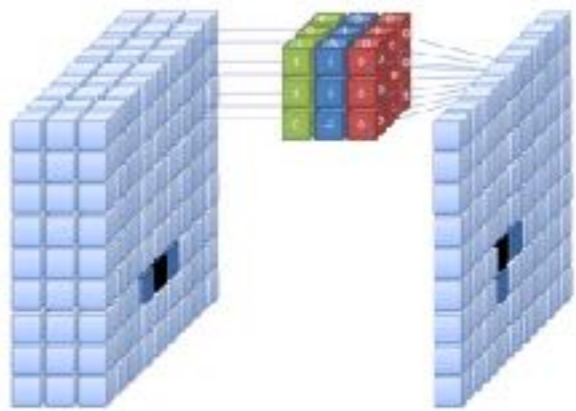
**Tensor Networks
in Deep Learning
Modeling**

Part III

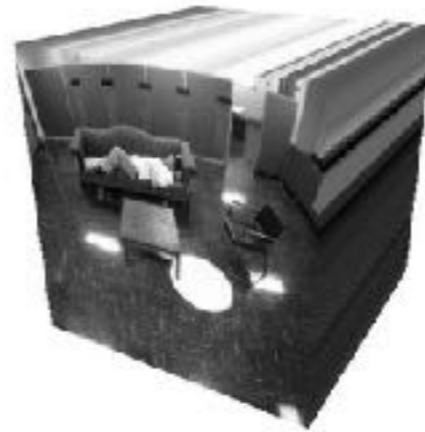
**Frontiers and
Future Trends**

Tensor Networks for Data Representation

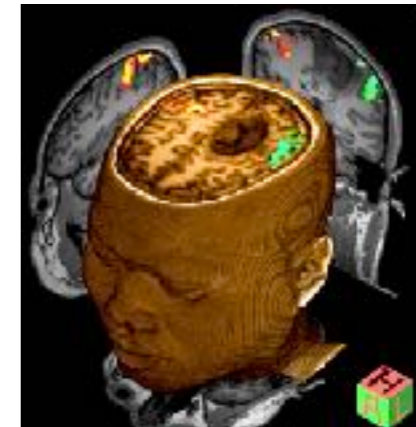
Multiway Structured Data



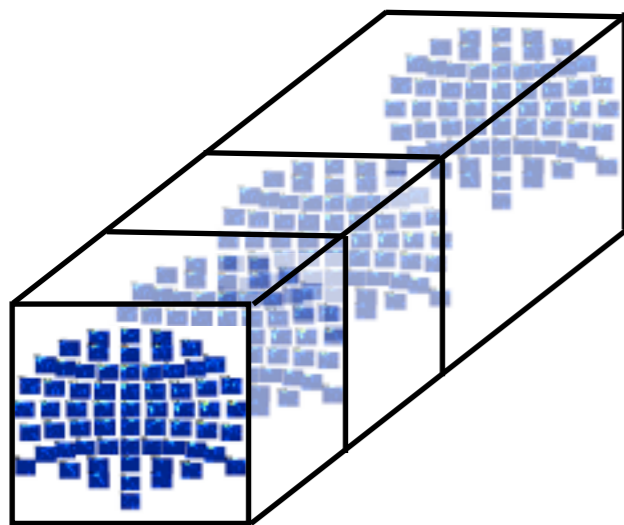
Feature Maps in CNN
(Spatial x Spatial x Filter)
WIKIMEDIA COMMONS



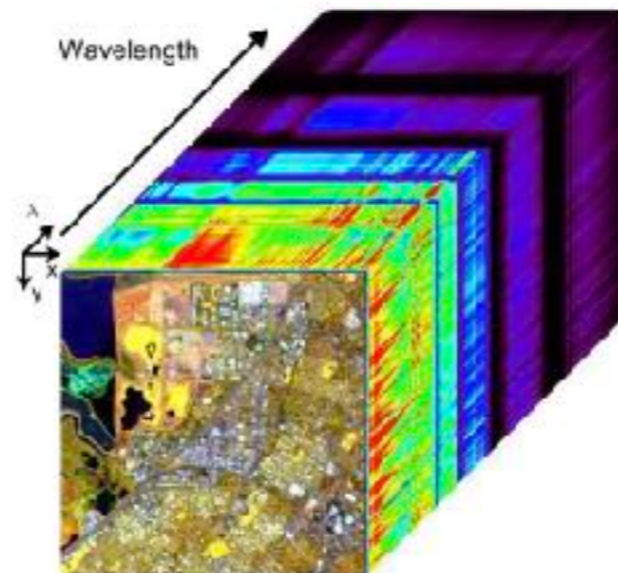
Video Data
(Spatial x Spatial x Time)



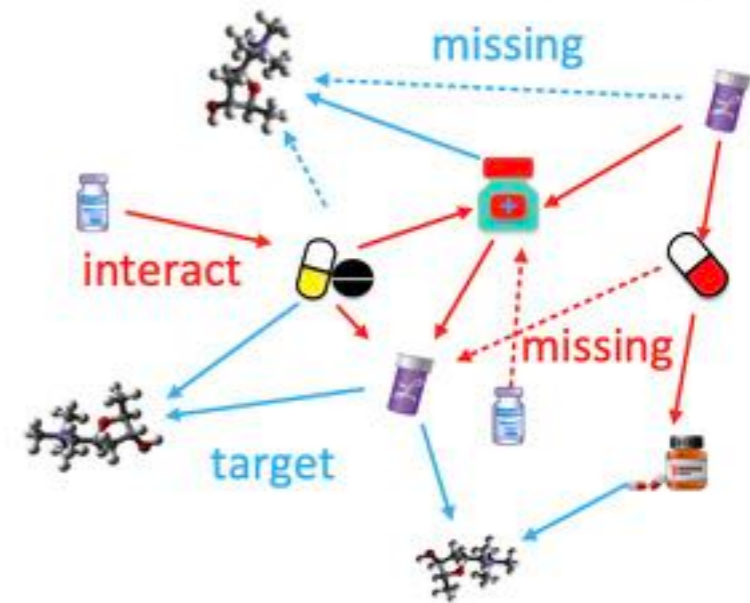
fMRI Data
(Spatial x Spatial x Spatial)



EEG
(Spatial x Time x Frequency)

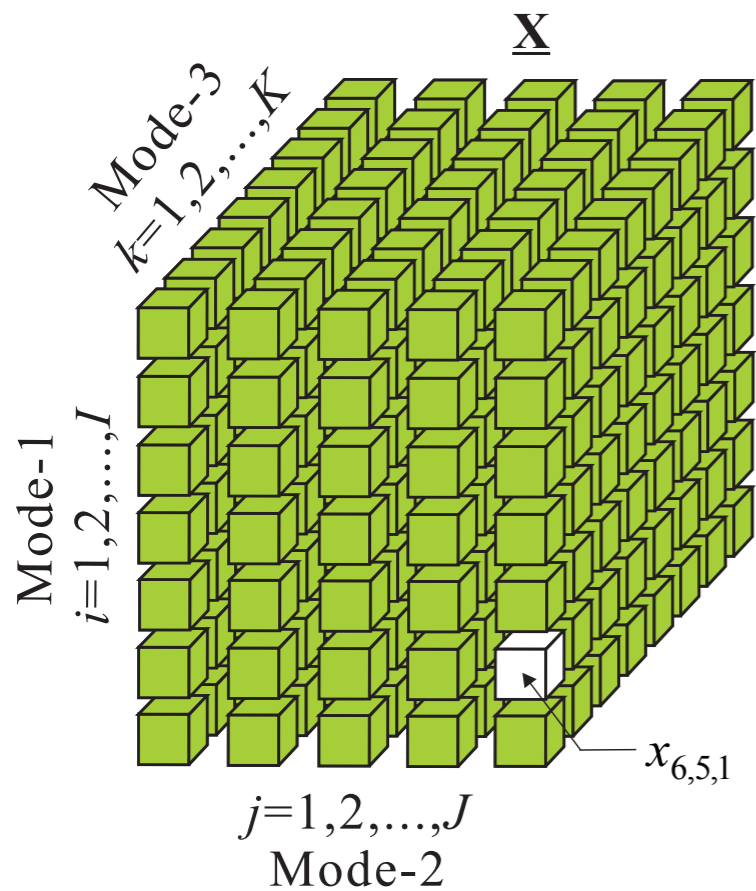


Hyperspectral Image
(Spatial x Spatial x Spectral)

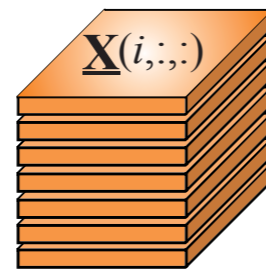


Medical Knowledge Graph
(Entity x Entity x Relation)
(Wang et al., 2017)

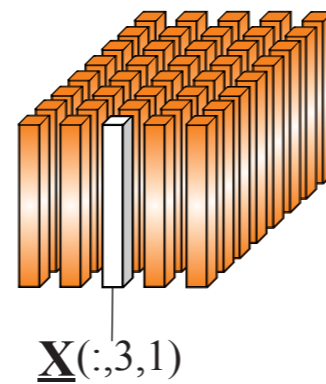
Tensor: Generalization of Vector and Matrix



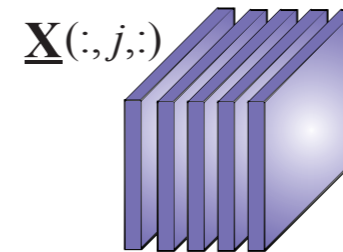
Horizontal Slices



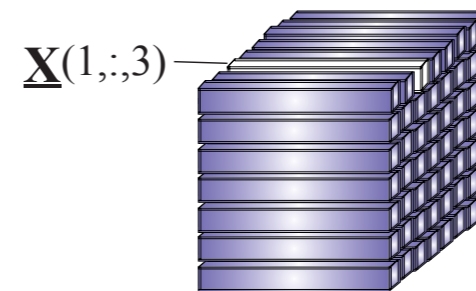
Column (Mode-1)
Fibers



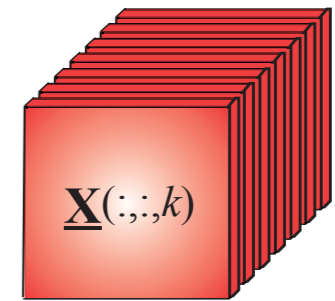
Lateral Slices



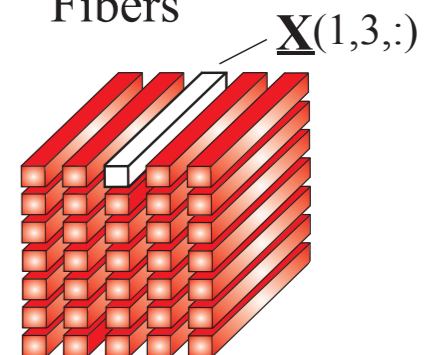
Row (Mode-2)
Fibers



Frontal Slices

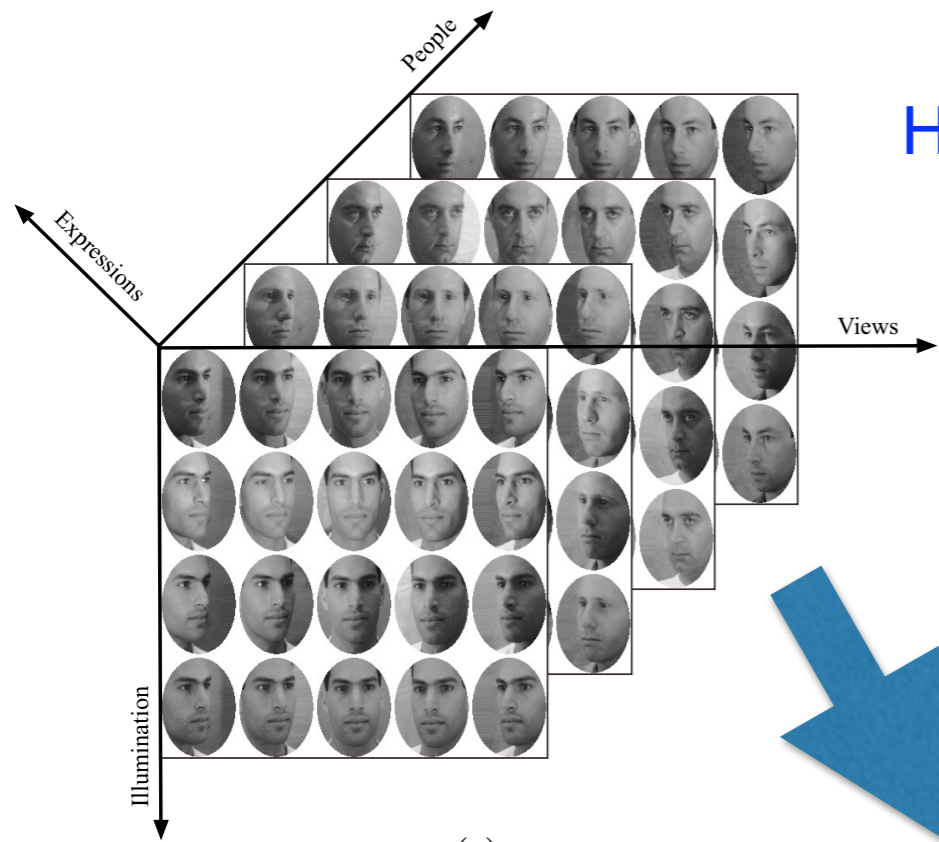


Tube (Mode-3)
Fibers



Multilinear Dimension Reduction

- ▶ Multi-linear extension of SVD, PCA and ICA

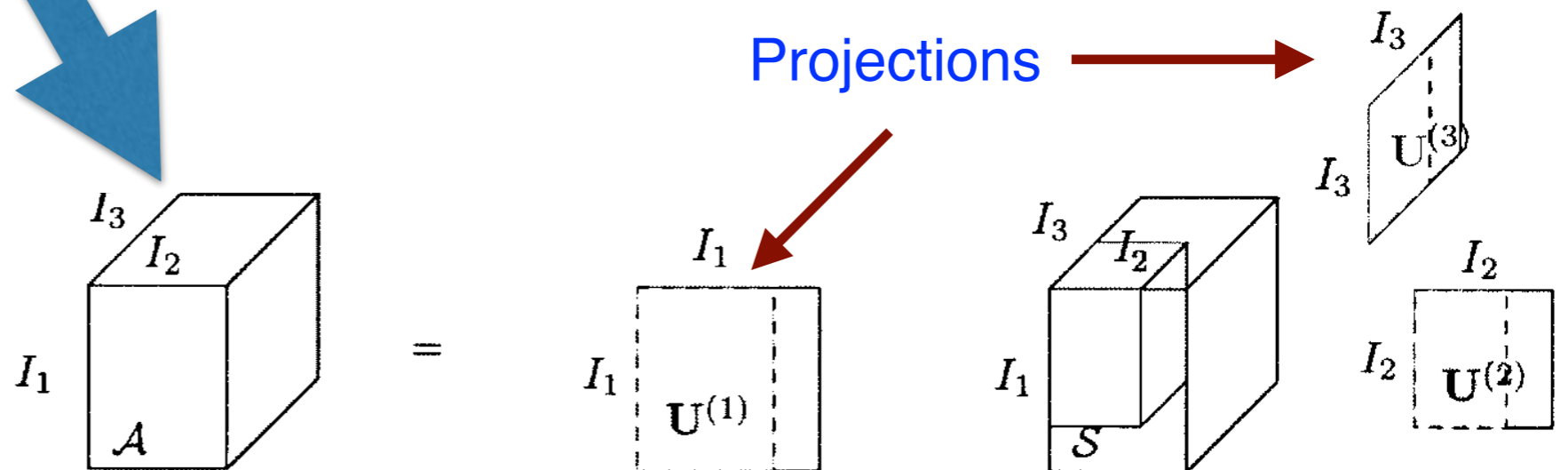


Higher-order singular value decomposition

$$\mathcal{A} = \mathcal{S} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \cdots \times_N \mathbf{U}^{(N)}$$

↓
Orthonormal

Projections →

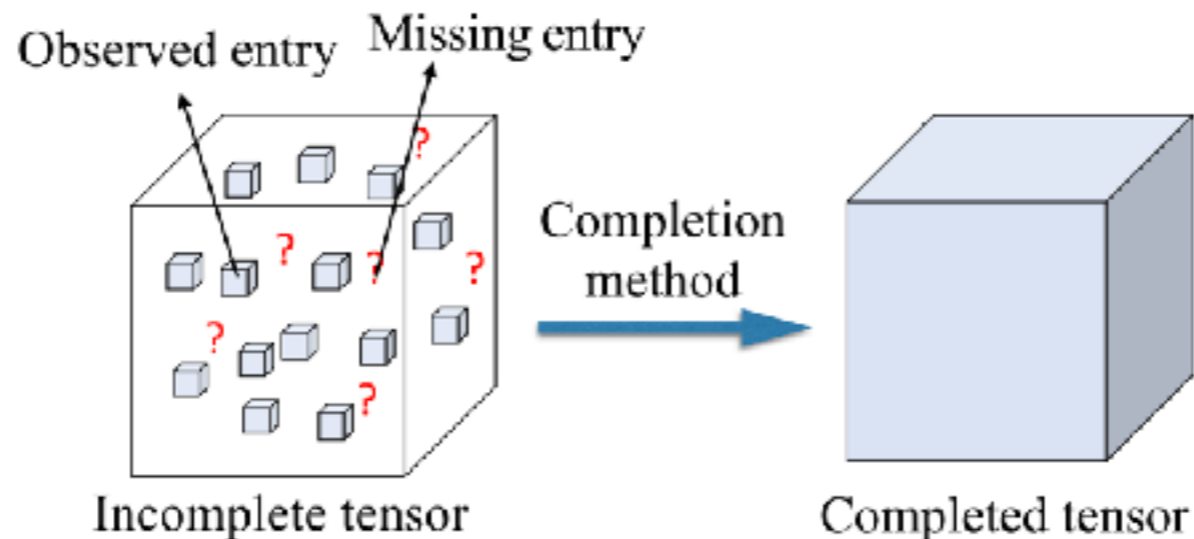
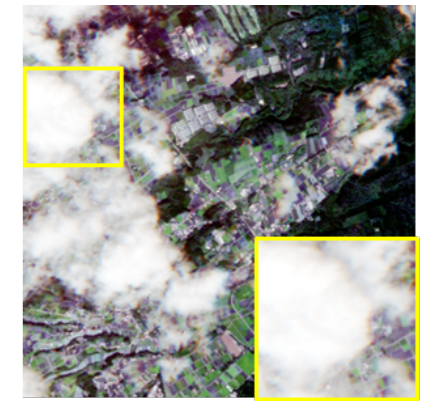
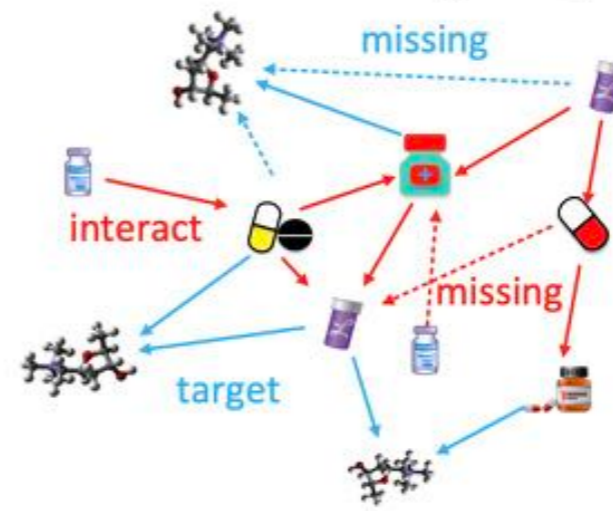
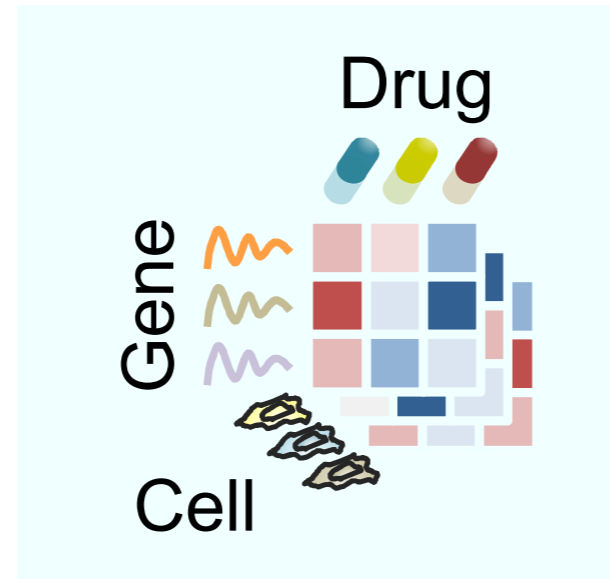
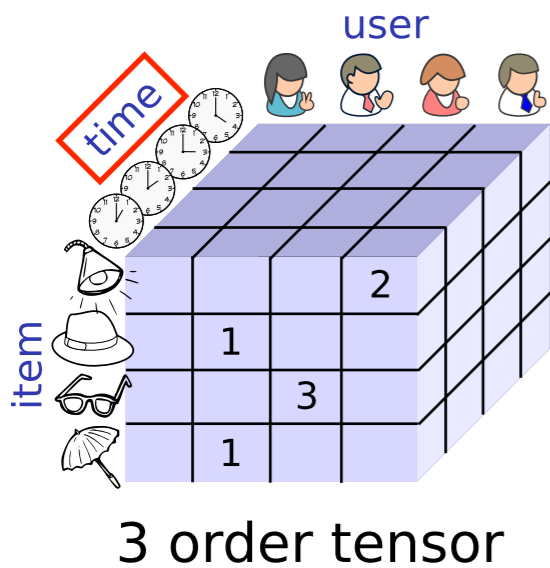


Multilinear Analysis of Image Ensembles: TensorFaces (Vasilescu et al., ECCV 2002)

A Multi-linear Singular Value Decomposition (Lathauwer et al, SIAM J. Matrix Anal. Appl., 2000)

Missing Values Problems in ML

- ▶ Recommender system (user x item x time), Knowledge graph prediction
- ▶ Image inpainting/denoising, drug repositioning



Tensor Completion

Problem formulation (low-rank approximation)

$$\min_{\mathcal{X}} \text{Rank}(\mathcal{X}), \quad s.t. \mathcal{Y}_{\Omega} = \mathcal{X}_{\Omega}$$

Main Approaches:

- ▶ Low-rank approximation via **convex optimization** (**high computation cost**)

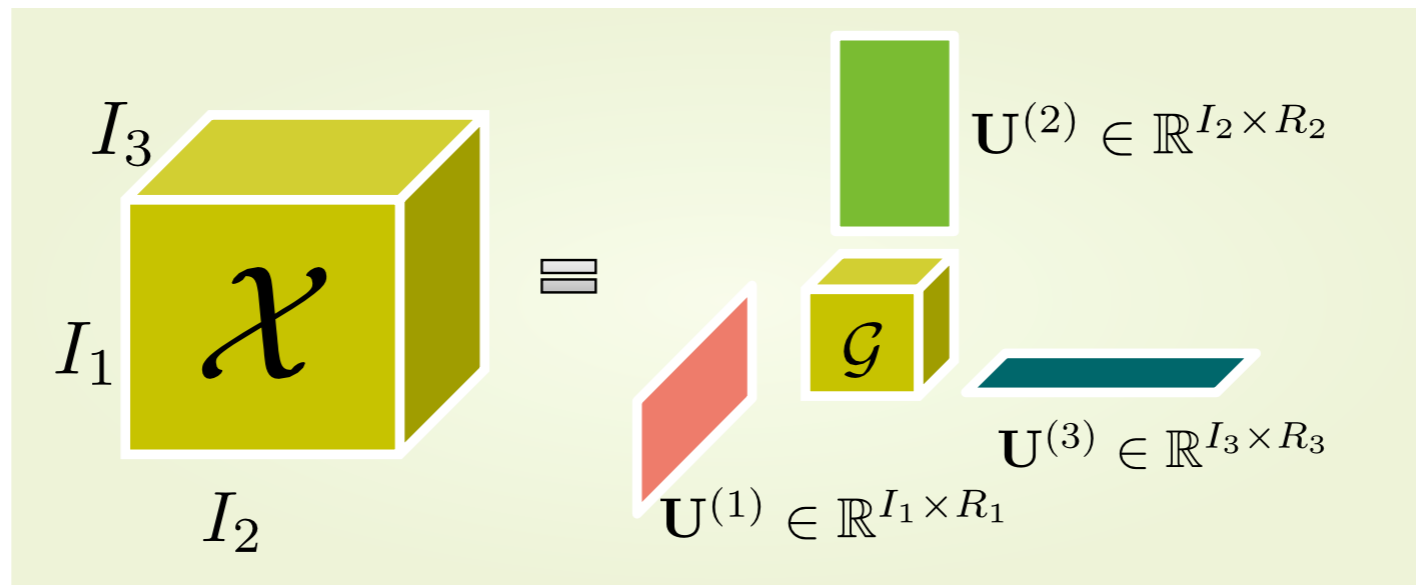
$$\min_{\mathcal{X}} \|\Omega * (\mathcal{Y} - \mathcal{X})\|_F^2 + \lambda \|\mathcal{X}\|_*$$

- ▶ **Decomposition** based approach (**rank selection problem**)

$$\min_{\mathcal{G}} \|\Omega * (\mathcal{Y} - \mathcal{X})\|_F^2, \quad s.t. \quad \mathcal{X} = \text{TN}(\mathcal{G}_1, \dots, \mathcal{G}_d).$$

- ▶ Combining convex optimization with decomposition approach
- ▶ Incorporating priori knowledge, such as auxiliary information from other data, smoothness, sparsity, non-negativity and etc.

