

# Convolutions Through the Lens of Tensor Networks

Felix Dangel

December 01, 2023



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

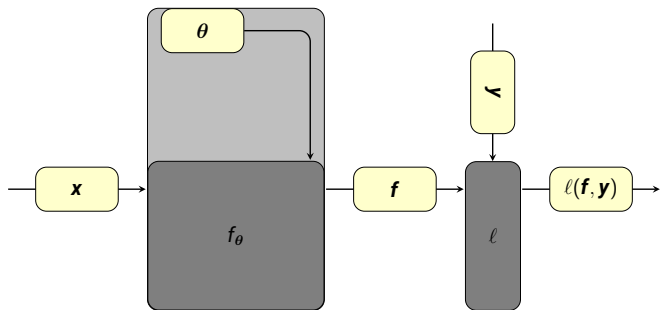
Felix Dangel

December 01, 2023



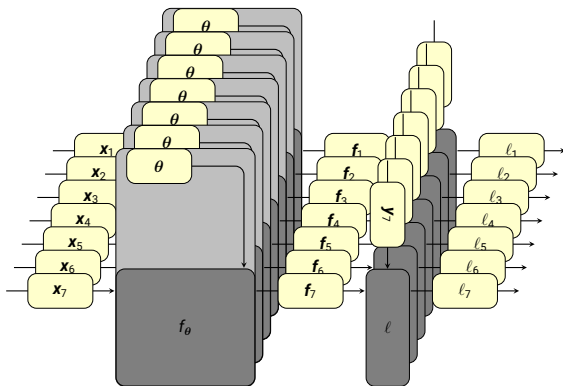
# **The Broader Picture**

# Supervised Deep Learning in a Nutshell (Math)



$$l(f_\theta(x), y)$$

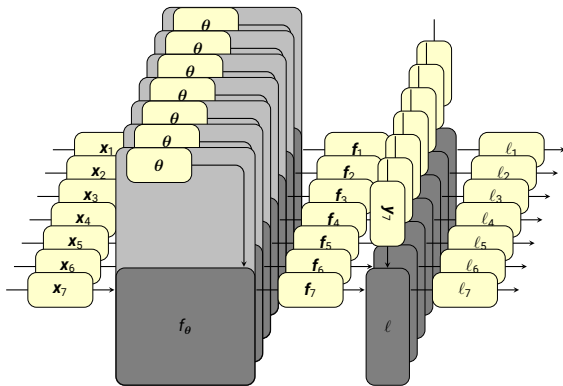
# Supervised Deep Learning in a Nutshell (Math)



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

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

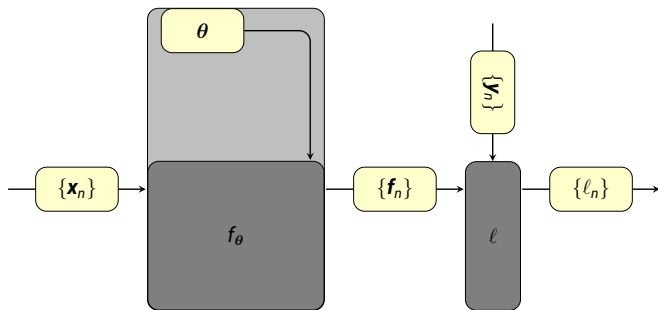
# Supervised Deep Learning in a Nutshell (Math)



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

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

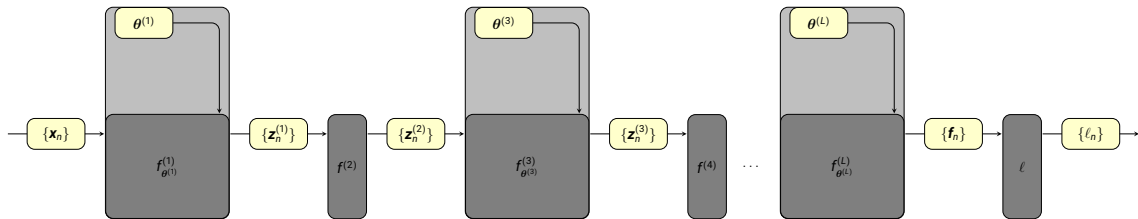
# Supervised Deep Learning in a Nutshell (Math)



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

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

# Supervised Deep Learning in a Nutshell (Math)



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

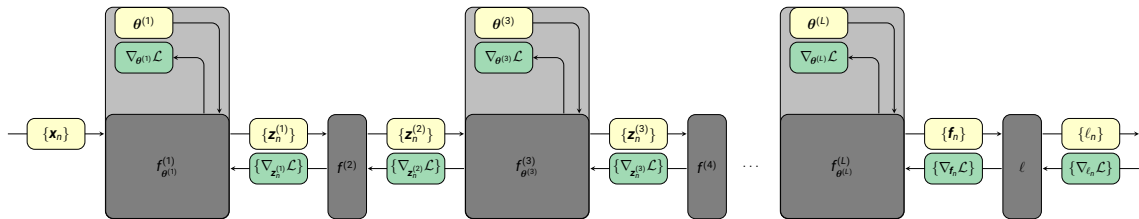
$$\mathcal{L}_{\mathbb{B}}(\theta) = \frac{1}{|\mathbb{B}|} \sum_{(\mathbf{x}_n, \mathbf{y}_n) \in \mathbb{B}} \ell(\mathbf{f}_{\theta}(\mathbf{x}_n), \mathbf{y}_n)$$

where

$$\mathbb{B} \sim \mathbb{D},$$
$$\mathbf{f}_{\theta} = \mathbf{f}_{\theta^{(L)}} \circ \mathbf{f}_{\theta^{(L-1)}} \circ \dots \circ \mathbf{f}_{\theta^{(1)}}.$$



# Supervised Deep Learning in a Nutshell (Math)



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

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

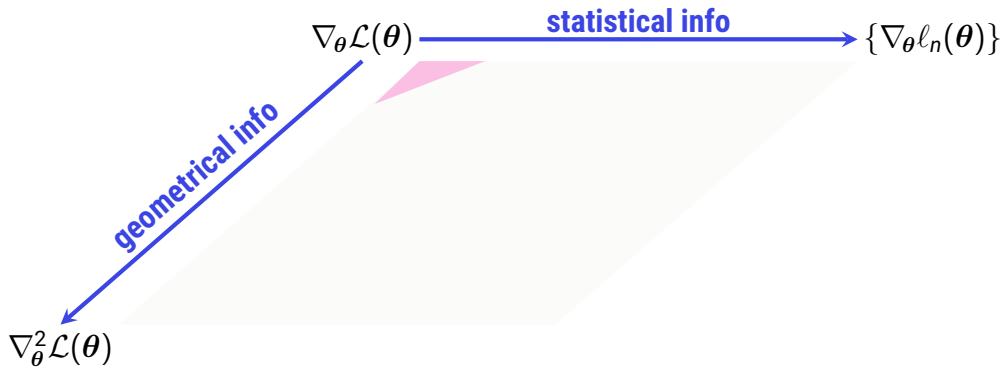
where

$$\mathbb{B} \sim \mathbb{D},$$
$$\mathbf{f}_{\boldsymbol{\theta}} = \mathbf{f}_{\theta^{(L)}}^{(L)} \circ \mathbf{f}_{\theta^{(L-1)}}^{(L-1)} \circ \dots \circ \mathbf{f}_{\theta^{(1)}}^{(1)}.$$

