# On the Relationship Between Self-Attention and Convolutional Layers

**Jean-Baptiste Cordonnier**
*with Andreas Loukas and Martin Jaggi*

Swiss Machine Learning Days
November 13th, 2019

# Conclusion

A Multi-Head Self-Attention Layer

can express any Convolutional Layer.

**Building block of state of the art NLP models**

Transformers
(Vaswani et al. 2017)
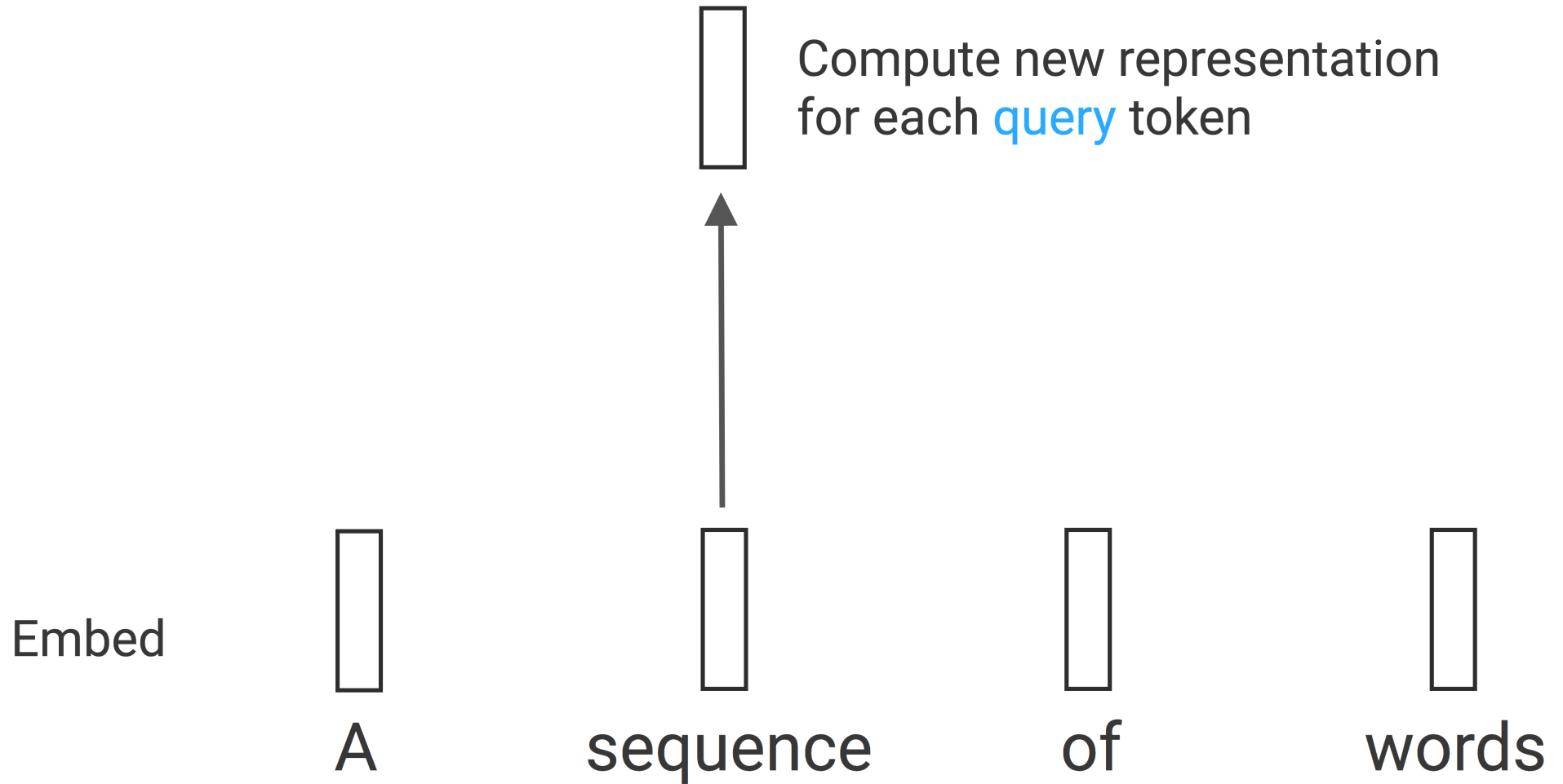
GPT2
(Radford et al. 2018)

BERT
(Devlin et al. 2019)

**are competitive with CNNs.**

applied to vision task,
achieve same performance,
at same computation cost.

