



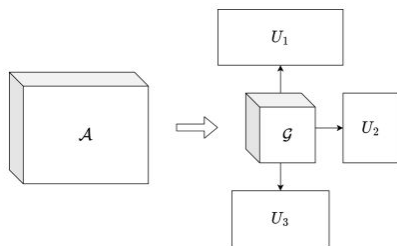
## MARS: Masked Automatic Ranks Selection in Tensor Decomposition

Paper Review

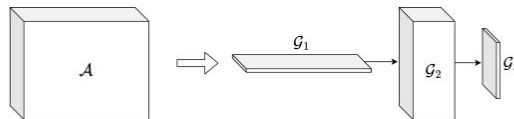
UFO Team

**Motivation&Problem:** Tensor decomposition methods are effective in compression and acceleration of neural networks. However, to achieve balance between compression and performance we need to carefully select tensor ranks.

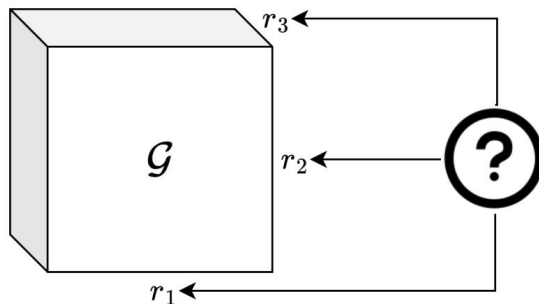
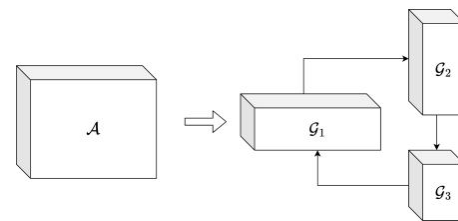
**Tucker**



**Tensor Train (TT)**



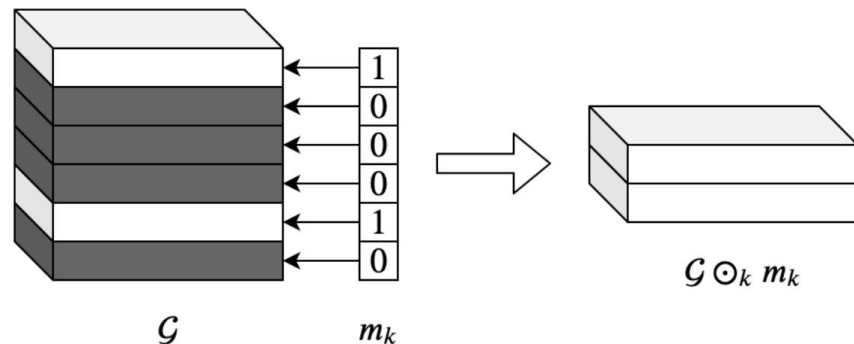
**Tensor Ring (TR)**



- Typical hyperparameter selection techniques, like cross-validation, are poorly suited for the choice of multiple tensor ranks.
- Existing approaches for core ranks selection are not general. They also could require significant computational overhead.

# Proposed method

- Ranks can be represented by binary masks over tensor dimensions



- Bayesian model for these binary masks is proposed, MAP estimate is considered
- Learned masks are binarized and applied to core tensors, resulting in compressed model

$$p(Y, \mathbf{m}, \mathbf{G} \mid X) = \prod_{i=1}^N p(y_i \mid x_i, \mathbf{G} \odot \mathbf{m}) p(\mathbf{m}) p(\mathbf{G})$$

$$\mathcal{G}_k \odot \mathbf{m} := \mathcal{G}_k \odot_{k_1} m_{k_1} \cdots \odot_{k_p} m_{k_p}$$

# Model

$$p(Y, \mathbf{m}, \mathbf{G} \mid X) = \prod_{i=1}^N p(y_i \mid x_i, \mathbf{G} \odot \mathbf{m}) p(\mathbf{m}) p(\mathbf{G})$$

Prior over masks is assumed to be Bernoulli with the success parameter  $\pi$

$$p(\mathbf{m}) = p(\mathbf{m} \mid \pi) = \prod_k \prod_{s=1}^{r_k} \pi^{m_k(s)} (1 - \pi)^{1 - m_k(s)}$$