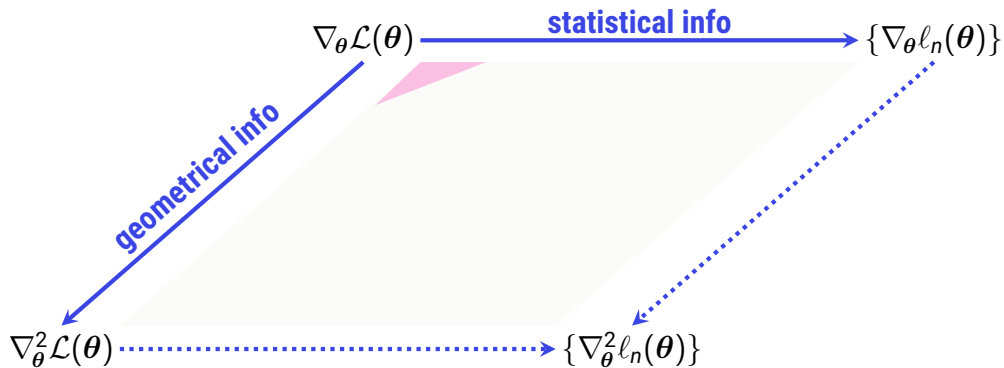
$$\nabla_{\theta} \mathcal{L}(\theta)$$

$\nabla_{\theta} \mathcal{L}(\theta)$  **statistical info**  $\rightarrow \{\nabla_{\theta} l_n(\theta)\}$

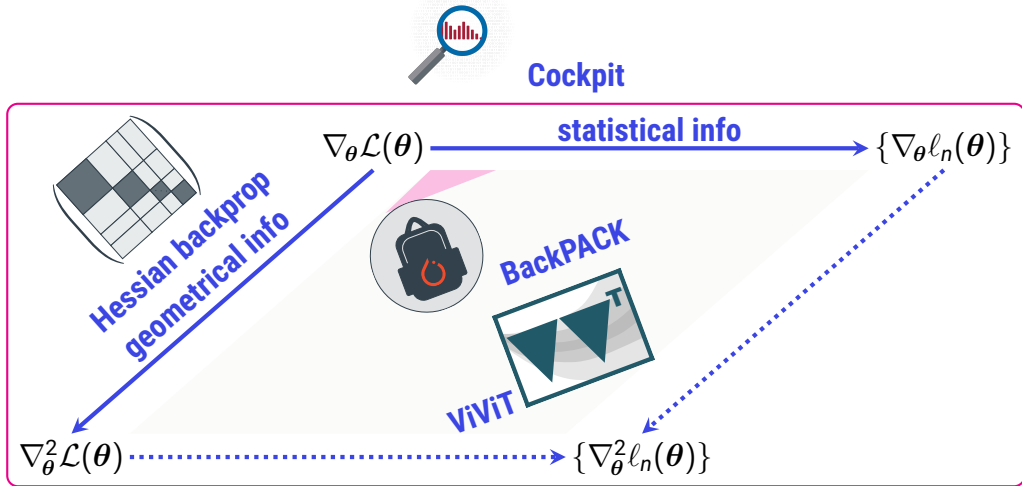
# Goal: Make Richer Information More Accessible



# Goal: Make Richer Information More Accessible



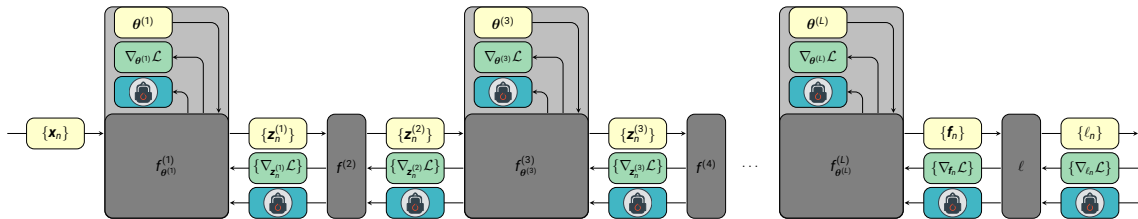
# Goal: Make Richer Information More Accessible



