

# Structured Adaptive and Random Spinners for Fast Machine Learning Computations

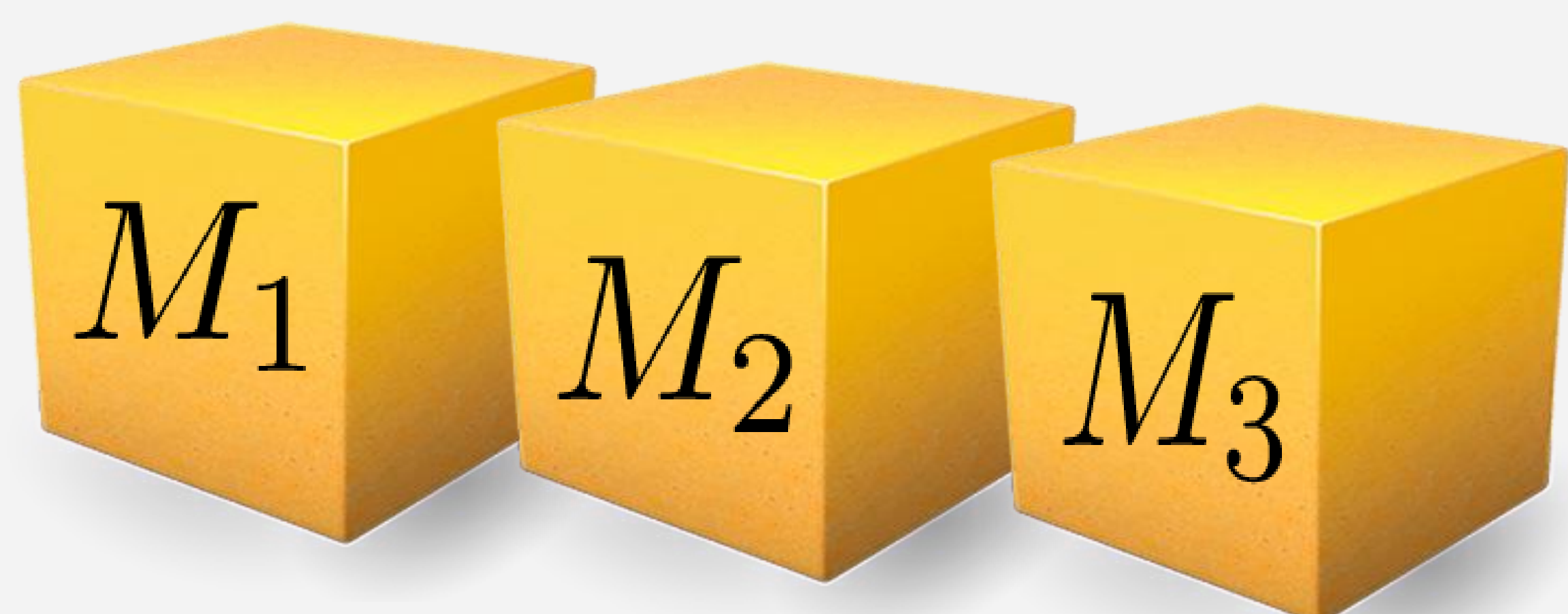
Mariusz Bojarski, Anna Choromanska, Krzysztof Choromanski, Francois Fagan, Cedric Gouy-Pailler, Anne Morvan, Nourhan Sakr, Tamas Sarlos, Jamal Atif

