Convolutions Through the Lens of Tensor Networks

Felix Dangel

June 26/July 02, 2024



Convolutions Through the Lens of Tensor Networks as einsum

Felix Dangel

June 26/July 02, 2024



The Broader Picture

The Crowded Valley of Deep Learning Optimizers [Schmidt et al., 2021, ICML]

 \mathbf{V}

AcceleGrad ACClip AdaAlter AdaBatch AdaBayes/AdaBayes-SS AdaBelief AdaBlock AdaBound/AMSBound AdaComp Adadelta Adafactor AdaFix AdaFom AdaFTRL Adagrad ADAHESSIAN Adai AdaLoss Adam Adam+ AdamAL AdaMax AdamBS AdamNC AdaMod AdamP AdamT AdamW AdamX ADAS AdaS AdaScale AdaSGD AdaShift AdaSqrt Adathm AdaX/AdaX-W AEGD ALI-G AMSGrad AngularGrad ArmijoLS ARSG ASAM AutoLRS AvaGrad BAdam BGAdam BPGrad BRMSProp BSGD C-ADAM CADA COCOB Cool Momentum CProp Curveball Dadam DeepMemory DGNOpt DiffGrad EAdam EKFAC Eve Expectigrad FastAdaBelief FRSGD G-AdaGrad GADAM Gadam GOALS GOLS-I Grad-Avg GRAPES Gravilon Gravity HAdam HyperAdam K-BFGS/K-BFGS(L) K-FAC KF-ON-CNN KFLR/KFRA L4Adam/L4Momentum LAMB LaProp LARS LHOPT LookAhead M-SVAG MADGRAD MAS MEKA MTAdam MVRC-1/MVRC-2 Nadam NAMSB/NAMSG ND-Adam Nero Nesterov Accelerated Momentum Noisy Adam/Noisy K-FAC NosAdam Novograd NT-SGD Padam PAGE PAL PolyAdam Polyak Momentum PowerSGD/PowerSGDM Probabilistic Polvak ProbLS PSTorm OHAdam/OHM RAdam Ranger RangerLars RecursiveOptimizer RescaledExp RMSProp RMSterov S-SGD SAdam Sadam/SAMSGrad SALR SAM SC-Adagrad/SC-RMSProp SDProp SGD SGD-BB SGD-G2 SGDEM SGDHess SGDM SGDP SGDR SHAdagrad Shampoo SignAdam++ SignSGD SKON/S4ON SM3 SMG SNGM SoftAdam SRSGD Step-Tuned SGD STORM SWATS SWNTS TAdam TEKFAC VAdam VR-SGD vSGD-b/vSGD-g/vSGD-l vSGD-fd WNGrad YellowFin Yoqi 🤉











Customized Backpropagation

- 1. What to backpropagate
- 2. How to backpropagate
- 3. How to extract target quantity



Convolutions Are More Tedious than Linear Layers

Personal example: Vector-Jacobian products (VJPs)

```
def _weight_jac_t_mat_prod(
   self.
   module: Linear.
  g inp: Tuple[Tensor]
  g out: Tuple[Tensor].
  mat: Tensor
  sum batch: int = True.
  subsampling: List[int] = None.
-> Tensor:
  """Batch-apply transposed Jacobian of the output w.r.t. the weight.
  Args:
       module: Linear layer.
      g_inp: Gradients w.r.t. module input. Not required by the implementation.
      g out: Gradients w.r.t. module output. Not required by the implementation.
      mat: Batch of "V" vectors of same shape as the laver output
          (``[N, *, out_features]``) to which the transposed output-input Jacobian
          is applied. Has shape \sum V. N. *, out features \sum  if subsampling is not
          used, otherwise ``N`` must be ``len(subsampling)`` instead.
      sum_batch: Sum the result's batch axis. Default: "True".
       subsampling: Indices of samples along the output's batch dimension that
          should be considered. Defaults to "None" (use all samples).
   Returns:
       Batched transposed Jacobian vector products. Has shape
      "[V. N. *module.weight.shape]" when ``sum batch`` is ``False``. With
      ``sum_batch=True``, has shape ``[V, *module.weight.shape]``. If sub-
      sampling is used. ``N`` must be ``len(subsampling)`` instead.
   d_weight = subsample(module.input0. subsampling=subsampling)
  equation = f''vn...o.n...i \rightarrow vf'' if sum batch else 'n'}oi"
   return einsum(equation, mat, d weight)
```

user manner transformett, nament, nament, sur-tagen(tamer), sur-tag I have there properties construction (a) and here (b) are in the second The state of the s I Restar to a read of the life of the second state of the second s promi i (interpretante mark, interpretation restant), i i i interpretation restante (i interpretation) 17788 1.71. 18864.071841 mer + higher, may have a

nager grad is anger grad and (r) anger grad is anger grad and (r)

Convolutions Are More Tedious than Linear Layers

- Personal example: Vector-Jacobian products (VJPs)
- + Other algorithmic & theoretical ideas

Concept	for MLPs	for CNNs
Approximate Hessian diagonal	1989	2023
Kronecker-factored curvature (KFAC, KFRA, KFLR)	2015, 2017, 2017	2016, 2020, 2020
Kronecker-factored quasi-Newton methods (KBFGS)	2021	2022
Neural tangent kernel (NTK)	2018	2019
Hessian rank	2021	2023
Gradient descent learning dynamics	2014	2023

Convolutions Are More Tedious than Linear Layers

- Personal example: Vector-Jacobian products (VJPs)
- + Other algorithmic & theoretical ideas

Concept	for MLPs	for CNNs
Approximate Hessian diagonal	1989	2023
Kronecker-factored curvature (KFAC, KFRA, KFLR)	2015, 2017, 2017	2016, 2020, 2020
Kronecker-factored quasi-Newton methods (KBFGS)	2021	2022
Neural tangent kernel (NTK)	2018	2019
Hessian rank	2021	2023
Gradient descent learning dynamics	2014	2023

