

Algorithmic and theoretical aspects of sparse deep neural networks

Quoc Tung Le

July 26th, 2023



Léon Zheng

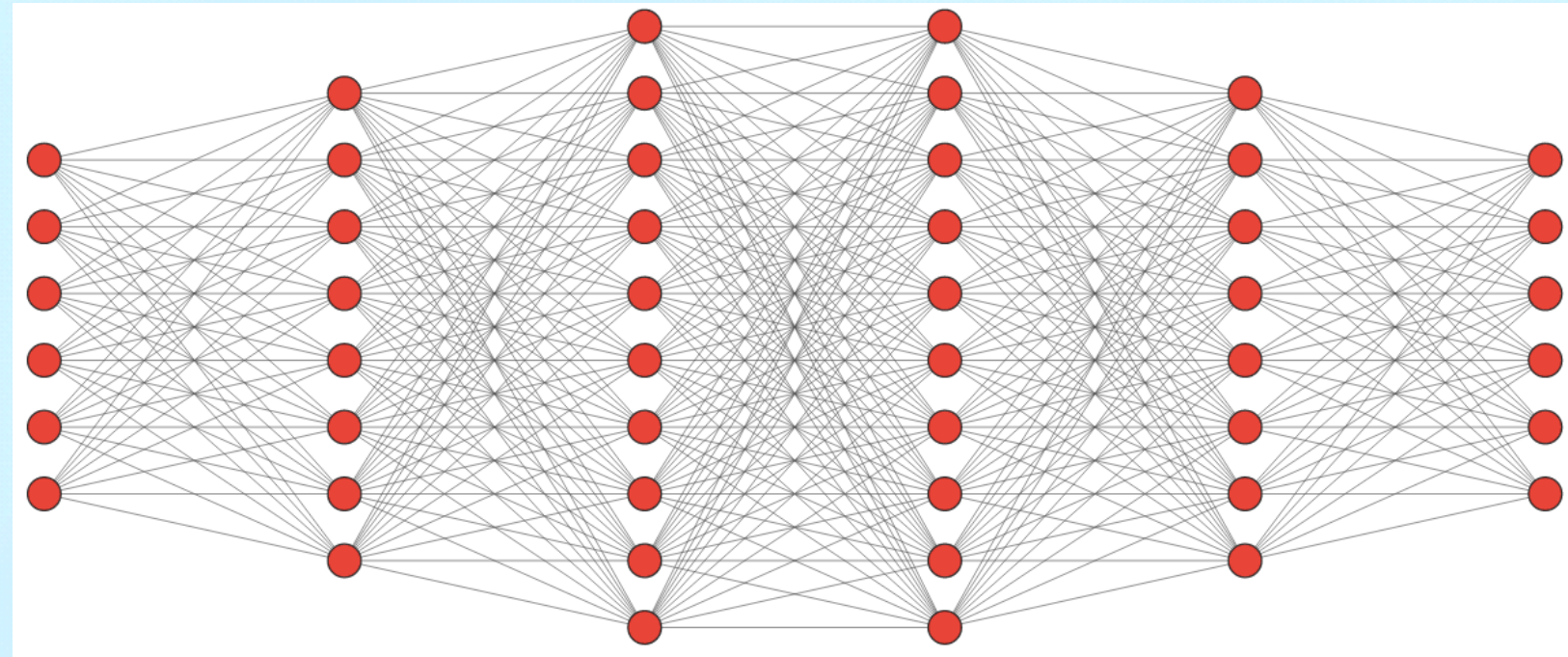


Elisa Riccietti

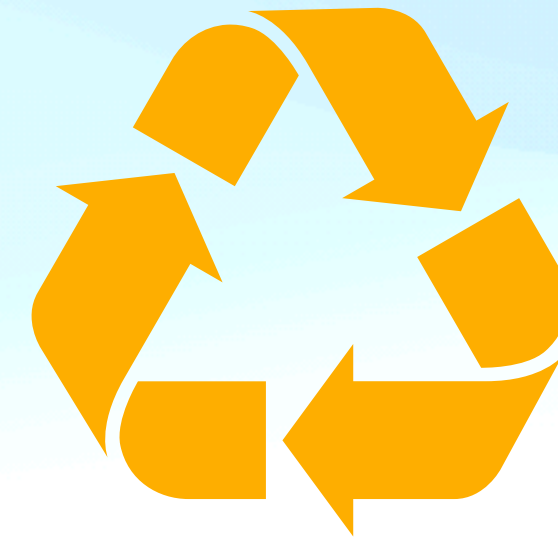


Rémi Gribonval

Deep learning and the recipe of success



Huge parametric models
(over-parameterization)



Massive data set

High computing
powers (GPUs)

Deep learning and the recipe of success

Huge parametric models

Massive dataset

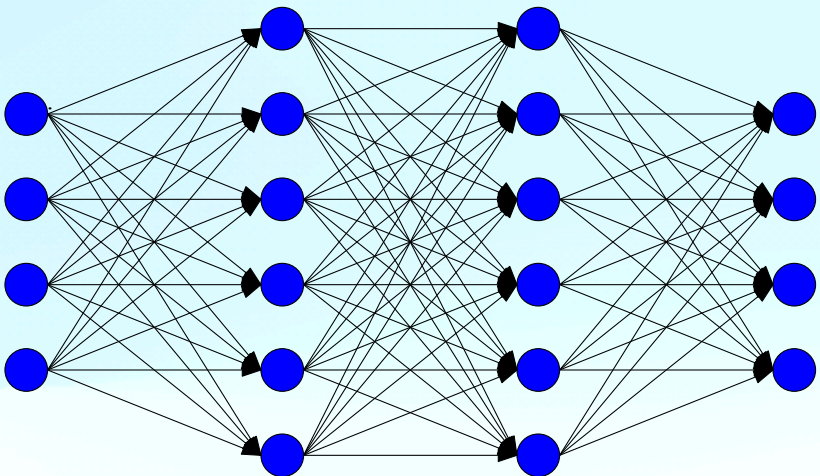
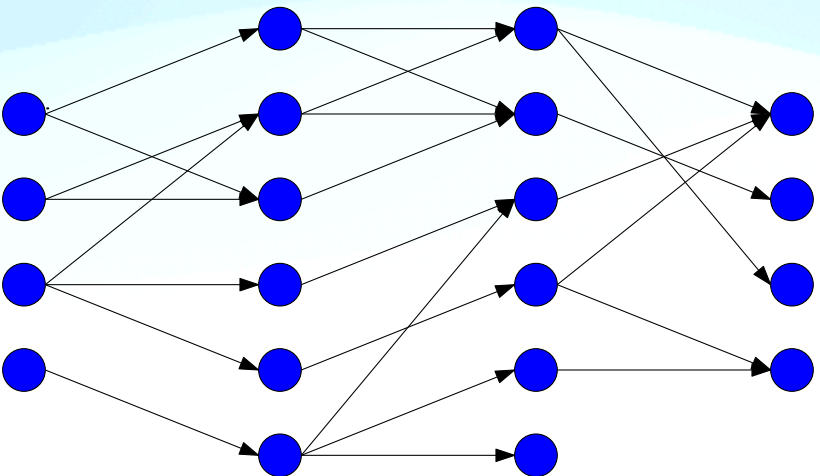
High computing powers (GPUs)



? Can we make these three components **scalable**?

Sparsity: a good old friend

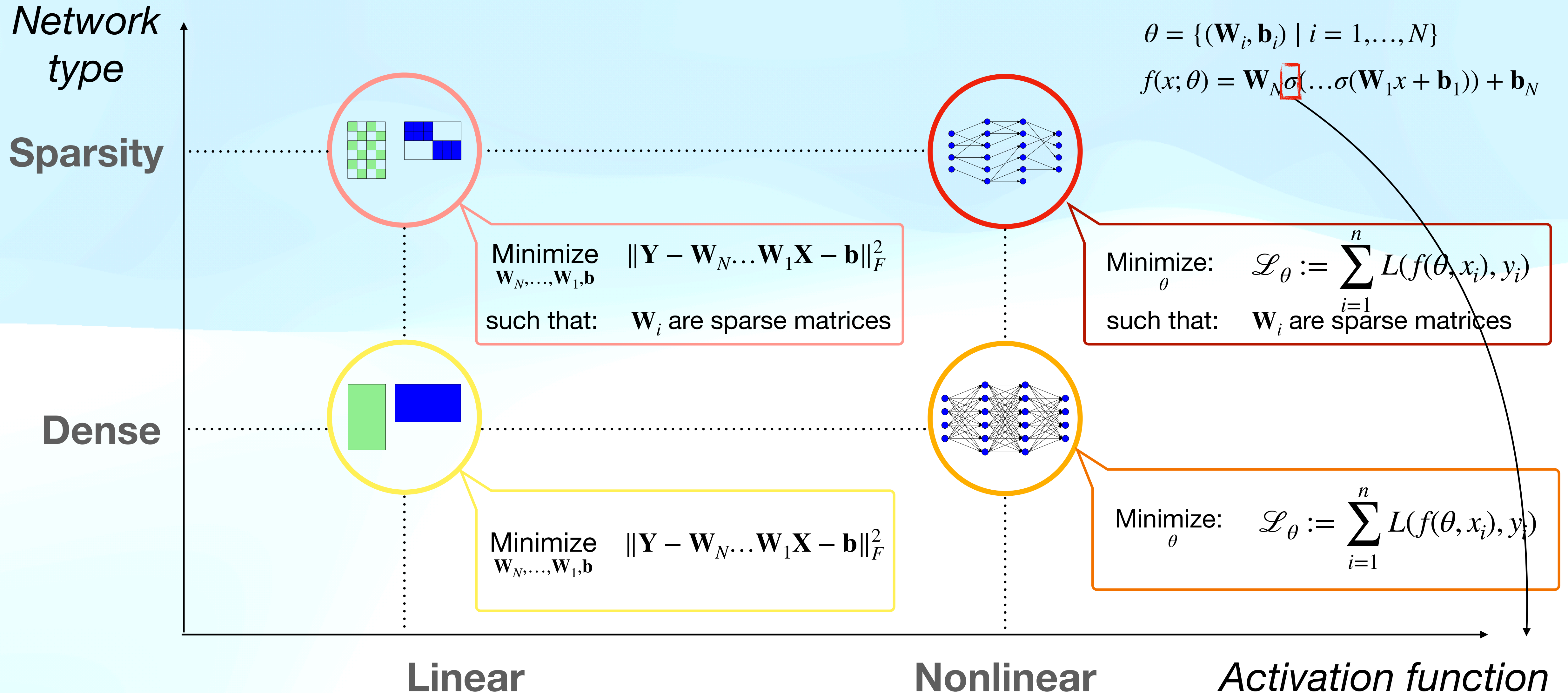
Sparse Deep Neural Networks

Usual Deep Neural Networks	Sparse Deep Neural Networks
	
<ul style="list-style-type: none">✗ Over-parameterization → storage problems✗ Extremely data hungry → large training sets✗ Memorization of training data → privacy leakage	<ul style="list-style-type: none">☑ Small parametric models (with high sparsity)☑ Fewer learned parameters → fewer training data☑ Reduce privacy issue



Sounds good. Then how can I train a sparse neural network?

Sparsity in Deep Neural Networks



Questions for Sparse Neural Networks

Existence of optimal solutions

? Does the training problem of sparse neural networks **always admit an optimal solution?**

Tractability

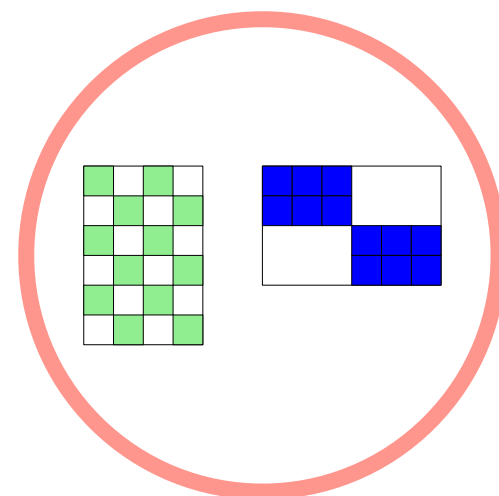
? Is it **polynomially tractable** to train sparse neural networks?