MAP estimate is considered.  
This discrete optimization problem can be reduced to continuous optimization

$$\sum_{i=1}^N \log p(y_i \mid x_i, \mathbf{G} \odot \mathbf{m}) + \log p(\mathbf{m}) + \log p(\mathbf{G}) \longrightarrow \max_{\mathbf{m}, \mathbf{G}}$$

# Model

Under the assumption that  $q(\mathbf{m})$  is a factorized Bernoulli distribution, problem from previous slide is equivalent to

$$\begin{aligned} \mathbb{E}_{\mathbf{m} \sim q_{\phi}(\mathbf{m})} & \left[ \sum_{i=1}^N \log p(y_i \mid x_i, \mathbf{G} \odot \mathbf{m}) \right] + \\ & + \sum_k \sum_{s=1}^{r_k} [\phi_k(s) \log \pi + (1 - \phi_k(s)) \log(1 - \pi)] + \\ & + \log p(\mathbf{G}) \longrightarrow \max_{\phi, \mathbf{G}}. \end{aligned}$$

After applying reparameterization trick, we can use stochastic gradient descent to get MAP estimate of  $\mathbf{G}$  (cores) and  $\phi$  (probability parameter from Bernoulli distribution of masks).

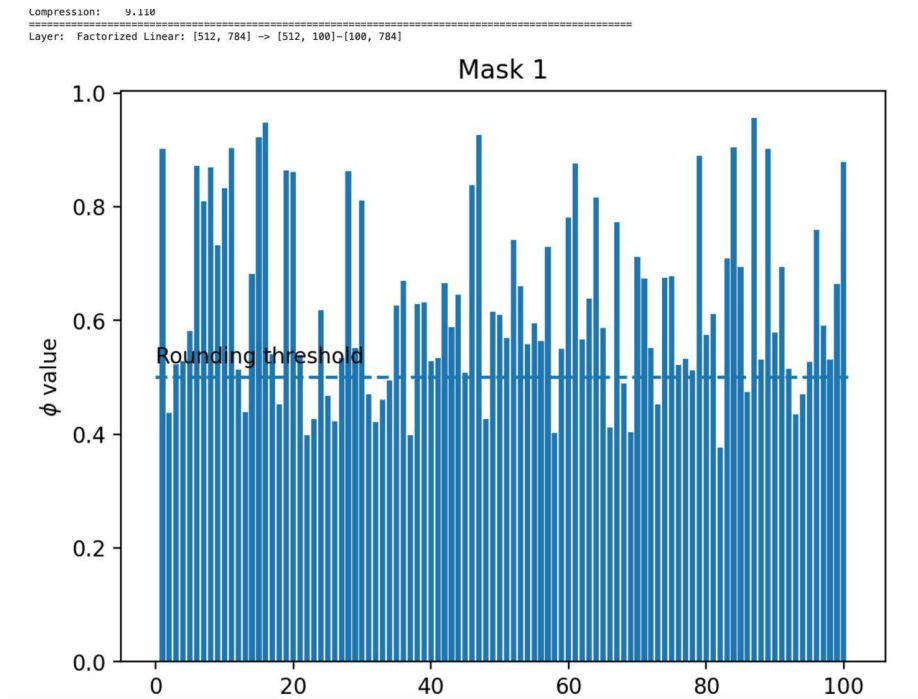
Finally, binary masks are obtained by rounding  $\phi$

Mode-k Hadamard product between cores and corresponding masks gives compressed tensors.

# Reproducing paper results

<https://github.com/xiyori/mars-reproducibility>

- The code provided by the authors only reproduces 1 experiment with a simple fully-connected model
- We implement 2 more experiments: LeNet-5 on MNIST and ResNet-110 on CIFAR10
- Our results align with the paper sufficiently well



# Reproducing paper results: LeNet-5

- Tucker decomposition for convolution
- Low-rank factorization for linear (Skeleton decomposition!)

$$A = UV^T, U \in \mathbb{R}^{m \times r}, V \in \mathbb{R}^{n \times r}$$

Model	Compression	Accuracy	Speed-up (inference)	Slow-down (training)
Original results				
Baseline	1×	99.2%	1×	-
MARS + Tucker	10 ± 0.8×	99.0 ± 0.07%	1.19 ± 0.01×	-
5-ensemble	2×	99.5%	faster in parallel	-
Our reproduction				
Baseline	1×	99.3%	1×	1×
MARS + Tucker	7.1×	99.0 ± 0.08%	1.32 ± 0.04×	1.04×
5-ensemble	1.4×	99.1%	faster in parallel	5.2×