This talk: A simplifying perspective of convolutions through tensor networks

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)

Convolutions – A Visual Walkthrough

A Simple Convolution



$\textbf{Y} = \textbf{X} \star \textbf{W}$ with

- + $\mathbf{X} \in \mathbb{R}'$
- + $\mathbf{W} \in \mathbb{R}^{K}$
- + $\mathbf{Y} \in \mathbb{R}^{O}$



A Simple Convolution

\mathbf{V}

$\textbf{Y} = \textbf{X} \star \textbf{W}$ with

- + $\mathbf{X} \in \mathbb{R}'$
- + $\mathbf{W} \in \mathbb{R}^{K}$
- + $\mathbf{Y} \in \mathbb{R}^{O}$

 \star

S is for Stride

7

$\textbf{Y} = \textbf{X} \star \textbf{W}$ with

- + $\mathbf{X} \in \mathbb{R}'$
- + $\mathbf{W} \in \mathbb{R}^{K}$
- + $\mathbf{Y} \in \mathbb{R}^{O}$



S is for Stride

V

$\mathbf{Y} = \mathbf{X} \star \mathbf{W}$ with

- + $\mathbf{X} \in \mathbb{R}'$
- + $\mathbf{W} \in \mathbb{R}^{K}$
- + $\mathbf{Y} \in \mathbb{R}^{O}$

 \star

P is for Padding



$\mathbf{Y} = \mathbf{X} \star \mathbf{W}$ with

- + $\mathbf{X} \in \mathbb{R}'$
- + $\mathbf{W} \in \mathbb{R}^{K}$
- + $\mathbf{Y} \in \mathbb{R}^{O}$



P is for Padding



$\boldsymbol{Y} = \boldsymbol{X} \star \boldsymbol{W}$ with

- + $\mathbf{X} \in \mathbb{R}'$
- + $\mathbf{W} \in \mathbb{R}^{K}$
- + $\mathbf{Y} \in \mathbb{R}^{O}$

 \star

D is for Dilation

7

$\textbf{Y} = \textbf{X} \star \textbf{W}$ with

- + $\mathbf{X} \in \mathbb{R}'$
- + $\mathbf{W} \in \mathbb{R}^{K}$
- + $\mathbf{Y} \in \mathbb{R}^{O}$



D is for Dilation

V

$\boldsymbol{Y} = \boldsymbol{X} \star \boldsymbol{W}$ with

- + $\mathbf{X} \in \mathbb{R}'$
- + $\mathbf{W} \in \mathbb{R}^{K}$
- + $\mathbf{Y} \in \mathbb{R}^{O}$



C_{in} is for Input Channels

$\mathbf{Y} = \mathbf{X} \star \mathbf{W}$ with

- + $\mathbf{X} \in \mathbb{R}^{C_{in} \times I}$ + $\mathbf{W} \in \mathbb{R}^{C_{in} \times K}$
- + $\mathbf{Y} \in \mathbb{R}^{O}$





C_{in} is for Input Channels

$\int_{\mathbf{v}}$

$\boldsymbol{Y} = \boldsymbol{X} \star \boldsymbol{W}$ with

- + $\mathbf{X} \in \mathbb{R}^{C_{in} \times I}$
- + $\mathbf{W} \in \mathbb{R}^{C_{\text{in}} \times K}$
- + $\mathbf{Y} \in \mathbb{R}^{O}$

 \star

Cout is for Output Channels

- $\mathbf{Y} = \mathbf{X} \star \mathbf{W}$ with
 - + $\mathbf{X} \in \mathbb{R}^{C_{in} \times I}$
 - + $\mathbf{W} \in \mathbb{R}^{C_{\text{out}} \times C_{\text{in}} \times K}$ + $\mathbf{Y} \in \mathbb{R}^{C_{\text{out}} \times O}$



C_{out} is for Output Channels

 \star

=

$\mathbf{Y} = \mathbf{X} \star \mathbf{W}$ with

- + $\mathbf{X} \in \mathbb{R}^{C_{in} imes I}$
- + $\mathbf{W} \in \mathbb{R}^{C_{\text{out}} \times C_{\text{in}} \times K}$
- + $\mathbf{Y} \in \mathbb{R}^{C_{\text{out}} \times O}$

G is for Groups

- $\mathbf{Y} = \mathbf{X} \star \mathbf{W}$ with
 - + $\mathbf{X} \in \mathbb{R}^{C_{in} \times I}$
 - + $\mathbf{W} \in \mathbb{R}^{C_{\text{out}} \times c_{\text{in/G}} \times K}$ + $\mathbf{Y} \in \mathbb{R}^{C_{\text{out}} \times O}$



G is for Groups

V

$\boldsymbol{Y} = \boldsymbol{X} \star \boldsymbol{W}$ with

- + $\mathbf{X} \in \mathbb{R}^{C_{in} \times I}$
- + $\mathbf{W} \in \mathbb{R}^{C_{\text{out}} \times C_{\text{in/G}} \times K}$
- + $\mathbf{Y} \in \mathbb{R}^{C_{\text{out}} \times O}$

 \star

N is for Batch Size

\mathbf{V}

$\mathbf{Y} = \mathbf{X} \star \mathbf{W}$ with

- + $\mathbf{X} \in \mathbb{R}^{N \times I}$
- + $\mathbf{W} \in \mathbb{R}^{K}$
- + $\mathbf{Y} \in \mathbb{R}^{N \times O}$



N is for Batch Size

V

$\boldsymbol{Y} = \boldsymbol{X} \star \boldsymbol{W}$ with

- + $\mathbf{X} \in \mathbb{R}^{N \times I}$
- + $\mathbf{W} \in \mathbb{R}^{K}$
- + $\mathbf{Y} \in \mathbb{R}^{N \times O}$

=

*

Putting Everything Together in One Dimension...