Landscape

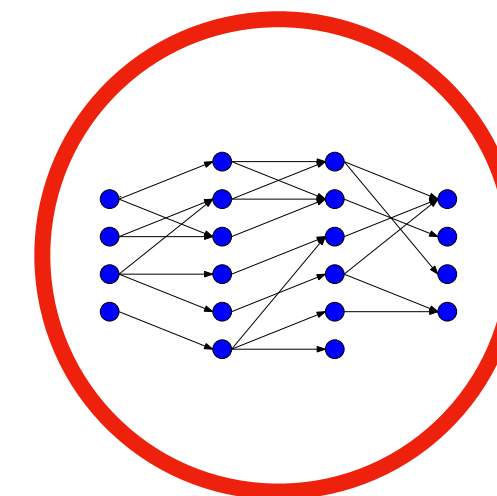
? What does the **landscape** of loss function look like? (e.g., does it have **local minima?**, etc.)

Strategy:

First



Then

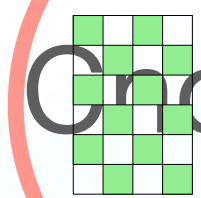


From sparsity in Deep Learning to Matrix Factorisation

Sparse Matrix Factorisation

Given A and \mathcal{E}_j some sets of *sparse* matrices, solve:

$$\min_{W^{(N)}, \dots, W^{(1)}} \left\| A - \prod_{j=1}^N W^{(j)} \right\|_F^2 \quad \text{subject to: } W^{(j)} \in \mathcal{E}_j, \forall j \in \{1, \dots, N\}$$



Choice of sparse matrices \mathcal{E}_j in Linear Sparse Neural Networks \approx Sparse Matrix Factorisation

- k -sparse per row,

- k -sparse per column

- k -sparse in total

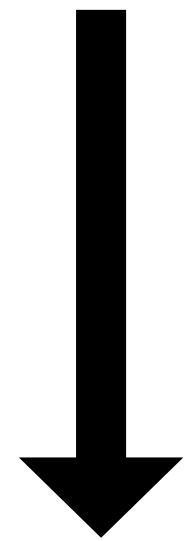
Linear Sparse Neural Network	Sparse Matrix Factorisation
Minimize $\ Y - W_N \dots W_1 X - b\ _F^2$ W_N, \dots, W_1, b	Minimize $\ A - W_N \dots W_1\ _F^2$ W_N, \dots, W_1

Sparse support factorization

Special case of sparse matrix factorization

SPARSE MATRIX
FACTORISATION

$$\min_{W^{(N)}, \dots, W^{(1)}} \|A - \prod_{j=1}^N W^{(j)}\|_F^2 \text{ subject to: } W^{(j)} \in \mathcal{E}_j, \forall j \in \{1, \dots, N\}$$



FIXED SUPPORT
MATRIX
FACTORISATION



- $N = 2$
- $(\mathcal{E}_1, \mathcal{E}_2)$: set of matrices whose **supports** are included in **given sets** I and J

$$\min_{X, Y} \|A - XY^T\|_F^2 \text{ subject to: } \text{supp}(X) \subseteq I, \text{supp}(Y) \subseteq J$$

Fixed support matrix factorization (FSMF)

$$\min_{X,Y} \|A - XY^T\|_F^2 \text{ subject to: } \text{supp}(X) \subseteq I, \text{supp}(Y) \subseteq J$$

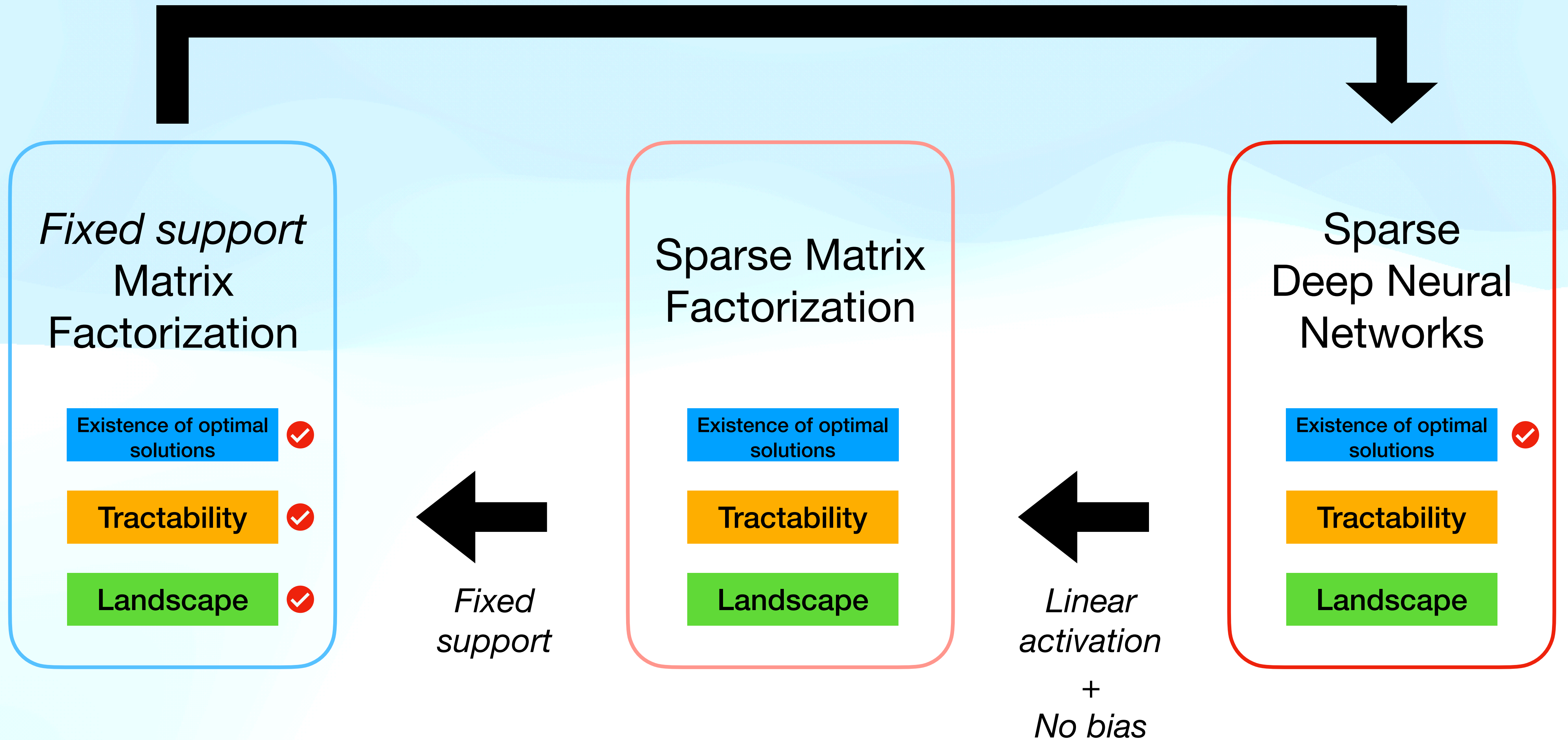
$$A \approx \begin{matrix} & \mathbf{X} & \\ \begin{matrix} ? & ? & 0 & ? \\ 0 & ? & ? & 0 \\ ? & 0 & ? & 0 \\ 0 & ? & 0 & ? \end{matrix} & \times & \begin{matrix} \mathbf{Y}^T \\ 0 & ? & 0 & ? & 0 \\ ? & ? & 0 & 0 & ? \\ ? & 0 & ? & ? & 0 \\ ? & 0 & ? & 0 & ? \end{matrix} \end{matrix}$$

inside support

outside support

SUPPORT CONSTRAINTS

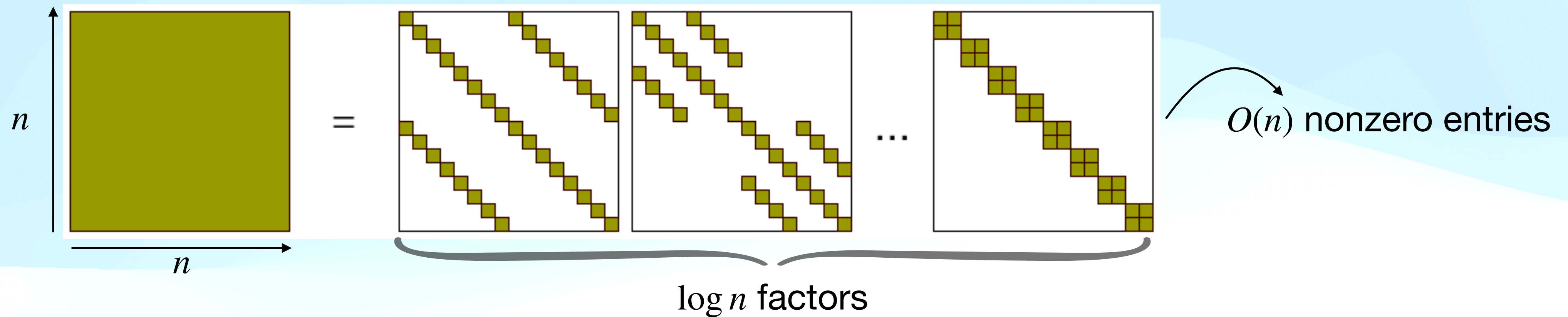
Overview of our reasoning



Further motivation for sparse matrix factorization

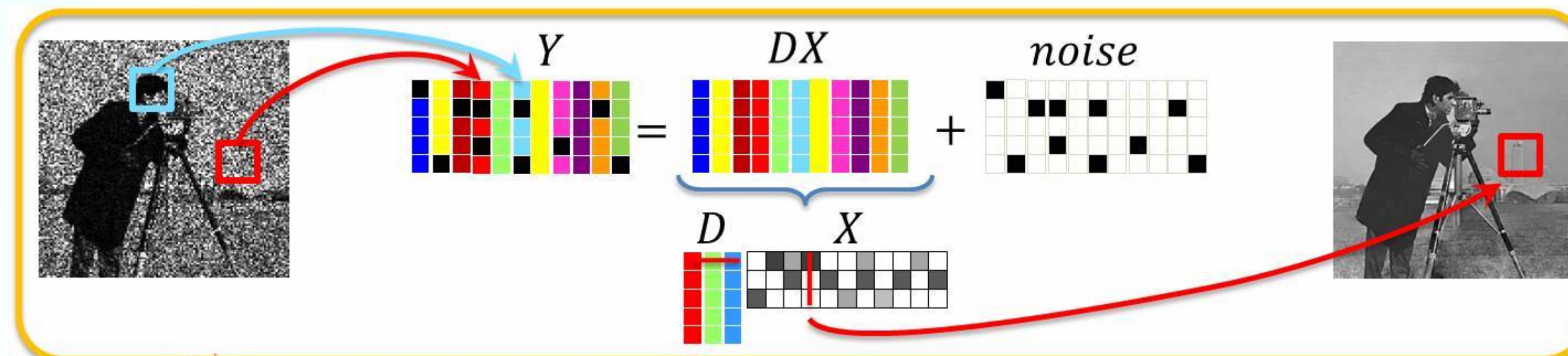
Why sparse matrix factorisation?