## Why structured projections ?

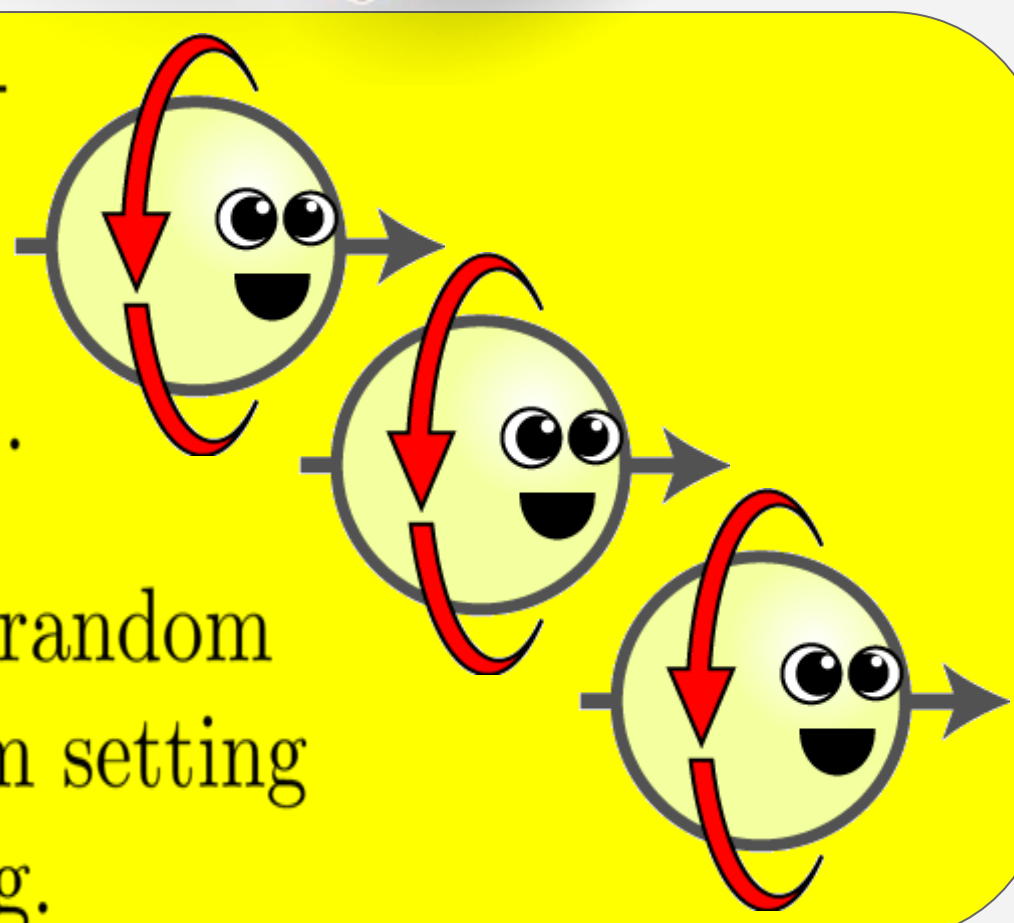
- Replace unstructured matrices in various algorithms
  - Kernel methods via random feature maps
  - Neural networks and dimensionality reduction techniques
  - Cross-polytope LSH methods and convex optimization
  - Random projection trees
  - Quasi-Monte Carlo techniques
  - Advantages:** Speed-ups. Storage compression. Almost no loss of accuracy.
- Provide tradeoff between required accuracy level and computational time/storage complexity