$\textbf{Y} = \textbf{X} \star \textbf{W}$ with

- + $\mathbf{X} \in \mathbb{R}^{N \times C_{in} \times I}$
- + $\mathbf{W} \in \mathbb{R}^{C_{\text{out}} \times C_{\text{in/G}} \times K}$
- + $\mathbf{Y} \in \mathbb{R}^{N \times C_{\text{out}} \times O}$



Putting Everything Together in One Dimension...

$\textbf{Y} = \textbf{X} \star \textbf{W}$ with

- + $\mathbf{X} \in \mathbb{R}^{N \times C_{in} \times I}$
- + $\mathbf{W} \in \mathbb{R}^{C_{\text{out}} \times C_{\text{in/G}} \times K}$
- + $\mathbf{Y} \in \mathbb{R}^{N \times C_{\text{out}} \times O}$

 \star

... and in Two Dimensions

V

- $\textbf{Y} = \textbf{X} \star \textbf{W}$ with
 - + $\mathbf{X} \in \mathbb{R}^{N \times C_{in} \times I_1 \times I_2}$
 - + $\mathbf{W} \in \mathbb{R}^{C_{\text{out}} \times C_{\text{in/G}} \times K_1 \times K_2}$
 - + $\mathbf{Y} \in \mathbb{R}^{N \times C_{\text{out}} \times O_1 \times O_2}$



*


... and in Two Dimensions

$\textbf{Y} = \textbf{X} \star \textbf{W}$ with

- + $\mathbf{X} \in \mathbb{R}^{N \times C_{in} \times I_1 \times I_2}$
- + $\mathbf{W} \in \mathbb{R}^{C_{\text{out}} \times c_{\text{in/G}} \times K_1 \times K_2}$
- + $\mathbf{Y} \in \mathbb{R}^{N \times C_{\text{out}} \times O_1 \times O_2}$

=

 \star

Tensor Multiplication & Tensor Networks



(Disclaimer: Tensor = multi-dimensional array; sorry)

(Disclaimer: Tensor = multi-dimensional array; sorry)

Consider two vectors **a**, **b**

Product	Notation
Inner Element-wise	$c = \mathbf{a}^{\top} \mathbf{b}$ $\mathbf{c} = \mathbf{a} \odot \mathbf{b}$
Outer/Kronecker	$oldsymbol{C} = oldsymbol{a}oldsymbol{b}^ op = oldsymbol{a}\otimesoldsymbol{b}^ op$



(Disclaimer: Tensor = multi-dimensional array; sorry)

Consider two vectors **a**, **b**

Product	Notation	Index notation
Inner Element-wise	$oldsymbol{c} = oldsymbol{a}^ op oldsymbol{b}$ $oldsymbol{c} = oldsymbol{a} \odot oldsymbol{b}$	$egin{array}{lll} \mathbf{c} = \sum_i \mathbf{a}_i \mathbf{b}_i \ \mathbf{c}_i = \mathbf{a}_i \mathbf{b}_i \end{array}$
Outer/Kronecker	$\mathbf{C} = \mathbf{a} \mathbf{b}^ op = \mathbf{a} \otimes \mathbf{b}^ op$	${\cal C}_{i,j}={\sf a}_i{\sf b}_j$

 \mathbf{V}

(Disclaimer: Tensor = multi-dimensional array; sorry)

Consider two vectors **a**, **b** and two matrices **A**, **B**

Product	Notation	Index notation
Inner Element-wise Outer/Kronecker	$m{c} = m{a}^{ op} m{b}$ $m{c} = m{a} \odot m{b}$ $m{C} = m{a} m{b}^{ op} = m{a} \otimes m{b}^{ op}$	$egin{aligned} \mathbf{c} &= \sum_i \mathbf{a}_i \mathbf{b}_i \ \mathbf{c}_i &= \mathbf{a}_i \mathbf{b}_i \ \mathbf{C}_{i,j} &= \mathbf{a}_i \mathbf{b}_j \end{aligned}$
Inner Element-wise Kronecker	$C = AB$ $C = A \odot B$ $C = A \otimes B$	

 \mathbf{V}

(Disclaimer: Tensor = multi-dimensional array; sorry)

Consider two vectors **a**, **b** and two matrices **A**, **B**

Product	Notation	Index notation
Inner Element-wise Outer/Kronecker	$c = a^{\top} b$ $c = a \odot b$ $C = ab^{\top} = a \otimes b^{\top}$	$egin{aligned} \mathbf{c} &= \sum_i \mathbf{a}_i \mathbf{b}_i \ \mathbf{c}_i &= \mathbf{a}_i \mathbf{b}_i \ \mathbf{C}_{i,j} &= \mathbf{a}_i \mathbf{b}_j \end{aligned}$
Inner Element-wise Kronecker	$C = AB$ $C = A \odot B$ $C = A \otimes B$	$egin{aligned} & m{C}_{i,k} = \sum_j m{A}_{i,j} m{B}_{j,k} \ & m{C}_{i,j} = m{A}_{i,j} m{B}_{i,j} \ & m{C}_{(i,k),(j,l)} = m{A}_{i,j} m{B}_{k,l} \end{aligned}$



(Disclaimer: Tensor = multi-dimensional array; sorry)

Consider two vectors **a**, **b** and two matrices **A**, **B**

Product	Notation	Index notation	In	Out
Inner Element-wise	$\mathbf{c} = \mathbf{a}^{\top} \mathbf{b}$ $\mathbf{c} = \mathbf{a} \odot \mathbf{b}$	$m{c} = \sum_i a_i m{b}_i \ m{c}_i = a_i m{b}_i$	(i), (i) (i), (i)	() (<i>i</i>)
Outer/Kronecker	$\mathbf{C} = \mathbf{a}\mathbf{b}^ op = \mathbf{a}\otimes\mathbf{b}^ op$	$C_{i,j} = a_i b_j$	(i), (j)	(<i>i</i> , <i>j</i>)
Inner Element-wise Kronecker	$C = AB$ $C = A \odot B$ $C = A \otimes B$	$egin{aligned} & C_{i,k} = \sum_j A_{i,j} B_{j,k} \ & C_{i,j} = A_{i,j} B_{i,j} \ & C_{(i,k),(j,l)} = A_{i,j} B_{k,l} \end{aligned}$	(i,j), (j,k) (i,j), (i,j) (i,j), (k,l)	(i, k) (i, j) ((i, k), (j, l))

