

CV

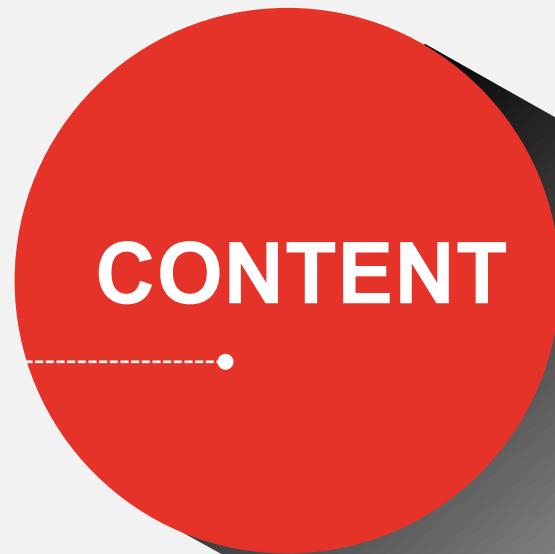
Few-Shot Learning

College of Information Science and Engineering
Ocean University of China

<http://vision.ouc.edu.cn>

Chenchen Qiu





01 Introduction to Few-shot Learning

**02 Popular methods to Few-shot
Learning**

03 Prior Knowledge

04 Multi-attention Network



01

Introduction to Few-shot Learning

In this part, I will briefly introduce the fundamental concept of few-shot learning, especially few-shot and one-shot image classification.

Comparison



Deep Learning

massive amounts of data
+
GPUs



Few-shot Learning

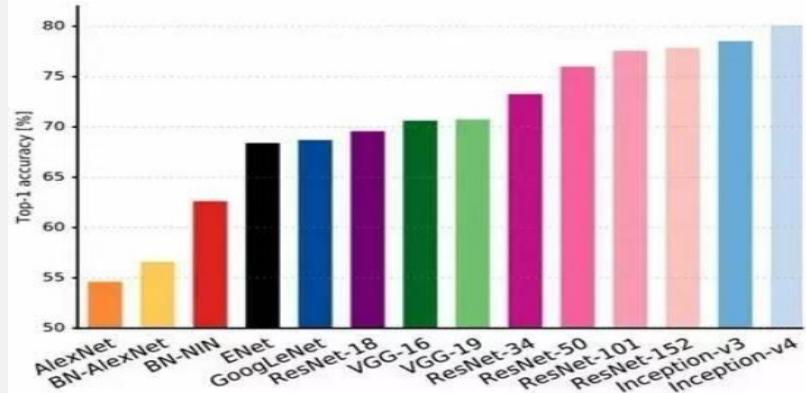
Training of machine learning
algorithms using a very small set of
training data



Zero-shot Learning

Solving a task despite not having
received any training examples of
that task

ImageNet ILSVRC



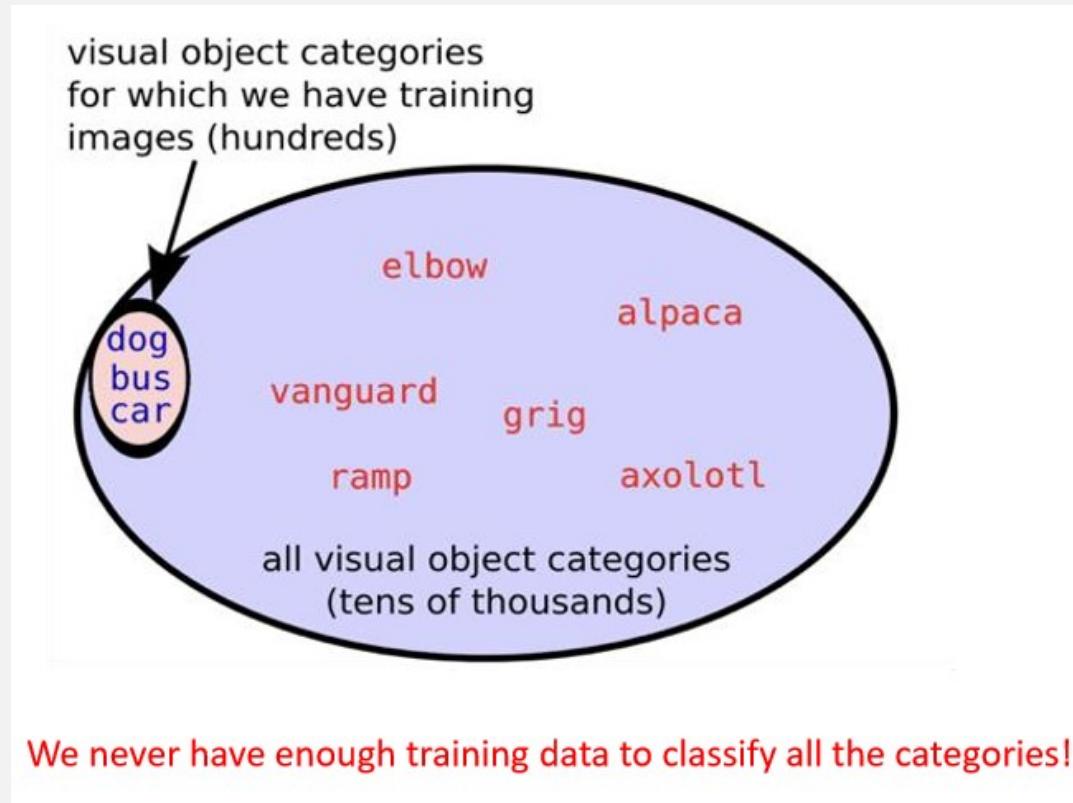
Voice Assistant



Face Verification



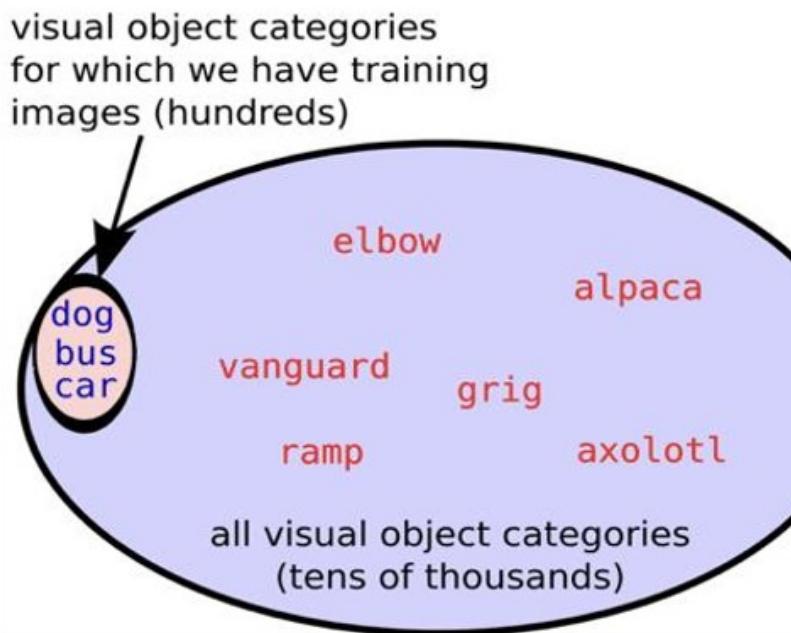
Dilemma of Large-scale Supervised Recognition



Problems :

- No memory: Knowledge learned is not retained
- Needs a large number of training examples

Dilemma of Large-scale Supervised Recognition



We never have enough training data to classify all the categories!

What we want? **Learn as humans do!**

- Humans have the ability to recognize without seeing examples (zero-shot learning);
- Retain learned knowledge from previous tasks & use it to help future learning(transfer learning);

Our Goal



One-shot learning aims to learn information about object categories from **one, or only a few**, training images.



Dataset



Omniglot

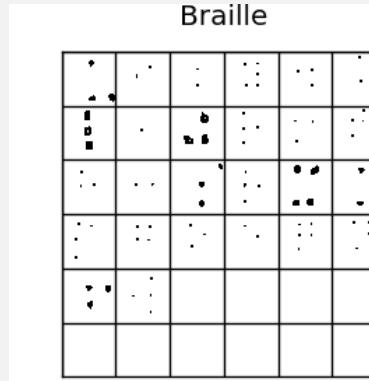


Minilmagenet



Animals with Attributes2

Omniglot



Bengali

ଫ୍ରେଜାନ୍ଟ	ଲ୍ୟାମ୍ବ	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍
କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍
କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍
କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍
କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍
କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍	କ୍ରୀପ୍ଟୋଗ୍ରାଫ୍

Sanskrit

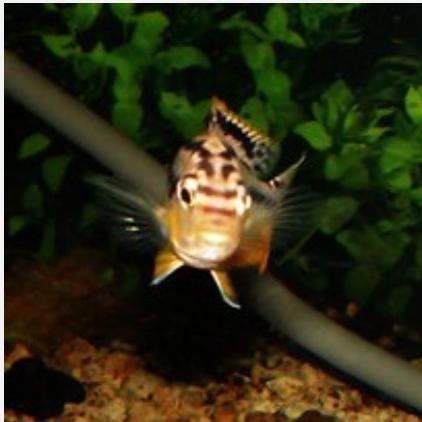
ପ	ଝ	ଙ୍ଗ	ଷ	ମ	ଳ	ଘ
ଟ	ଠ	କ	ବ୍ର	ଫ	ଅ	ଵ
ତ	ଦ	ଏ	ନ	ଜ୍ଞ	ଧ	ସ
ଦ	ଦୀ	ଆ	ଭୀ	ଓୟ	ତ୍ତ	ତ
ର	ଛ	ଣ	ଣ	ଲ	ଥ୍ରଦ	
କୁ	ଚି	କୁ	ଶ	ହ	ଶର୍କୁ	
ଚି	ଚି	ନ	ତ୍ତୁ	ତ୍ତୁ		
ରୁ	ଫି	କ୍ଷ	ବ			

Greek

Φ	Λ	Β	Σ	Ζ
Μ	Α	Κ	Χ	Ϝ
Β	Θ	Υ	Τ	Ω
Ω	Π	η	Ω	ε
Ρ	ξ	ς	ψ	

- This dataset contains **1623** different handwritten characters from **50** different alphabets.
- Each of the 1623 characters was drawn online via Amazon's Mechanical Turk by 20 different people.