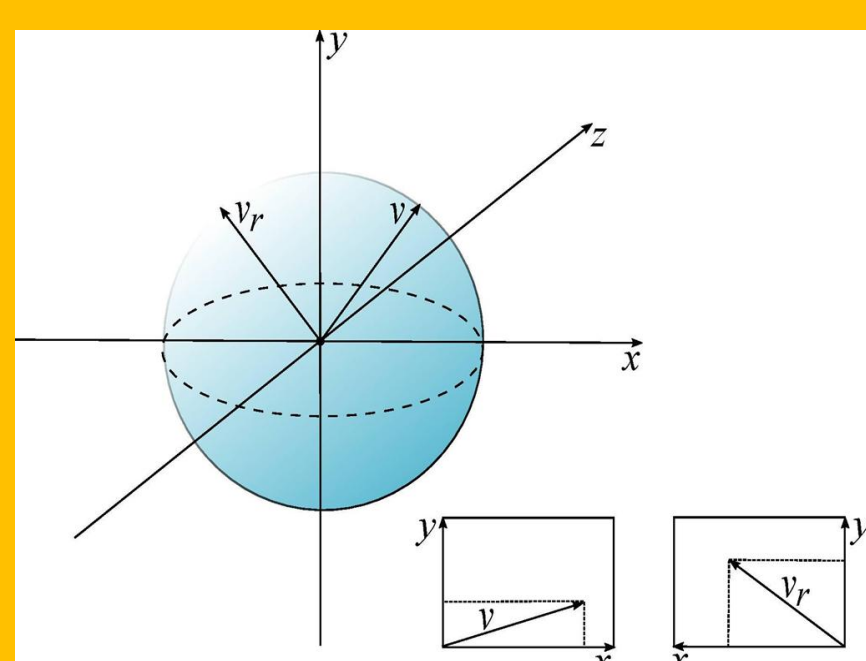
## Structured Spinners



- **Balanceness:**  $M_1$  and  $M_2 M_1$  are  $(\delta(n), p(n))$ -balanced isometries.
- **Decorrelation:**  $M_2 = V(W^1, \dots, W^n) D_{\rho_1, \dots, \rho_n}$ .
- **Budget:**  $M_3 = C(r, n)$  for  $r \in \mathbb{R}^k$ , where  $r$  is random Rademacher/Gaussian in the random setting and is learned in the adaptive setting.