# 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 ``[N, *]`` 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 ``[N, *]`` 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)
```



# Convolutions Are More Tedious than Linear Layers



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

<b>Concept</b>	<b>for MLPs</b>	<b>for CNNs</b>
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

# Convolutions Are More Tedious than Linear Layers



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

<b>Concept</b>	<b>for MLPs</b>	<b>for CNNs</b>
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**

# Convolutions – A Visual Walkthrough

# A Simple Convolution

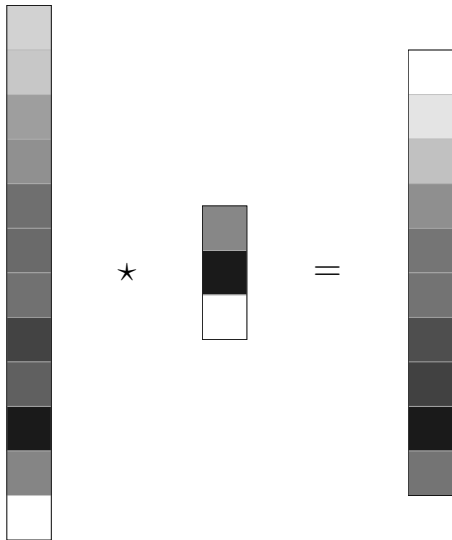


$\mathbf{Y} = \mathbf{X} \star \mathbf{W}$  with

★  $\mathbf{X} \in \mathbb{R}^l$

★  $\mathbf{W} \in \mathbb{R}^k$

★  $\mathbf{Y} \in \mathbb{R}^o$



# A Simple Convolution



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

★  $\mathbf{X} \in \mathbb{R}^I$

★  $\mathbf{W} \in \mathbb{R}^K$

★  $\mathbf{Y} \in \mathbb{R}^O$

★

=

# S is for Stride

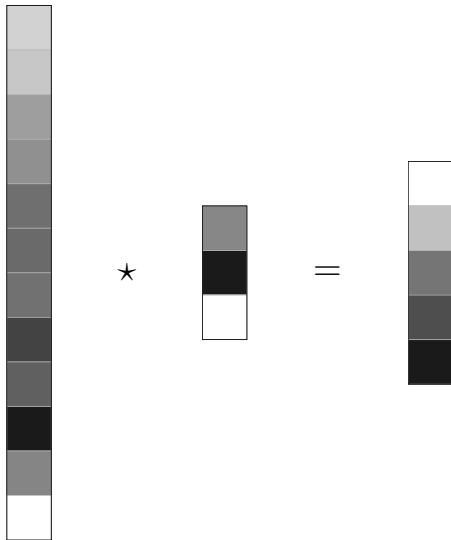


$\mathbf{Y} = \mathbf{X} \star \mathbf{W}$  with

★  $\mathbf{X} \in \mathbb{R}^l$

★  $\mathbf{W} \in \mathbb{R}^k$

★  $\mathbf{Y} \in \mathbb{R}^o$



# S is for Stride



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

★  $\mathbf{X} \in \mathbb{R}^I$

★  $\mathbf{W} \in \mathbb{R}^K$

★  $\mathbf{Y} \in \mathbb{R}^O$

★

=

# $P$ is for Padding

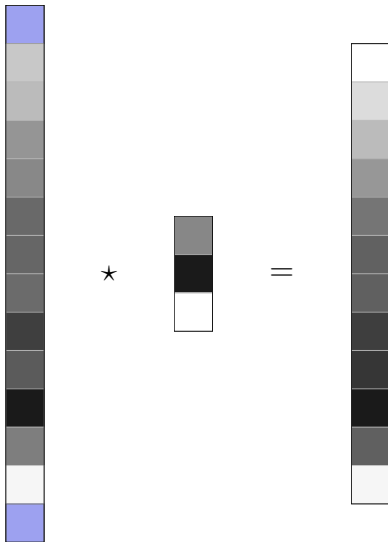


$\mathbf{Y} = \mathbf{X} \star \mathbf{W}$  with

★  $\mathbf{X} \in \mathbb{R}^l$

★  $\mathbf{W} \in \mathbb{R}^k$

★  $\mathbf{Y} \in \mathbb{R}^o$





# $P$ is for Padding



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

★  $\mathbf{X} \in \mathbb{R}^l$

★  $\mathbf{W} \in \mathbb{R}^k$

★  $\mathbf{Y} \in \mathbb{R}^o$

★

=

# D is for Dilation

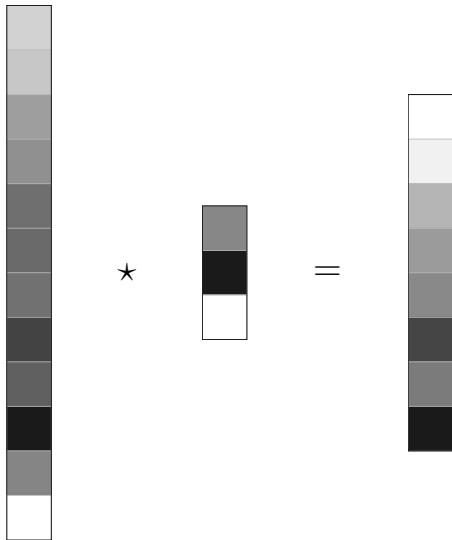


$\mathbf{Y} = \mathbf{X} \star \mathbf{W}$  with

★  $\mathbf{X} \in \mathbb{R}^l$

★  $\mathbf{W} \in \mathbb{R}^k$

★  $\mathbf{Y} \in \mathbb{R}^o$



# $D$ is for Dilation



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

★  $\mathbf{X} \in \mathbb{R}^l$

★  $\mathbf{W} \in \mathbb{R}^k$

★  $\mathbf{Y} \in \mathbb{R}^o$

★

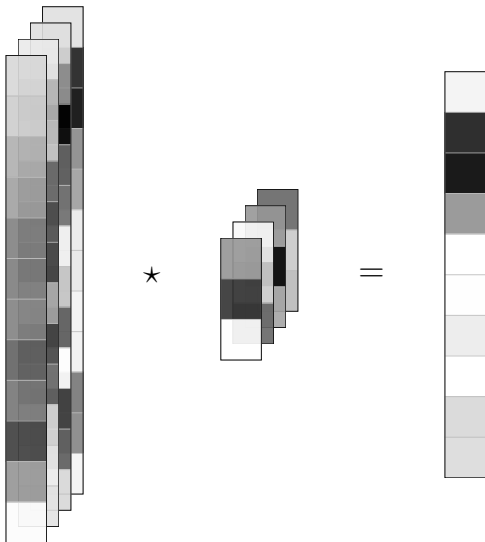
=

# $C_{in}$ 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



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

★  $\mathbf{X} \in \mathbb{R}^{C_{in} \times I}$

★  $\mathbf{W} \in \mathbb{R}^{C_{in} \times K}$

★  $\mathbf{Y} \in \mathbb{R}^O$

★

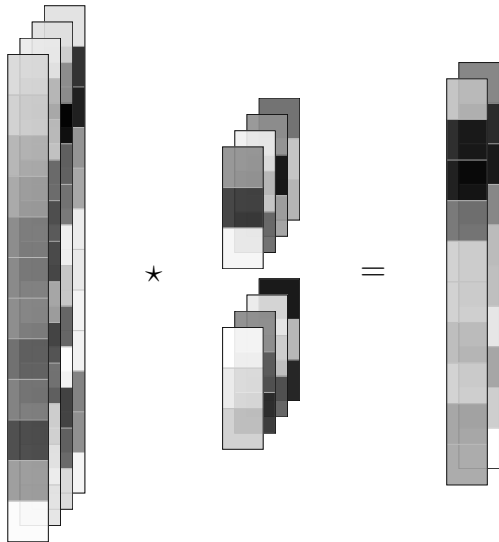
=

# $C_{out}$ is for Output Channels



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

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



# $C_{\text{out}}$ is for Output Channels



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

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

★

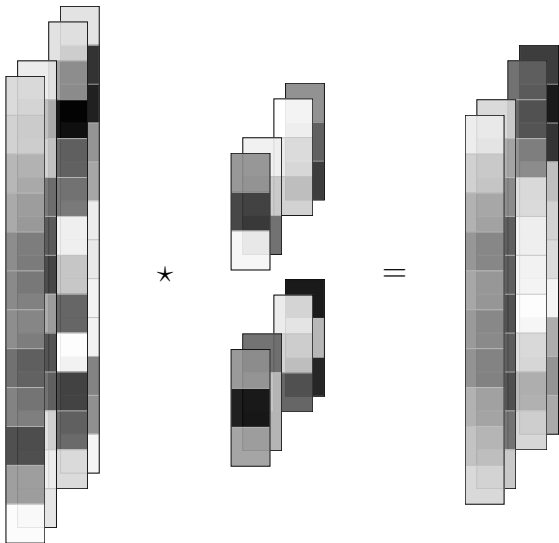
=

# G is for Groups



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

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





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

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

★

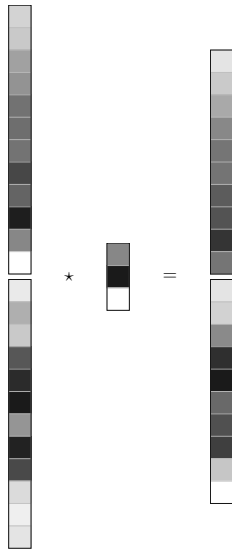
=

# $N$ is for Batch Size



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



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

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

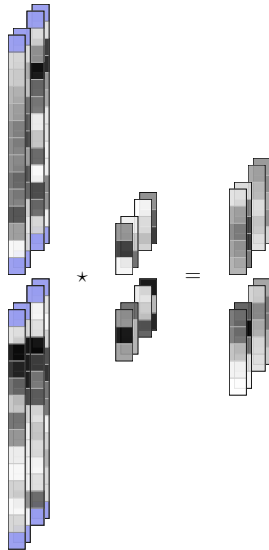
★ =

# Putting Everything Together in One Dimension...



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

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



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

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

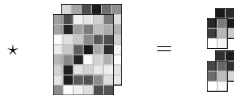
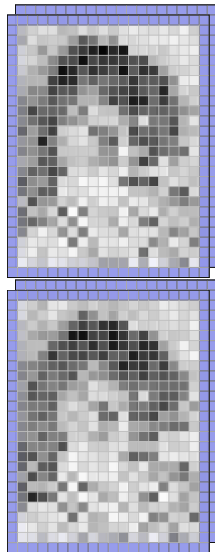
★ =

# ...and in Two Dimensions



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

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



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

★  $\mathbf{X} \in \mathbb{R}^{N \times C_{in} \times I_1 \times I_2}$

★  $\mathbf{W} \in \mathbb{R}^{C_{out} \times C_{in}/G \times K_1 \times K_2}$

★  $\mathbf{Y} \in \mathbb{R}^{N \times C_{out} \times O_1 \times O_2}$

★

=

# Tensor Multiplication & Tensor Networks



# Tensor Multiplication Unifies Vector & Matrix Multiplications



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

# Tensor Multiplication Unifies Vector & Matrix Multiplications



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

Consider two vectors  $\mathbf{a}, \mathbf{b}$

---

<b>Product</b>	<b>Notation</b>
Inner	$\mathbf{c} = \mathbf{a}^\top \mathbf{b}$
Element-wise	$\mathbf{c} = \mathbf{a} \odot \mathbf{b}$
Outer/Kronecker	$\mathbf{C} = \mathbf{a}\mathbf{b}^\top = \mathbf{a} \otimes \mathbf{b}^\top$