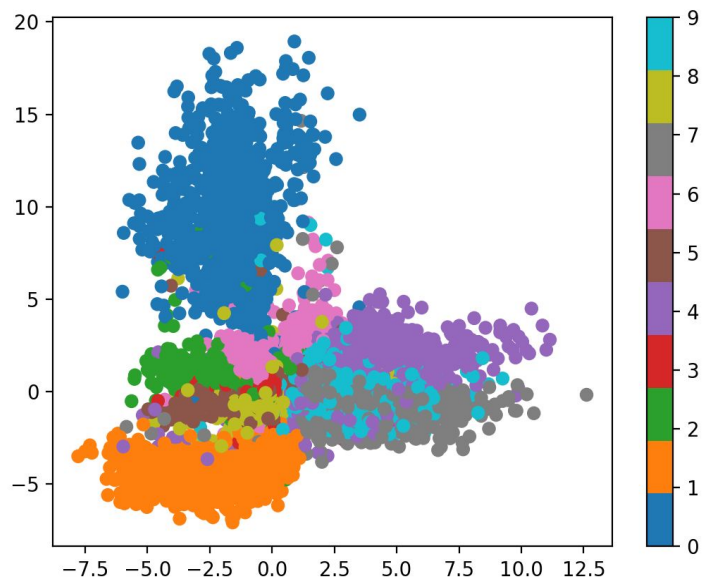
# Reproducing paper results: ResNet-110

- Tensor Train for convolutions in 2nd and 3rd blocks
- Naive:  $2d + 1$  tensors  $r_{k-1} \times n_k \times r_k$
- Proper: similar number of tensors  $r_{k-1} \times n_k \times m_k \times r_k$ , inspired by Garipov et al. (2016)

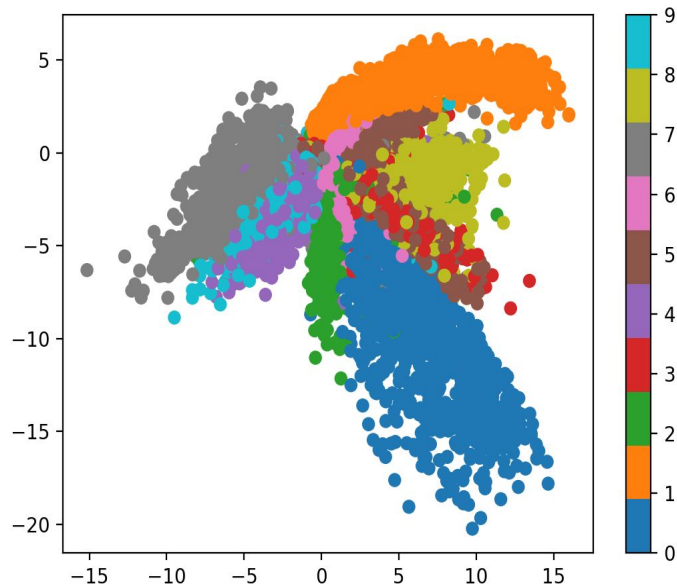
Model	Compression	Accuracy	Slow-down (inference)	Slow-down (training)
Original results				
Baseline	$1\times$	92.6%	-	-
MARS (naive)	$7.0\times$	90.7%	-	-
MARS (proper)	$5.5\times$	91.1%	-	-
Our reproduction				
Baseline	$1\times$	92.3%	$1\times$	$1\times$
MARS (naive)	at least $2.7\times$	at least 89%	$1.60\times$	$2.9\times$
MARS (proper)	at least $2.3\times$	?%	$1.54\times$	$2.8\times$