We can infer the summations given the index signature

Given two tensors A, B with index tuples S_A, S_B

$$\boldsymbol{\mathsf{C}} := \ast_{(S_{\boldsymbol{\mathsf{A}}},S_{\boldsymbol{\mathsf{B}}},S_{\boldsymbol{\mathsf{C}}})}(\boldsymbol{\mathsf{A}},\boldsymbol{\mathsf{B}}) \quad \Leftrightarrow \quad [\boldsymbol{\mathsf{C}}]_{S_{\boldsymbol{\mathsf{C}}}} = \sum_{(S_{\boldsymbol{\mathsf{A}}}\cup S_{\boldsymbol{\mathsf{B}}})\setminus S_{\boldsymbol{\mathsf{C}}}} [\boldsymbol{\mathsf{A}}]_{S_{\boldsymbol{\mathsf{A}}}} [\boldsymbol{\mathsf{B}}]_{S_{\boldsymbol{\mathsf{B}}}} \,,$$

indices not present in the output are summed out; S_C satisfies $S_C \subseteq S_A \cup S_B$.

Given two tensors A, B with index tuples S_A, S_B

$$\mathbf{C} := \ast_{(S_{\mathbf{A}},S_{\mathbf{B}},S_{\mathbf{C}})}(\mathbf{A},\mathbf{B}) \quad \Leftrightarrow \quad [\mathbf{C}]_{S_{\mathbf{C}}} = \sum_{(S_{\mathbf{A}}\cup S_{\mathbf{B}})\setminus S_{\mathbf{C}}} [\mathbf{A}]_{S_{\mathbf{A}}} [\mathbf{B}]_{S_{\mathbf{B}}} \,,$$

indices not present in the output are summed out; S_C satisfies $S_C \subseteq S_A \cup S_B$.

[Example] Matrix-matrix multiplication $\mathbf{C} = \mathbf{A}\mathbf{B}$ as tensor multiplication

$$\begin{array}{rcl} S_{A} & = & (i,j) \\ S_{B} & = & (j,k) \\ S_{C} & = & (i,k) \end{array}$$

Given two tensors A, B with index tuples S_A, S_B

$$\mathbf{C} := \ast_{(S_{\mathbf{A}},S_{\mathbf{B}},S_{\mathbf{C}})}(\mathbf{A},\mathbf{B}) \quad \Leftrightarrow \quad [\mathbf{C}]_{S_{\mathbf{C}}} = \sum_{(S_{\mathbf{A}}\cup S_{\mathbf{B}})\setminus S_{\mathbf{C}}} [\mathbf{A}]_{S_{\mathbf{A}}}[\mathbf{B}]_{S_{\mathbf{B}}} \,,$$

indices not present in the output are summed out; S_C satisfies $S_C \subseteq S_A \cup S_B$.

[Example] Matrix-matrix multiplication **C** = **AB** as tensor multiplication

$$\begin{array}{rcl} S_{\boldsymbol{A}} &=& (i,j) \\ S_{\boldsymbol{B}} &=& (j,k) \\ S_{\boldsymbol{C}} &=& (i,k) \end{array} \implies (S_{\boldsymbol{A}} \cup S_{\boldsymbol{B}}) \setminus S_{\boldsymbol{C}} = (i,j,k) \setminus (i,k) = (j) \end{array}$$

Given two tensors \mathbf{A}, \mathbf{B} with index tuples $S_{\mathbf{A}}, S_{\mathbf{B}}$

$$\mathbf{C} := \ast_{(S_{\mathbf{A}},S_{\mathbf{B}},S_{\mathbf{C}})}(\mathbf{A},\mathbf{B}) \quad \Leftrightarrow \quad [\mathbf{C}]_{S_{\mathbf{C}}} = \sum_{(S_{\mathbf{A}}\cup S_{\mathbf{B}})\setminus S_{\mathbf{C}}} [\mathbf{A}]_{S_{\mathbf{A}}} [\mathbf{B}]_{S_{\mathbf{B}}} \,,$$

indices not present in the output are summed out; S_C satisfies $S_C \subseteq S_A \cup S_B$.

[Example] Matrix-matrix multiplication $\mathbf{C} = \mathbf{A}\mathbf{B}$ as tensor multiplication

$$\begin{array}{rcl} S_{\boldsymbol{A}} &=& (i,j) \\ S_{\boldsymbol{B}} &=& (j,k) \\ S_{\boldsymbol{C}} &=& (i,k) \end{array} \implies (S_{\boldsymbol{A}} \cup S_{\boldsymbol{B}}) \setminus S_{\boldsymbol{C}} = (i,j,k) \setminus (i,k) = (j) \end{array}$$

 $\implies \boldsymbol{C} = *_{((i,j),(j,k),(i,k))}(\boldsymbol{A},\boldsymbol{B}), \qquad \text{Or} \qquad \text{C} = \text{einsum}(\texttt{``ij,jk->ik''}, A, B).$

Many libraries provide $*_{(...)}(...)$ as einsum.