---

# Tensor Multiplication Unifies Vector & Matrix Multiplications



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

Consider two vectors  $\mathbf{a}, \mathbf{b}$

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	$\mathbf{C} = \mathbf{a} \mathbf{b}^\top = \mathbf{a} \otimes \mathbf{b}^\top$	$C_{i,j} = a_i b_j$

# 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	$\mathbf{C} = \mathbf{A} \mathbf{B}$	
Element-wise	$\mathbf{C} = \mathbf{A} \odot \mathbf{B}$	
Kronecker	$\mathbf{C} = \mathbf{A} \otimes \mathbf{B}$	

# 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	$\mathbf{C} = \mathbf{A} \otimes \mathbf{B}$	$C_{(i,k),(j,l)} = A_{i,j} B_{k,l}$

# 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	$\mathbf{C} = \mathbf{a} \mathbf{b}^\top = \mathbf{a} \otimes \mathbf{b}^\top$	$C_{i,j} = a_i b_j$	$(i), (j)$	$(i, j)$
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	$\mathbf{C} = \mathbf{A} \odot \mathbf{B}$	$C_{i,j} = A_{i,j} B_{i,j}$	$(i, j), (i, j)$	$(i, j)$
Kronecker	$\mathbf{C} = \mathbf{A} \otimes \mathbf{B}$	$C_{(i,k),(j,l)} = A_{i,j} B_{k,l}$	$(i, j), (k, l)$	$((i, k), (j, l))$

**We can infer the summations given the index signature**

Given two tensors  $\mathbf{A}$ ,  $\mathbf{B}$  with index tuples  $\mathbf{S}_A$ ,  $\mathbf{S}_B$

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

indices not present in the output are summed out;  $\mathbf{S}_C$  satisfies  $\mathbf{S}_C \subseteq \mathbf{S}_A \cup \mathbf{S}_B$ .



Given two tensors  $\mathbf{A}$ ,  $\mathbf{B}$  with index tuples  $\mathbf{S}_A$ ,  $\mathbf{S}_B$

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

indices not present in the output are summed out;  $\mathbf{S}_C$  satisfies  $\mathbf{S}_C \subseteq \mathbf{S}_A \cup \mathbf{S}_B$ .

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

$$\begin{aligned} \mathbf{S}_A &= (i, j) \\ \mathbf{S}_B &= (j, k) \\ \mathbf{S}_C &= (i, k) \end{aligned}$$





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

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

Given two tensors  $\mathbf{A}$ ,  $\mathbf{B}$  with index tuples  $\mathbf{S}_A$ ,  $\mathbf{S}_B$

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

indices not present in the output are summed out;  $\mathbf{S}_C$  satisfies  $\mathbf{S}_C \subseteq \mathbf{S}_A \cup \mathbf{S}_B$ .

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

$$\begin{aligned} \mathbf{S}_A &= (i, j) \\ \mathbf{S}_B &= (j, k) \\ \mathbf{S}_C &= (i, k) \end{aligned} \quad \Rightarrow \quad (\mathbf{S}_A \cup \mathbf{S}_B) \setminus \mathbf{S}_C = (i, j, k) \setminus (i, k) = (j)$$

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

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



**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  $\mathbf{C} = \mathbf{AB}$  as tensor multiplication


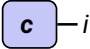
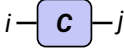

$$\begin{aligned} \mathbf{S}_{\mathbf{A}} &= (i, j) \\ \mathbf{S}_{\mathbf{B}} &= (j, k) \\ \mathbf{S}_{\mathbf{C}} &= (i, k) \end{aligned} \quad \Rightarrow \quad (\mathbf{S}_{\mathbf{A}} \cup \mathbf{S}_{\mathbf{B}}) \setminus \mathbf{S}_{\mathbf{C}} = (i, j, k) \setminus (i, k) = (j)$$

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

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

# Tensor Networks: Graphical Notation for Tensor Multiplications



Operation	Diagram
Scalar	
Vector	
Matrix	
Tensor	

# Tensor Networks: Graphical Notation for Tensor Multiplications



Operation	Diagram
Scalar	
Vector	
Matrix	
Tensor	
Matricize	
Flatten	

# Tensor Networks: Graphical Notation for Tensor Multiplications



Operation	Diagram
Scalar	
Vector	
Matrix	
Tensor	
Matricize	
Flatten	
Trace	

# Tensor Networks: Graphical Notation for Tensor Multiplications



Operation	Diagram
Scalar	
Vector	
Matrix	
Tensor	
Matricize	
Flatten	
Trace	

Operation	Diagram

# Tensor Networks: Graphical Notation for Tensor Multiplications



Operation	Diagram
Scalar	
Vector	
Matrix	
Tensor	
Matricize	
Flatten	
Trace	

Operation	Diagram
$i \rightarrow AB \rightarrow k$	
$i \rightarrow A \odot B \rightarrow j$	



# Tensor Networks: Graphical Notation for Tensor Multiplications



Operation	Diagram
Scalar	
Vector	
Matrix	
Tensor	
Matricize	
Flatten	
Trace	

Operation	Diagram

# Tensor Networks: Graphical Notation for Tensor Multiplications



Operation	Diagram
Scalar	
Vector	
Matrix	
Tensor	
Matricize	
Flatten	
Trace	

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

# Benefits of TN/einsum Abstraction



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



## 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  $\{\mathbf{C}_n\}_{n=1}^N = \{\mathbf{A}_n \mathbf{B}_n\}_{n=1}^N$



## 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")
```



