

# Deep Learning: models and algorithms

Christian Wolf

July 8<sup>th</sup>, 2021

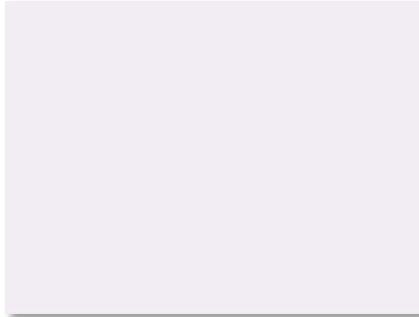
# Our group

Christian Wolf  
INSA-Lyon, LIRIS UMR CNRS 5205  
[liris.cnrs.fr/christian.wolf](http://liris.cnrs.fr/christian.wolf)

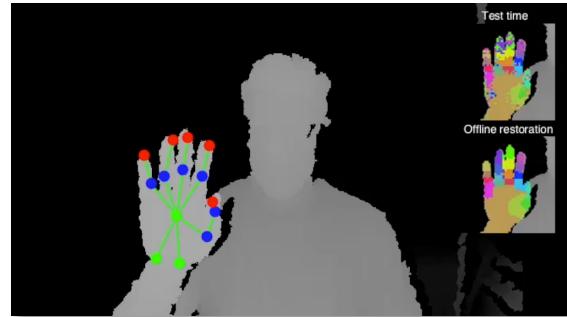


*The group in Feb. 2020: Corentin Kervadec, Steeven Janny, Edward Beeching, Fabien Baradel, Théo Jaunet, Quentin Possamaï.*

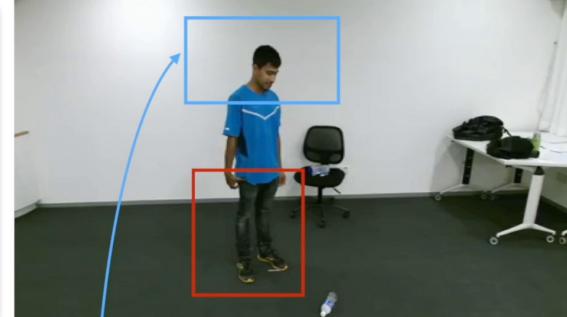
# Our group: learning vision & robotics



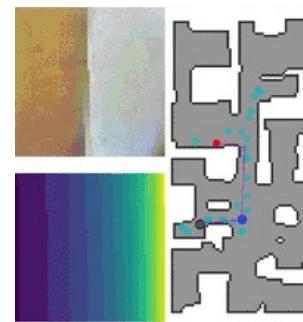
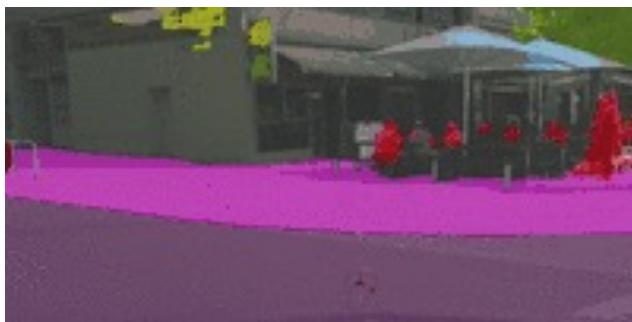
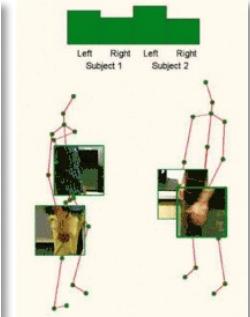
Gesture  
recognition



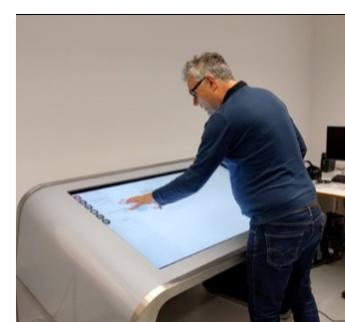
Pose estimation



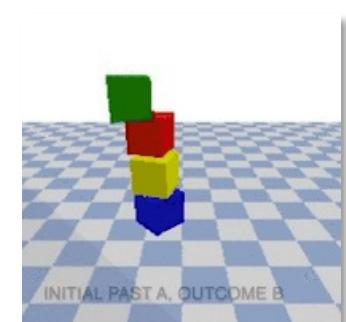
Activity Recognition



Robot Perception and Navigation

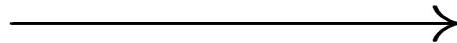


H-C Interaction

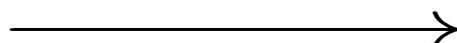


Physics

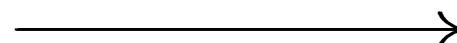
# Taking decisions



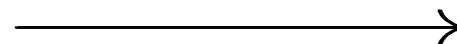
{dog, **cat**, avocado, chair, ...}



{0, 1, ... 26, 27, **28**..., 98, 99, ...}

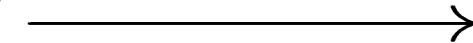


{Left, right, forward, backward, ...}



Motor control

“A blue parrot with a yellow belly  
sitting on a branch in a forest”



# The 3 problems of Machine Learning

## 1. Expressivity

- What is the complexity of the functions my model can represent?

## 2. Trainability

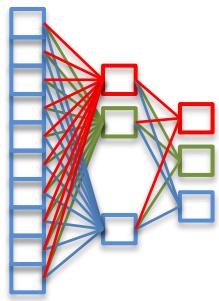
- How easy is training of my model (i.e. solving the optimization problem)?

## 3. Generalization

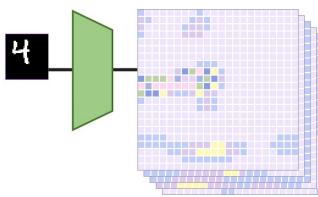
- How does my model behave on unseen data?
- In presence of a shift in distributions?

(D'après Eric Jang & Jascha Sohl-Dickstein)

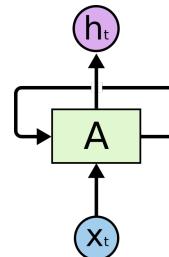
# The Deep Toolbox



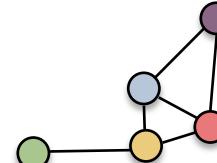
MLP



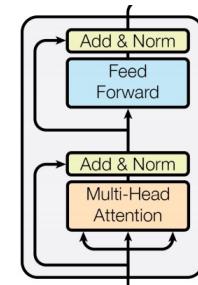
CNN /  
Convolutions



RNN /  
Recurrence



GN, GCN /  
Graphs, geometry



Transformers /  
Self-attention

*What do I know about the data and the task?*

*Nothing  
(vector space)*

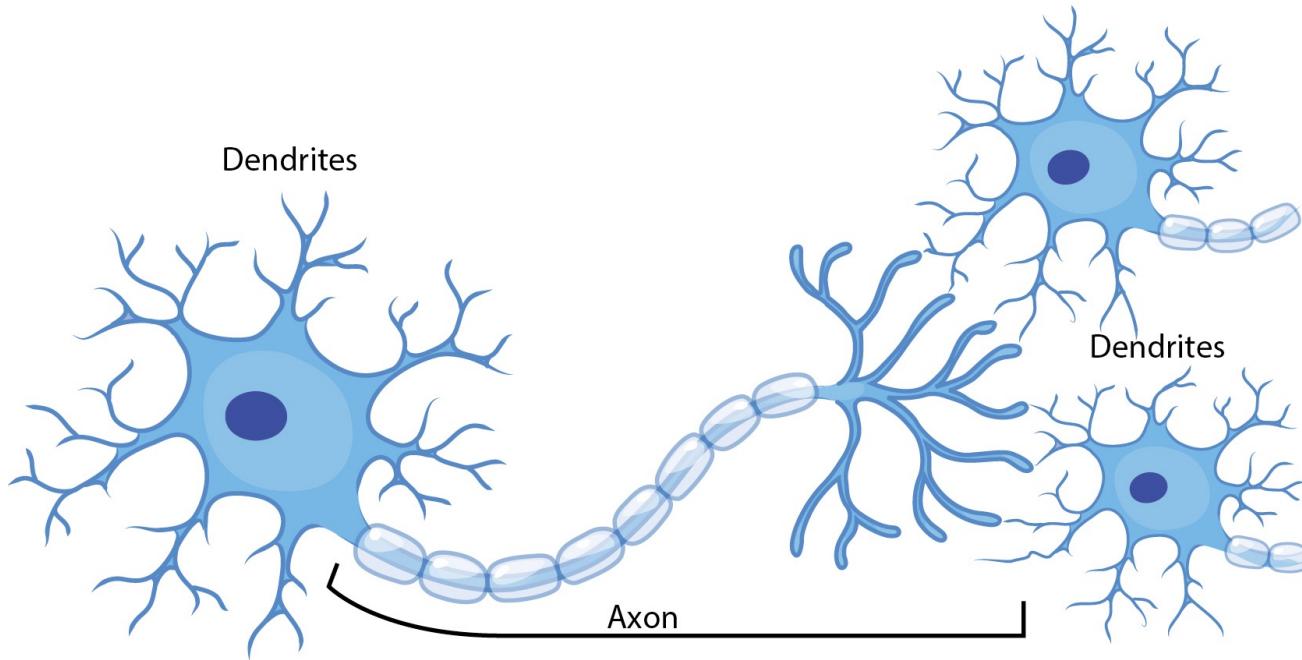
*Translation  
equivariance*

*Sequential data,  
Markov property*

*Graph structured  
data*

*Permutation  
equivariance*

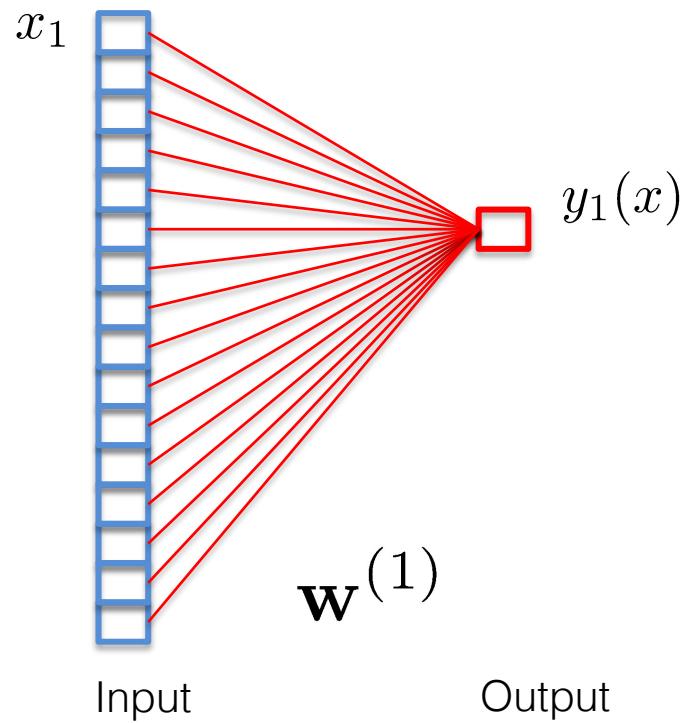
# Biological neurons



Devin K. Phillips

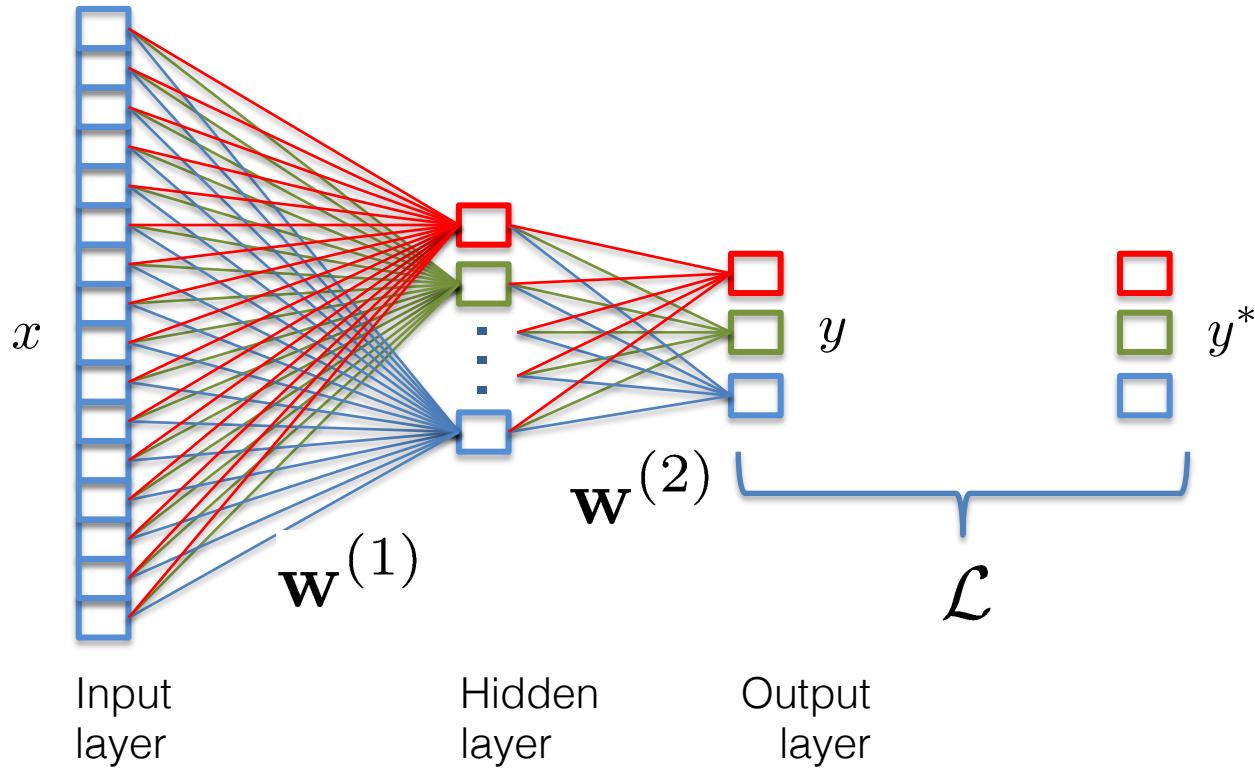
# Neural networks

## « Perceptron »



$$y(\mathbf{x}, \mathbf{w}) = \sum_{i=0}^D \mathbf{w}_i \mathbf{x}_i$$

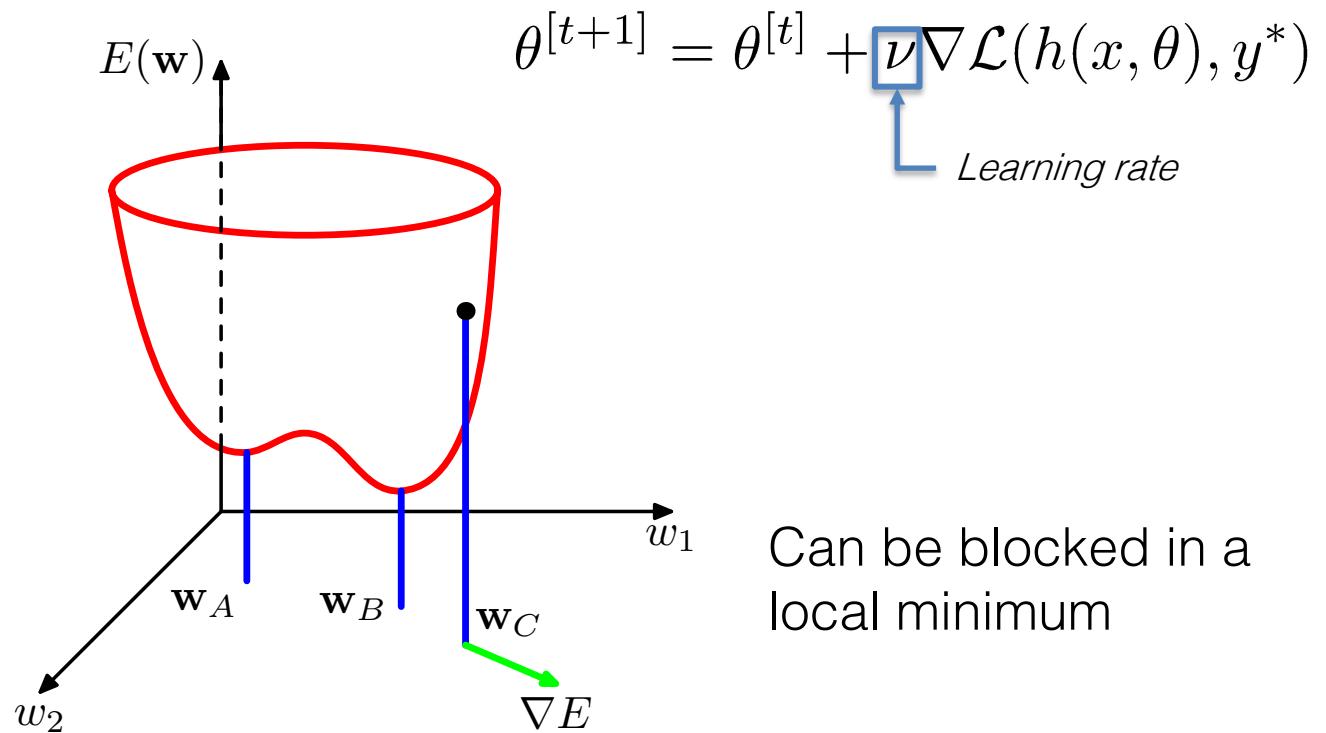
# Deep neural networks



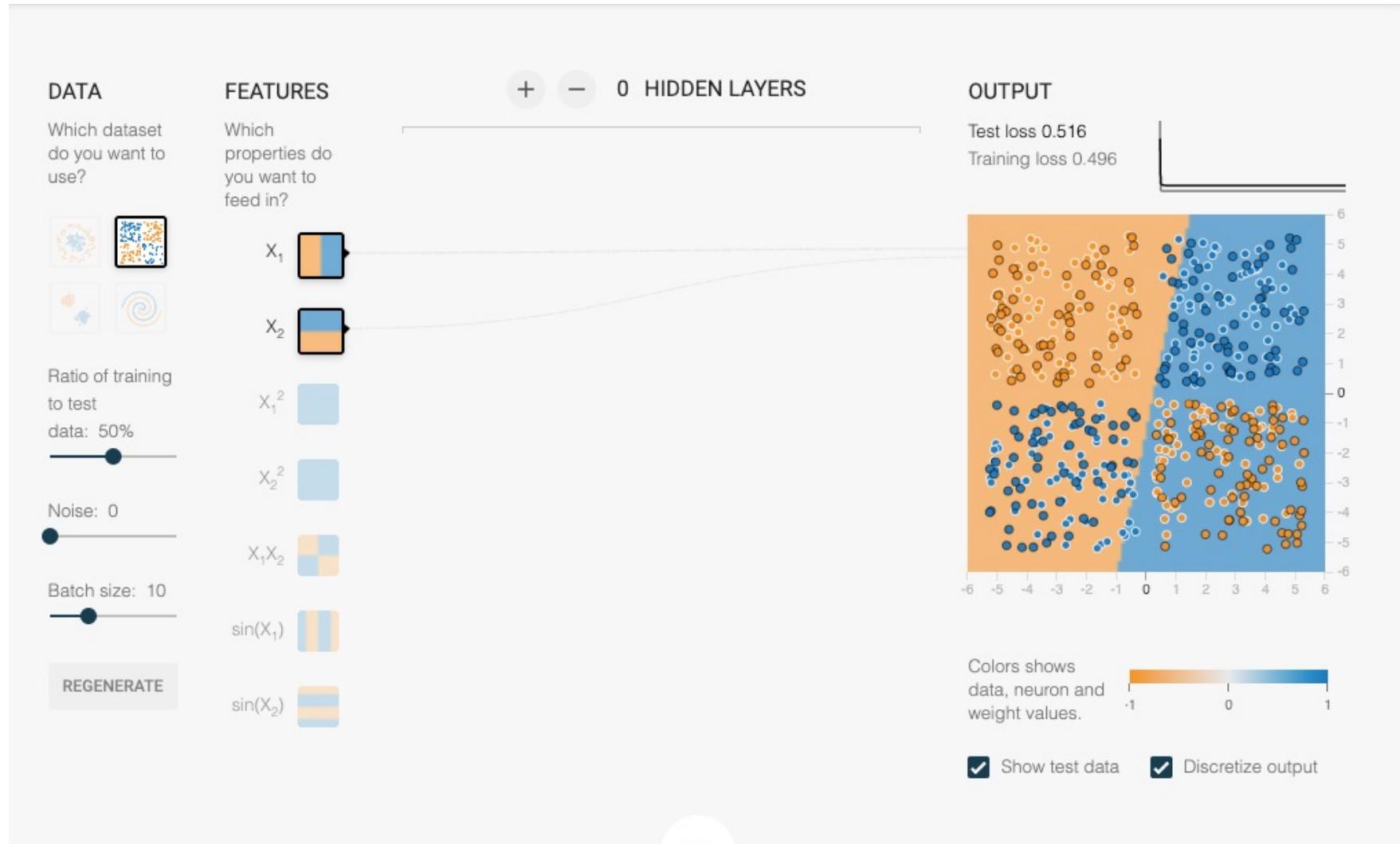
# Gradient descent

Minimize the error on known data

"Empirical Risk Minimization"

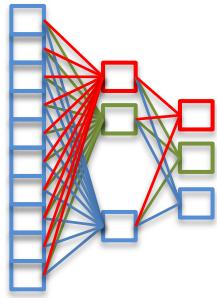


# Tensorflow Playground

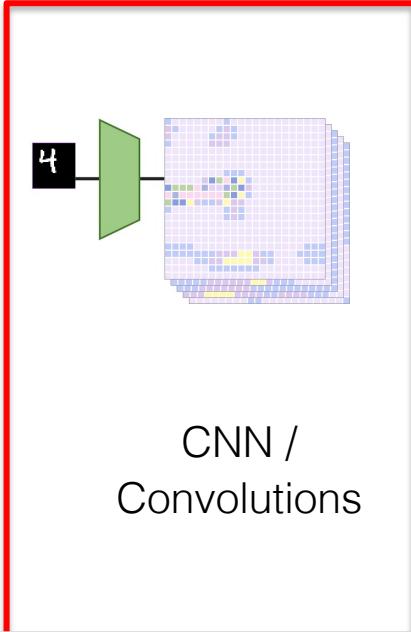


<https://playground.tensorflow.org>

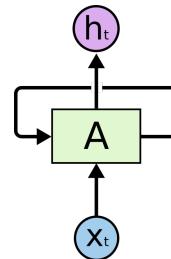
# The Deep Toolbox



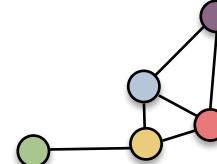
MLP



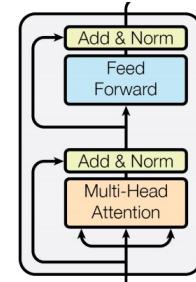
CNN /  
Convolutions



RNN /  
Recurrence



GN, GCN /  
Graphs, geometry



Transformers /  
Self-attention

*What do I know about the data and the task?*

*Nothing  
(vector space)*

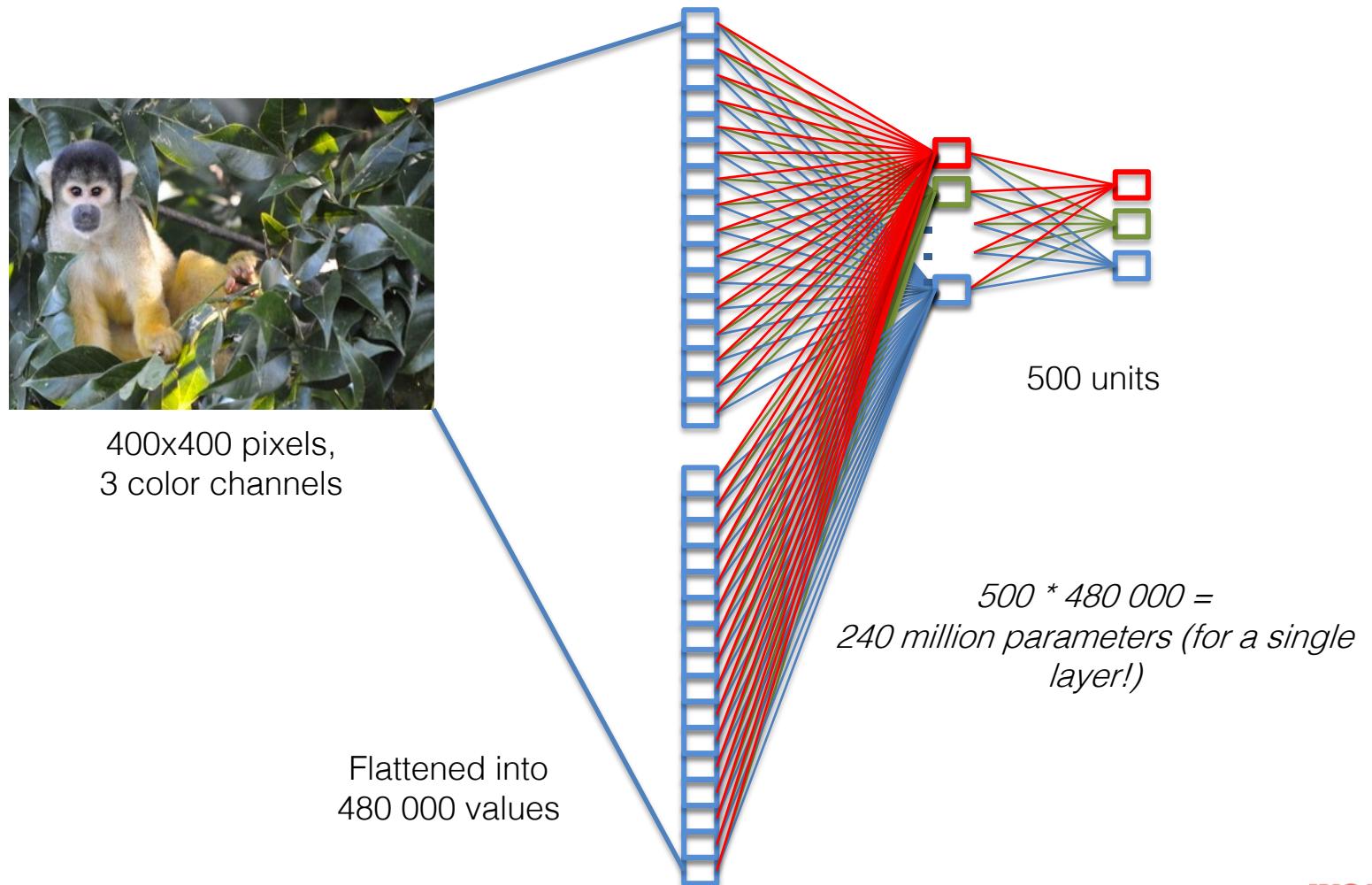
*Translation  
equivariance*

*Sequential data,  
Markov property*

*Graph structured  
data*

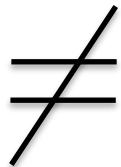
*Permutation  
equivariance*

# Fully connected layers might be harmful



# Fully connected layers might be harmful

Parameters learned for one part of the image do not generalize to other parts of the image.