(Bello et al. 2019)
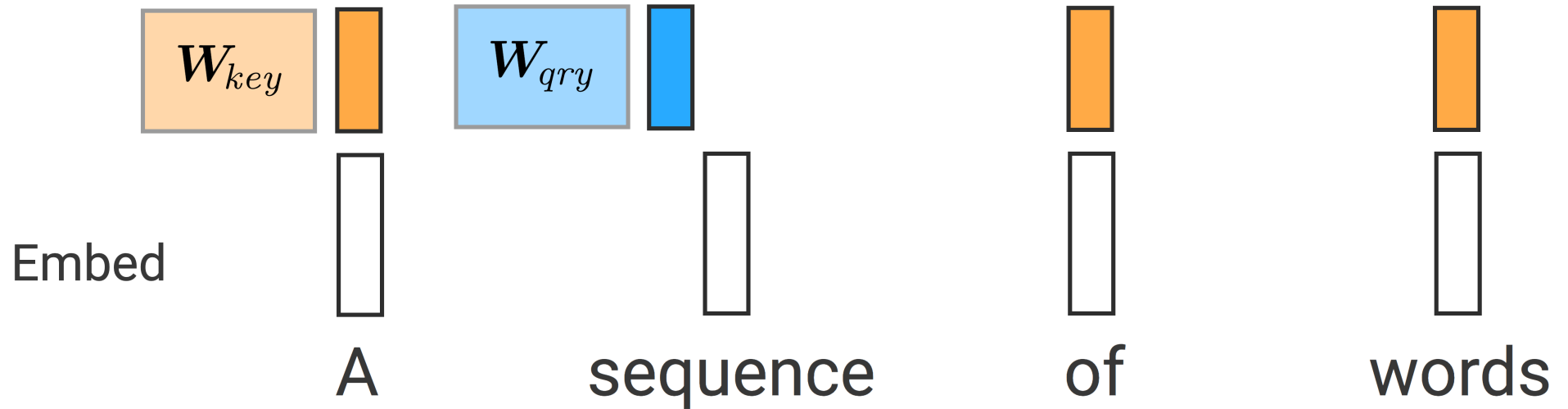(Ramachandran et al. 2019)
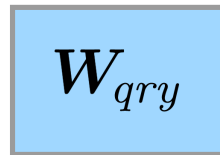
# Self-Attention (Vaswani et al. 2017)

Compute new representation for each query token

Embed

A    sequence    of    words

# Self-Attention (Vaswani et al. 2017)



Compute new representation for each query token

$W_{key}$

$W_{qry}$

Embed

A      sequence      of      words

# Self-Attention (Vaswani et al. 2017)
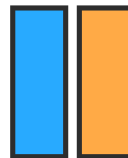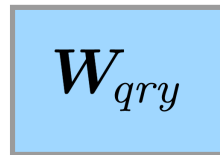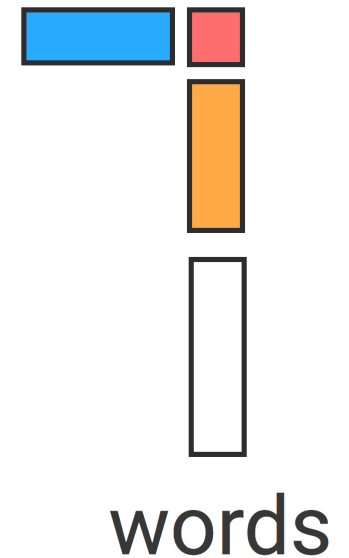
Parameters