MinilmageNet



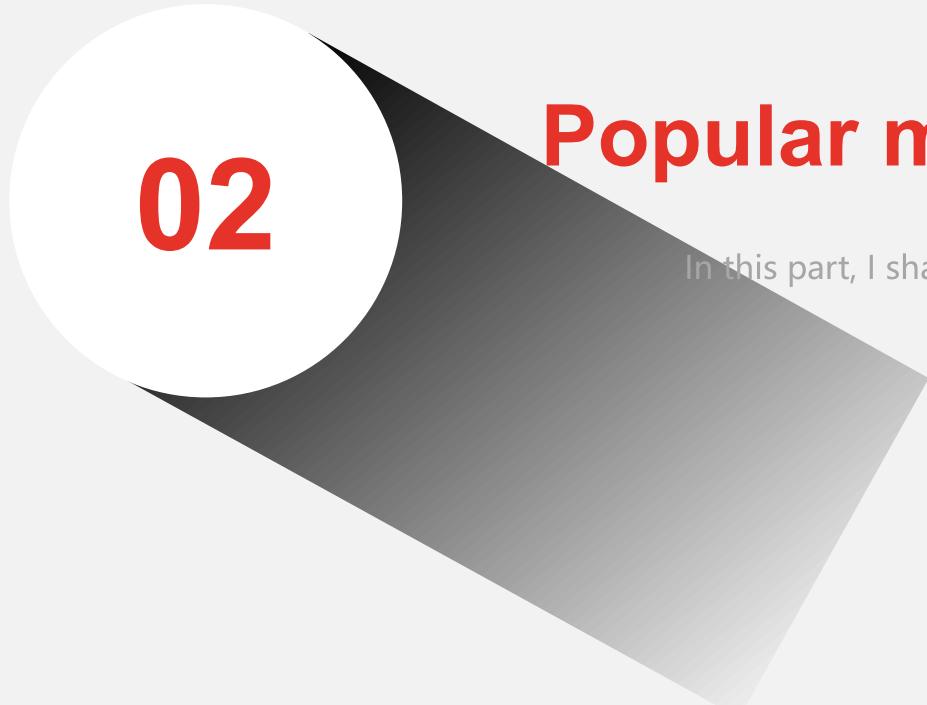
- The minilmageNet test used 60,000 images with 100 classes, each having 600 examples.
- 80 classes were used for training, and testing was done on the other 20.

Animals with Attributes2

<u>polar bear</u>	
black:	no
white:	yes
brown:	no
stripes:	no
water:	yes
eats fish:	yes



- This dataset provides a platform to benchmark transfer-learning algorithms, in particular attribute base classification and **zero-shot learning** .
- It consists of 37322 images of 50 animals classes.

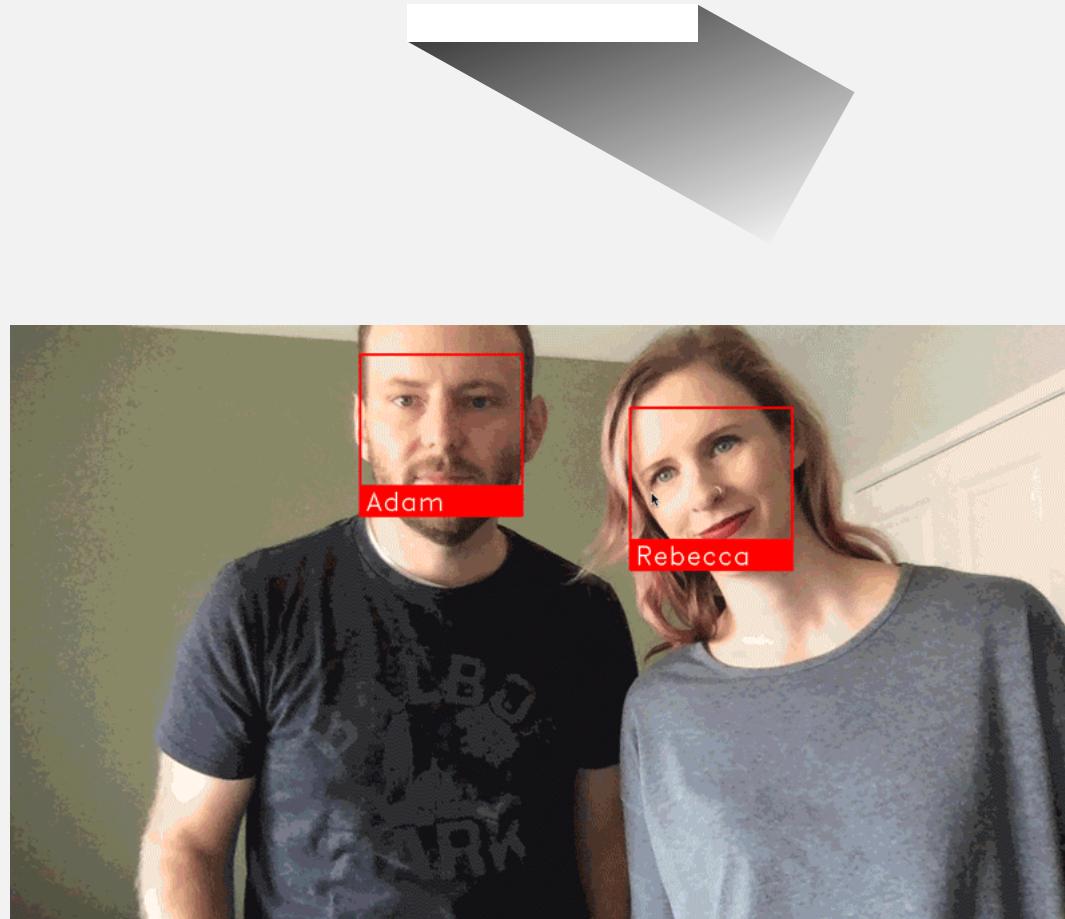
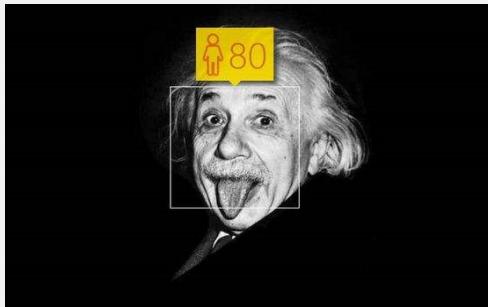


02

Popular methods to Few-shot Learning

In this part, I shall describe some of the papers that I gathered.

Face Recognition



Face Recognition



Bruce Banner



Ironman



Scarlet Witch



Black Widow



Learning from **one example** to recognize the person again.