# Translation invariance

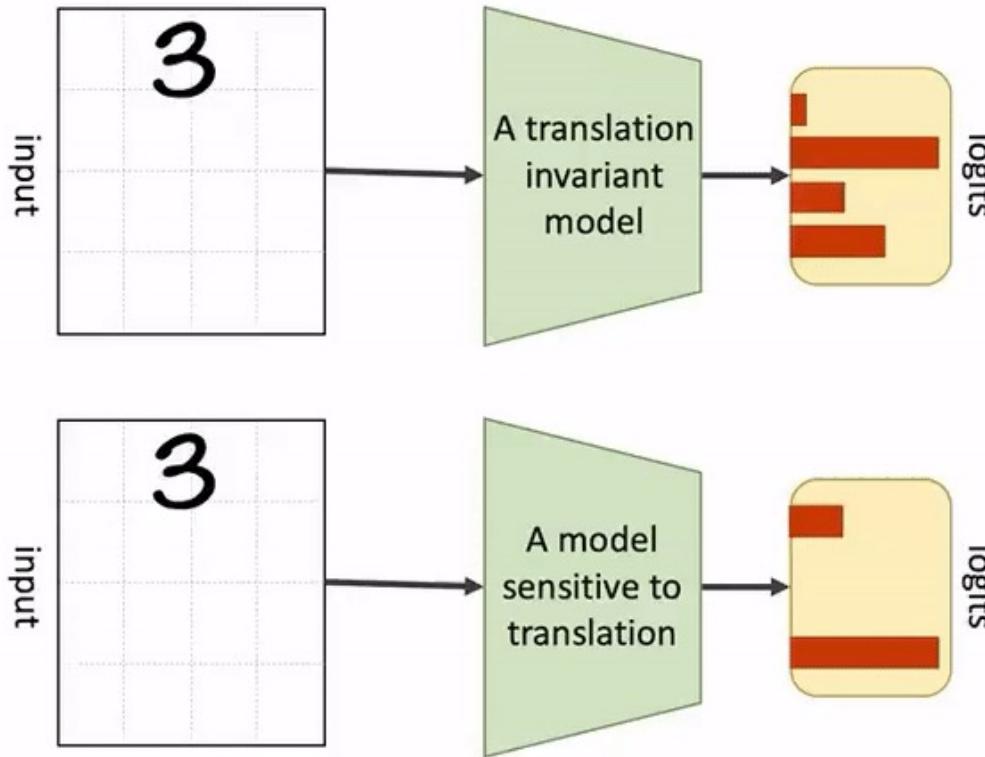
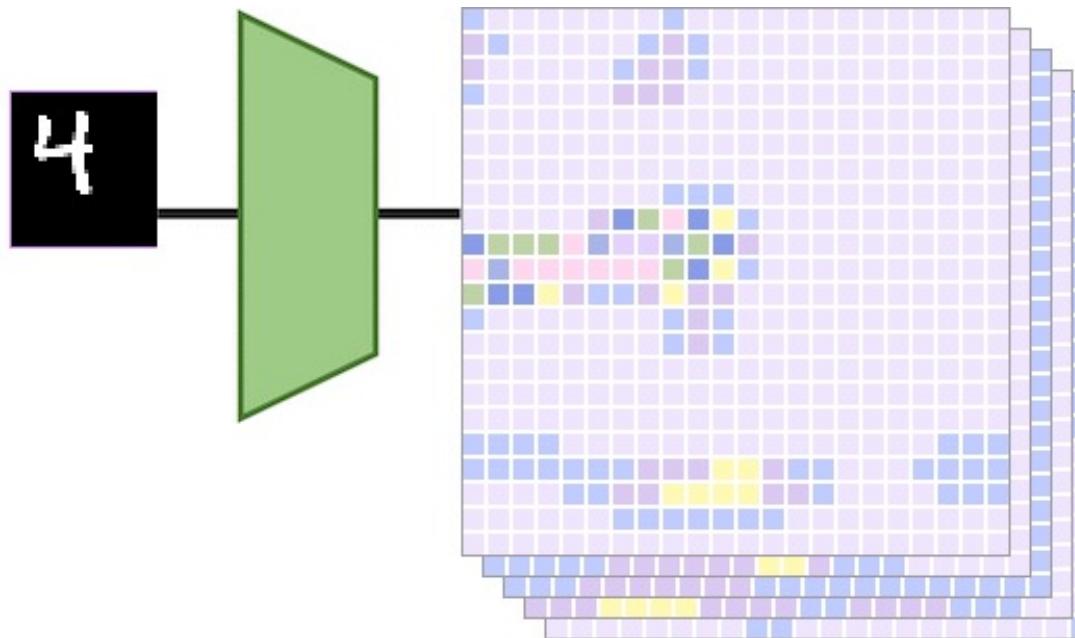


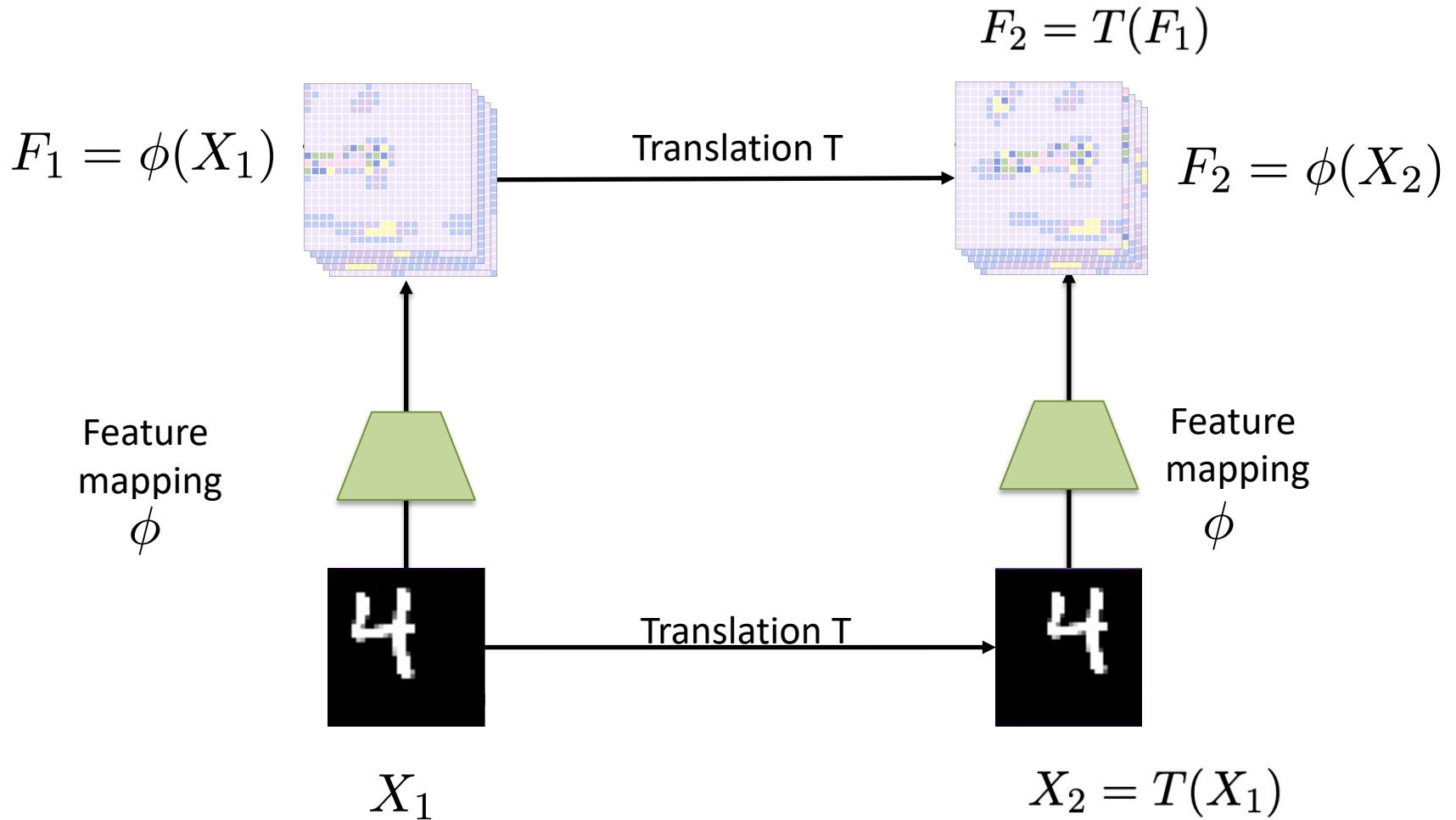
Figure by Samira Abnar (U. of Amsterdam)  
<https://samiraabnar.github.io/>

# Translation equivariance

An equivariant mapping preserves the algebraic structure of the transformation.



# Translation Equivariance



# Shift-equivariant linear operations

Classical MLPs are sequences of linear layers followed by point-wise non-linearities. What kind of linear and shift-equivariant operators can there exist?

Let's consider an operator  $\emptyset$  with following properties:

**Linearity:**

$$\phi(\alpha f + \beta f') = \alpha\phi(f) + \beta\phi(f')$$

**Shift-equivariance:**

$$\phi({}^{m,n}S(f)) = {}^{m,n}S(\phi(f))$$

where  ${}^{m,n}S(f)$  shifts a signal by  ${}^{m,n}$ .

**Impulse response**  $h$ :

$$h = \phi({}^{0,0}p)$$

where  ${}^{0,0}p$  is a Dirac impulse centered at 0, 0.

# Shift-equivariant linear operations

We decompose the signal  $f$  into a series of Diracs  ${}^{m,n}p$ :

$$[\phi(f)](x, y) = \left[ \phi \left( \sum_{m=-M/2}^{M/2} \sum_{n=-N/2}^{N/2} f(m, n) {}^{m,n}p \right) \right] (x, y)$$

*Dirac at position  $m, n$*

We use **linearity**:

$$= \left[ \sum_{m=-M/2}^{M/2} \sum_{n=-N/2}^{N/2} f(m, n) \phi({}^{m,n}p) \right] (x, y)$$

We use **shift-equivariance**:

$$= \left[ \sum_{m=-M/2}^{M/2} \sum_{n=-N/2}^{N/2} f(m, n) \phi({}^{m,n}S({}^{0,0}p)) \right] (x, y)$$

We use **linearity** again:

$$= \left[ \sum_{m=-M/2}^{M/2} \sum_{n=-N/2}^{N/2} f(m, n) {}^{m,n}S(\phi({}^{0,0}p)) \right] (x, y)$$

# Shift-equivariant linear operations

Change in notation: we replace  $\phi(0,0)p$  by  $h$  (through definition):

$$[\phi(f)](x, y) = \left[ \sum_{m=-M/2}^{M/2} \sum_{n=-N/2}^{N/2} f(m, n)^{m,n} S(h) \right] (x, y)$$

Change in notation: make the shift-operator explicit:

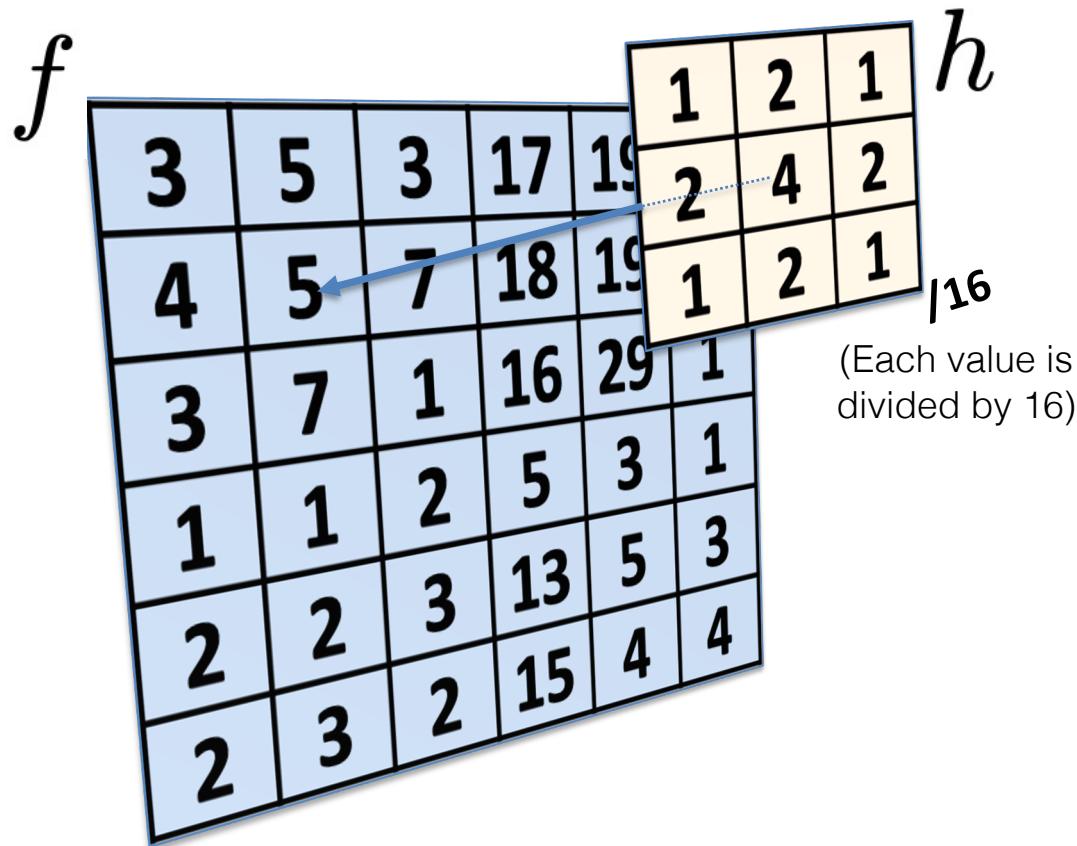
$$= \sum_{m=-M/2}^{M/2} \sum_{n=-N/2}^{N/2} f(m, n) h(x - m, y - n)$$

Change of variables:

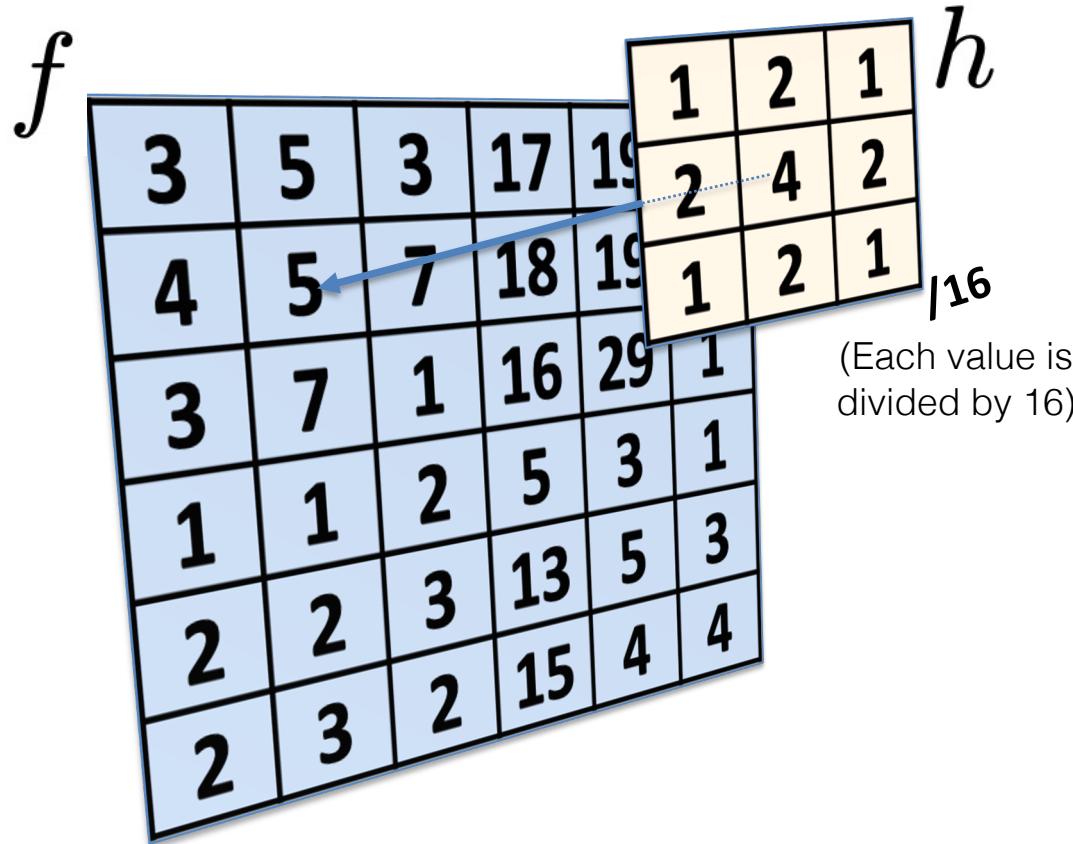
$$= \sum_{m'=-M/2}^{M/2} \sum_{n'=-N/2}^{N/2} f(x - m', y - n') h(m', n')$$

$(m' = x - m, n' = y - n) \Rightarrow$  we get a convolution!

# Convolutions



# Convolutions



A diagram showing the calculation of the output value 5. A circled asterisk (\*) connects the input feature map  $f$  and the kernel  $h$ . Below them are their respective 3x3 matrices. The formula for calculating the output value is provided:  $(3+2*5+3+2*4+5*5+2*7+3+2*7+1)/16 = 5$ .

3	5	3
4	5	7
3	7	1

\*

1	2	1
2	4	2
1	2	1

/16

$$(3+2*5+3+2*4+5*5+2*7+3+2*7+1)/16 = 5$$

# Convolutions

$f$

3	5	3	17	19
4	5	7	18	19
3	7	1	16	29
1	1	2	5	3
2	2	3	13	5
2	3	2	15	4
2	3	1		

1	2	1
2	4	2
1	2	1

$h$

(Each value is divided by 16)

$\phi(f)$


5

$$\begin{array}{c}
 * \\
 \begin{array}{|c|c|c|} \hline 3 & 5 & 3 \\ \hline 4 & 5 & 7 \\ \hline 3 & 7 & 1 \\ \hline \end{array} \quad \begin{array}{|c|c|c|} \hline 1 & 2 & 1 \\ \hline 2 & 4 & 2 \\ \hline 1 & 2 & 1 \\ \hline \end{array} /16
 \end{array}$$

$$\begin{aligned}
 &(3+2*5+3+2*4+5*5+2*7 \\
 &+3+2*7+1)/16 = 5
 \end{aligned}$$

# Convolutions

f

3	5	3	17	19	15
4	5	7	18	19	14
3	7	1	16	29	1
1	1	2	5	3	1
2	2	3	13	5	3
2	3	2	15	4	4

h

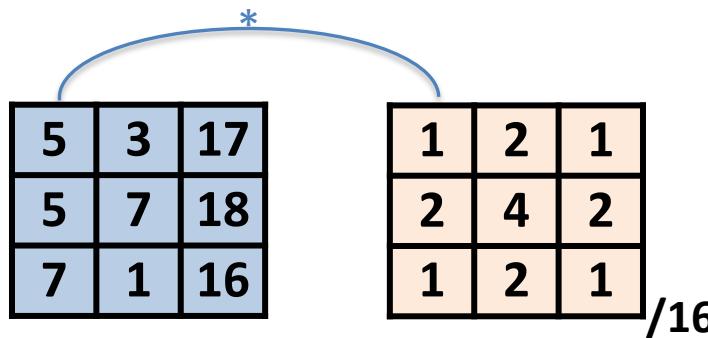
$$\phi(f)$$

| 16

(Each value divided by 1)

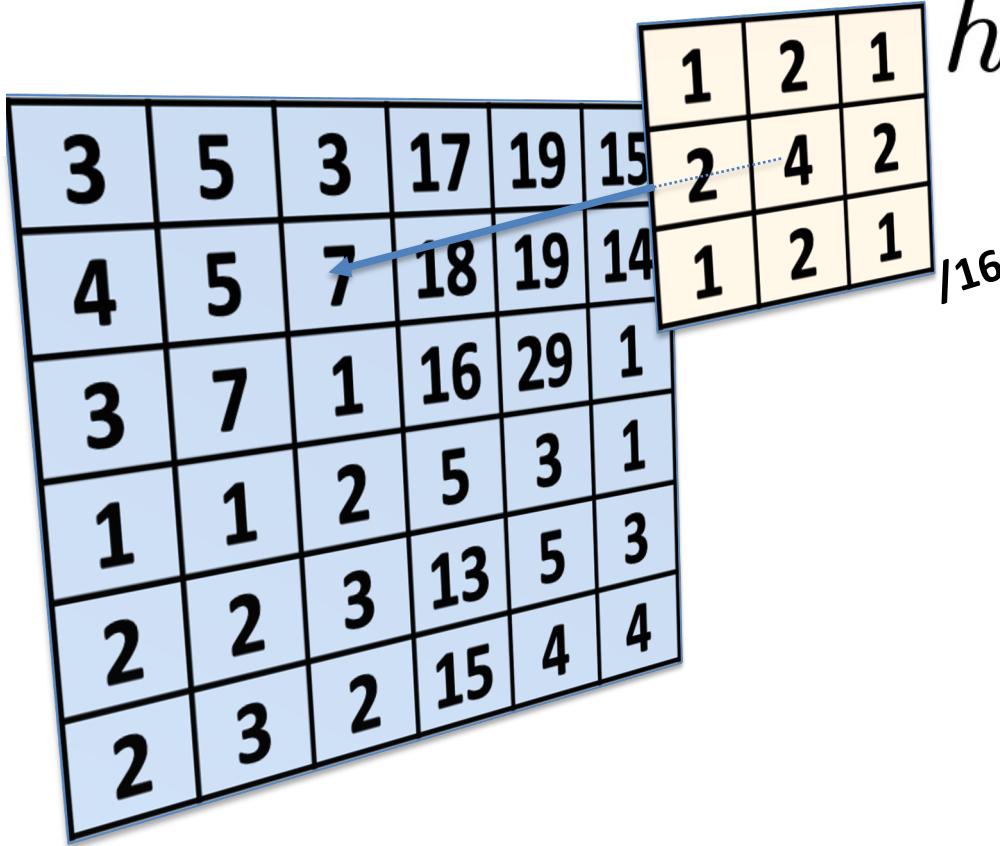
A 5x5 grid of light purple squares with black outlines. The top-left square contains the large black number '5'. The grid is tilted diagonally to the right.