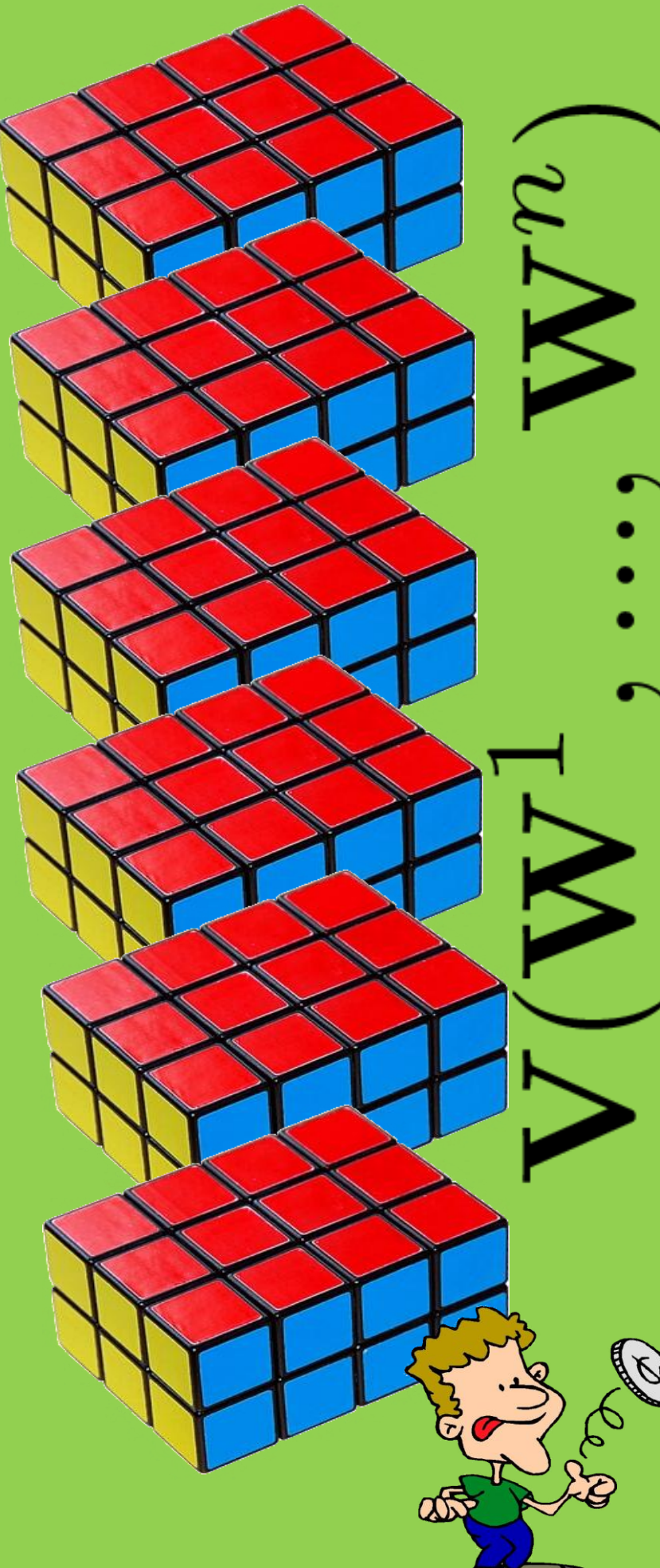
### ➤ Balanceness



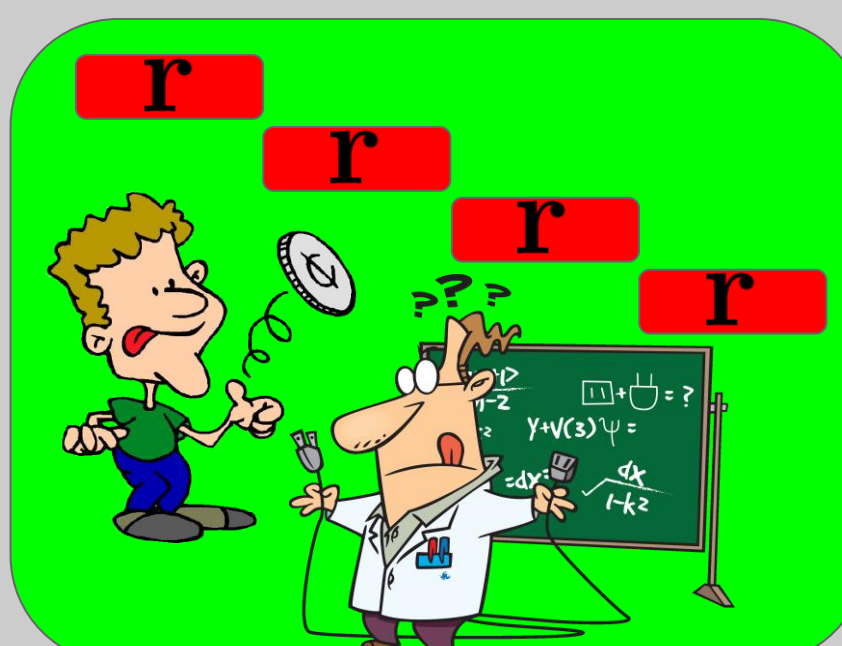
1	1	1	1	1	1	1	1
1	1	1	1	-1	-1	-1	-1
1	1	-1	-1	-1	-1	1	1
1	1	-1	-1	1	1	-1	-1
1	-1	-1	1	1	-1	-1	1
1	-1	-1	1	-1	1	1	-1
1	-1	1	-1	-1	1	-1	1
1	-1	1	-1	1	-1	1	-1



### ➤ Decorrelation



### ➤ Budget or randomness/learnable parameters



- Defines the capacity of the model
- From quadratic with no computational and storage gains to linear with  $O(n \log(n))$  time complexity

## Smooth Sets of Matrices

**Definition 1.**  $((\Delta_F, \Delta_2)$ -smooth sets):

A deterministic set of matrices  $W^1, \dots, W^n \in \mathbb{R}^{k \times n}$  is  $(\Delta_F, \Delta_2)$ -smooth if:

- $\|W^i\|_2 = \dots = \|W^n\|_2$  for  $i = 1, \dots, n$ , where  $W^j$  is the  $j^{\text{th}}$  column of  $W^i$ ,
- for  $i \neq j$  and  $l = 1, \dots, n$  we have:  $(W^i)^T \cdot W^j = 0$ ,
- $\max_{i,j} \| (W^j)^T W^i \|_F \leq \Delta_F$  and  $\max_{i,j} \| (W^j)^T W^i \|_2 \leq \Delta_2$ .

Permutation matrices

$$A_{\text{circ}} = \begin{pmatrix} g_0 & g_1 & g_2 & g_3 & g_4 \\ g_4 & g_0 & g_1 & g_2 & g_3 \\ g_3 & g_4 & g_0 & g_1 & g_2 \\ g_2 & g_3 & g_4 & g_0 & g_1 \\ g_1 & g_2 & g_3 & g_4 & g_0 \end{pmatrix}$$

$$P_1 = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

$$P_2 = \begin{pmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 \end{pmatrix}$$

$$P_3 = \begin{pmatrix} 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \end{pmatrix}$$

$$P_4 = \begin{pmatrix} 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \end{pmatrix}$$

$$P_5 = \begin{pmatrix} 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{pmatrix}$$

