

Title: Deep Learning Convolutions Through the Lens of Tensor Networks

Speakers: Felix Dangel

Series: Machine Learning Initiative

Date: December 01, 2023 - 2:30 PM

URL: <https://pirsa.org/23120027>

Abstract: Despite their simple intuition, convolutions are more tedious to analyze than dense layers, which complicates the transfer of theoretical and algorithmic ideas. We provide a simplifying perspective onto convolutions through tensor networks (TNs) which allow reasoning about the underlying tensor multiplications by drawing diagrams, and manipulating them to perform function transformations and sub-tensor access. We demonstrate this expressive power by deriving the diagrams of various autodiff operations and popular approximations of second-order information with full hyper-parameter support, batching, channel groups, and generalization to arbitrary convolution dimensions. Further, we provide convolution-specific transformations based on the connectivity pattern which allow to re-wire and simplify diagrams before evaluation. Finally, we probe computational performance, relying on established machinery for efficient TN contraction. Our TN implementation speeds up a recently-proposed KFAC variant up to 4.5x and enables new hardware-efficient tensor dropout for approximate backpropagation.

Zoom link <https://pitp.zoom.us/j/99090845943?pwd=NHBNVbVhbnRNSOGNSVzNGS21xc1lFdz09>

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 layer output
            ([N, *, out_features]) to which the transposed output-input Jacobian
            is applied. Has shape [V, N, *, out_features] 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->v{' if sum_batch else 'n'}oi"
    return einsum(equation, mat, d_weight)
```

```
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 layer output
            ([N, *, out_features]) to which the transposed output-input Jacobian
            is applied. Has shape [V, N, *, out_features] 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)