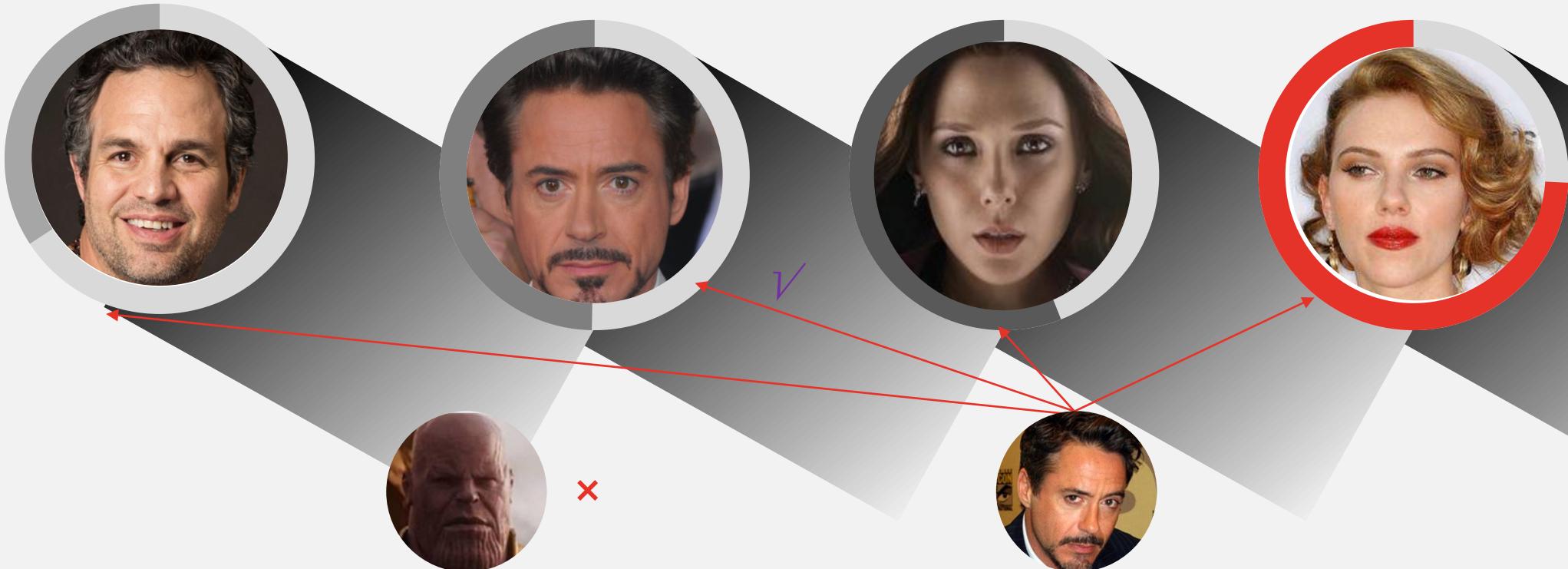
Learning a “similarity” function

$d(\text{img1}, \text{img2})$ = degree of difference between images

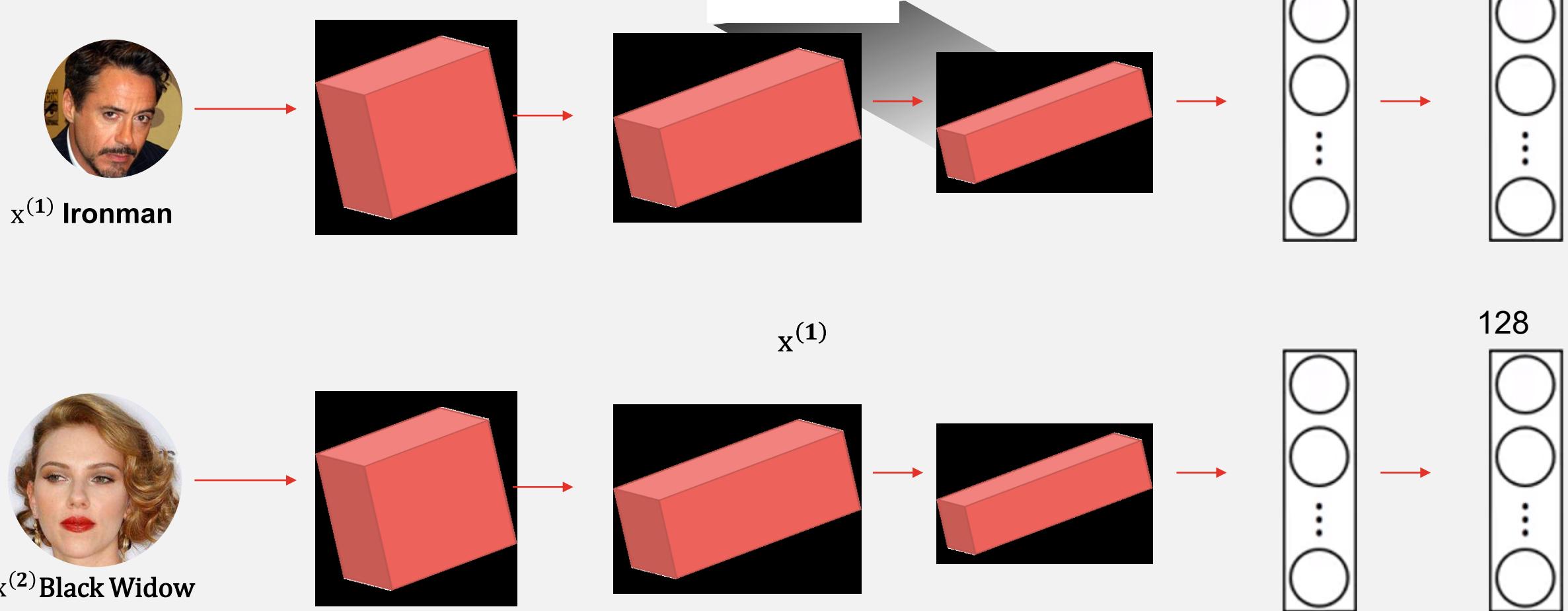
If $d(\text{img1}, \text{img2}) \leq \tau$ same
 $> \tau$ different



Face verification

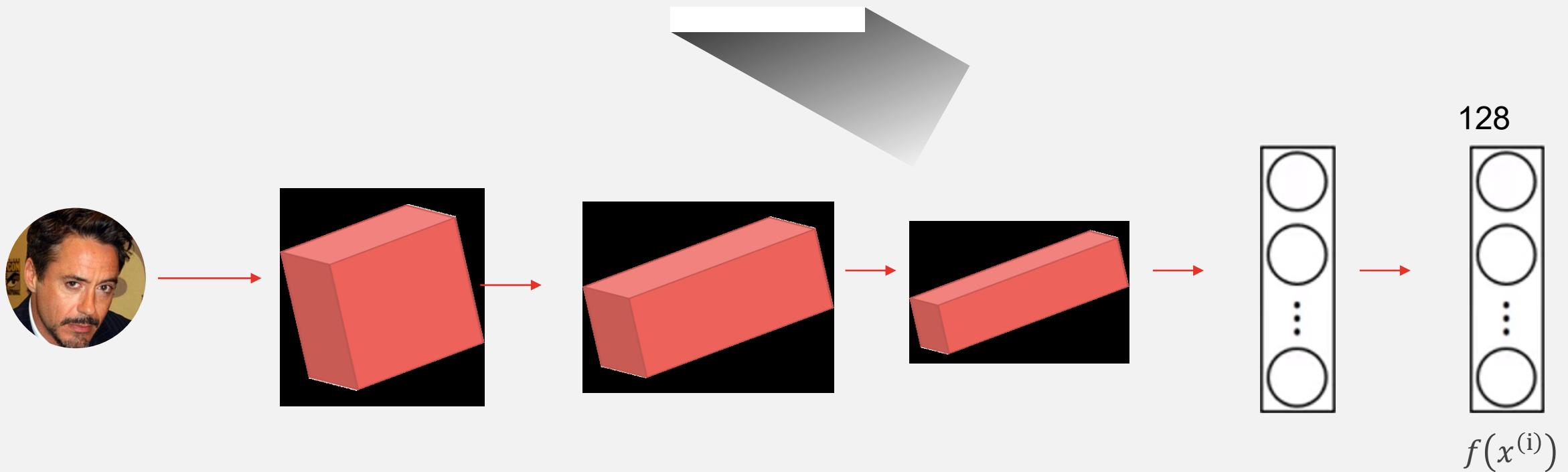


Siamese Network



$$d(x^{(1)}, x^{(2)}) = \|f(x^{(1)}) - f(x^{(2)})\|_d^2$$

Goal of learning



- Parameters of NN define an **encoding** $f(x^{(i)})$
- Learn parameters so that:

If $x^{(i)}, x^{(j)}$ are the person, $\|f(x^{(i)}) - f(x^{(j)})\|^2$ is small.

If $x^{(i)}, x^{(j)}$ are different persons, $\|f(x^{(i)}) - f(x^{(j)})\|^2$ is large.

Face Recognition



Anchor



Positive



Anchor



Negative

Goal: $\|f(A) - f(P)\|^2 \leq \|f(A) - f(N)\|^2$

$$\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 \leq 0$$

$$\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 \leq \alpha$$

Loss Function

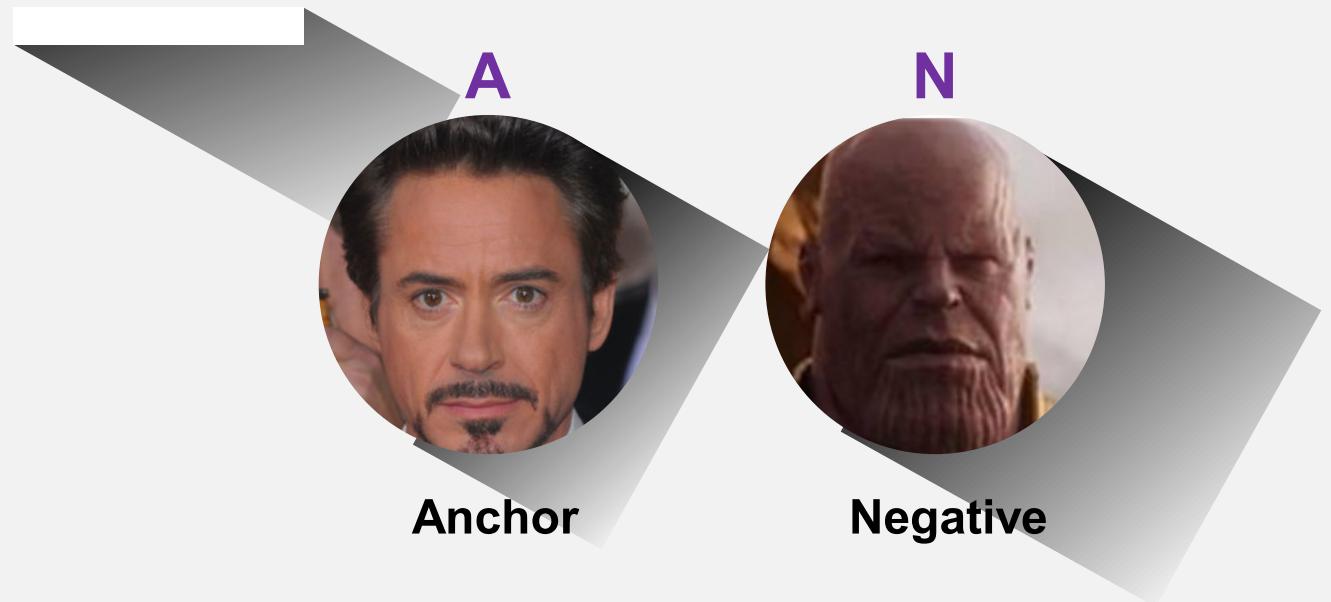


$$L(A, P, N) = \max(\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha, 0)$$

$$J = \sum_{i=1}^m L(A^{(i)}, P^{(i)}, N^{(i)})$$

Training set: 10k pictures of 1k persons

Choosing the triplets A,P,N



During training, if A, P, N are chosen randomly,
 $\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha \leq 0$ is easily satisfied.

$$d(A, P) + \alpha \leq d(A, N)$$

Training set using triplet loss

Anchor



Positive



Anchor



$$J = \sum_{i=1}^m L(A^{(i)}, P^{(i)}, N^{(i)})$$

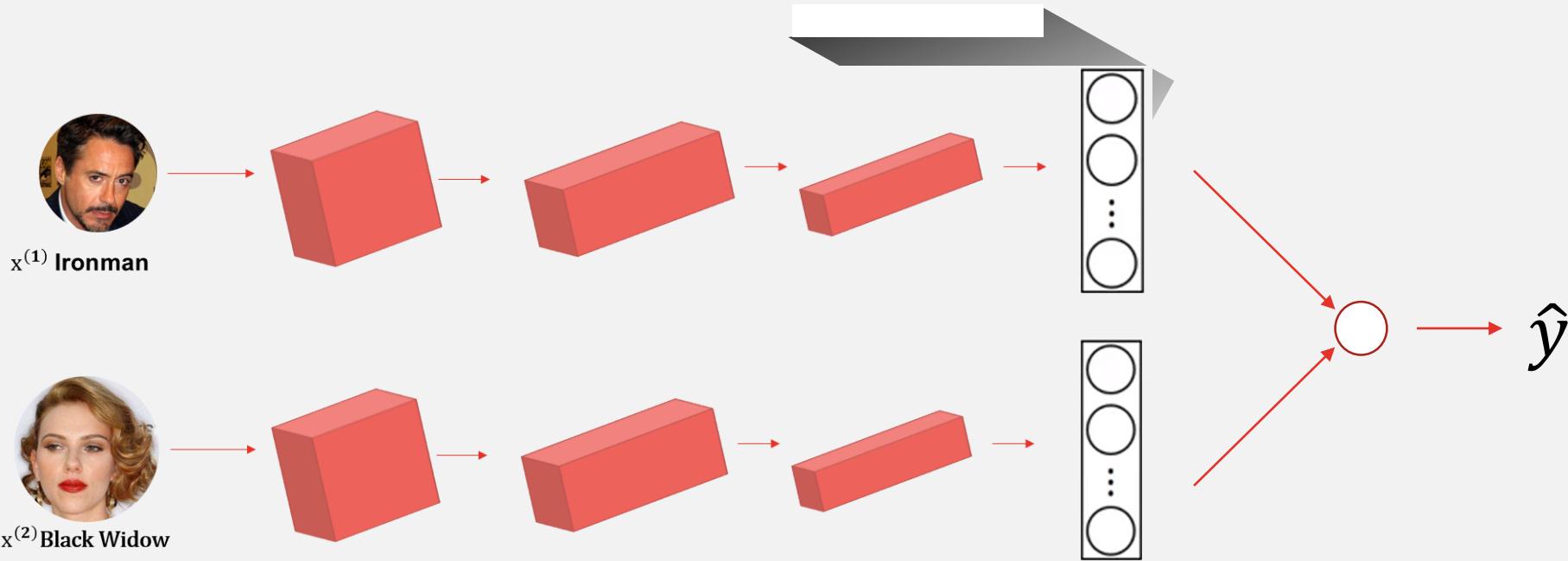
:

:

:

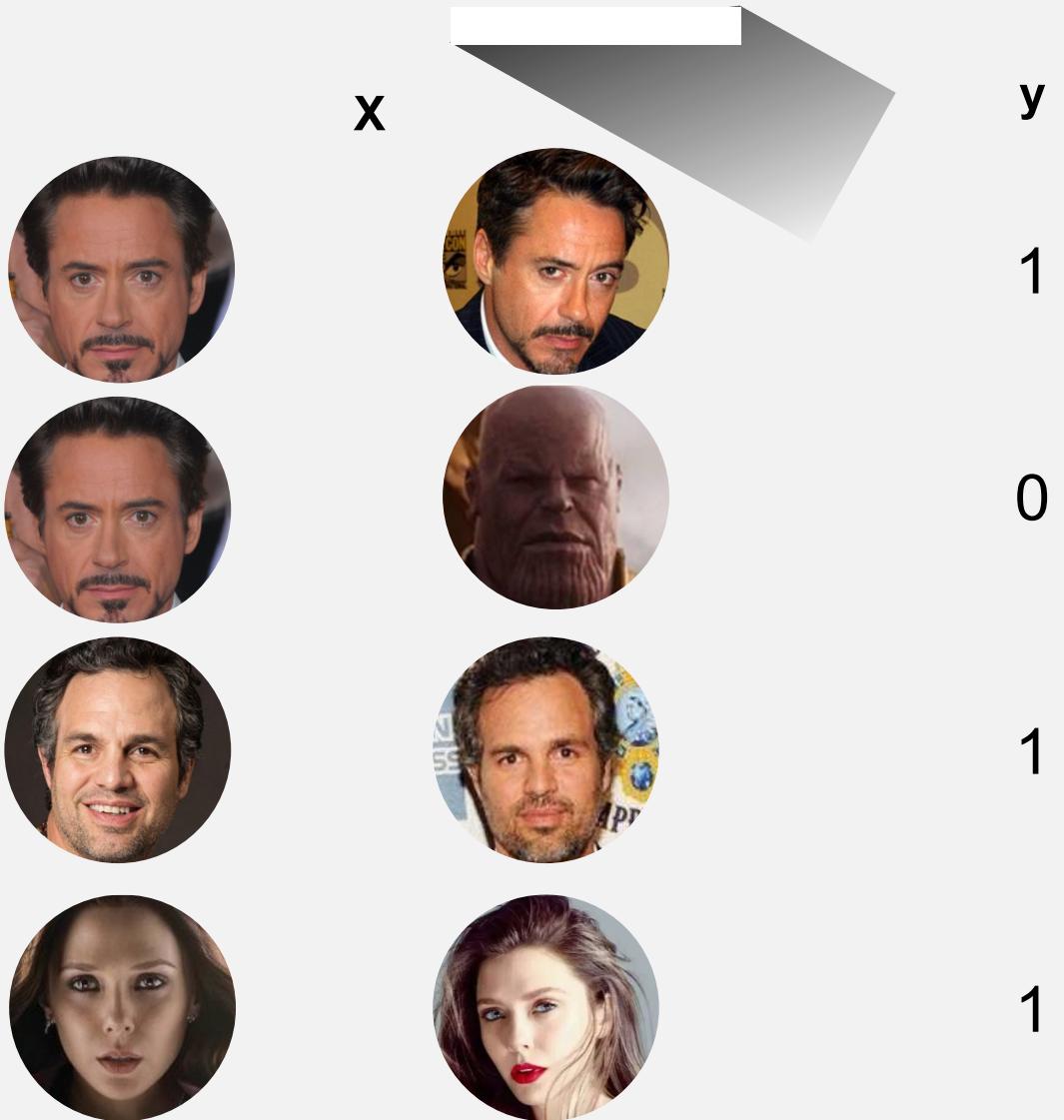


Learning the similarity function

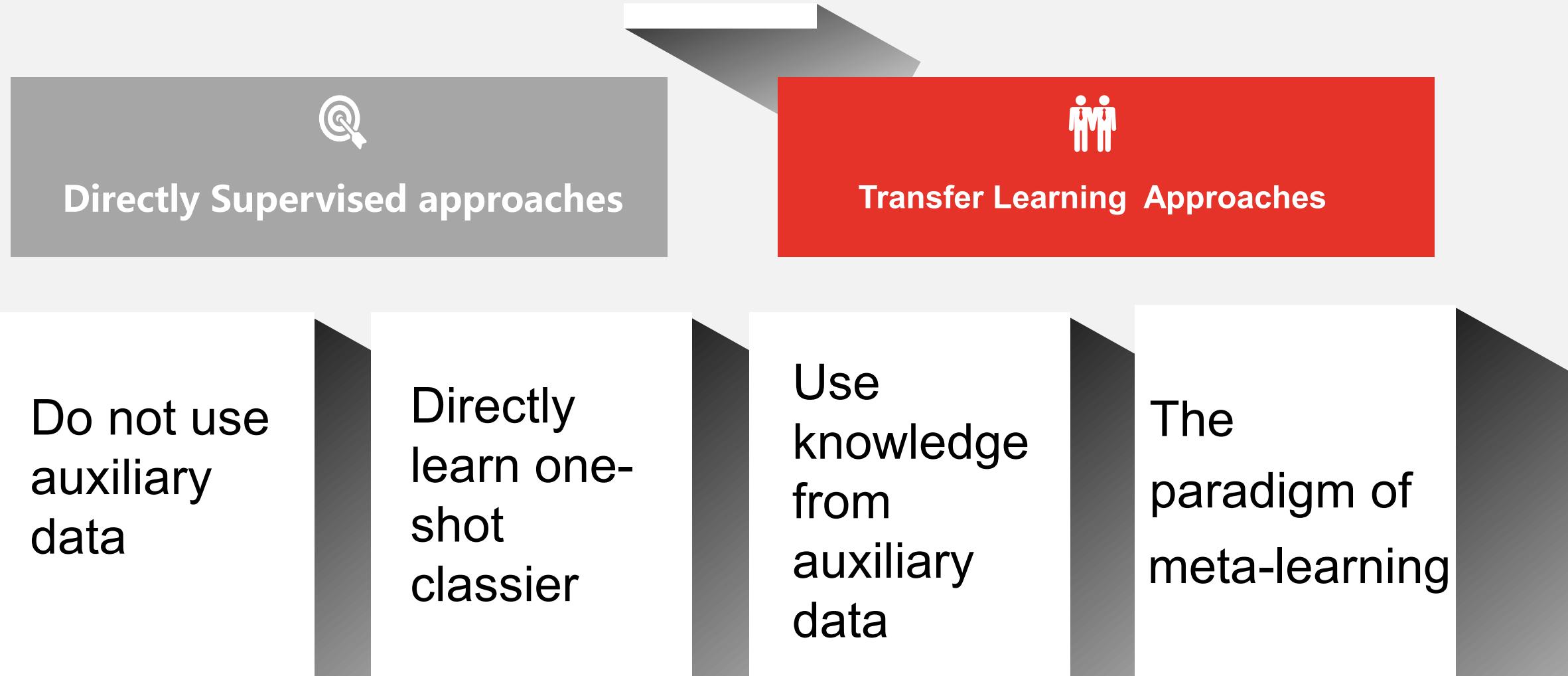


$$y = \sigma \left(\sum_{k=1}^{128} w_i \left| f(x^{(i)})_k - f(x^{(j)})_k \right| + b \right)$$

Face verification supervised learning



Outlines of One-shot Learning Methods



Directly supervised learning-based approaches



Instance-based learning

- K-nearest neighbor

Non-parametric methods

- Fei-Fei et al. A Bayesian approach to unsupervised one-shot learning of object categories. In *CVPR*, 2003.
- Fei-Fei et al. One-shot learning of object categories. *TPAMI*, 2006.