Fast linear operator: if $\mathbf{A} \approx \mathbf{W}_1 \dots \mathbf{W}_J$ then $\mathbf{A}\mathbf{x} = \mathbf{W}_1 \dots \mathbf{W}_J \mathbf{x}$

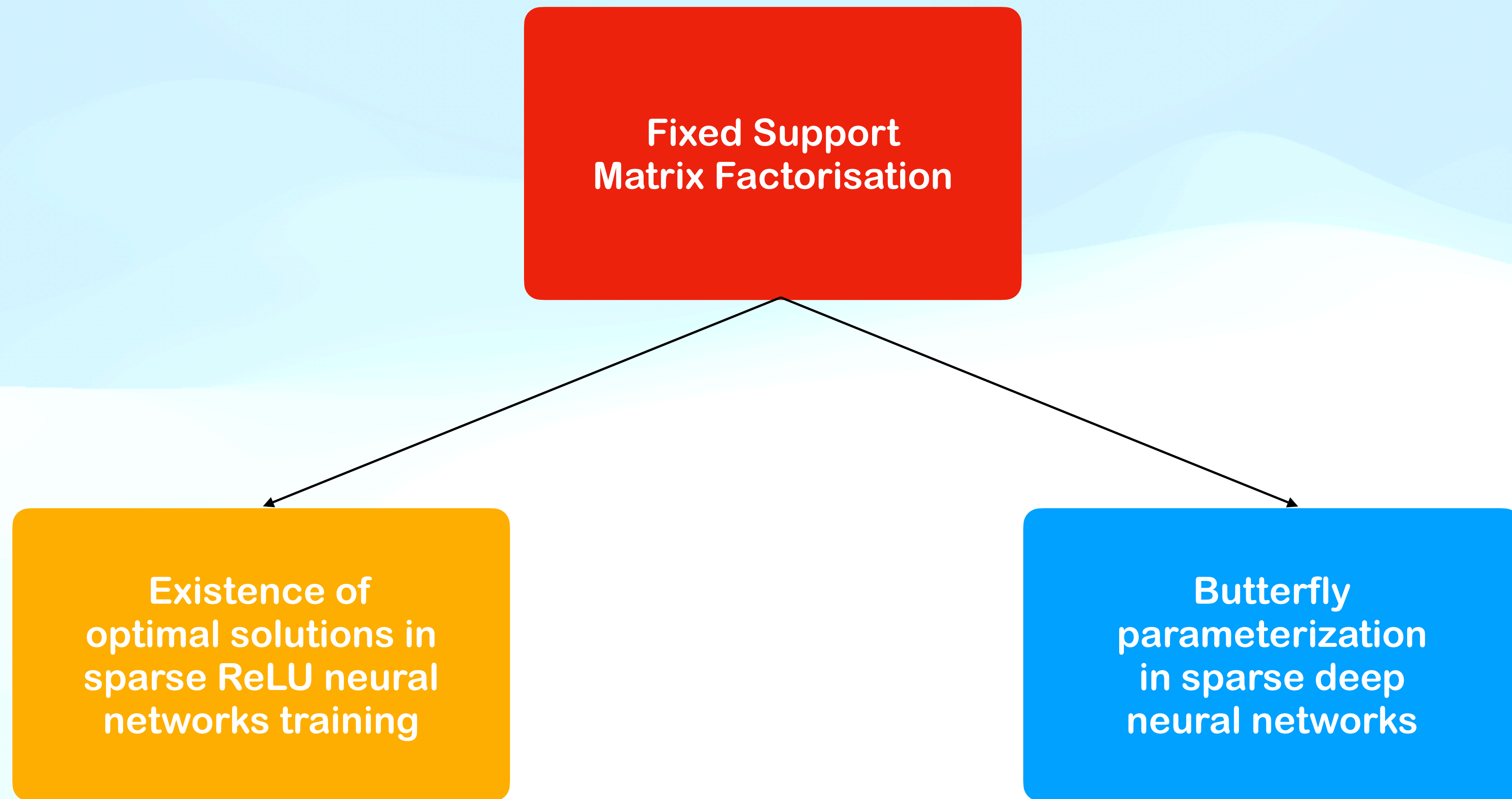


The Discrete Fourier Transformation sparse factorisation and its $O(n \log n)$ fast algorithm

Dictionary learning: given a dataset \mathbf{Y} , find atoms \mathbf{D} and look-up table \mathbf{X}



Plan of the talk





Fixed support matrix factorization

Results on (FSMF)

$$\min_{X,Y} L(X, Y) = \|A - XY^T\|_F^2 \text{ subject to: } \text{supp}(X) \subseteq I, \text{supp}(Y) \subseteq J$$

Existence of optimal solutions

👉 Does (FSMF) **always admit an optimal solution?**

Tractability

👉 Is (FSMF) **polynomially tractable?**

Landscape

👉 What does the **landscape** of $L(X, Y)$ look like? (e.g., does it have **local minima?**, etc.)

Non-existence of optimal solutions

(FSMF) does not always admit an optimal solution

$$A = \begin{matrix} & \overset{n}{\longleftarrow} \\ \begin{matrix} \uparrow \\ n \\ \downarrow \end{matrix} & \begin{matrix} \text{Grid 1} \\ \times \\ \text{Grid 2} \\ \downarrow \\ n \end{matrix} \end{matrix}$$

$$A = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

$$n = 2$$

👉 Infimum is **zero**: $X_k = \begin{pmatrix} \frac{1}{k} & 0 \\ 1 & -1 \end{pmatrix}$, $Y_k^\top = \begin{pmatrix} 1 & k \\ 0 & k \end{pmatrix}$, $\lim X_k Y_k^\top = A$.

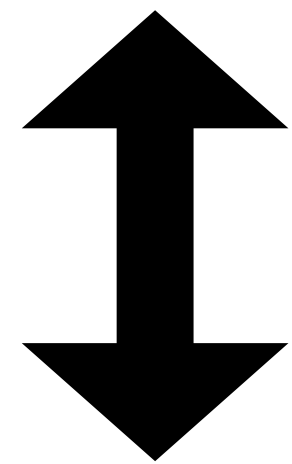
👉 Infimum is **not** attained: There is no feasible (X, Y) such that $A = XY^\top$.

(Gene H. Golub and Charles F. Van Loan, *Matrix Computations*)

Equivalence between existence - closedness

ORIGINAL
FORMULATION

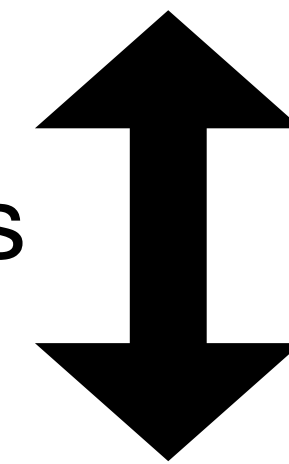
$$\min_{X,Y} \|A - XY^T\|_F^2 \text{ subject to: } \text{supp}(X) \subseteq I, \text{supp}(Y) \subseteq J$$



EQUIVALENT
FORMULATION

$$\min_{B \in \mathcal{P}_{I,J}} \|A - B\|_F^2 \text{ where } \mathcal{P}_{I,J} := \{XY^T \mid \text{supp}(X) \subseteq I, \text{supp}(Y) \subseteq J\}$$

Change of variables



PROJECTION OF A ONTO THE SET $\mathcal{P}_{I,J}$

Optimal solutions exist if and only if $\mathcal{P}_{I,J}$ is closed



Deciding the closedness of $\mathcal{P}_{I,J}$

Given (I, J) , decide the **closedness** of $\mathcal{P}_{I,J}$.

REMINDER: $\mathcal{P}_{I,J} := \{XY^T \mid \text{supp}(X) \subseteq I, \text{supp}(Y) \subseteq J\}$

$\mathcal{P}_{I,J}$ is
closed?



$\overline{\mathcal{P}_{I,J}} \setminus \mathcal{P}_{I,J}$
is empty?

closure of $\mathcal{P}_{I,J}$

semi-algebraic
set

☞ The emptyness of a semi-algebraic set is decidable (using Quantifier Elimination Algorithm).

(S. Basu, R. Pollack, M-F Roy, *Algorithms in Real Algebraic Geometry*)

☞ The complexity is doubly exponential (w.r.t. the sizes of I, J and the matrix).

NP-hardness

THEOREM 1

For arbitrary support constraint (I, J) , (FSMF) is NP-hard.

PROOF: Rank-one matrix completion is reducible to (FSMF).

$$\begin{array}{|c|c|c|c|c|} \hline 2 & ? & 1 & ? & ? \\ \hline ? & 2 & ? & 4 & 0 \\ \hline ? & 9 & ? & ? & 1 \\ \hline 5 & ? & 2 & 9 & ? \\ \hline \end{array} = \begin{array}{|c|} \hline \text{blue} \\ \hline \text{blue} \\ \hline \text{blue} \\ \hline \text{blue} \\ \hline \end{array} \times \begin{array}{|c|c|c|c|c|} \hline \text{blue} \\ \hline \text{blue} \\ \hline \text{blue} \\ \hline \text{blue} \\ \hline \text{blue} \\ \hline \end{array}$$

This problem is NP-hard.

(N.Gillis, F. Glineur, SIAM Journal on Matrix Analysis and Applications)

Tractability with structured supports

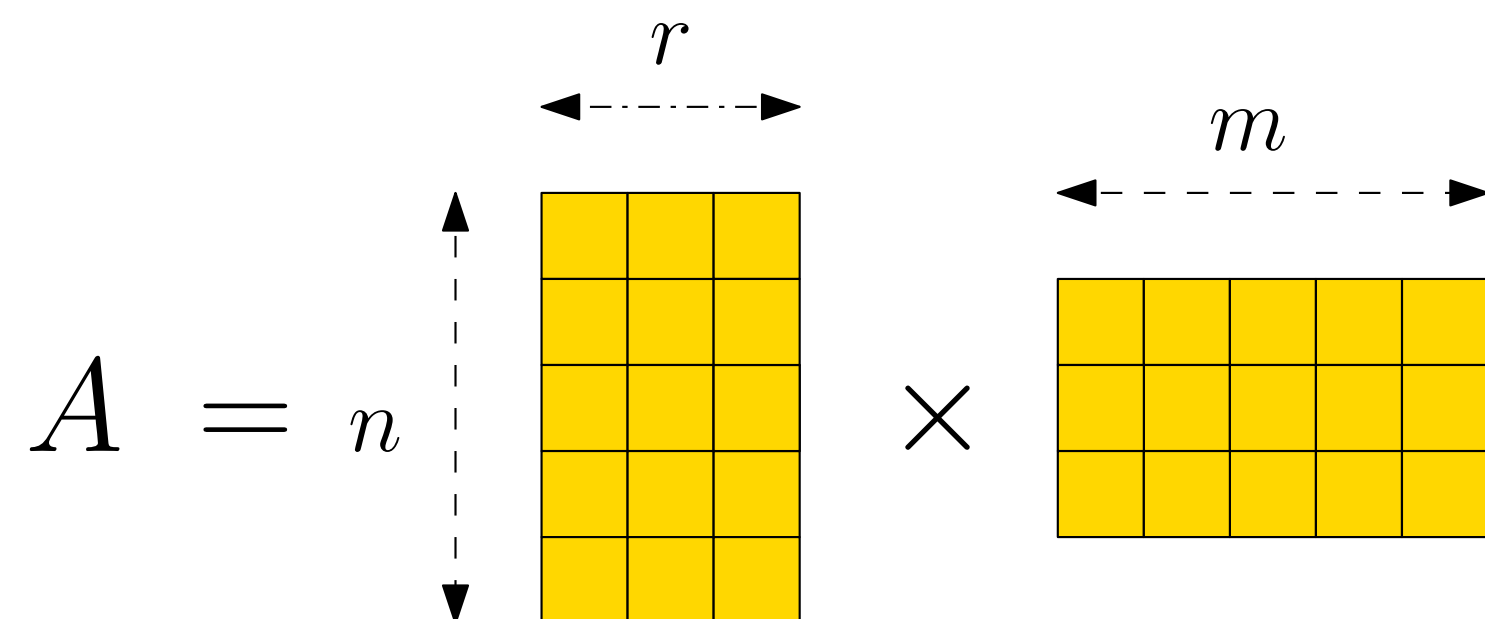
An example of tractable instances

Unconstrained Matrix Factorisation

When there is **no** constraint on the supports of (X, Y)

$$\text{Minimize } L(X, Y) = \|A - XY^T\|_F^2$$

$$X \in \mathbb{R}^{n \times r}, Y \in \mathbb{R}^{m \times r}$$



Best rank r approximation
of the matrix A .

(S. Burer, R. D.C. Monteiro, *Mathematical Programming*)

👉 Algorithm: Using (Truncated) Singular Value Decomposition.

? Can (Truncated) Singular Value Decomposition still work in constrained cases?