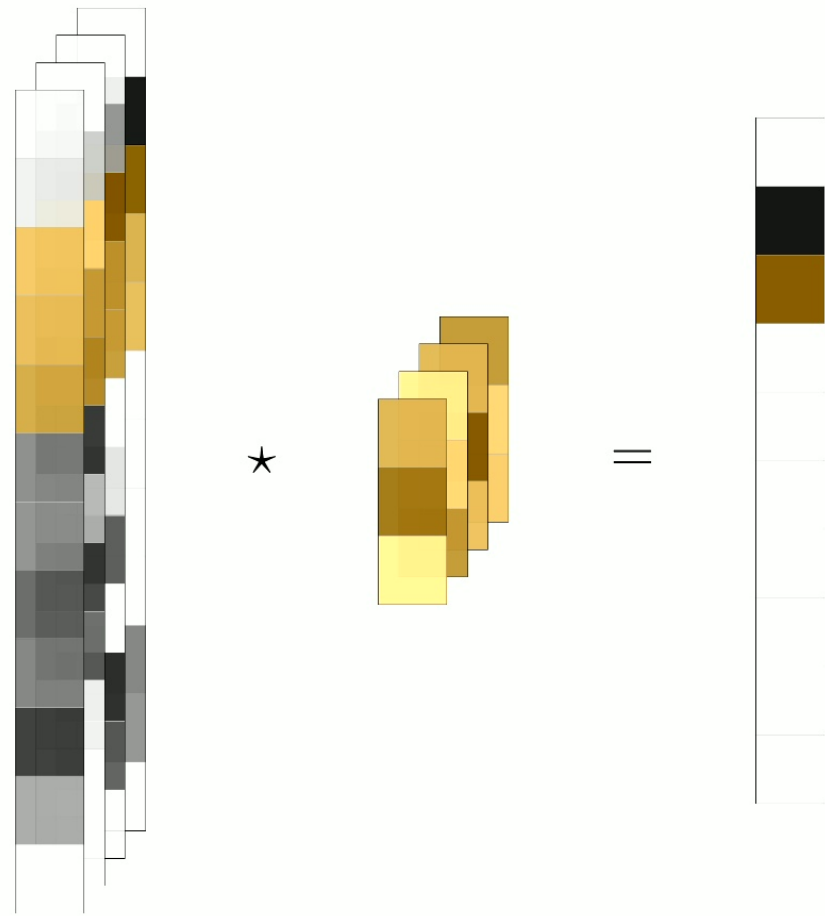
    equation = f"vn...o,n...i->v{' if sum_batch else 'n'}oi"
    return einsum(equation, mat, d_weight)
```

C_{in} is for Input Channels



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

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



A Mathematical Definition of Tensor Multiplication [e.g. Laue et al., 2020]



General case: Given N tensors $\mathbf{A}_1, \dots, \mathbf{A}_N$ with index tuples $\mathbf{S}_{\mathbf{A}_1}, \dots, \mathbf{S}_{\mathbf{A}_N}$

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

indices not present in the output are summed out; $\mathbf{S}_{\mathbf{C}}$ satisfies $\mathbf{S}_{\mathbf{C}} \subseteq \mathbf{S}_{\mathbf{A}_1} \cup \dots \cup \mathbf{S}_{\mathbf{A}_N}$.

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

$$\begin{aligned} \mathbf{S}_{\mathbf{A}} &= (i, j) \\ \mathbf{S}_{\mathbf{B}} &= (j, k) \\ \mathbf{S}_{\mathbf{C}} &= (i, k) \end{aligned} \quad \Rightarrow \quad (\mathbf{S}_{\mathbf{A}} \cup \mathbf{S}_{\mathbf{B}}) \setminus \mathbf{S}_{\mathbf{C}} = (i, j, k) \setminus (i, k) = (j)$$

$$\Rightarrow \mathbf{C} = *_{((i,j),(j,k),(i,k))}(\mathbf{A}, \mathbf{B}), \quad \text{or} \quad \mathbf{C} = \text{einsum}(\text{"ij,jk->ik"}, \mathbf{A}, \mathbf{B}).$$

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

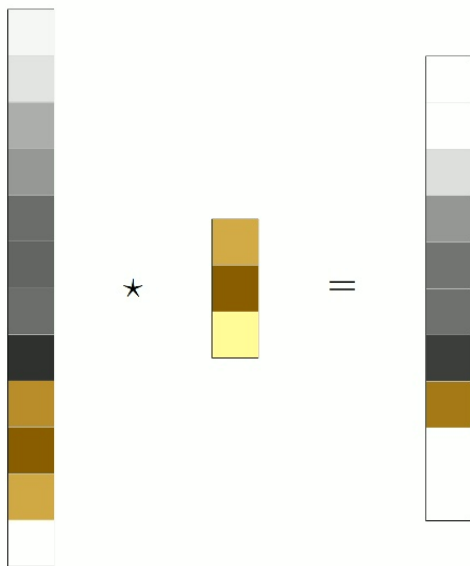
Convolution as Matrix Multiplication



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

$$\begin{aligned} [\mathbf{X}] &\in \mathbb{R}^{C_{in} K \times O} \\ \mathbf{Y} &\in \mathbb{R}^{C_{out} \times O} \\ \mathbf{W} &\in \mathbb{R}^{C_{out} \times C_{in} K} \end{aligned}$$

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



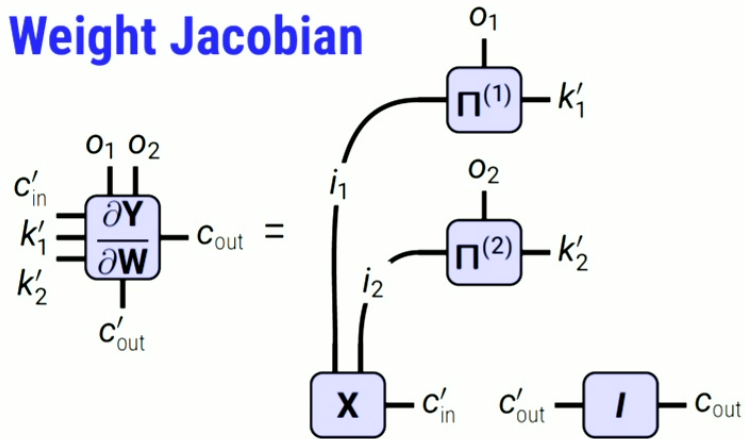
$$\mathbf{Y} = \mathbf{W} [\mathbf{X}]$$



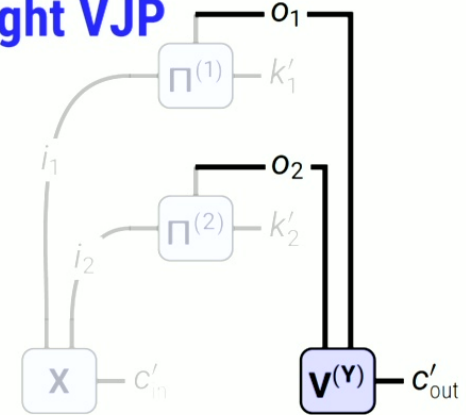
Derivatives & Autodiff Routines



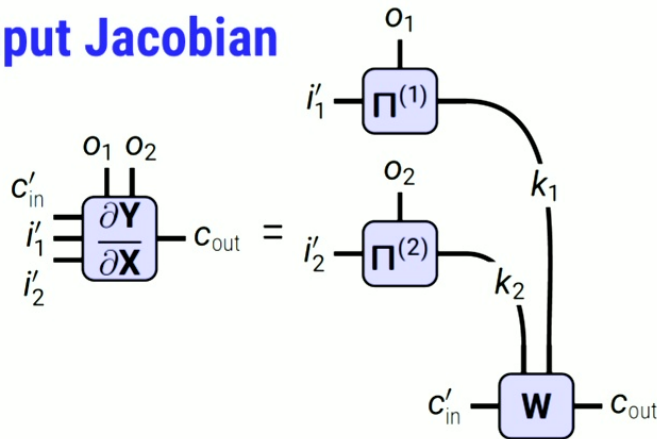
Weight Jacobian



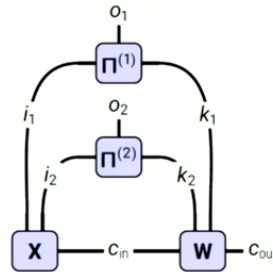
Weight VJP



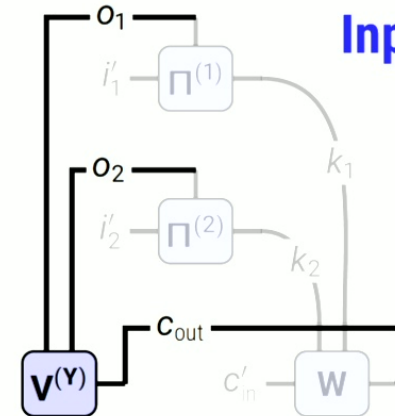
Input Jacobian



Conv.



Input VJP



Convolutions Through the Lens of Tensor Networks

Felix Dangel

December 01, 2023



Convolutions Are More Tedious than Linear Layers