Tensor Decomposition



Tucker Decomposition (HOSVD):

$$\mathcal{X} = \mathcal{G} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \cdots \times_M \mathbf{U}^{(M)},$$

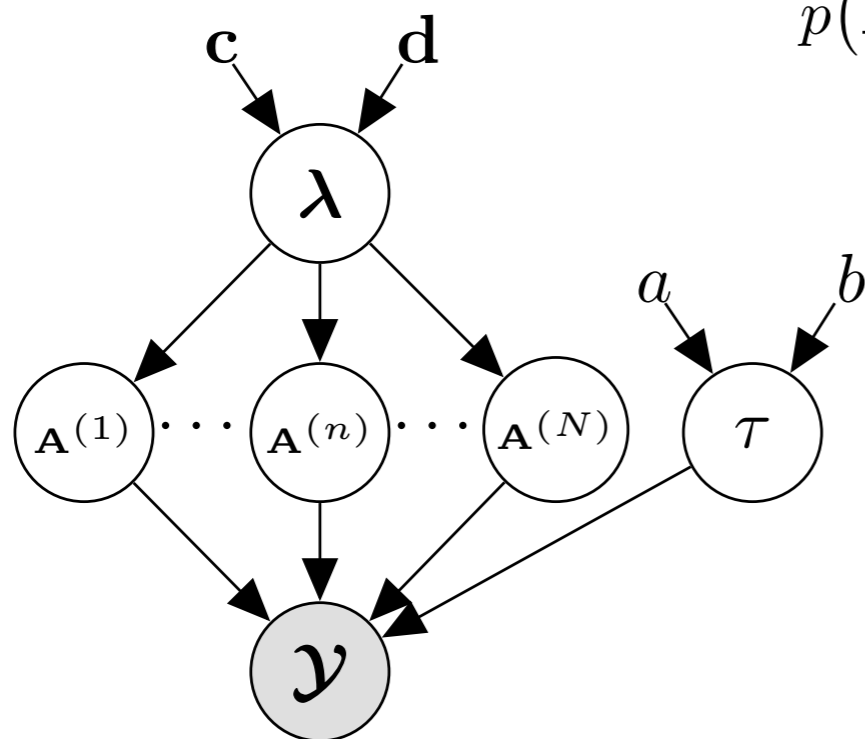
$$\text{rank}_{ML}(\underline{\mathbf{X}}) = \{\text{rank}(\mathbf{X}_{(1)}), \text{rank}(\mathbf{X}_{(2)}), \dots, \text{rank}(\mathbf{X}_{(N)})\},$$

Canonical Polyadic Decomposition (CPD) :

$$\mathcal{X} = \sum_{r=1}^R \lambda_r \mathbf{u}_r^{(1)} \circ \mathbf{u}_r^{(2)} \circ \cdots \circ \mathbf{u}_r^{(M)}, \quad R_1 \leq \text{rank}_{CP}(\underline{\mathbf{X}}) \leq R_2 R_3 \cdots R_N.$$

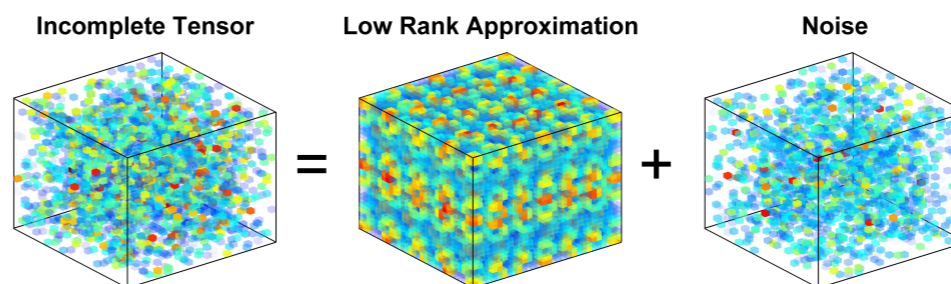
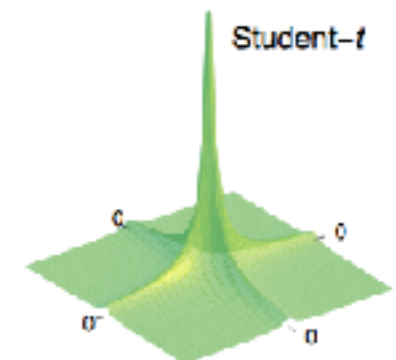
Learning Optimal Tensor Rank

- ▶ Probabilistic modeling of tensor decomposition
- ▶ **Group sparsity prior** imposed on factor matrices
- ▶ Bayesian inference for posteriors of model parameters



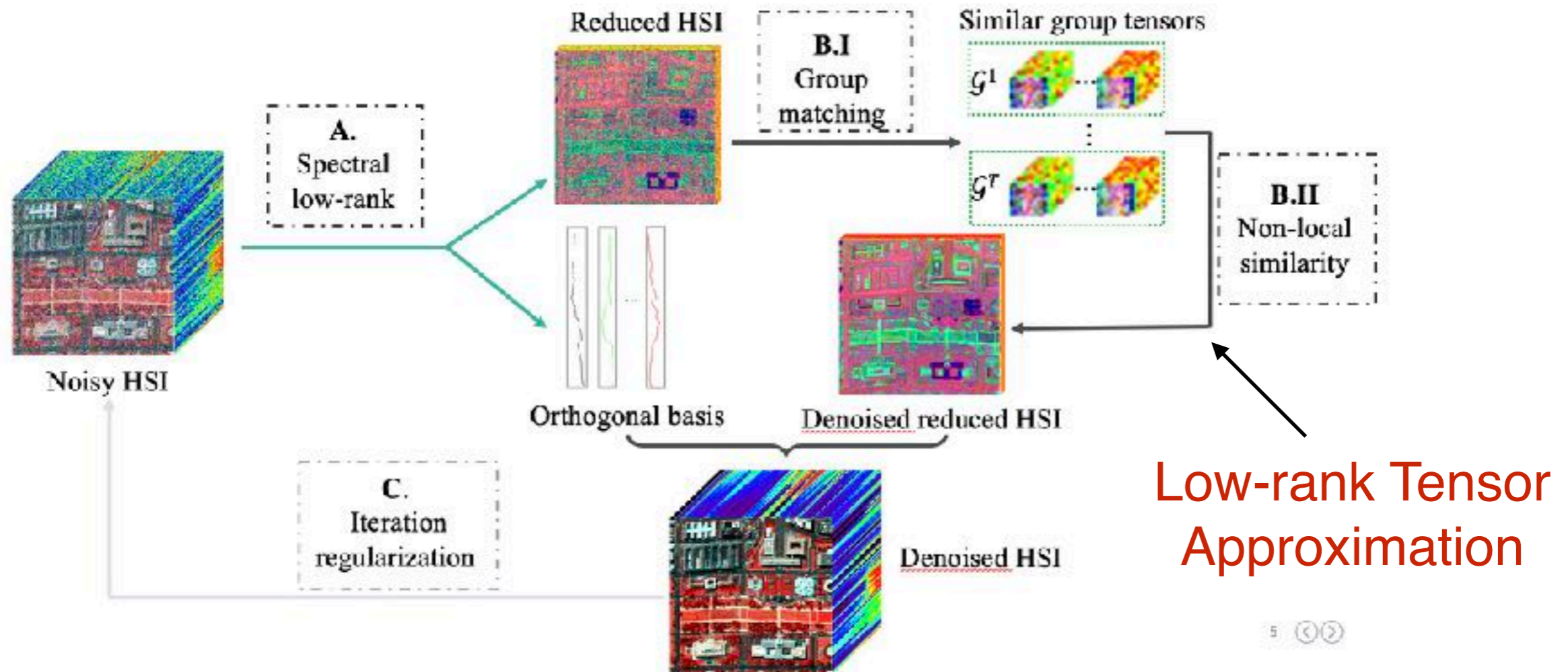
$$p(\mathbf{A}^{(n)} | \boldsymbol{\lambda}) = \prod_{i_n=1}^{I_n} \mathcal{N}(\mathbf{a}_{i_n}^{(n)} | \mathbf{0}, \boldsymbol{\Lambda}^{-1}), \forall n \in [1, N],$$

$$p(\boldsymbol{\lambda}) = \prod_{r=1}^R \text{Ga}(\lambda_r | c_0^r, d_0^r),$$



(Zhao et al, TPAMI 2015)

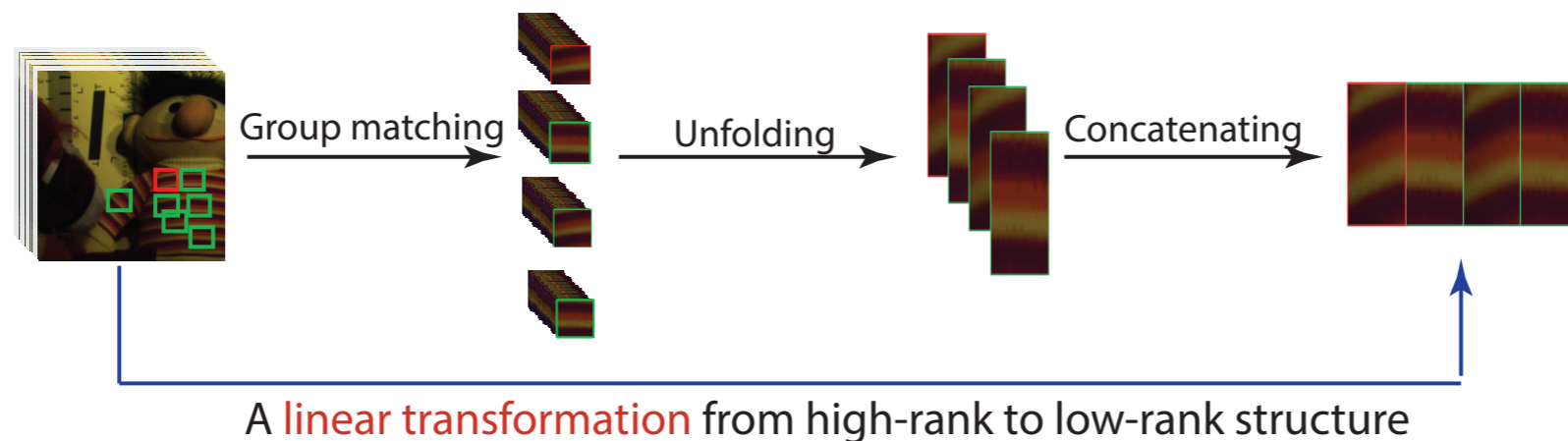
Applications to Hyperspectral Image Denoising



Non-Local Meets Global: An Integrated Paradigm for Hyperspectral Denoising (He et al., CVPR 2019)

Tensor Completion Under Multiple Transformation

- ▶ Image is not always globally low-rank
- ▶ Non-local similar patches are often low-rank
- ▶ Lack of theoretical analysis



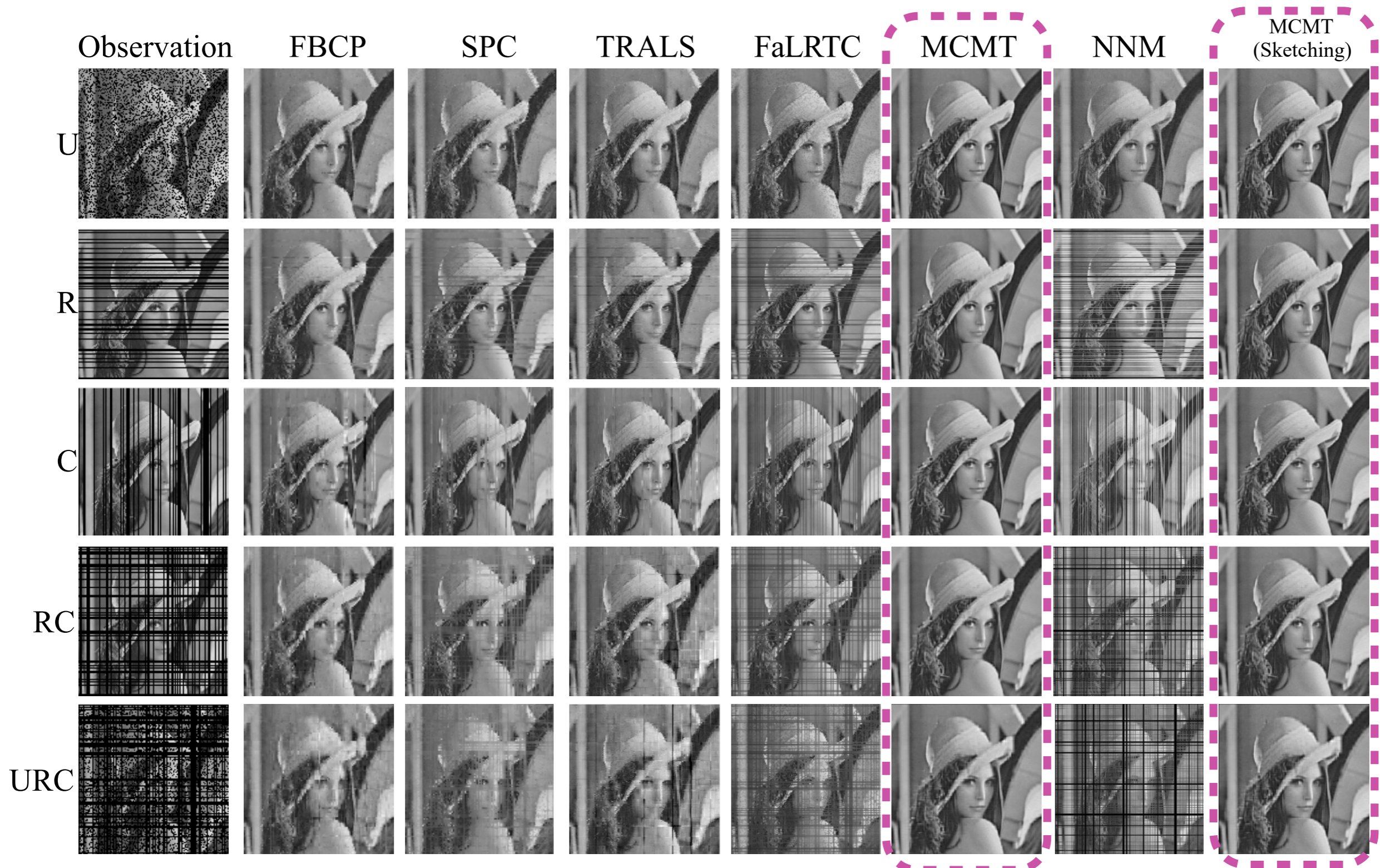
Low-rankness under linear transformation

$$\min_{\mathbf{X} \in \mathbb{R}^{m_1 \times m_2}} \|\mathcal{Q}(\mathbf{X})\|_* \quad s.t. \quad \|\mathcal{P}_\Omega(\mathbf{X}) - \mathcal{P}_\Omega(\mathbf{Y})\|_F \leq \delta,$$

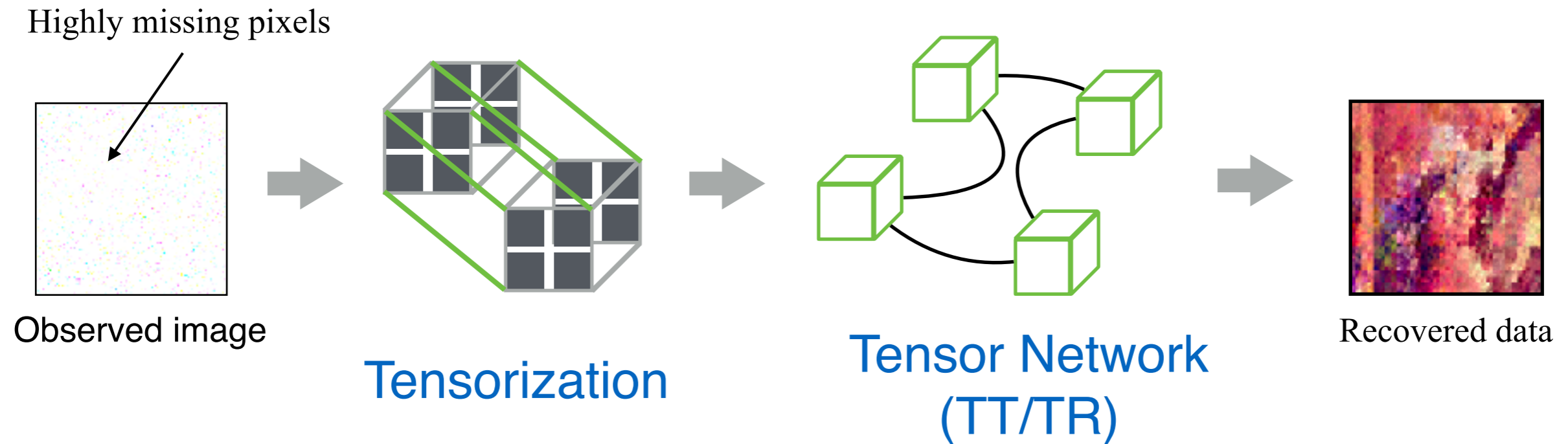
↓
Linear transformation

Guaranteed Matrix Completion under Multiple Linear Transformations (Li et al, CVPR 2019)

Illustrative Experiment



Tensor Ring with Low-rank Cores



$$\min_{\mathcal{G}} \left\| \Omega * (\mathcal{Y} - \hat{\mathcal{Y}}) \right\|_F^2 + \lambda \sum_{n=1}^d \sum_{i=1}^3 \left\| \mathcal{G}_{(i)}^{(n)} \right\|_*, \quad s.t. \quad \hat{\mathcal{Y}} = \text{TR}(\mathcal{G}^{(1)}, \dots, \mathcal{G}^{(d)}).$$

Fitting error
Nuclear norm on core tensor
TT/TR decomposition

- ▶ Combined decomposition and nuclear norm minimization
- ▶ Tensorization + TN improves performance of TD on original tensor

L. Yuan et al. , “Tensor ring decomposition with rank minimization on latent space: An efficient approach for tensor completion ”, (AAAI 2019)

Tensor Network Diagram



$$(1 \times 1)$$

Scalar



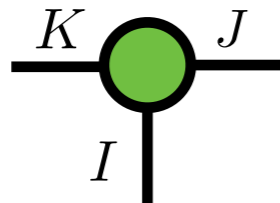
$$(I \times 1)$$

Vector



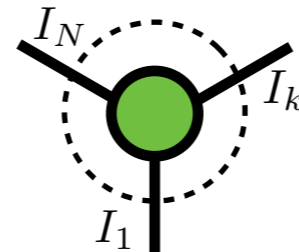
$$(I \times J)$$

Matrix



$$(I \times J \times K)$$

3-order Tensor



$$(I_1 \times I_2 \times \cdots \times I_N)$$

N-order Tensor

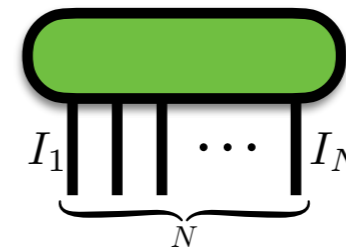
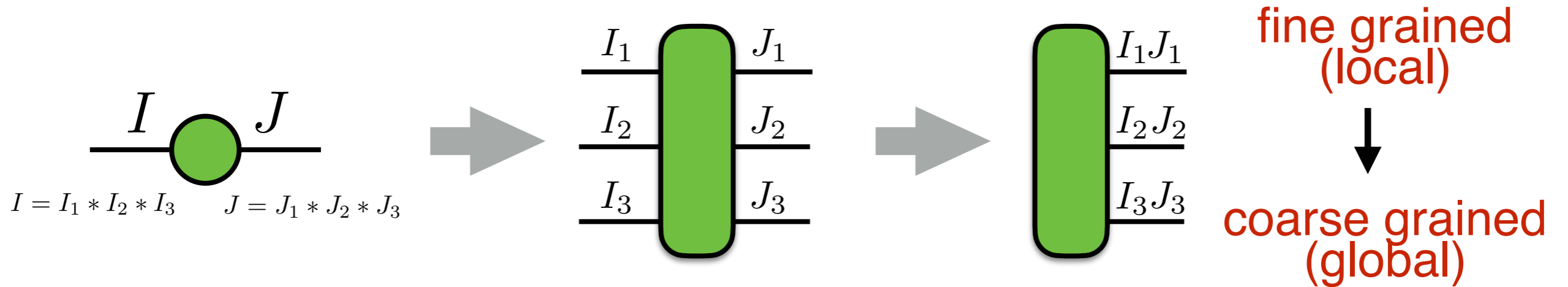
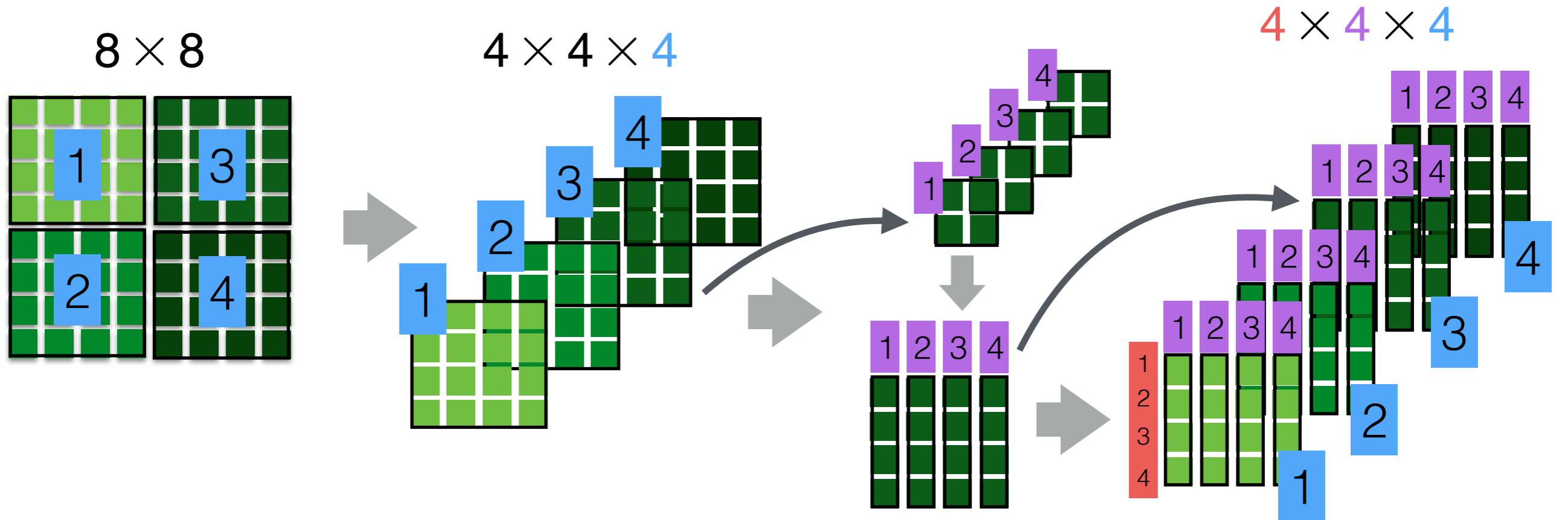


Image Tensorization



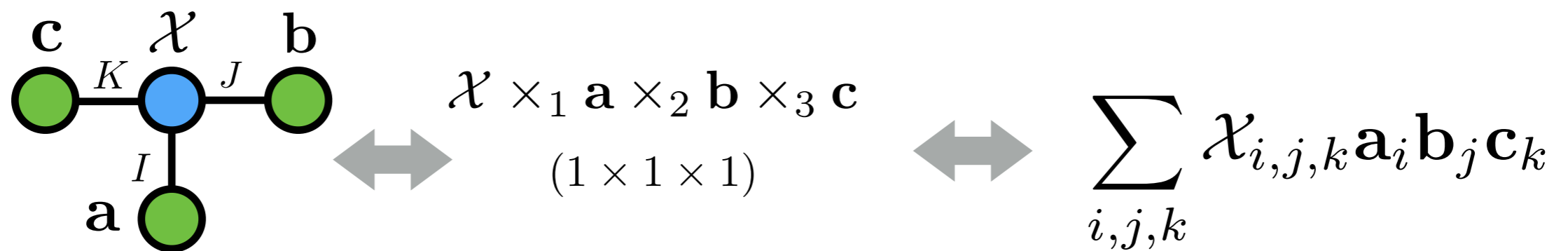
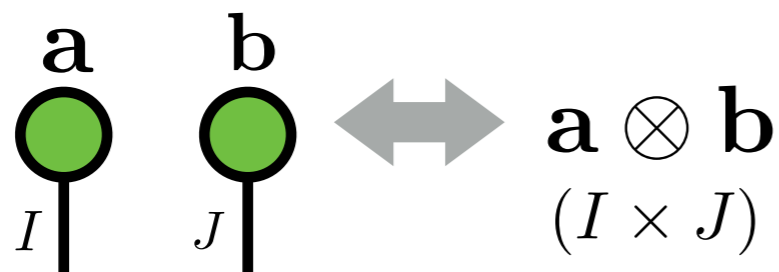
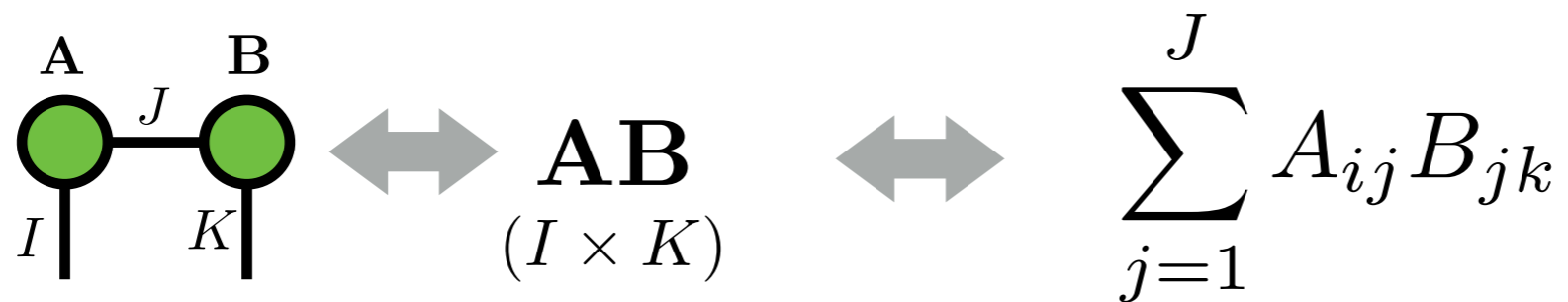
- Explore correlations from global to local

What Is Tensor Network?