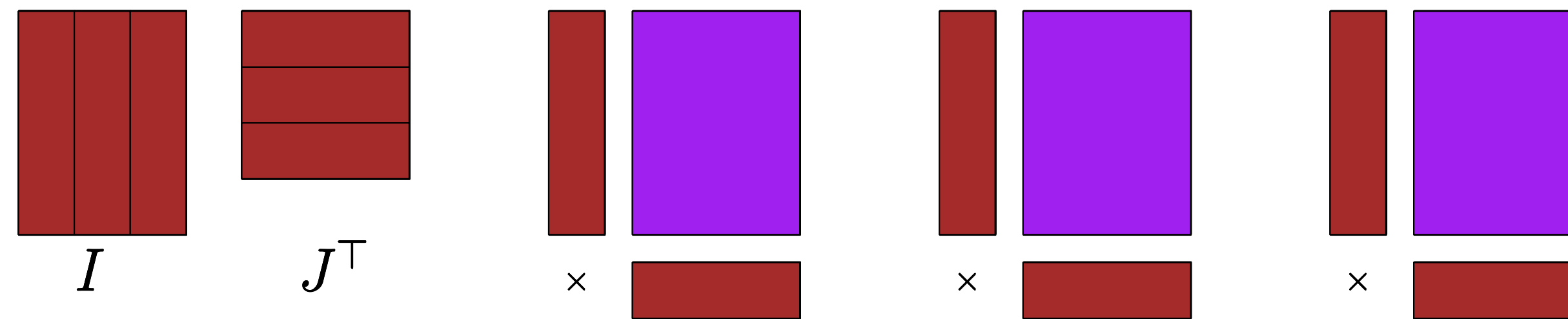
Tractability with structured supports (cont)

Rank one contribution supports

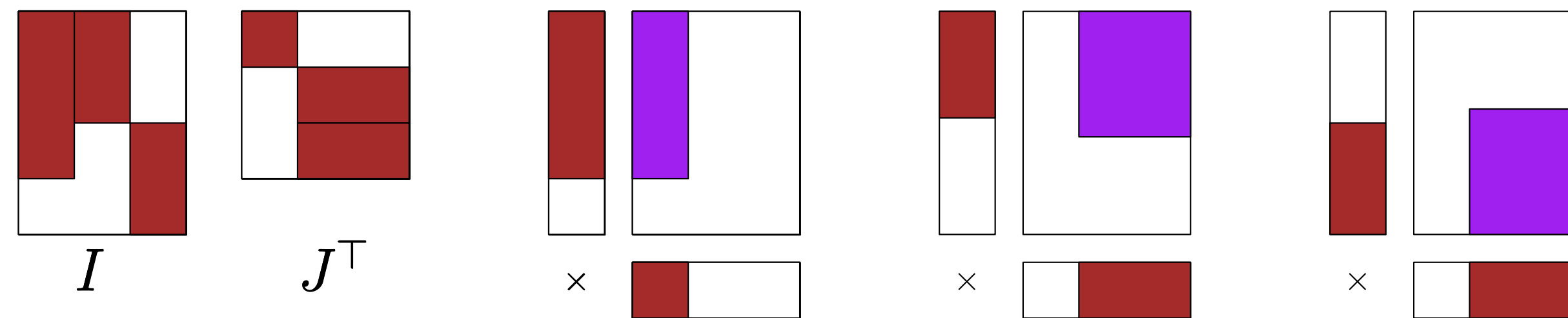
$$XY^T = \boxed{X_1 Y_1^T} + \boxed{X_2 Y_2^T} + \boxed{X_3 Y_3^T}$$

k th column of X rank at most 1 of Y

No support constraint



With support constraint



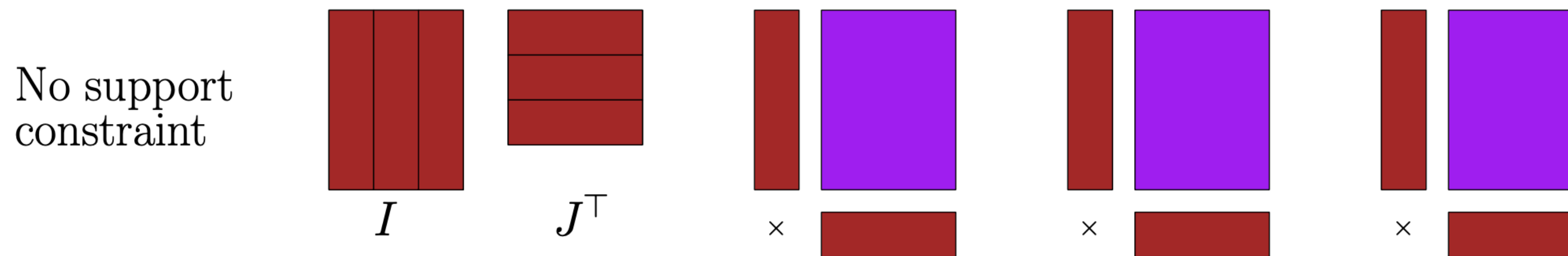
zeros

rank-one support

Tractability with structured supports (cont)

What is special about unconstrained matrix factorization?

$$XY^{\top} = X_1Y_1^{\top} + X_2Y_2^{\top} + X_3Y_3^{\top}$$



All rank-one supports are **identical**

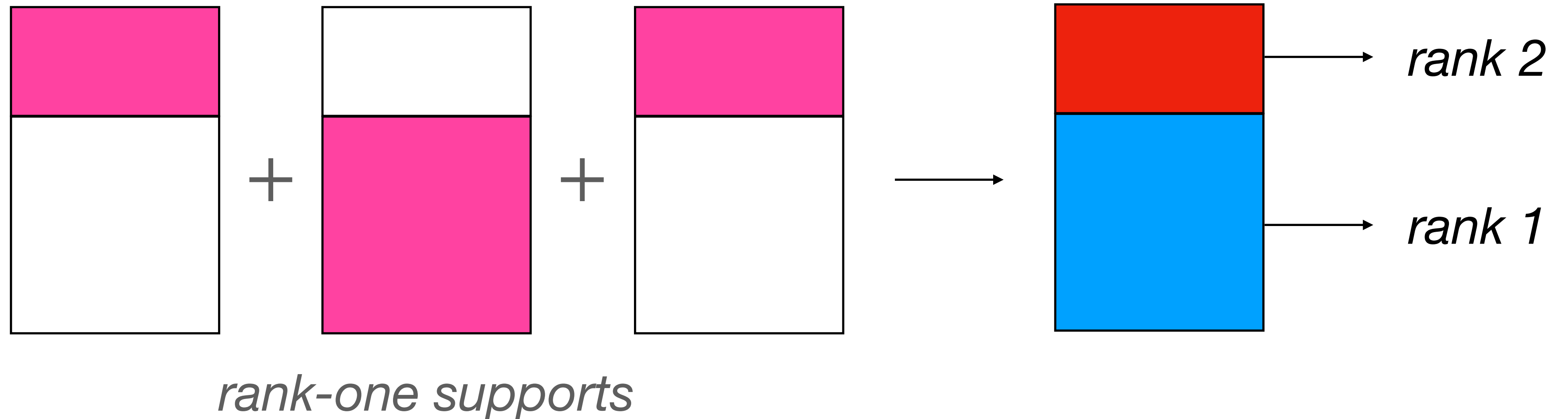
THEOREM II

If all rank-one supports are **pairwise disjoint** or **identical**, then the corresponding instance of (FSMF) is polynomially tractable.

Tractability with structured supports (cont)

THEOREM II

If all rank-one supports are **pairwise disjoint** or **identical**, then the corresponding instance of (FSMF) is polynomially tractable.



👉 Algorithm: Using (Truncated) Singular Value Decomposition for **submatrices** of the target matrix.

Litteratures on the landscape of $L(X, Y)$

$$L(X, Y) = \|A - XY^T\|_F^2$$

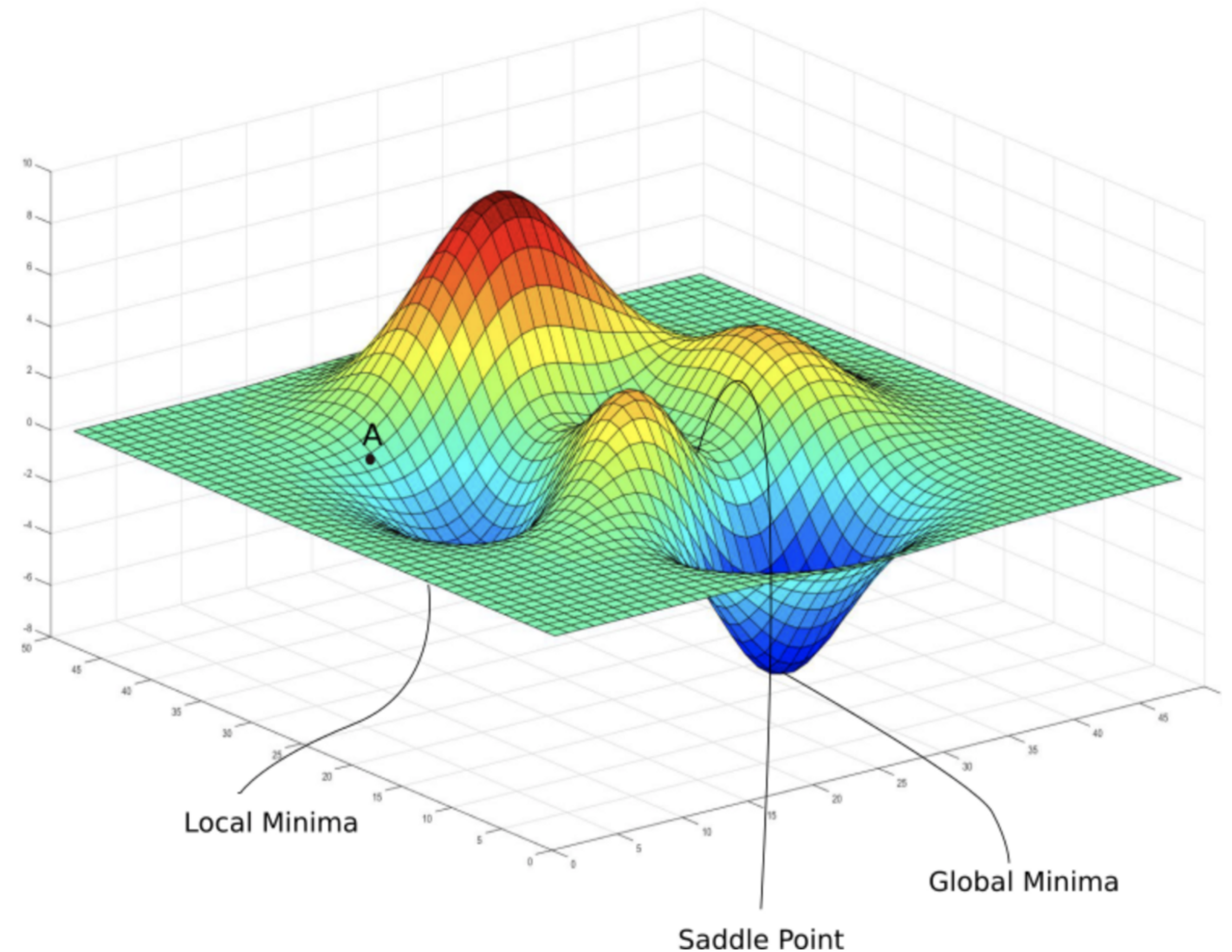
Has been studied for:

- 👉 Linear and shallow neural networks
- 👉 Matrix sensing, phase retrieval, matrix completion.

(Q. Li, Z. Zhu, G. Tang, The nonconvex geometry of low-rank matrix optimization, Information and Inference, 2018)

(Z. Zhu, D. Soudry, Y.C. Eldar, M.B. Wakin, The global optimization geometry of shallow linear neural networks, JMIV, 2019)

(L. Venturi, A. S. Bandeira, J. Bruna, Spurious valleys in one-hidden-layer neural network optimization landscapes, JMLR, 2019)