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

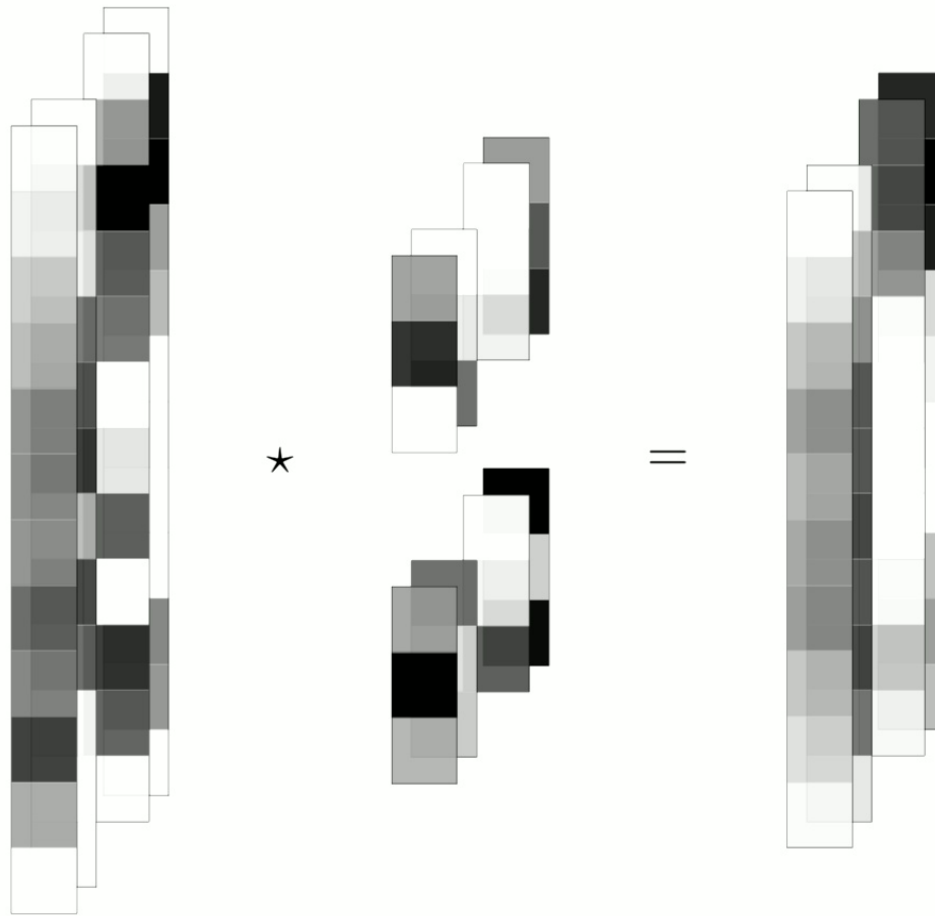
Concept	for MLPs	for CNNs
Approximate Hessian diagonal	1989	2023
Kronecker-factored curvature (KFAC, KFRA, KFLR)	2015, 2017, 2017	2016, 2020, 2020
Kronecker-factored quasi-Newton methods (KBFGS)	2021	2022
Neural tangent kernel (NTK)	2018	2019
Hessian rank	2021	2023

G is for Groups



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

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



Tensor Networks: Graphical Notation for Tensor Multiplications



Operation	Diagram
Scalar	
Vector	
Matrix	
Tensor	
Matricize	
Flatten	
Trace	

Operation	Diagram

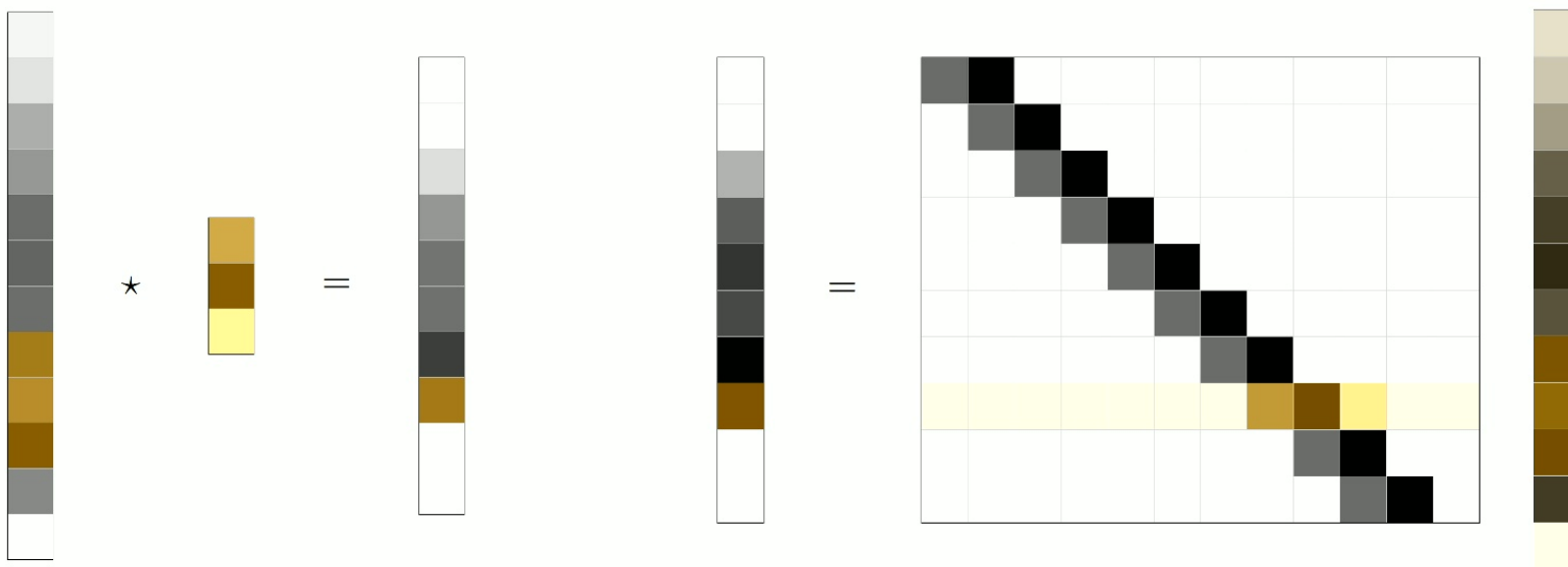
Convolution as Structured Matrix Multiplication



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

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

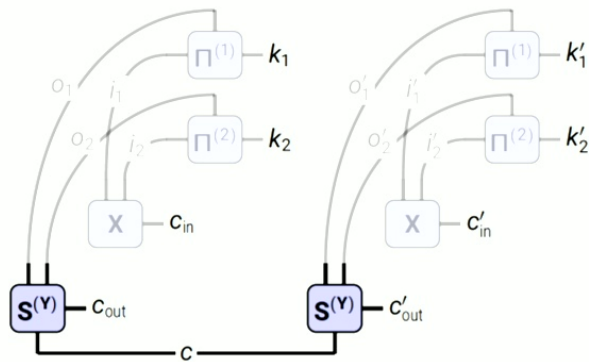
$$\mathbf{y} = \mathbf{A}(\mathbf{W})\mathbf{x}$$



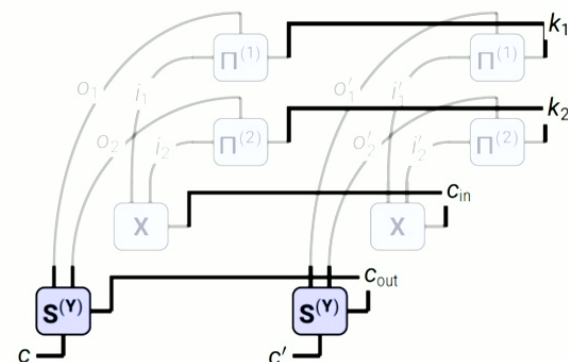
Higher-order Derivatives & Approximations



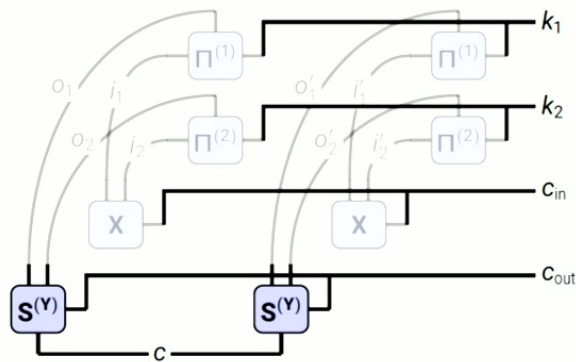
Fisher/GGN block



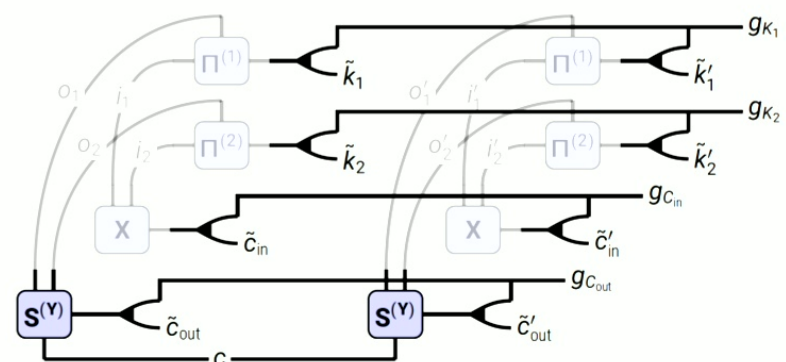
GGN Gram/Empirical NTK matrix



Fisher/GGN diagonal



Fisher/GGN mini-block diagonal



Convolutions ~~Through the Lens of~~ ~~Tensor Networks~~ as einsum

Felix Dangel

December 01, 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

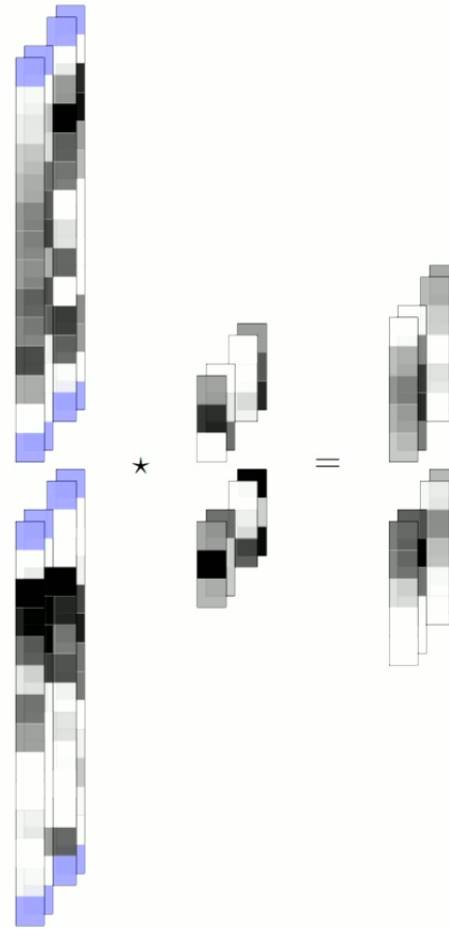
This talk: A simplifying perspective of convolutions through tensor networks

Putting Everything Together in One Dimension...



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

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



Tensor Networks: Graphical Notation for Tensor Multiplications



Operation	Diagram
Scalar	
Vector	
Matrix	
Tensor	
Matricize	
Flatten	
Trace	

Operation	Diagram
$i \rightarrow AB \rightarrow k$	
$i \rightarrow A \odot B \rightarrow j$	
$(i, k) \rightarrow A \otimes B \rightarrow (j, l)$	
$\text{diag}(\mathbf{A}) \rightarrow i$	
$i \rightarrow \text{diag}(\mathbf{a}) \rightarrow i$	

Convolution as Structured Matrix Multiplication

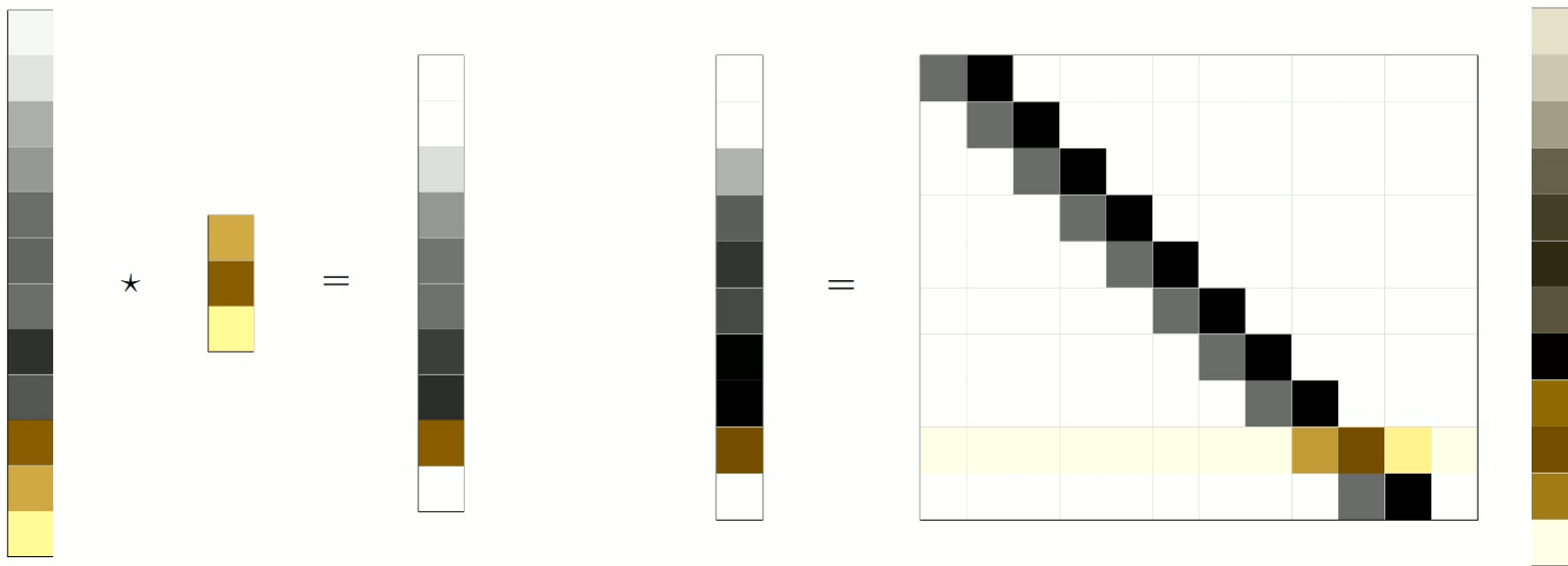