- ▶ A powerful tool to describe strongly entangled quantum **many-body systems** in physics
- ▶ Approximation of **N-order tensor** as contractions of $O(N)$ **smaller tensors**



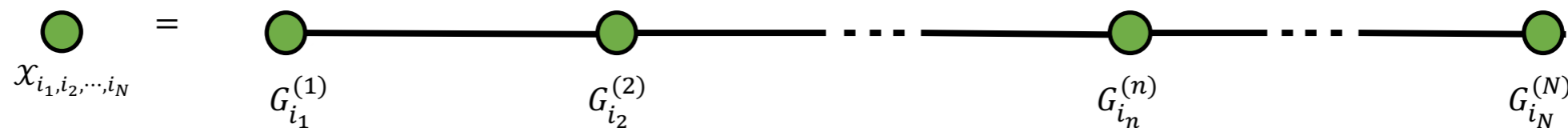
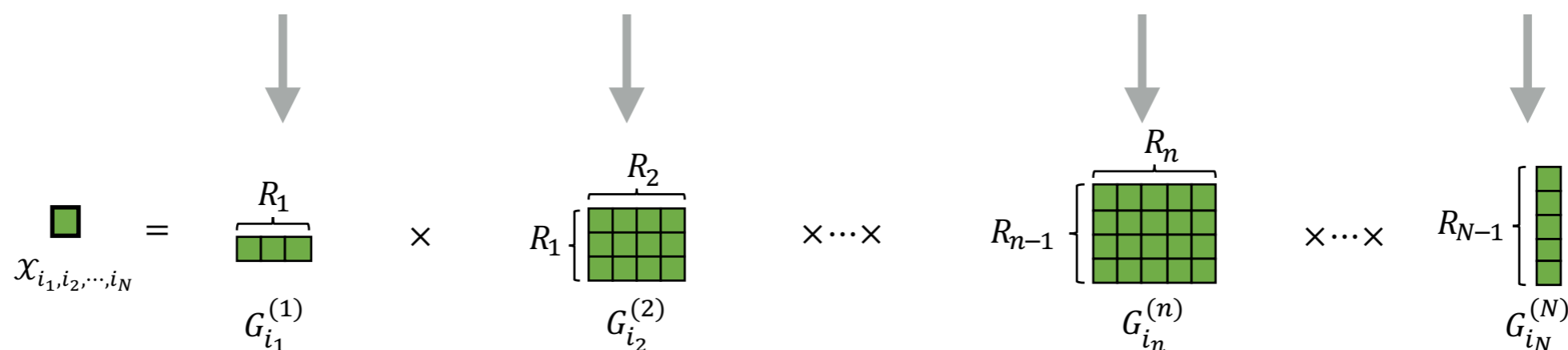
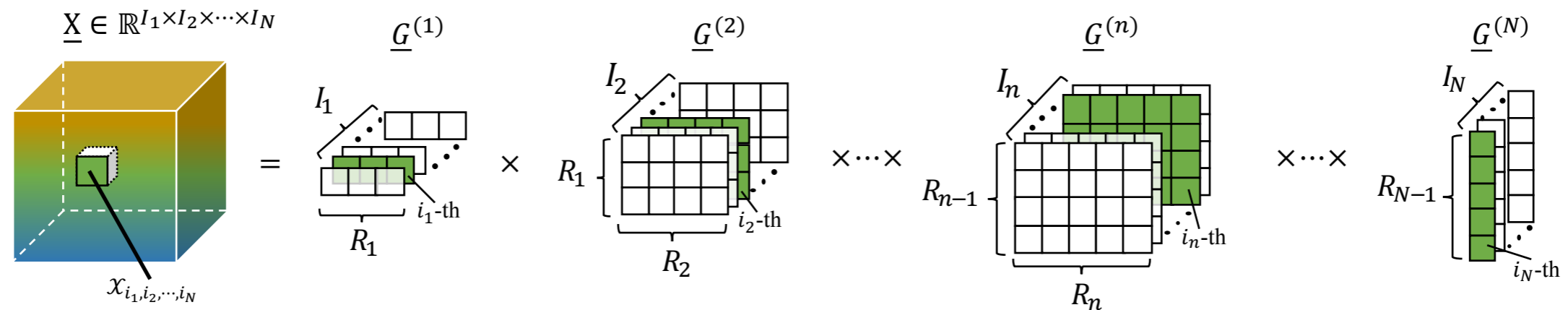
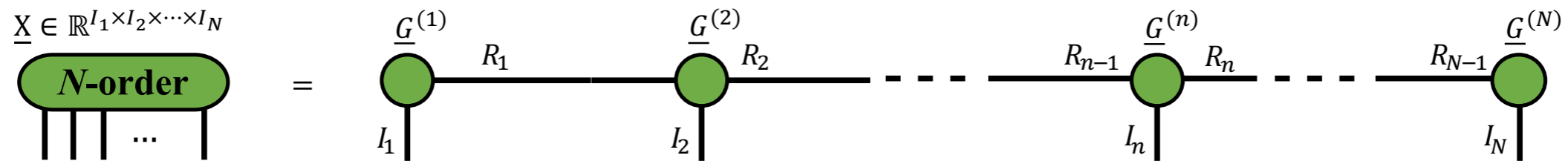
Tensor Network Operations



- Convenient to denote complex operations

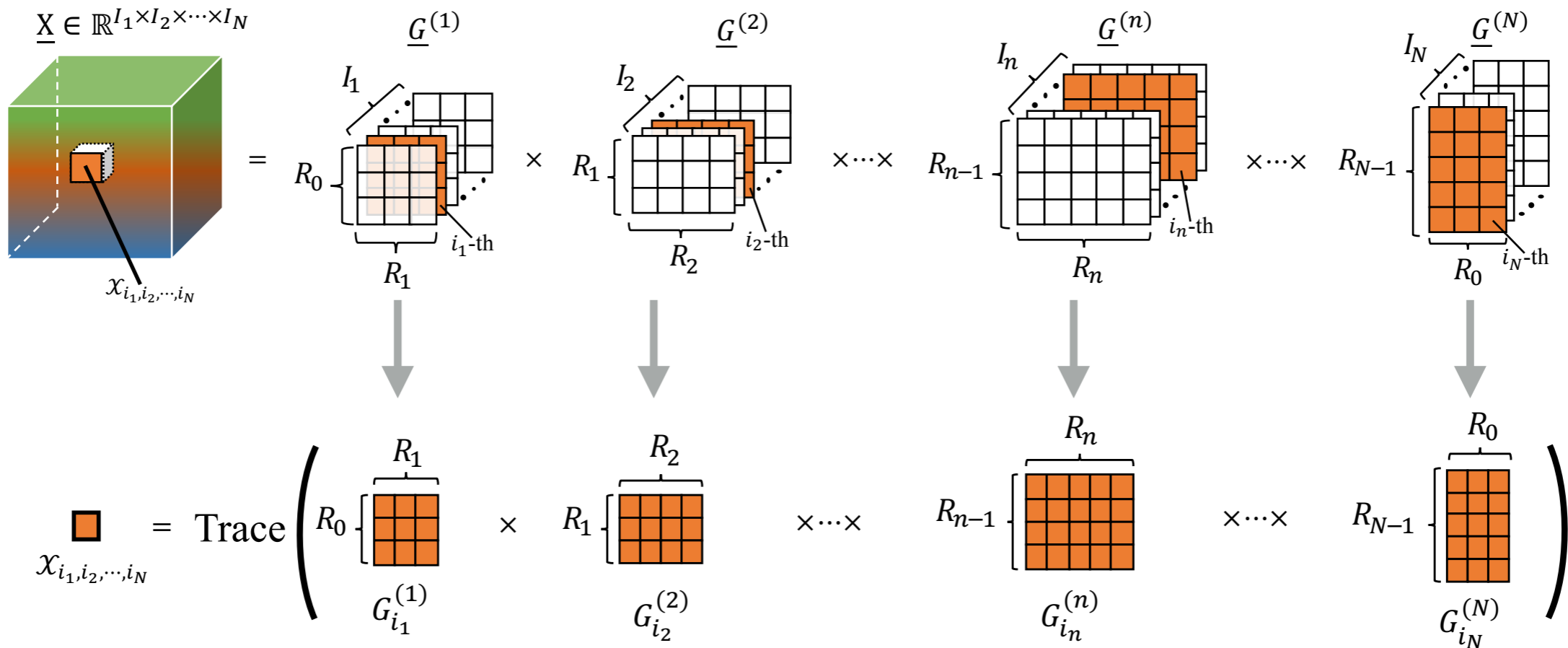
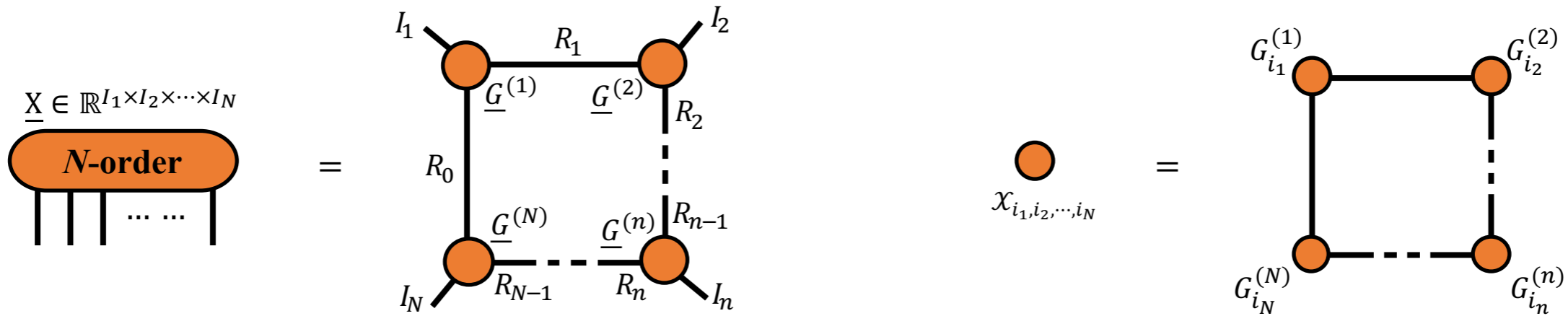
Tensor Train Decomposition

(Oseledets, SIAM J. Sci. Comput. 2011)

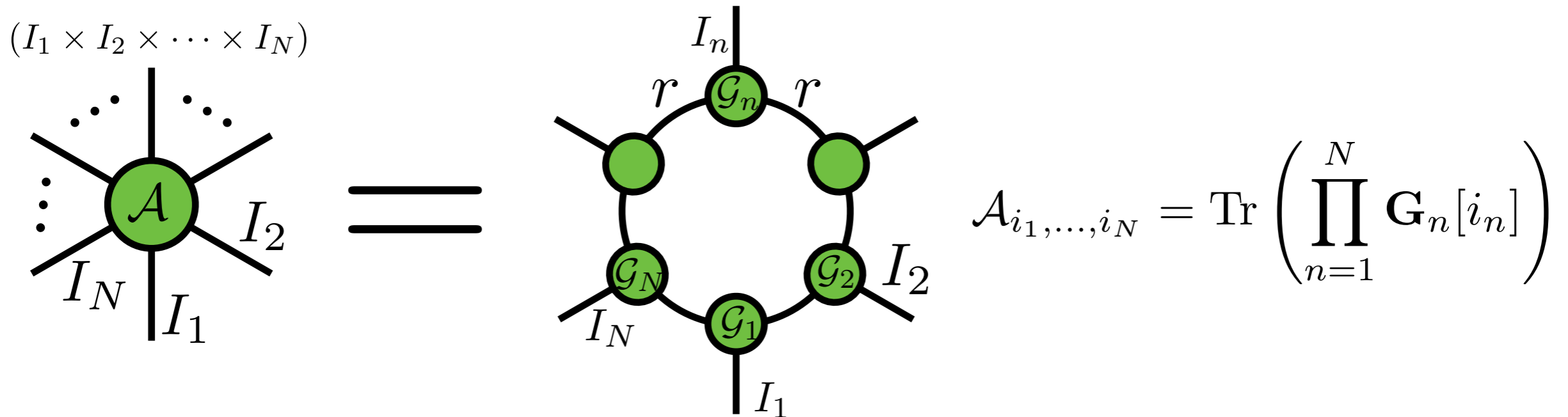


Tensor Ring Decomposition

(Zhao et al., arXiv 2016, ICASSP 2019)



Tensor Ring Decomposition



- ▶ Highly compact representation

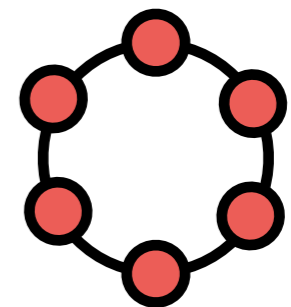
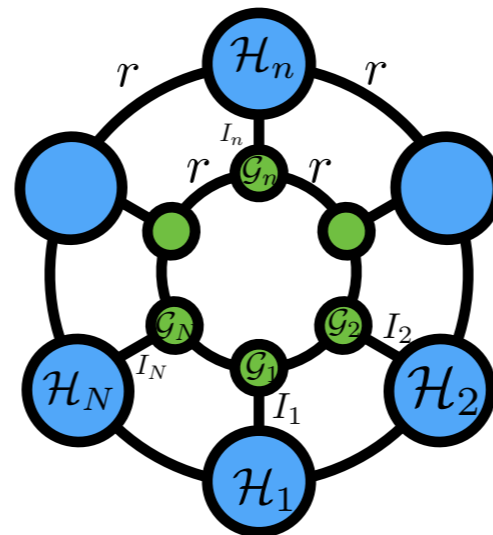
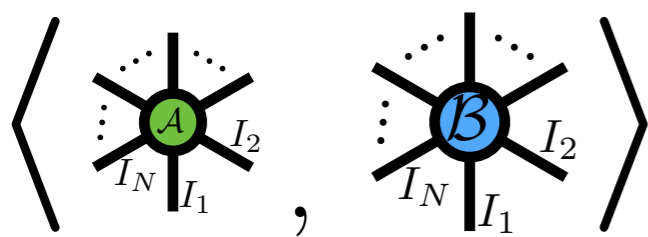
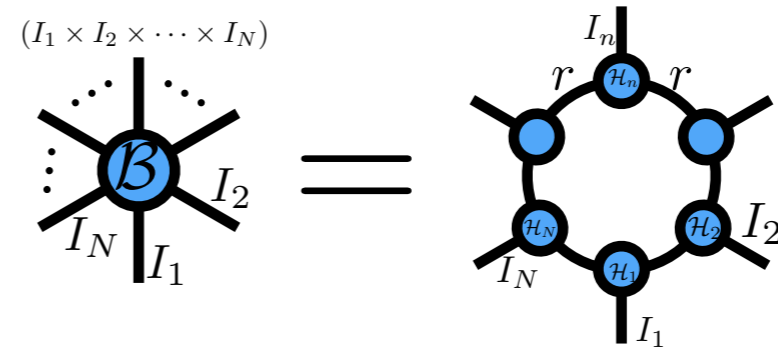
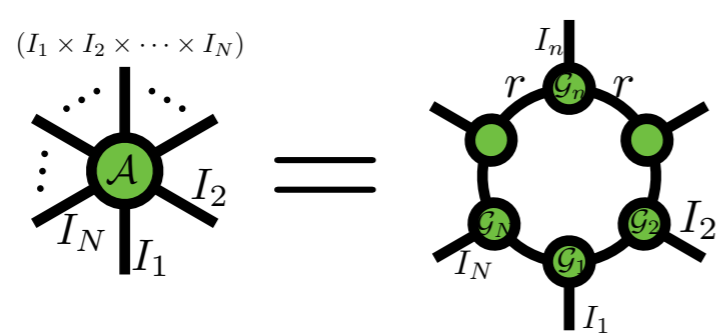
$$\mathcal{O}(I^N) \Rightarrow \mathcal{O}(N I r^2)$$

- ▶ Circular permutation invariance

$$\tilde{\mathcal{A}}(I_n \times \dots \times I_N \times I_1 \times \dots \times I_{n-1}) \Rightarrow \{\mathcal{G}_n, \dots, \mathcal{G}_N, \mathcal{G}_1, \dots, \mathcal{G}_{n-1}\}$$

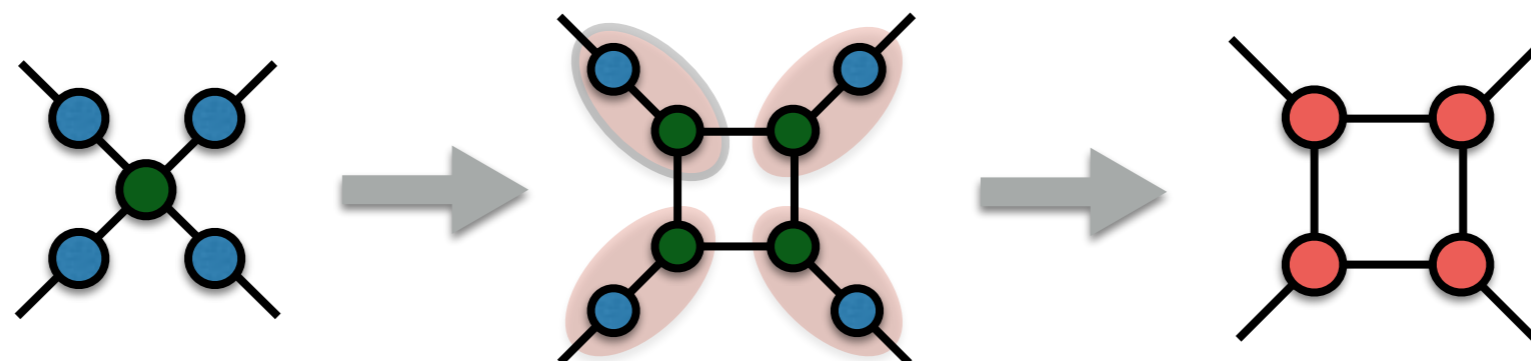
- ▶ If any $r = 1$, TR becomes TT.

Efficient Computation



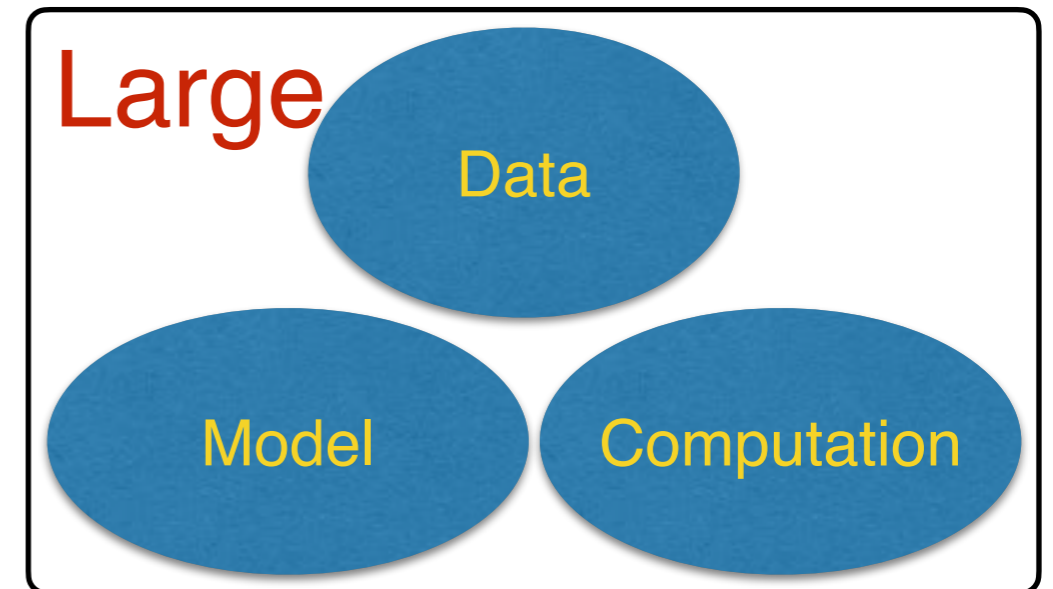
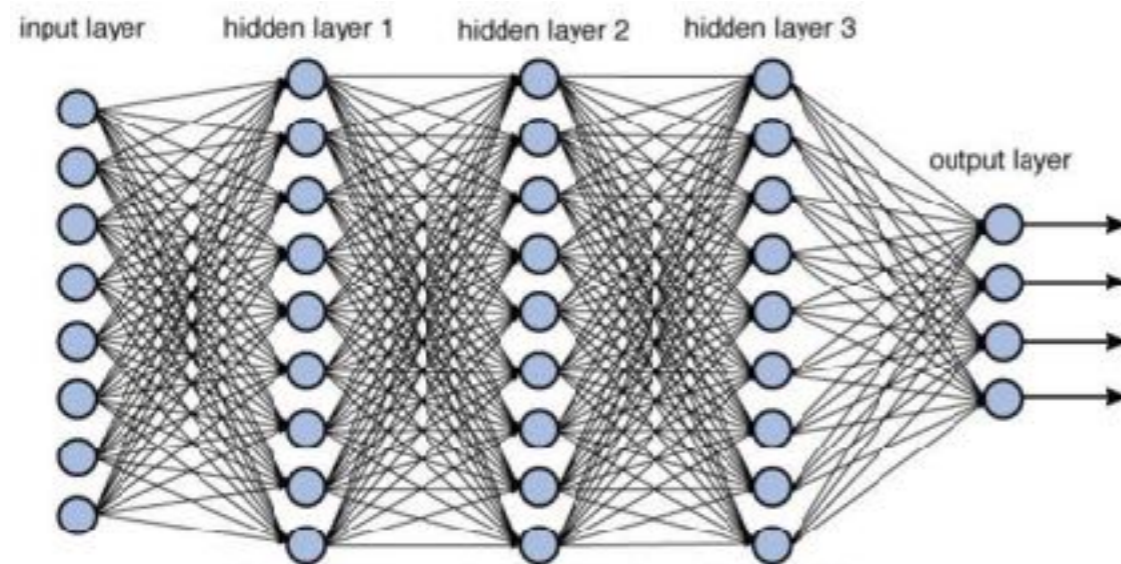
Computation
 $\mathcal{O}(I^N)$

Computation
 $\mathcal{O}(N I r^4)$



Tensor Networks in Deep Learning

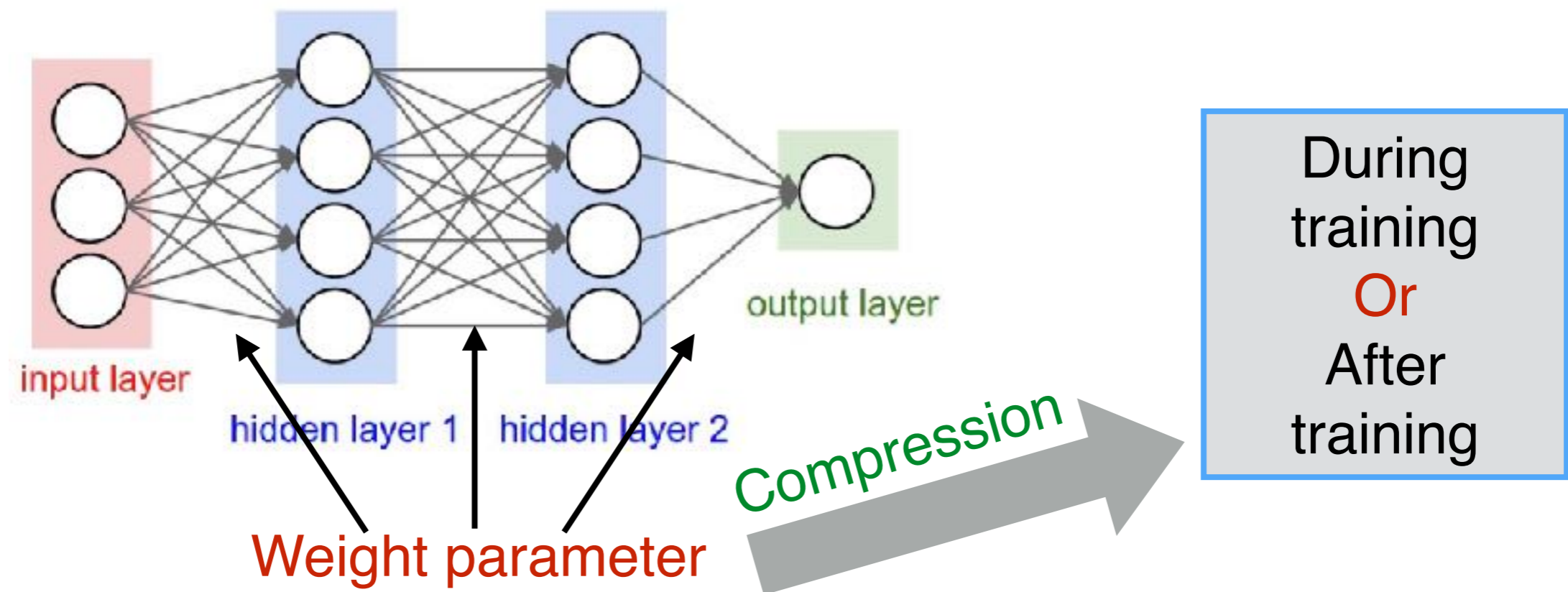
Challenges in DNNs



- ▶ **Dimension of input** increases for more complex data and task
- ▶ Wider layer improves model capacity but **parameters increase** dramatically
- ▶ Large number of **training samples**, large number of **model parameters**, large **computation cost**
- ▶ Huge storage (memory, disk requirement)
- ▶ Slow inference and energy costly
- ▶ Hard to deploy models on small devices (e.g. smart phone)

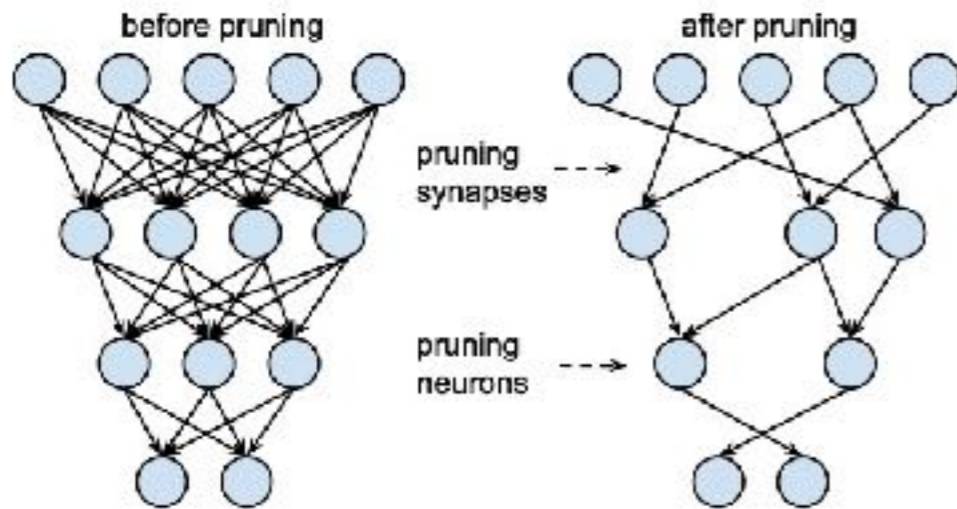
Model Compression

Goal: Make a lightweight model that is fast, memory-efficient and energy-efficient