/16



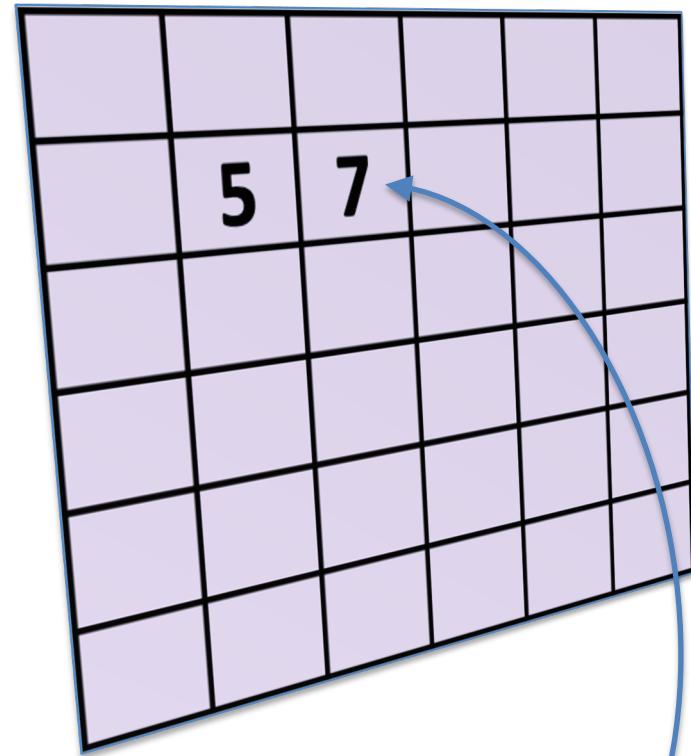
# Convolutions

$f$



$h$

$\phi(f)$



\*

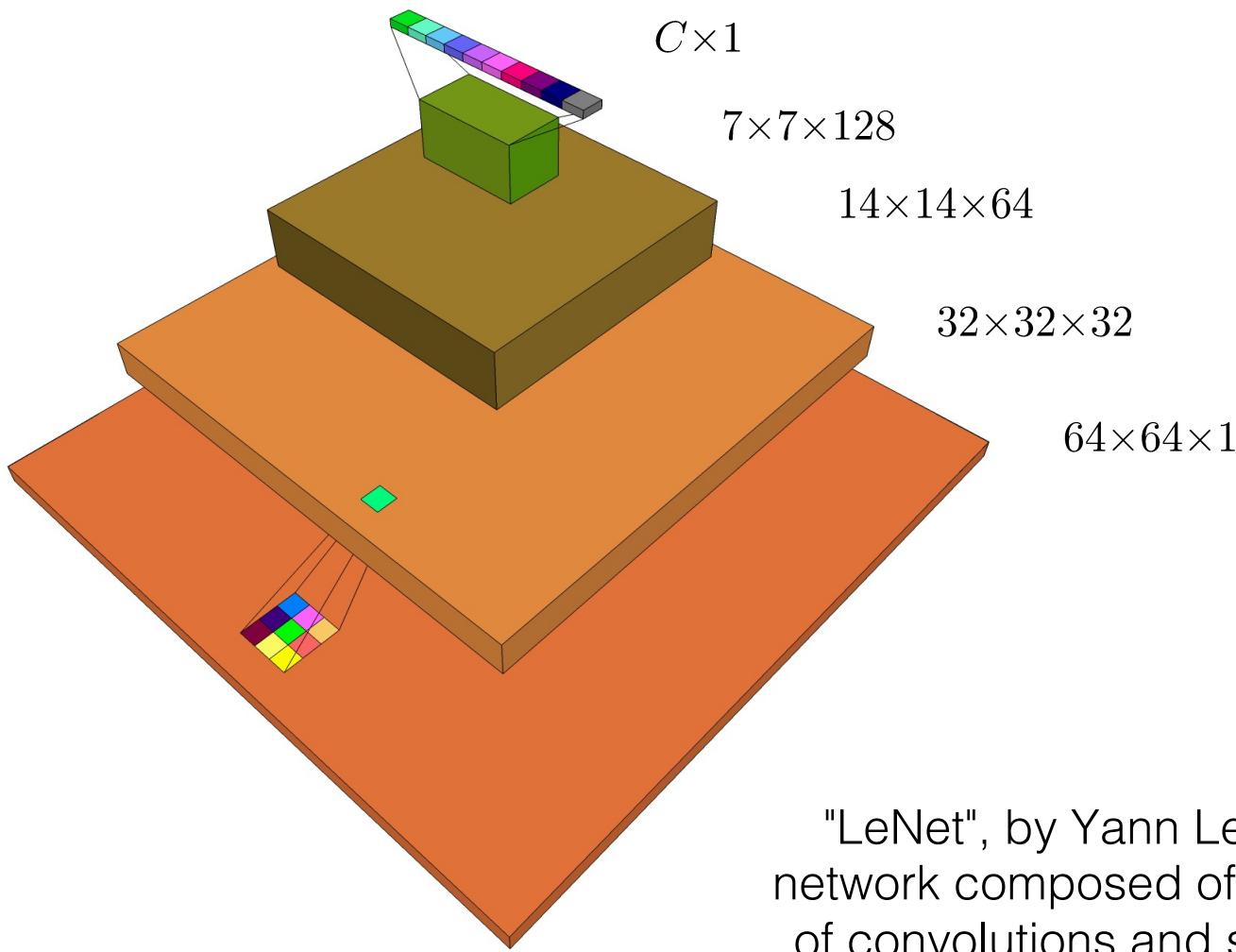
5	3	17
5	7	18
7	1	16

1	2	1
2	4	2
1	2	1

$/16$

$$(5+2*3+17+2*5+4*7+2*18+7+2*1+16)/16 = 7$$

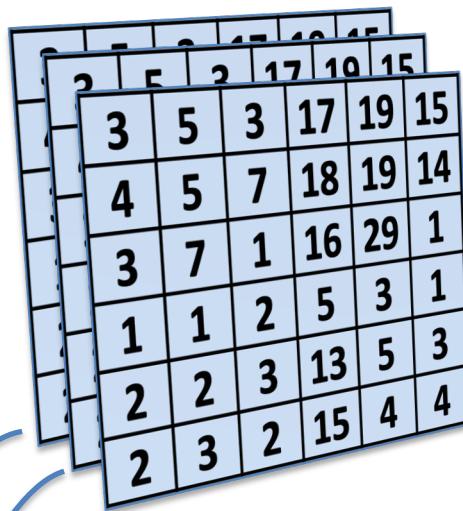
# « LeNet »



"LeNet", by Yann LeCun, is a network composed of sequences of convolutions and spatial size reductions (« pooling »).

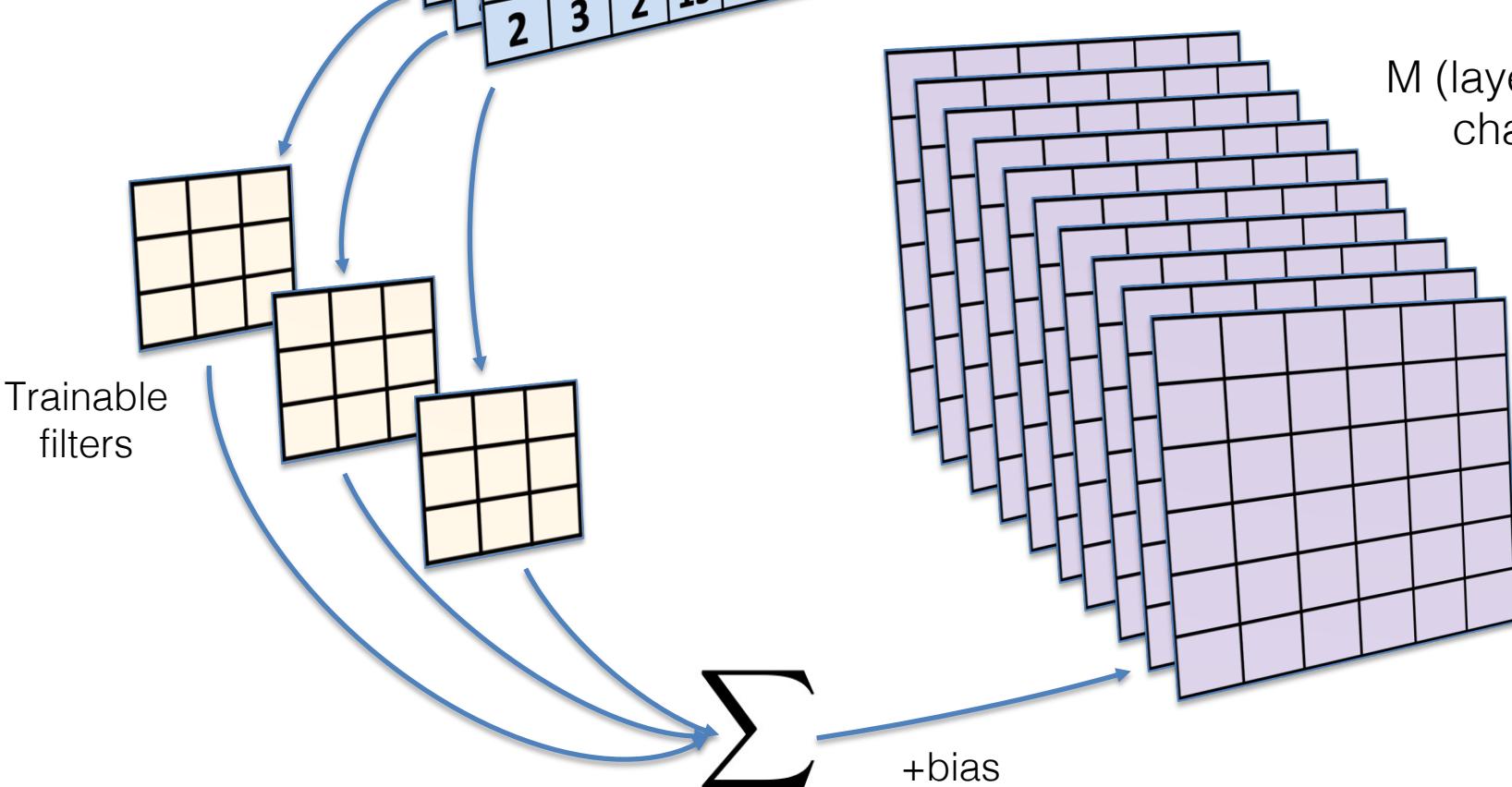
[LeCun et al., 1998]

$N$  (layer) input channels



Each of the  $M$  output channels is a linear combination of the  $N$  filtered input images.

We need  $N \times M$  filters!



# Template matching

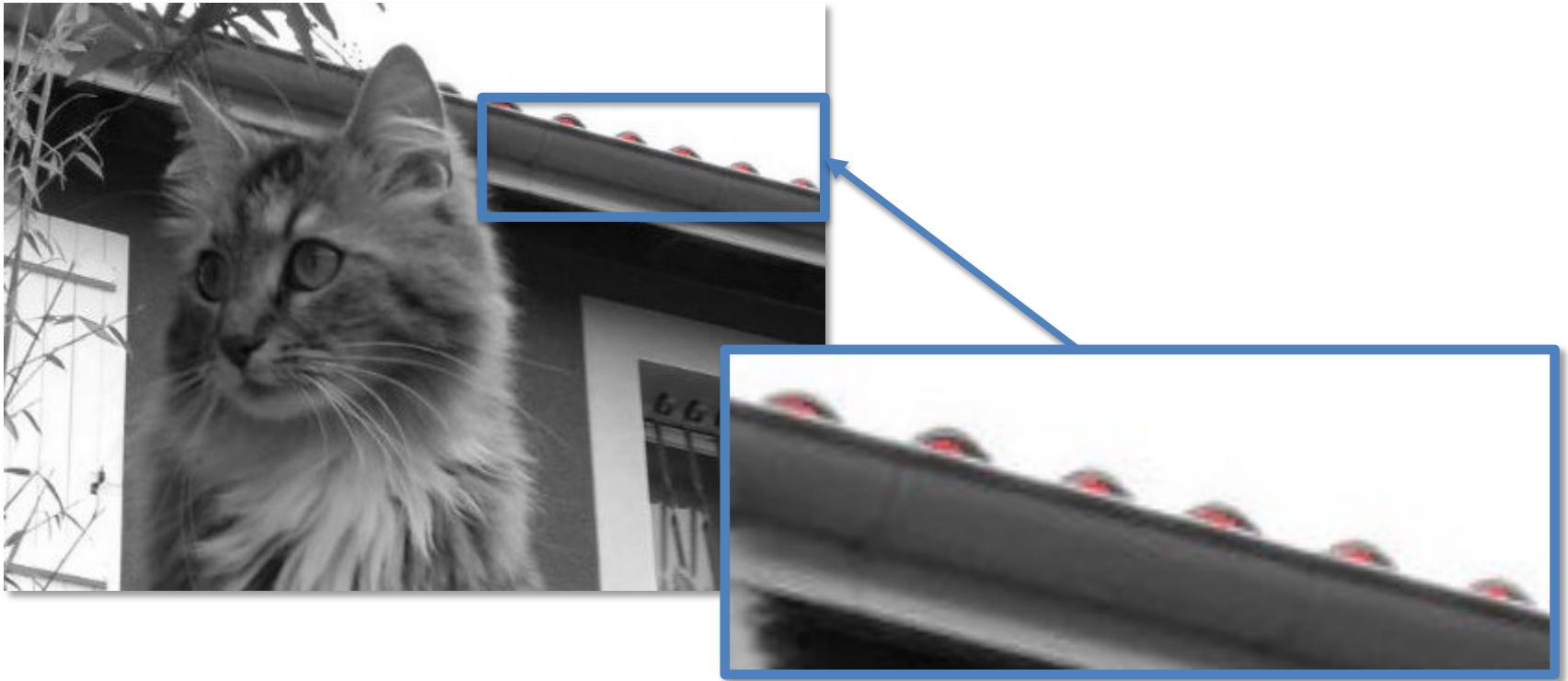
Input image:



Filter kernel

- We normalize input and kernel (subtract mean, divide by standard deviation).
- We convolve the input with the kernel
- We threshold response and superimpose it on the red channel.

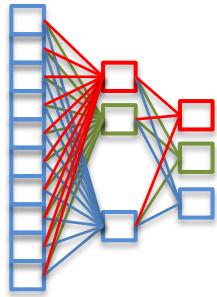
# Template matching



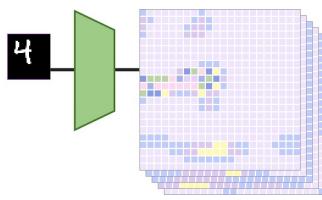
Convolutions can match patterns:

- either in input space
- or in intermediate network layers (which still have a spatial organization).

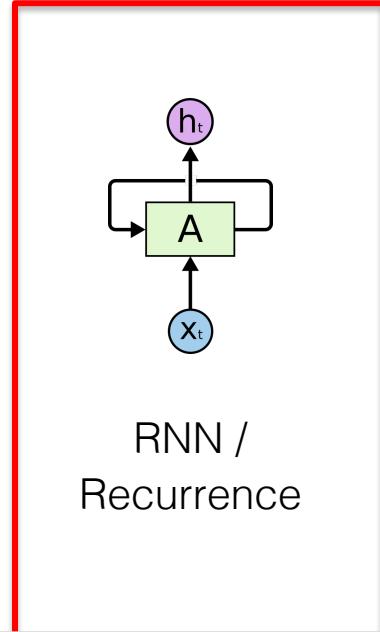
# The Deep Toolbox



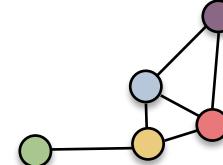
MLP



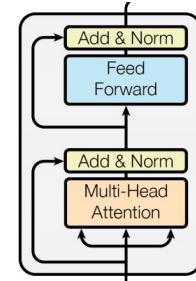
CNN /  
Convolutions



RNN /  
Recurrence



GN, GCN /  
Graphs, geometry



Transformers /  
Self-attention

*What do I know about the data and the task?*

*Nothing  
(vector space)*

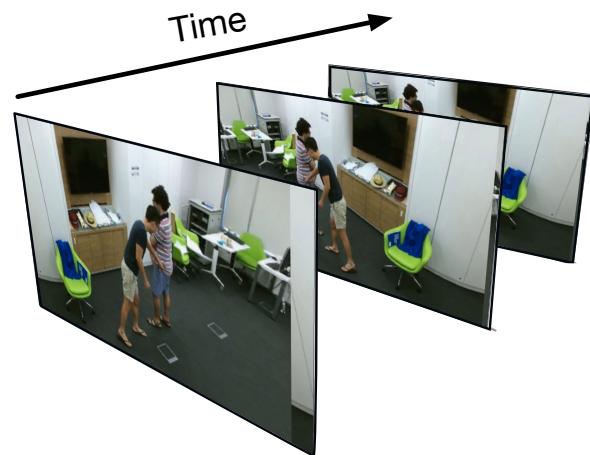
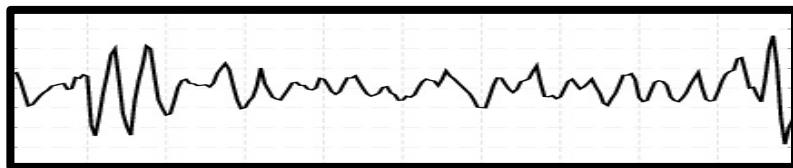
*Translation  
equivariance*

*Sequential data,  
Markov property*

*Graph structured  
data*

*Permutation  
equivariance*

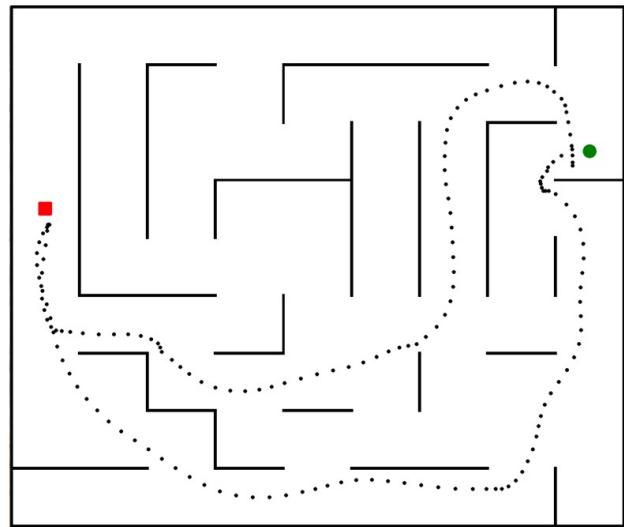
# Problem: dealing with sequences



Toutes les familles heureuses se ressemblent. Chaque famille malheureuse, au contraire, l'est à sa façon.



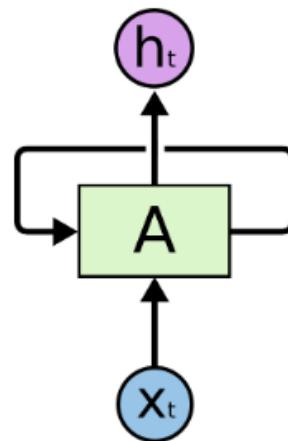
Happy families are all alike. Every unhappy family is unhappy in its own way.



# Recurrent Neural Networks (RNNs)

As feed-forward networks, Recurrent Neural Networks (RNNs) predict some output from a given input.

However, they also pass information over time, from instant  $(t-1)$  to  $(t)$ :



Here, we write  $h_t$  for the output, since these networks can be stacked into multiple layers, i.e.  $h_t$  is input into a new layer.

Chris Olah, <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

# From state $h$ to output $y$

$$\mathbf{h}^{(t)} = \psi \left[ \mathbf{U} \cdot \mathbf{h}^{(t-1)} + \mathbf{W} \cdot \mathbf{x}^{(t)} \right]$$

$$\mathbf{y}^{(t)} = \phi \left[ \mathbf{V} \cdot \mathbf{h}^{(t)} \right]$$

# Toy example with handcrafted parameters

In a sequential problem, we surveil a farm and watch for the appearance of objects. At each instant  $t$  we observe a vector  $\mathbf{z}$  which can indicate the appearance on the scene of

- a wolf 
- a farmer 

The objective is to output an estimate of danger, i.e.

- presence of the wolf w/o the farmer, or
- presence of both, arrival of the wolf before the farmer.