General case: Given *N* tensors A_1, \ldots, A_N with index tuples S_{A_1}, \ldots, S_{A_N}

$$\mathbf{C} := \ast_{(S_{\mathbf{A}_1}, \dots, S_{\mathbf{A}_N}, S_{\mathbf{C}})}(\mathbf{A}_1, \dots, \mathbf{A}_N) \quad \Leftrightarrow \quad [\mathbf{C}]_{S_{\mathbf{C}}} = \sum_{(S_{\mathbf{A}_1} \cup \dots \cup S_{\mathbf{A}_N}) \setminus S_{\mathbf{C}}} [\mathbf{A}_1]_{S_{\mathbf{A}_1}} \cdots [\mathbf{A}_N]_{S_{\mathbf{A}_N}} \,,$$

indices not present in the output are summed out; S_C satisfies $S_C \subseteq S_{A_1} \cup \cdots \cup S_{A_N}$.

[Example] Matrix-matrix multiplication $\mathbf{C} = \mathbf{AB}$ as tensor multiplication

$$\begin{array}{rcl} S_{\mathbf{A}} &=& (i,j) \\ S_{\mathbf{B}} &=& (j,k) \\ S_{\mathbf{C}} &=& (i,k) \end{array} \implies (S_{\mathbf{A}} \cup S_{\mathbf{B}}) \setminus S_{\mathbf{C}} = (i,j,k) \setminus (i,k) = (j) \end{array}$$

 $\implies \boldsymbol{C} = *_{((i,j),(j,k),(i,k))}(\boldsymbol{A},\boldsymbol{B}), \qquad \text{Or} \qquad \text{C} = \text{einsum}(\texttt{``ij,jk->ik''}, A, B).$

Many libraries provide $*_{(...)}(...)$ as einsum.

Tensor Networks: Graphical Notation for Tensor Multiplications ${ m V}$



Tensor Networks: Graphical Notation for Tensor Multiplications ${ m V}$



Tensor Networks: Graphical Notation for Tensor Multiplications ${ m V}$



Tensor Networks: Graphical Notation for Tensor Multiplications igvee



Tensor Networks: Graphical Notation for Tensor Multiplications igvee



Tensor Networks: Graphical Notation for Tensor Multiplications igvee



Tensor Networks: Graphical Notation for Tensor Multiplications old V



Benefits of TN/einsum Abstraction

 $\mathbf{V}_{\mathbf{v}}$

Powerful syntax extensions: multi-letter indices, index (un-)grouping, ...



Powerful syntax extensions: multi-letter indices, index (un-)grouping, ...

[Example] Batched matrix-matrix multiplication

Given $\{A_n\}_{n=1}^N, \{B_n\}_{n=1}^N$, compute $\{C_n\}_{n=1}^N = \{A_nB_n\}_{n=1}^N$



Powerful syntax extensions: multi-letter indices, index (un-)grouping, ...

[Example] Batched matrix-matrix multiplication

Given $\{A_n\}_{n=1}^N, \{B_n\}_{n=1}^N$, compute $\{C_n\}_{n=1}^N = \{A_nB_n\}_{n=1}^N$

C = einsum(A, B, "batch i j, batch j k -> batch i k")

 $\mathbf{V}_{\mathbf{v}}$

Powerful syntax extensions: multi-letter indices, index (un-)grouping, ...

[Example] Batched matrix-matrix multiplication

Given $\{A_n\}_{n=1}^N$, $\{B_n\}_{n=1}^N$, compute $\{C_n\}_{n=1}^N = \{A_nB_n\}_{n=1}^N$

C = einsum(A, B, "batch i j, batch j k -> batch i k")

[Example] Improved readability (from an ICML paper Lin et al. [2024])



Powerful syntax extensions: multi-letter indices, index (un-)grouping, ...

[Example] Batched matrix-matrix multiplication

Given $\{A_n\}_{n=1}^N$, $\{B_n\}_{n=1}^N$, compute $\{C_n\}_{n=1}^N = \{A_nB_n\}_{n=1}^N$

C = einsum(A, B, "batch i j, batch j k -> batch i k")

[Example] Improved readability (from an ICML paper Lin et al. [2024])

```
# average channel groups
x = rearrange(x, "b (g c_in) i1 i2 -> b g c_in i1 i2", g=groups)
x = reduce(x, "b g c_in i1 i2 -> b c_in i1 i2", "mean")
x_unfold = F.unfold(x, kernel_size, dilation=dilation, padding=padding, stride=stride)
return rearrange(x_unfold, "b c_in_k1_k2 o1_o2 -> b o1_o2 c_in_k1_k2")
```

[Example] Matrix-matrix-vector-product

Consider $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{3.000 \times 3.000}$ and $\mathbf{v} \in \mathbb{R}^{3.000}$.





[Example] Matrix-matrix-vector-product

Consider $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{3.000 \times 3.000}$ and $\mathbf{v} \in \mathbb{R}^{3.000}$.

$$i - A - j - B - k - \mathbf{v}$$

$$\rightarrow$$
 i-AB-k-V \rightarrow i-ABV



Schedule 1:

A @ B @ v

[Example] Matrix-matrix-vector-product

Consider $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{3.000 \times 3.000}$ and $\mathbf{v} \in \mathbb{R}^{3.000}$. $i - \mathbf{A} - j - \mathbf{B} - k - \mathbf{v}$ $\rightarrow i - \mathbf{A} - \mathbf{B} - k - \mathbf{v}$ $\rightarrow i - \mathbf{A} - \mathbf{B} - \mathbf{v}$ $\rightarrow i - \mathbf{A} - \mathbf{B} - \mathbf{v}$ $\rightarrow i - \mathbf{A} - \mathbf{B} \mathbf{v}$ Schedule 1: $\mathbf{A} \in \mathbf{B} \in \mathbf{v}$ $\mathbf{A} \in (\mathbf{B} \in \mathbf{v})$

