# Week 14
## Master Thesis 2020

Tobias Engelhardt Rasmussen (s153057)

DTU Compute

December 17, 2020

# Outline

# Since Last

- FLOPs (Floating Point Operations - $+$ - * %) for calculating theoretical speed-up
- Timing - Found a timing function that seems to give stable results
- Writing - Methodology and appendices

# Outline

# FLOPs given a network layer

Stolen:

Convolutional FLOPs maybe a bit too simple... Does not take padding or stride into account.

Made my own calculations based on the size of the output tensor with very similar results...

# FLOPs given a network

My calculations for convolutions based on the number of output values:

$$FLOPs = \underbrace{T \cdot F' \cdot H' \cdot W'}_{\text{Size of output}} \cdot \left( 2 \cdot \underbrace{S \cdot K^3}_{\text{Size of filter}} - 1 \right)$$

Which would correspond to: (for 2D)

$$FLOPs = C_{out} \cdot H' \cdot W' \cdot \left( 2 \cdot C_{in} K^2 - 1 \right)$$

While theirs is (rearranged):

$$FLOPs = C_{out} \cdot H \cdot W \cdot 2 \cdot \left( C_{in} K^2 + 1 \right)$$

# Outline

# Timing results

| Number of parameters | 3D conv 1 | 3D conv 2 | Linear 1 | Linear 2 | Linear 3 | Total |
|---|---|---|---|---|---|---|
| Original | 14,520 | 58,080 | 358,400 | 10,752 | 168 | 441,920 |
| Compressed | 4,852 | 3,690 | 14,266 | 212 | 168 | 23,188 |
| **Ratio** | 0.334 | 0.064 | 0.039 | 0.020 | 1 | 0.052 |

The timing is performed on a CPU

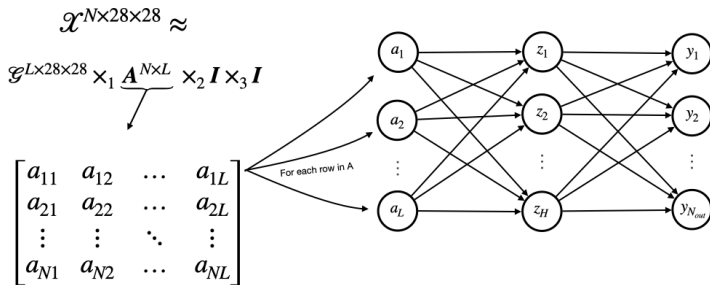| | Timing | 3D conv 1 | 3D conv 2 | Linear 1 | Linear 2 | Linear 3 | Total |
|---|---|---|---|---|---|---|---|
| **FLOPs (1000)** | Original | 15,618,892 | 697,158 | 716.8 | 21.5 | 0.34 | 16,316,789.4 |
| | Compressed | 1,943,628 | 27,490.1 | 28.53 | 0.4 | 0.34 | 1,971,147.4 |
| | **Speed-up** | 8.0359 | 25.3603 | 25.1288 | 50.8369 | 1 | **8.2778** |
| **Time (s)** | Original | 2.6807 | 0.01623 | $1.3 \cdot 10^{-5}$ | $5.6 \cdot 10^{-7}$ | $1.1 \cdot 10^{-7}$ | 2.9594 |
| | Compressed | 1.1688 | 0.007246 | $2.1 \cdot 10^{-8}$ | $2.1 \cdot 10^{-8}$ | $9.8 \cdot 10^{-8}$ | 1.0805 |
| | **Speed-up** | 2.2936 | 2.2434 | 6.4310 | 2.5622 | 1.1249 | **2.7390** |
| Accounts for approx. (%) | | 95.7 | 4.27 | $4.4 \cdot 10^{-3}$ | $1 \cdot 10^{-4}$ | $2 \cdot 10^{-6}$ | 100 |

# Outline

# Methodology Section

Overall outline:

- ▶ Decomposing the Input
  - ▶ Estimating the loadings for the testing data
- ▶ Compressing a pre-trained Network Using Tucker
  - ▶ The linear / dense layer
  - ▶ The convolutional layer
  - ▶ Rank selection
  - ▶ One-shot Compression of an entire CNN using Tucker

