# **Convolutions and More as Einsum**

Felix Dangel

NeurIPS 2024



# Conceptual Overview



### Convolution

\*

=

# Conceptual Overview

### Convolution

\*

=

#### **Tensor networks**



# Conceptual Overview

#### Convolution

\*

=

#### **Tensor networks**



#### **Evaluation**

einsum("c\_in i1 i2, k1 o1 i1, k2 o2 i2, c\_out c\_in k1 k2 -> c\_out o1 o2", X, Pi1, Pi2, W)

















Conv.

w

- Court





Conv.

w - cout





Χ



 $old X \ o [\![ old X ]\!]$ 



$$\begin{split} \mathbf{X} & \\ \to \llbracket \mathbf{X} \rrbracket & \\ \to \mathbf{1}^{\top}\llbracket \mathbf{X} \rrbracket & \end{split}$$



$$\begin{split} & \textbf{X} \\ & \rightarrow \llbracket \textbf{X} \rrbracket \\ & \rightarrow \textbf{1}^{\top}\llbracket \textbf{X} \rrbracket \\ & \rightarrow \left( \textbf{1}^{\top}\llbracket \textbf{X} \rrbracket \right)^{\top} \left( \textbf{1}^{\top}\llbracket \textbf{X} \rrbracket \right) \end{split}$$



$$\begin{split} & \textbf{X} \\ & \rightarrow \llbracket \textbf{X} \rrbracket \\ & \rightarrow \textbf{1}^{\top}\llbracket \textbf{X} \rrbracket \\ & \rightarrow \left( \textbf{1}^{\top}\llbracket \textbf{X} \rrbracket \right)^{\top} \left( \textbf{1}^{\top}\llbracket \textbf{X} \rrbracket \right) \end{split}$$

#### **Tensor Network**





 $(1) - k'_1$ 

k'



Time: **9.87 ms** 

Time: 2.69 ms (3.7 x)

(features.1.0.block.0 convolution of ConvNeXt-base with (32, 3, 256, 256) input)

 $\mathbf{V}$ 

Non-conventional operations can be much cheaper with tensor networks

### State of the art

$$\begin{split} & \textbf{X} \\ & \rightarrow \llbracket \textbf{X} \rrbracket \\ & \rightarrow \textbf{1}^{\top} \llbracket \textbf{X} \rrbracket \\ & \rightarrow \left( \textbf{1}^{\top} \llbracket \textbf{X} \rrbracket \right)^{\top} \left( \textbf{1}^{\top} \llbracket \textbf{X} \rrbracket \right) \end{split}$$

#### Tensor Network



Time: 9.87 ms Extra memory: 3.07 GiB Time: **2.69 ms (3.7 x)** Extra memory: **0 MiB** 

(features.1.0.block.0 convolution of ConvNeXt-base with (32, 3, 256, 256) input)



#### ImageNet GPU benchmark (NVIDIA A40) with (128, 3, 256, 256) inputs

Net	Optimizer	Per-iteration [s]	Peak memory [GiB]
ResNet18	SGD	1.17 · 10 <sup>−1</sup> (1.00 x)	3.62 (1.00 x)
VGG19	SGD	6.90 · 10 <sup>-1</sup> (1.00 x)	14.1 (1.00 x)



#### ImageNet GPU benchmark (NVIDIA A40) with (128, 3, 256, 256) inputs

Net	Optimizer	Per-iteration [s]	Peak memory [GiB]
ResNet18	SGD <b>SINGD</b>	$ \begin{array}{c c} 1.17 \cdot 10^{-1} & (1.00 \text{ x}) \\ 1.94 \cdot 10^{-1} & (1.67 \text{ x}) \end{array} \end{array} $	3.62 (1.00 x) 4.54 <b>(</b> 1.25 <b>x)</b>
VGG19	SGD SINGD	6.90 · 10 <sup>-1</sup> (1.00 x) 1.02 (1.48 x)	14.1 (1.00 x) 32.1 (2.28 x)

(SINGD uses KFAC-reduce and diagonal pre-conditioners which are updated every step)



ImageNet GPU benchmark (NVIDIA A40) with (128, 3, 256, 256) inputs

Net	Optimizer	Per-iteration [s]	Peak memory [GiB]
ResNet18	SGD	$1.17 \cdot 10^{-1} (1.00 \text{ x})$	3.62(1.00  x)
	SINGD+TN	$1.56 \cdot 10^{-1} (1.33 \text{ x})$	3.63 (1.00 x)
VGG19	SGD	6.90 · 10 <sup>-1</sup> (1.00 x)	14.1 (1.00 x)
	SINGD	1.02 <b>(</b> 1.48 <b>x</b> )	32.1 <mark>(2.28 x)</mark>
	SINGD+TN	8.01 · 10 <sup>−1</sup> (1.16 x)	16.1 <mark>(</mark> 1.14 <b>x)</b>

(SINGD uses KFAC-reduce and diagonal pre-conditioners which are updated every step)

Significantly reduces the computational gap of second-order methods.

# Convolutions and More as einsum

 $\sim$ 

- + TN perspective simplifies the transfer of algorithmic ideas
- + Enables flexible/faster implementations of black box routines
- + Relies on automatically efficient evaluation inside einsum

# Convolutions and More as einsum

- + TN perspective simplifies the transfer of algorithmic ideas
- + Enables flexible/faster implementations of black box routines
- + Relies on automatically efficient evaluation inside einsum

#### Try it out!

```
from einconv.expressions import kfac_reduce
from torch import einsum
```

```
# create the tensor network
equation, operands, shape = kfac_reduce.
    einsum_expression(..., simplify=True)
# evaluate it
einsum(equation, *operands).reshape(shape)
```



#### pip install einconv

# Convolutions and More as einsum

- + TN perspective simplifies the transfer of algorithmic ideas
- + Enables flexible/faster implementations of black box routines
- + Relies on automatically efficient evaluation inside einsum

#### Try it out!

```
from einconv.expressions import kfac_reduce
from torch import einsum
```

```
# create the tensor network
equation, operands, shape = kfac_reduce.
    einsum_expression(..., simplify=True)
# evaluate it
einsum(equation, *operands).reshape(shape)
```



#### pip install einconv

Paper: arxiv/2307.02275
Code: github.com/f-dangel/einconv





- Runa Eschenhagen, Alexander Immer, Richard E. Turner, Frank Schneider, and Philipp Hennig. Kronecker-factored approximate curvature for modern neural network architectures. In Advances in Neural Information Processing Systems (NeurIPS), 2023.
- Wu Lin, Felix Dangel, Runa Eschenhagen, Kirill Neklyudov, Agustinus Kristiadi, Richard E. Turner, and Alireza Makhzani. Structured inverse-free natural gradient descent: Memory-efficient & numerically-stable KFAC. In International Conference on Machine Learning (ICML), 2024.