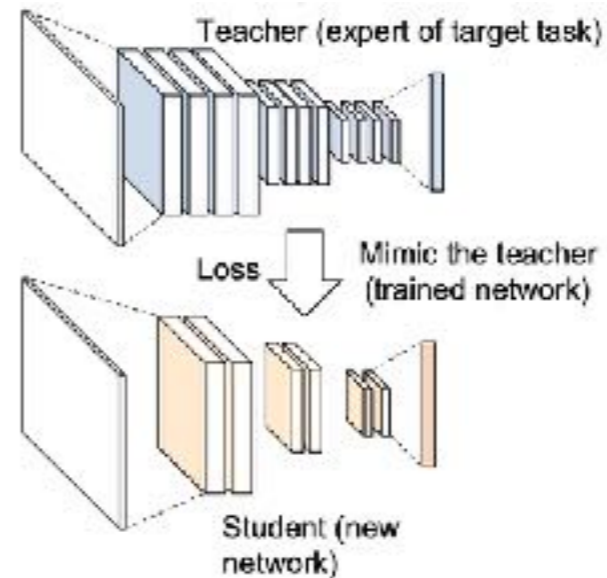


- ▶ To **compress** model parameters with comparable performance
- ▶ To **increase capacity** without increasing model parameters
- ▶ To improve **computation efficiency**

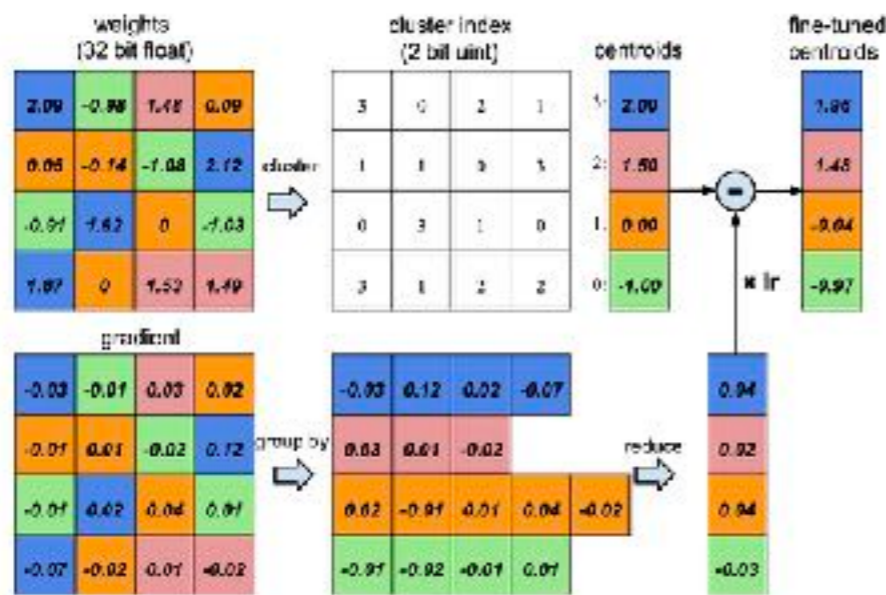
Taxonomy of Model Compression



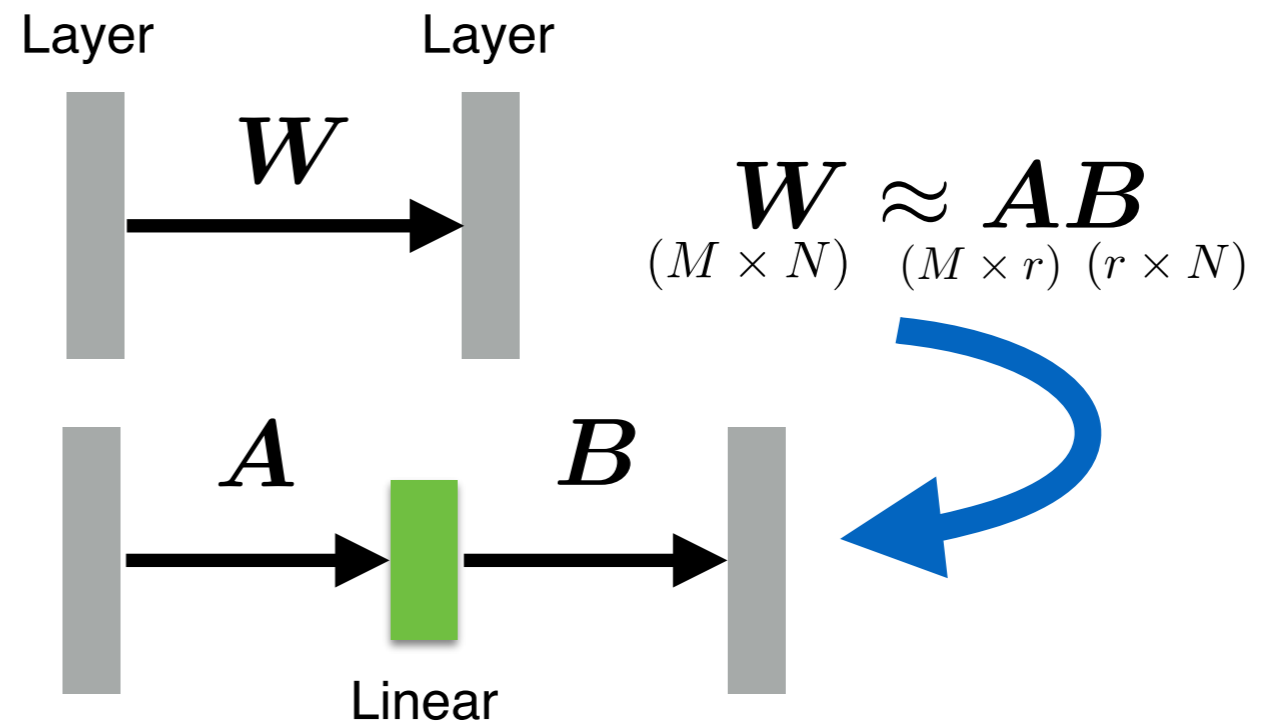
Weight Pruning
(Han et al., NIPS 2015)



Knowledge Distillation
(Hinton et al., NIPS 2014 Workshop)



Quantization and Sharing of Weight
(Han et al., ICLR 2016)



Low-rank Matrix/Tensor Factorization
(Novikov et al., NIPS 2015)

Compression

Low-rank **matrix factorization of weights** in fully connected layer

$$\underset{(M \times N)}{\mathbf{W}} \approx \underset{(M \times r)}{\mathbf{A}} \underset{(r \times N)}{\mathbf{B}} \quad \text{Compression: } \mathcal{O}(MN) \rightarrow \mathcal{O}(r(M + N))$$

Low-rank **tensor network factorization of weights**

- ▶ Step 1: $\mathbf{W} \rightarrow \mathcal{W}$ (**Matrix to d -order tensor**)
- ▶ Step 2: $\mathcal{W} \approx \text{TT}(\mathcal{G}_1 \cdots \mathcal{G}_d)$ (**Tensor network representation**)

Loss function:

$$\mathcal{L}(\mathbf{W}, \mathbf{x}, \mathbf{y}) \rightarrow \mathcal{L}(\{\mathcal{G}_1, \cdots, \mathcal{G}_d\}, \mathbf{x}, \mathbf{y})$$

Compression:

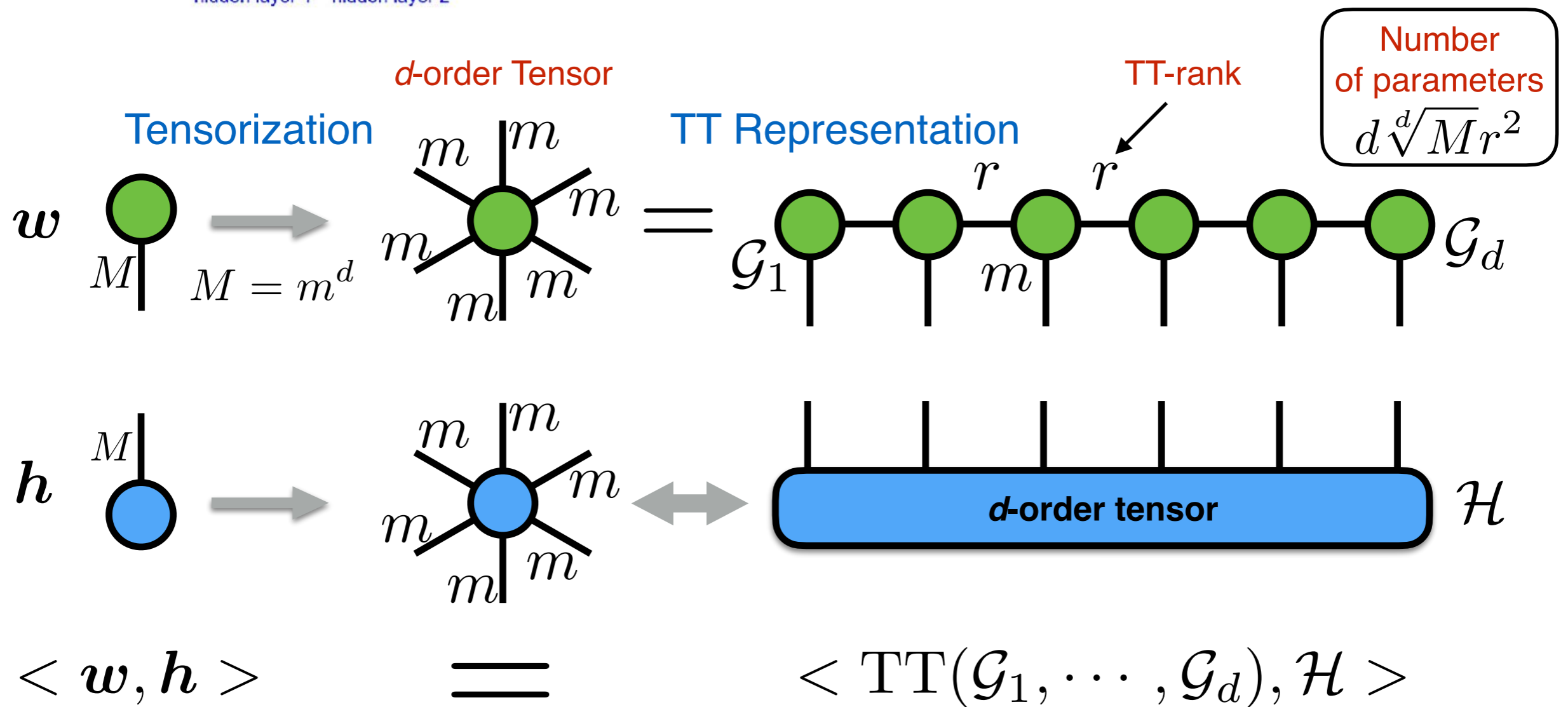
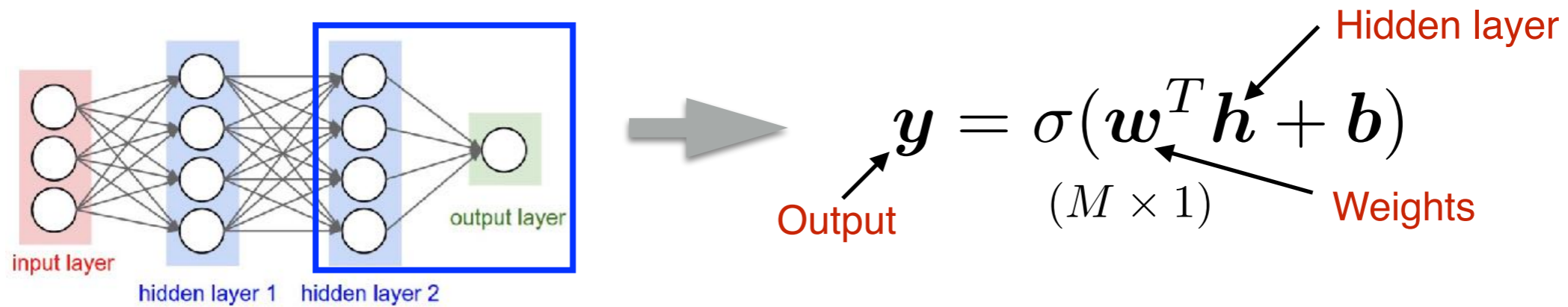
$$\mathcal{O}(MN) \rightarrow \mathcal{O}(dr^2 \sqrt[d]{M} \sqrt[d]{N})$$

Order of tensor TT-rank

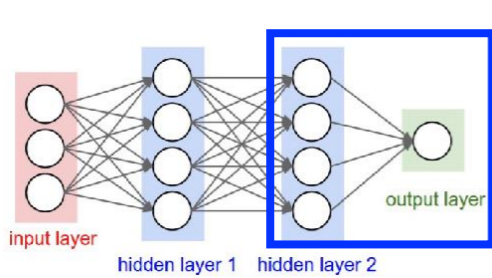
Train very “wide” model

Tensorizing Neural Networks (Novikov et al., NIPS 2015)

Vector to Scalar Transformation via TN

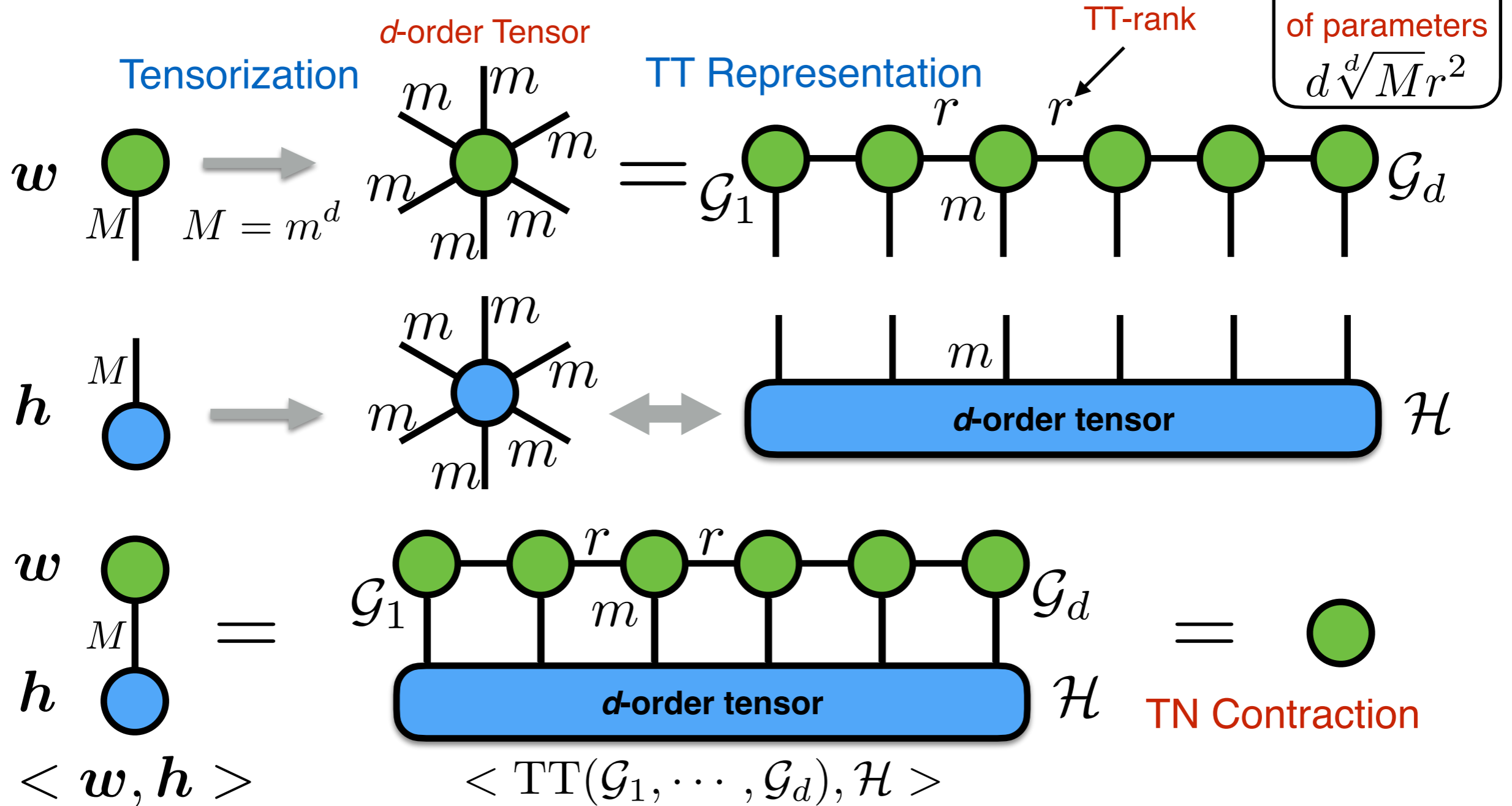


Vector to Scalar Transformation via TN

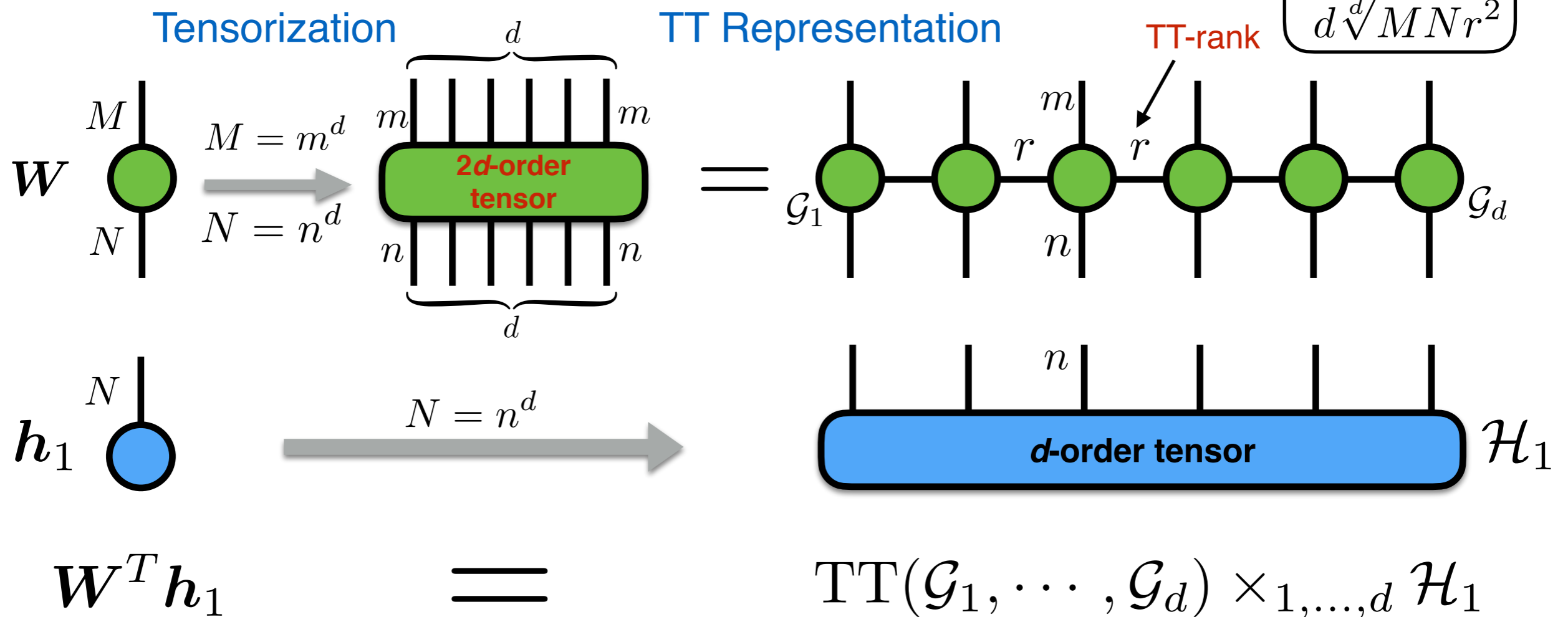
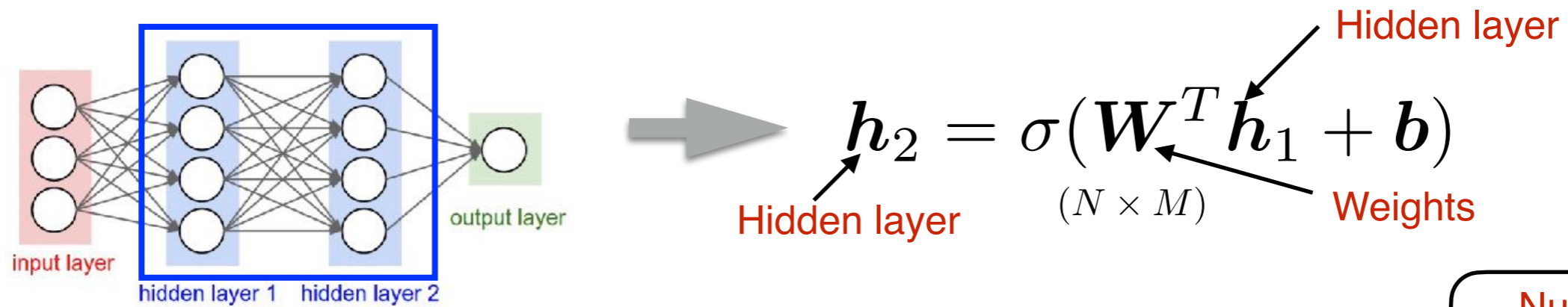


$$y = \sigma(\mathbf{w}^T \mathbf{h} + b)$$

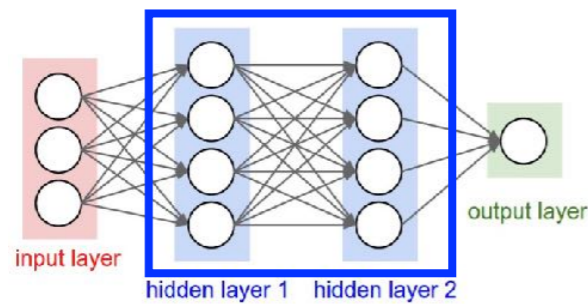
Labels: **Output** (pointing to y), **Weights** (pointing to \mathbf{w}), **Hidden layer** (pointing to \mathbf{h}), $(M \times 1)$ (under \mathbf{w}).



Vector to Vector Transformation via TN



Vector to Vector Transformation via TN



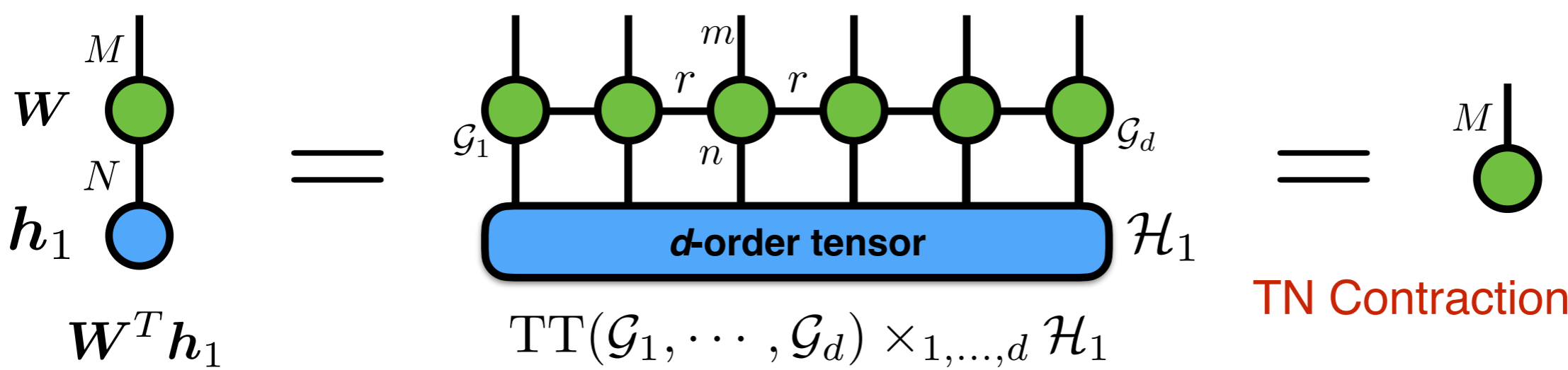
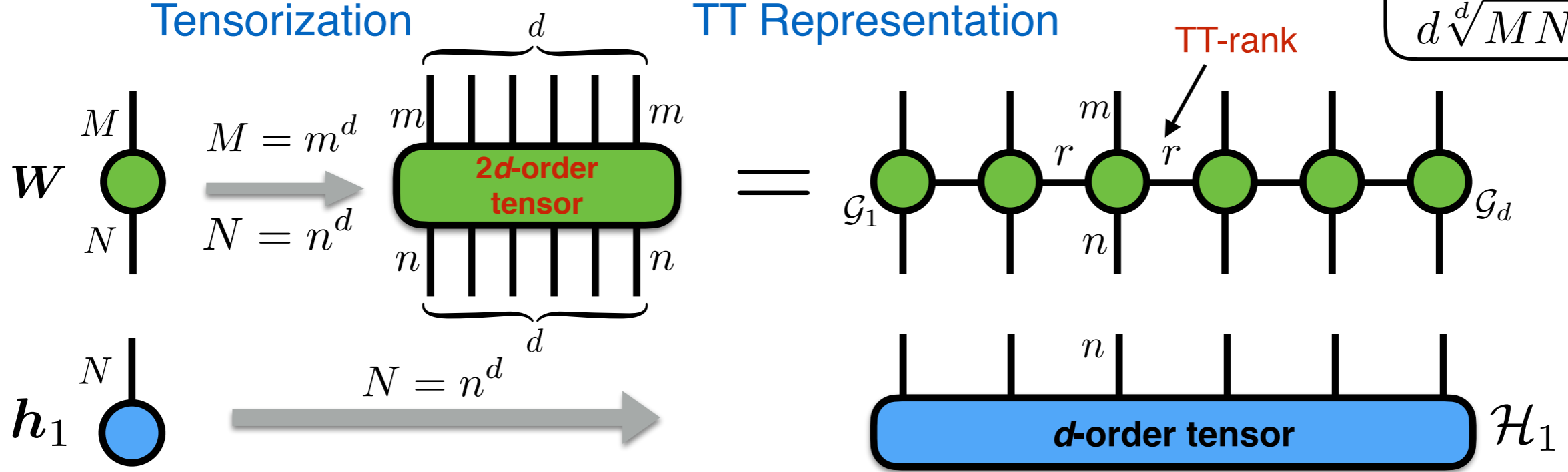
$$h_2 = \sigma(\mathbf{W}^T h_1 + b)$$

Labels: **Hidden layer** (pointing to h_1), **Weights** (pointing to \mathbf{W}), **Hidden layer** (pointing to h_2). Dimension: $(N \times M)$.