Examples:

- Circulant
- Toeplitz
- Hankel
- Low-displacement rank
- FastFood

- HDHDHD
- General Kronecker products

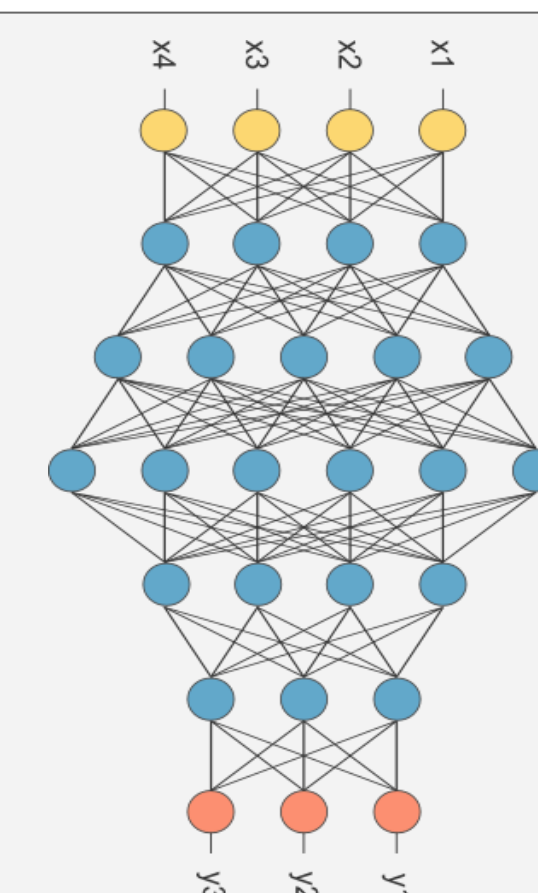
P-model

## Theoretical results

**Theorem 1** Consider a matrix  $M \in \mathbb{R}^{m \times n}$  encoding the weights of connections between a layer  $l_0$  of size  $n$  and a layer  $l_1$  of size  $m$  in some learned unstructured neural network model. Assume that the input to layer  $l_0$  is taken from the  $d$ -dimensional space  $\mathcal{L}$  (although potentially embedded in a much higher dimensional space). Then with probability at least

$$1 - 2p(n)d - 2 \binom{md}{2} e^{-\Omega(\min(\frac{t^2 n^2}{K^4 \Lambda_F^2 \delta^4(n)}, \frac{tn}{K^2 \Lambda_2 \delta^2(n)}))}$$

for  $t = \frac{1}{md}$  and with respect to random choices of  $M_1$  and  $M_2$ , there exists a vector  $r$  defining  $M_3$  such that the structured spinner  $M^{\text{struct}} = M_3 M_2 M_1$  equals to  $M$  on  $\mathcal{L}$ .



**Lemma 1 (structured random setting theorem)** Let  $A$  be a randomized algorithm using unstructured Gaussian matrices  $G$  and let  $A^{TS}$  be its structured version obtained by replacing the unstructured matrix  $G$  by TripleSpinner with blocks of  $m$  rows each. Denote by  $d$  the dimensionality of the space on which  $A$  acts. Then for  $n$  large enough and  $\epsilon = o_{md}(1)$  with probability  $p_{\text{succ}}$  at least:

$$1 - 2p(n)d - 2 \binom{md}{2} e^{-\Omega(\min(\frac{\epsilon^2 n^2}{K^4 \Lambda_F^2 \delta^4(n)}, \frac{\epsilon n}{K^2 \Lambda_2 \delta^2(n)}))}$$

with respect to the random choices of  $M_1$  and  $M_2$  the following holds for any  $S$  such that  $A^{-1}(S)$  is measurable and  $b$ -convex:

$$|\mathbb{P}[A(q) \in S] - \mathbb{P}[A^{TS}(q) \in S]| \leq b\eta,$$

where the the probabilities in the last formula are with respect to the random choice of  $M_3$  and  $\eta = \frac{\delta^3(n)}{n^{\frac{5}{6}}}$ .

## Experiments