# Toy example with handcrafted parameters

$$\mathbf{h}^{(t)} = \psi \left[ \mathbf{U} \cdot \mathbf{h}^{(t-1)} + \begin{array}{l} \text{Wolf arrived}^* \\ \text{Farmer arrived}^* \\ \text{Time since wolf arrived} \\ \text{Time since farmer arrived} \end{array} \right]$$

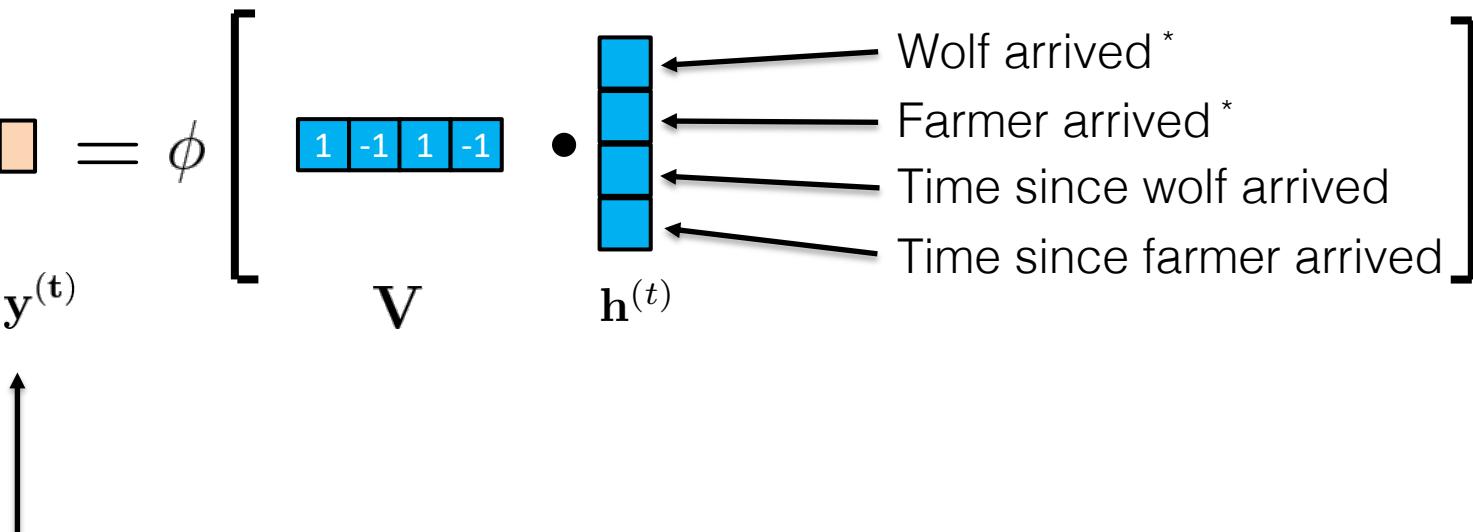
$$\begin{array}{l} \text{Appearance template "wolf"} \\ \text{Appearance template "farmer"} \end{array} \rightarrow \left[ \begin{array}{c} \mathbf{W} \cdot \mathbf{x}^{(t)} \\ + \end{array} \right]$$

\*Appropriate activation functions required to normalize state values

# Toy example with handcrafted parameters

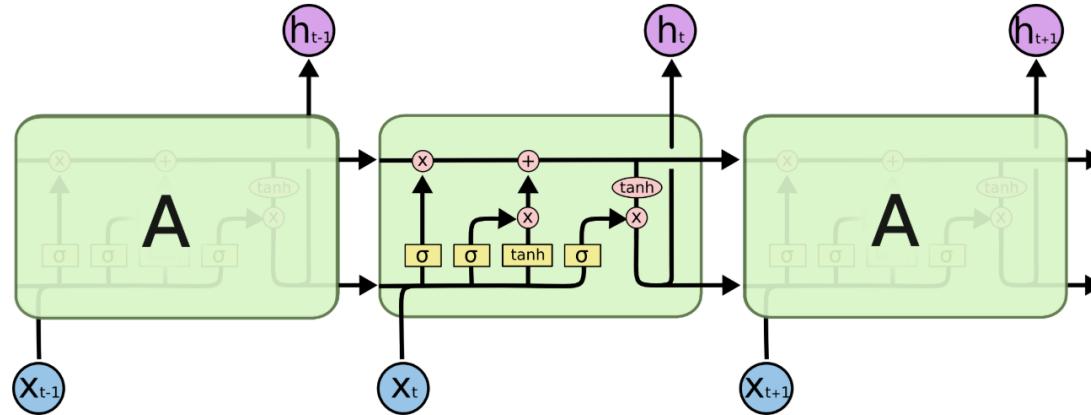
$$\boxed{y^{(t)}} = \phi \left[ \begin{matrix} V & \bullet & h^{(t)} \end{matrix} \right] \cdot \begin{matrix} \text{Wolf arrived}^* \\ \text{Farmer arrived}^* \\ \text{Time since wolf arrived} \\ \text{Time since farmer arrived} \end{matrix}$$

Output: danger!



# LSTM Networks

LSTM (=Long-short term Memory) networks use gating mechanisms, which handle information flow in a fully trainable way.



$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

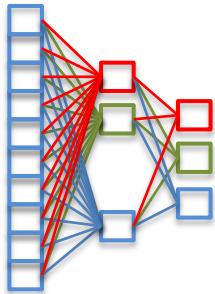
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

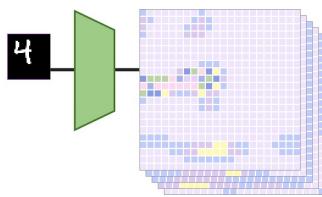
$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

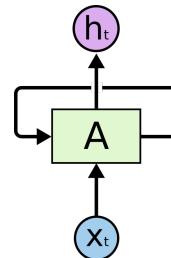
# The Deep Toolbox



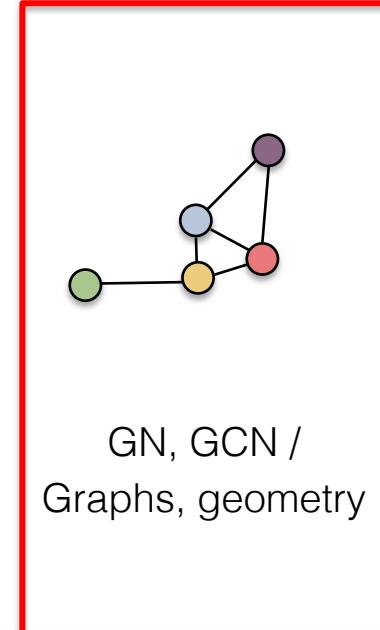
MLP



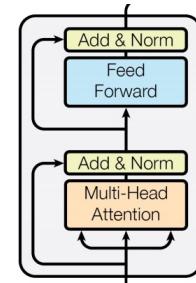
CNN /  
Convolutions



RNN /  
Recurrence



GN, GCN /  
Graphs, geometry



Transformers /  
Self-attention

*What do I know about the data and the task?*

*Nothing  
(vector space)*

*Translation  
equivariance*

*Sequential data,  
Markov property*

*Graph structured  
data*

*Permutation  
equivariance*

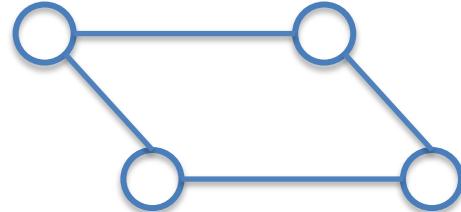
# Structured Input and/or structured Output

Predicting for multiple inter-dependent variables

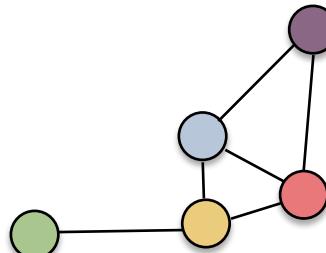
- Sequences
- Images and other 2D grids
- Multi-label Problems
- General graphs
- Kinematic trees



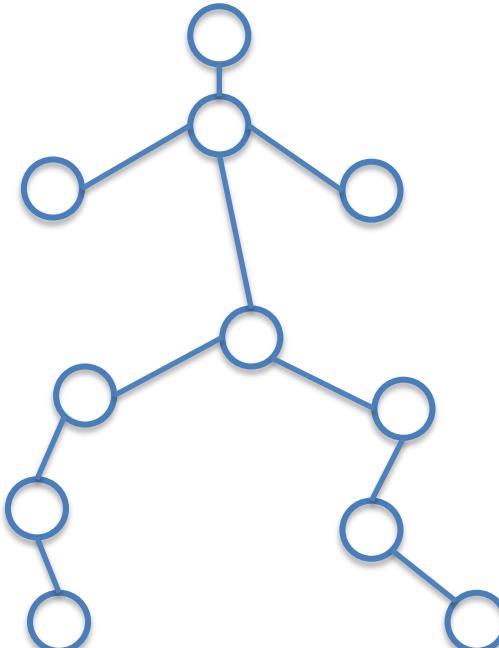
Sequences



Images and  
other 2D grids



General graphs

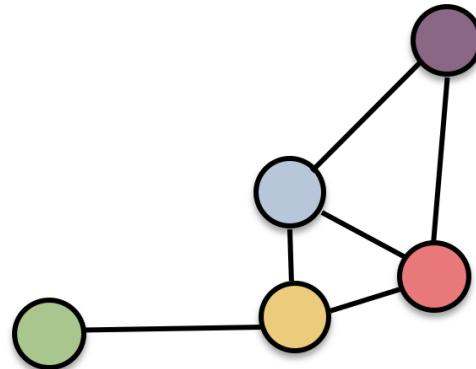


Kinematic trees

# Graphs: definition

A graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  consists of:

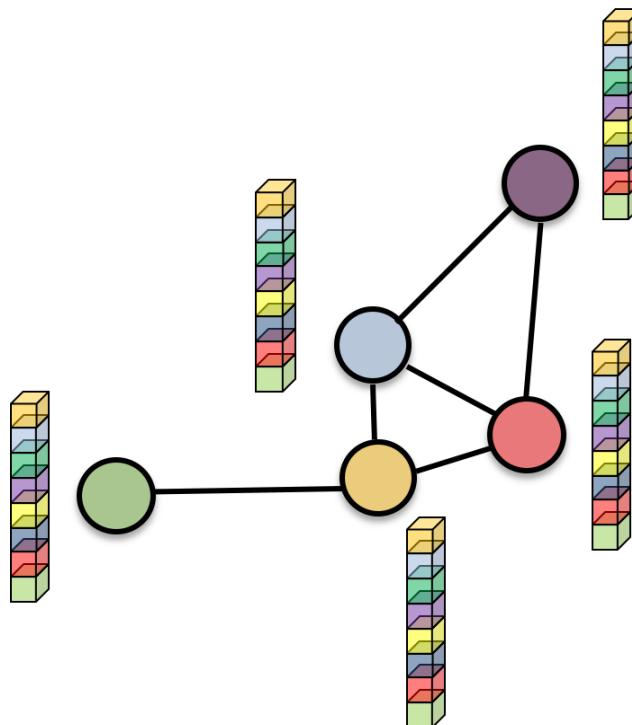
- a set  $\mathcal{V}$  of nodes, and
- a set  $\mathcal{E} \in \mathcal{V} \times \mathcal{V}$  of edges



# Attributed Graphs

Attributed graphs also have values for each node:  $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{X})$ :

- a set  $\mathcal{V}$  of nodes, and
- a set  $\mathcal{E} \in \mathcal{V} \times \mathcal{V}$  of edges
- a set  $\mathcal{X} = \{x_0, \dots, x_N\}$  of values, each value being associated to a node. In our case we find **embeddings**  $x_i \in \mathbb{R}^d$  for each node.



# Relational Reasoning

$$g(x_1, x_2, \dots, x_N) = \max(h(x_1), h(x_2), \dots, h(x_N))$$

“PointNet” [Qi, Su, Mo, Guibas, CVPR 2017]

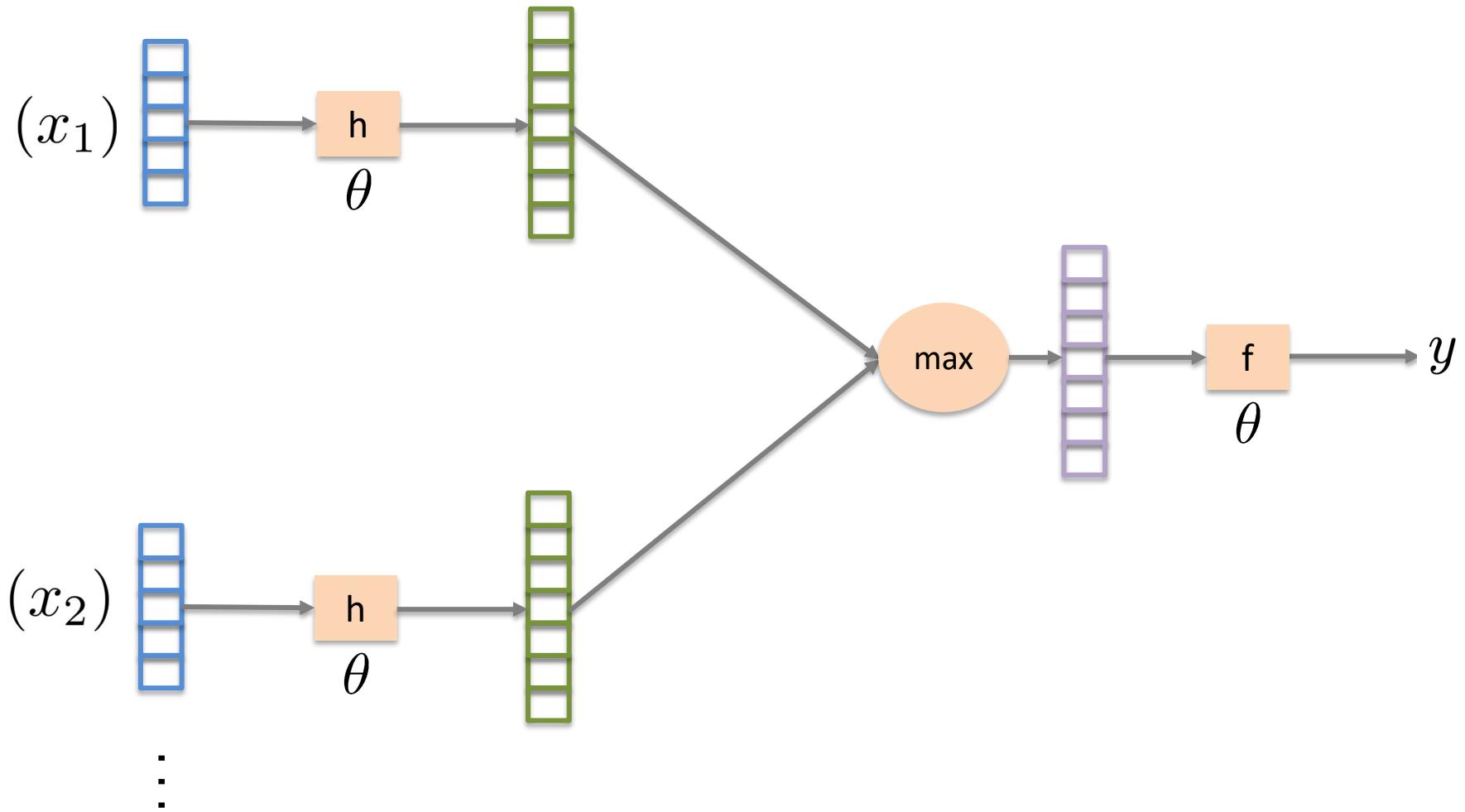
Defined over points of a point cloud

$$g(x_1, x_2, \dots, x_N) = \sum_{i,j} h(x_i, x_j)$$

“Relational Reasoning” [Santoro et al., NIPS 2017]

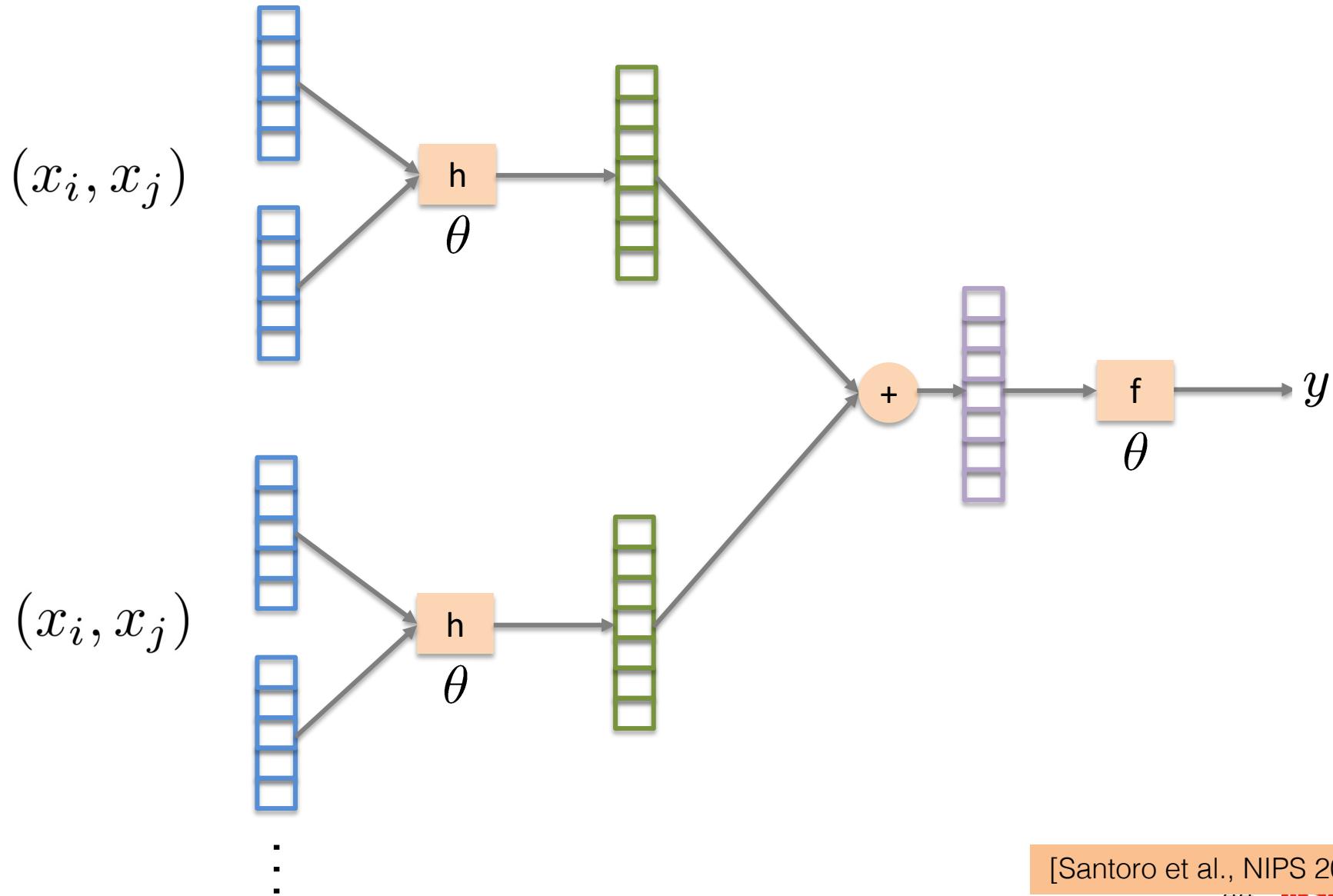
Defined over feature map cells

# PointNet



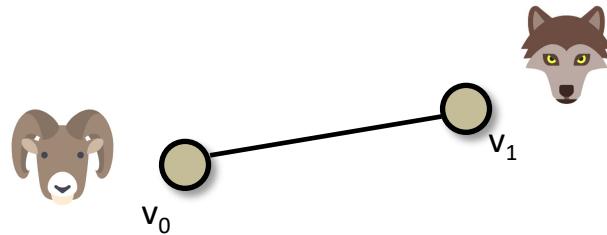
[Qi, Su, Mo, Guibas, CVPR 2017]

# Relational Reasoning: pairwise terms

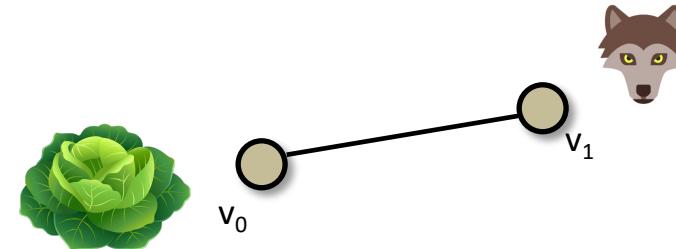


[Santoro et al., NIPS 2017]

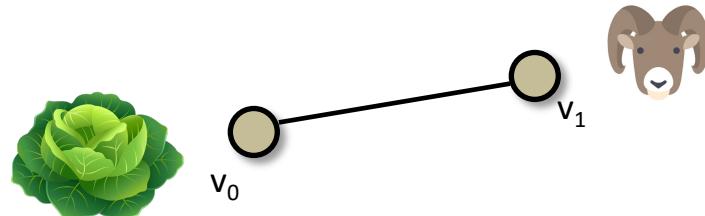
# A toy problem: will somebody get eaten?