Transfer learning-based approaches



Attribute-based
Algorithms

- M2LATM
- TMV-HLP



Meta-learning
Algorithms

- MAML
- Meta-learn LSTM
- Meta-Net



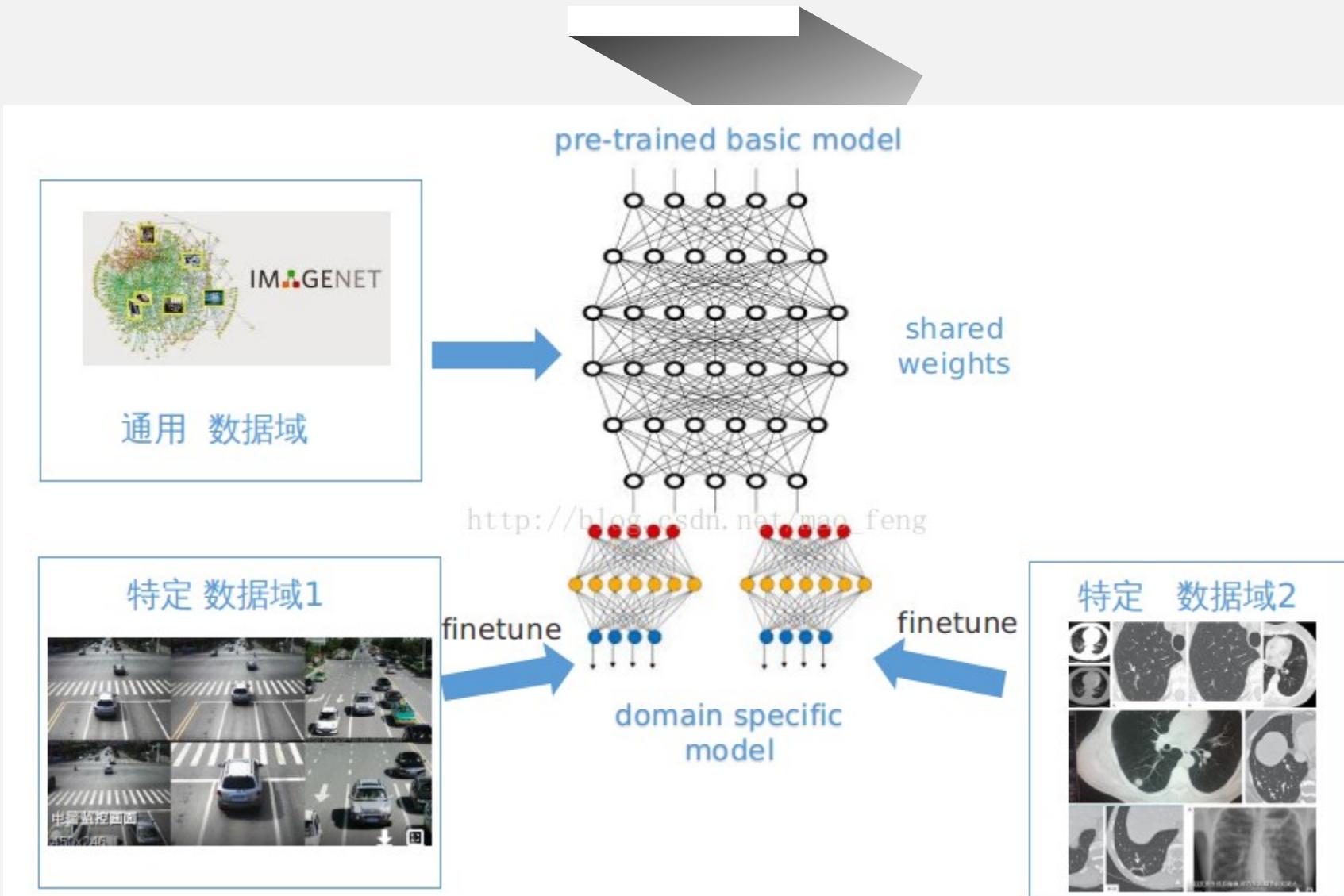
Metric-learning
Algorithms

- Matching Nets
- Prototypical
Nets
- Relation Net

Data Augmentation for One-shot Learning

-
- 01 Learning one-shot models by utilizing the **manifold information** of large amount of unlabeled data in a **semi-supervised** or transductive setting
 - 02 Adaptively learning the one-shot classifiers from off-shelf trained models
 - 03 Borrowing examples from **relevant categories** or semantic vocabularies to augment training set
 - 04 Synthesizing **new** labelled training **data** by rendering virtual examples or composing synthesized representations or distorting existing training examples;
 - 05 Generating **new examples** using Generative Adversarial Networks (**GANs**);
 - 06 Attribute-guided augmentation (**AGA**) to **synthesize samples** at desired values or strength

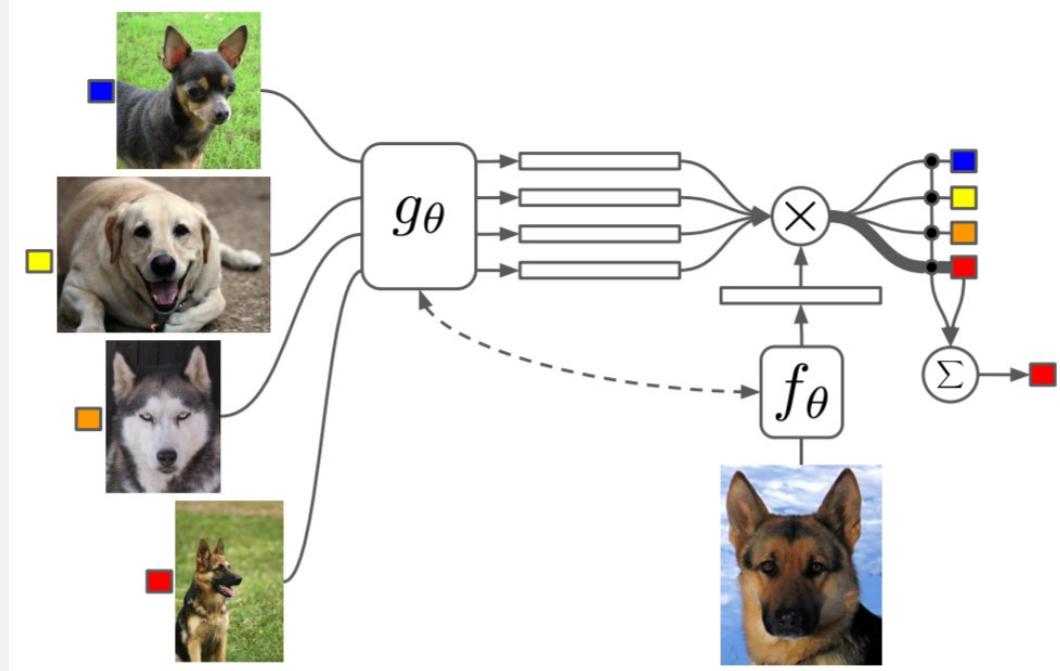
Finetune



Metric-learning Algorithms

- **Siamese Neural Network**
- **Matching Networks**
- **Prototypical Networks**
- **Relation Network**

Matching Networks



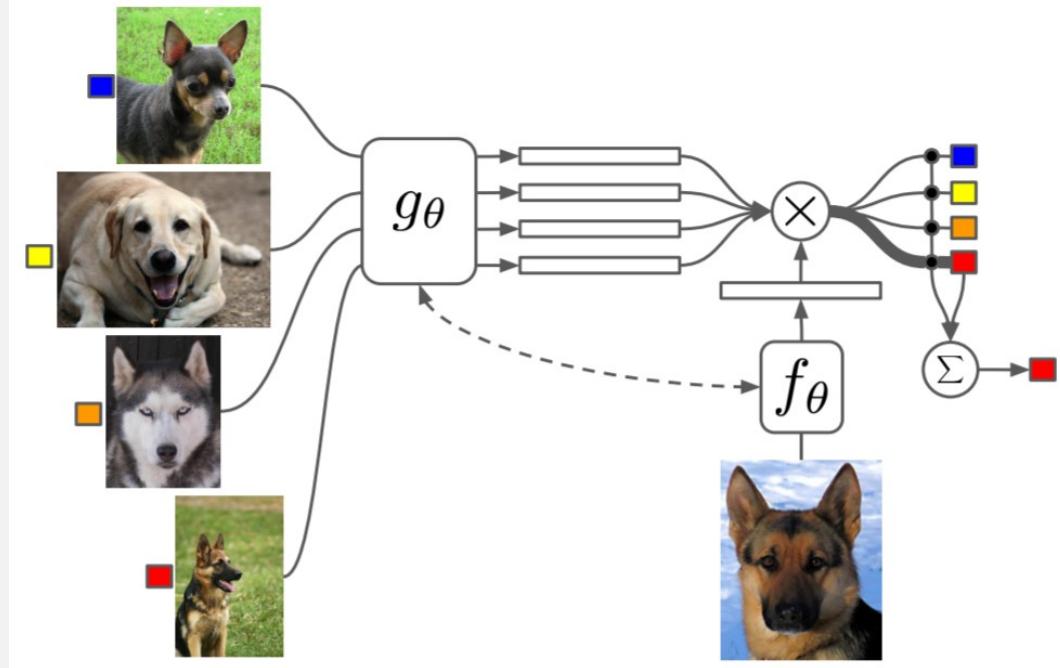
Novelty of their work

Modeling level

Training procedure



Matching Networks



Novelty of Modeling level

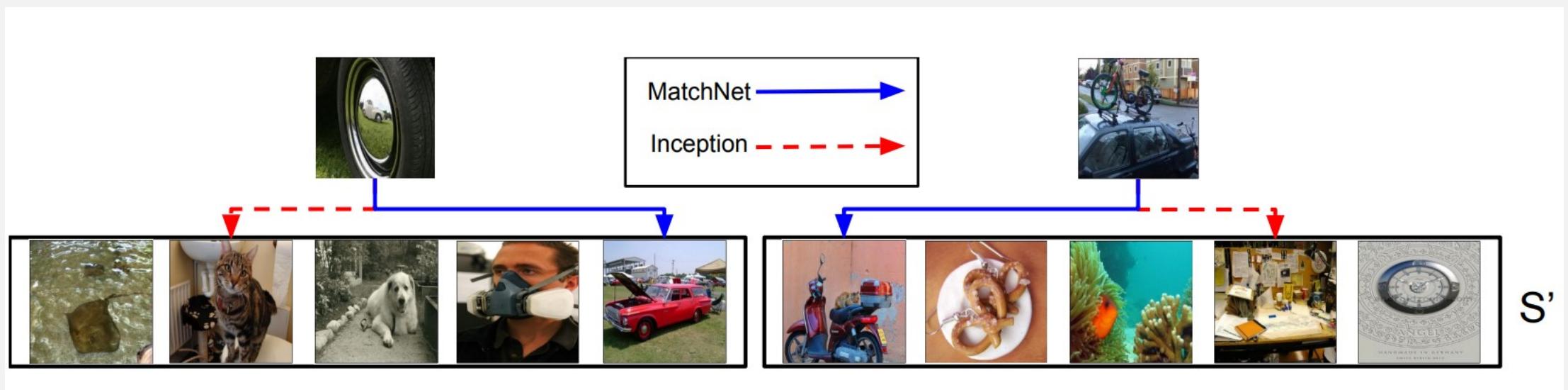
$$\arg \max_y P(y|\hat{x}, S),$$

$$\hat{y} = \sum_{i=1}^k a(\hat{x}, x_i) y_i$$

$$a(\hat{x}, x_i) = \frac{e^{c(f(\hat{x}), g(x_i))}}{\sum_{i=1}^k e^{c(f(\hat{x}), g(x_i))}}$$