$W_{key}$

$W_{qry}$

Compute new representation
for each query token

Embed

A          sequence          of          words

# Self-Attention (Vaswani et al. 2017)

Parameters

$W_{key}$

$W_{qry}$

Compute new representation for each query token

Embed

A          sequence          of          words

# Self-Attention (Vaswani et al. 2017)



Parameters

$W_{key}$

$W_{qry}$

$W_{val}$

Embed

A          sequence          of          words

# Self-Attention (Vaswani et al. 2017)

Parameters

$\boldsymbol{W}_{key}$

$\boldsymbol{W}_{qry}$

$\boldsymbol{W}_{val}$

Embed

A          sequence          of          words
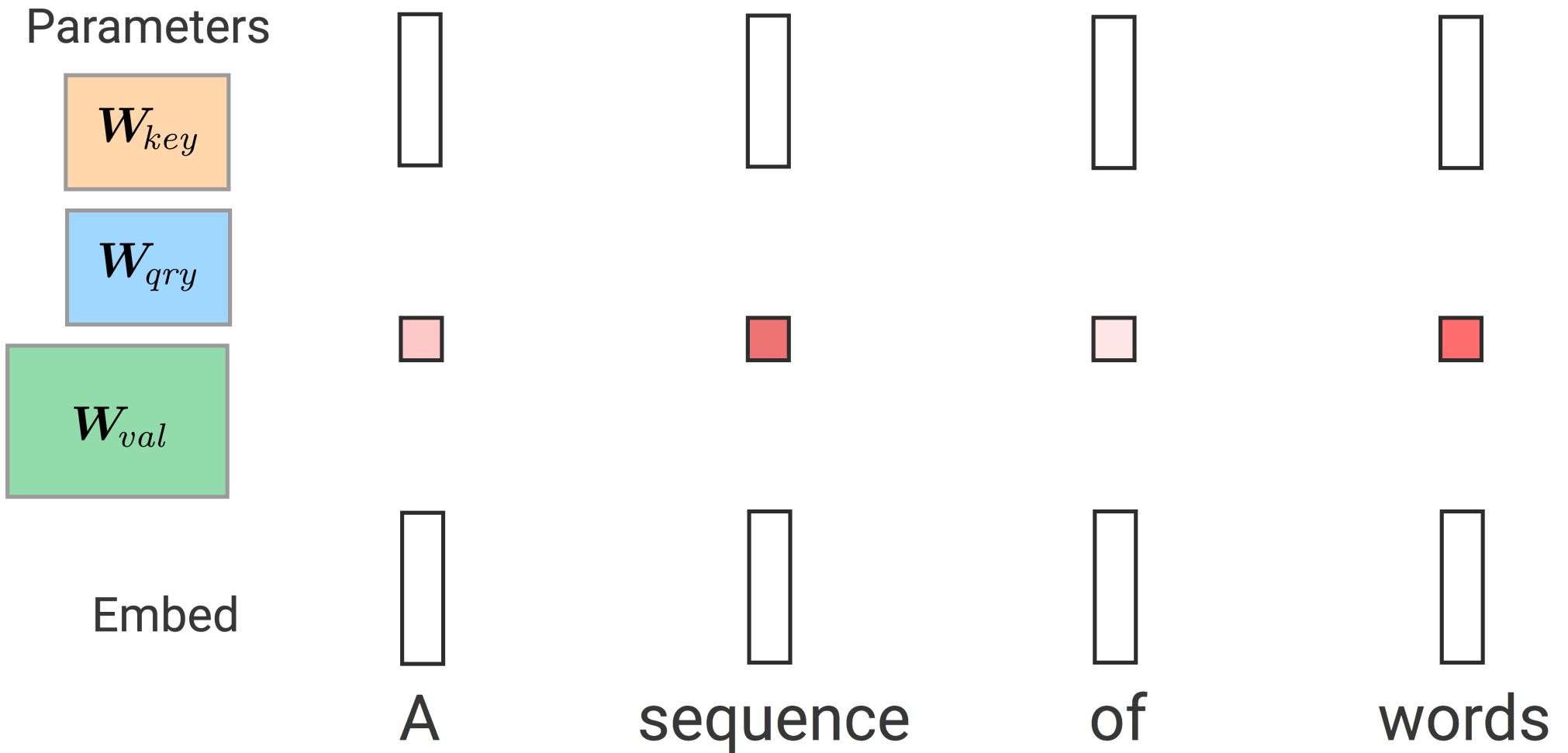
# Self-Attention (Vaswani et al. 2017)