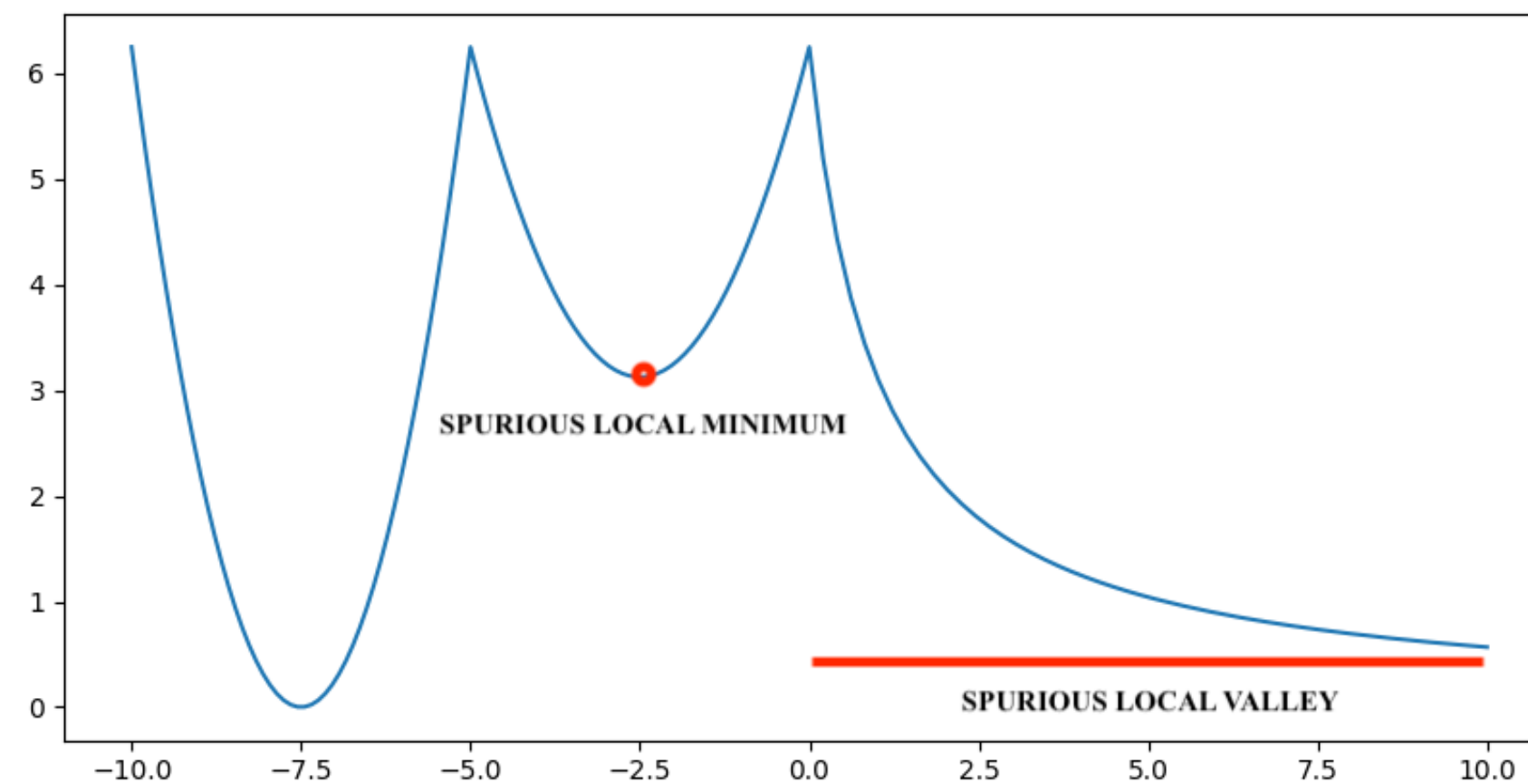


? What does the landscape look like in the constrained cases?

Unconstrained matrix factorization

$$L(X, Y) = \|A - XY^T\|_F^2$$

- 👉 There is no spurious local minimum for any A .
- 👉 There is no spurious local valley for any A .



? Do these properties still hold in constrained cases?

Benign landscape of tractable instances

Reminder: Fixed Support Matrix Factorization

$$\min_{X,Y} \|A - XY^T\|_F^2 \text{ subject to: } \text{supp}(X) \subseteq I, \text{supp}(Y) \subseteq J$$

THEOREM III

If (I, J) satisfies the condition of **Theorem II**, then there is no spurious local minima and spurious local valleys.

Summary on (FSMF)

Existence of optimal solutions

- There are instances (A, I, J) which (FSMF) admits no optimal solution.

NP-hardness

- For arbitrary (I, J) , (FSMF) is *NP-hard* to solve.

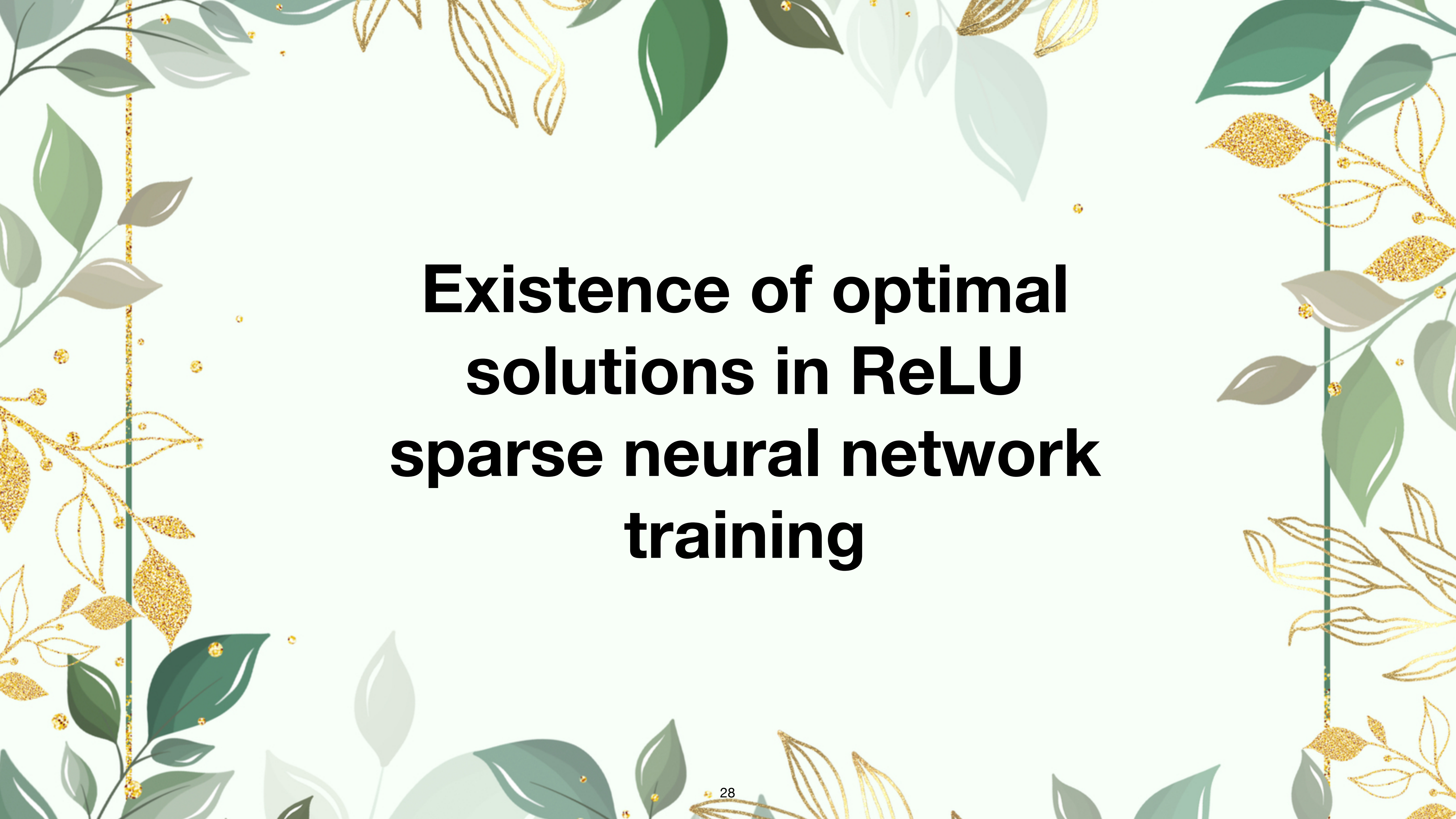
Tractability

- For certain structured (I, J) , (FSMF) has a polynomial algorithm.

Benign landscape

- With the same family of structured (I, J) , loss function of (FSMF) has no local minima.

(Q-T. Le, E. Riccietti, R. Gribonval, *SIAM Journal of Matrix Analysis and Applications*, 2023)



Existence of optimal solutions in ReLU sparse neural network training

Sparse ReLU Neural Networks Training

Optimization problem of Sparse Neural Networks

Given data set $\mathcal{D} := (X, Y)$ and \mathcal{E}_j some sets of *sparse* matrices, solve:

$$\text{Minimize}_{W^{(j)}, b^{(j)}} \quad \|Y - W^{(N)} \sigma(\dots \sigma(W^{(1)}X + b^{(1)}) + \dots) + b^{(N)}\|_F^2$$

$$\text{subject to:} \quad W^{(j)} \in \mathcal{E}_j, \forall j \in \{1, \dots, N\}$$

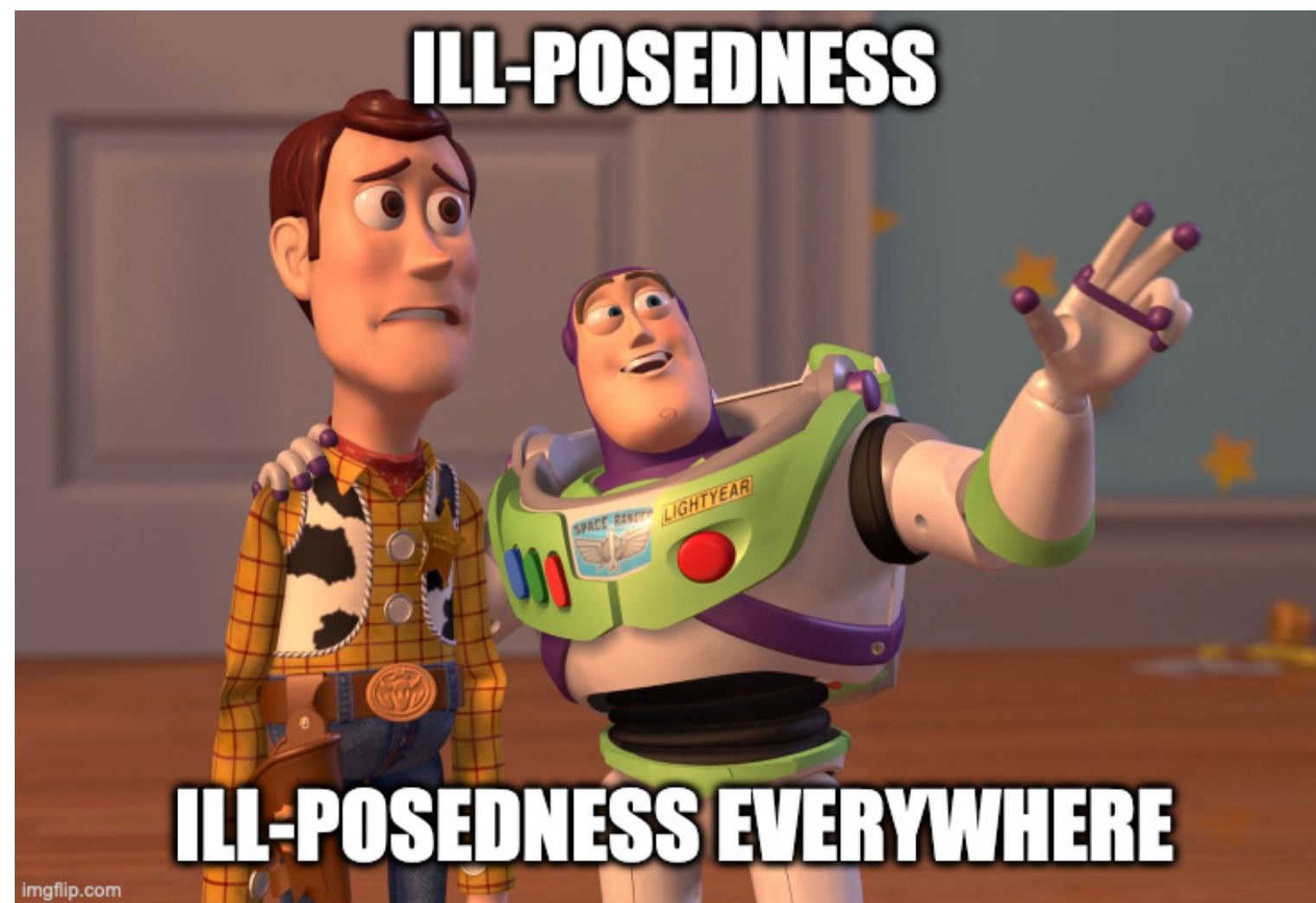
👉 σ is the **ReLU** activation function: $\sigma(x) = \max(x, 0)$.

👉 In practice, \mathcal{E}_j is usually chosen as the set of **k-sparse matrices**.

(J. Frankle, M. Carbin, ICLR 2019), (S. Han, H. Mao, W-J. Dally, ICLR 2016)

👉 We consider **quadratic loss function** for simplification. Our argument works for any **coercive** loss function.