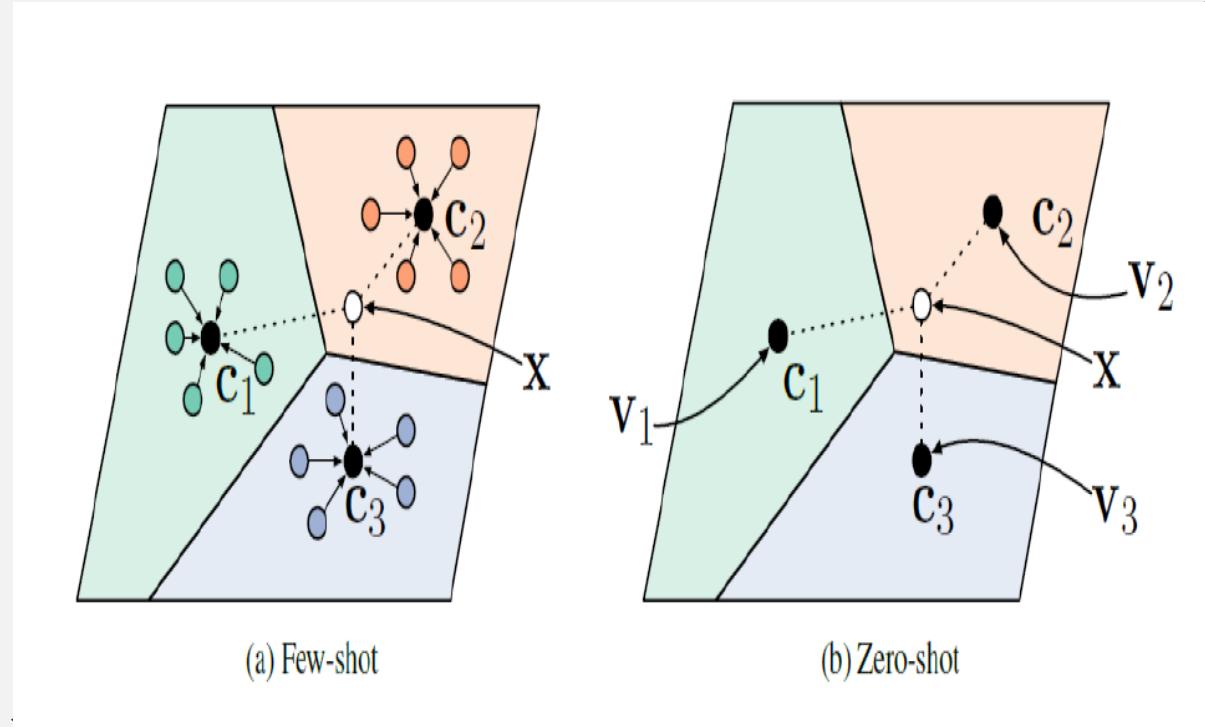
Matching Networks

Training procedure



$$\theta = \arg \max_{\theta} E_{L \sim T} \left[E_{S \sim L, B \sim L} \left[\sum_{(x,y) \in B} \log P_{\theta} (y|x, S) \right] \right]$$

Prototypical Networks for Few-shot Learning



Novelty of their work



Prototypical Networks for Few-shot Learning

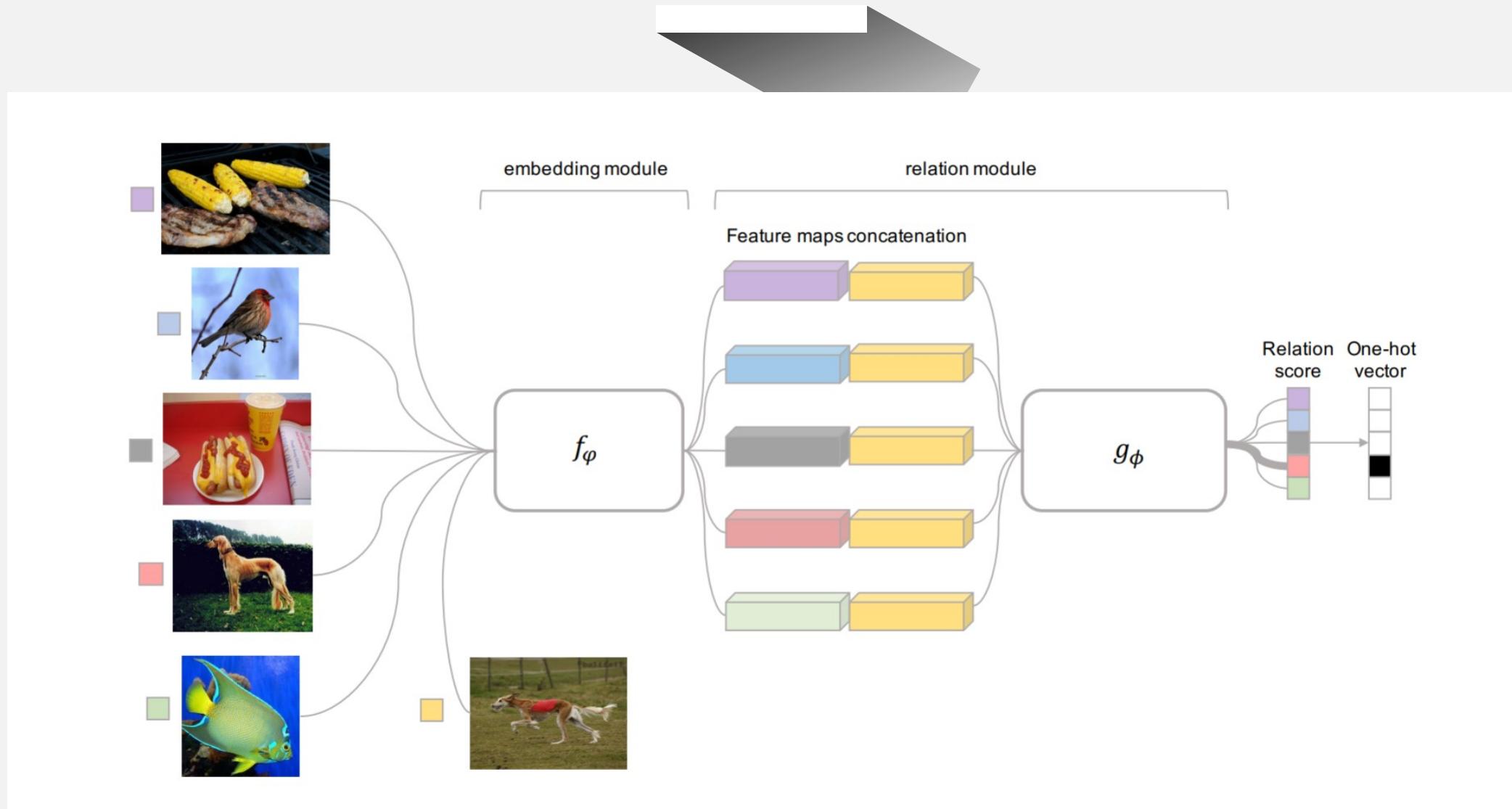
Table 1: Few-shot classification accuracies on Omniglot.

Model	Dist.	Fine Tune	5-way Acc.		20-way Acc.	
			1-shot	5-shot	1-shot	5-shot
MATCHING NETWORKS [29]	Cosine	N	98.1%	98.9%	93.8%	98.5%
MATCHING NETWORKS [29]	Cosine	Y	97.9%	98.7%	93.5%	98.7%
NEURAL STATISTICIAN [6]	-	N	98.1%	99.5%	93.2%	98.1%
PROTOTYPICAL NETWORKS (OURS)	Euclid.	N	98.8%	99.7%	96.0%	98.9%

Table 2: Few-shot classification accuracies on *miniImageNet*. All accuracy results are averaged over 600 test episodes and are reported with 95% confidence intervals. *Results reported by [22].

Model	Dist.	Fine Tune	5-way Acc.	
			1-shot	5-shot
BASELINE NEAREST NEIGHBORS*	Cosine	N	$28.86 \pm 0.54\%$	$49.79 \pm 0.79\%$
MATCHING NETWORKS [29]*	Cosine	N	$43.40 \pm 0.78\%$	$51.09 \pm 0.71\%$
MATCHING NETWORKS FCE [29]*	Cosine	N	$43.56 \pm 0.84\%$	$55.31 \pm 0.73\%$
META-LEARNER LSTM [22]*	-	N	$43.44 \pm 0.77\%$	$60.60 \pm 0.71\%$
PROTOTYPICAL NETWORKS (OURS)	Euclid.	N	$49.42 \pm 0.78\%$	$68.20 \pm 0.66\%$

Relation network for few-shot learning



Relation network for few-shot learning

motivation

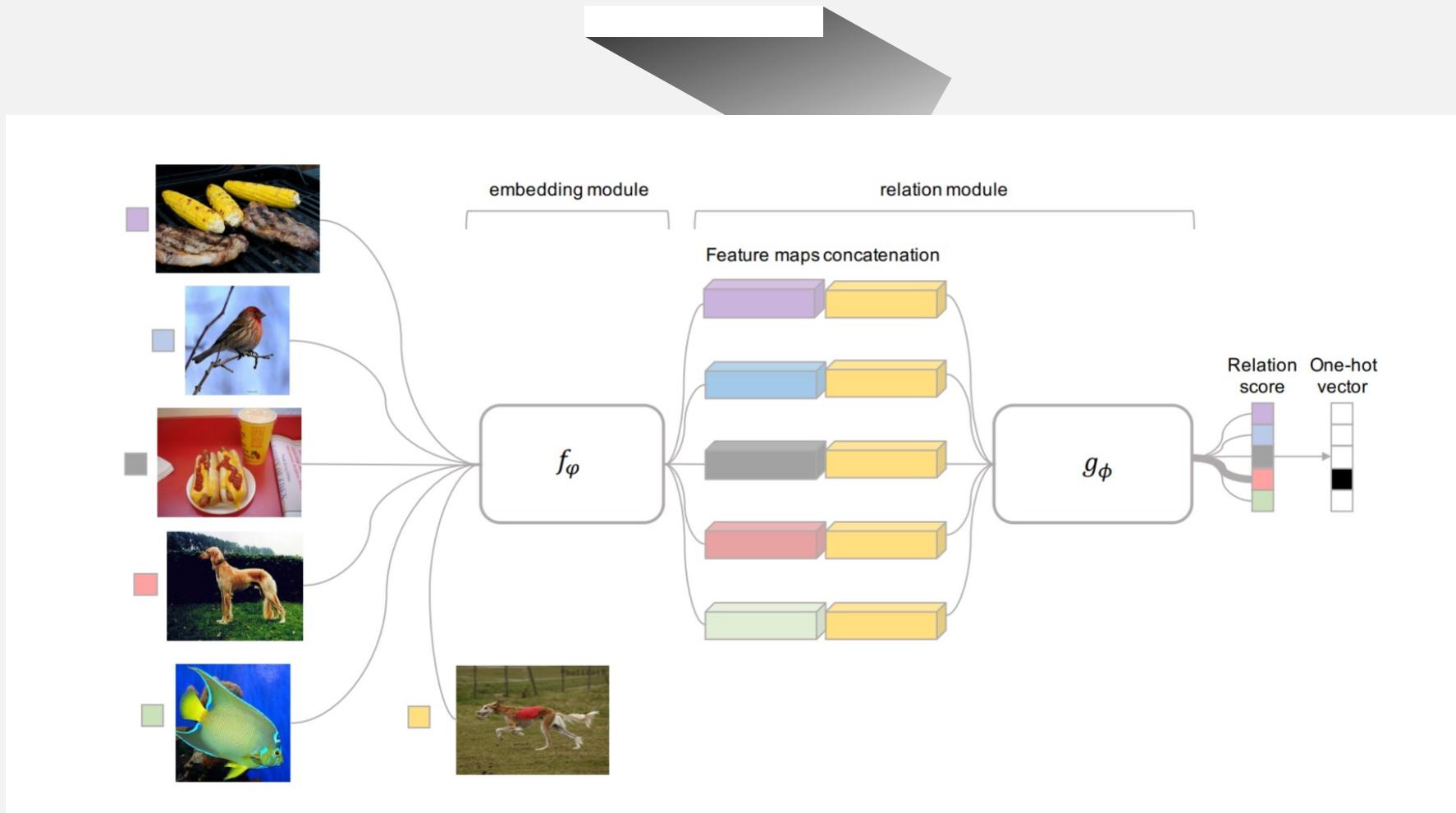


• H4 懶

卤蛋精!