Number of parameters $d \sqrt[3]{MNr^2}$

Tensorization

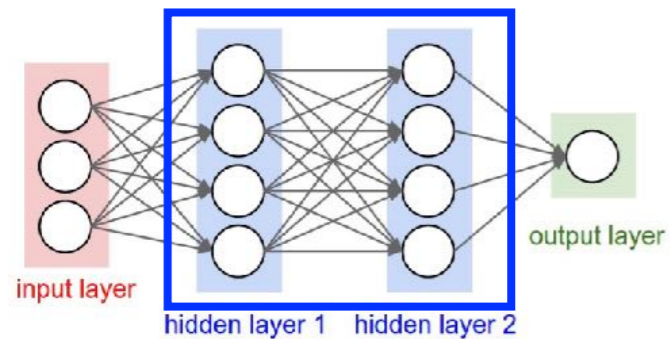
TT Representation



TT Contraction

$$\text{TT}(\mathcal{G}_1, \dots, \mathcal{G}_d) \times_{1, \dots, d} \mathcal{H}_1$$

Learning of TN Weights

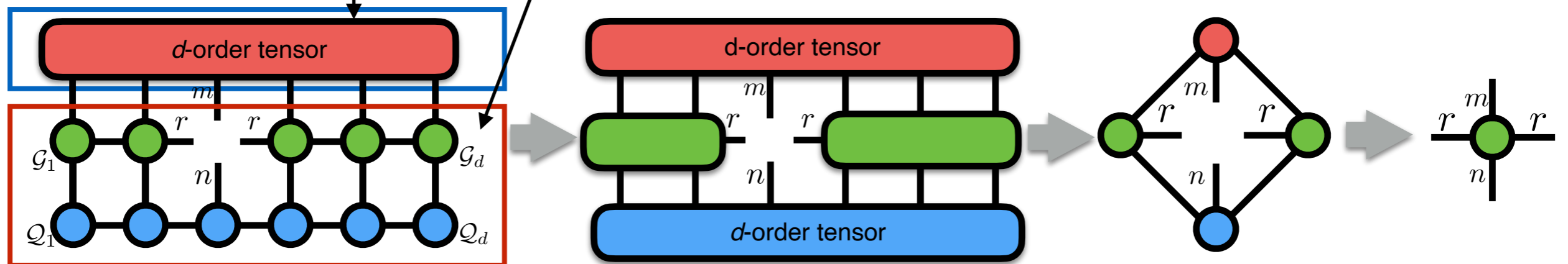
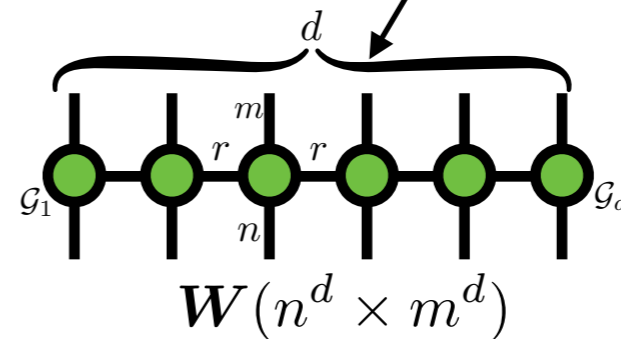


$$h_2 = \sigma(\mathbf{W}^T h_1 + \mathbf{b})$$

\mathbf{o}_2

- ▶ Loss: $\min_{\mathbf{W}} \mathcal{L}(\mathbf{W}, \mathbf{x}, y) \Rightarrow \min_{\mathcal{G}_1, \dots, \mathcal{G}_d} \mathcal{L}(\{\mathcal{G}_1, \dots, \mathcal{G}_d\}, \mathbf{x}, y)$
- ▶ Gradients over TN core tensors

$$\frac{\partial \mathcal{L}}{\partial \mathcal{G}_k} = \frac{\partial \mathcal{L}}{\partial h_2} \frac{\partial h_2}{\partial \mathbf{o}_2} \frac{\partial \mathbf{o}_2}{\partial \mathcal{G}_k}$$



Compression Performance

► Model compression on ImageNet

Architecture	Matrices compr.	Network compr.	Error
FC FC FC	1	1	11.2
TT4 FC FC	50 972	3.9	11.2
TT2 FC FC	194 622	3.9	11.5
TT1 FC FC	713 614	3.9	12.8
TT4 TT4 FC	37 732	7.4	12.3
LR1 FC FC	3 521	3.9	97.6
LR5 FC FC	704	3.9	53.9
LR50 FC FC	70	3.7	14.9

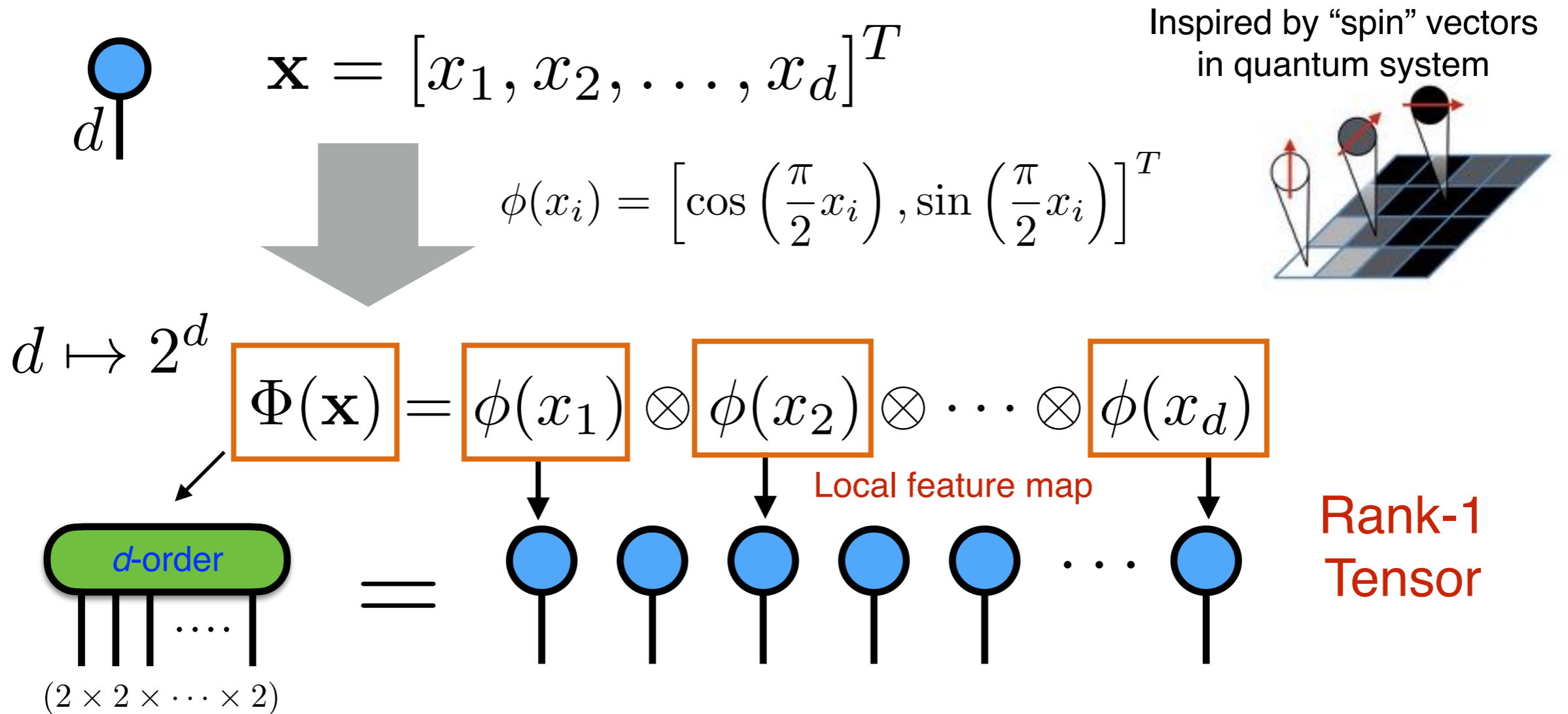
► Speed

Operation	Time	Memory	Type	1 im. time (ms)	100 im. time (ms)
FC forward pass	$O(MN)$	$O(MN)$	CPU fully-connected layer	16.09	97.2
TT forward pass	$O(dr^2 m \max\{M, N\})$	$O(r \max\{M, N\})$	CPU TT-layer	1.24	94.7
FC backward pass	$O(MN)$	$O(MN)$	GPU fully-connected layer	2.7	33
TT backward pass	$O(d^2 r^4 m \max\{M, N\})$	$O(r^3 \max\{M, N\})$	GPU TT-layer	1.92	12.86

Tensorizing Neural Networks (Novikov et al., NIPS 2015)

Feature Mapping

- ▶ Mapping input data into TN representation



- ▶ Different mapping for each feature is possible
- ▶ $\phi(x_i)$ can be generalized to I -dimensional, e.g., $\Phi(\mathbf{x}) : d \mapsto I^d$

Exponential Machines

- ▶ Mapping of input data

$$\mathbf{x} = [x_1, x_2, \dots, x_d]^T$$

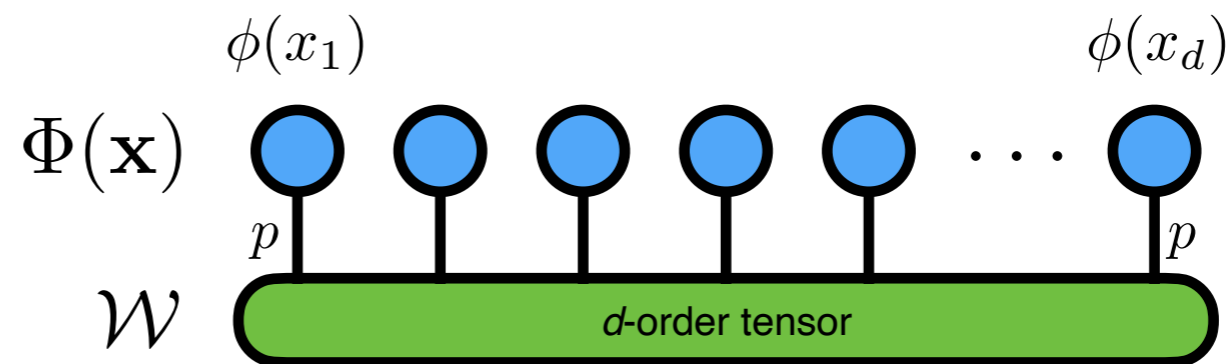
$$\phi(x_i) = [1, x_i, x_i^2, \dots, x_i^{p-1}]^T$$



$$d \mapsto p^d \quad \Phi(\mathbf{x}) = \phi(x_1) \otimes \phi(x_2) \otimes \dots \otimes \phi(x_d)$$

- ▶ p -degree polynomial function

$$f(\mathbf{x}) = \langle \mathcal{W}, \Phi(\mathbf{x}) \rangle$$



Example:

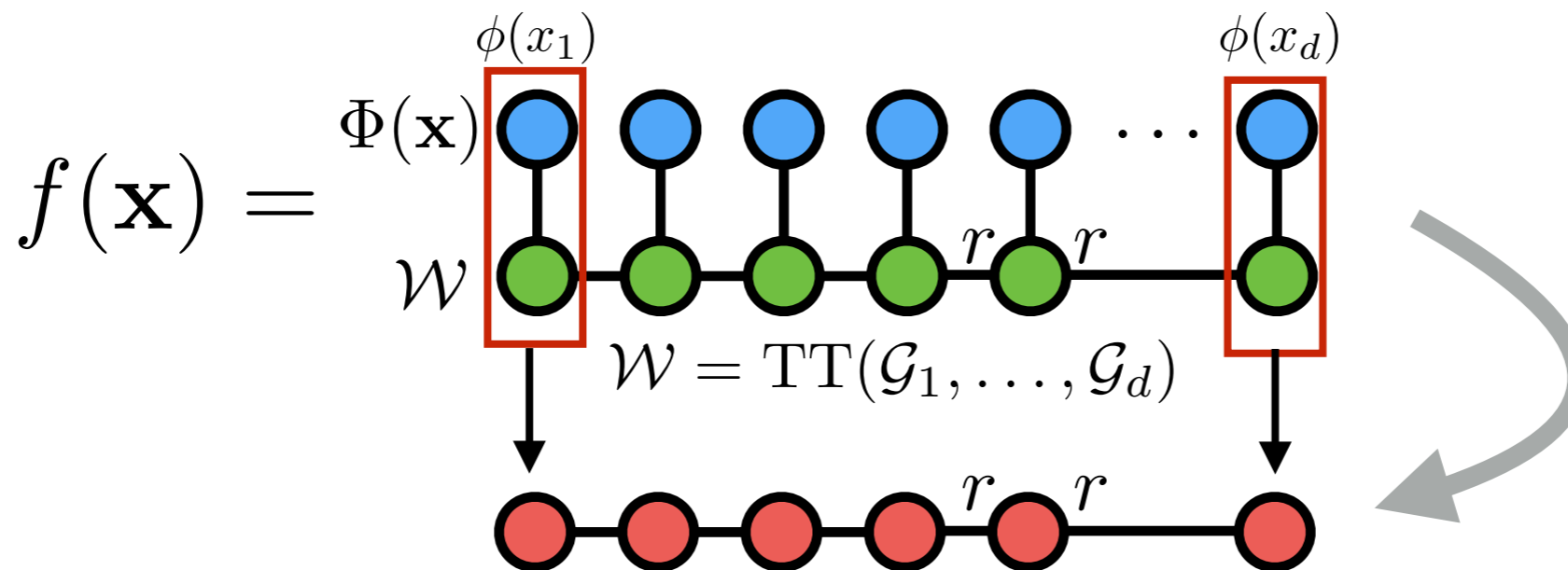
$$d = 3, p = 2 \quad f(\mathbf{x}) = \mathcal{W}_{000} + \mathcal{W}_{100}x_1 + \mathcal{W}_{010}x_2 + \mathcal{W}_{001}x_3 \\ + \mathcal{W}_{110}x_1x_2 + \mathcal{W}_{101}x_1x_3 + \mathcal{W}_{011}x_2x_3 \\ + \mathcal{W}_{111}x_1x_2x_3$$

- ▶ Interactions of features are fully captured
- ▶ \mathcal{W} contains weights of interactions and increases exponentially with the dimension of input

Exponential Machines

$$f(\mathbf{x}) = \sum_{i_1=1}^p \cdots \sum_{i_d=1}^p \mathcal{W}_{i_1, \dots, i_d} \prod_{k=1}^d x_k^{i_k-1}$$

	Memory	Computation
Normal	p^d	$\mathcal{O}(dp^d)$
With TN	dpr^2	$\mathcal{O}(dpr^2)$



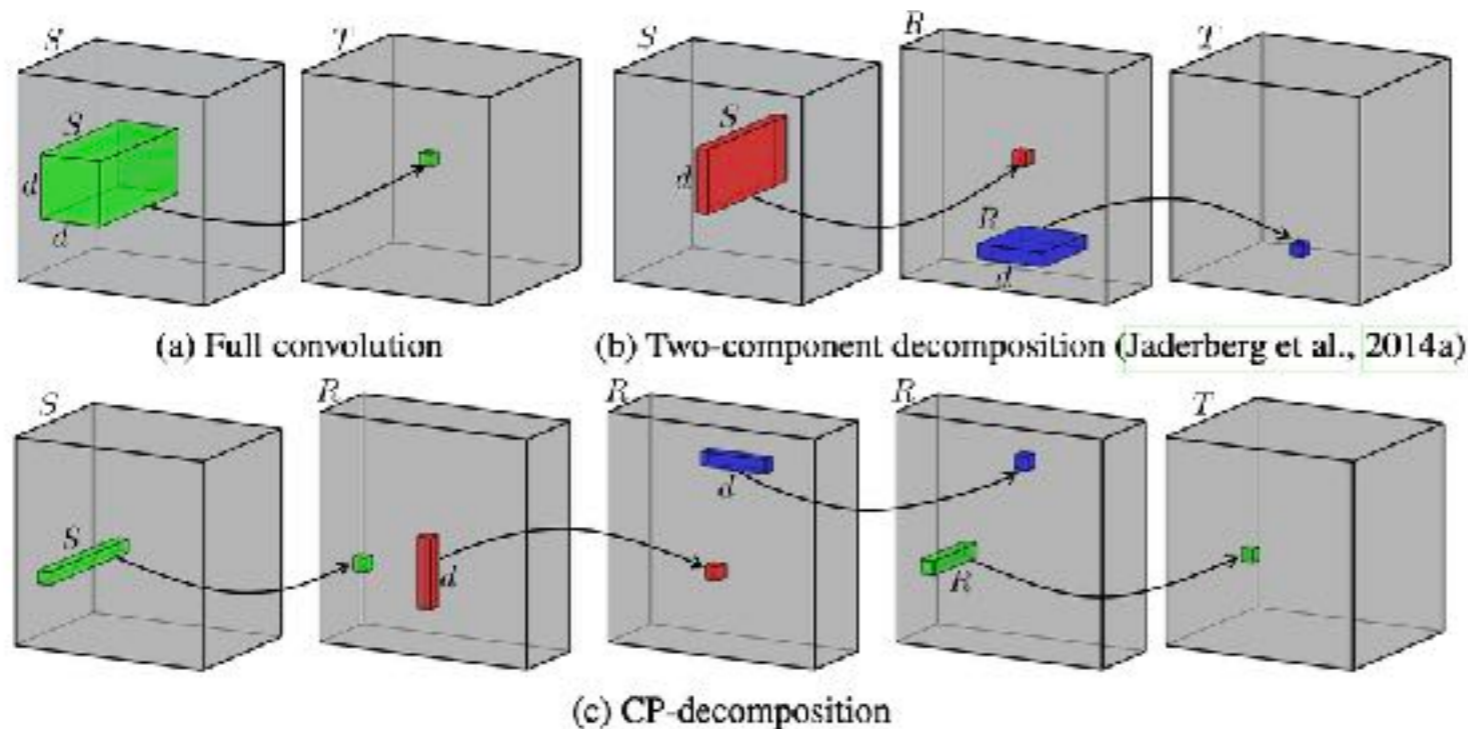
$$f(\mathbf{x}) = \left(\sum_{i_1=1}^p x_i^{i_1-1} \mathcal{G}_1[i_1] \right) \cdots \left(\sum_{i_d=1}^p x_d^{i_d-1} \mathcal{G}_d[i_d] \right) = \tilde{\mathbf{G}}_1 \tilde{\mathbf{G}}_2 \cdots \tilde{\mathbf{G}}_d$$

- ▶ TN representation is highly efficient for **very large** d, p
- ▶ Low-rank structure of TN (i.e., TT-rank) controls expressive power

Speed-up and Compression of CNNs

- ▶ 4-order kernel tensor is represented by CPD

$$K(i, j, s, t) = \sum_{r=1}^R K^x(i - x + \delta, r) K^y(j - y + \delta, r) K^s(s, r) K^t(t, r).$$



ImageNet (AlexNet)
4x speedup
1% accuracy drop

Kernel Size

Input channels

Number of Parameters
& Computations

$(d \times d \times S \times T)$
 Spatial (points to $d \times d$)
 Output channels (points to $S \times T$)

$STd^2 \rightarrow R(S + 2d + T)$

Speeding-up convolutional neural networks using fine-tuned cp-decomposition (Lebedev et al. ICLR 2015)

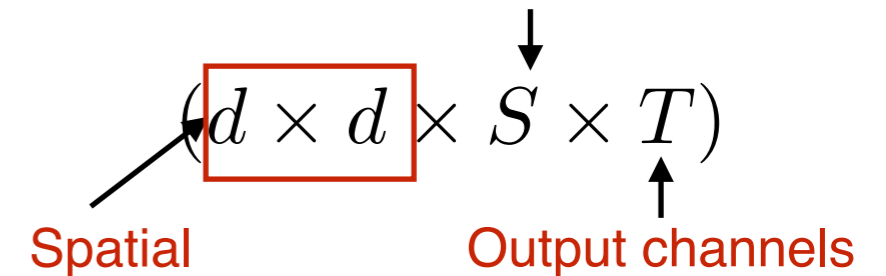
CNN Compression by Tucker

Convolution Layer

$$\mathcal{Y}_t = \mathcal{X} * \mathcal{K}_t \quad t \in [1, T]$$

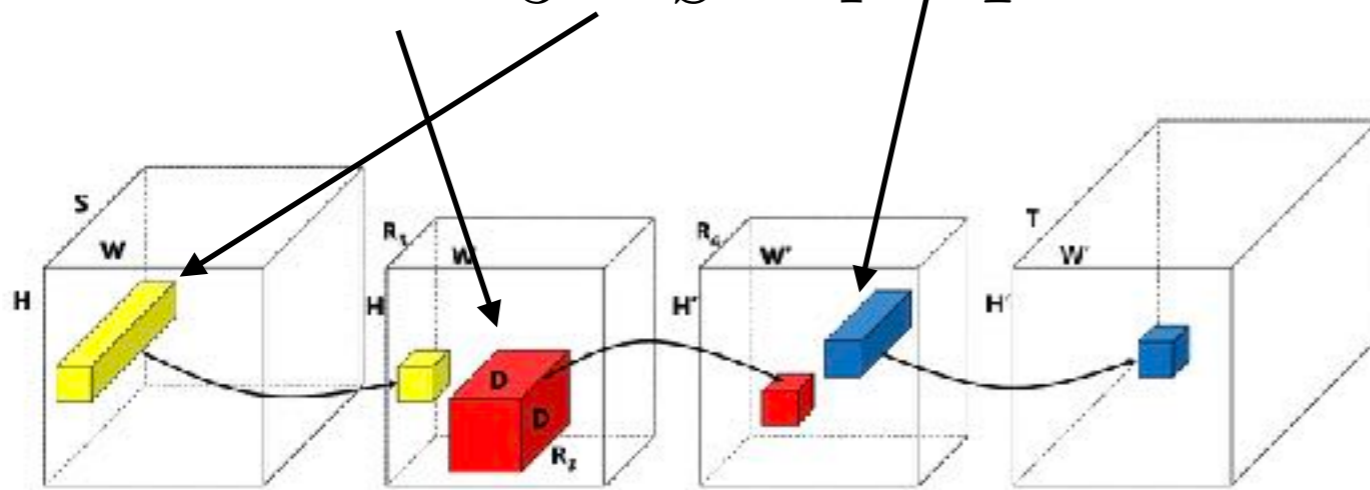
$(H \times W \times T)$ $(H \times W \times S)$ $(d \times d \times S)$

Size of Kernel \mathcal{K} Input channels



Tucker 2 decomposition

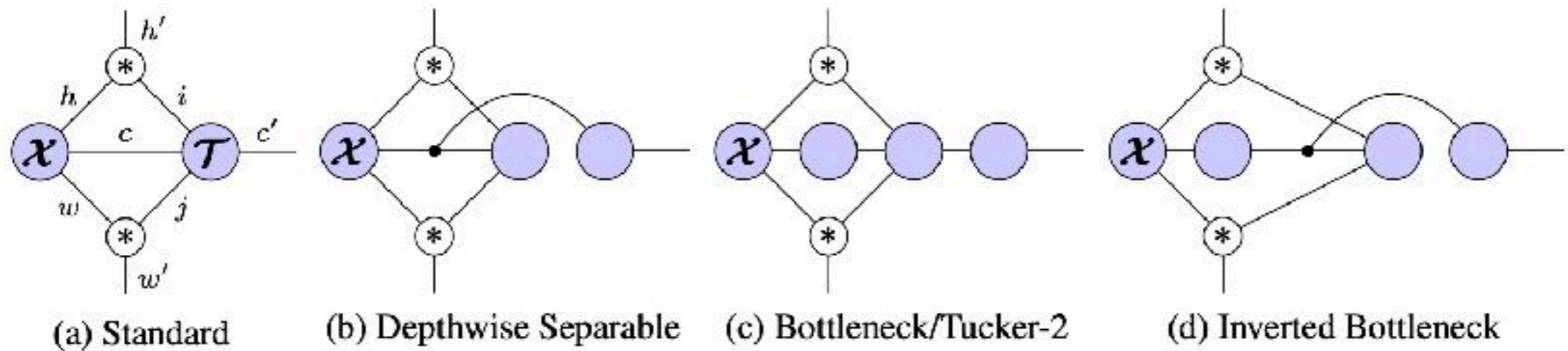
$$\mathcal{K} = \mathcal{G} \times_3 \mathbf{V}_S \times_4 \mathbf{V}_T$$



3.68x speedup for VGG

- ▶ Tensor Ranks are **hyper-parameters**
- ▶ Control trade-off between **kernel size and accuracy loss**
- ▶ Kernel size corresponds to **speed and model size**

TN Decompositions for CNN Kernels

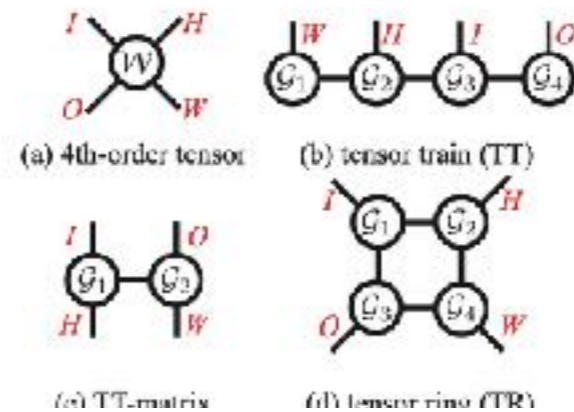
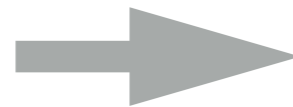


- ▶ Interpret algorithms as different TN decompositions
- ▶ Explore unexplored TN decompositions

Is Weight Kernel Low-rank?