Parameters

$W_{key}$

$W_{qry}$

$W_{val}$

Attention probabilities

Embed

A          sequence          of          words

# Self-Attention on Images



$W$

$D_{in}$

$H$

a key pixel    the query pixel

**Attention probabilities**

$\mathbf{A}_{\boldsymbol{q},:}^{(1)}$

$\mathbf{A}_{\boldsymbol{q},:}^{(2)}$

$\mathbf{A}_{\boldsymbol{q},:}^{(3)}$

**Parameters**

$D_k$

$D_{in}$ $\boldsymbol{W}_{qry}^{(1)}$

$\boldsymbol{W}_{qry}^{(2)}$

$\boldsymbol{W}_{qry}^{(3)}$

$N_h$

$D_k$

$D_{in}$ $\boldsymbol{W}_{key}^{(1)}$

$\boldsymbol{W}_{key}^{(2)}$

$\boldsymbol{W}_{key}^{(3)}$

$D_h$

$D_{in}$ $\boldsymbol{W}_{val}^{(1)}$

$\boldsymbol{W}_{val}^{(2)}$

$\boldsymbol{W}_{val}^{(3)}$

# Self-Attention on Images

**Attention probabilities**

$\mathbf{A}_{q,:}^{(1)}$

$\mathbf{A}_{q,:}^{(2)}$

$\mathbf{A}_{q,:}^{(3)}$

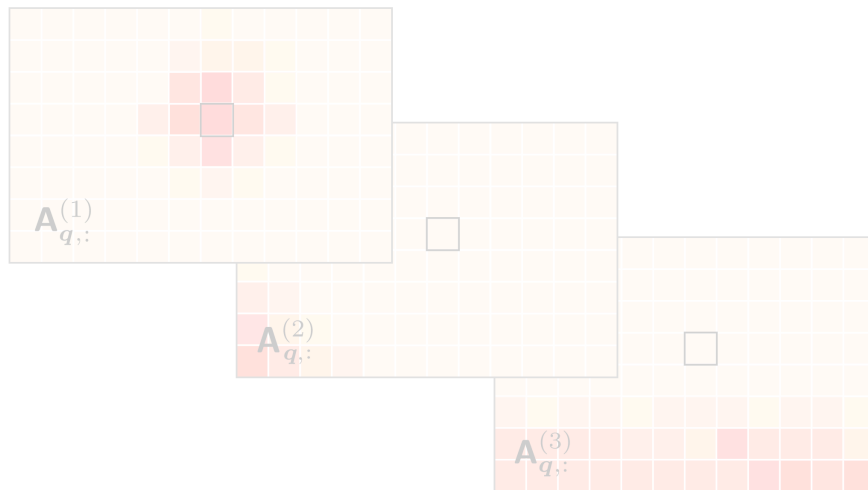**Same performance as ResNet on ImageNet (Ramachandran, 2019)**