Non-existence of optimal solutions - ill-posedness



Tensor decomposition (order at least three)	<p>TENSOR RANK AND THE ILL-POSEDNESS OF THE BEST LOW-RANK APPROXIMATION PROBLEM</p> <p>VIN DE SILVA* AND LEK-HENG LIM†</p>
Matrix Completion	<p>Low-Rank Matrix Approximation with Weights or Missing Data is NP-hard</p> <p>Nicolas Gillis¹ and François Glineur¹</p>
Robust Principle Component Analysis	<p>Matrix rigidity and the ill-posedness of Robust PCA and matrix completion*</p> <p>Jared Tanner^{†‡} Andrew Thompson[§] Simon Vary[†]</p>
(Classical) Neural Network Training	<p>Best k-Layer Neural Network Approximations</p> <p>Lek-Heng Lim¹ · Mateusz Michałek^{2,3} · Yang Qi⁴</p>

? How about the training problem of sparse ReLU neural networks?

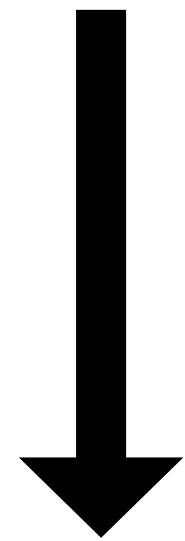
Fixed support sparse ReLU neural networks

Given data set $\mathcal{D} := (X, Y)$, solve:

GENERAL

$$\min_{W^{(j)}, b^{(j)}} \|Y - W^{(L)} \sigma(\dots \sigma(W^{(1)}X + b^{(1)}) + \dots) + b^{(L)}\|_F^2$$

subject to: $W^{(j)} \in \mathcal{E}_j, \forall j \in \{1, \dots, N\}$



FIXED SUPPORT

$$\min_{W^{(j)}, b^{(j)}} \|Y - W^{(N)} \sigma(\dots \sigma(W^{(1)}X + b^{(1)}) + \dots) + b^{(N)}\|_F^2$$

subject to: $\text{supp}(W^{(j)}) \in I_j, \forall j \in \{1, \dots, N\}$

DÉJÀ VU: closedness vs existence of optimal solutions



Given a support constraint (I_1, \dots, I_N) , **does** optimal solutions always exist **for all data set** \mathcal{D} for the corresponding training problem ?

The support constraint (I_1, \dots, I_N) makes the training problem **always admit optimal solutions** if and only if **for all** input sets X , the image of the function $\theta := \{(W^i, b^i)\} \mapsto W^{(N)}\sigma(\dots\sigma(W^{(1)}X + b^{(1)}) + \dots) + b^{(N)}$ is **closed**.



Sufficient condition for the existence of optimal solutions

THEOREM IV

For **two-layer** neural networks ($N = 2$) with **output dimension** equal to **one**, any support constraint makes the training problem **always admit optimal solutions**.

(Q-T. Le, E. Riccietti, R. Gribonval, preprint, 2023)

COROLLARY I

For **two-layer** neural networks ($N = 2$) with **output dimension** equal to **one**, the constraints $\mathcal{E}_j := \{X \mid \|X\|_0 \leq k_j\}, j = 1, 2$ makes the training problem **always admit optimal solutions**.

Necessary condition for the existence of optimal solutions

THEOREM V

For **two-layer** neural networks ($N = 2$) with support constraint (I, J) , if the training problem always admits optimal solutions, then $\mathcal{P}_{I,J}$ is **closed**.

(Q-T. Le, E. Riccietti, R. Gribonval, preprint, 2023)

$$\mathcal{P}_{I,J} := \{XY^T \mid \text{supp}(X) \subseteq I, \text{supp}(Y) \subseteq J\}$$

this is decidable



THEOREM VI

For fixed support neural networks with support constraint (I_1, \dots, I_N) , if the training problem always admits optimal solutions, then $\mathcal{P}_{I_1, \dots, I_N}$ is **closed**.

Necessary condition for the existence of optimal solutions

THEOREM

For **two-layer** neural networks ($N = 2$) with support constraint (I, J) , if the training problem always admits optimal solutions, then $\mathcal{P}_{I,J}$ is **closed**.

The condition is just necessary because when there is **no constraint** on the support, the training problem is ill-posed for certain data set.

(L-H. Lim, M. Michalek, Y. Qi, *Constructive Approximation* 2019)



Butterfly parameterization in sparse neural networks

Introduction to butterfly parameterization

Discrete Fourier Transform (DFT): $y_k = \sum_{n=0}^{N-1} e^{-\frac{i2\pi}{N}kn} x_n, \quad k = 0, \dots, N-1$

$$y = F_N x \quad F_N = \left(e^{-\frac{i2\pi}{N}kn} \right)_{k,n \in \{0, \dots, N-1\}}$$

Algebraic properties of the Fourier kernel:

$$y_k = \sum_{m=0}^{N/2-1} e^{-\frac{i2\pi}{N/2}km} x_{2m} + e^{-\frac{i2\pi}{N}k} \sum_{m=0}^{N/2-1} e^{-\frac{i2\pi}{N/2}km} x_{2m+1}$$

$$y_{k+N/2} = \sum_{m=0}^{N/2-1} e^{-\frac{i2\pi}{N/2}km} x_{2m} - e^{-\frac{i2\pi}{N}k} \sum_{m=0}^{N/2-1} e^{-\frac{i2\pi}{N/2}km} x_{2m+1}$$

DFT on **even** indices

DFT on **odd** indices

The Cooley - Tukey algorithm (radix-2)

Matrix factorization of DFT

$$\begin{aligned}
 F_N x &= \begin{pmatrix} F_{N/2} x_e + D_{N/2} F_{N/2} x_o \\ F_{N/2} x_e - D_{N/2} F_{N/2} x_o \end{pmatrix} \\
 &= \begin{pmatrix} I_{N/2} & D_{N/2} \\ I_{N/2} & -D_{N/2} \end{pmatrix} \begin{pmatrix} F_{N/2} & \mathbf{0} \\ \mathbf{0} & F_{N/2} \end{pmatrix} \boxed{P_N} x
 \end{aligned}$$

$D_{N/2} = \begin{pmatrix} 1 & & & \\ & e^{-\frac{i2\pi}{N}} & & \\ & & e^{-\frac{i2\pi}{N}2} & \\ & & & \dots \\ & & & & e^{-\frac{i2\pi}{N}(\frac{N}{2}-1)} \end{pmatrix}$

Unrolling the recursion:

Permutation matrix

$$\begin{aligned}
 F_N &= B_N \begin{pmatrix} F_{N/2} & \mathbf{0} \\ \mathbf{0} & F_{N/2} \end{pmatrix} P_N \\
 &= B_N \begin{pmatrix} B_{N/2} & \mathbf{0} \\ \mathbf{0} & B_{N/2} \end{pmatrix} \begin{pmatrix} B_{N/4} & & & \\ & B_{N/4} & & \\ & & B_{N/4} & \\ & & & B_{N/4} \end{pmatrix} \begin{pmatrix} P_{N/2} & \mathbf{0} \\ \mathbf{0} & P_{N/2} \end{pmatrix} P_N \\
 &= \dots
 \end{aligned}$$

Matrix factorization of DFT

Assume that $N = 2^L$:

$$F_N = X^{(1)} \dots X^{(L)} P$$

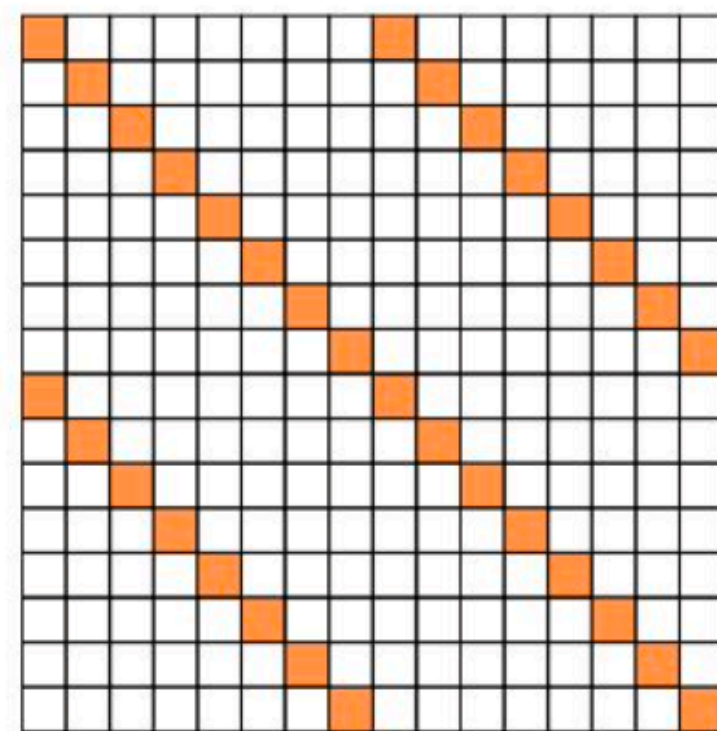
$$\text{supp}(X^{(\ell)}) \subseteq \mathbf{I}_{2^{\ell-1}} \otimes \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} \otimes \mathbf{I}_{\frac{N}{2^\ell}}$$

support constraint in binary matrix form

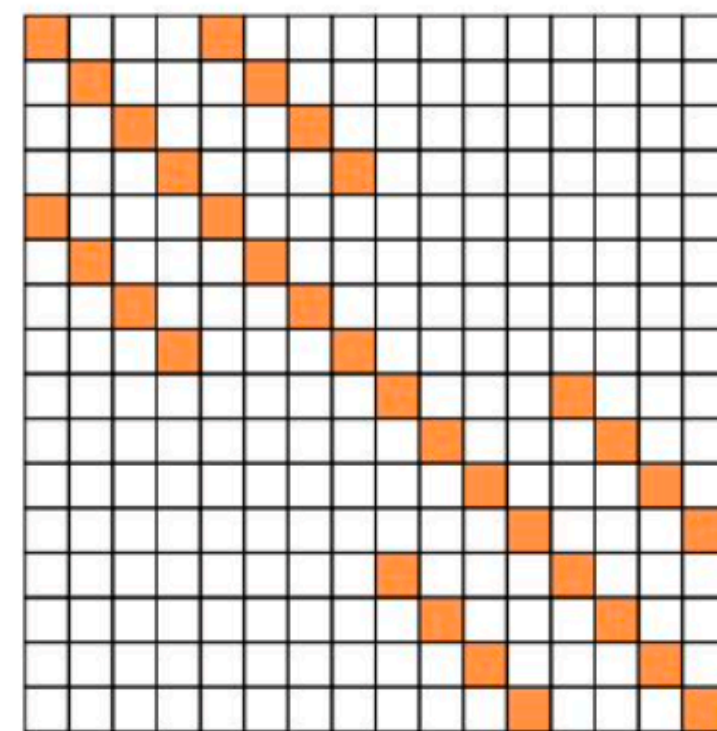
$O(N^2)$

$O(N \log N)$

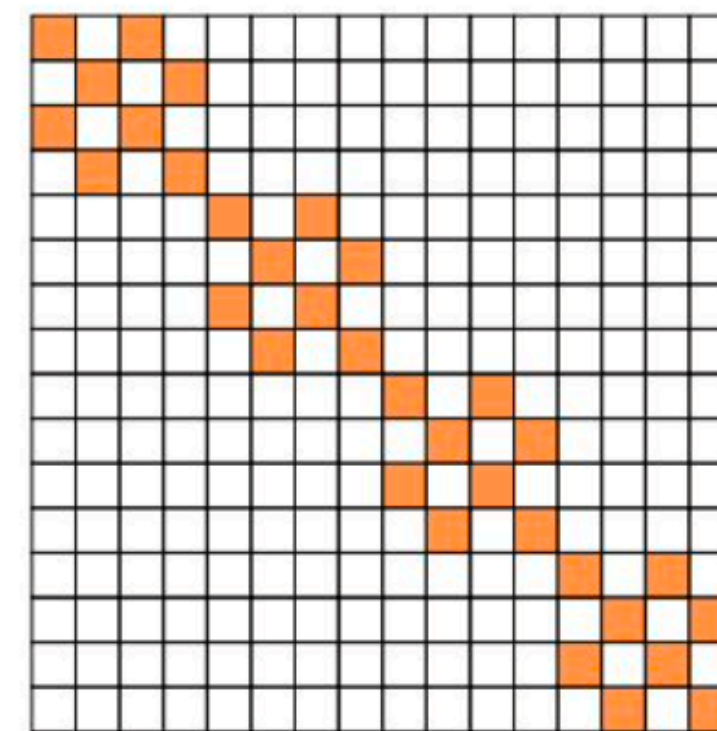
For $N = 16$:



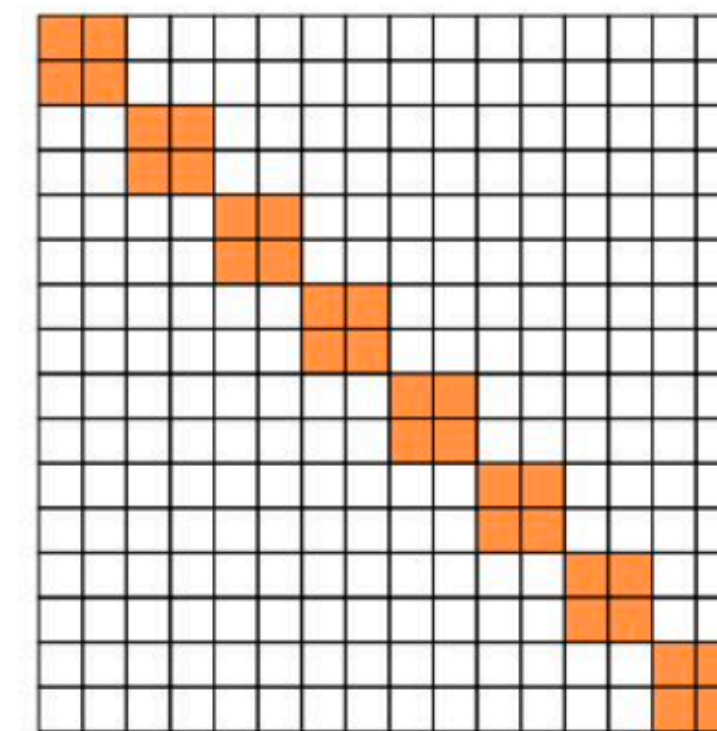
$\text{supp}(X^{(1)})$



$\text{supp}(X^{(2)})$




$\text{supp}(X^{(3)})$



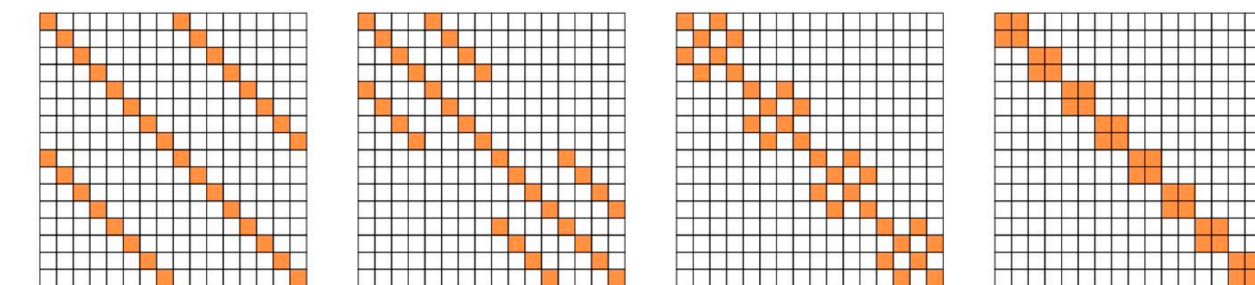
$\text{supp}(X^{(4)})$

Butterfly parameterization in Sparse Neural Networks

Classical neural networks \rightarrow

 $x + b$
 $O(N^2)$

Butterfly sparse neural networks

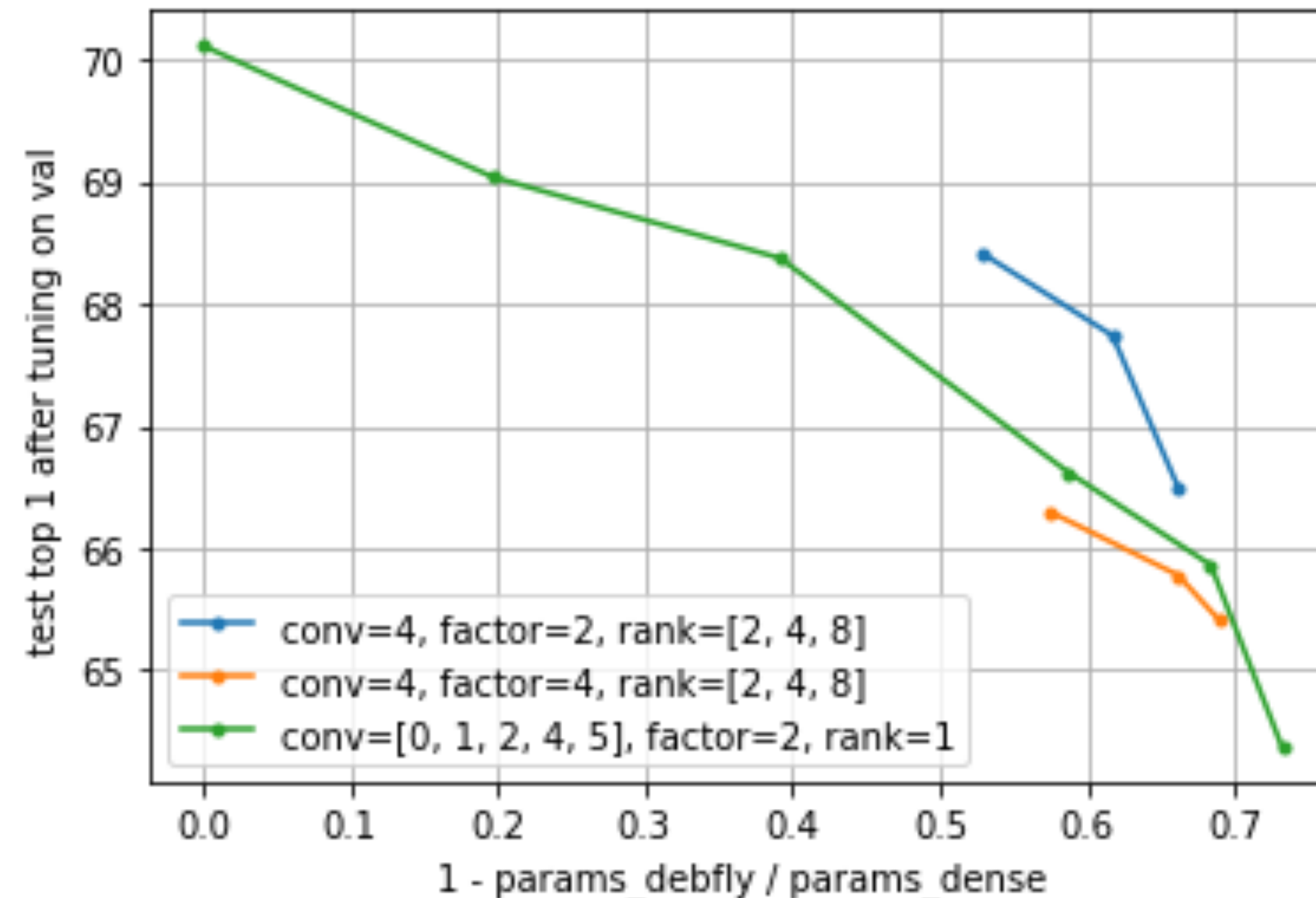
 $x + b$
 $O(N \log N)$

Parameterization	Number of factors	Matrix size	Introduced by
Butterfly	L	$2^L \times 2^L$	T. Dao et. al., 2019
Kaleidoscope	2L	$2^L \times 2^L$	T. Dao et. al., 2020
Monarch	2	$m \times n$	T. Dao et. al., 2022
Deformable butterfly	flexible	$m \times n$	R. Lin et. al., 2022

!! Supports of factors are **fixed, sparse** and very **structured**.

Interpretation of butterfly parameterization

? Among all existing parameterization, which one should we choose?



Trade-off between performance and compression

Approximation a matrix by butterfly parameterization

$$\min_{W^{(1)}, \dots, W^{(L)}} \|A - W^{(1)} \dots W^{(L)}\|_F \quad \text{s.t. } W^{(\ell)} \text{ is butterfly} \quad (1)$$

👉 Generalized version of (FSMF) with *structured* supports.

👉 Existing algorithm: hierarchical factorization - **butterfly algorithm**.

(Michielssen & Boag, 1996); (O'Neil, Woolfe & Rokhlin, 2010); (Liu et. al. 2021)

⚠️ No theoretical guarantee yet.

!! Hypothesis class of matrix:

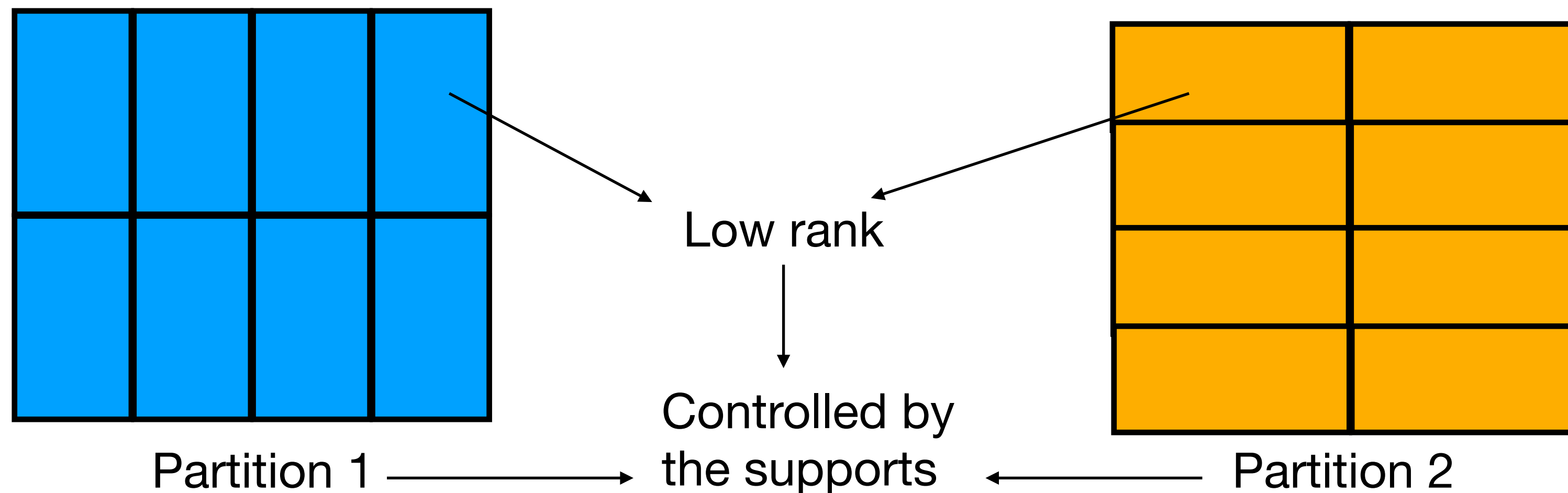
$$\mathcal{B} := \{ W^{(1)} \dots W^{(L)} \mid \text{Infimum of (1)} = 0 \}$$

Analysis of butterfly parameterization

THEOREM VII

If E^* is the best error approximation of (1), the butterfly algorithm yields a solution whose distance to A is smaller than $(2^{L-1} - 1)E^*$

Algebraic description of \mathcal{B} :



👉 Also known as *complementary low-rank matrices* in the literature.

Contribution and future works

TAKE AWAY MESSAGE

- Link between sparse matrix factorization and its variant (FSMF) with sparse ReLU neural networks.
- Necessary/Sufficient condition for the existence of optimal solutions sparse ReLU neural networks.
- Butterfly parameterization in sparse deep neural networks

POSSIBLE IMPROVEMENT?

- Better algorithms to decide the ill-posedness of (FSMF).
- A full characterization of ill-posedness of sparse ReLU neural networks.

