

Ultimate tensorization: compressing convolutional and FC layers alike

Timur Garipov¹ Dmitry Podoprikin^{1,2}
Alexander Novikov³ Dmitry Vetrov^{2,3}



¹Moscow State University,
Moscow, Russia



²Yandex LLC,
Moscow, Russia



³Higher School of Economics,
Moscow, Russia

Why compress convolutions?

Large neural networks:

- are expensive to download to smartphones;
- drain battery.

Why compress convolutions?

Large neural networks:

- are expensive to download to smartphones;
- drain battery.

We focus on compressing **convolutional** layers because:

- several modern architectures (Inception, ResNet) lack fully-connected layers;
- we can **already** compress fully-connected layers to move the bottleneck into convolutional layers¹.

¹Novikov et al. *Tensorizing Neural Networks*, NIPS-15

How to compress convolutions?