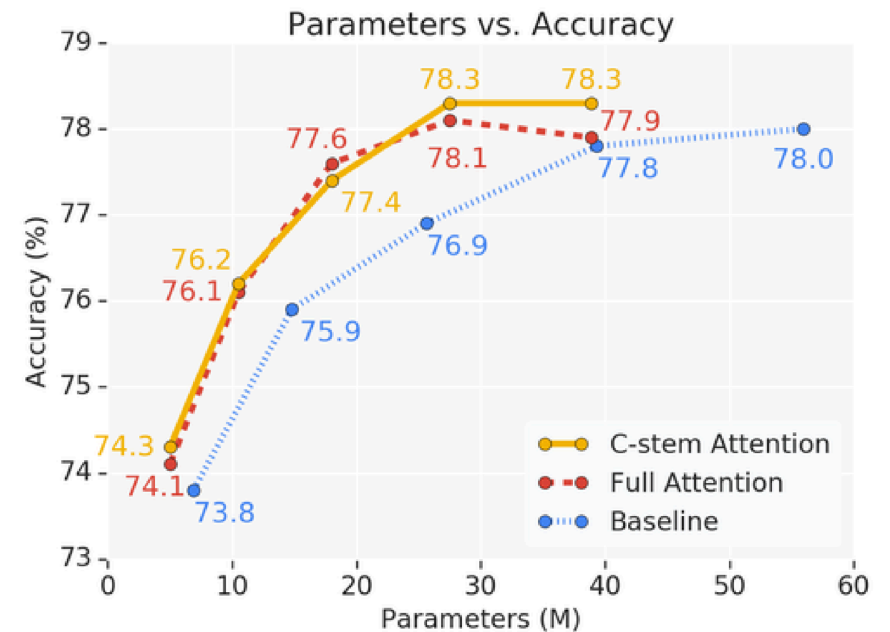
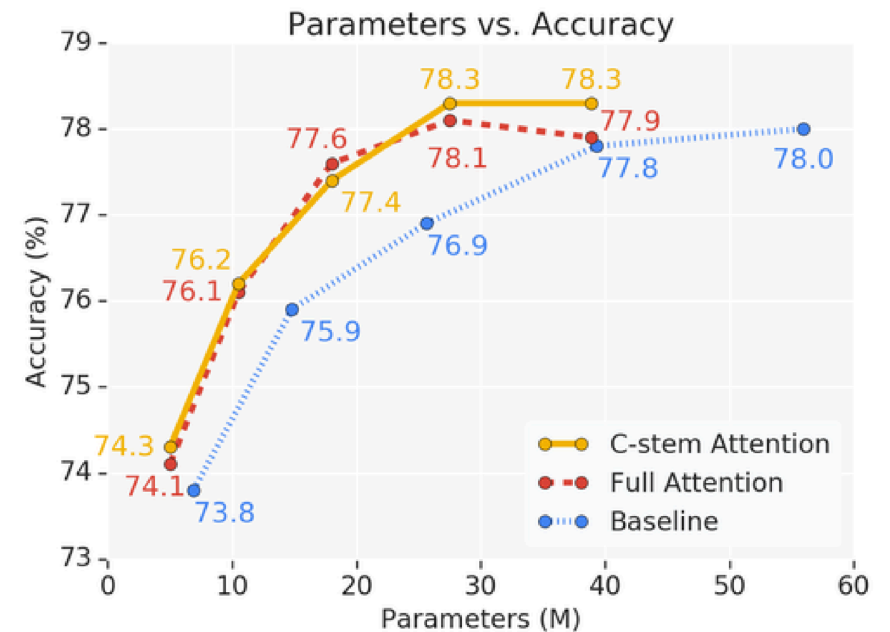Parameters vs. Accuracy