# MARS for AutoEncoders

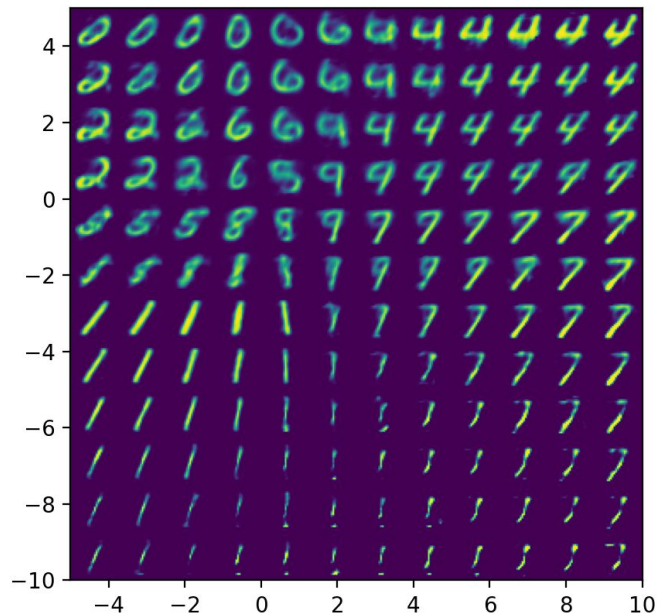


Factorized autoencoders latent space  
Compression: 3.819



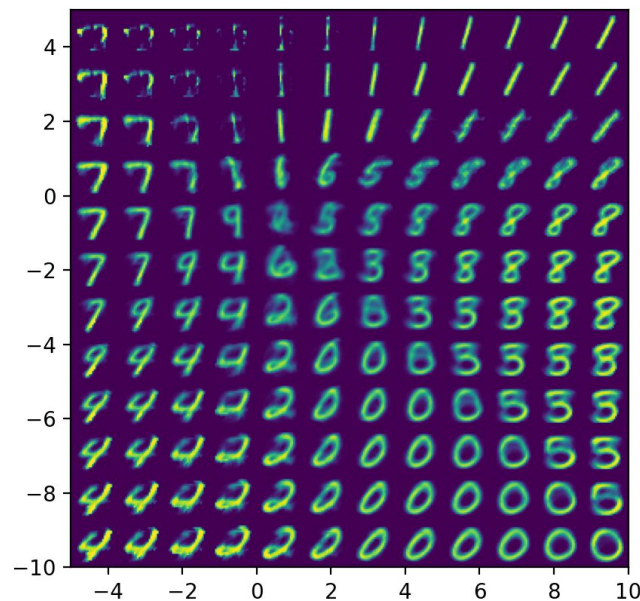
Base autoencoders latent space's latent space