## 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)
```



## 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")
```

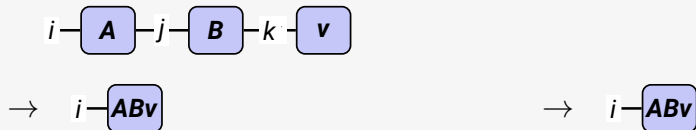
[Example] Improved readability (from an ICLR 2024 submission)

```
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")
```



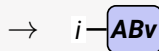
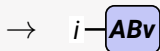
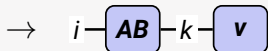
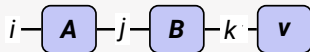
## [Example] Matrix-matrix-vector-product

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  $\mathbf{v} \in \mathbb{R}^{3.000}$ .

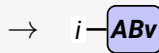
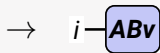
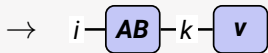
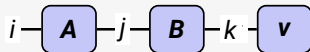


Schedule 1:

A @ B @ v

## [Example] Matrix-matrix-vector-product

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



Schedule 1:

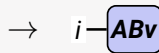
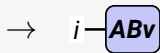
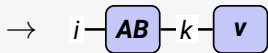
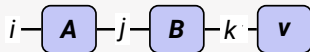
A @ B @ v

Schedule 2:

A @ (B @ v)

## [Example] Matrix-matrix-vector-product

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



Schedule 1: **1.1 s**

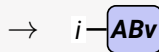
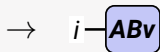
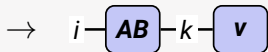
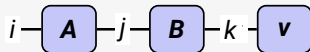
A @ B @ v

Schedule 2: **2.2 ms**

A @ (B @ v)

## [Example] Matrix-matrix-vector-product

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



Schedule 1: **1.1 s**

A @ B @ v

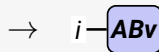
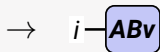
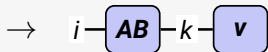
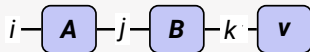
Schedule 2: **2.2 ms**

A @ (B @ 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  $\mathbf{v} \in \mathbb{R}^{3.000}$ .



Schedule 1: **1.1 s**

`A @ B @ v`

Schedule 2: **2.2 ms**

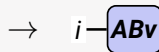
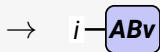
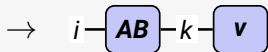
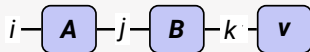
`A @ (B @ v)`

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  $\mathbf{v} \in \mathbb{R}^{3.000}$ .



Schedule 1: **1.1 s**

`A @ B @ v`

Schedule 2: **2.2 ms**

`A @ (B @ v)`

PyTorch: **1.1 s**

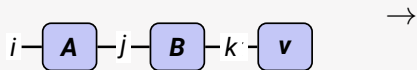
with update + pip install `opt_einsum`, **2.3 ms**

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

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

[Example] Batching/vmap-ing: adding legs

$$(\mathbf{A}, \mathbf{B}, \mathbf{v}) \mapsto \mathbf{ABv}$$

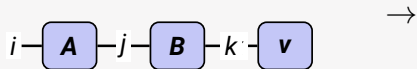




[Example] Batching/vmap-ing: adding legs

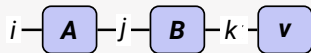
$$(\mathbf{A}, \mathbf{B}, \mathbf{v}) \mapsto \mathbf{ABv}$$

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



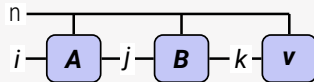
[Example] Batching/vmap-ing: adding legs

$$(\mathbf{A}, \mathbf{B}, \mathbf{v}) \mapsto \mathbf{ABv}$$



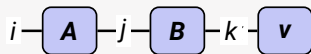
→

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



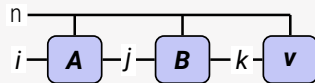
[Example] Batching/vmap-ing: adding legs

$$(\mathbf{A}, \mathbf{B}, \mathbf{v}) \mapsto \mathbf{ABv}$$



→

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

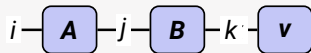


[Example] Differentiation: removing arguments

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

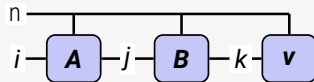
[Example] Batching/vmap-ing: adding legs

$$(\mathbf{A}, \mathbf{B}, \mathbf{v}) \mapsto \mathbf{ABv}$$



→

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

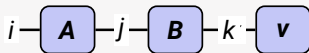


[Example] Differentiation: removing arguments

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

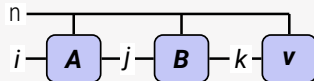
[Example] Batching/vmap-ing: adding legs

$$(\mathbf{A}, \mathbf{B}, \mathbf{v}) \mapsto \mathbf{ABv}$$



→

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



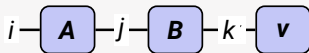
[Example] Differentiation: removing arguments

$$i \rightarrow \frac{\partial(\mathbf{ABv})}{\partial \mathbf{B}} \leftarrow j' \leftarrow k' = \frac{\partial \left( i \rightarrow \mathbf{A} \rightarrow j \rightarrow \mathbf{B} \rightarrow k \rightarrow \mathbf{v} \right)}{\partial \left( j' \rightarrow \mathbf{B} \rightarrow k' \right)} = i \rightarrow \mathbf{A} \rightarrow j' \quad k' \rightarrow \mathbf{v}$$

$$i \rightarrow \frac{\partial(\mathbf{ABv})}{\partial \mathbf{A}} \leftarrow j'$$

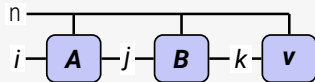
[Example] Batching/vmap-ing: adding legs

$$(\mathbf{A}, \mathbf{B}, \mathbf{v}) \mapsto \mathbf{ABv}$$



→

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



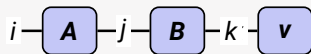
[Example] Differentiation: removing arguments

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

$$i - \frac{\partial(\mathbf{ABv})}{\partial \mathbf{A}} - j' = \frac{\partial \left( i - \mathbf{I} - l - \mathbf{A} - j - \mathbf{Bv} \right)}{\partial \left( i' - \mathbf{A} - j' \right)}$$

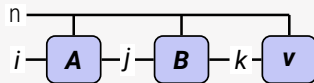
[Example] Batching/vmap-ing: adding legs

$$(\mathbf{A}, \mathbf{B}, \mathbf{v}) \mapsto \mathbf{ABv}$$



→

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



[Example] Differentiation: removing arguments

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

$$i - \frac{\partial(\mathbf{ABv})}{\partial \mathbf{A}} - j' = \frac{\partial \left( i - \mathbf{I} - l - \mathbf{A} - j - \mathbf{Bv} \right)}{\partial \left( i' - \mathbf{A} - j' \right)} = i - \mathbf{I} - i' - j' - \mathbf{Bv}$$