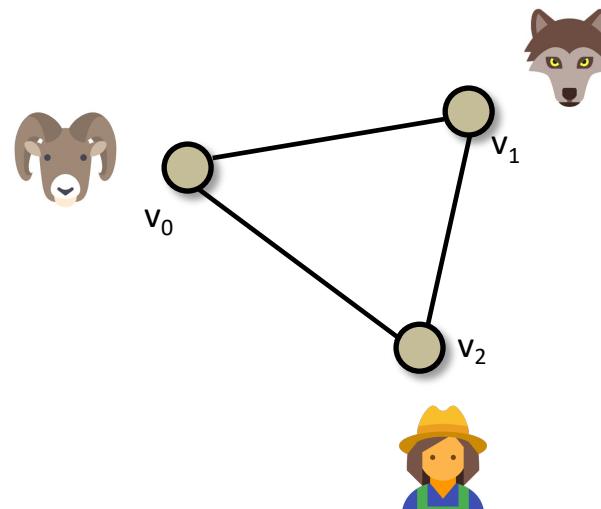
Yes



No

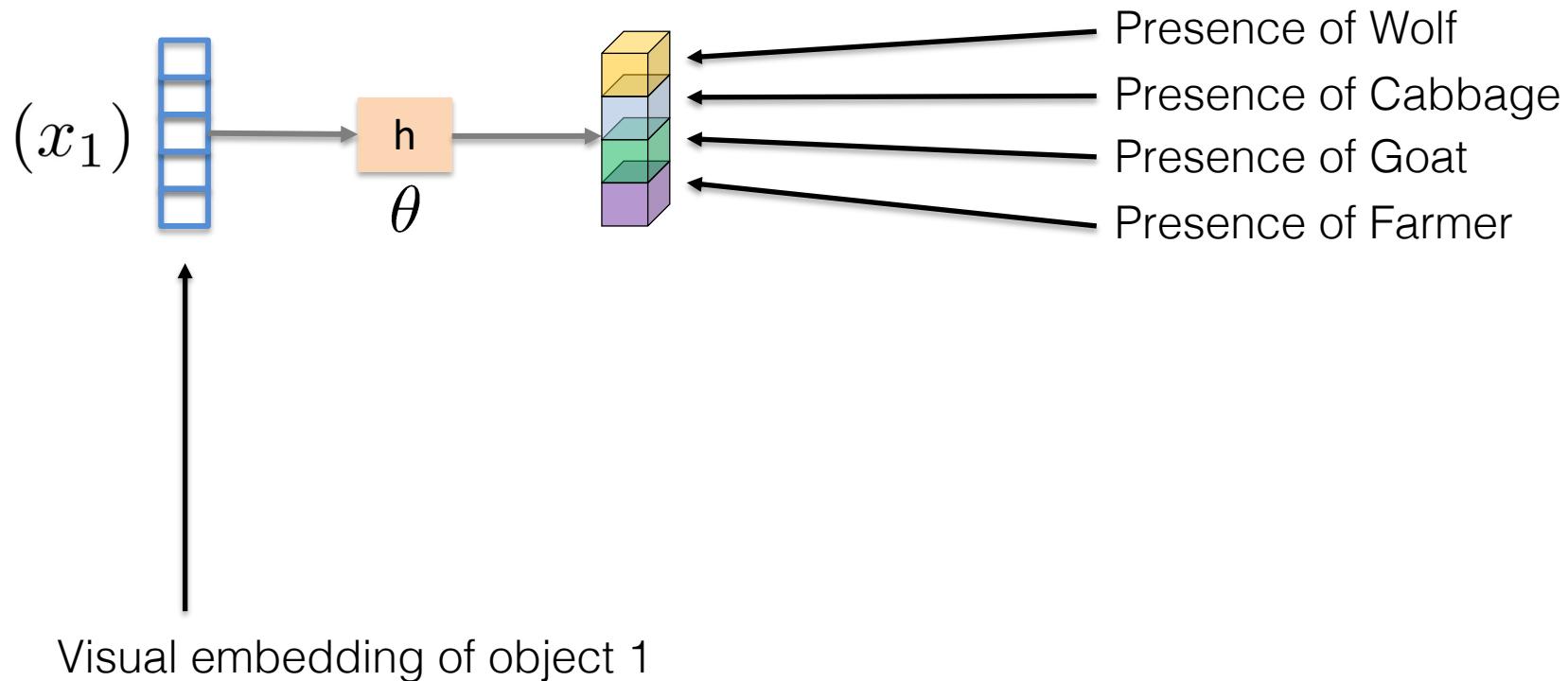


Yes

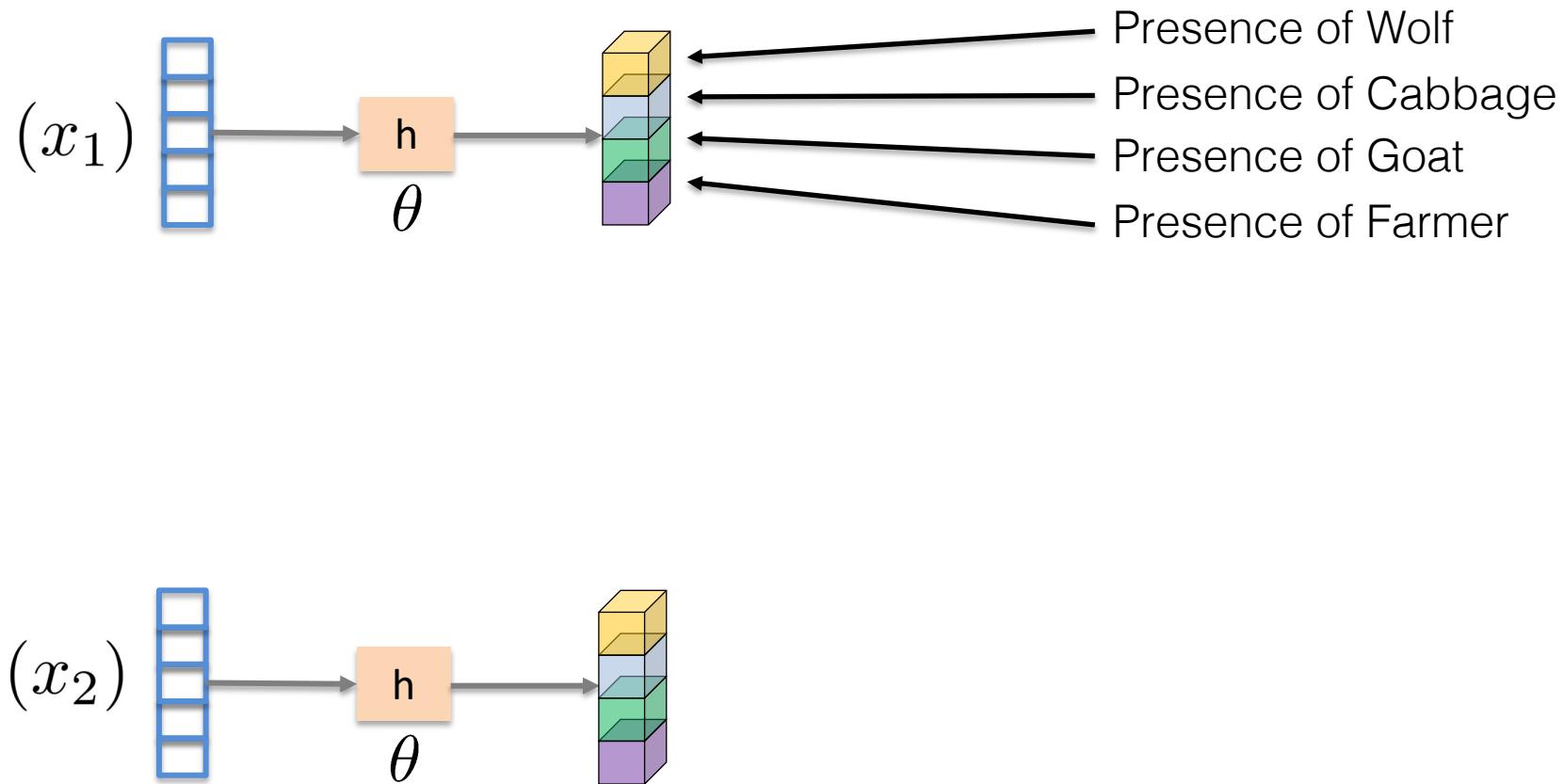


No

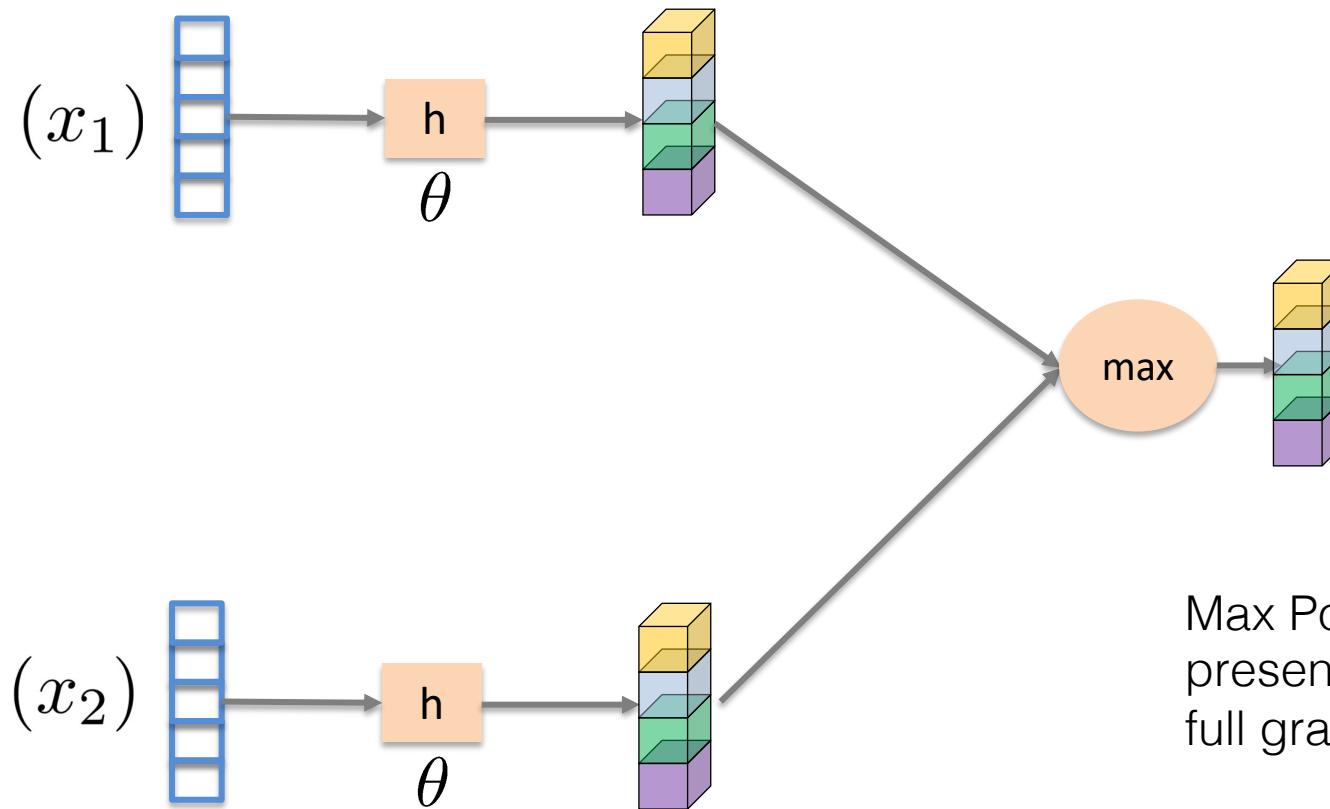
# Representation by PointNet



# Representation by PointNet

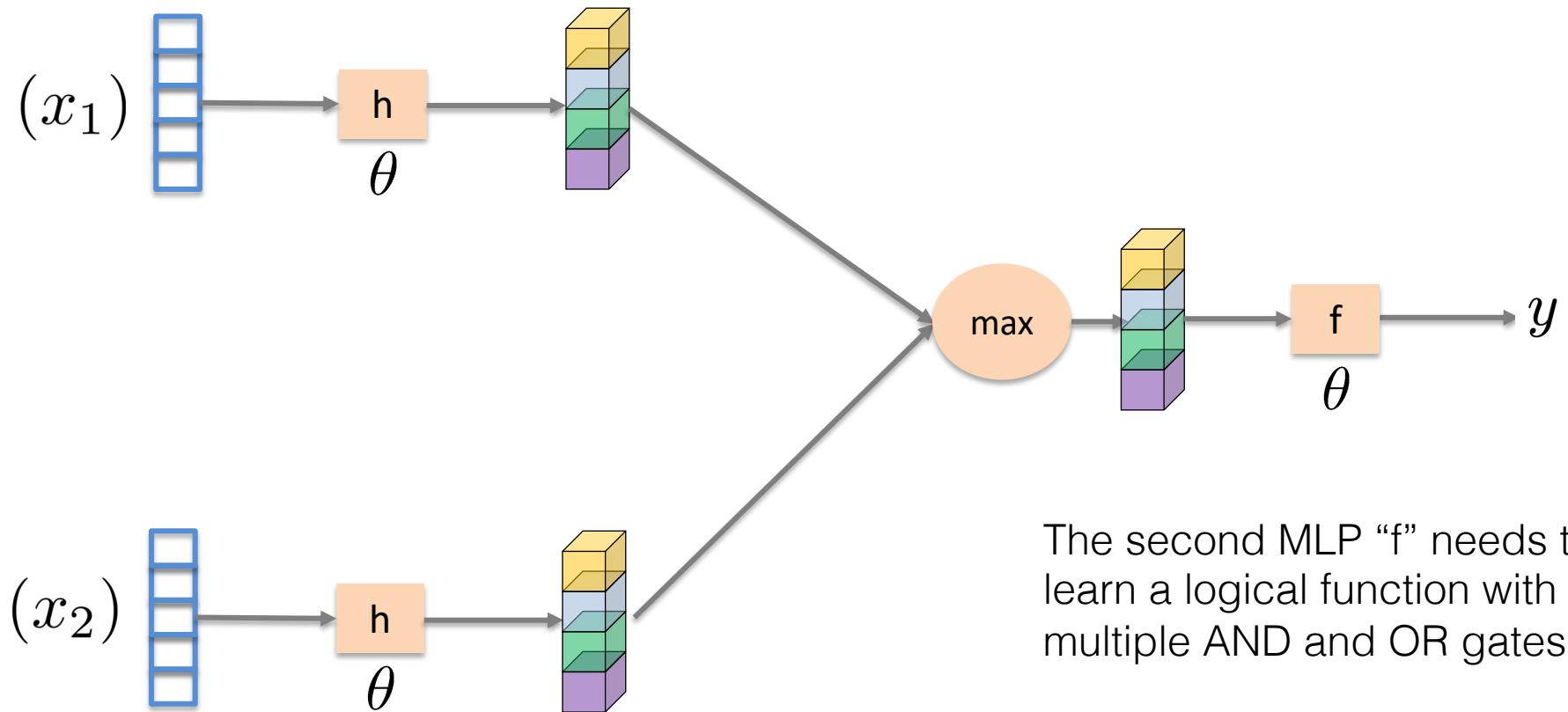


# Representation by PointNet

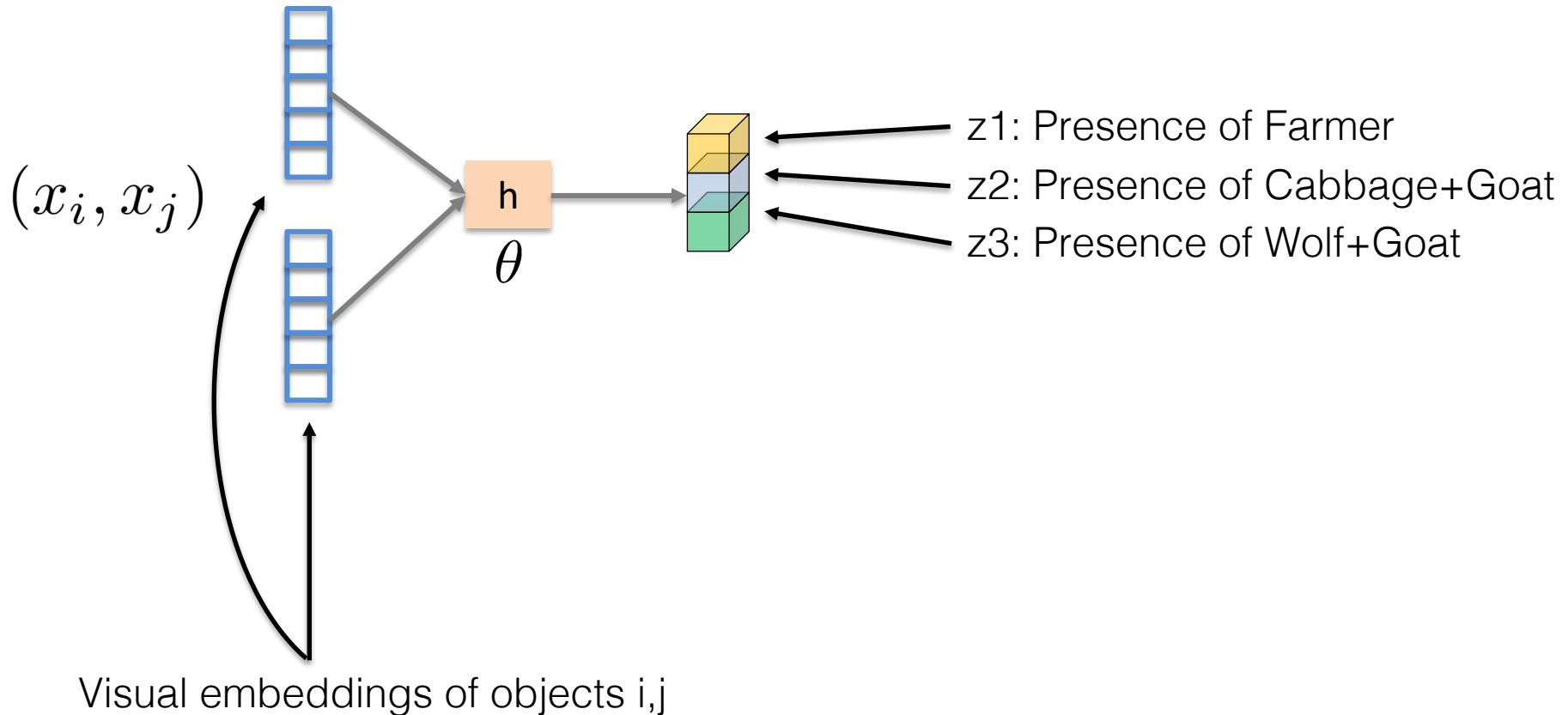


Max Pooling will cumulate the presence information from the full graph into a single vector