<https://faust.inria.fr/>

<https://arxiv.org/abs/2112.00386>

<https://arxiv.org/abs/2306.02666>

THANK YOU

Analysis of butterfly parameterization

 The supports of all existing factors have the form:

$$\text{supp}(W^{(\ell)}) \subseteq \mathbf{I}_a \otimes \mathbf{1}_{b \times c} \otimes \mathbf{I}_d$$

EXAMPLE:

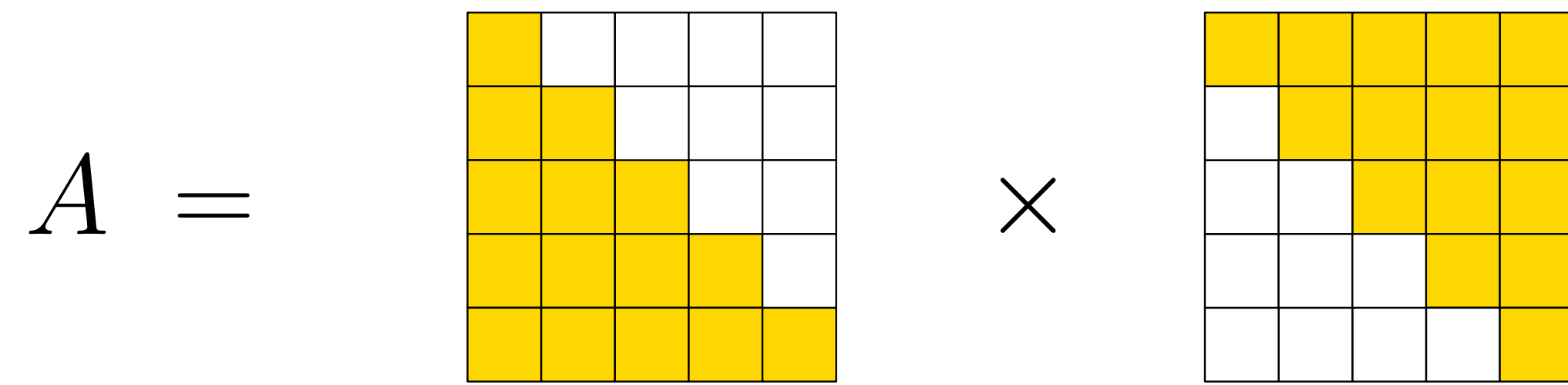
Parameterization	Support forms
Butterfly	$\mathbf{I}_{2^{\ell-1}} \otimes \mathbf{1}_{2 \times 2} \otimes \mathbf{I}_{\frac{N}{2^\ell}}$
Kaleidoscope	
Monarch	$\mathbf{1}_{a \times b} \otimes \mathbf{I}_c$ and $\mathbf{I}_b \otimes \mathbf{1}_{c \times d}$
Deformable butterfly	$\mathbf{I}_a \otimes \mathbf{1}_{b \times c} \otimes \mathbf{I}_d$

 The product of two consecutive factors remains *butterfly* .

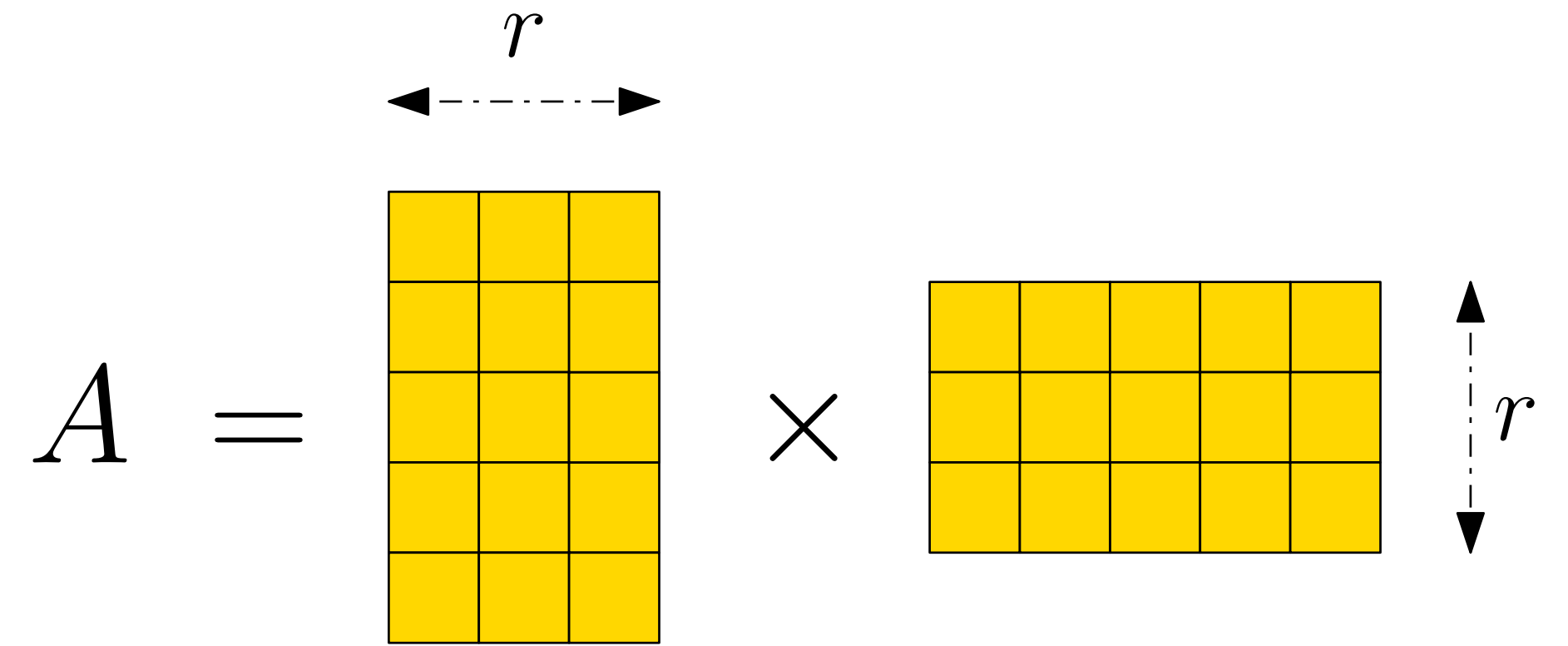
$$\text{supp}(W^{(\ell)} W^{(\ell+1)}) \subseteq \mathbf{I}_{a'} \otimes \mathbf{1}_{b' \times c'} \otimes \mathbf{I}_{d'}$$

 This does not include the **Kaleidoscope** parameterization.

And Fixed Support Matrix Factorization



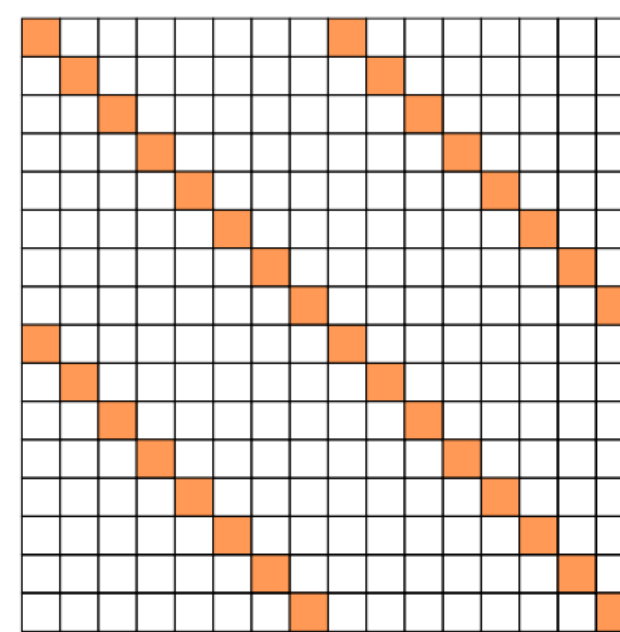
LU decomposition



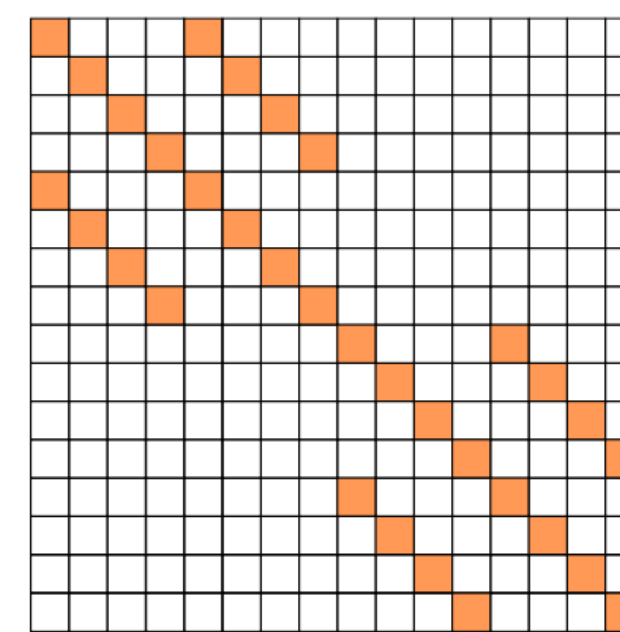
Low rank approximation



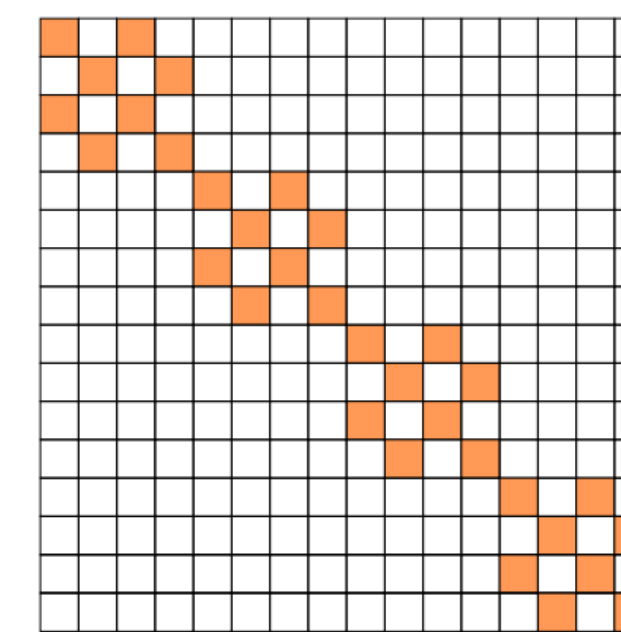
Hierarchical matrix



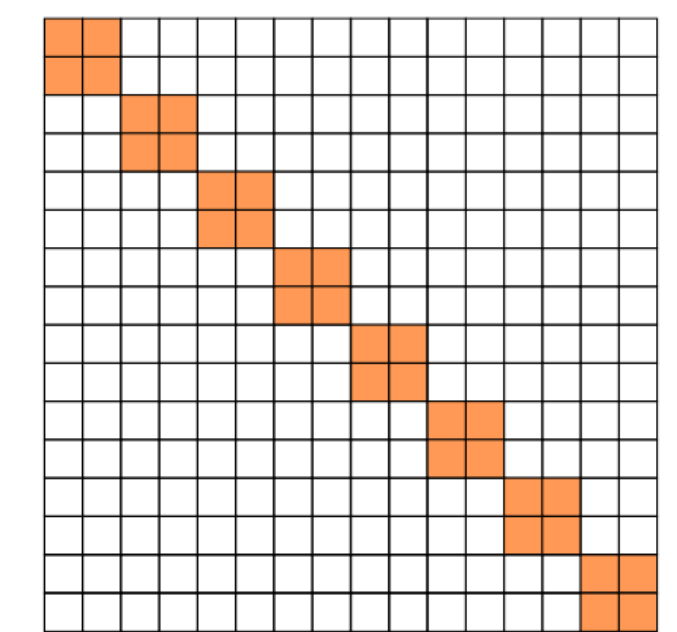
(a) $S_{bf}^{(4)}$



(b) $S_{bf}^{(3)}$



(c) $S_{bf}^{(2)}$



(d) $S_{bf}^{(1)}$

Butterfly matrix factorization

Question for sparse neural networks

Problem formulation

Feed forward networks:

$$\theta = \{(\mathbf{W}_i, \mathbf{b}_i) \mid i = 1, \dots, N\}$$

$$f(x; \theta) = \mathbf{W}_N \sigma(\dots \sigma(\mathbf{W}_1 x + \mathbf{b}_1)) + \mathbf{b}_N$$

Training:

Conventional Deep Neural Networks	Sparse Deep Neural Networks
Minimize: $\mathcal{L}_\theta := \sum_{i=1}^n L(f(\theta, x_i), y_i)$	Minimize: $\mathcal{L}_\theta := \sum_{i=1}^n L(f(\theta, x_i), y_i)$ such that: \mathbf{W}_i are sparse matrices

Existing algorithms / approaches for Sparse Deep Neural Networks training:

- Pruning & Retraining, Lottery Ticket Hypothesis (Han et al., IPL 2015), (Zhu et al., 2017), (Jonathan et al., 2019)
- Regularisation l_0 or l_1 (Bengio et al., 2013), (Yu et al., 2017), (Collins et al., 2014), (Liu et al., 2015)
- Bayesian/ Variational approaches (Neklyudov et al., 2017), (Ullrich et al., 2017), (Louizos et al., 2017)