$$\begin{aligned} \mathbf{X} &\in \mathbb{R}^{C_{in} \times I} \\ \mathbf{Y} &\in \mathbb{R}^{C_{out} \times O} \\ \mathbf{W} &\in \mathbb{R}^{C_{out} \times C_{in} \times K} \end{aligned}$$

$$\begin{aligned} \mathbf{x} &\in \mathbb{R}^{C_{in} I} \\ \mathbf{y} &\in \mathbb{R}^{C_{out} O} \\ \mathbf{A}(\mathbf{W}) &\in \mathbb{R}^{C_{out} O \times C_{in} I} \end{aligned}$$

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

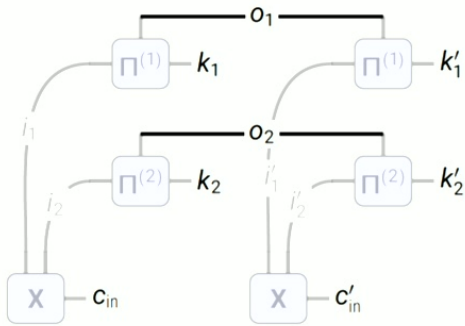
$$\mathbf{y} = \mathbf{A}(\mathbf{W})\mathbf{x}$$



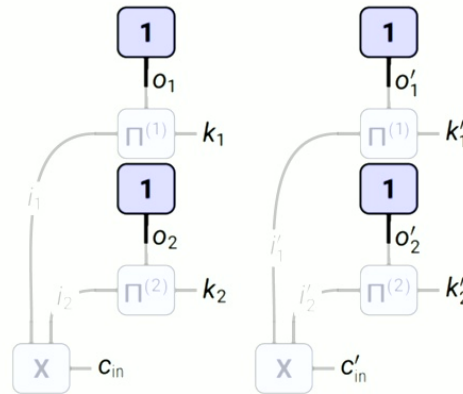
... and More!



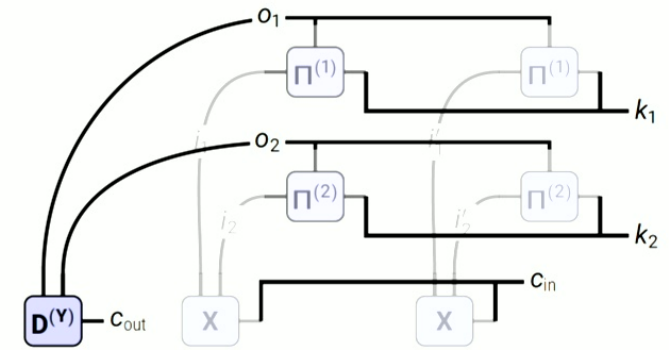
KFC/KFAC-expand



KFAC-reduce



Approximate Hessian diagonal



Check out the table in the paper's appendix!

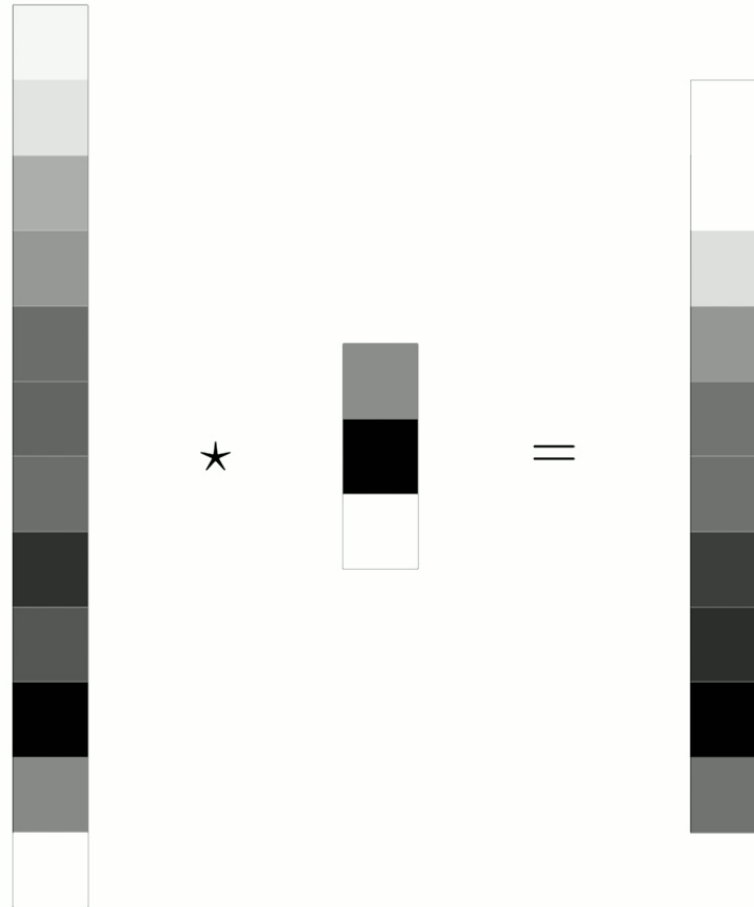
The Broader Picture

A Simple Convolution



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

- ★ $\mathbf{X} \in \mathbb{R}^l$
- ★ $\mathbf{W} \in \mathbb{R}^k$
- ★ $\mathbf{Y} \in \mathbb{R}^o$

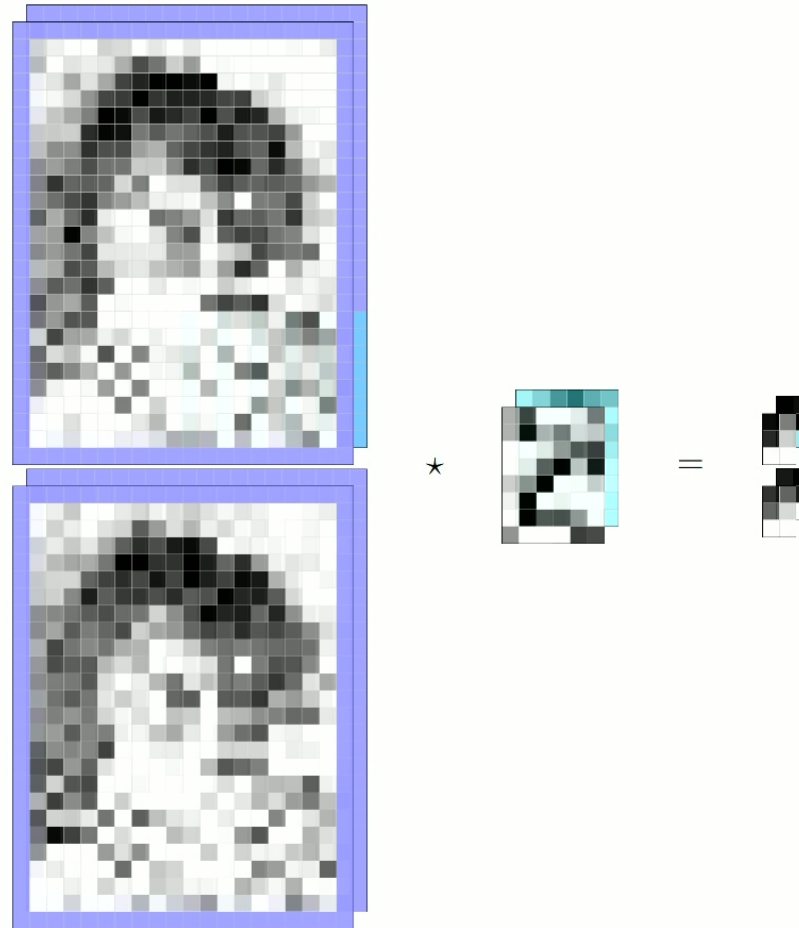


...and in Two Dimensions



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

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



Benefits of TN/einsum Abstraction

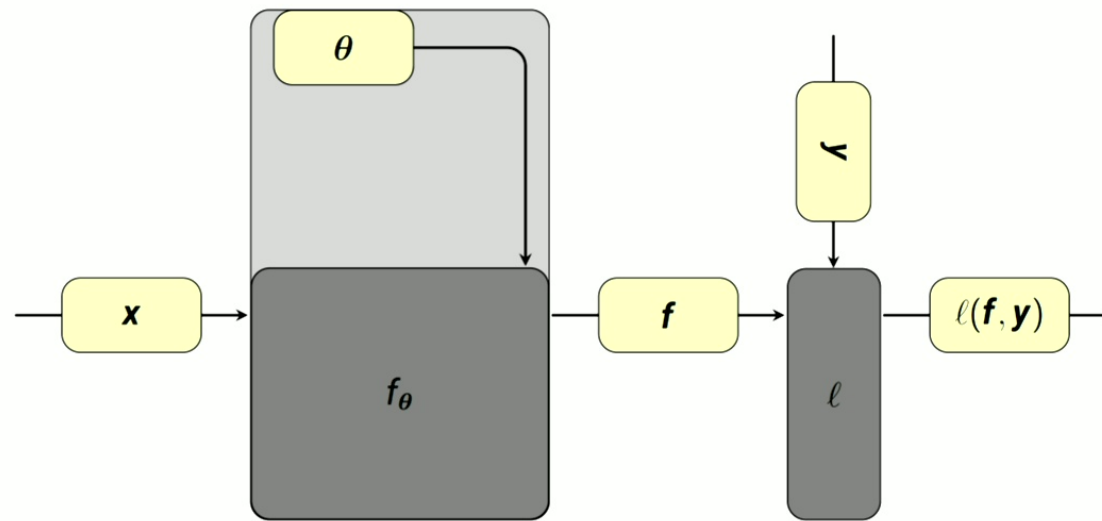
Convolution as Tensor Network [Hayashi et al., 2019]



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

Example & Conclusion

Supervised Deep Learning in a Nutshell (Math)



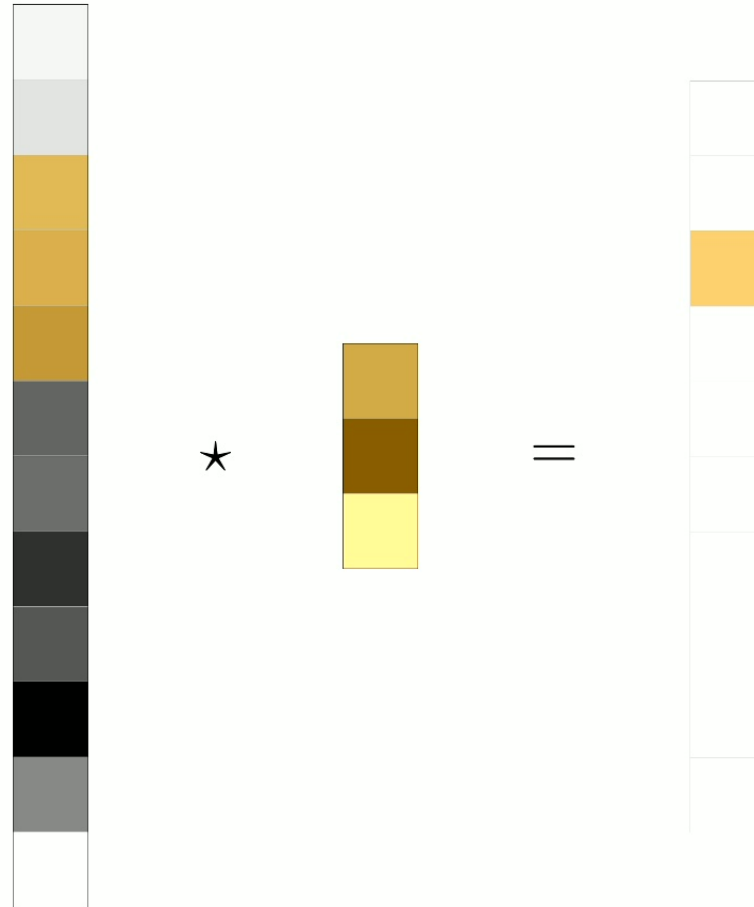
$$l(f_\theta(\mathbf{x}), \mathbf{y})$$

A Simple Convolution



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

- ★ $\mathbf{X} \in \mathbb{R}^l$
- ★ $\mathbf{W} \in \mathbb{R}^k$
- ★ $\mathbf{Y} \in \mathbb{R}^o$

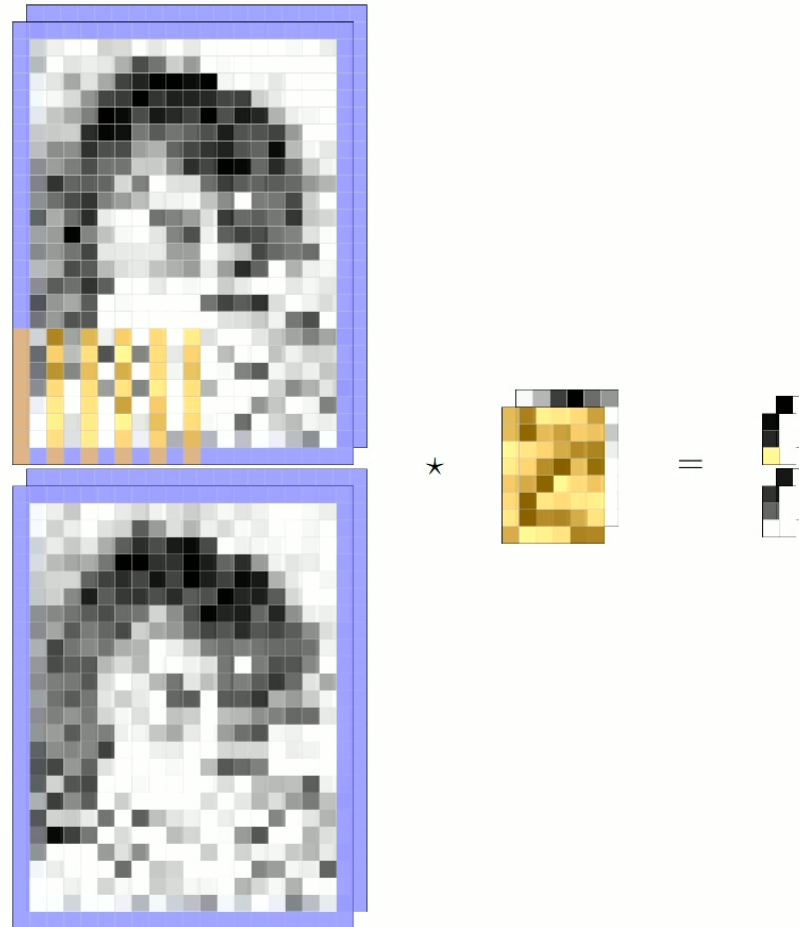


...and in Two Dimensions



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

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



Improved einsum Syntax: einops [Rogozhnikov, 2022, ICLR]



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

[Example] Batched matrix-matrix multiplication

Given $\{\mathbf{A}_n\}_{n=1}^N$, $\{\mathbf{B}_n\}_{n=1}^N$, compute $\{\mathbf{C}_n\}_{n=1}^N = \{\mathbf{A}_n \mathbf{B}_n\}_{n=1}^N$

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

[Example] Improved readability (from an ICLR 2024 submission)