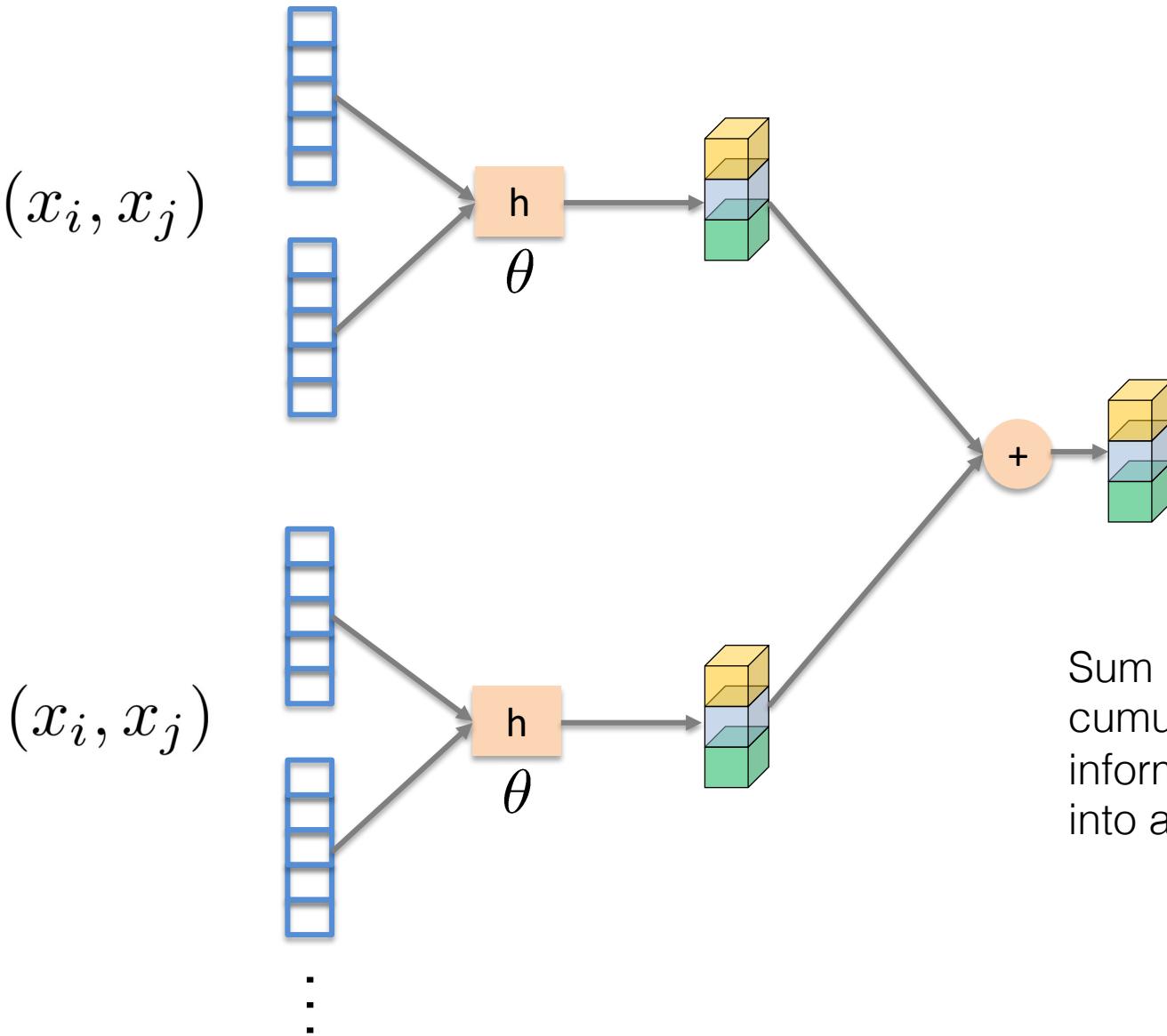
# Representation by PointNet



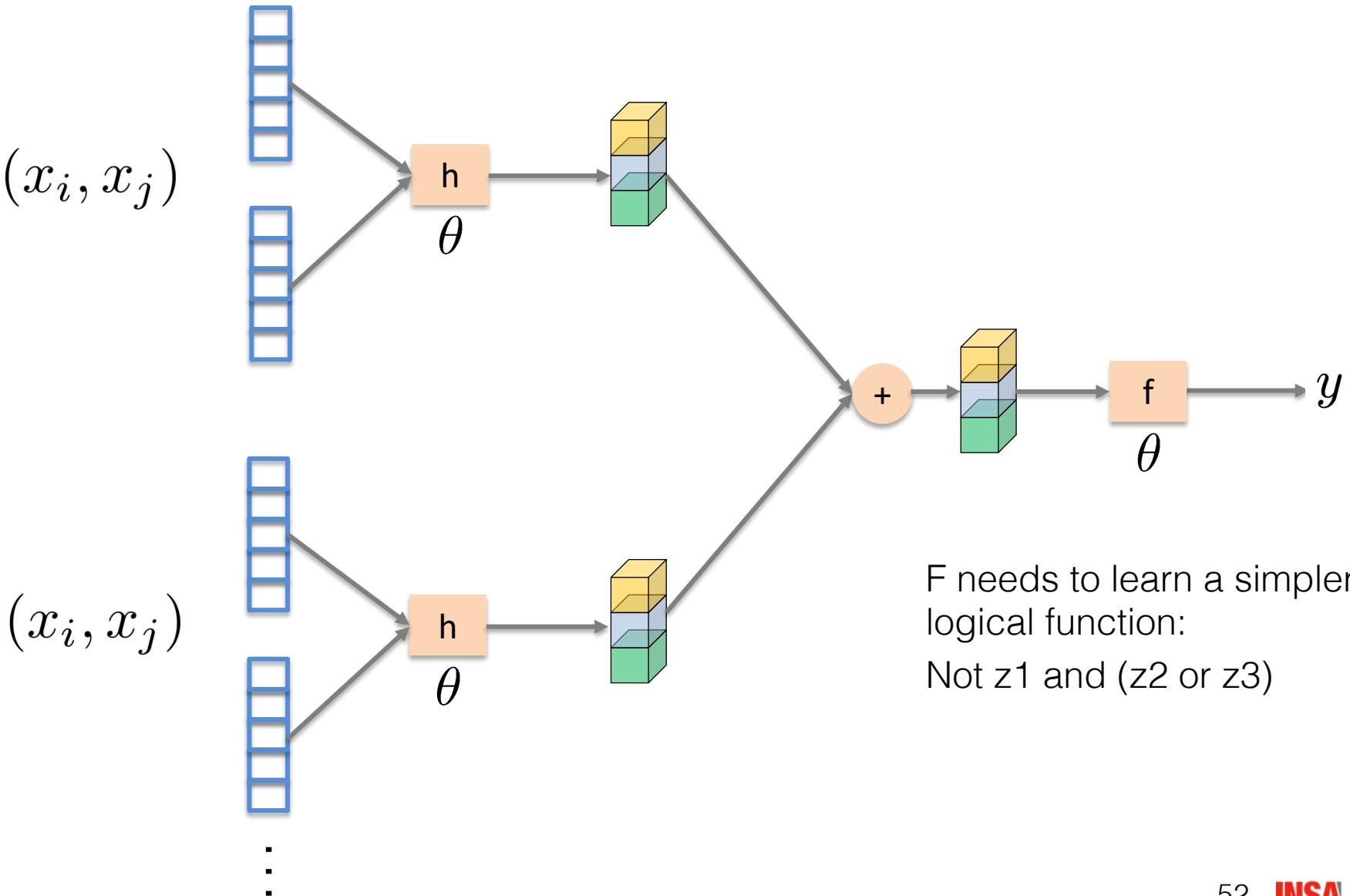
# Representation with pairwise terms



# Representation with pairwise terms

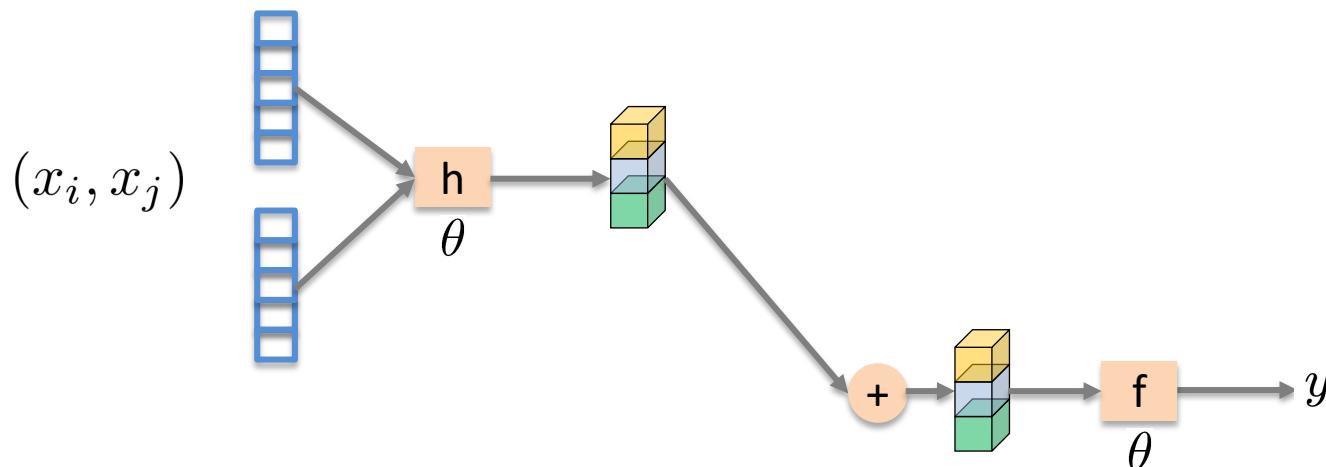


# Representation with pairwise terms



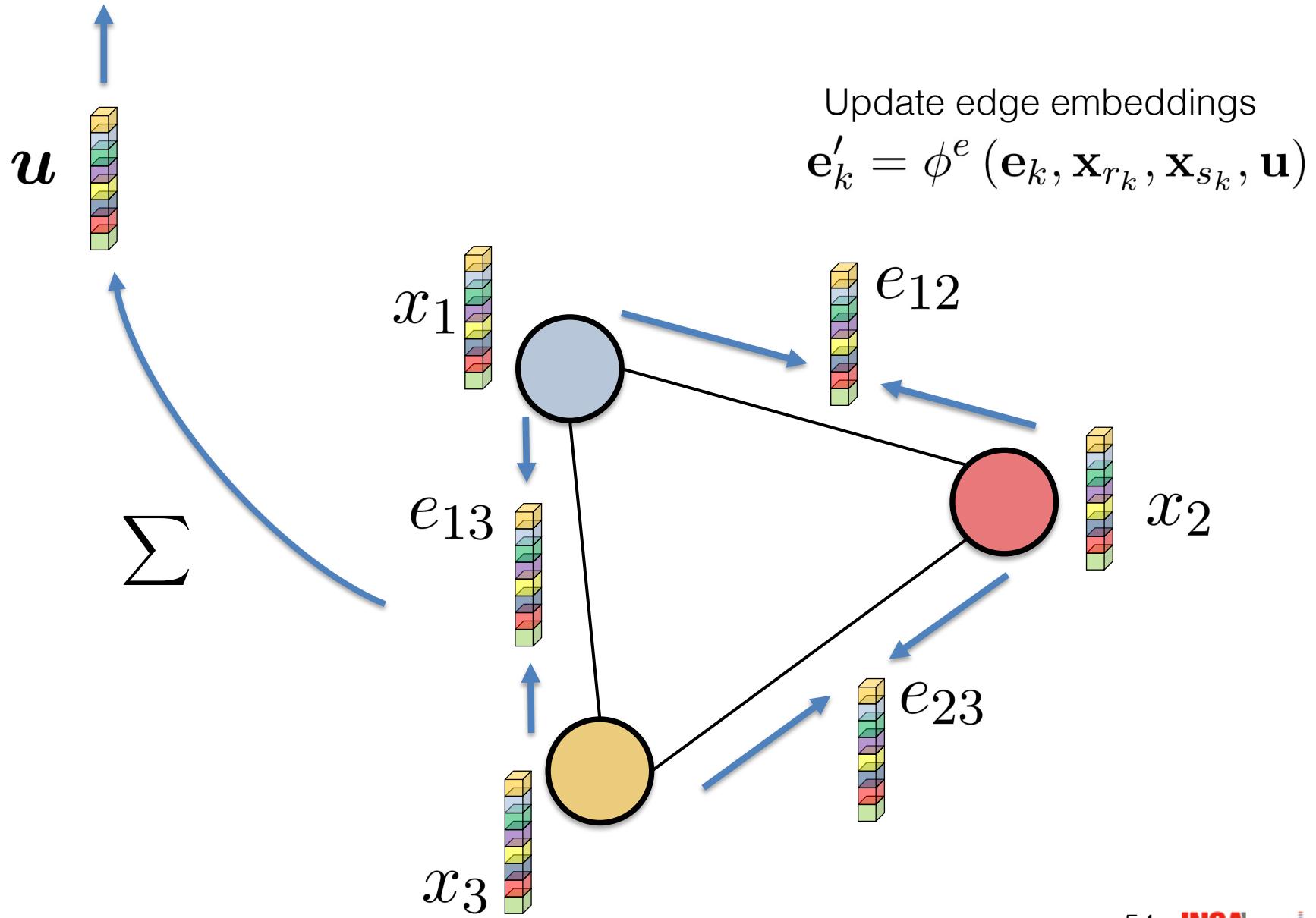
# Comparison

- Increasing the complexity of  $h$  may allow to decrease the complexity of  $f$ .
- There is no known rule which determines the best trade-off between  $h$  and  $f$  for a given problem.
- Example: there are problems dominated by pairwise relationships in the data where models without pairwise terms work better.



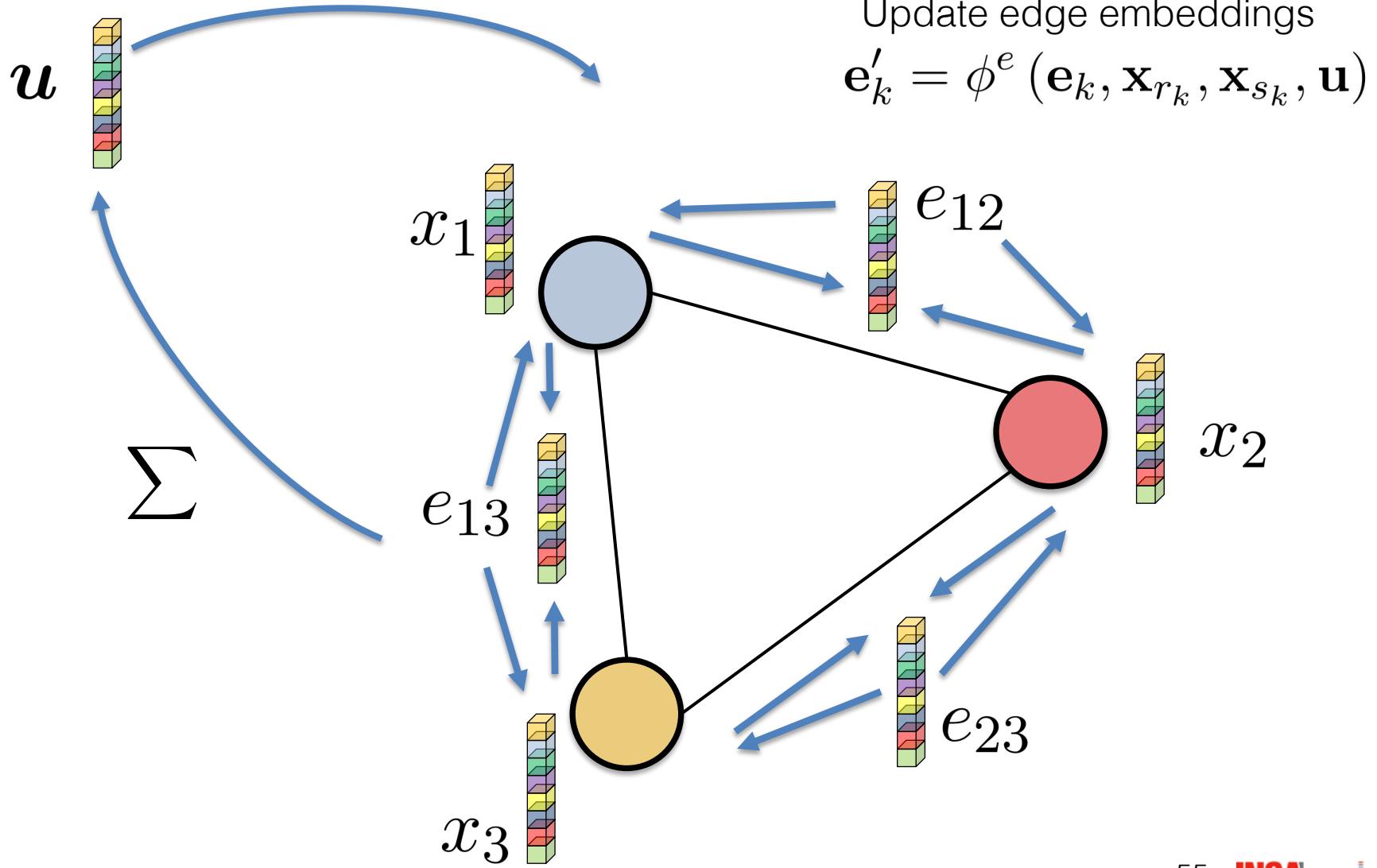
Decision

# Relational reasoning

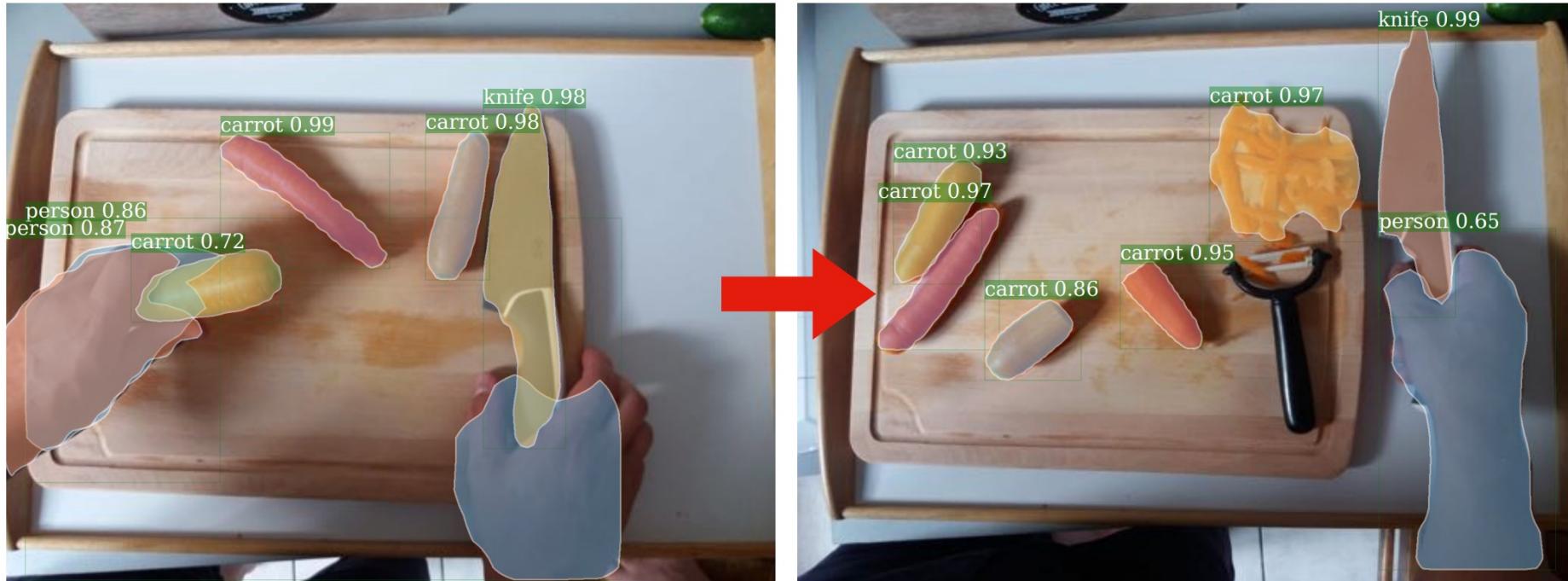


# Graph Networks

All messages are influenced by the global (context) embedding



# Object level Visual Reasoning



[Baradel, Neverova, Wolf, Mille, Mori, ECCV 2018]



Fabien Baradel  
Phd @ LIRIS,  
INSA-Lyon



Natalia Neverova  
Facebook AI  
Research, Paris



Christian Wolf  
Insa-Lyon,  
LIRIS  
INRIA Chroma

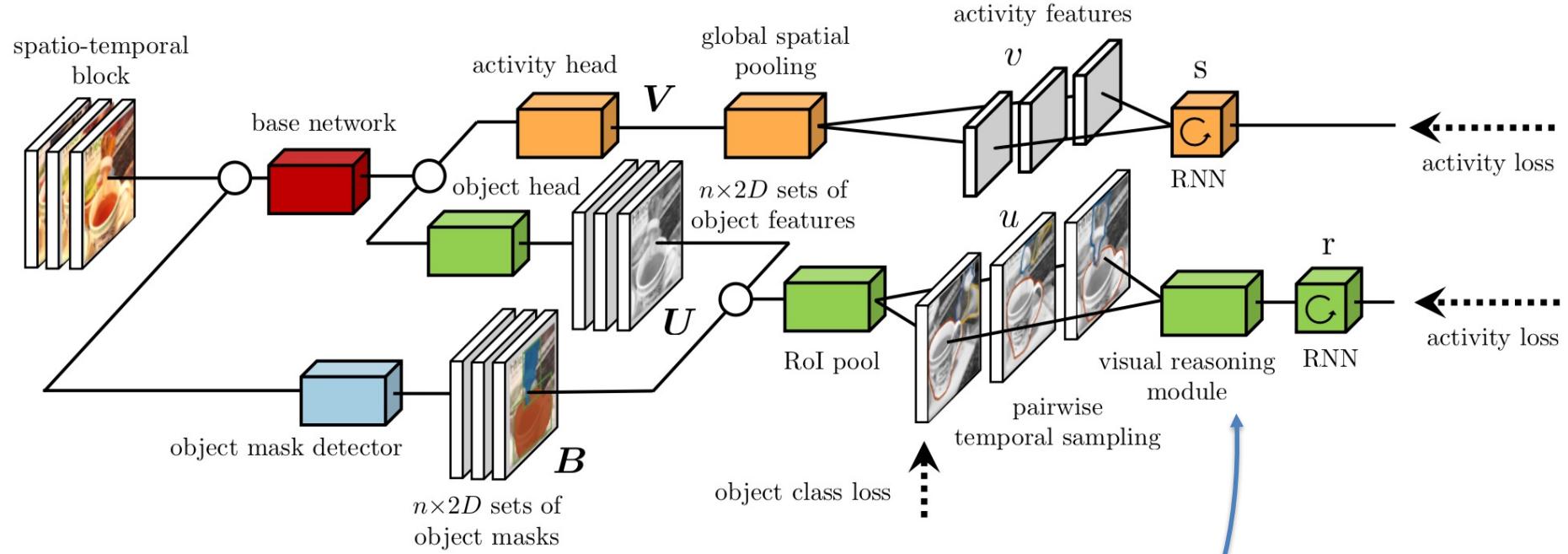


Julien Mille  
LIFAT,  
INSA VdL



Greg Mori  
Simon  
Fraser  
University,  
Canada

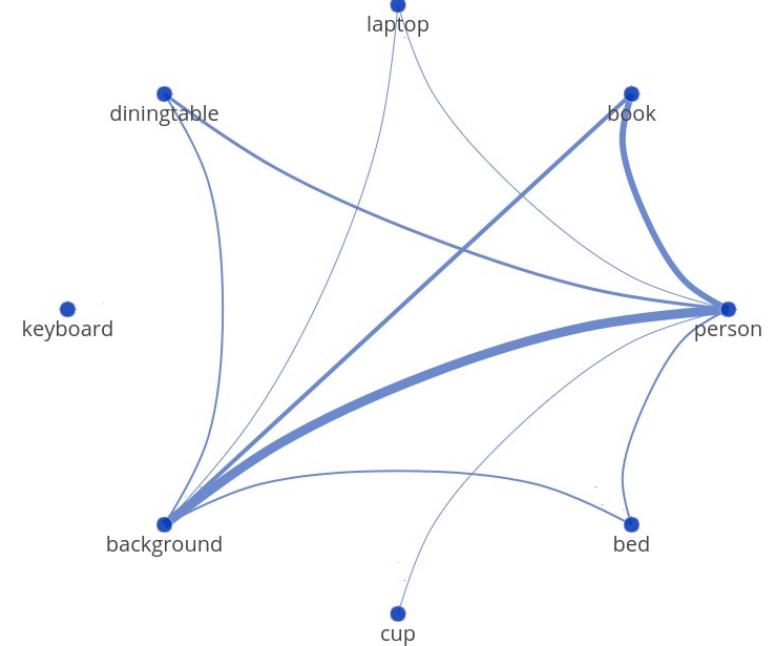
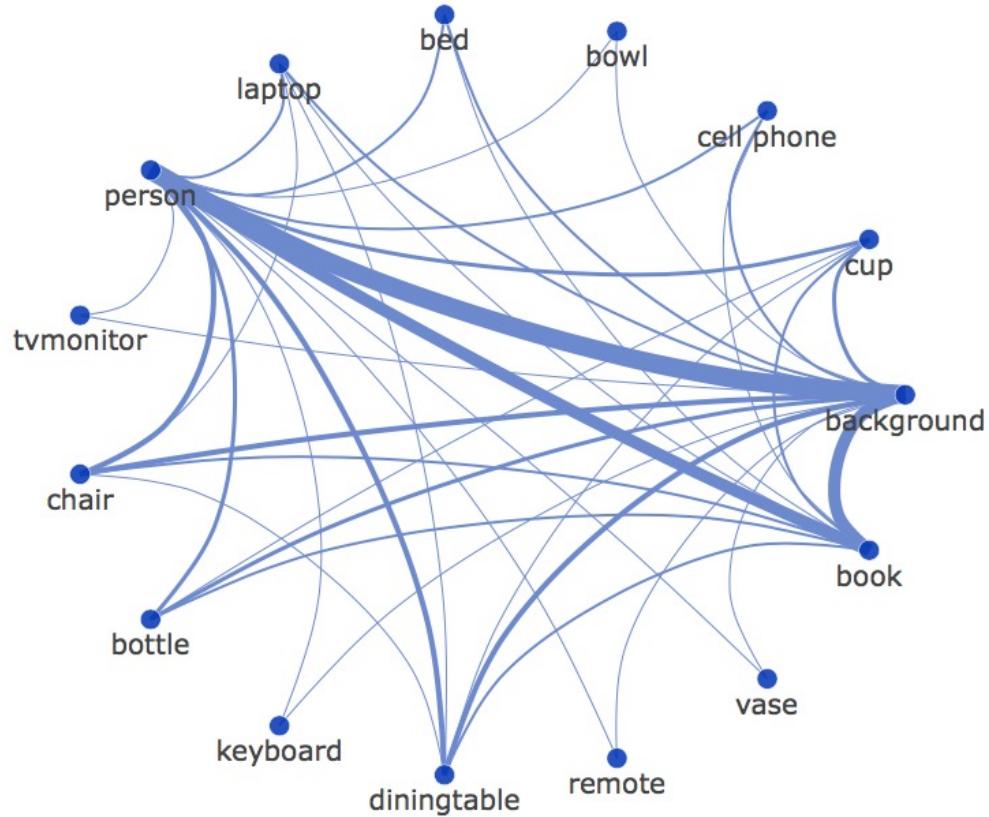
# Object level Visual Reasoning



$$\mathbf{g}_t = \sum_{j,k} h_\theta(\mathbf{o}_{t'}^j, \mathbf{o}_t^k)$$

[Baradel, Neverova, Wolf, Mille, Mori, ECCV 2018]

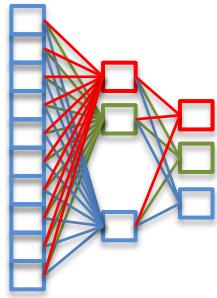
# Learned interactions



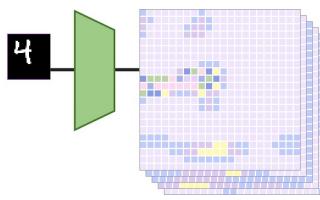
## Class: person-book interaction

[Baradel, Neverova, Wolf, Mille, Mori, ECCV 2018]

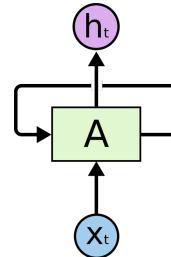
# The Deep Toolbox



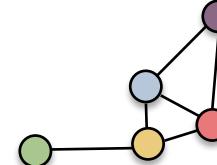
MLP



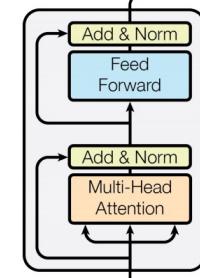
CNN /  
Convolutions



RNN /  
Recurrence



GN, GCN /  
Graphs, geometry



Transformers /  
Self-attention

*What do I know about the data and the task?*

*Nothing  
(vector space)*

*Translation  
equivariance*

*Sequential data,  
Markov property*

*Graph structured  
data*

*Permutation  
equivariance*

# Learning language reasoning

Все счастливые семьи похожи друг на друга, каждая несчастливая семья несчастлива по-своему..

Toutes les familles heureuses se ressemblent. Chaque famille malheureuse, au contraire, l'est à sa façon.

Happy families are all alike. Every unhappy family is unhappy in its own way.

Alle glücklichen Familien gleichen einander, jede unglückliche Familie ist auf ihre eigene Weise unglücklich

[L. Tolstoy, 1873]

# Examples

If you really want to hear about it, the first thing you'll probably want to know VO  
where I was born, and what my lousy childhood was like, and how my parents were  
occupied and all before they had me, and all that David Copperfield kind of crap, but I  
don't feel like going into it, if you want to know the truth. In the first place, that stuff  
bores me, and in the second place, my parents would have two hemorrhages apiece if  
I told anything pretty personal about them.

Si vous voulez vraiment en entendre parler, la première chose que vous voudrez probablement savoir est où je suis né, et ce que mon enfance moche était, et comment mes parents étaient occupés et tout ce qu'ils avaient avant moi, et tout ce que David Copperfield, mais je n'ai pas envie d'y aller, si tu veux savoir la vérité. En premier lieu, cela m'ennuie, et en second lieu, mes parents auraient deux hémorragies si je leur racontais quelque chose de très personnel. VF Google

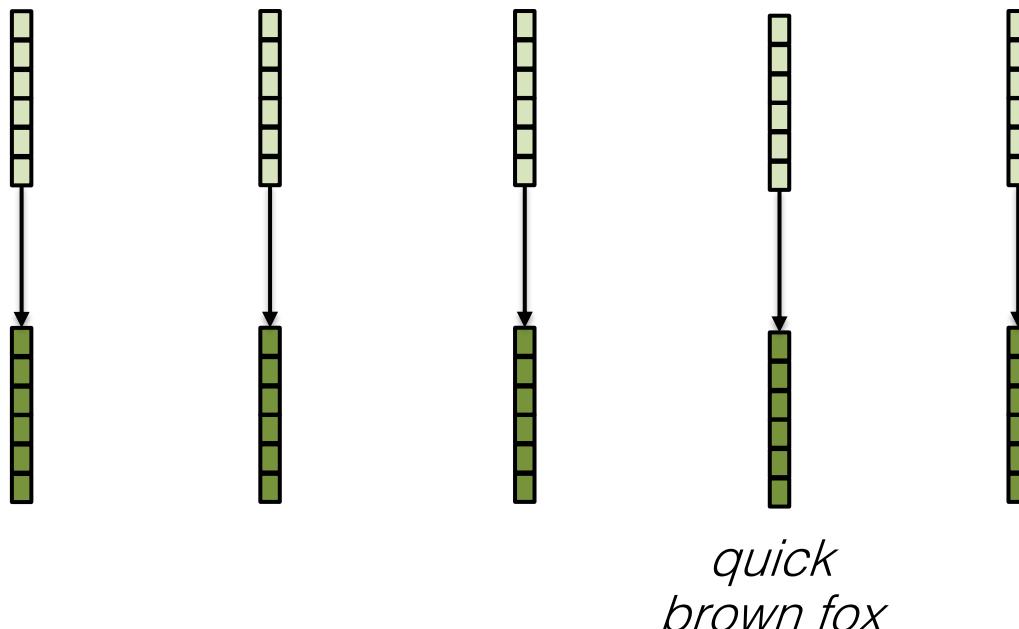
Si vous voulez vraiment que je vous dise, alors sûrement la première chose que je vais aller demander c'est où je suis né, et à quoi ça ressemblait, ma saloperie d'enfance et c'est que faisaient mes parents avant de m'avoir, et toutes ces conneries à la David Copperfield, mais j'ai pas envie de raconter ça et tout. Primo, ce genre de trucs ça me rase, et secundo, mes parents ils auraient chacun une attaque ou même deux chacun, si je me mettais à baratiner sur leur compte quelque chose d'un peu personnel. VF officielle

# Contextualization

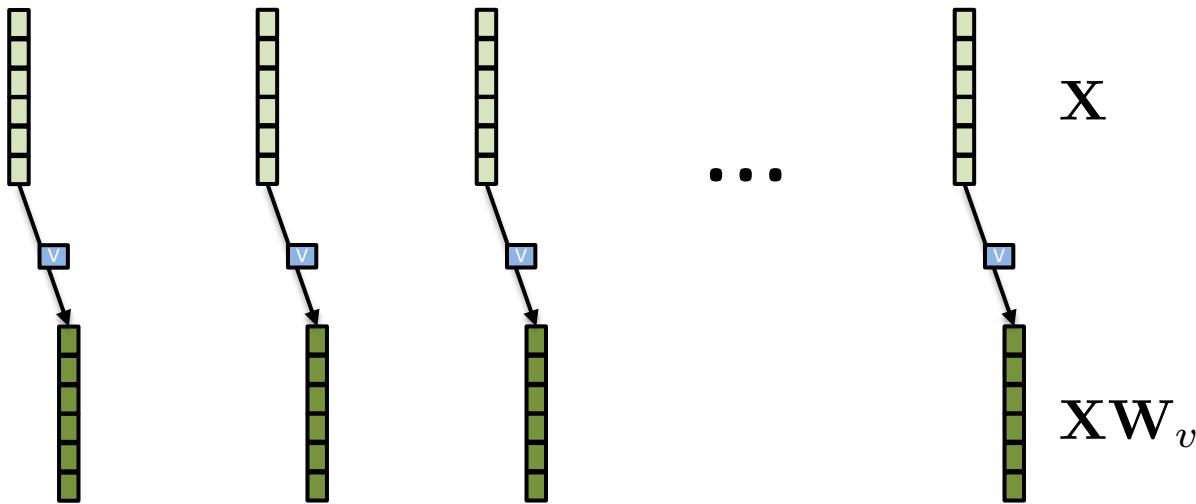
We suppose a set  $\mathbf{X} = \{\mathbf{x}_i\}$  of items (vectors).