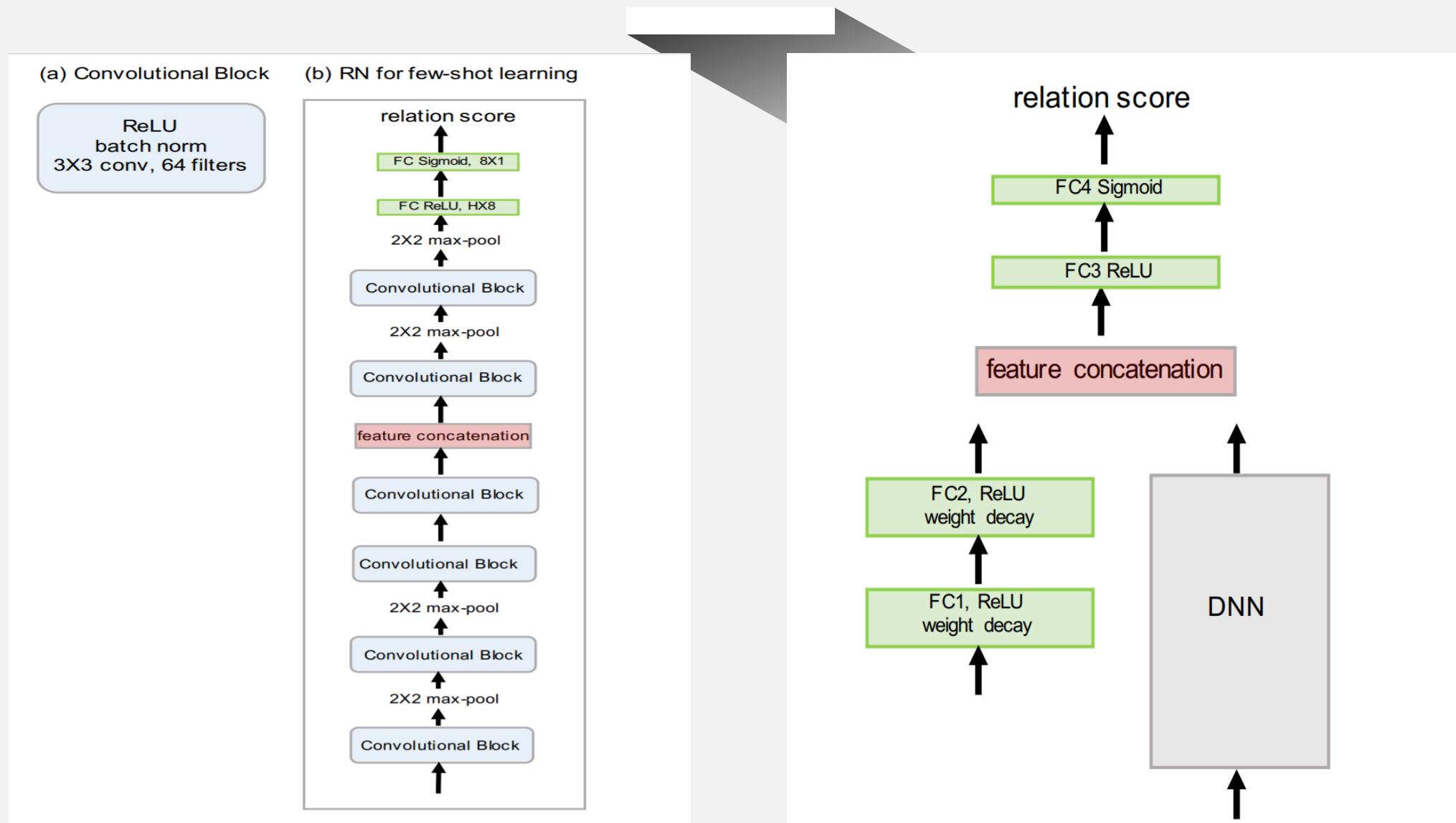


Relation network for few-shot learning



$$r_{i,j} = g_\phi \left(C \left(f_\phi(v_c), f_\phi(x_j) \right) \right), \quad i = 1, 2, 3, \dots, C$$

Relation network for few-shot learning

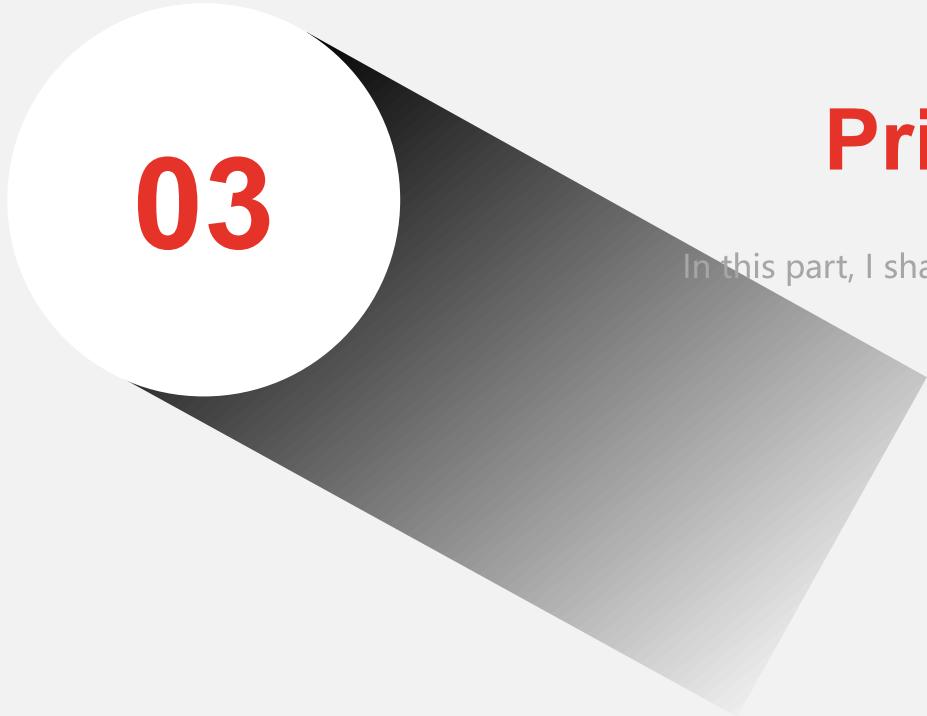


Relation network for few-shot learning

Fitting
everything!

Model	Fine Tune	5-way Acc.		20-way Acc.	
		1-shot	5-shot	1-shot	5-shot
MANN [32]	N	82.8%	94.9%	-	-
CONVOLUTIONAL SIAMESE NETS [20]	N	96.7%	98.4%	88.0%	96.5%
CONVOLUTIONAL SIAMESE NETS [20]	Y	97.3%	98.4%	88.1%	97.0%
MATCHING NETS [39]	N	98.1%	98.9%	93.8%	98.5%
MATCHING NETS [39]	Y	97.9%	98.7%	93.5%	98.7%
SIAMESE NETS WITH MEMORY [18]	N	98.4%	99.6%	95.0%	98.6%
NEURAL STATISTICIAN [8]	N	98.1%	99.5%	93.2%	98.1%
META NETS [27]	N	99.0%	-	97.0%	-
PROTOTYPICAL NETS [36]	N	98.8%	99.7%	96.0%	98.9%
MAML [10]	Y	98.7 ± 0.4%	99.9 ± 0.1%	95.8 ± 0.3%	98.9 ± 0.2%
RELATION NET	N	99.6 ± 0.2%	99.8 ± 0.1%	97.6 ± 0.2%	99.1 ± 0.1%

Model	FT	5-way Acc.	
		1-shot	5-shot
MATCHING NETS [39]	N	43.56 ± 0.84%	55.31 ± 0.73%
META NETS [27]	N	49.21 ± 0.96%	-
META-LEARN LSTM [29]	N	43.44 ± 0.77%	60.60 ± 0.71%
MAML [10]	Y	48.70 ± 1.84%	63.11 ± 0.92%
PROTOTYPICAL NETS [36]	N	49.42 ± 0.78%	68.20 ± 0.66%
RELATION NET	N	50.44 ± 0.82%	65.32 ± 0.70%



03

Prior Knowledge

In this part, I shall introduce the RNN, attention mechanism, VQA problem.

Prior Knowledge



RNN

I'll introduce the recurrent neural network, including LSTM, GRU and word2vec.



Attention mechanism

Summary of attention mechanism.