```
1 x = x.unfold(2, kernel_size[0], stride[0])
2 x = x.unfold(3, kernel_size[1], stride[1])
3 x = x.transpose_(1, 2).transpose_(2, 3)
4 return torch.mean(
5     x.reshape((x.size(0), x.size(1), x.size(2), groups, -1, x.size(4), x.size(5))),
6     3,
7 ).view(x.size(0), x.size(1), x.size(2), -1)
```

Convolution as Tensor Network [Hayashi et al., 2019]



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

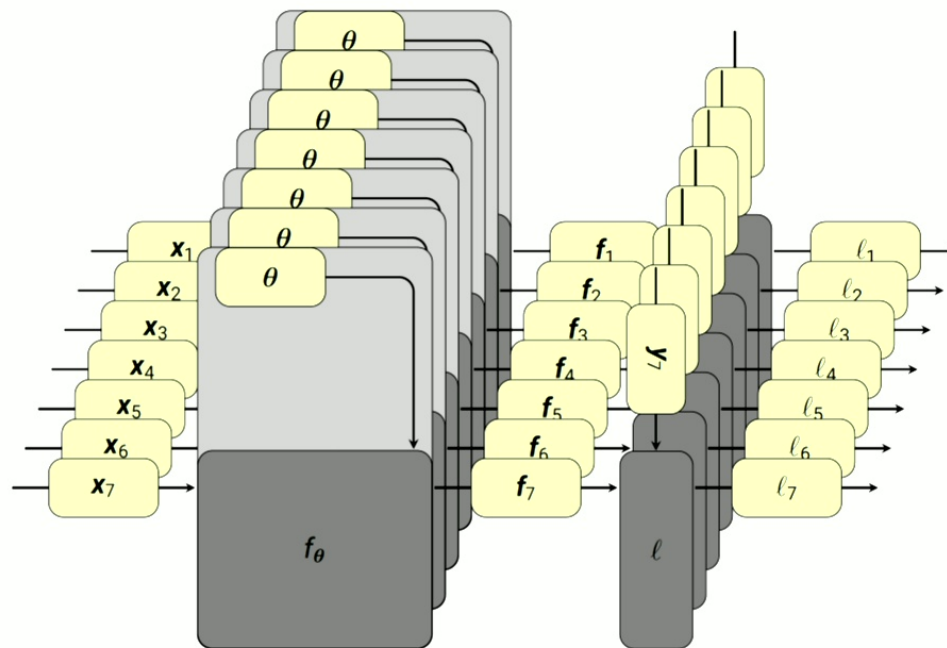
$$\begin{aligned} Y_{c_{\text{out}},o} &= \sum_{c_{\text{in}},k} X_{c_{\text{in}},i(k,o)} W_{c_{\text{out}},c_{\text{in}},k} \\ &= \sum_{c_{\text{in}},k} \sum_i X_{c_{\text{in}},i} \Pi_{i,o,k} W_{c_{\text{out}},c_{\text{in}},k} \end{aligned}$$



State of the art

X

Supervised Deep Learning in a Nutshell (Math)



Minimize $\mathcal{L}_{\mathbb{D}}(\boldsymbol{\theta})$ with

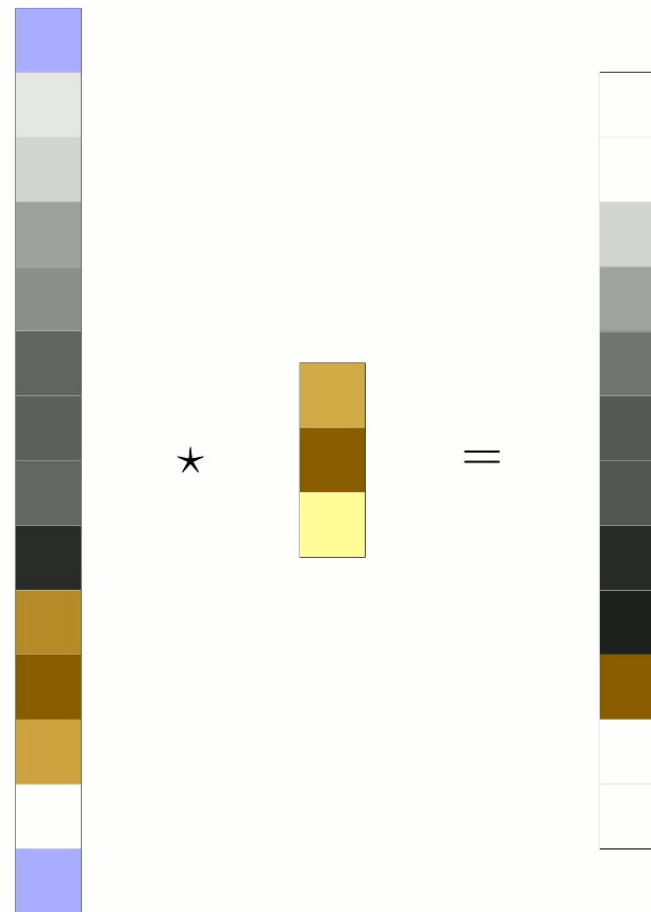
$$\mathcal{L}_{\mathbb{D}}(\boldsymbol{\theta}) = \frac{1}{|\mathbb{D}|} \sum_{(\mathbf{x}_n, \mathbf{y}_n) \in \mathbb{D}} \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_n), \mathbf{y}_n)$$

P is for Padding



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

- ★ $\mathbf{X} \in \mathbb{R}^l$
- ★ $\mathbf{W} \in \mathbb{R}^k$
- ★ $\mathbf{Y} \in \mathbb{R}^o$

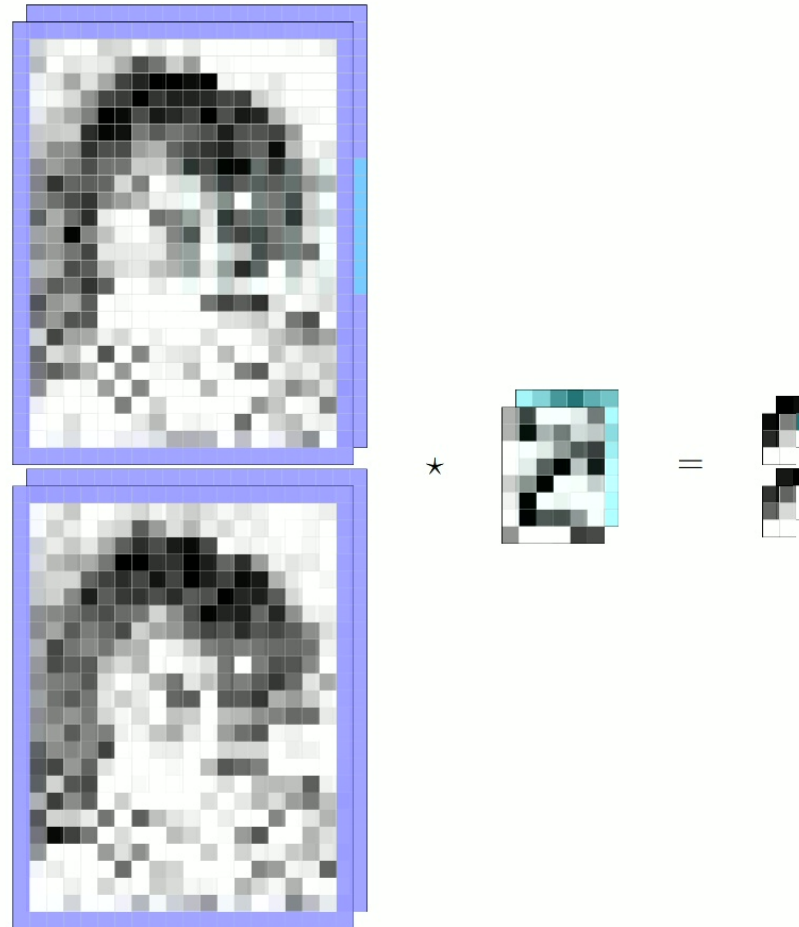


...and in Two Dimensions



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

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





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

[Example] Batched matrix-matrix multiplication

Given $\{\mathbf{A}_n\}_{n=1}^N$, $\{\mathbf{B}_n\}_{n=1}^N$, compute $\{\mathbf{C}_n\}_{n=1}^N = \{\mathbf{A}_n \mathbf{B}_n\}_{n=1}^N$

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

[Example] Improved readability (from an ICLR 2024 submission)

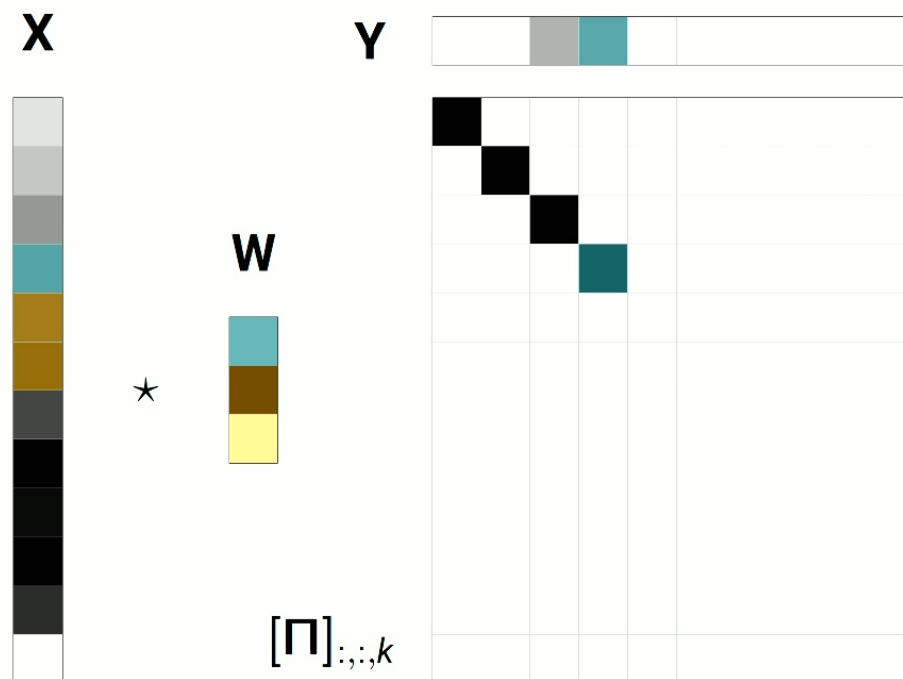
```
1 # average channel groups
2 x = rearrange(x, "b (g c_in) i1 i2 -> b g c_in i1 i2", g=groups)
3 x = reduce(x, "b g c_in i1 i2 -> b c_in i1 i2", "mean")
4
5 x_unfold = F.unfold(x, kernel_size, dilation=dilation, padding=padding, stride=stride)
6 return rearrange(x_unfold, "b c_in_k1_k2 o1_o2 -> b o1_o2 c_in_k1_k2")
```

Convolution as Tensor Network [Hayashi et al., 2019]



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

$$\begin{aligned}
 Y_{c_{out},o} &= \sum_{c_{in},k} X_{c_{in},i(k,o)} W_{c_{out},c_{in},k} \\
 &= \sum_{c_{in},k} \sum_i X_{c_{in},i} \Pi_{i,o,k} W_{c_{out},c_{in},k}
 \end{aligned}$$



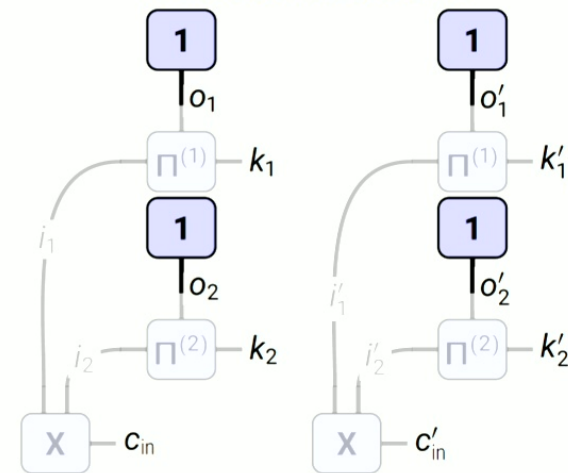


State of the art

$$\begin{aligned} \mathbf{X} &\rightarrow \llbracket \mathbf{X} \rrbracket \\ &\rightarrow \mathbf{1}^\top \llbracket \mathbf{X} \rrbracket \\ &\rightarrow (\mathbf{1}^\top \llbracket \mathbf{X} \rrbracket)^\top (\mathbf{1}^\top \llbracket \mathbf{X} \rrbracket) \end{aligned}$$

Time: **9.87 ms**

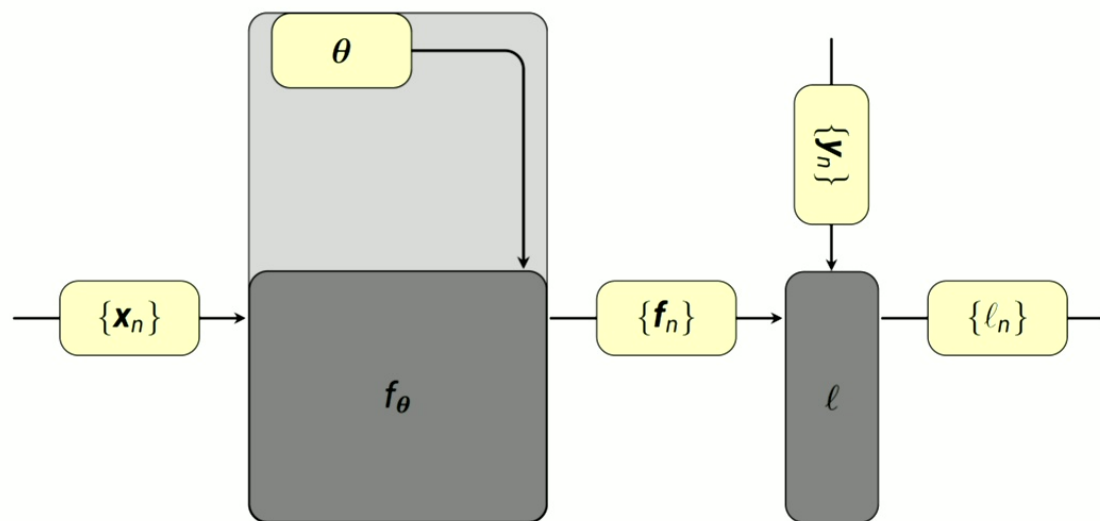
Tensor Network



Time: **2.69 ms (3.7 x)**

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

Supervised Deep Learning in a Nutshell (Math)



Minimize $\mathcal{L}_{\mathbb{D}}(\theta)$ with

$$\mathcal{L}_{\mathbb{B}}(\theta) = \frac{1}{|\mathbb{B}|} \sum_{(\mathbf{x}_n, \mathbf{y}_n) \in \mathbb{B}} \ell(f_\theta(\mathbf{x}_n), \mathbf{y}_n) \quad \text{where} \quad \mathbb{B} \sim \mathbb{D},$$

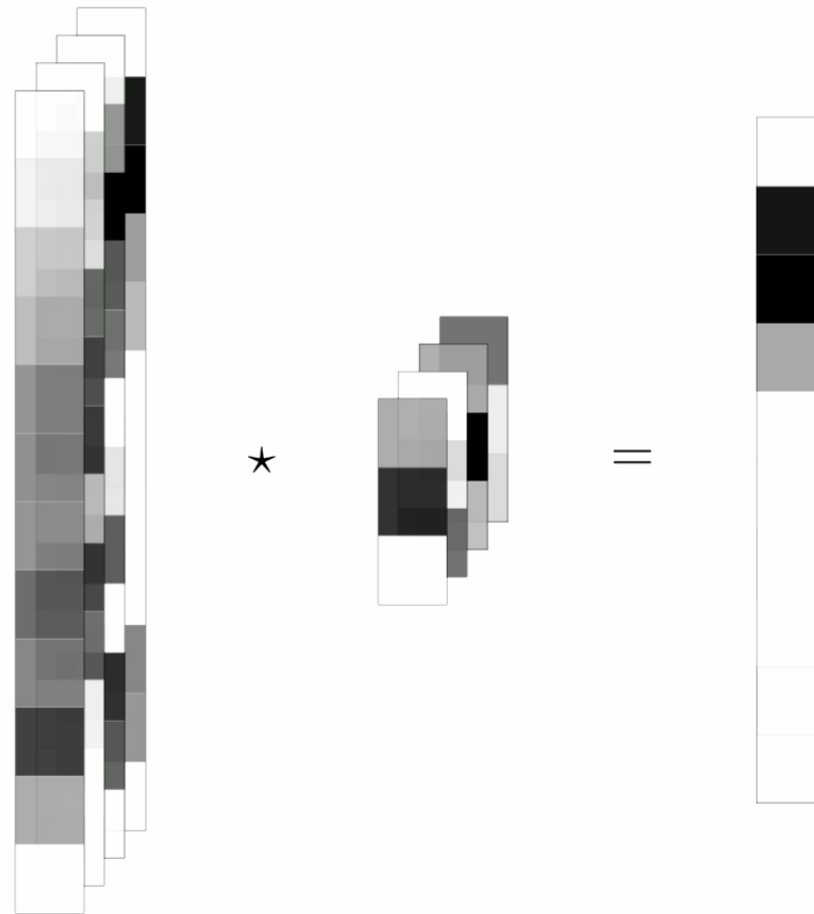