Iteratively "enrich" each item by providing context from the other items

The      quick      brown      fox      jumps    ....

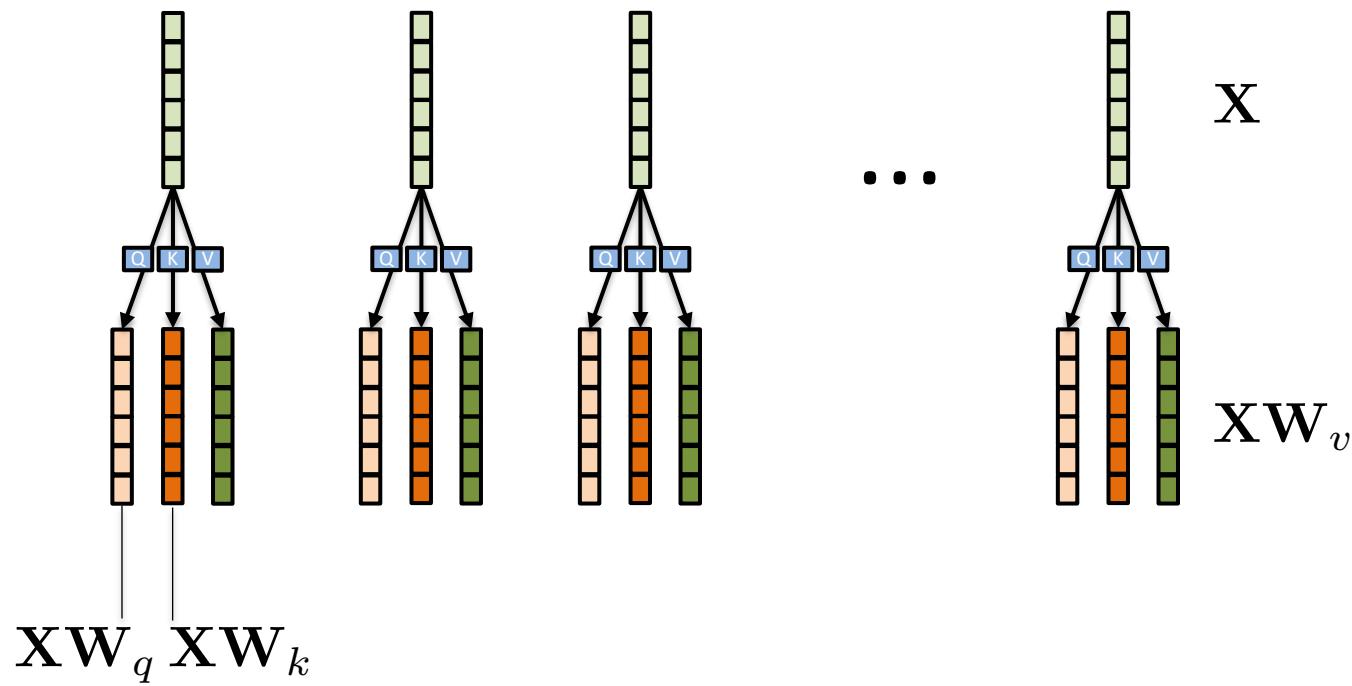


# Transformers: attention is all you need

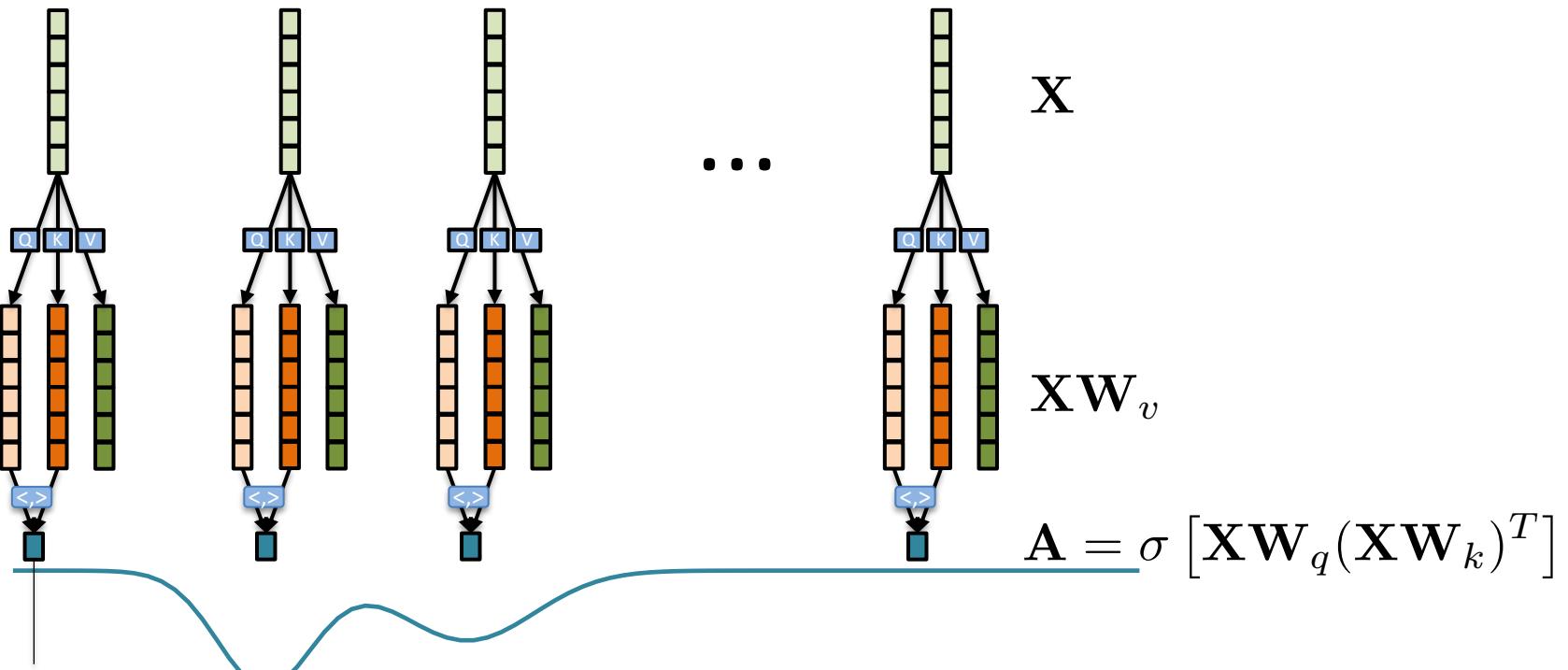


We suppose unordered data (a set)  $\mathbf{X} = \{\mathbf{x}_i\}$ . Each input item is "transformed" by a linear transform weighted by attention.

# Transformers: attention is all you need

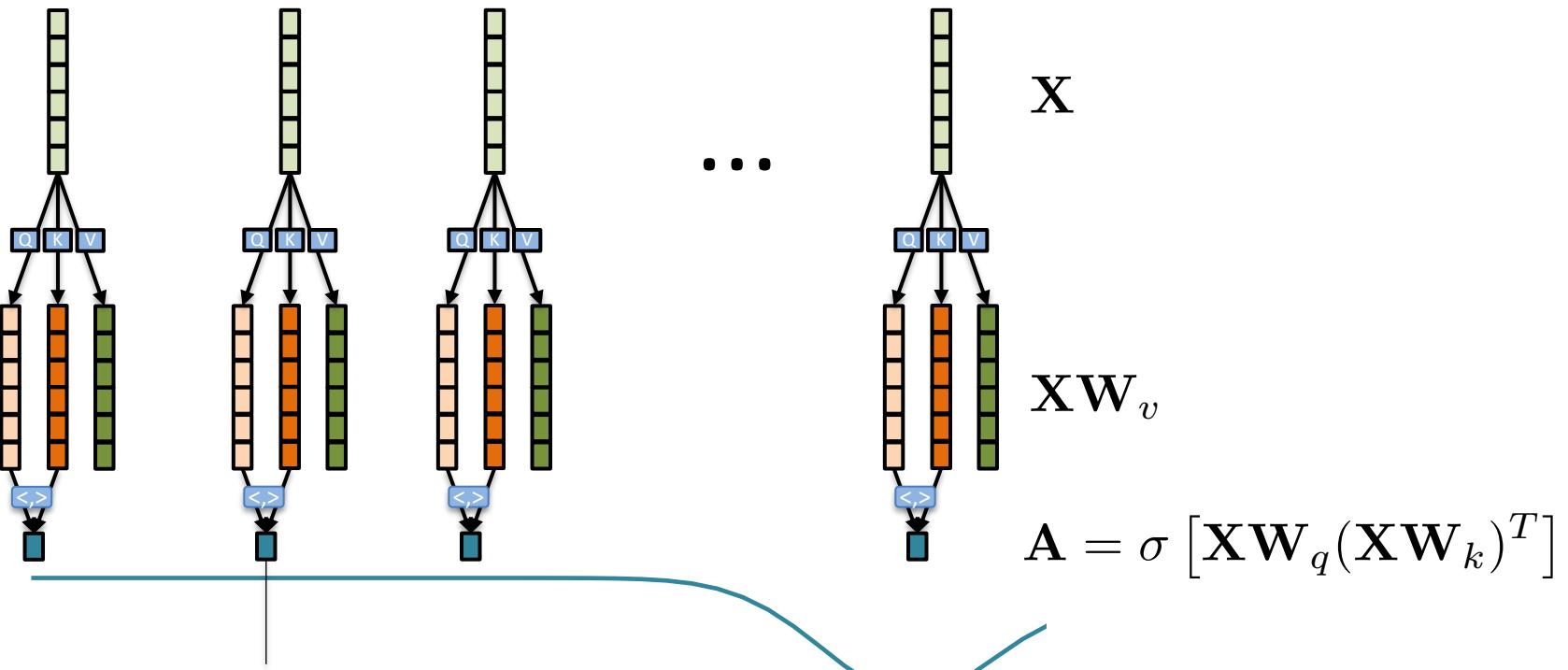


# Transformers: attention is all you need



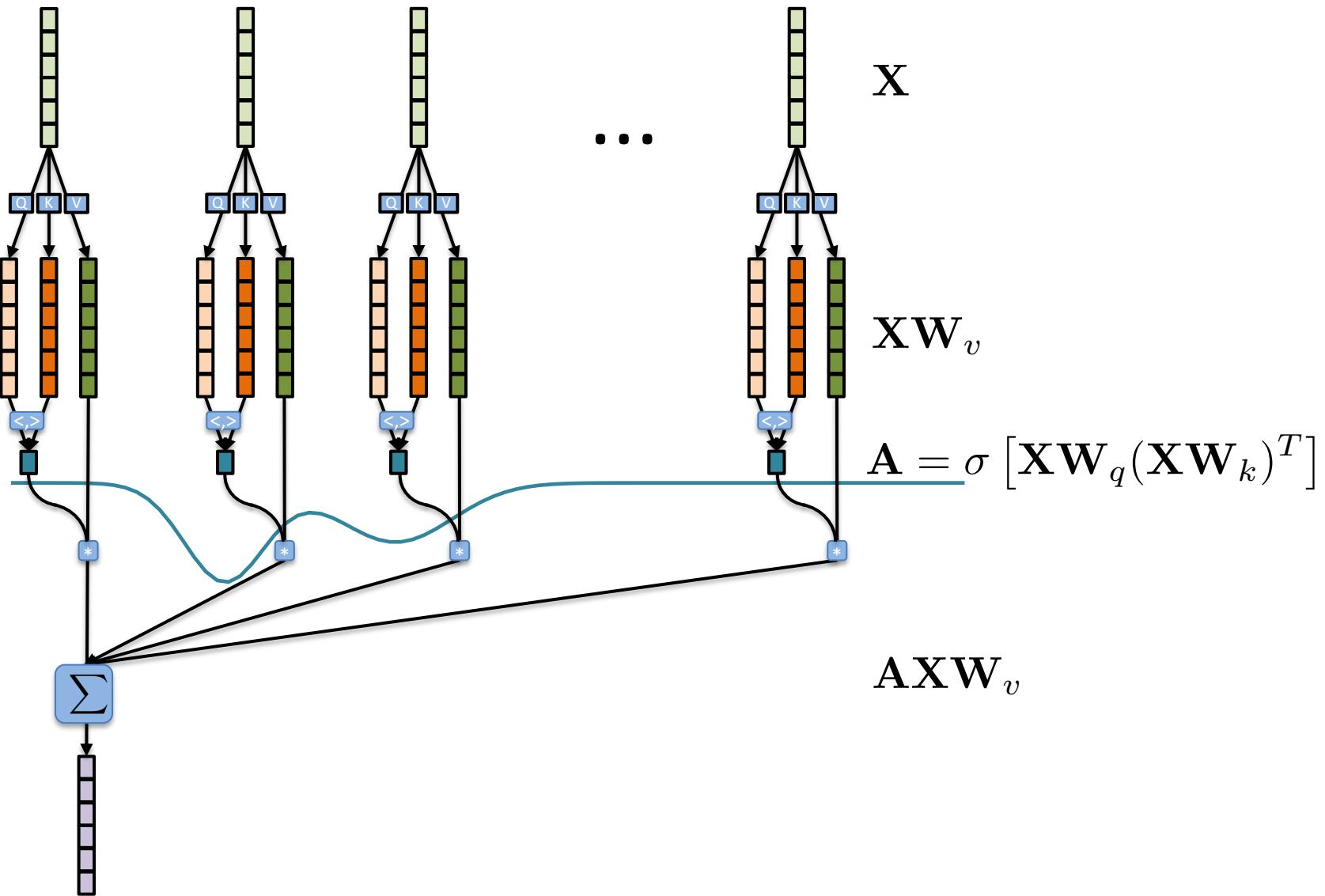
Each item has its own distribution,  
associated with its query vector!

# Transformers: attention is all you need

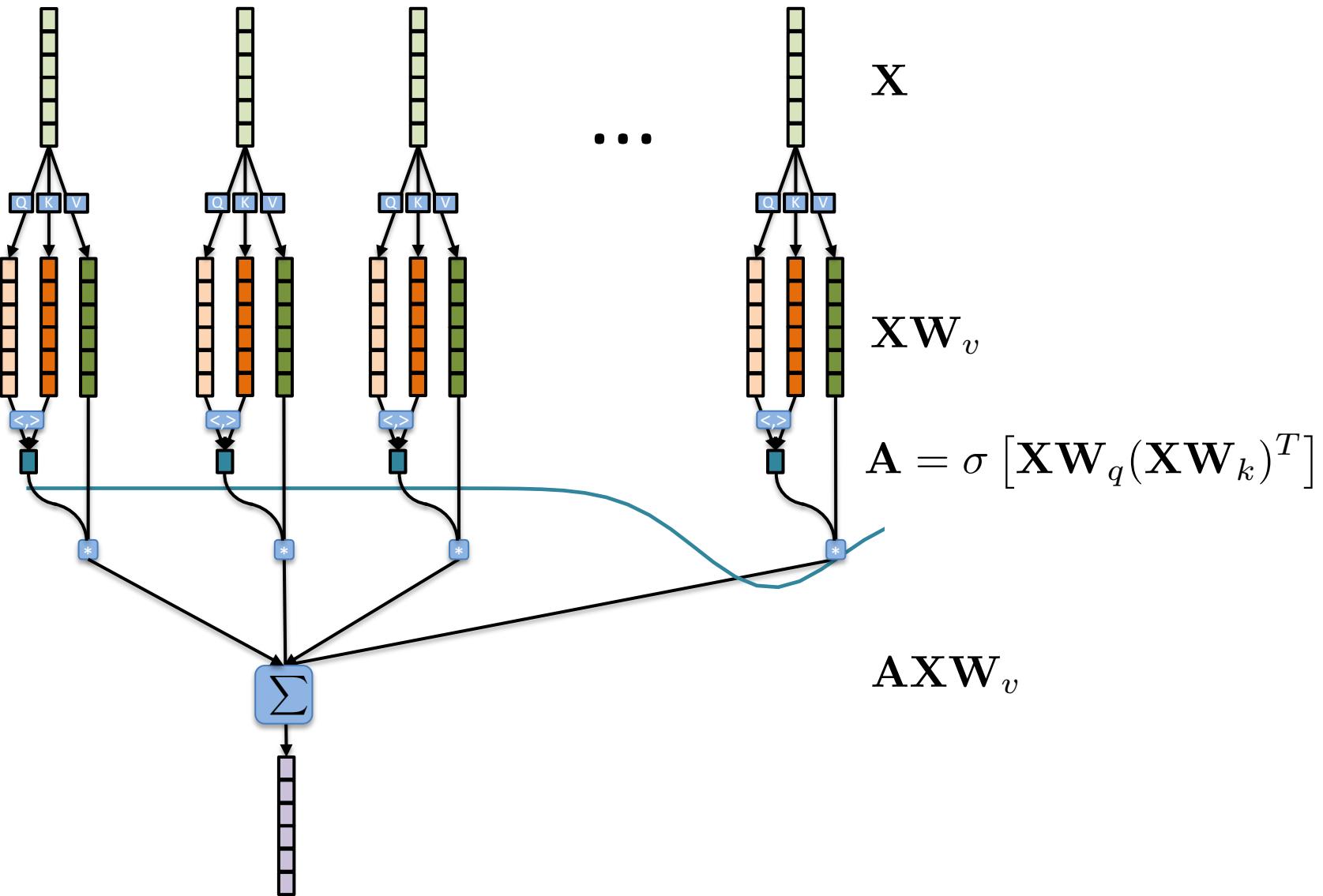


Each item has its own distribution,  
associated with its query vector!

# Transformers: attention is all you need



# Transformers: attention is all you need

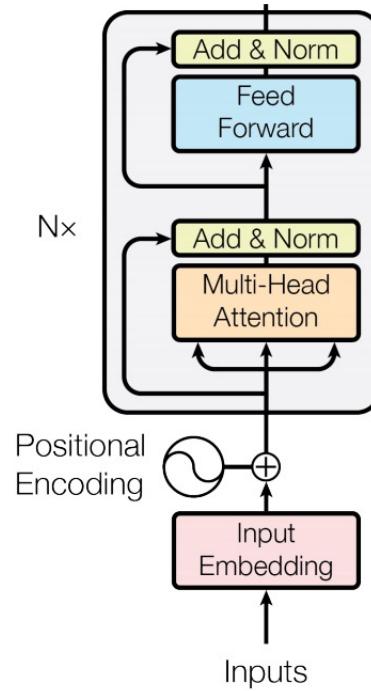


# Integration ...

Toutes les familles heureuses se ressemblent. Chaque famille malheureuse, au contraire, l'est à sa façon.



Happy families are all alike. Every unhappy family is unhappy in its own way.



# A concrete large-scale application

# Vision and Language Reasoning

*"How much money do I have in my hand?"*



*"What is in this jar?"*



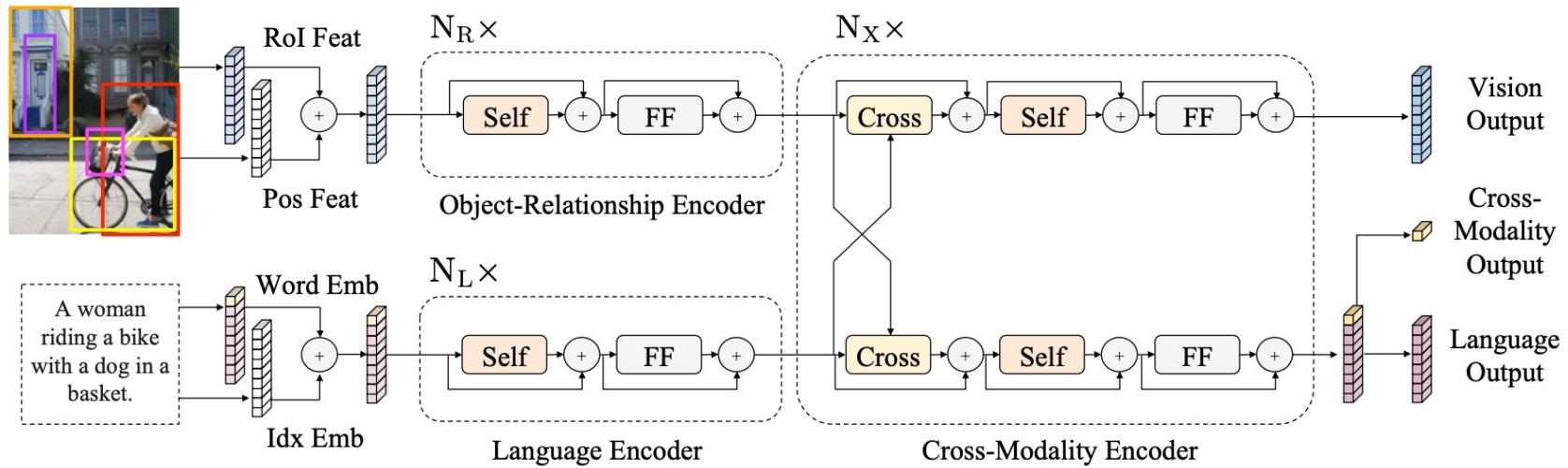
*"Did I leave the door open?"*



*"Did I leave the lights on?"*

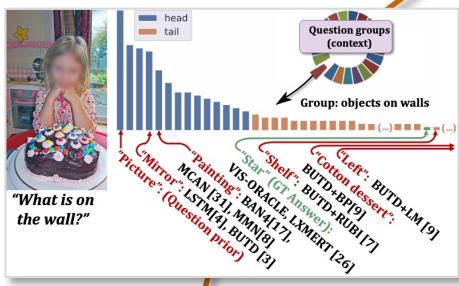
# LXMERT

A vision and language encoder with self-attention and cross-attention.