## Recommendation: Use einsum in your code!

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



## 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!)

## 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

# Convolutions as Tensor Networks

# Convolution as Matrix Multiplication



$$\mathbf{X} \in \mathbb{R}^{C_{in} \times I}$$

$$\mathbf{Y} \in \mathbb{R}^{C_{out} \times O}$$

$$\mathbf{W} \in \mathbb{R}^{C_{out} \times C_{in} \times K}$$

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

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

$\star$  =

=

# Convolution as Matrix Multiplication



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

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

$\star$  =

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

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

=

# Convolution as Structured Matrix Multiplication



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

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

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

$$\star \quad = \quad =$$

# Convolution as Structured Matrix Multiplication



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

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

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

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

$\star$                       =                      =



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





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

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

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

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

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

**X**

**Y**

**W**

★

$[\Pi]_{:, :, k}$

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

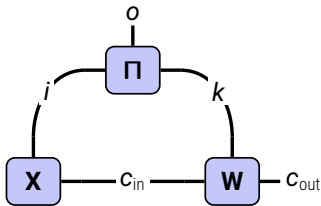
$$\begin{aligned}
 Y_{c_{out},o} &= \sum_{c_{in},k} X_{c_{in},i(k,o)} W_{c_{out},c_{in},k} \\
 &= \sum_{c_{in},k} \sum_i X_{c_{in},i} \Pi_{i,o,k} W_{c_{out},c_{in},k}
 \end{aligned}$$

**X**

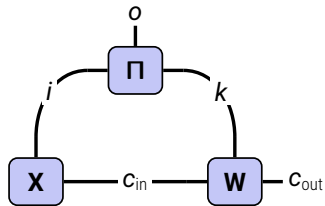
**Y**

**W**

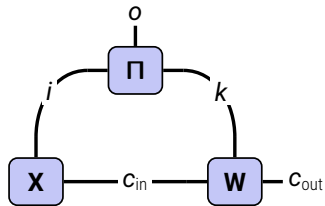
\*



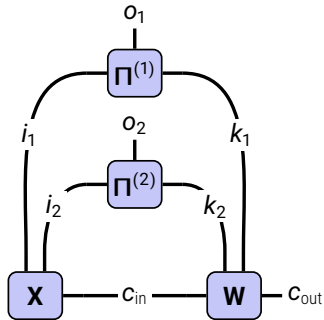
$[\Pi]_{:, :, k}$



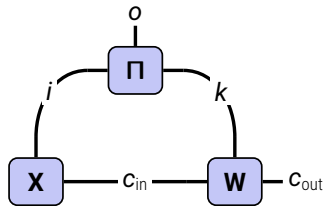
1d



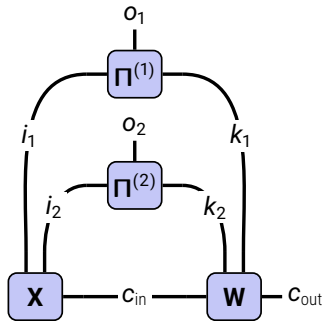
1d



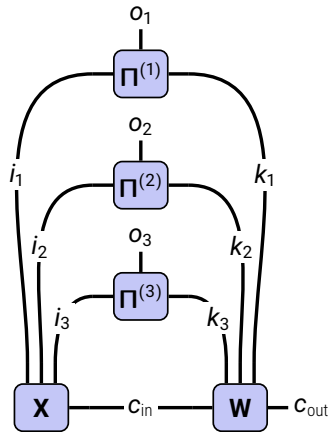
2d



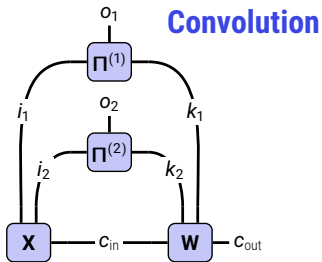
1d



2d

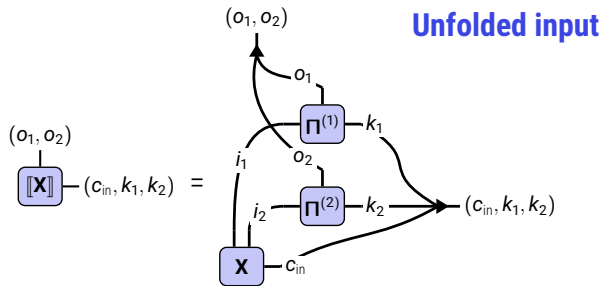
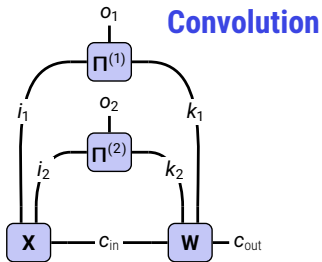


3d

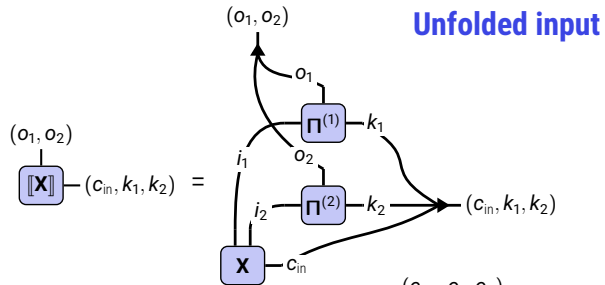
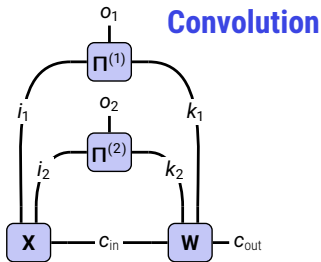




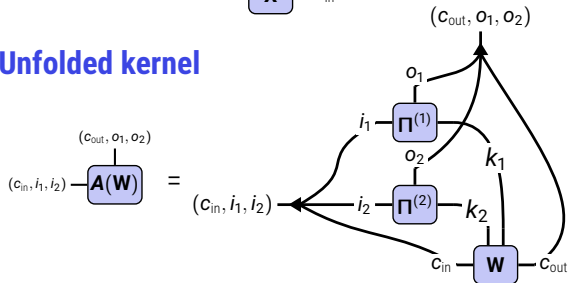
# Connection to Matrix-multiplication Perspective



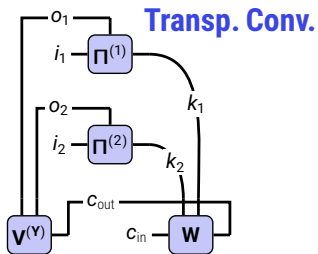
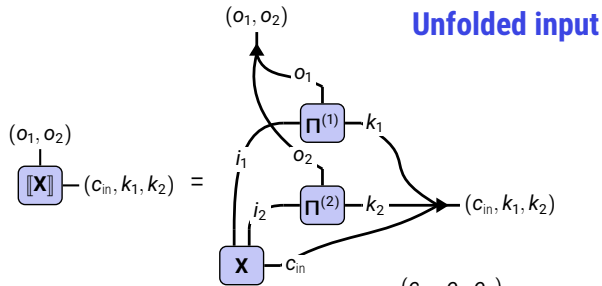
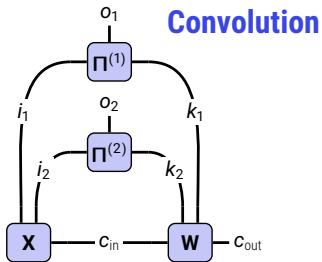
# Connection to Matrix-multiplication Perspective