1

C_{in} is for Input Channels



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

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



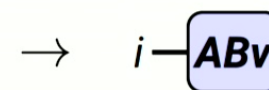
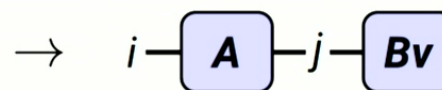
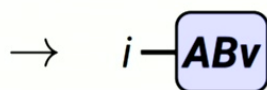
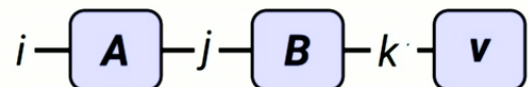
Tensor Multiplication & Tensor Networks

Contraction Path Optimization (e.g. `opt_einsum` [Smith and Gray, 2018])



[Example] Matrix-matrix-vector-product

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



Schedule 1: **1.1 s**

`A @ B @ v`

Schedule 2: **2.2 ms**

`A @ (B @ v)`

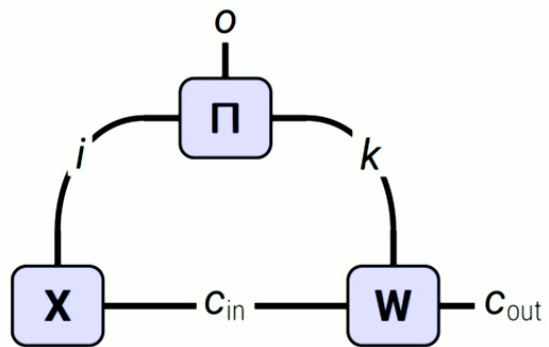
PyTorch: **1.1 s**

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.

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



1d

Non-conventional operations can be much cheaper with tensor networks

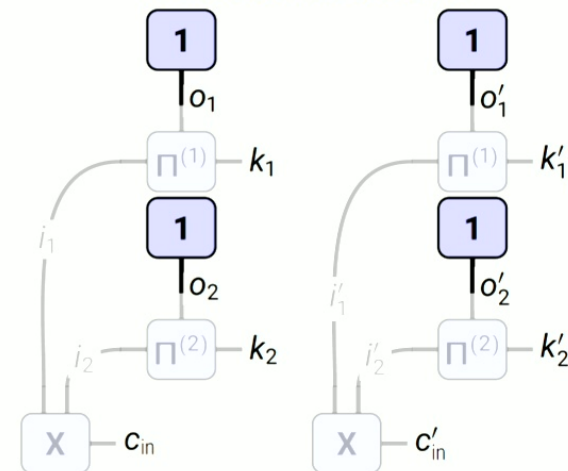
State of the art

$$\begin{aligned} \mathbf{X} &\rightarrow \llbracket \mathbf{X} \rrbracket \\ &\rightarrow \mathbf{1}^\top \llbracket \mathbf{X} \rrbracket \\ &\rightarrow (\mathbf{1}^\top \llbracket \mathbf{X} \rrbracket)^\top (\mathbf{1}^\top \llbracket \mathbf{X} \rrbracket) \end{aligned}$$

Time: **9.87 ms**

Extra memory: **3.07 GiB**

Tensor Network

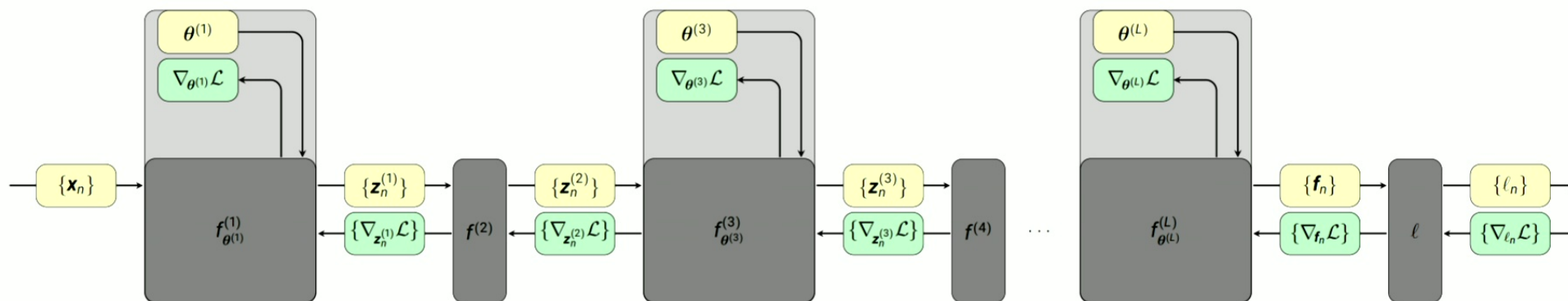


Time: **2.69 ms (3.7 x)**

Extra memory: **0 MiB**

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

Supervised Deep Learning in a Nutshell (Math)



Minimize $\mathcal{L}_{\mathbb{D}}(\theta)$ with

$$\mathcal{L}_{\mathbb{B}}(\theta) = \frac{1}{|\mathbb{B}|} \sum_{(\mathbf{x}_n, \mathbf{y}_n) \in \mathbb{B}} \ell(f_{\theta}(\mathbf{x}_n), \mathbf{y}_n)$$

where

$$\mathbb{B} \sim \mathbb{D},$$

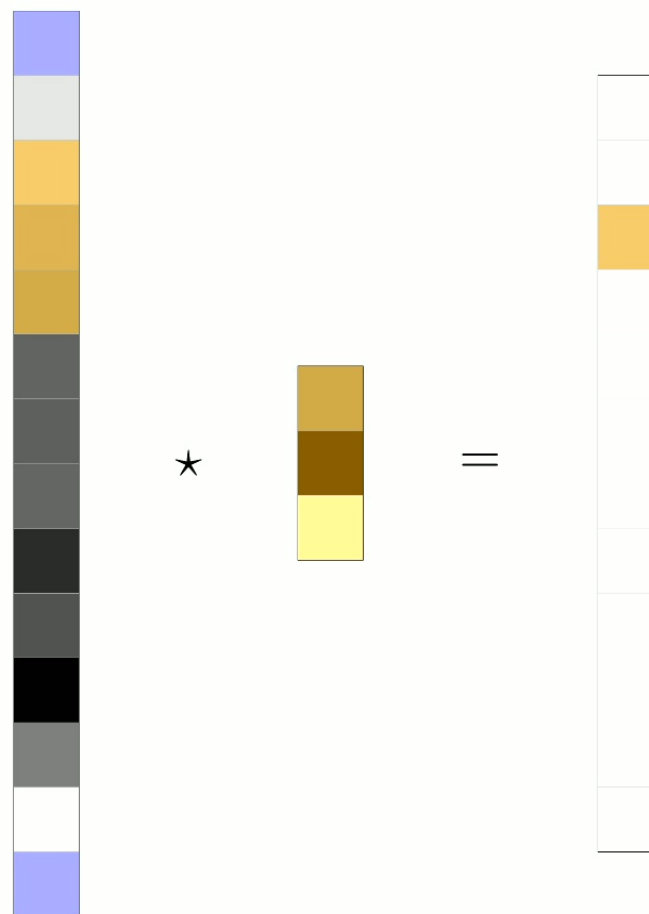
$$f_{\theta} = f_{\theta^{(L)}} \circ f_{\theta^{(L-1)}} \circ \dots \circ f_{\theta^{(1)}}.$$

P is for Padding



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

- ★ $\mathbf{X} \in \mathbb{R}^l$
- ★ $\mathbf{W} \in \mathbb{R}^k$
- ★ $\mathbf{Y} \in \mathbb{R}^o$



Tensor Multiplication Unifies Vector & Matrix Multiplications



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

Consider two vectors \mathbf{a}, \mathbf{b} and two matrices \mathbf{A}, \mathbf{B}

Product	Notation	Index notation
Inner	$\mathbf{c} = \mathbf{a}^\top \mathbf{b}$	$c = \sum_i a_i b_i$
Element-wise	$\mathbf{c} = \mathbf{a} \odot \mathbf{b}$	$c_i = a_i b_i$
Outer/Kronecker	$\mathbf{C} = \mathbf{a} \mathbf{b}^\top = \mathbf{a} \otimes \mathbf{b}^\top$	$C_{i,j} = a_i b_j$
Inner	$\mathbf{C} = \mathbf{A} \mathbf{B}$	$C_{i,k} = \sum_j A_{i,j} B_{j,k}$
Element-wise	$\mathbf{C} = \mathbf{A} \odot \mathbf{B}$	$C_{i,j} = A_{i,j} B_{i,j}$
Kronecker	$\mathbf{C} = \mathbf{A} \otimes \mathbf{B}$	$C_{(i,k),(j,l)} = A_{i,j} B_{k,l}$

Function Transformations with Tensor Networks



[Example] Batching/vmap-ing: adding legs

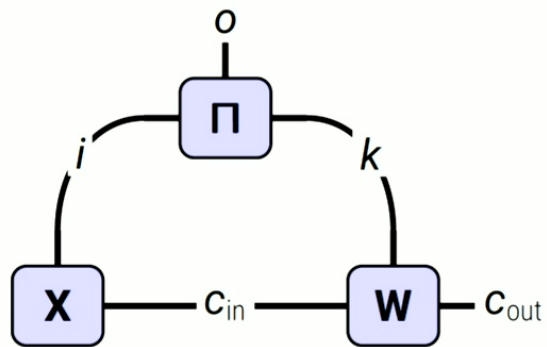
$$(\mathbf{A}, \mathbf{B}, \mathbf{v}) \mapsto \mathbf{ABv} \qquad \{(\mathbf{A}_n, \mathbf{B}_n, \mathbf{v}_n)\}_{n=1}^N \mapsto \{\mathbf{A}_n \mathbf{B}_n \mathbf{v}_n\}_{n=1}^N$$

[Example] Differentiation: removing arguments

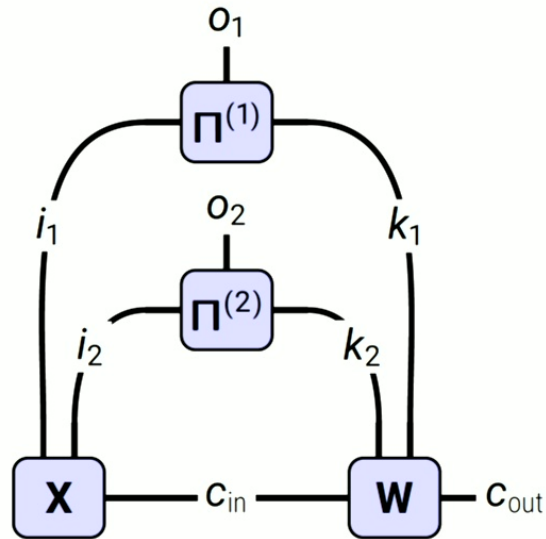
$$i - \left[\frac{\partial(\mathbf{ABv})}{\partial \mathbf{B}} \right] - \begin{matrix} j' \\ k' \end{matrix} = \frac{\partial \left(i - \mathbf{A} - j - \mathbf{B} - k - \mathbf{v} \right)}{\partial \left(j' - \mathbf{B} - k' \right)} = i - \mathbf{A} - j' \quad k' - \mathbf{v}$$