# Self-Attention on Images

**Attention probabilities**

**Same performance as ResNet on ImageNet (Ramachandran, 2019)**
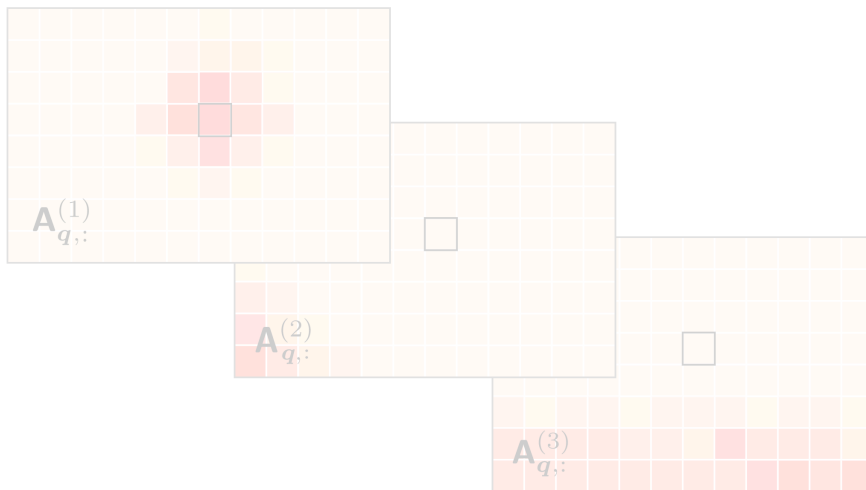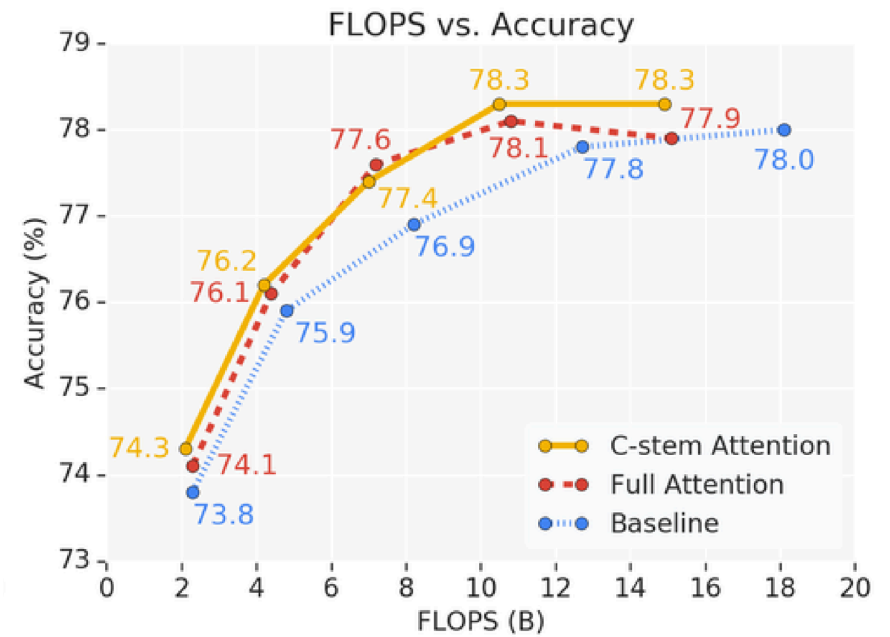
# Self-Attention on Images

**Attention probabilities**



$\mathbf{A}_{q,:}^{(1)}$

$\mathbf{A}_{q,:}^{(2)}$

$\mathbf{A}_{q,:}^{(3)}$

**Same performance as ResNet on ImageNet (Ramachandran, 2019)**



FLOPS vs. Accuracy

- C-stem Attention
- Full Attention
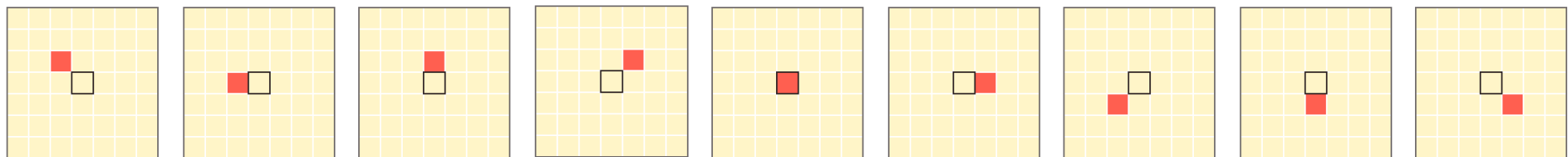- Baseline