[Example] Matrix-matrix-vector-product

Consider $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{3.000 \times 3.000}$ and $\mathbf{v} \in \mathbb{R}^{3.000}$. $i - \mathbf{A} - j - \mathbf{B} - k - \mathbf{v}$ $\rightarrow i - \mathbf{AB} - k - \mathbf{v} \rightarrow i - \mathbf{ABv} \rightarrow i - \mathbf{A} - j - \mathbf{Bv} \rightarrow i - \mathbf{ABv}$ Schedule 1: 1.1 s $A \in B \in \mathbf{v}$ $A \in (B \in \mathbf{v})$

[Example] Matrix-matrix-vector-product

Consider $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{3.000 \times 3.000}$ and $\mathbf{v} \in \mathbb{R}^{3.000}$. $i - \mathbf{A} - j - \mathbf{B} - k - \mathbf{v}$ $\rightarrow i - \mathbf{AB} - k - \mathbf{v} \rightarrow i - \mathbf{ABv} \rightarrow i - \mathbf{A} - j - \mathbf{Bv} \rightarrow i - \mathbf{ABv}$ Schedule 1: 1.1 s $A \in B \in \mathbf{v}$ Schedule 2: 2.2 ms $A \in (B \in \mathbf{v})$

einsum(**"ij,jk,k->i"**, A, B, v)

[Example] Matrix-matrix-vector-product

Consider $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{3.000 \times 3.000}$ and $\mathbf{v} \in \mathbb{R}^{3.000}$. $i - \mathbf{A} - j - \mathbf{B} - \mathbf{k} - \mathbf{v}$ $\rightarrow i - \mathbf{A} \mathbf{B} - \mathbf{k} - \mathbf{v} \rightarrow i - \mathbf{A} \mathbf{B} \mathbf{v} \rightarrow i - \mathbf{A} - j - \mathbf{B} \mathbf{v} \rightarrow i - \mathbf{A} \mathbf{B} \mathbf{v}$ Schedule 1: 1.1 s $A \in B \in \mathbf{v}$ $A \in (B \in \mathbf{v})$

PyTorch: 1.1 s

einsum(**"ij,jk,k->i"**, A, B, v)

[Example] Matrix-matrix-vector-product

Consider $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{3.000 \times 3.000}$ and $\mathbf{v} \in \mathbb{R}^{3.000}$. $i - \mathbf{A} - j - \mathbf{B} - k - \mathbf{v}$ $\rightarrow i - \mathbf{AB} - k - \mathbf{v} \rightarrow i - \mathbf{ABv} \rightarrow i - \mathbf{A} - j - \mathbf{Bv} \rightarrow i - \mathbf{ABv}$ Schedule 1: 1.1 s $A \in B \in \mathbf{v}$ $A \in (B \in \mathbf{v})$ Schedule 2: 2.2 ms $A \in (B \in \mathbf{v})$

PyTorch: 1.1s with update + pip install opt_einsum, 2.3 ms

einsum(**"ij,jk,k->i"**, A, B, v)

Order matters; einsum automatically finds a 'good' schedule.

Function Transformations with Tensor Networks

[Example] Batching/vmap-ing: adding legs

 $(\boldsymbol{A}, \boldsymbol{B}, \boldsymbol{v}) \mapsto \boldsymbol{A} \boldsymbol{B} \boldsymbol{v}$

$$i - A - j - B - k - v$$
 \rightarrow

Function Transformations with Tensor Networks



Function Transformations with Tensor Networks




$$i - \underbrace{\frac{\partial (\mathbf{ABv})}{\partial \mathbf{B}}}_{\mathbf{B}'} = \frac{\partial \left(i - \mathbf{A} - j - \mathbf{B} - k - \mathbf{v}\right)}{\partial \left(j' - \mathbf{B} - k'\right)}$$



$$i - \underbrace{\frac{\partial (\mathbf{A}\mathbf{B}\mathbf{v})}{\partial \mathbf{B}}}_{\mathbf{k}'} = \frac{\partial \left(i - \mathbf{A} - j - \mathbf{B} - k - \mathbf{v}\right)}{\partial \left(j' - \mathbf{B} - k'\right)} = i - \mathbf{A} - j' \qquad k' - \mathbf{v}$$



$$i - \underbrace{\frac{\partial (ABV)}{\partial B}}_{k'} = \frac{\partial \left(i - A - j - B - k - V\right)}{\partial \left(j' - B - k'\right)} = i - A - j' \qquad k' - V$$





$$i - \frac{\partial (\mathbf{ABV})}{\partial \mathbf{B}} = \frac{j'}{k'} = \frac{\partial \left(i - \mathbf{A} - j - \mathbf{B} - k - \mathbf{V}\right)}{\partial \left(j' - \mathbf{B} - k'\right)} = i - \mathbf{A} - j' \quad k' - \mathbf{V}$$
$$i - \frac{\partial (\mathbf{ABV})}{\partial \mathbf{A}} = \frac{j'}{j'} = \frac{\partial \left(i - \mathbf{I} - \mathbf{I} - \mathbf{A} - j - \mathbf{BV}\right)}{\partial \left(i' - \mathbf{A} - j'\right)}$$





einsum is Probably All You Need

\mathbf{V}

Recommendation: Use einsum in your code!

- Better readability (e.g. einops [Rogozhnikov, 2022])
- Automatic optimization (e.g. opt_einsum [Smith and Gray, 2018])

einsum is Probably All You Need

V

Recommendation: Use einsum in your code!

- Better readability (e.g. einops [Rogozhnikov, 2022])
- Automatic optimization (e.g. opt_einsum [Smith and Gray, 2018])

Also (not discussed)

- Automatic distribution (e.g. cotengra [Gray and Kourtis, 2021])
- Randomized/Approximate evaluation (e.g. [Ma et al., 2024])

einsum is Probably All You Need

\mathbf{V}

Recommendation: Use einsum in your code!