$$i - \left[\frac{\partial(\mathbf{ABv})}{\partial \mathbf{A}} \right] - \begin{matrix} i' \\ j' \end{matrix} = \frac{\partial \left(i - \mathbf{I} - l - \mathbf{A} - j - \mathbf{Bv} \right)}{\partial \left(i' - \mathbf{A} - j' \right)} = i - \mathbf{I} - i' \quad j' - \mathbf{Bv}$$

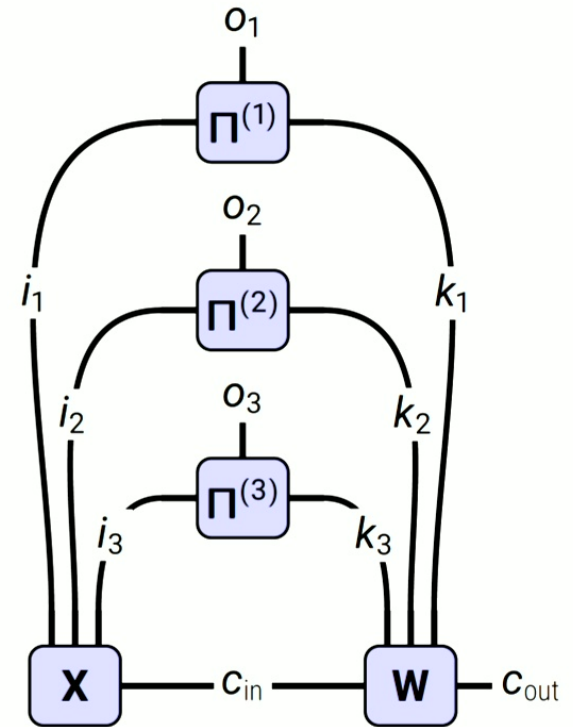
Higher-dimensional Convolutions [Hayashi et al., 2019]



1d



2d



3d

Non-conventional operations can be much cheaper with tensor networks

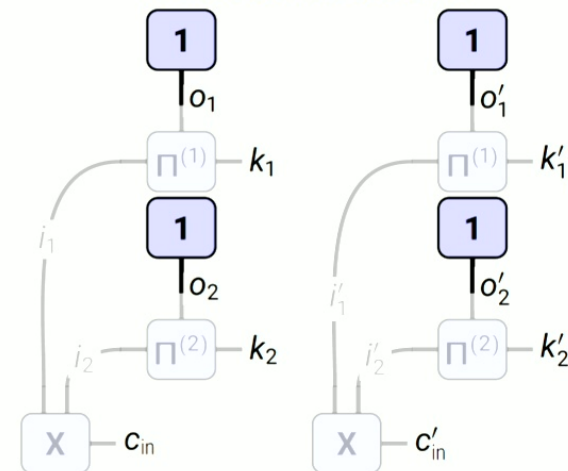
State of the art

$$\begin{aligned}
 &\mathbf{X} \\
 &\rightarrow \llbracket \mathbf{X} \rrbracket \\
 &\rightarrow \mathbf{1}^\top \llbracket \mathbf{X} \rrbracket \\
 &\rightarrow (\mathbf{1}^\top \llbracket \mathbf{X} \rrbracket)^\top (\mathbf{1}^\top \llbracket \mathbf{X} \rrbracket)
 \end{aligned}$$

Time: **9.87 ms**

Extra memory: **3.07 GiB**

Tensor Network

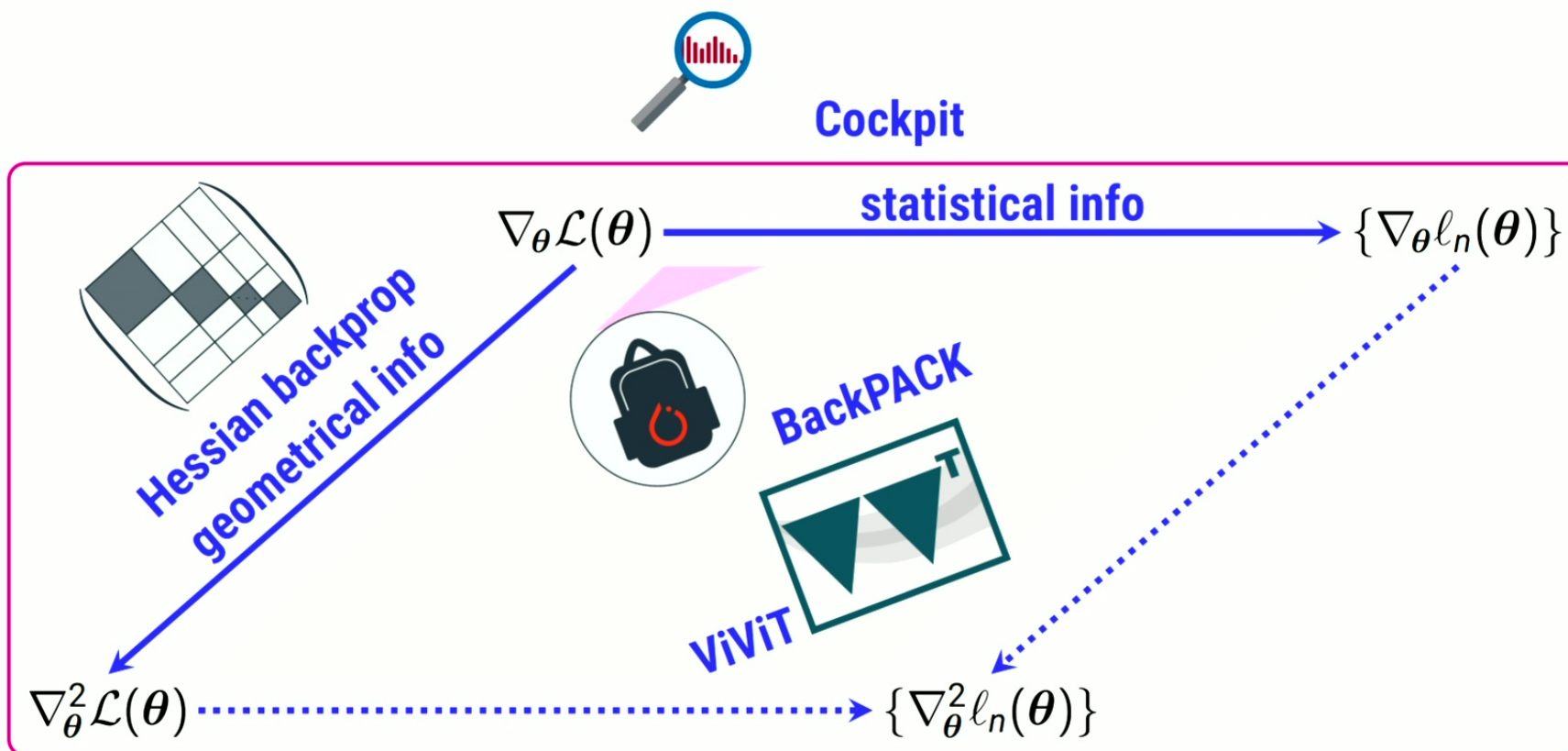


Time: **2.69 ms (3.7 x)**

Extra memory: **0 MiB**

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

Goal: Make Richer Information More Accessible

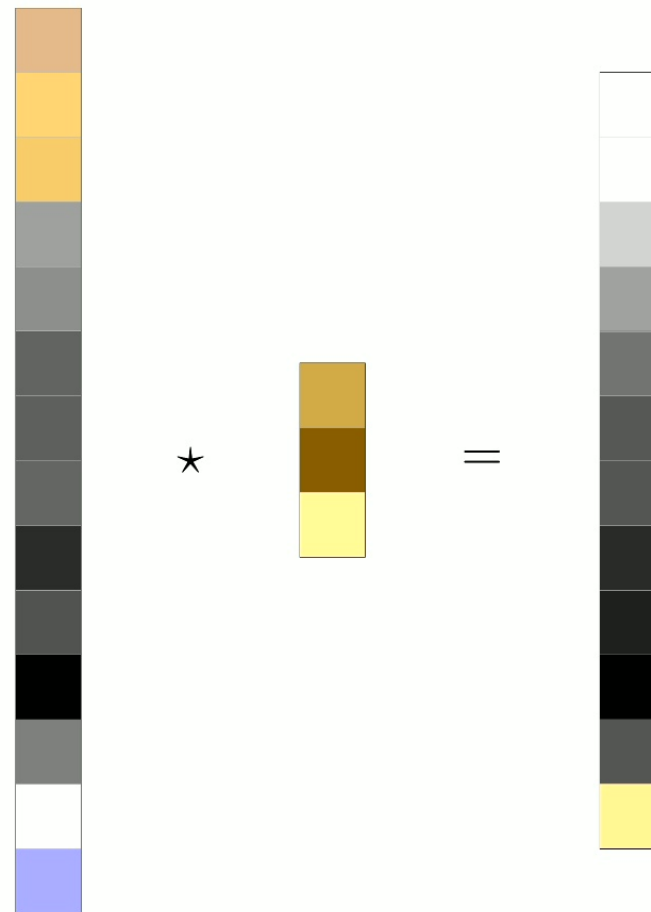


P is for Padding



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

- ★ $\mathbf{X} \in \mathbb{R}^l$
- ★ $\mathbf{W} \in \mathbb{R}^k$
- ★ $\mathbf{Y} \in \mathbb{R}^o$



Tensor Multiplication Unifies Vector & Matrix Multiplications



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

Consider two vectors \mathbf{a}, \mathbf{b} and two matrices \mathbf{A}, \mathbf{B}

Product	Notation	Index notation	In	Out
Inner	$\mathbf{c} = \mathbf{a}^\top \mathbf{b}$	$c = \sum_i a_i b_i$	$(i), (i)$	$()$
Element-wise	$\mathbf{c} = \mathbf{a} \odot \mathbf{b}$	$c_i = a_i b_i$	$(i), (i)$	(i)
Outer/Kronecker	$\mathbf{C} = \mathbf{a} \mathbf{b}^\top = \mathbf{a} \otimes \mathbf{b}^\top$	$C_{ij} = a_i b_j$	$(i), (j)$	(i, j)
Inner	$\mathbf{C} = \mathbf{A} \mathbf{B}$	$C_{i,k} = \sum_j A_{i,j} B_{j,k}$	$(i, j), (j, k)$	(i, k)
Element-wise	$\mathbf{C} = \mathbf{A} \odot \mathbf{B}$	$C_{ij} = A_{ij} B_{ij}$	$(i, j), (i, j)$	(i, j)
Kronecker	$\mathbf{C} = \mathbf{A} \otimes \mathbf{B}$	$C_{(i,k),(j,l)} = A_{ij} B_{k,l}$	$(i, j), (k, l)$	$((i, k), (j, l))$

We can infer the summations given the index signature

einsum is Probably All You Need



Recommendation: Use einsum in your code!

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

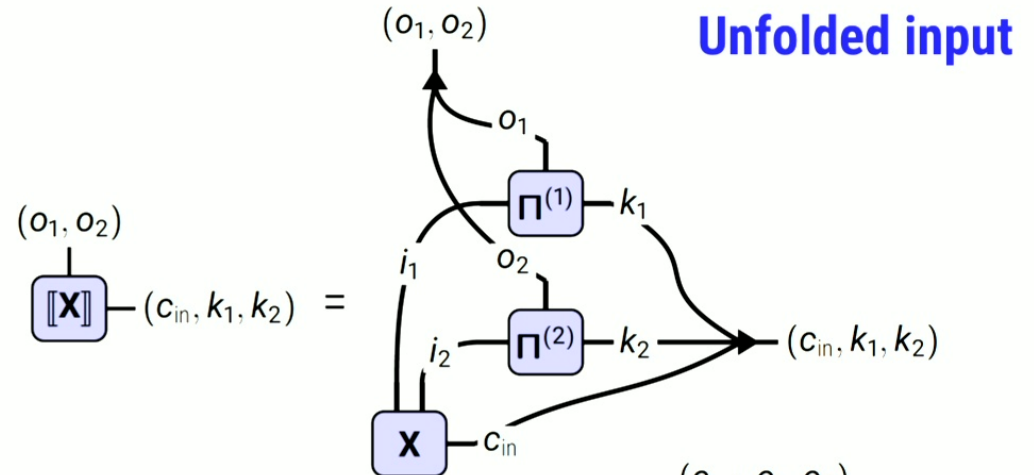
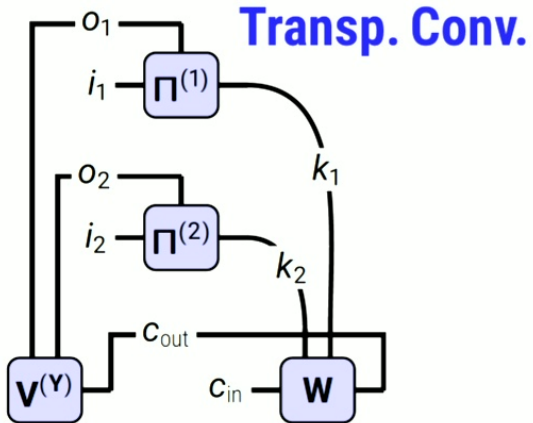
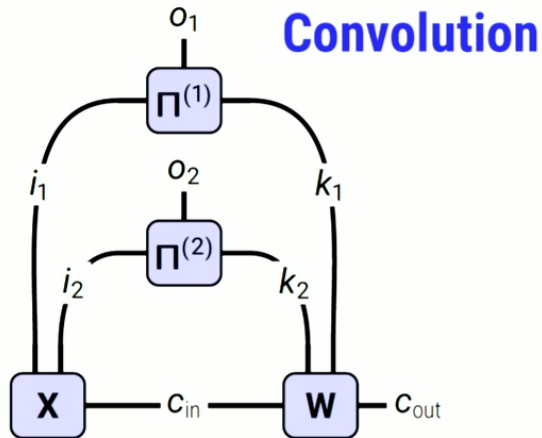
Also (not discussed)

- 😊 Automatic distribution (e.g. `cotengra` [Gray and Kourtis, 2021])
- 😊 Randomized/Approximate evaluation (somebody should do this!)

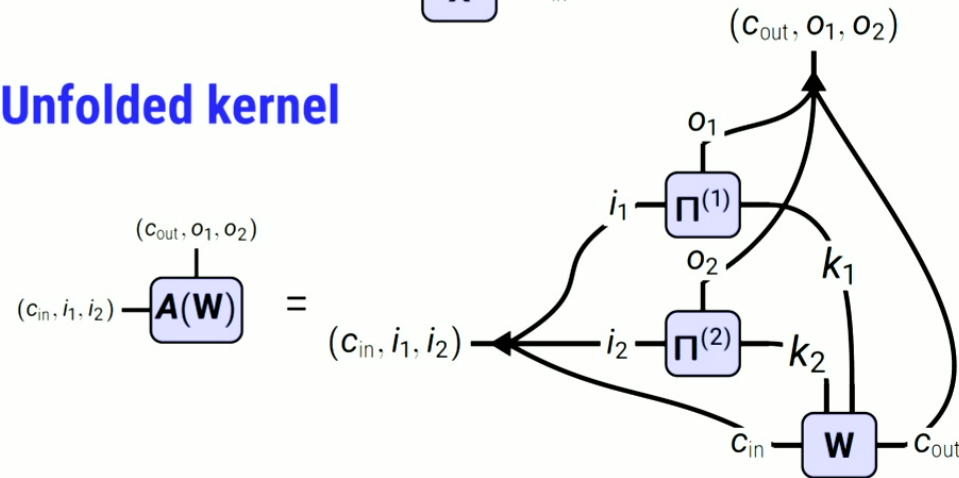
Also (regarding tensor networks)

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

Connection to Matrix-multiplication Perspective



Unfolded kernel



Convolutions and More as einsum

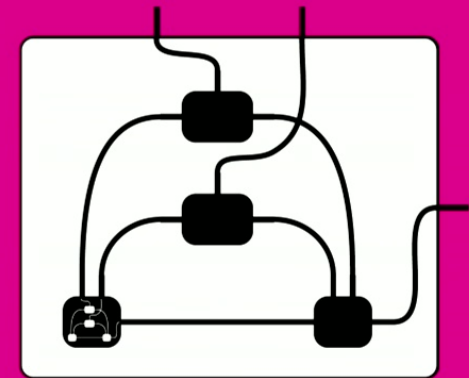


- ✦ TN perspective simplifies the transfer of algorithmic ideas
- ✦ Enables flexible/faster implementations of black box routines
- ✦ Relies on automatically efficient evaluation inside einsum

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)
```