# Self-Attention on Images
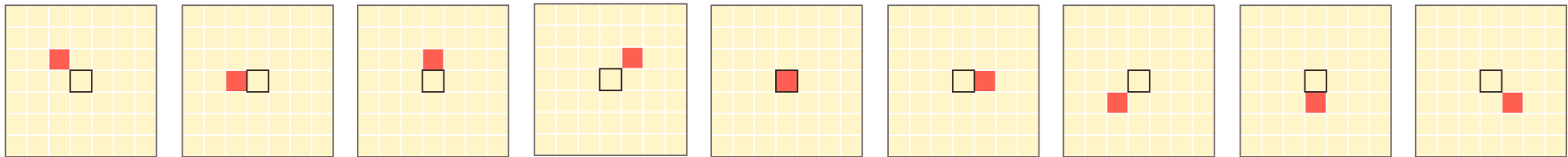
## Attention probabilities



## What if attention probabilities could look like this?

# Self-Attention on Images

**What if attention probabilities could look like this?**



**Then the multi-head self-attention could express any 3x3 convolution**

**Theorem 1.** *A multi-head self-attention layer with $N_h$ heads of dimension $D_h$, output dimension $D_{out}$ and a relative positional encoding of dimension $D_p \geq 3$ can express any convolutional layer of kernel size $\sqrt{N_h} \times \sqrt{N_h}$ and $\min(D_h, D_{out})$ output channels.*

# Self-Attention on Images

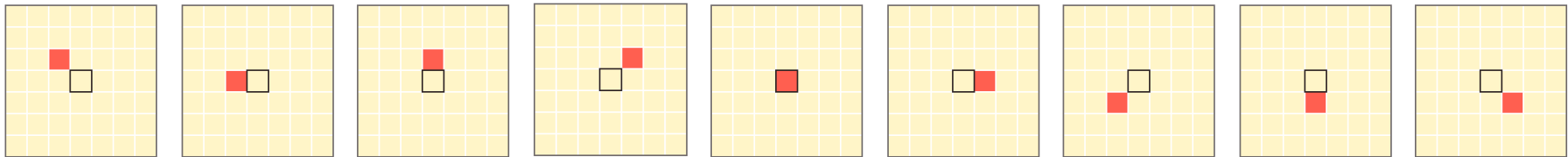**What if attention probabilities could look like this?**



**Then the multi-head self-attention could express any 3x3 convolution**

**Theorem 1.** *A multi-head self-attention layer with $N_h$ heads of dimension $D_h$, output dimension $D_{out}$ and a relative positional encoding of dimension $D_p \geq 3$ can express any convolutional layer of kernel size $\sqrt{N_h} \times \sqrt{N_h}$ and $\min(D_h, D_{out})$ output channels.*

# Self-Attention on Images

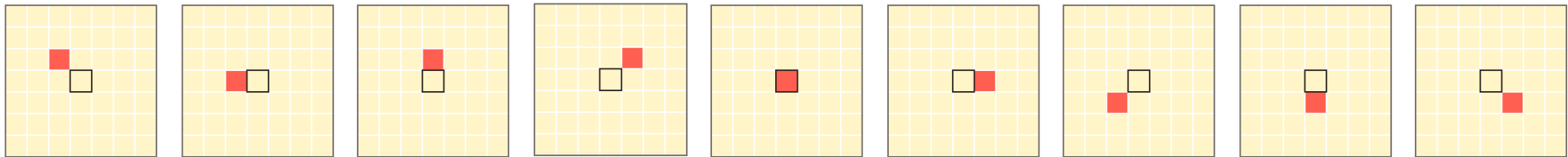**What if attention probabilities could look like this?**



**Then the multi-head self-attention could express any 3x3 convolution**

**Theorem 1.** *A multi-head self-attention layer with $N_h$ heads of dimension $D_h$, output dimension $D_{out}$ and a* relative positional encoding *of dimension $D_p \geq 3$ can express any convolutional layer of kernel size $\sqrt{N_h} \times \sqrt{N_h}$ and $\min(D_h, D_{out})$ output channels.*

# Conclusion

A Multi-Head Self-Attention Layer

can express any Convolutional Layer.

**Building block of state of the art NLP models**

Transformers                    GPT2
(Vaswani et al. 2017)     (Radford et al. 2018)

BERT
(Devlin et al. 2019)

**are competitive with CNNs.**

applied to vision task,
achieve same performance,
at same computation cost.

(Bello et al. 2019)
(Ramachandran et al. 2019)