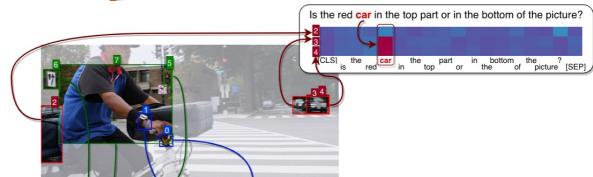


[Tan et Bansal, EMNLP 2019]

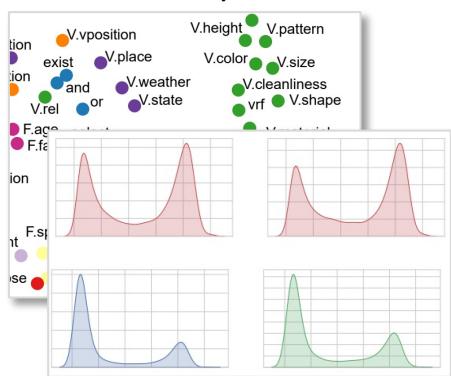
How can we evaluate biases  
in learning? (CVPR 2021a)



Can we ground object  
detection through language  
(in writing)



How can we visualize and  
transfer reasoning? (CVPR  
2021b)



VQA

Can we weakly supervise word-object alignment? (ECAI 2020)



Can we supervise reasoning  
programs? (arxiv, under  
review)



Corentin  
Kervadec



Grigory  
Antipov



Moez  
Baccouche

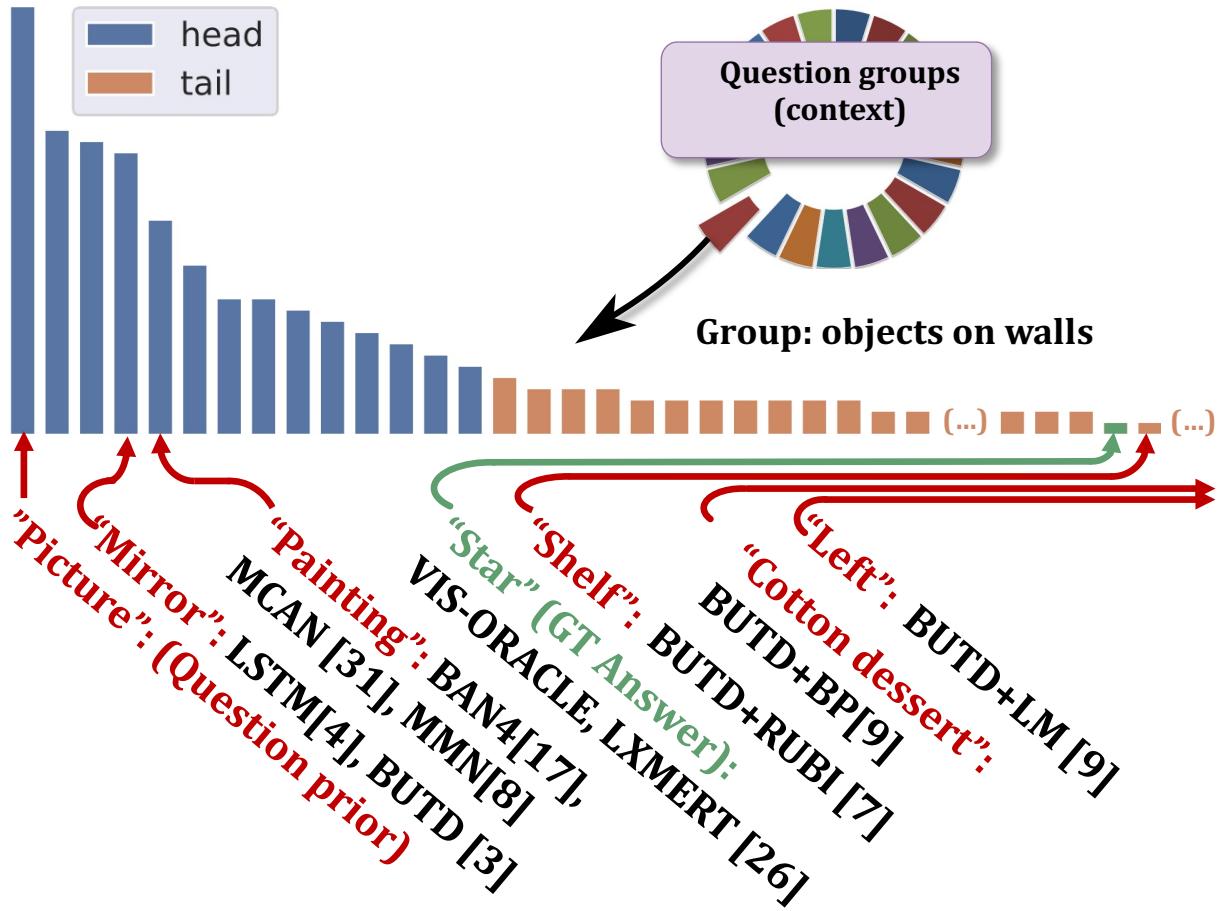


Christian  
Wolf

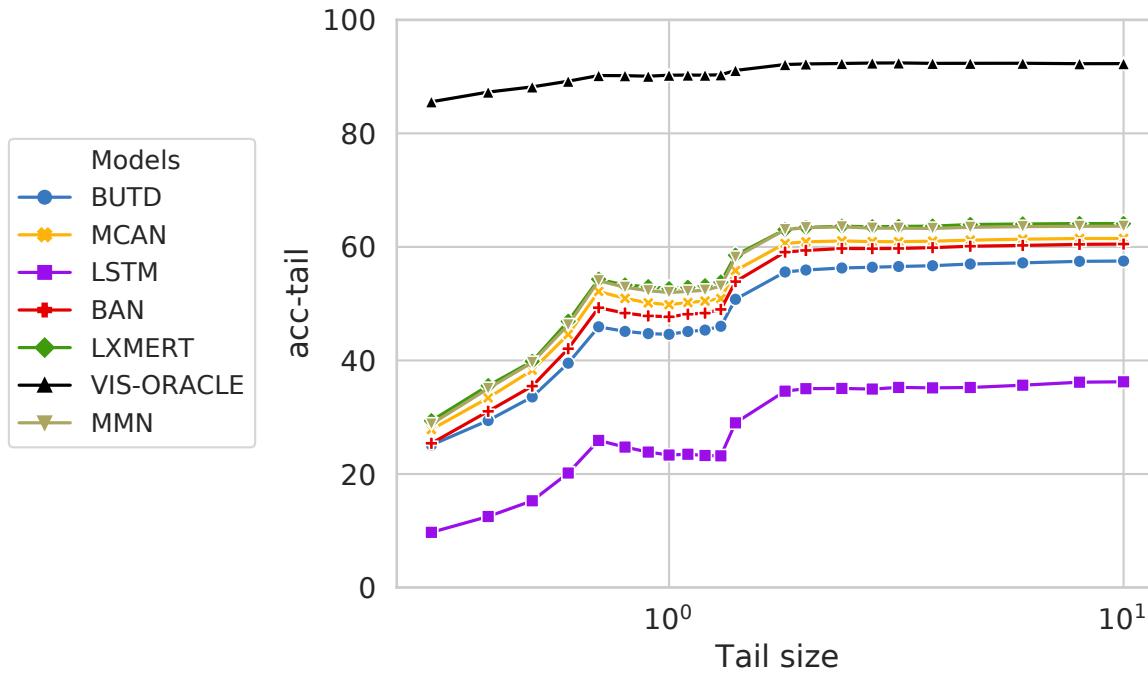
# Roses are Red, Violets are blue ... But should VQA expect them to?



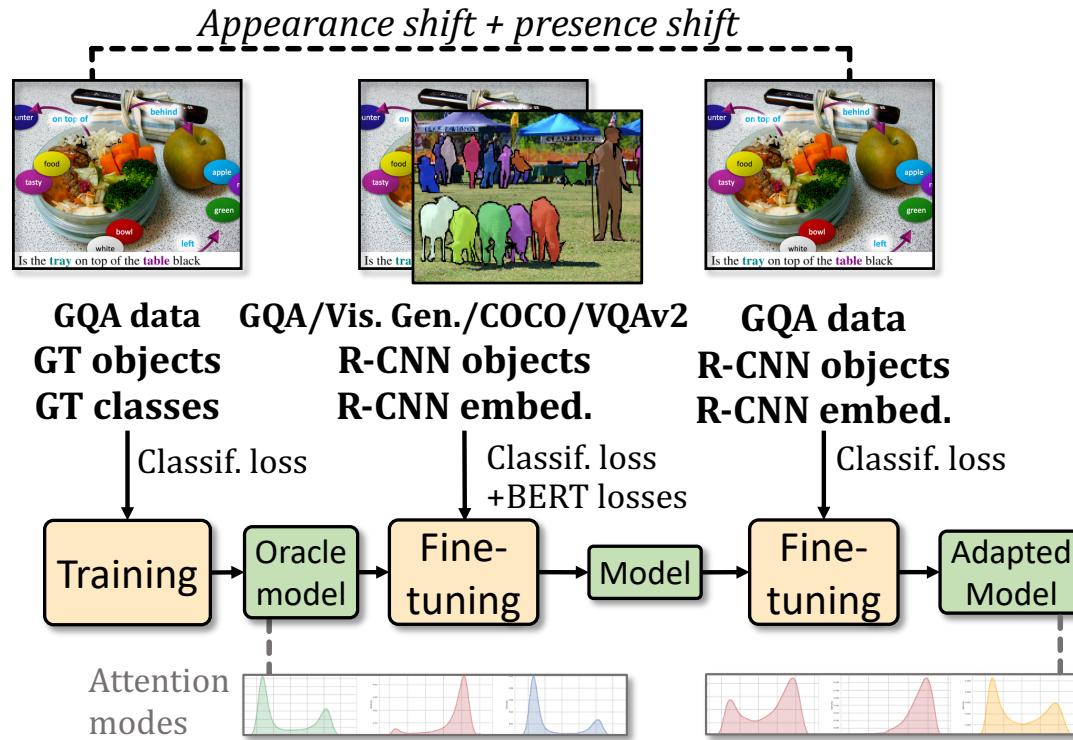
**"What is on  
the wall?"**



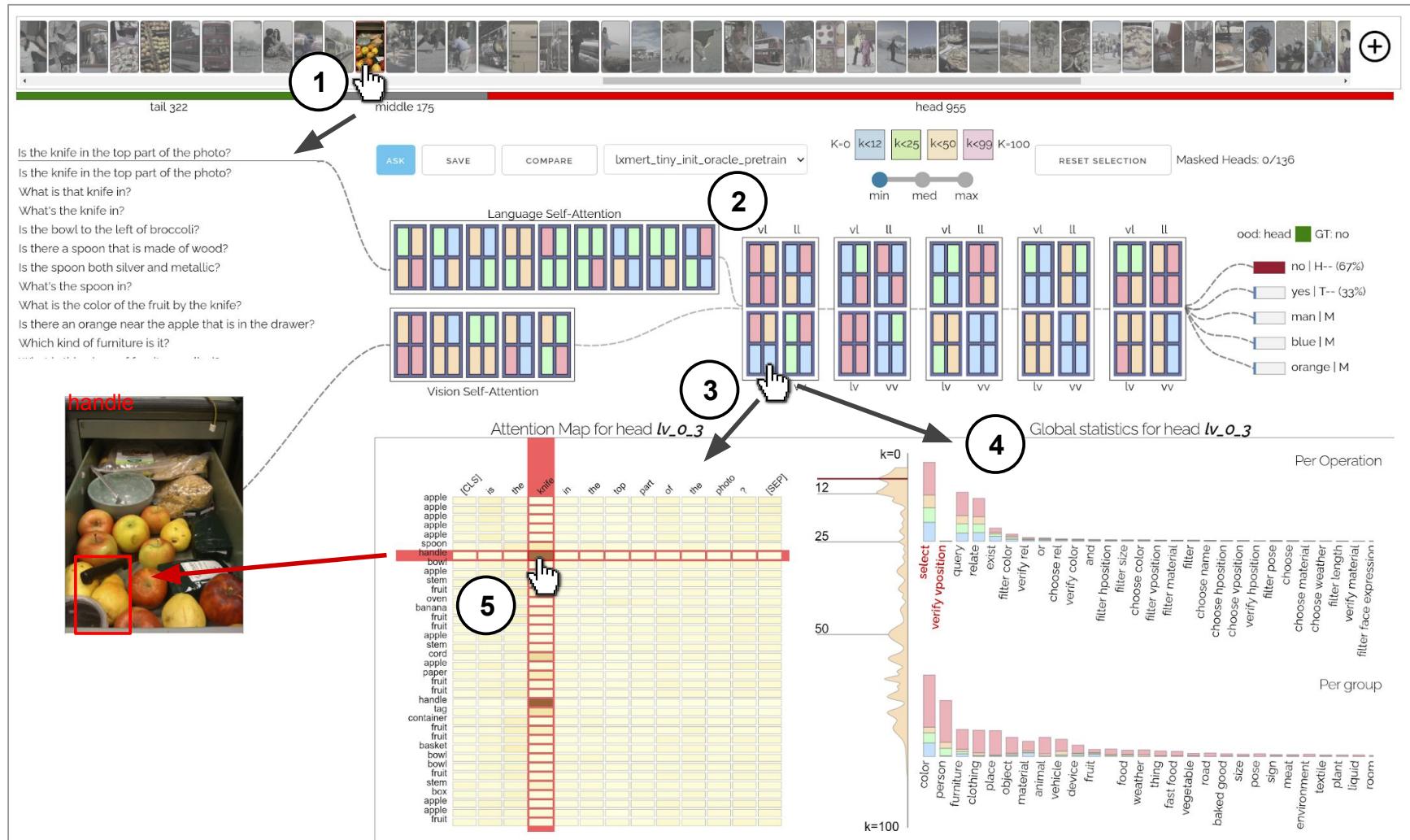
# Reasoning vs. bias exploitation



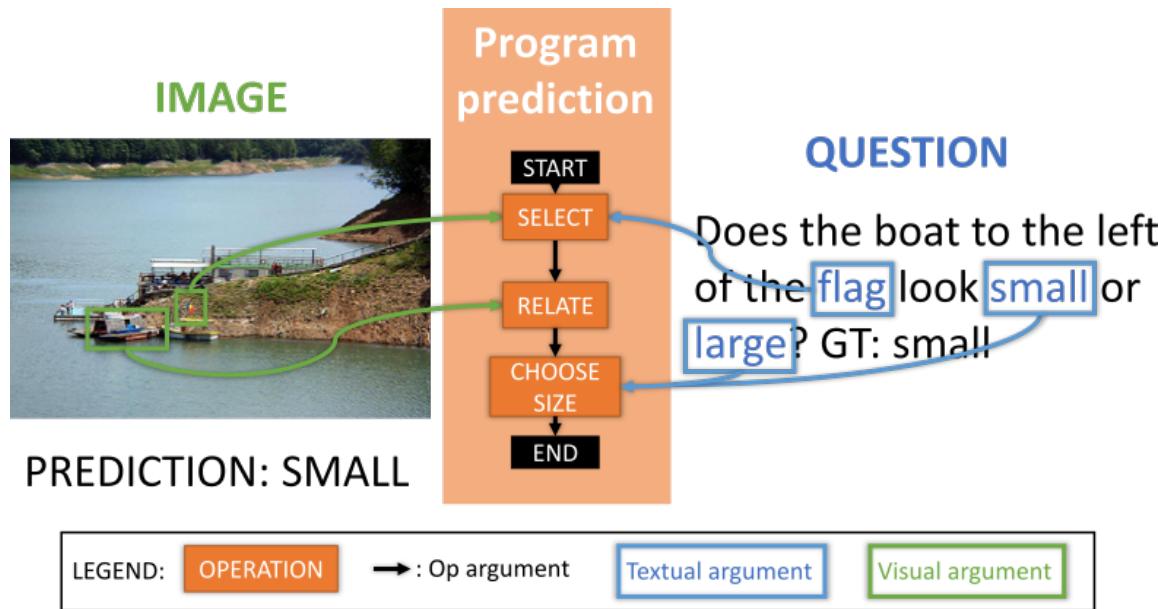
# Oracle transfer



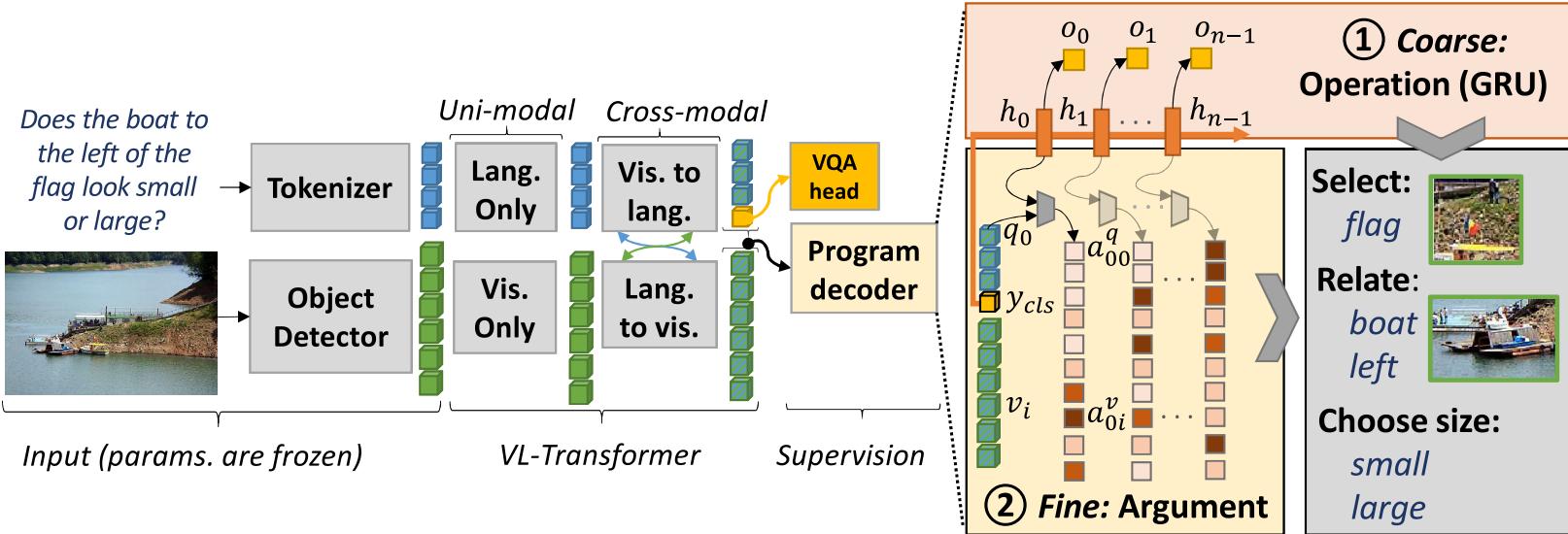
# Visualizing Reasoning Patterns



# GT reasoning programs

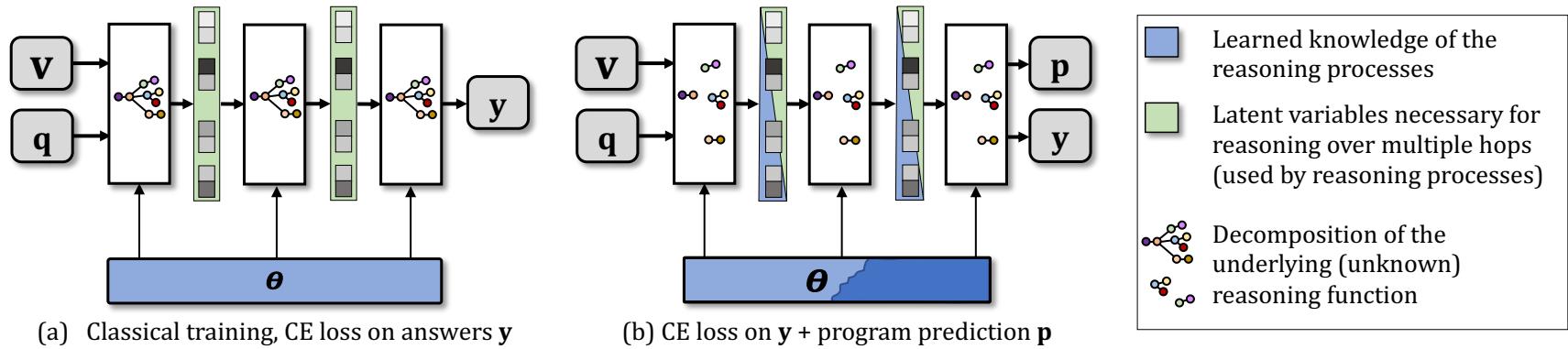


# Program prediction



# Program supervision

We theoretically show through PAC-bounds that program supervision as auxiliary loss can decrease sample complexity under some hypotheses.

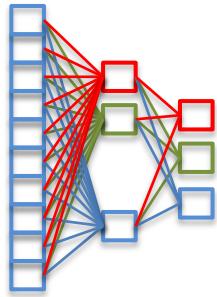


# Oracle + program supervision

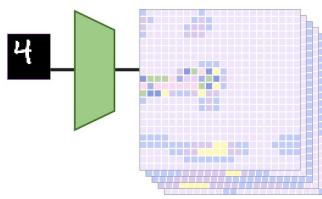
Model	Oracle transf.	Prog. sup.	GQA-OOD [20]		test-dev	GQA [17]			AUC <sup>†</sup> prog.
			acc-tail	acc-head		binary*	open*	test-std	
scratch	(a) Baseline		42.9	49.5	52.4	-	-	-	/
	(b) Oracle transfer	✓	$48.2 \pm 0.3$	$54.6 \pm 1.1$	$57.0 \pm 0.3$	74.5	42.1	57.3	/
	(c) Ours	✓	$48.8 \pm 0.1$	$56.1 \pm 0.3$	$57.8 \pm 0.2$	<b>75.4</b>	<b>43.0</b>	<b>58.2</b>	97.1
+ Lxmert	(d) Baseline		47.5	55.2	58.5	-	-	-	/
	(e) Oracle transfer	✓	47.1	54.8	58.4	77.1	42.6	58.8	/
	(f) Ours	✓	$48.0 \pm 0.6$	$56.6 \pm 0.6$	$59.3 \pm 0.3$	<b>77.3</b>	<b>44.1</b>	<b>59.7</b>	96.4

Table 1: Impact of program supervision on *Oracle transfer* [23] for vision-language transformers. LXMERT [36] pre-training is done on the GQA unbalanced training set. We report scores on GQA [17] (*test-dev* and *test-std*) and GQA-OOD (*test*). \* binary and open scores are computed on the test-std; <sup>†</sup> we evaluate visual argument prediction by computing AUC@0.66 on GQA-val.

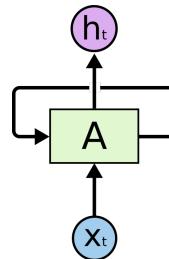
# Conclusion (1)



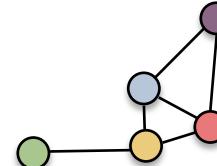
MLP



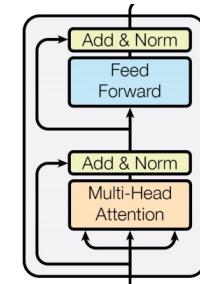
CNN /  
Convolutions



RNN /  
Recurrence



GN, GCN /  
Graphs, geometry



Transformers /  
Self-attention

*What do I know about the data and the task?*

*Nothing  
(vector space)*

*Translation  
equivariance*

*Sequential data,  
Markov property*

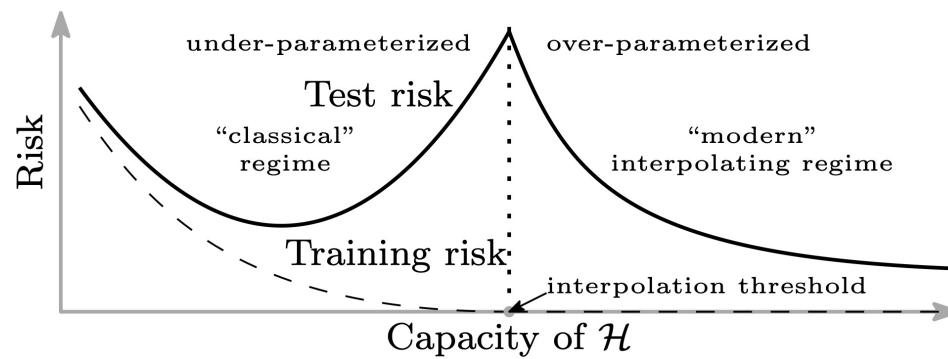
*Graph structured  
data*

*Permutation  
equivariance*

# Conclusion (2)

What is missing?

- Inductive biases for images
  - Convolutions vs. Transformers vs. MLPs
- ML formulations in deep learning
  - Self-supervised learning
  - Reinforcement Learning
- ML + physics (what do we model?)
- Theoretical models of Deep Learning



[Belkin et al., 2019]