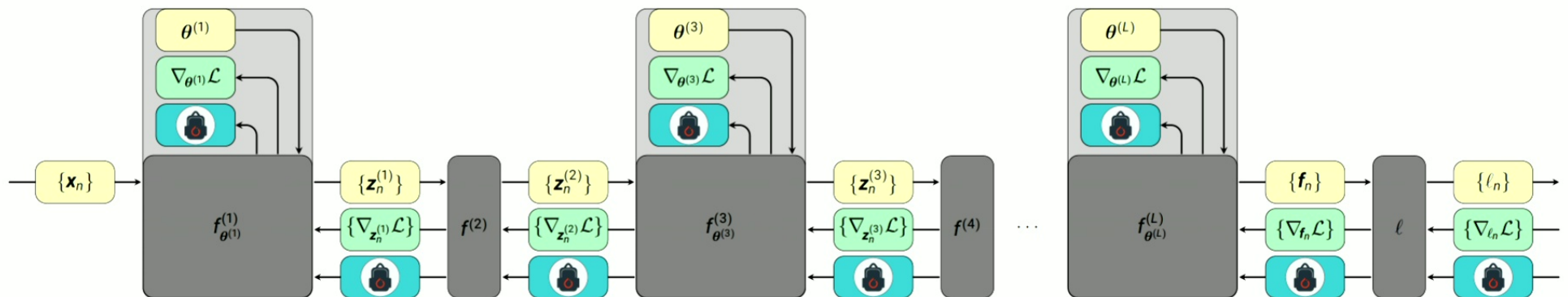


pip install einconv

Customized Backpropagation



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

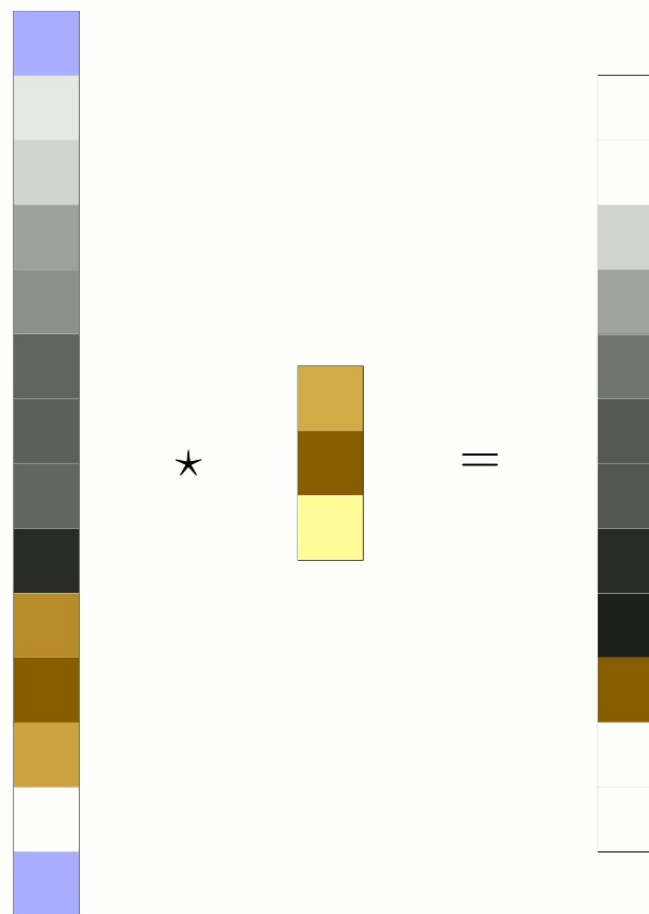


P is for Padding



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

- ★ $\mathbf{X} \in \mathbb{R}^l$
- ★ $\mathbf{W} \in \mathbb{R}^k$
- ★ $\mathbf{Y} \in \mathbb{R}^o$



A Mathematical Definition of Tensor Multiplication [e.g. Laue et al., 2020]



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

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

indices not present in the output are summed out; $S_{\mathbf{C}}$ satisfies $S_{\mathbf{C}} \subseteq S_{\mathbf{A}} \cup S_{\mathbf{B}}$.

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

$$\begin{aligned} S_{\mathbf{A}} &= (i, j) \\ S_{\mathbf{B}} &= (j, k) \\ S_{\mathbf{C}} &= (i, k) \end{aligned} \quad \Rightarrow \quad (S_{\mathbf{A}} \cup S_{\mathbf{B}}) \setminus S_{\mathbf{C}} = (i, j, k) \setminus (i, k) = (j)$$

$$\Rightarrow \mathbf{C} = *_{((i,j),(j,k),(i,k))}(\mathbf{A}, \mathbf{B}), \quad \text{or} \quad \mathbf{C} = \text{einsum}(\text{"ij,jk->ik"}, \mathbf{A}, \mathbf{B}).$$

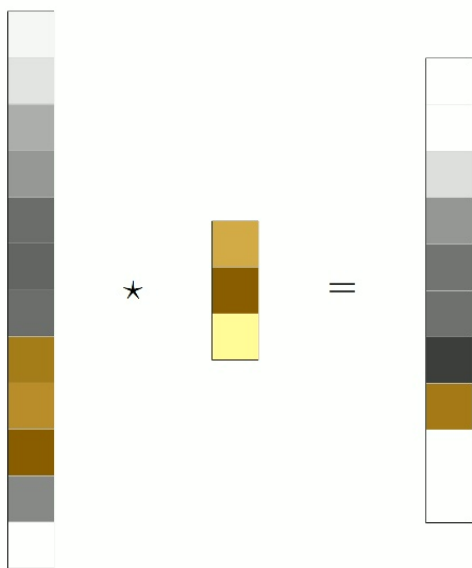
Many libraries provide $*_{(\dots)}(\dots)$ as `einsum`.

Convolution as Matrix Multiplication



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

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



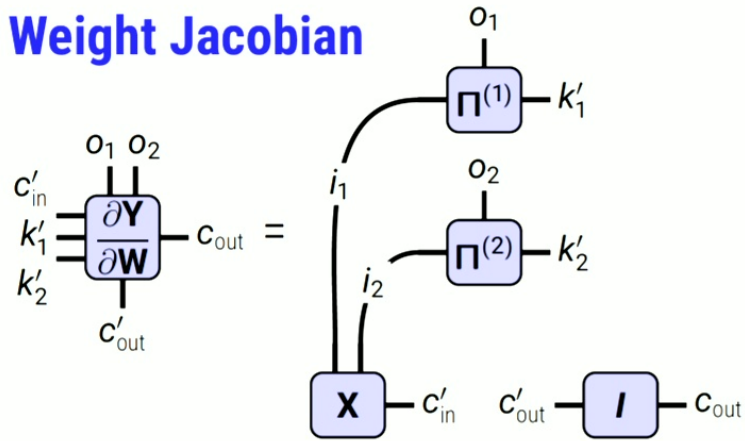
$$\mathbf{Y} = \mathbf{W} [\mathbf{X}]$$



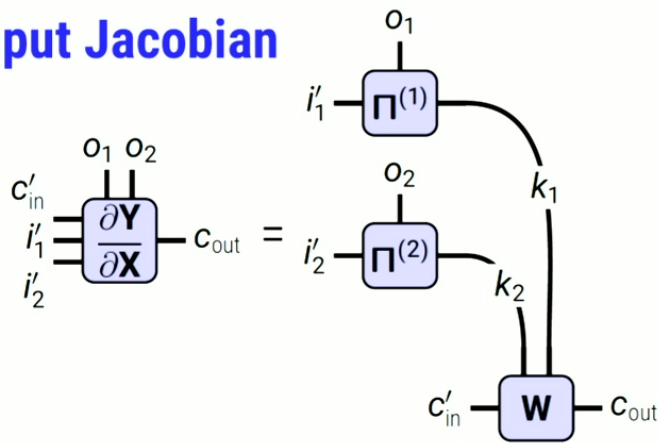
Derivatives & Autodiff Routines



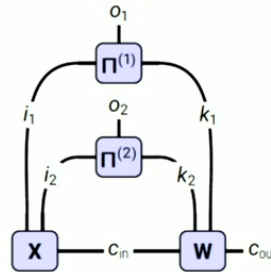
Weight Jacobian



Input Jacobian



Conv.



Convolutions and More as einsum

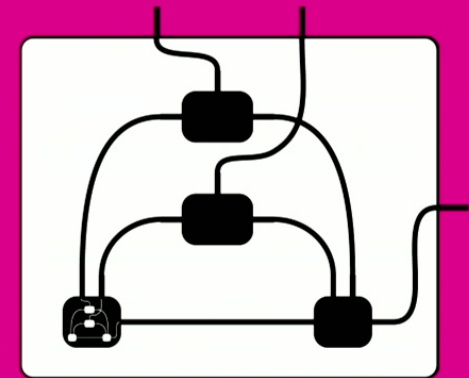


- ✦ TN perspective simplifies the transfer of algorithmic ideas
- ✦ Enables flexible/faster implementations of black box routines
- ✦ Relies on automatically efficient evaluation inside einsum

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](https://arxiv.org/abs/2307.02275)

Code: github.com/f-dangel/einconv

Thank you 🙏 questions?