Low-rank regularizer

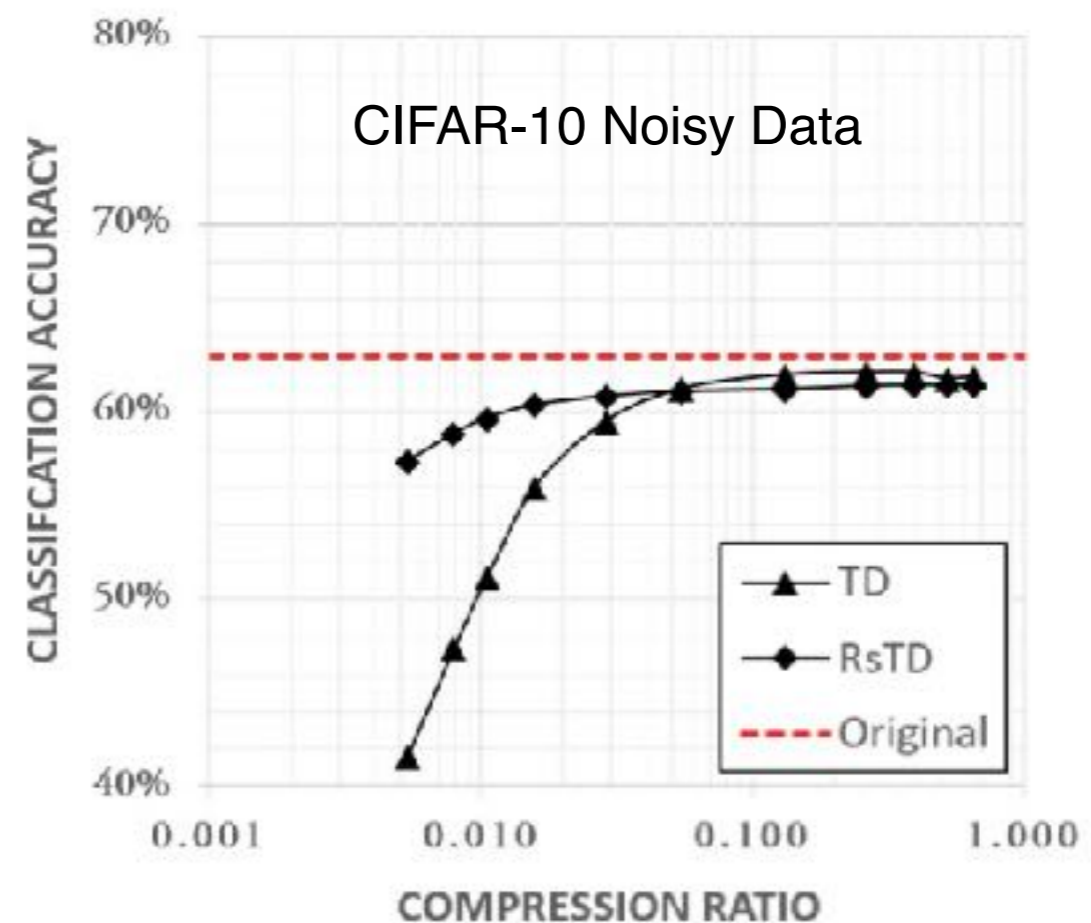
$$\mathcal{W} = \text{TD}(\mathcal{G}_1, \dots, \mathcal{G}_N)$$



Random shuffling (Rs) low-rank

$$\mathcal{W} = R(\text{TD}(\mathcal{G}_1, \dots, \mathcal{G}_N))$$

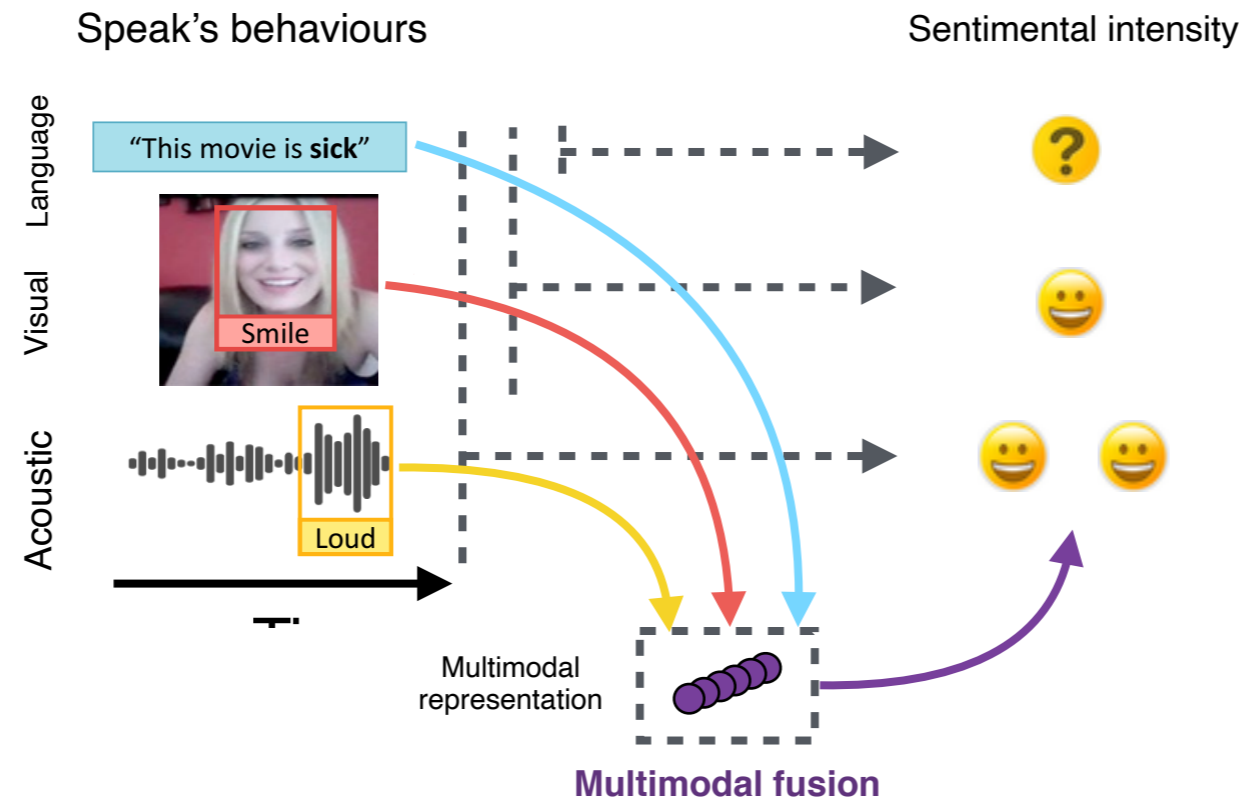
- ▶ RsTD performs better in higher compression rate
- ▶ More robust to noisy data
- ▶ CNN kernels have **inherent general low-rank** structure



Low-rank embedding of kernels in convolutional neural networks under random shuffling (Li et al., ICASSP 2019)

Multimodal Learning

- ▶ Multimodal sentimental classification (Acoustic, Visual, Language)



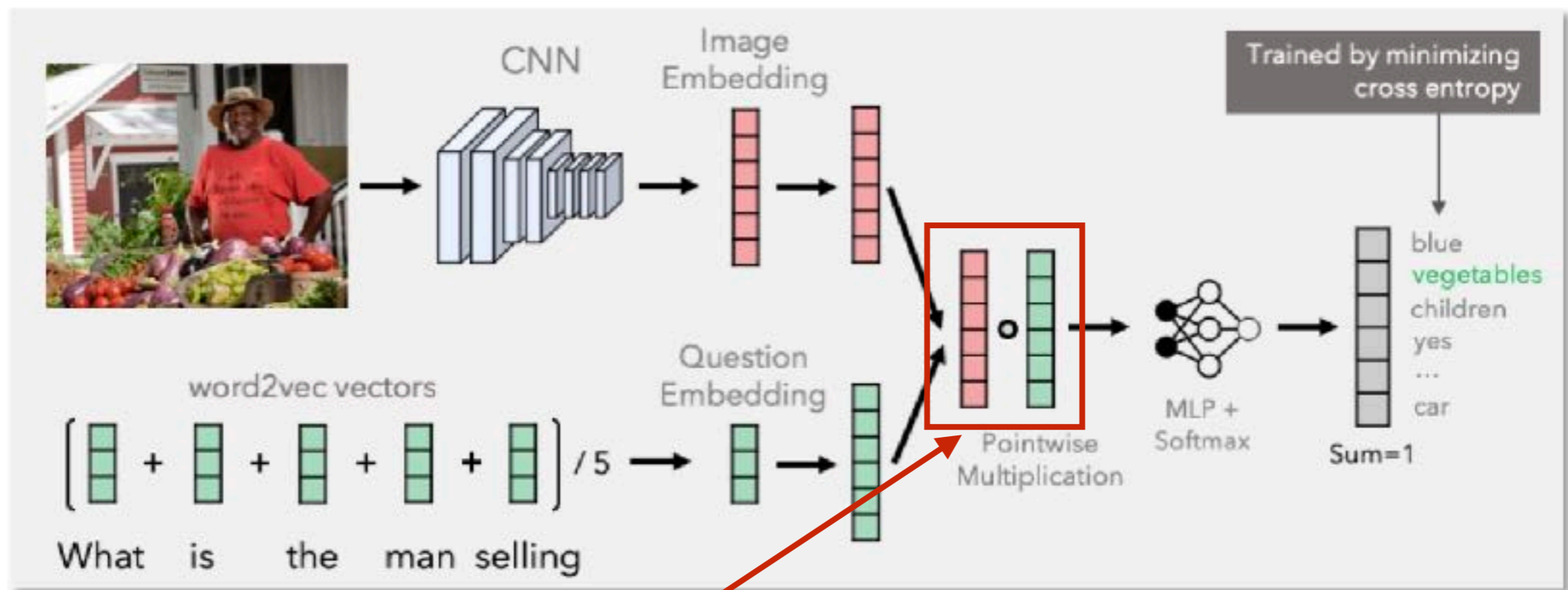
- ▶ Visual question answering (Image + Language)



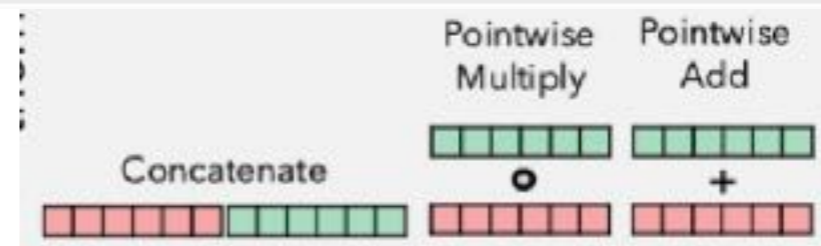
Q : "What do you see?" (Ground Truth : a₃)
a₁ : "A courtyard with flowers"
a₂ : "A restaurant kitchen"
a₃ : "A family with a stroller, tables for dining"
a₄ : "People waiting on a train"

VQA: Visual Question Answering

- ▶ Given an **image** and a **natural language question** about the image, the task is to provide an accurate **natural language answer**.



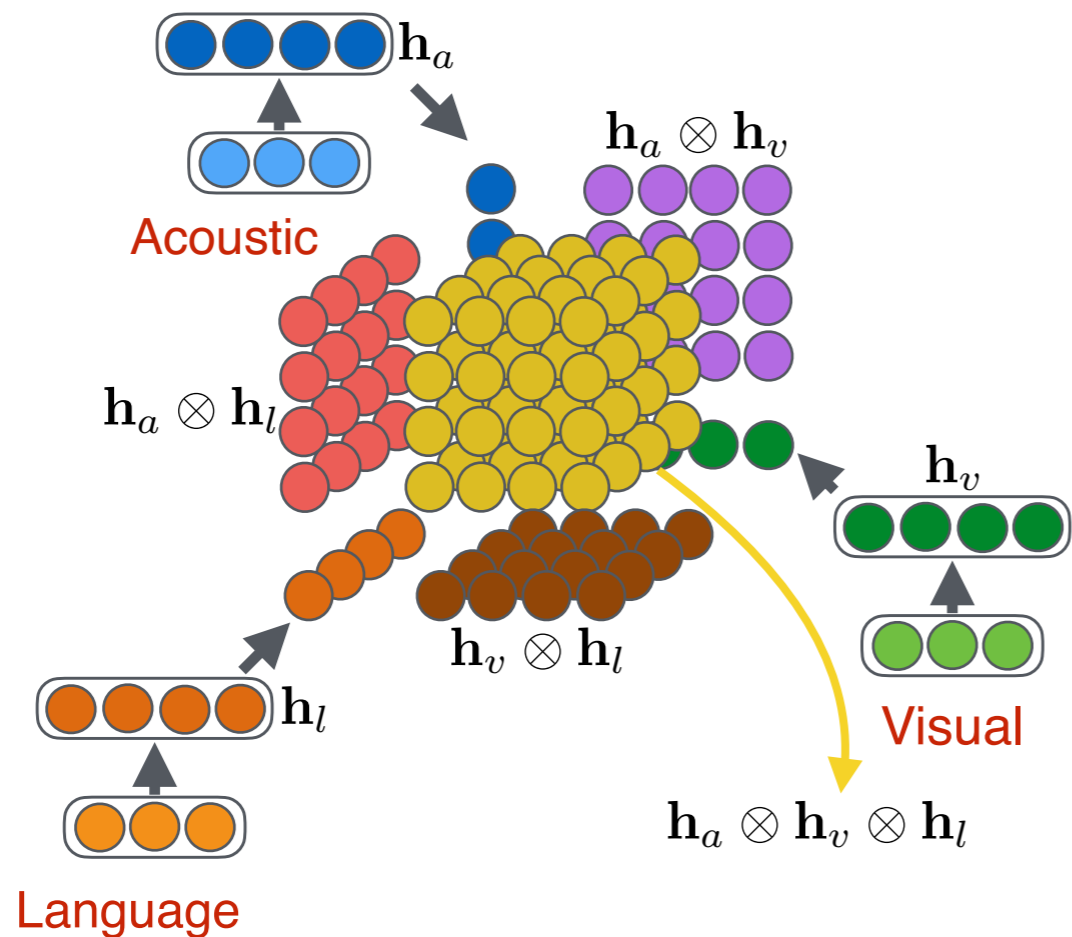
Joint Representation



VQA: Visual Question Answering (Agrawal et al., ICCV 2015)

Tensor Fusion Network for Multimodal Learning

- ▶ Trilinear fusion: linear, bilinear and trilinear interactions
- ▶ Capture cross-modality nonlinear interaction information



$$\mathbf{z}^m = \begin{bmatrix} \mathbf{z}^l \\ 1 \end{bmatrix} \otimes \begin{bmatrix} \mathbf{z}^v \\ 1 \end{bmatrix} \otimes \begin{bmatrix} \mathbf{z}^a \\ 1 \end{bmatrix}$$

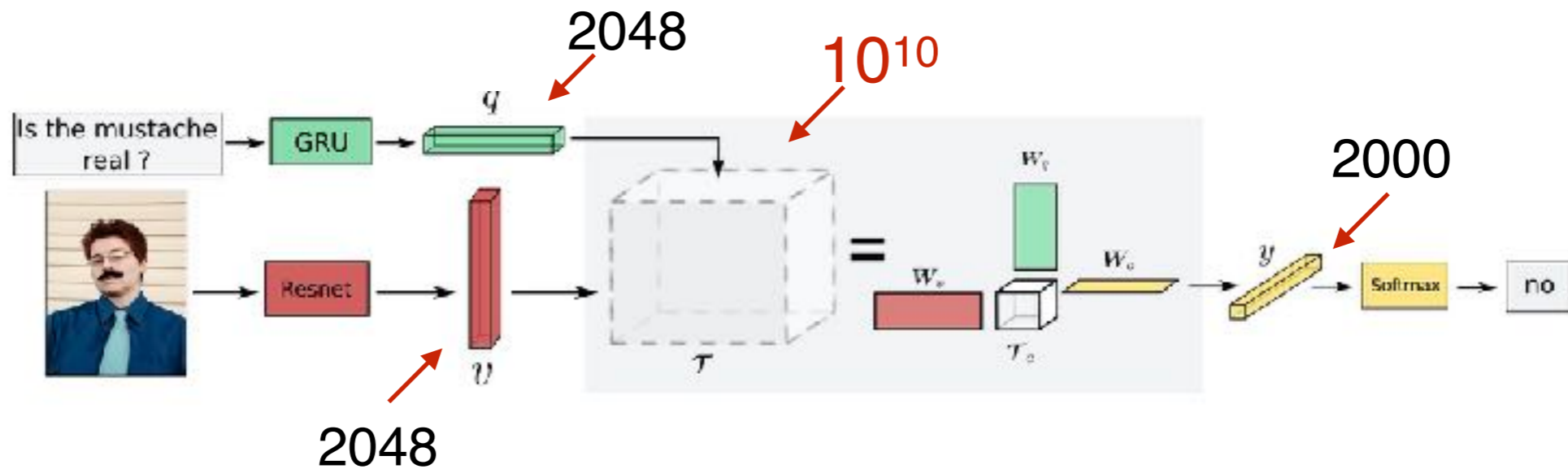
Baseline	Binary		5-class	Regression	
	Acc(%)	F1	Acc(%)	MAE	r
TFN _{language}	74.8	75.6	38.5	0.99	0.61
TFN _{visual}	66.8	70.4	30.4	1.13	0.48
TFN _{acoustic}	65.1	67.3	27.5	1.23	0.36
TFN _{bimodal}	75.2	76.0	39.6	0.92	0.65
TFN _{trimodal}	74.5	75.0	38.9	0.93	0.65
TFN _{notrimodal}	75.3	76.2	39.7	0.919	0.66
TFN	77.1	77.9	42.0	0.87	0.70
TFN _{early}	75.2	76.2	39.0	0.96	0.63

- ▶ Exponential increase of dimensionality and complexity

Tensor Fusion Network for Multimodal Sentiment Analysis (Zadeh et al., EMNLP 2017)

Multimodal Tucker Fusion

- ▶ Bilinear fusion suffer from **huge dimensionality** issue



- ▶ Tensor decomposition of model parameter \mathcal{T}

$$y = (\mathcal{T} \times_1 \mathbf{q}) \times_2 \mathbf{v} \quad \longrightarrow \quad y = ((\mathcal{T}_c \times_1 (\mathbf{q}^\top \mathbf{W}_q)) \times_2 (\mathbf{v}^\top \mathbf{W}_v)) \times_3 \mathbf{W}_o$$

Tucker
Decomposition: $\mathcal{T} = ((\mathcal{T}_c \times_1 \mathbf{W}_q) \times_2 \mathbf{W}_v) \times_3 \mathbf{W}_o$

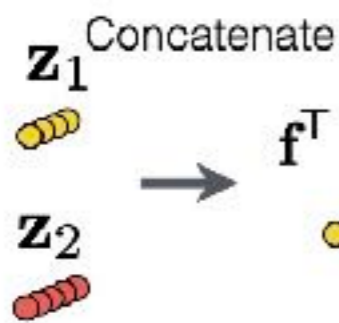
Tensor Rank is key to balance expressivity and complexity

MUTAN: Multimodal Tucker Fusion for Visual Question Answering (Ben-younes et al., ICCV 2017)

High-order Tensor Fusion

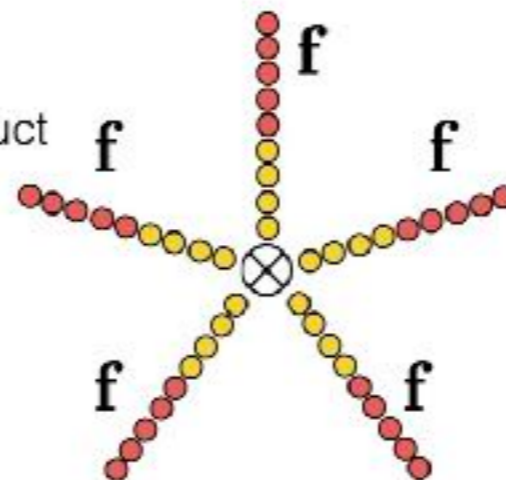
- ▶ Expressive power of tensor fusion is limited
- ▶ From first-order to **high-order** intra-modal and cross-modal feature interactions
- ▶ **Tensor Polynomial Pooling (PTP)**

Modality 1



Modality 2

P-order
tensor product



Feature Interactions

- ▶ Linear
- ▶ Bilinear
- ▶ Trilinear
- ▶ Intra-modal
- ▶ High-order

$$\mathcal{F} = \underbrace{\mathbf{f} \otimes \mathbf{f} \otimes \dots \otimes \mathbf{f}}_{P\text{-order}}$$

Dimensionality increases exponentially with P

Example: $\mathcal{F}_{1,1,1} = f_1^3$, $\mathcal{F}_{1,2,1} = f_1^2 f_2$,
 $\mathcal{F}_{1,2,3} = f_1 f_2 f_3$, $\mathcal{F}_{1,2,2} = f_1 f_2^2$

$$(md)^P$$

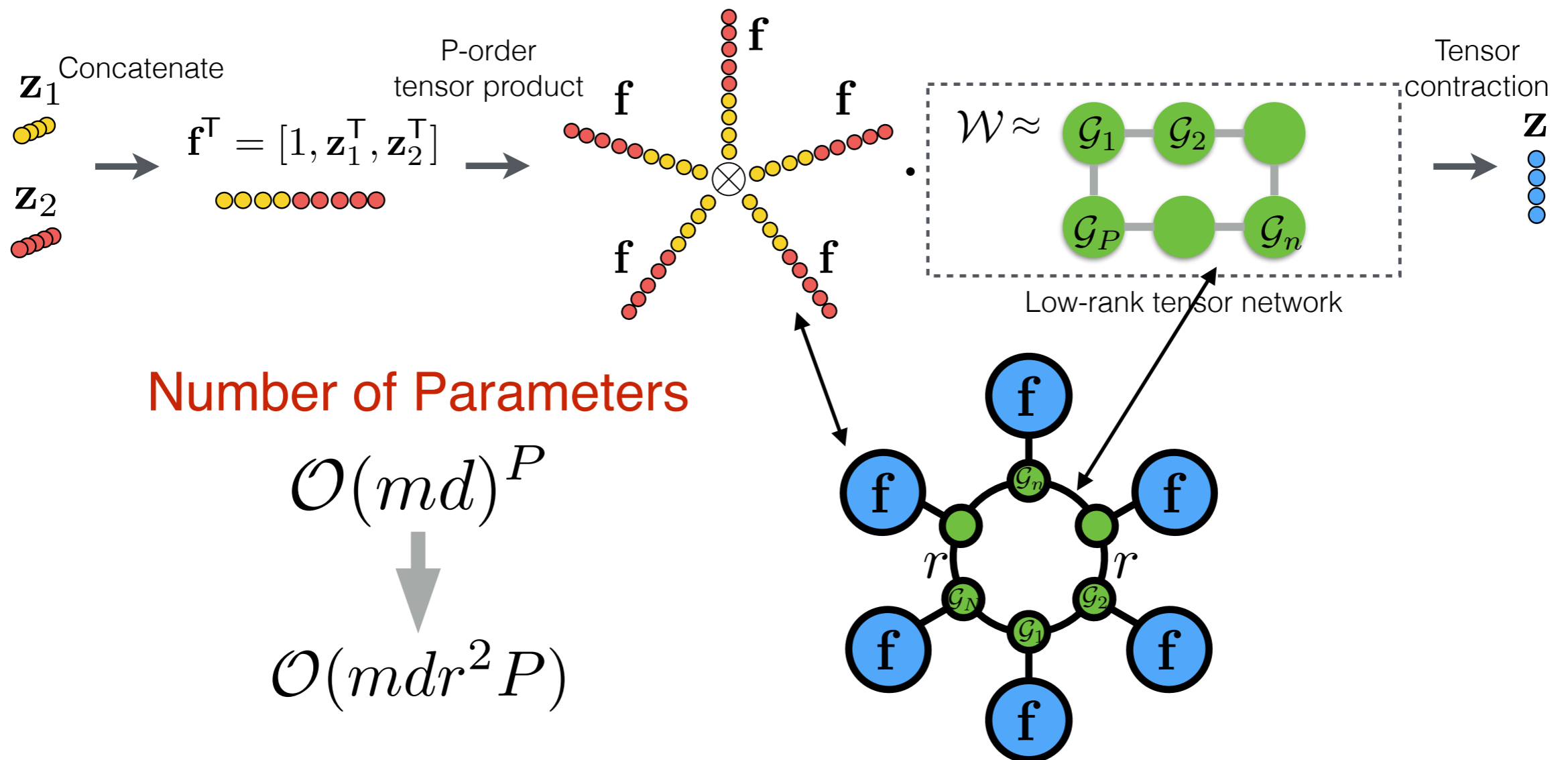
Number of modality

Dimension of feature

Order

Deep Multimodal Multilinear Fusion with High-order Polynomial Pooling (Hou et al., NeurIPS 2019)

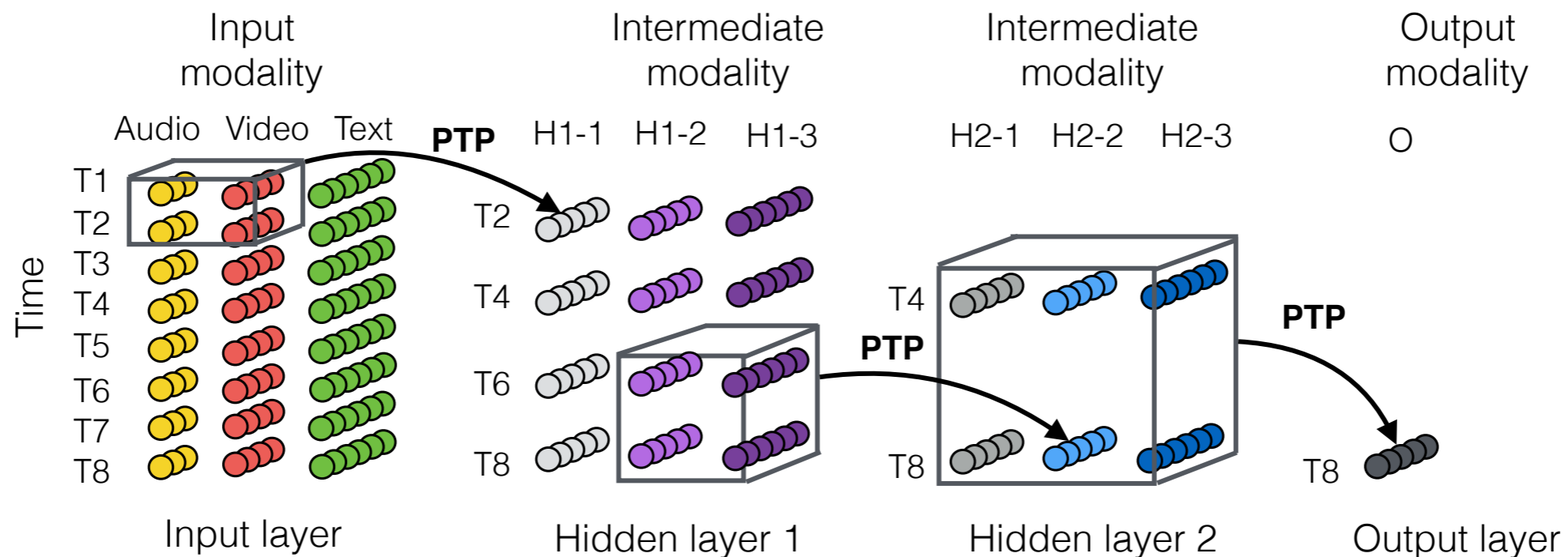
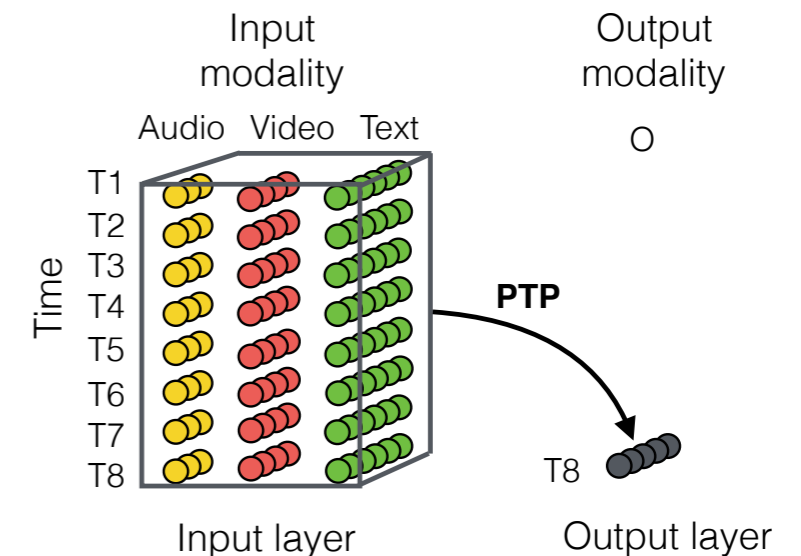
Tensor Polynomial Pooling (PTP)



- Tensor network representation of \mathcal{W} is crucial

Hierarchical Polynomial Fusion Network (HPFN)

- ▶ Interactions between **modality** and **time steps** are modeled explicitly
- ▶ **Local** temporal-modality correlations are fused and **multi-layers** are designed
- ▶ PTP resembles **convolutional filter**



Tensor Networks for Theoretical Study of DNNs

Theoretical understanding of DNNs is limited

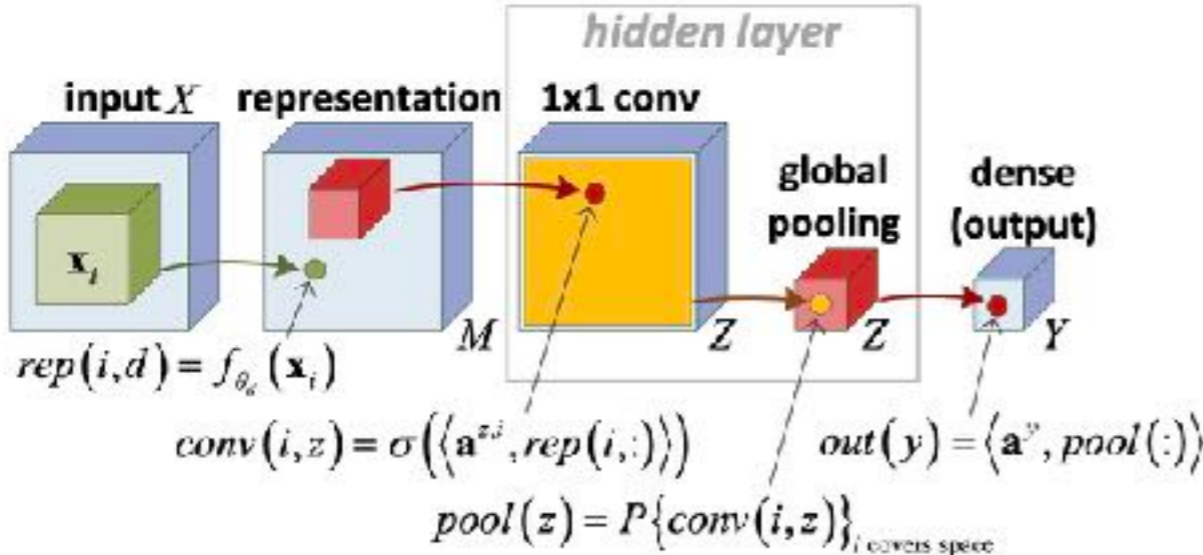
- ▶ How to analyze expressive power of DNNs?
- ▶ How is depth efficiency of DNNs?
- ▶ Generalization error bound
- ▶ ...