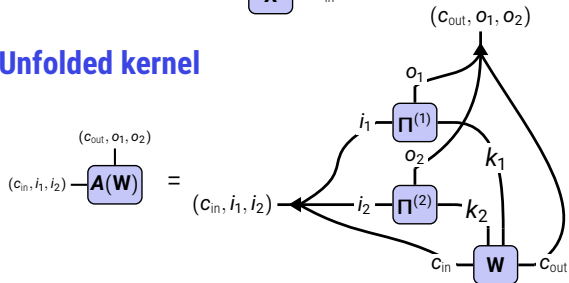
## Unfolded kernel



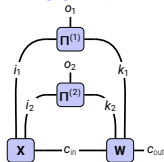
# Connection to Matrix-multiplication Perspective



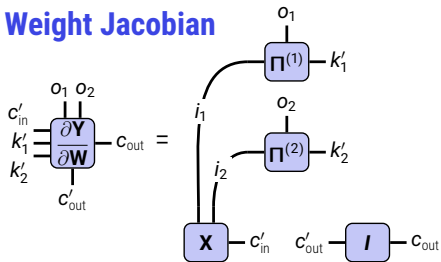
**Unfolded kernel**



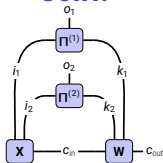
Conv.



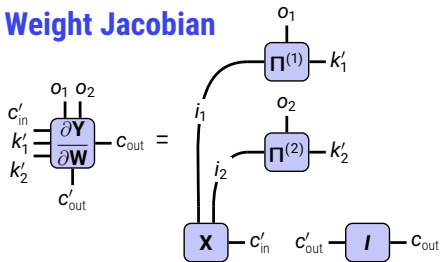
## Weight Jacobian



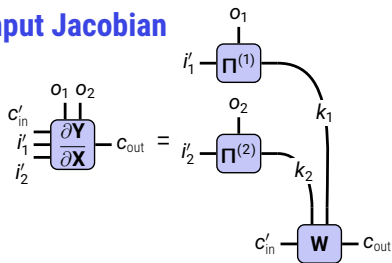
## Conv.