- Better readability (e.g. einops [Rogozhnikov, 2022])
- Automatic optimization (e.g. opt_einsum [Smith and Gray, 2018])

Also (not discussed)

- Automatic distribution (e.g. cotengra [Gray and Kourtis, 2021])
- Randomized/Approximate evaluation (e.g. [Ma et al., 2024])

Also (regarding tensor networks)

- Drawing diagrams is more fun
- Simplifications/Transformations via graphical manipulations

Convolutions as Tensor Networks

 \mathbf{V}

- $\begin{array}{rcl} \textbf{X} & \in & \mathbb{R}^{C_{\text{in}} \times I} \\ \textbf{Y} & \in & \mathbb{R}^{C_{\text{out}} \times O} \\ \textbf{W} & \in & \mathbb{R}^{C_{\text{out}} \times C_{\text{in}} \times K} \end{array}$
 - $\mathbf{Y} = \mathbf{X} \star \mathbf{W}$

 $\pmb{Y}=\pmb{W}[\![\pmb{X}]\!]$

* =

=

Convolution as Matrix Multiplication

 \mathbf{V}

 $\begin{array}{rcl} \textbf{X} & \in & \mathbb{R}^{C_{\text{in}} \times I} \\ \textbf{Y} & \in & \mathbb{R}^{C_{\text{out}} \times O} \\ \textbf{W} & \in & \mathbb{R}^{C_{\text{out}} \times C_{\text{in}} \times K} \end{array}$

 $\begin{bmatrix} \mathbf{X} \end{bmatrix} \in \mathbb{R}^{C_{in}K \times O} \\ \mathbf{Y} \in \mathbb{R}^{C_{out} \times O} \\ \mathbf{W} \in \mathbb{R}^{C_{out} \times C_{in}K}$

 $\mathbf{Y} = \mathbf{X} \star \mathbf{W}$

 $oldsymbol{Y} = oldsymbol{W}[\![oldsymbol{X}]\!]$

* =

=

Convolution as Structured Matrix Multiplication

- $\begin{array}{rcl} \mathbf{X} & \in & \mathbb{R}^{C_{\text{in}} \times I} \\ \mathbf{Y} & \in & \mathbb{R}^{C_{\text{out}} \times O} \\ \mathbf{W} & \in & \mathbb{R}^{C_{\text{out}} \times C_{\text{in}} \times K} \end{array}$
 - $\mathbf{Y} = \mathbf{X} \star \mathbf{W}$

 $\mathbf{y} = \mathbf{A}(\mathbf{W})\mathbf{x}$

* = =

Convolution as Structured Matrix Multiplication

7

 $\begin{array}{rcl} \mathbf{X} & \in & \mathbb{R}^{C_{\text{in}} \times I} \\ \mathbf{Y} & \in & \mathbb{R}^{C_{\text{out}} \times O} \\ \mathbf{W} & \in & \mathbb{R}^{C_{\text{out}} \times C_{\text{in}} \times K} \end{array}$

 $\begin{array}{rcl} \textbf{X} & \in & \mathbb{R}^{C_{\text{in}} I} \\ \textbf{y} & \in & \mathbb{R}^{C_{\text{out}} O} \\ \textbf{A}(\textbf{W}) & \in & \mathbb{R}^{C_{\text{out}} O \times C_{\text{in}} I} \end{array}$

 $\mathbf{Y} = \mathbf{X} \star \mathbf{W}$

y = A(W)x

* = =

$$\mathbf{Y}_{\mathbf{c}_{\mathsf{out}},\mathbf{o}} = \sum_{\mathbf{c}_{\mathsf{in}},k} \mathbf{X}_{\mathbf{c}_{\mathsf{in}},i} \quad \mathbf{W}_{\mathbf{c}_{\mathsf{out}},\mathbf{c}_{\mathsf{in}},k}$$

$$Y_{c_{\text{out}},o} = \sum_{c_{\text{in}},k} X_{c_{\text{in}},i(k,o)} W_{c_{\text{out}},c_{\text{in}},k}$$

Index pattern $\Pi(I, K, S, P, D)$ captures the convolution's connectivity

$$\begin{split} \mathbf{Y}_{\mathbf{c}_{\text{out}},\mathbf{o}} &= \sum_{\mathbf{c}_{\text{in}},k} \mathbf{X}_{\mathbf{c}_{\text{in}},i(k,o)} \mathbf{W}_{\mathbf{c}_{\text{out}},\mathbf{c}_{\text{in}},k} \\ &= \sum_{\mathbf{c}_{\text{in}},k} \sum_{i} \mathbf{X}_{\mathbf{c}_{\text{in}},i} \Pi_{i,o,k} \mathbf{W}_{\mathbf{c}_{\text{out}},\mathbf{c}_{\text{in}},k} \end{split}$$

Index pattern $\Pi(I, K, S, P, D)$ captures the convolution's connectivity

Х

$$\begin{split} Y_{c_{\text{out}},o} &= \sum_{\mathbf{c}_{\text{in}},k} X_{c_{\text{in}},i(k,o)} W_{c_{\text{out}},c_{\text{in}},k} \\ &= \sum_{\mathbf{c}_{\text{in}},k} \sum_{i} X_{c_{\text{in}},i} \Pi_{i,o,k} W_{c_{\text{out}},c_{\text{in}},k} \end{split}$$

 \star

W

[Π]_{:,:,k}

γ

Index pattern $\Pi(I, K, S, P, D)$ captures the convolution's connectivity

Х

$$Y_{c_{\text{out}},o} = \sum_{c_{\text{in}},k} X_{c_{\text{in}},i(k,o)} W_{c_{\text{out}},c_{\text{in}},k}$$
$$= \sum_{c_{\text{in}},k} \sum_{i} X_{c_{\text{in}},i} \prod_{i,o,k} W_{c_{\text{out}},c_{\text{in}},k}$$



*

W

[**П**]_{:,:,*k*}

γ

Higher-dimensional Convolutions [Hayashi et al., 2019]



1d

Higher-dimensional Convolutions [Hayashi et al., 2019]



1d

2d

cout

Higher-dimensional Convolutions [Hayashi et al., 2019]



1d

2d

3d









Х







01



C_{in} – W

Cout

















Conv.

w

- Court





Conv.

w - cout



Fisher/GGN block



Fisher/GGN block



Fisher/GGN diagonal



Fisher/GGN block



Fisher/GGN diagonal



GGN Gram/Empirical NTK matrix



Fisher/GGN block



Fisher/GGN diagonal



GGN Gram/Empirical NTK matrix



Fisher/GGN mini-block diagonal



... and More!