Tensor Networks can be used as a tool for theoretical analysis of DNNs

- ▶ **Relation and equivalence** between TNs and DNNs
- ▶ Expressive power of network can be measured by **rank of tensor networks**
- ▶ **Expressive analysis of DNNs** by analyzing TN ranks
- ▶ Provide valuable insights on how they can be improved

Equivalence between CNNs and TNs

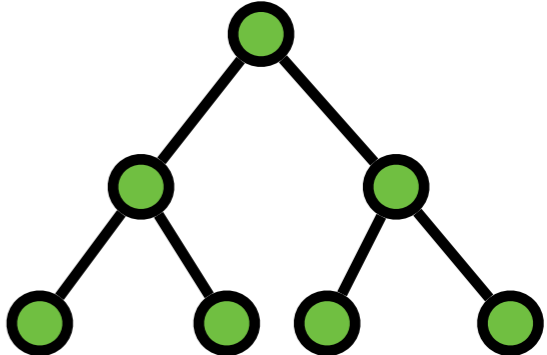
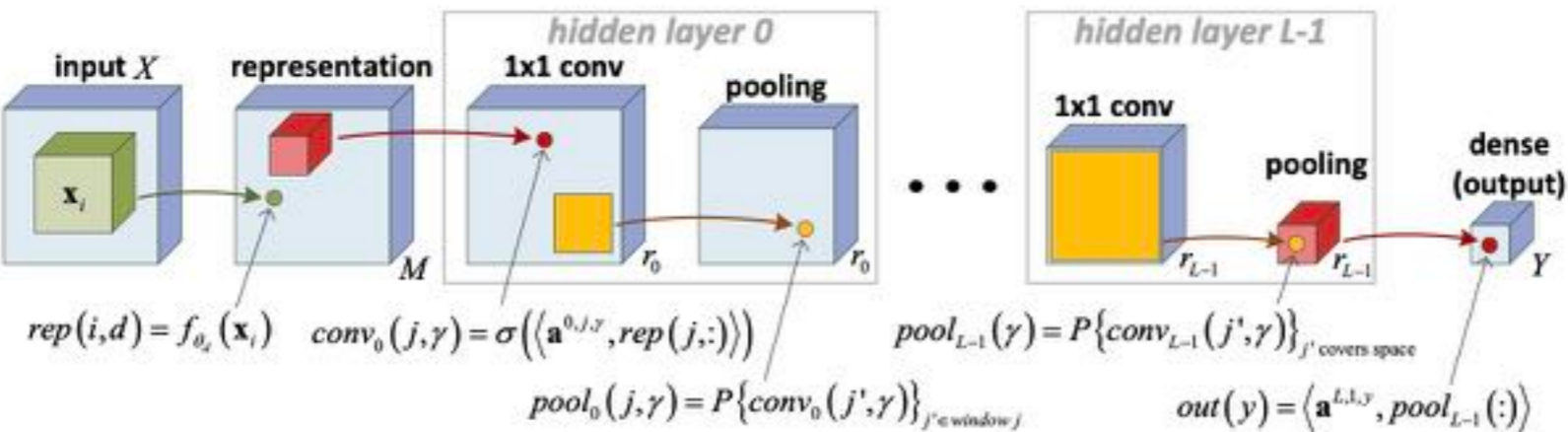
- ▶ Shallow ConvNet corresponds to CP decomposition



Score function:

$$A(h_y^S) = \sum_{z=1}^Z a_z^y \cdot (F \mathbf{a}^{z,1}) \otimes_g \cdots \otimes_g (F \mathbf{a}^{z,N})$$

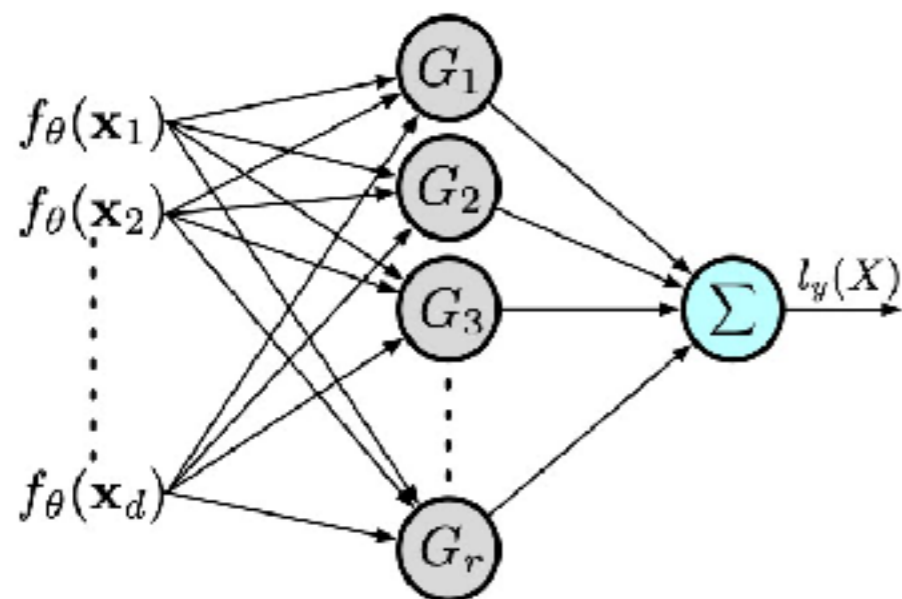
- ▶ Deep CNN network corresponds to hierarchical Tucker decomposition



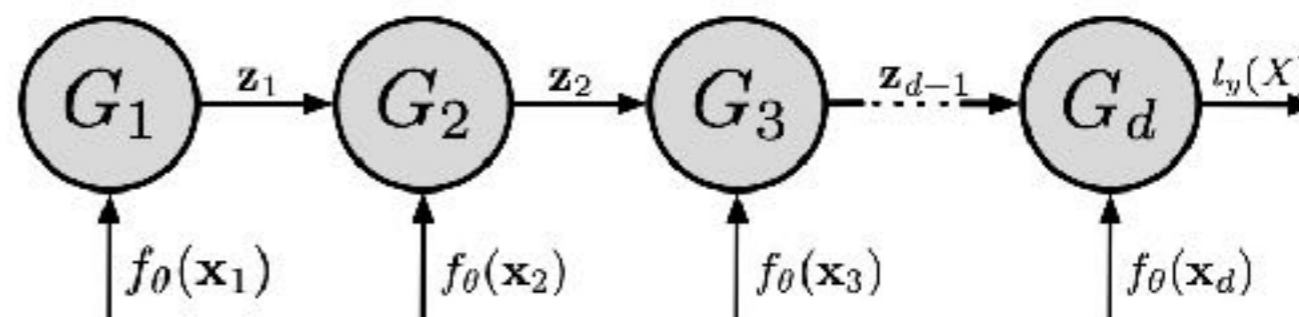
hierarchical Tucker

Convolutional Rectifier Networks as Generalized Tensor Decompositions (Cohen et al., ICML 2016)

Relations between TNs and DNNs

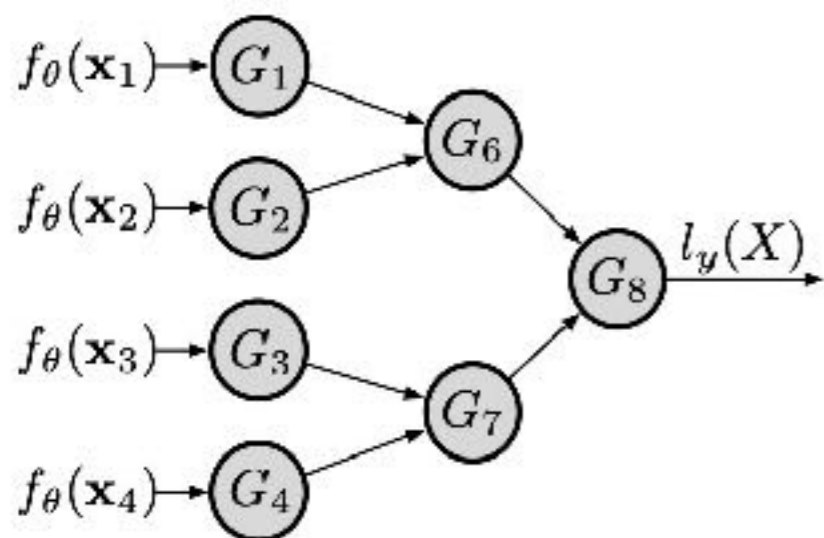


(a) CP-Network



TT-Network

Score function: $l_y(X) = \langle \mathcal{W}_y, \Phi(X) \rangle$



(b) HT-Network

Tensor Decompositions

CP-decomposition
TT-decomposition
HT-decomposition
rank of the decomposition

Deep Learning

shallow network
RNN
CNN
width of the network

Expressive Power of Recurrent Neural Networks (Khrulkov et al., ICLR 2018)

Theoretical Results

	TT-Network	HT-Network	CP-Network
TT-Network	r	$r^{\log_2(d)/2}$	r
HT-Network	r^2	r	r
CP-Network	$\geq r^{\frac{d}{2}}$	$\geq r^{\frac{d}{2}}$	r

Given a network of width r shown in a column, what is the upper bound on the width of equivalent networks shown in rows

- ▶ TT is exponentially more expressive than CP
- ▶ A shallow network of exponentially large width is required to mimic a recurrent neural network

Software for Tensor Networks

- ▶ **Google TensorNetwork** (<https://github.com/google/TensorNetwork>)
- ▶ **Tensorly** (<http://tensorly.org/>)
- ▶ **ITensor** (<https://itensor.org>)
- ▶ **TeNPy** (<https://tenpy.readthedocs.io/>)
- ▶ **Tntorch** (<https://tntorch.readthedocs.io/>)
- ▶ **TT-Toolbox** (<https://github.com/oseledets/TT-Toolbox>)
- ▶ Tensor network machine learning (**TNML**) (<https://github.com/emstoudenmire/TNML>)
- ▶ **TT_RNN** (https://github.com/Tuyki/TT_RNN)
- ▶ **TensorNet** (<https://github.com/Bihaqo/TensorNet>)
- ▶ More ...

Frontier and Future Direction

TNs for Function Approximation

Supervised learning

- ▶ Identify a nonlinear function from training dataset $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_M, y_M)\}$

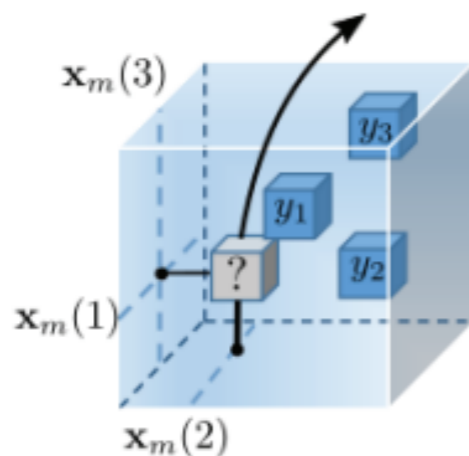
$$y = f(x_1, \dots, x_N)$$



N discrete inputs and an output

- ▶ The function space is modeled as tensor

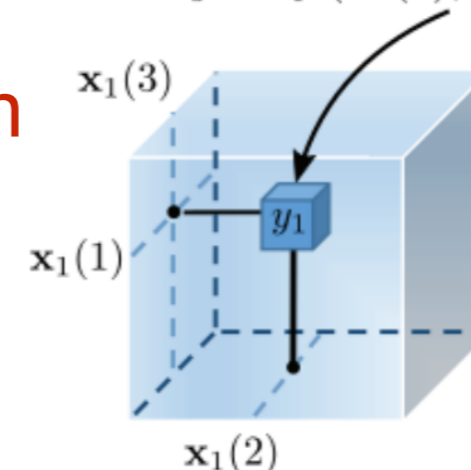
$$y_m = f(\mathbf{x}_m(1), \mathbf{x}_m(2), \mathbf{x}_m(3))$$



Tensor Completion



$$y_1 = f(\mathbf{x}_1(1), \mathbf{x}_1(2), \mathbf{x}_1(3))$$



Inference = Missing value prediction

Training points are observed entries

- ▶ Supervised learning to tensor completion

Nonlinear System Identification via Tensor Completion (Kargas et al. AAAI 2020)

A Simple Example

$h(x_2)$



Highly nonlinear

$h(x_1)$



Rank = 1

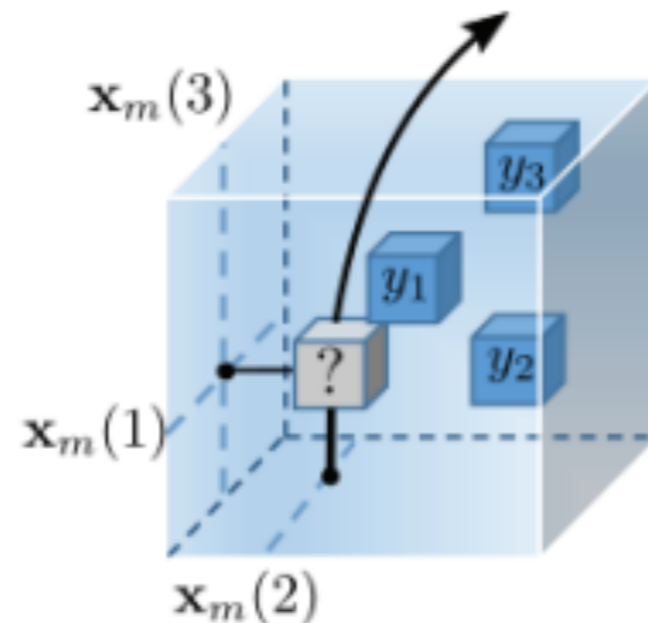
$$y = f(x_1, x_2)$$

$$f(x_1, x_2) = h(x_1)h(x_2)$$

Challenges

- ▶ Discretization of features is necessary
- ▶ Sample complexity of TNs? Is it comparable with DNNs
- ▶ Number of parameters (latent cores in TNs vs. weight parameters in DNNs)

$$y_m = f(\mathbf{x}_m(1), \mathbf{x}_m(2), \mathbf{x}_m(3))$$

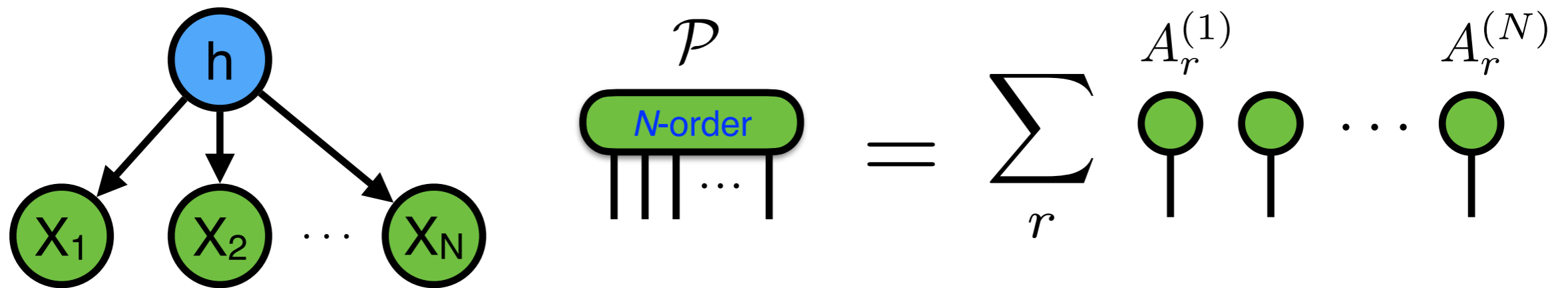


TNs for Probability Graphical Model

- ▶ Probability mass function = Non-negative Tensor

$$P(X_1, X_2, \dots, X_N) \quad X_n \in [1 : d]$$

- ▶ Tensor decomposition (CPD)



Joint Probability: $P(X_1, X_2, \dots, X_N) = \sum_h P(h)P(X_1|h)P(X_2|h) \cdots P(X_N|h)$

CPD: $P_{x_1, x_2, \dots, x_N} = \sum_r \lambda_r A_{x_1, r}^{(1)} A_{x_2, r}^{(2)} \cdots A_{x_N, r}^{(N)}$

- ▶ Conditional dependency controlled by tensor rank

Expressive power of tensor-network factorizations for probabilistic modeling (Glasser et al. NuerIPS, 2019)

Why use TN for Graphical Models

- ▶ **Efficient computation in:**
 - ▶ Marginalization, normalization (partition function)
 - ▶ Maximum log-likelihood, MAP estimate
 - ▶ Conditional distributions and sampling
 - ▶ ...

Example

- ▶ Factorize probability tensor in MPS/TR format

$$p(\mathbf{z}) = \prod_{k=1}^N \text{tr}(Q^{(k)}[z_k]),$$

- ▶ Unconstrained $p(\mathbf{z})$: I^N .

- ▶ MPS format $p(\mathbf{z})$: INr^2 .

where $Q^{(k)}[z_k] \in \mathbb{R}_+^{r_k \times r_{k+1}}$, r_k are MPS-ranks (or, boundary conditions).

- ▶ The marginal distribution is

$$p(z_{1:k-1}, z_{k+1:N}) = \sum_{i_k=1}^{I_k} p(z_1, \dots, z_k = i_k, \dots, z_N)$$

$$\tilde{Q}_k = \sum_{i_k=1}^{I_k} Q^{(k)}[i_k]$$

$$= \sum_{i_k=1}^{I_k} \text{tr} \left(\prod_{d=1}^{k-1} Q^{(d)}[z_d] \cdot Q^{(k)}[i_k] \cdot \prod_{d=k+1}^N Q^{(d)}[z_d] \right)$$

$$= \text{tr} \left(\prod_{d=1}^{k-1} Q^{(d)}[z_d] \cdot \tilde{Q}_k \cdot \prod_{d=k+1}^N Q^{(d)}[z_d] \right),$$

Summation simply performed on core tensor

Example

- ▶ Normalization

$$Z = \text{tr} \left(\prod_{k=1}^N \tilde{Q}_k \right)$$

Summation over all core tensors

- ▶ Conditional distribution

$$p(z_{1:k-1}, z_{k+1:N} \mid z_k) = \frac{p(z_{1:N})}{p(z_k)},$$

$$p(z_k) = \text{tr} \left(\prod_{d=1}^{k-1} \tilde{Q}_d \cdot Q^{(k)}[z_k] \cdot \prod_{d=k+1}^N \tilde{Q}_d \right)$$

Summation over all core tensors except k-th

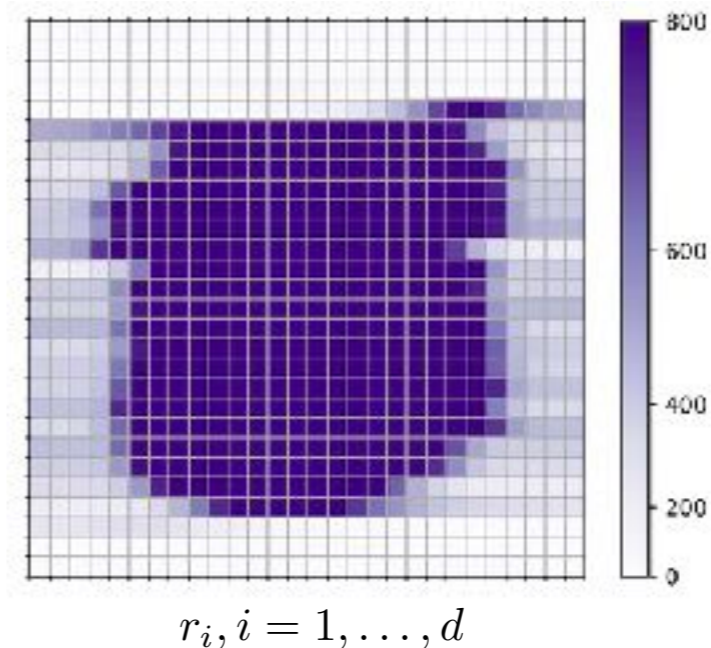
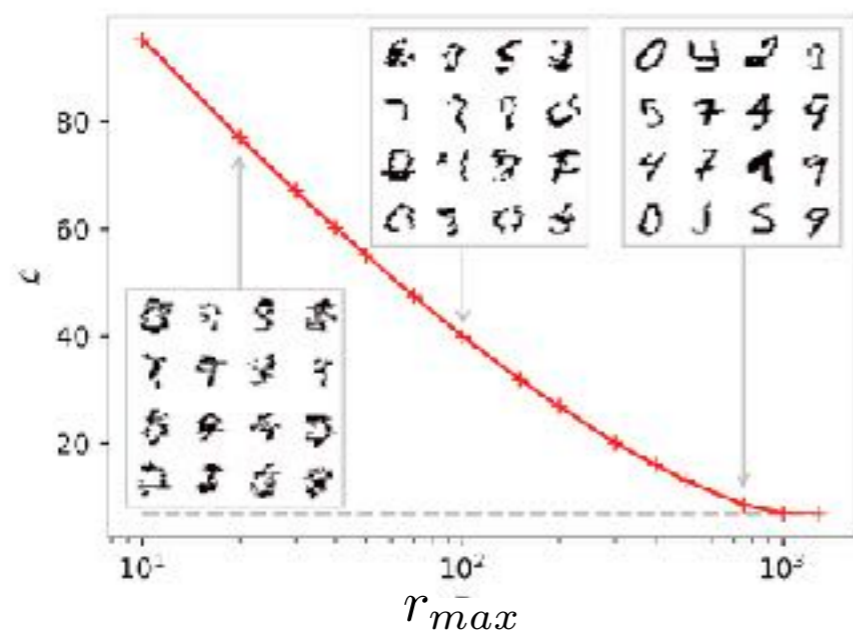
Unsupervised Generative Modeling

Han et al., Phys. Rev. X 2018

- ▶ Modeling joint probability distribution of binary variables

$$p(\mathbf{v}) = \left| \begin{array}{cccccc} \mathcal{G}_1 & & & & & \mathcal{G}_d \\ \bullet & \bullet & \bullet & \bullet & \bullet & \dots & \bullet \\ | & | & | & | & | & \dots & | \\ \bullet & \bullet & \bullet & \bullet & \bullet & \dots & \bullet \\ v_1 & & & & & & v_d \end{array} \right|^2 / \mathcal{Z}$$

- ▶ Minimizing negative log-likelihood to learn parameters $\{\mathcal{G}_i\}_{i=1}^d$
- ▶ Efficient direct sampling method to generate a sample bit by bit



Continuous Distribution by TNs

- ▶ Gaussian mixture distribution

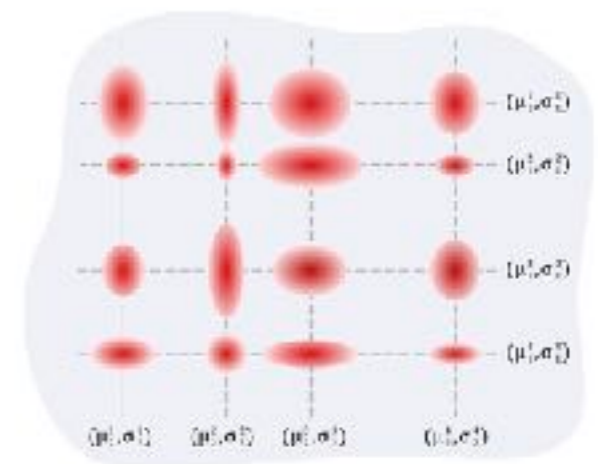
$$p(\mathbf{x}) = \sum_z p(z) \prod_{k=1}^N p(x_k | z),$$