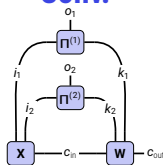
## Weight Jacobian



## Input Jacobian

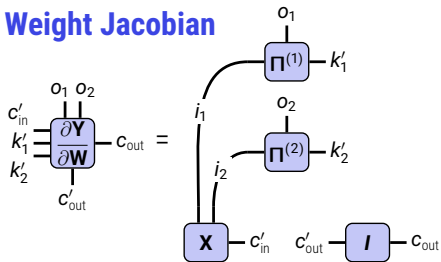


## Conv.

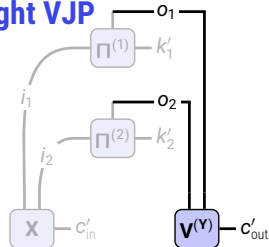




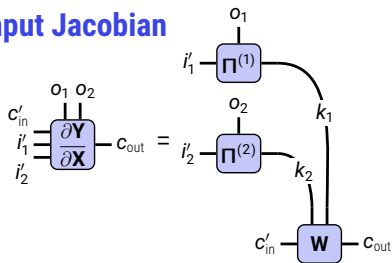
## Weight Jacobian



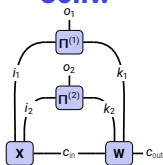
## Weight VJP



## Input Jacobian

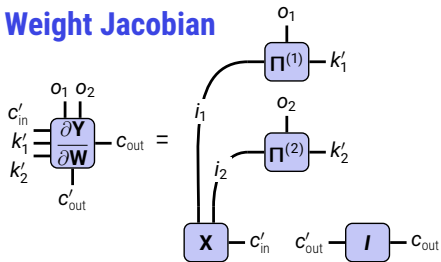


## Conv.

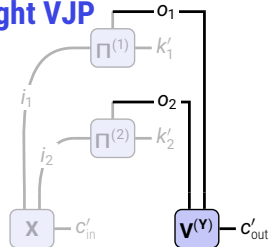




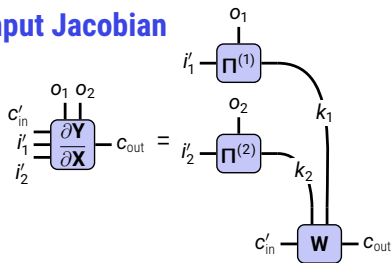
## Weight Jacobian



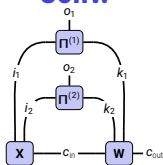
## Weight VJP



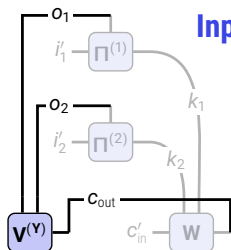
## Input Jacobian



## Conv.

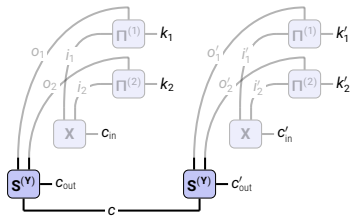


## Input VJP

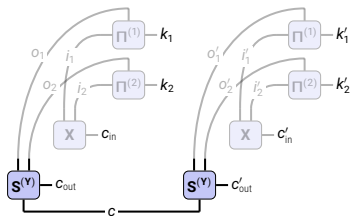




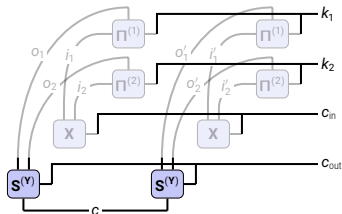
## Fisher/GGN block



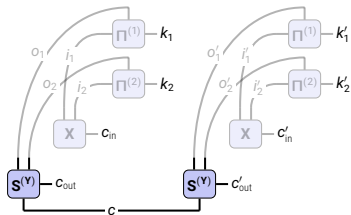
## Fisher/GGN block



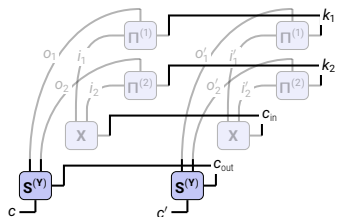
## Fisher/GGN diagonal



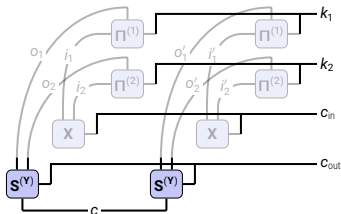
## Fisher/GGN block



## GGN Gram/Empirical NTK matrix

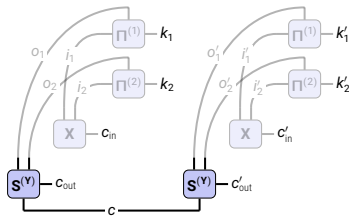


## Fisher/GGN diagonal

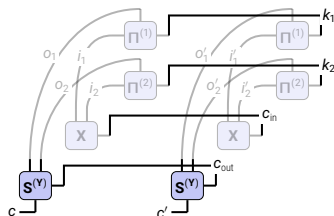




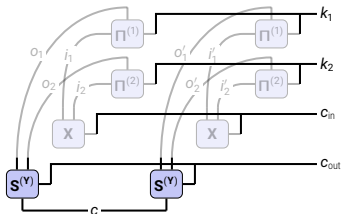
## Fisher/GGN block



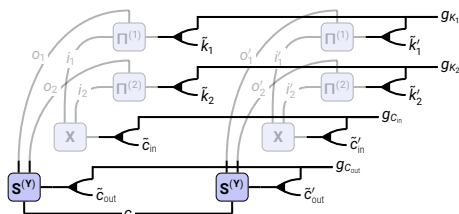
## GGN Gram/Empirical NTK matrix



## Fisher/GGN diagonal

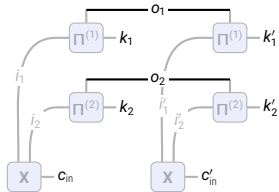


## Fisher/GGN mini-block diagonal

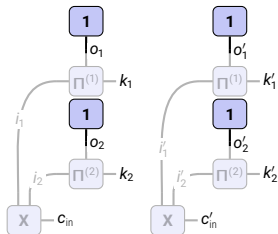




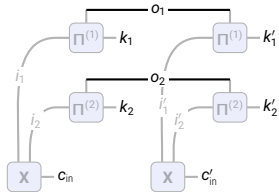
### KFC/KFAC-expand



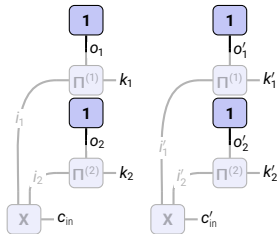
### KFAC-reduce



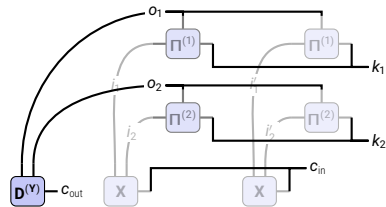
### KFC/KFAC-expand



### KFAC-reduce



### Approximate Hessian diagonal



**Check out the table in the paper's appendix!**

# Example & Conclusion



## State of the art

**x**



## State of the art

**x**

→ **[[x]]**

## State of the art

**X**

$$\rightarrow \llbracket \mathbf{X} \rrbracket$$

$$\rightarrow \mathbf{1}^\top \llbracket \mathbf{X} \rrbracket$$

## State of the art

**X**

$$\rightarrow \llbracket \mathbf{X} \rrbracket$$

$$\rightarrow \mathbf{1}^\top \llbracket \mathbf{X} \rrbracket$$

$$\rightarrow (\mathbf{1}^\top \llbracket \mathbf{X} \rrbracket)^\top (\mathbf{1}^\top \llbracket \mathbf{X} \rrbracket)$$

## State of the art

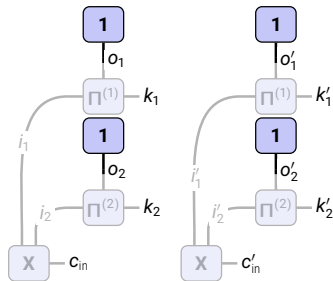
$\mathbf{X}$

$$\rightarrow \llbracket \mathbf{X} \rrbracket$$

$$\rightarrow \mathbf{1}^\top \llbracket \mathbf{X} \rrbracket$$

$$\rightarrow (\mathbf{1}^\top \llbracket \mathbf{X} \rrbracket)^\top (\mathbf{1}^\top \llbracket \mathbf{X} \rrbracket)$$

## Tensor Network



## State of the art

$\mathbf{X}$

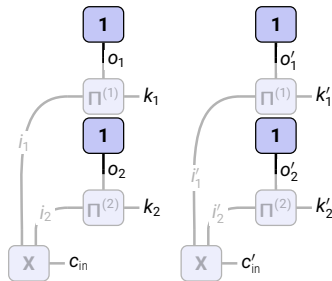
$$\rightarrow \llbracket \mathbf{X} \rrbracket$$

$$\rightarrow \mathbf{1}^\top \llbracket \mathbf{X} \rrbracket$$

$$\rightarrow (\mathbf{1}^\top \llbracket \mathbf{X} \rrbracket)^\top (\mathbf{1}^\top \llbracket \mathbf{X} \rrbracket)$$

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)

## Non-conventional operations can be much cheaper with tensor networks

### State of the art

$\mathbf{X}$

$$\rightarrow \llbracket \mathbf{X} \rrbracket$$

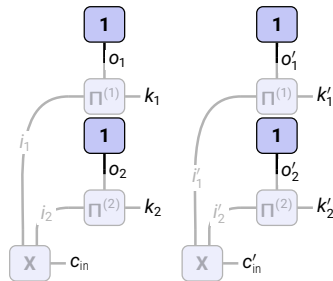
$$\rightarrow \mathbf{1}^\top \llbracket \mathbf{X} \rrbracket$$

$$\rightarrow (\mathbf{1}^\top \llbracket \mathbf{X} \rrbracket)^\top (\mathbf{1}^\top \llbracket \mathbf{X} \rrbracket)$$

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)



- ✦ 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

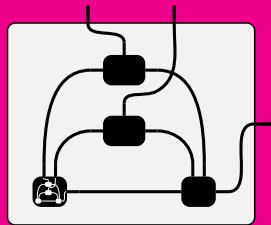


- ✦ TN perspective simplifies the transfer of algorithmic ideas
- ✦ Enables flexible/faster implementations of black box routines
- ✦ Relies on automatically efficient evaluation inside einsum

## 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

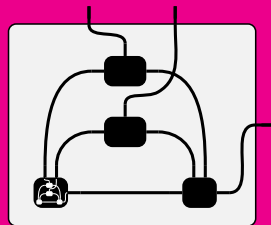


- ✦ TN perspective simplifies the transfer of algorithmic ideas
- ✦ Enables flexible/faster implementations of black box routines
- ✦ Relies on automatically efficient evaluation inside einsum

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

**Code:** [github.com/f-dangel/einconv](https://github.com/f-dangel/einconv)

Thank you 🙏 questions?

- 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.
- 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.
- James Martens and Roger Grosse. Optimizing neural networks with Kronecker-factored approximate curvature. In **International Conference on Machine Learning (ICML)**, 2015.
- 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.
- 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.



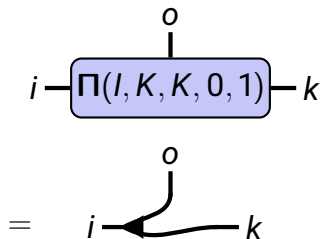
**For structured patterns, we can re-wire the TN before evaluation  
→ even better performance**

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

**Dense convolution ( $K = S$ )**

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

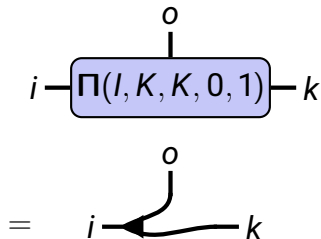
Dense convolution ( $K = S$ )



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

Dense convolution ( $K = S$ )

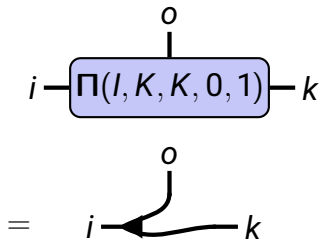
Down-sampling convolution ( $S > K$ )





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

Dense convolution ( $K = S$ )



Down-sampling convolution ( $S > K$ )