KFC/KFAC-expand

KFAC-reduce





... and More!



KFC/KFAC-expand

KFAC-reduce





Approximate Hessian diagonal

[Elsayed, Farrahi, Dangel, and Mahmood, 2024, ICML]



Check out the table in the paper's appendix!
Example & Conclusion



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)

 $\mathbf{V}_{\mathbf{v}}$

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 ^{−1} (1.00 x)	3.62 (1.00 x)
VGG19	SGD	6.90 · 10 ⁻¹ (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 SINGD	$ \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 (1.25 x)
VGG19	SGD SINGD	6.90 · 10 ⁻¹ (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 Singd Singd+Tn	$\begin{array}{c c} 1.17 \cdot 10^{-1} & (1.00 \text{ x}) \\ 1.94 \cdot 10^{-1} & (1.67 \text{ x}) \\ 1.56 \cdot 10^{-1} & (1.33 \text{ x}) \end{array}$	3.62 (1.00 x) 4.54 (1.25 x) 3.63 (1.00 x)
VGG19	SGD Singd Singd+TN	6.90 · 10 ⁻¹ (1.00 x) 1.02 (1.48 x) 8.01 · 10 ⁻¹ (1.16 x)	14.1 (1.00 x) 32.1 (2.28 x) 16.1 (1.14 x)

(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

_____****

- + 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





- Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. Advances in neural information processing systems (NeurIPS), 2019.
- Suzanna Becker and Yann Lecun. Improving the convergence of back-propagation learning with second-order methods. 1989.
- Aleksandar Botev, Hippolyt Ritter, and David Barber. Practical Gauss-Newton optimisation for deep learning. In International Conference on Machine Learning (ICML), 2017.
- Felix Dangel, Frederik Kunstner, and Philipp Hennig. BackPACK: Packing more into backprop. In International Conference on Learning Representations (ICLR), 2020.

Mohamed Elsayed and A. Rupam Mahmood. HesScale: Scalable computation of hessian diagonals. 2023.

Mohamed Elsayed, Homayoon Farrahi, Felix Dangel, and A. Rupam Mahmood. Revisiting scalable hessian diagonal approximations for applications in reinforcement learning. In International Conference on Machine Learning (ICML), 2024.

References II



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.

Donald Goldfarb, Yi Ren, and Achraf Bahamou. Practical quasi-newton methods for training deep neural networks, 2021.

Johnnie Gray and Stefanos Kourtis. Hyper-optimized tensor network contraction. Quantum, 2021.

- Roger Grosse and James Martens. A kronecker-factored approximate Fisher matrix for convolution layers. In **International Conference on Machine Learning (ICML)**, 2016.
- Kohei Hayashi, Taiki Yamaguchi, Yohei Sugawara, and Shin-ichi Maeda. Exploring unexplored tensor network decompositions for convolutional neural networks. In Advances in Neural Information Processing Systems (NeurIPS), 2019.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks, 2020.
- Sören Laue, Matthias Mitterreiter, and Joachim Giesen. A simple and efficient tensor calculus. In AAAI Conference on Artificial Intelligence, 2020.

References III

- \mathbf{V}
- 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.
- Linjian Ma, Matthew Fishman, Miles Stoudenmire, and Edgar Solomonik. Approximate contraction of arbitrary tensor networks with a flexible and efficient density matrix algorithm. 2024.
- James Martens and Roger Grosse. Optimizing neural networks with Kronecker-factored approximate curvature. In International Conference on Machine Learning (ICML), 2015.
- Hannah Pinson, Joeri Lenaerts, and Vincent Ginis. Linear cnns discover the statistical structure of the dataset using only the most dominant frequencies. In **International Conference on Machine Learning** (ICML), 2023.
- Yi Ren, Achraf Bahamou, and Donald Goldfarb. Kronecker-factored quasi-newton methods for deep learning, 2022.
- Alex Rogozhnikov. Einops: Clear and reliable tensor manipulations with einstein-like notation. In International Conference on Learning Representations (ICLR), 2022.



- Andrew M. Saxe, James L. McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks, 2014.
- R. M. Schmidt, F. Schneider, and P. Hennig. Descending through a crowded valley benchmarking deep learning optimizers. In International Conference on Machine Learning (ICML), 2021.
- Sidak Pal Singh, Gregor Bachmann, and Thomas Hofmann. Analytic insights into structure and rank of neural network hessian maps. Advances in Neural Information Processing Systems (NeurIPS), 2021.
- Sidak Pal Singh, Thomas Hofmann, and Bernhard Schölkopf. The hessian perspective into the nature of convolutional neural networks. 2023.
- Daniel G. A. Smith and Johnnie Gray. opt_einsum A python package for optimizing contraction order for einsum-like expressions. Journal of Open Source Software (JOSS), 2018.



For structured patterns, we can re-wire the TN before evaluation \longrightarrow even better performance



For structured patterns, we can re-wire the TN before evaluation \longrightarrow even better performance

Dense convolution (K = S)



For structured patterns, we can re-wire the TN before evaluation \longrightarrow even better performance

Dense convolution (K = S**)**





For structured patterns, we can re-wire the TN before evaluation \longrightarrow even better performance

Dense convolution (K = S)

Down-sampling convolution (S > K**)**





For structured patterns, we can re-wire the TN before evaluation \longrightarrow even better performance

Dense convolution (K = S)

Down-sampling convolution (S > K**)**