Mixture coefficient
discrete variable

- ▶ Multi-dimensional Gaussian mixture distribution

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{z}) \prod_{k=1}^N p(\mathbf{x}_k | z_k),$$

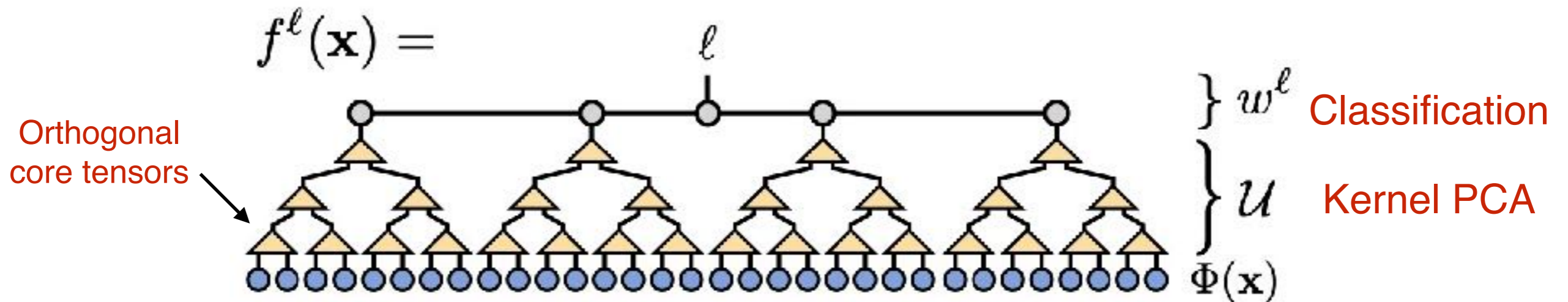
Mixture coefficient is a discrete vector
Probability is N-dimensional tensor



A prior of a googol gaussians: a tensor ring induced prior for generative models (Kuznetsov et al., NeurIPS 2019)

Learning Relevant Features of Data

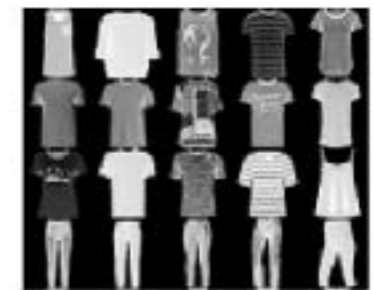
- ▶ Unsupervised learning with **reduced order of TN** representation
- ▶ Supervised learning for the top classification layer



TN Representation: $\Phi(\mathbf{x}) = \phi(x_1) \otimes \cdots \otimes \phi(x_d)$

Local feature mapping: $\phi(x_i) = [1, x_i]^T, i \in [1, d]$

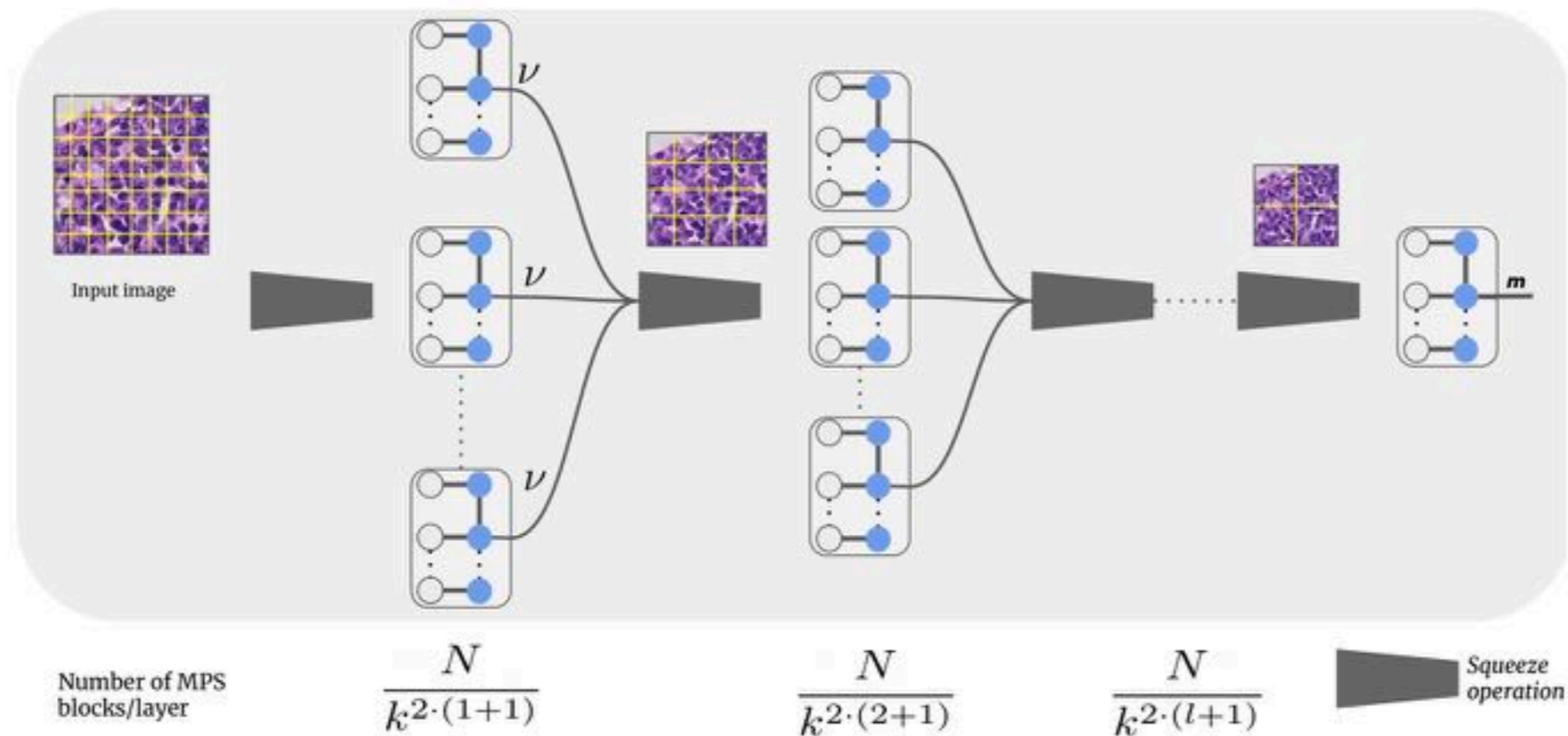
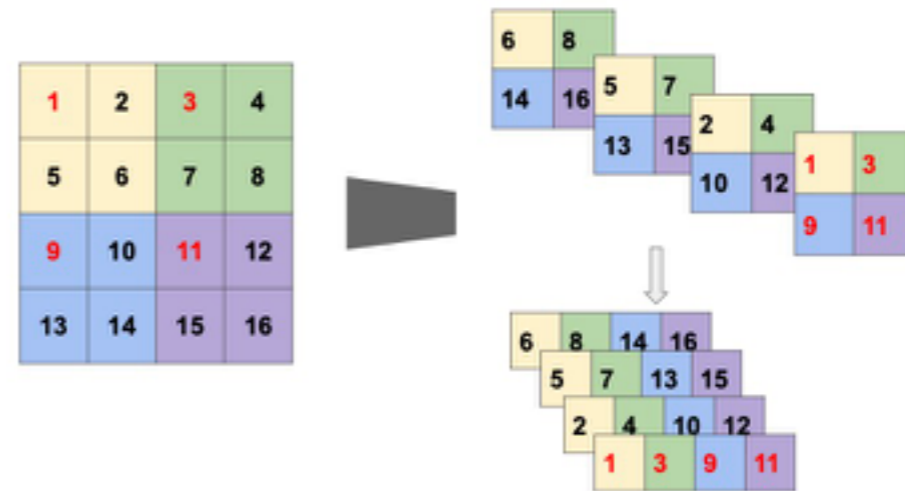
Input data: $\mathbf{x} = [x_1, \dots, x_d]^T$



89% accuracy on Fashion MNIST data set

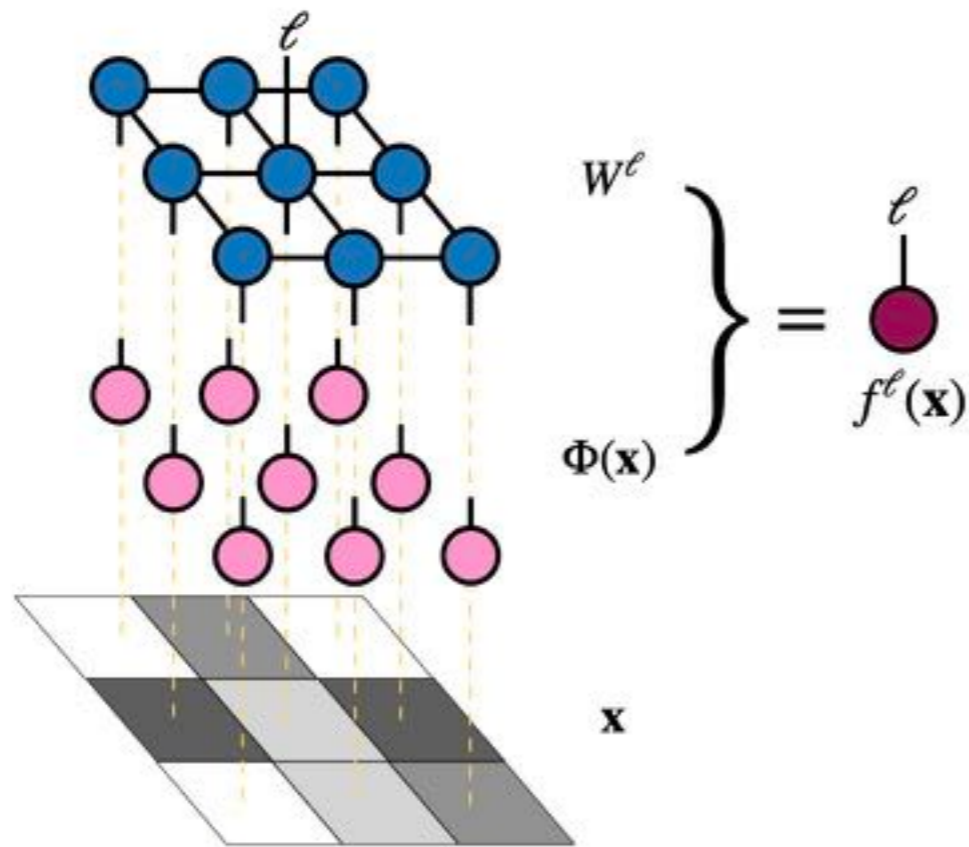
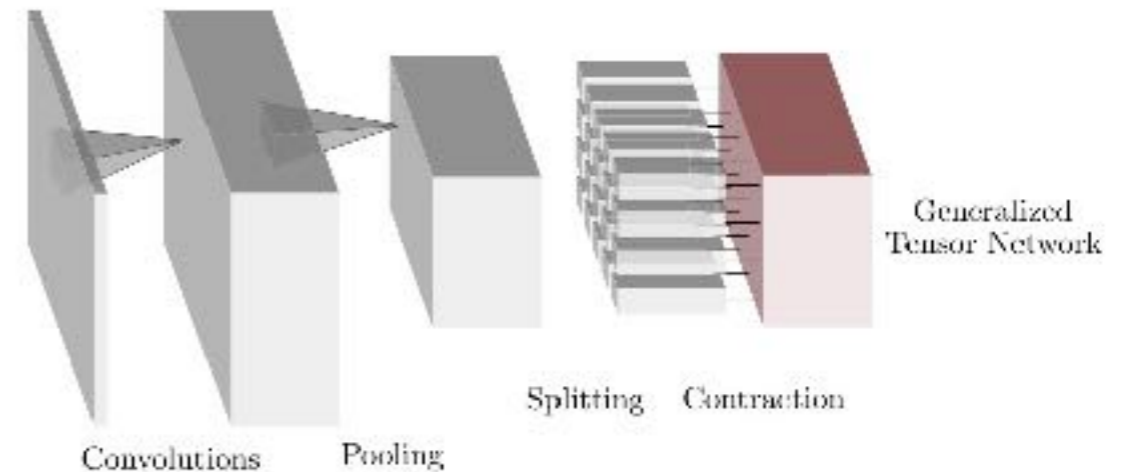
TN for Image Classification

- ▶ Capture global structure by squeeze operation
- ▶ Single layer to multi-layer TNs
- ▶ Comparable performance with DenseNet but fewer parameters

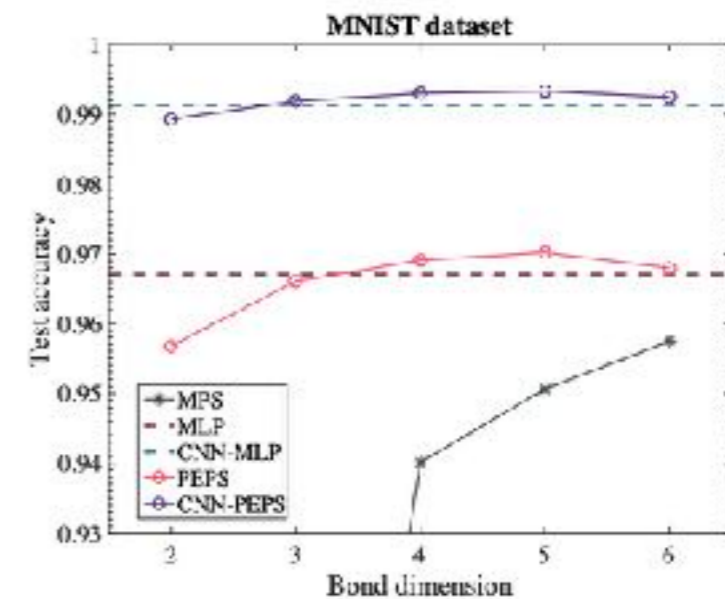


Supervised Learning with Projected Entangled Pair States

- ▶ One layer PEPS for supervised learning
- ▶ CNN + PEPS as learning model



7% parameters of CNN+MLP

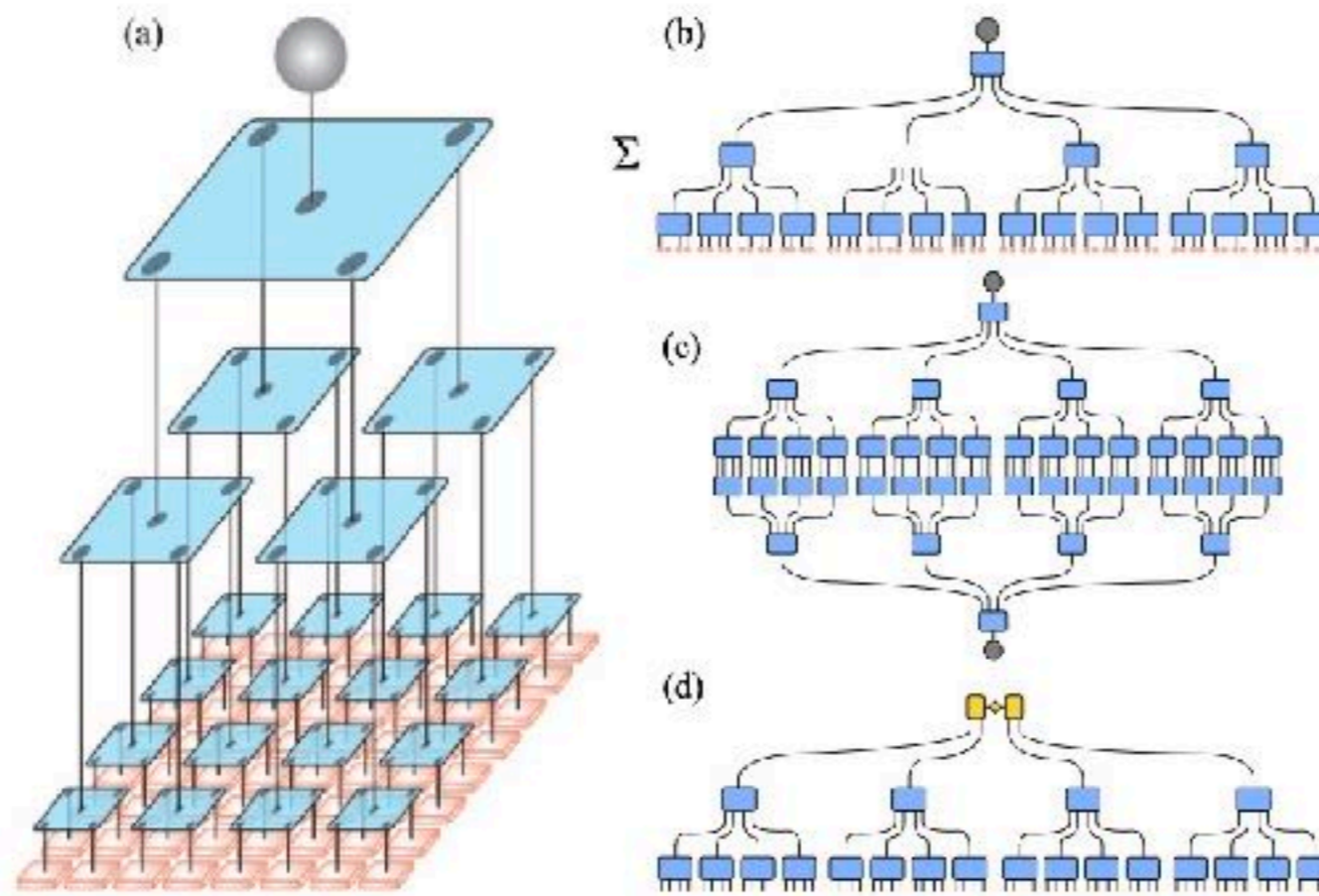


From Probabilistic Graphical Models to Generalized Tensor Networks for Supervised Learning (Glasser, arXiv 2019)

Supervised Learning with Projected Entangled Pair States (Chen et al., arXiv 2020)

Machine Learning by Hierarchical Tensor Networks

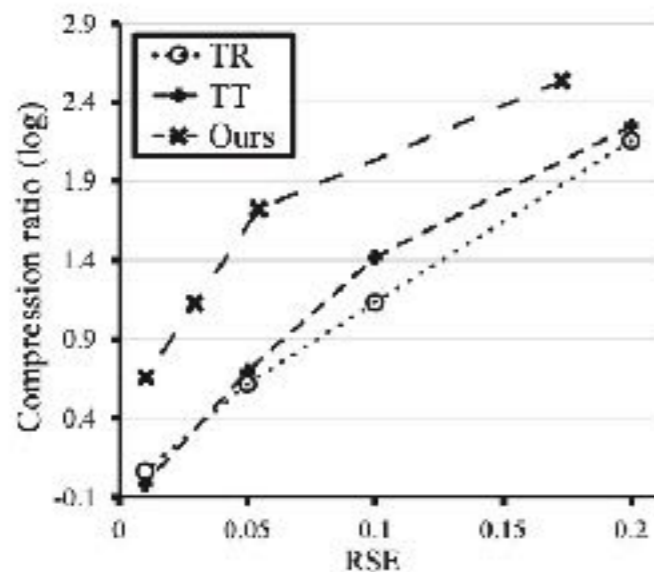
- ▶ Supervised learning with tree tensor networks



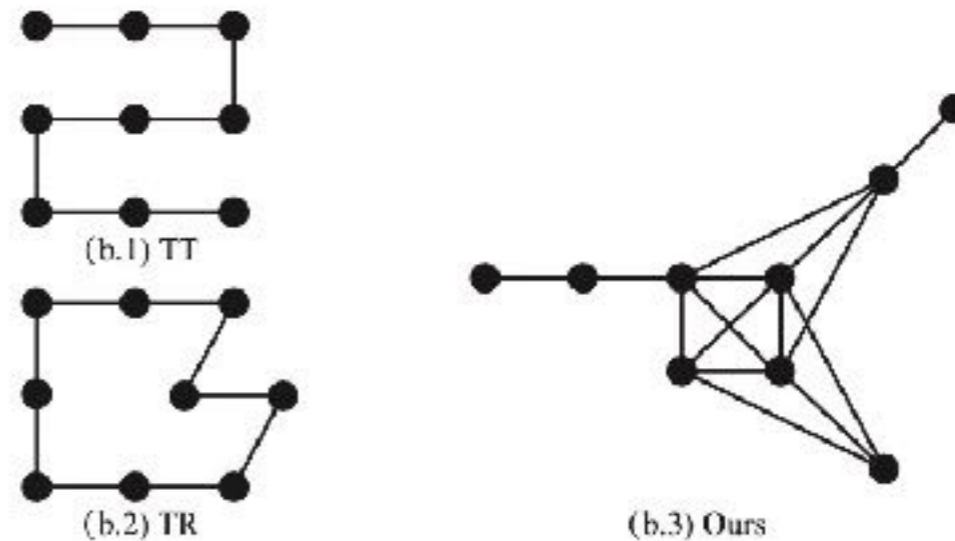
Model Architecture

Machine Learning by Unitary Tensor Network of Hierarchical Tree Structure (Liu et al., 2019)

Learning Tensor Network Structure



(a) RSE vs. CR



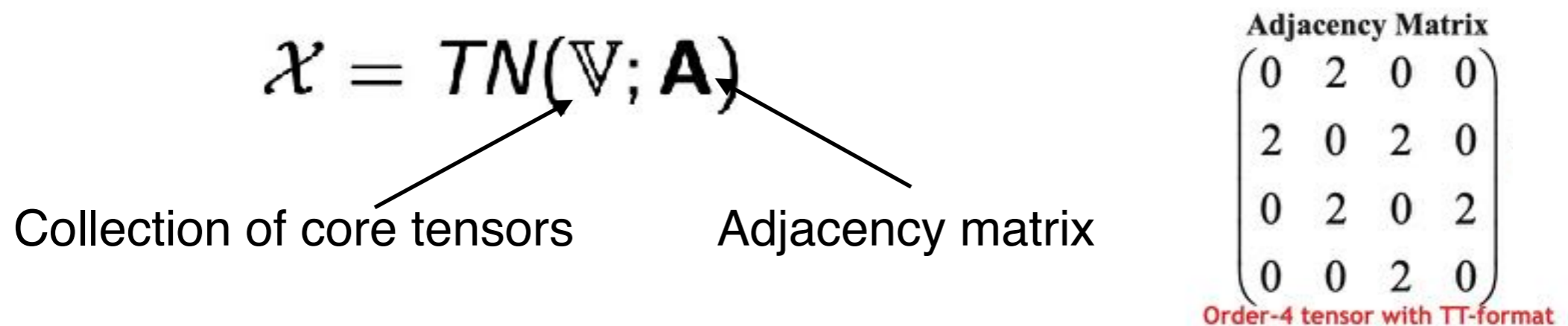
(b) Graphical structures of TN

Standard models may not be the most compressive one

- ▶ How to choose an optimal **tensor network structure** is necessary
- ▶ Can we **learn an optimal TN structure** by data?
- ▶ Challenge: given an 9-order tensor, there are more than **68 BILLION** candidates

Optimization Problem

- ▶ The TN structures can be fully described by **its adjacency matrix**



- ▶ Find an optimal **A** such that

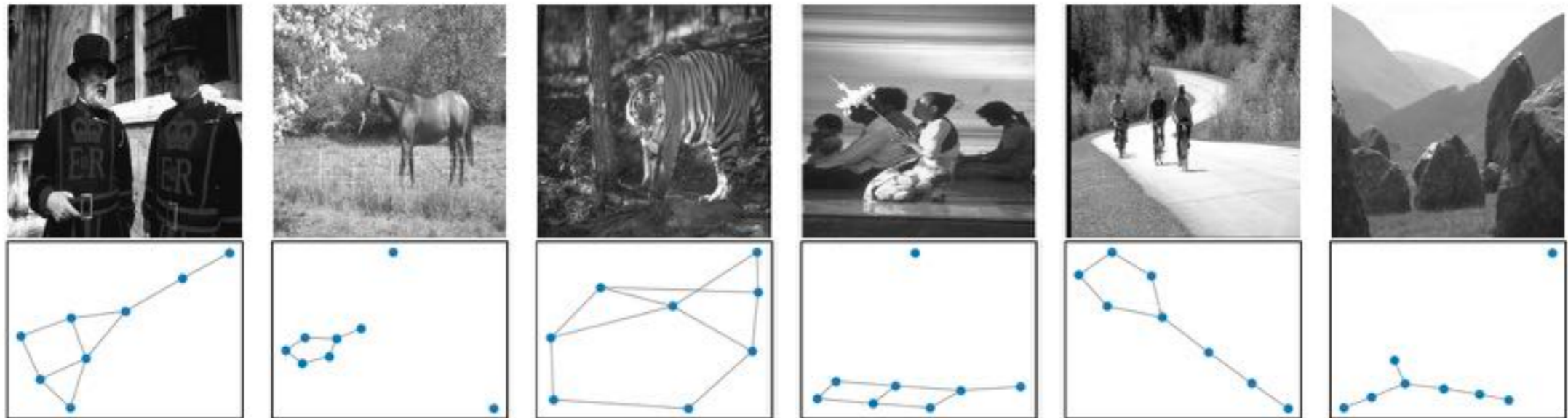
$$\min_{\mathbf{A} \in \mathbb{A}} \frac{1}{\epsilon(\mathbf{A})}, \quad s.t. \exists \hat{\mathbb{V}} \text{ which satisfies } \|\mathcal{X} - TN(\hat{\mathbb{V}}; \mathbf{A})\|_F^2 \leq \delta,$$

$$\epsilon(\mathbf{A}) = \frac{\text{Uncompressed size of } \mathcal{X}}{\text{Parameter size of } \mathbb{V} \text{ under } \mathbf{A}}.$$

- ▶ Maximize **compression ratio** such that approximation error is upper-bounded

Experiment and Discussion

- ▶ 10 natural images are random selected from LIVE dataset [Sheikh et al., 2006], and reshaped as order-8 tensors
- ▶ The learned topology have **more complex structures** than simple ones like line, tree or cycles.



- ▶ Optimal structures could significantly boost the expressive power of TN decomposition
- ▶ Near-optimal structures are sufficient to outperform the simple ones

Evolutionary Topology Search for Tensor Network Decomposition (Li et al., ICML 2020)

Discussions

- ▶ TNs are tools for data modeling, representing model parameters and function approximation
- ▶ Theory shows TNs have expressive power similar to DNNs
- ▶ Practically, performance of TNs for ML tasks is less attractive
- ▶ If TNs can be developed as a general ML model?
- ▶ Are there any advantages of TNs from perspective of robustness and interpretability?