# Decomposing the Input

For MNIST:

$$\mathcal{X}^{N\times 28\times 28} \approx$$

$$\mathcal{G}^{L\times 28\times 28} \times_1 \underbrace{\boldsymbol{A}^{N\times L}} \times_2 \boldsymbol{I} \times_3 \boldsymbol{I}$$

$$\begin{bmatrix} a_{11} & a_{12} & \dots & a_{1L} \\ a_{21} & a_{22} & \dots & a_{2L} \\ \vdots & \vdots & \ddots & \vdots \\ a_{N1} & a_{N2} & \dots & a_{NL} \end{bmatrix}$$

For each row in A

# One-Shot Compression of an Entire CNN Using Tucker

**Algorithm 1** One-Shot Tucker Compression of CNN

1: **net** ← define an appropriate CNN
2: train **net** using the given training data
3: **net_dcmp** ← make a copy of **net**
4: **for** each **convolutional** layer; layer in net_dcmp **do**
5:     $\mathcal{K}_{\text{layer}}$ ← take out weight tensor
6:     $R_4, R_5$ ← choose appropriate rank(s)
7:     $\mathcal{G}, \boldsymbol{U}^{(4)}, \boldsymbol{U}^{(5)}$ ← Decompose $\mathcal{K}_{\text{layer}}$ using relevant algorithm  ▷ Or just one $\boldsymbol{U}$
8:     **layer_dcmp_1** ← define new $1 \times 1$ convolution  ▷ If applicable
9:     **layer_dcmp_2** ← define new 3D convolution
10:    **layer_dcmp_3** ← define new $1 \times 1$ convolution  ▷ If applicable
11:    $\mathcal{K}_{\text{layer\_dcmp\_1}}$ ← $\boldsymbol{U}^{(4)}$  ▷ If applicable
12:    $\mathcal{K}_{\text{layer\_dcmp\_2}}$ ← $\mathcal{G}$
13:    $\mathcal{K}_{\text{layer\_dcmp\_3}}$ ← $\boldsymbol{U}^{(5)}$  ▷ If applicable
14:    $\boldsymbol{b}_{\text{layer\_dcmp\_3}}$ ← $\boldsymbol{b}_{\text{layer}}$ add the bias to the last layer  ▷ Or in the line above
15:    layer ← sequence(layer_dcmp_1, layer_dcmp_2, layer_dcmp_3)
16: **end for**
17: **for** each **linear** layer; layer in net_dcmp **do**
18:    $\boldsymbol{W}_{\text{layer}}$ ← take out weight matrix of size $N_{out} \times N_{in}$
19:    $R_A, R_B$ ← choose appropriate rank(s)
20:    $\boldsymbol{G}, \boldsymbol{A}, \boldsymbol{B}$ ← decompose $\boldsymbol{W}_{\text{layer}}$ using relevant algorithm  ▷ Or just $\boldsymbol{A}$ or $\boldsymbol{B}$
21:    **layer_dcmp_1** ← define new $R_B \times N_{in}$ linear layer  ▷ If applicable
22:    **layer_dcmp_2** ← define new $R_A \times R_B$ linear layer  ▷ Or $R_B \times N_{in}$ or $N_{out} \times R_A$
23:    **layer_dcmp_3** ← define new $N_{out} \times R_A$ linear layer  ▷ If applicable
24:    $\boldsymbol{W}_{\text{layer\_dcmp\_1}}$ ← $\boldsymbol{B}$  ▷ If applicable
25:    $\boldsymbol{W}_{\text{layer\_dcmp\_2}}$ ← $\boldsymbol{G}$
26:    $\boldsymbol{W}_{\text{layer\_dcmp\_3}}$ ← $\boldsymbol{A}$  ▷ If applicable
27:    $\boldsymbol{b}_{\text{layer\_dcmp\_3}}$ ← $\boldsymbol{b}_{\text{layer}}$ add the bias to the last layer  ▷ Or in the line above
28:    layer ← sequence(layer_dcmp_1, layer_dcmp_2, layer_dcmp_3)
29: **end for**
30: train **net_dcmp** using the given training data  ▷ fine-tuning