# MARS for AutoEncoders



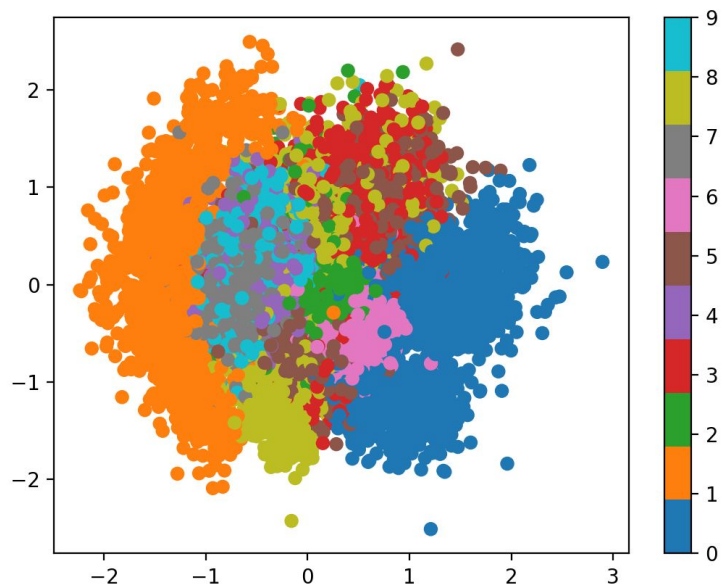
Factorized autoencoder reconstruction

Inference speedup: 1.11

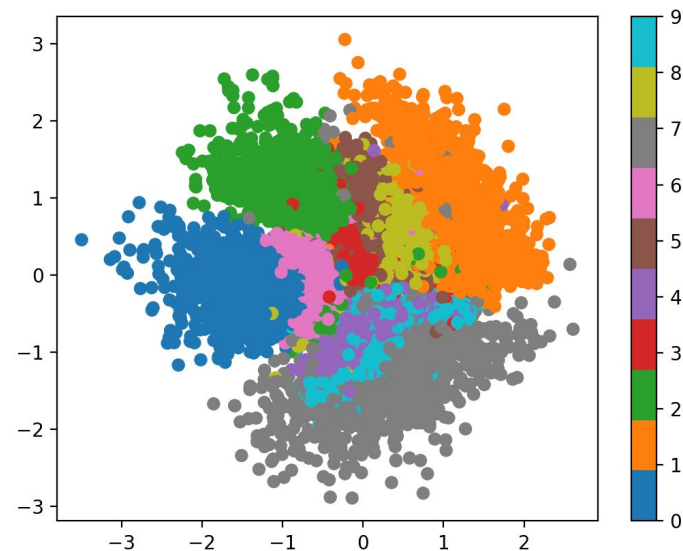


Base autoencoder reconstruction

# MARS for VAE

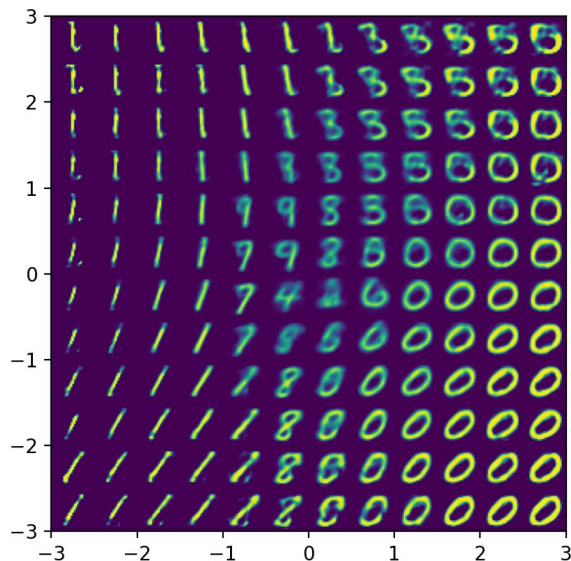


Factorized VAE latent space  
Compression: 5.383



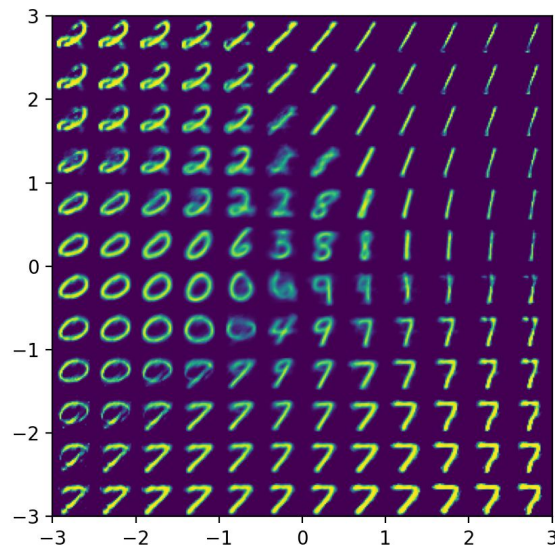
Base VAE latent space

# MARS for VAE



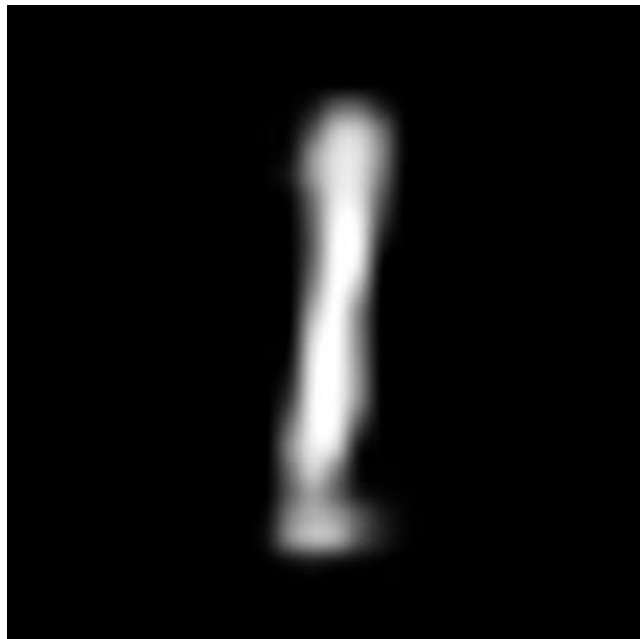
Factorized VAE reconstruction

Inference speedup: 1.12

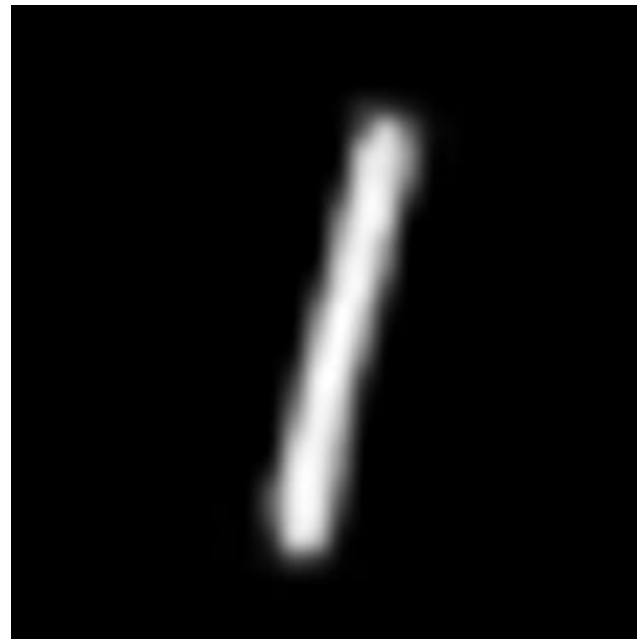


Base VAE reconstruction

# MARS for VAE



Factorized VAE transition



Base VAE transition

# MARS for U-Net

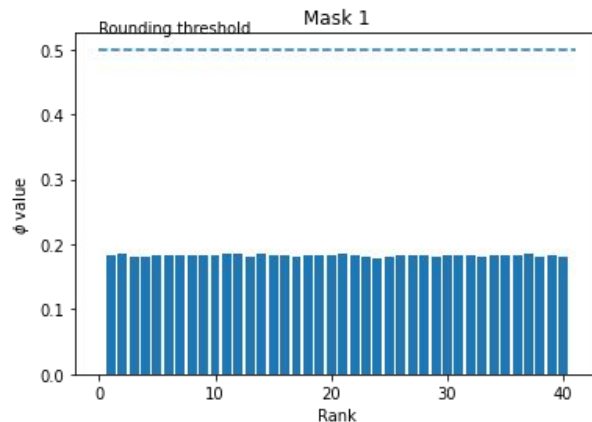
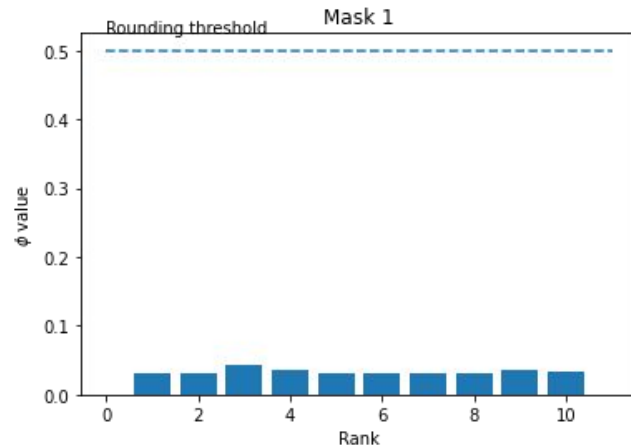
Approach for U-Net:

- Replace ordinary convolution blocks with either ones based on TT or Tucker decomposition (similar to famous approach with MobileNet convolutions)

As a result, size of model decreases from the start

Challenges:

- Hard to make large model to “take off”
- During training MARS more likely to not find any structure inside of models at all and give a zero mask
- Training process take 2x-3x more time, than for ordinary U-Net, because requires some time after convergence to find a mask



# Conclusion

## Model Advantages

- **Universal** - applicable to various tensorized models
- **Sensible** - closely rank approximation
- **Effective** - much better than manual selection of ranks and no worse than specialized rank selection schemes
- **Efficient** - no extra computational cost
- **Scalable** - easy tensorization of ResNet-110
- **Consistent** - learned masks probabilities are close to hard values  $\{0, 1\}$

## Model Disadvantages

- In practice it is difficult to choose hyperparameters
- Ensemble learning approach is quite long
- Model is working long for TensorTrain rank selection



## Our Team



**Elfat Sabitov**  
Experimental  
applications



**Petr Sychev**  
Theory  
understanding,  
preso



**Foma Shipilov**  
Experiments  
reproduction



**Petr Kushnir**  
Theory  
understanding,  
preso



**Sergey Kushneryuk**  
Experiments with  
CNN