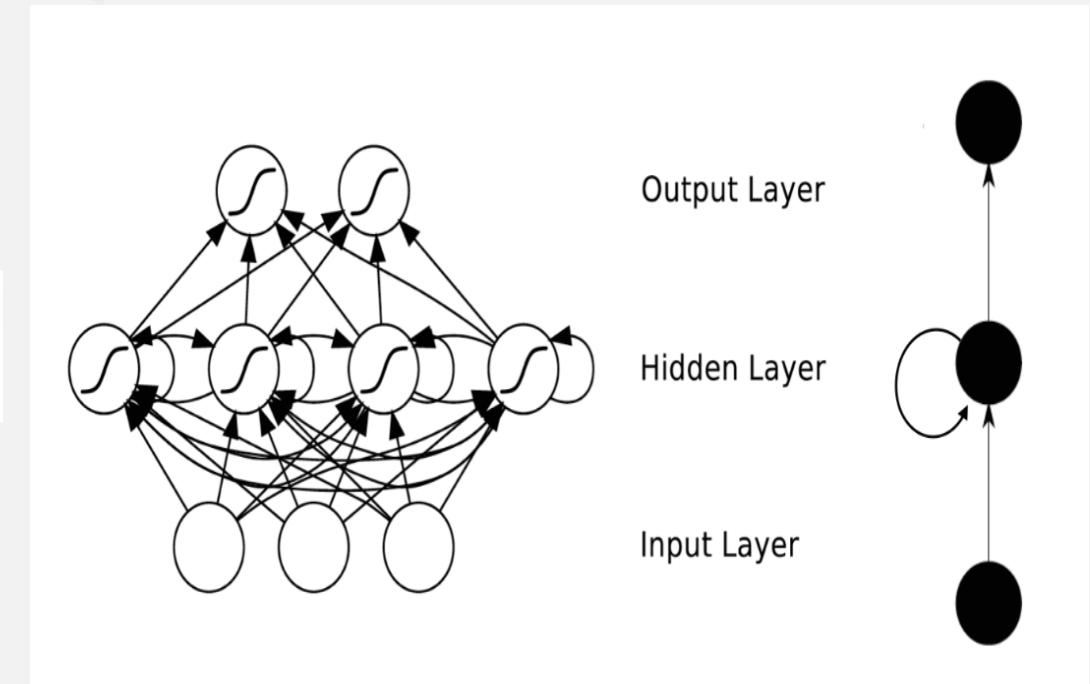
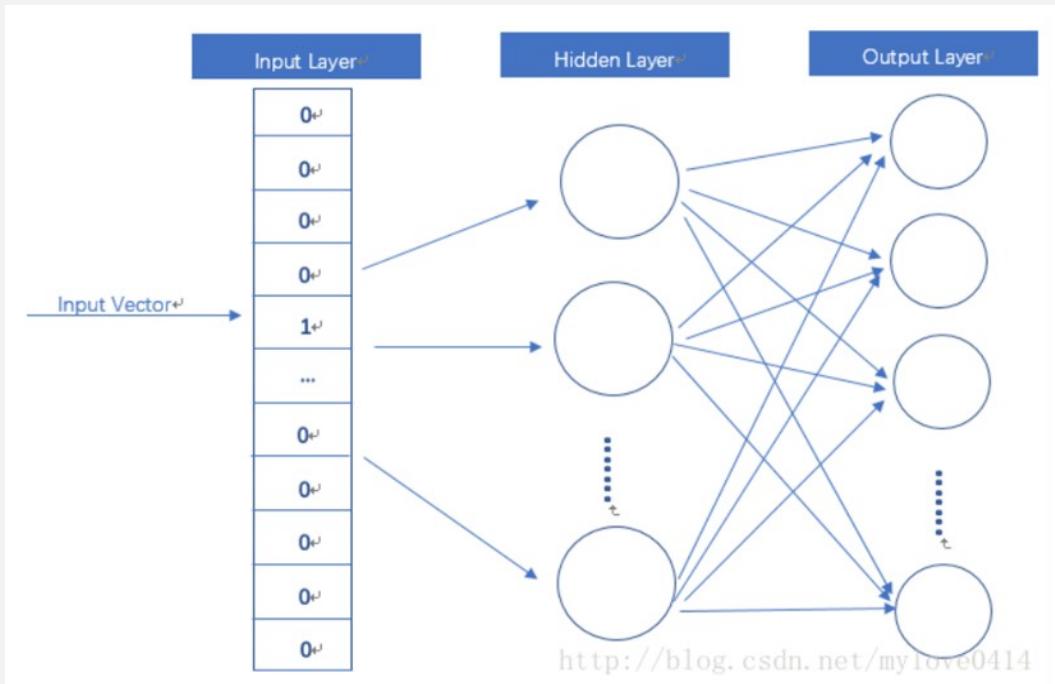
VQA

I'll describe the visual question answering problem.

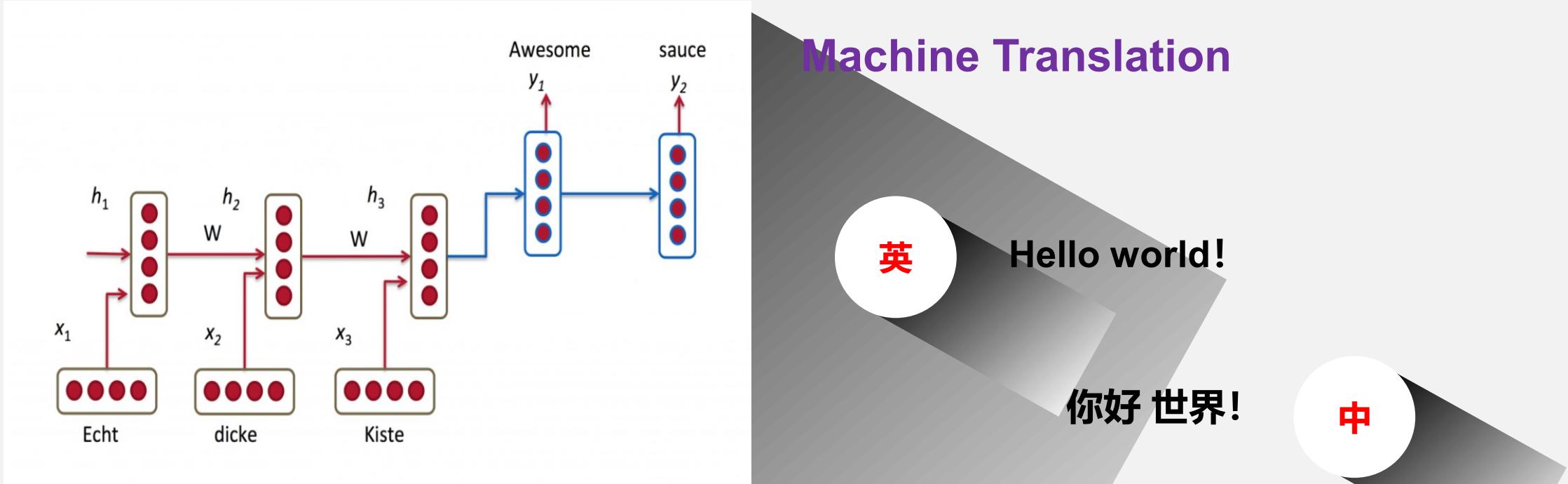


- 01 **What is Recurrent Neural Network?**
- 02 **Several Units**
- 03 **Representing words**
- 04 **Application of Recurrent Neural Network**

What are RNNs?



What can RNNs do?



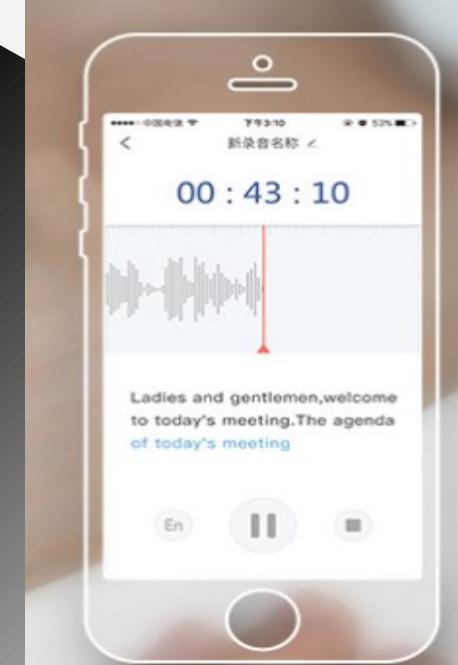
Speech Recognition



Smart speaker



Voice assistant

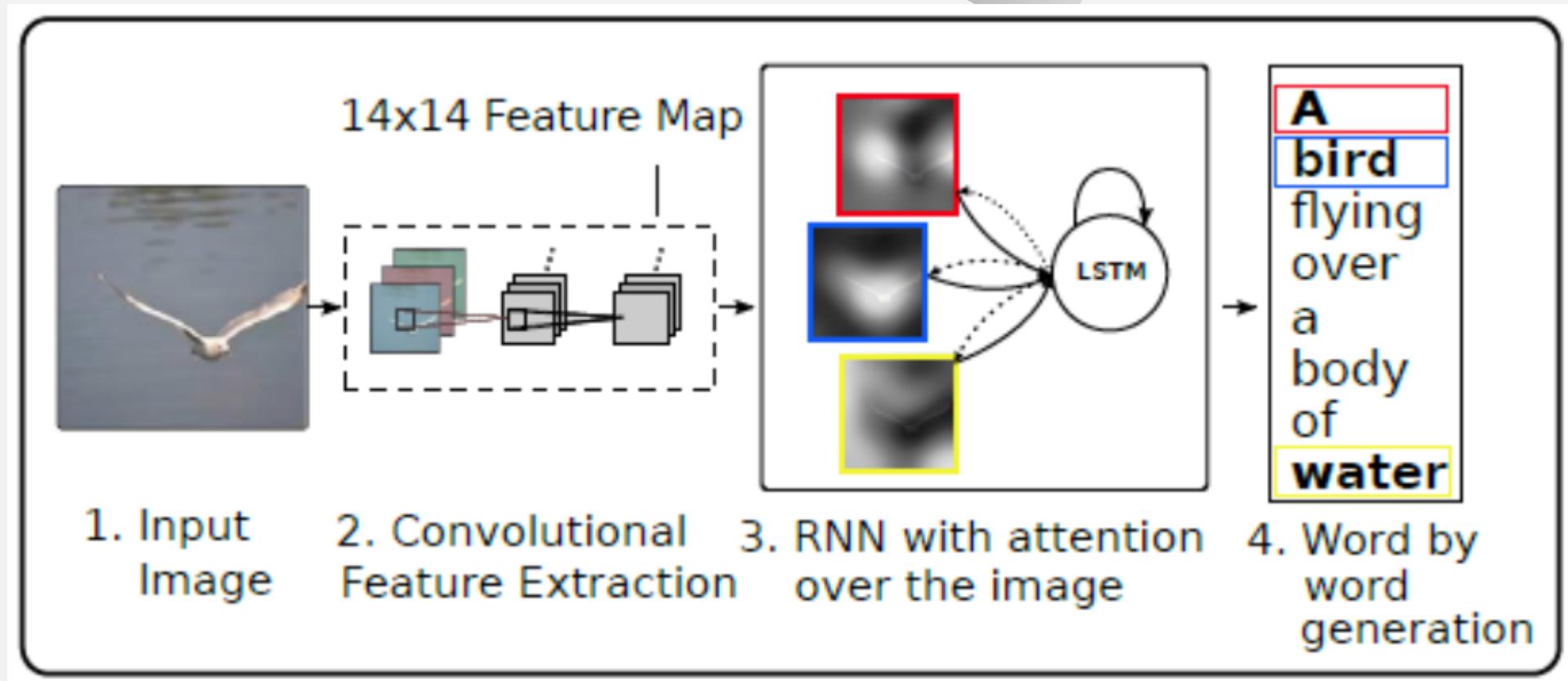


Real-time ASR

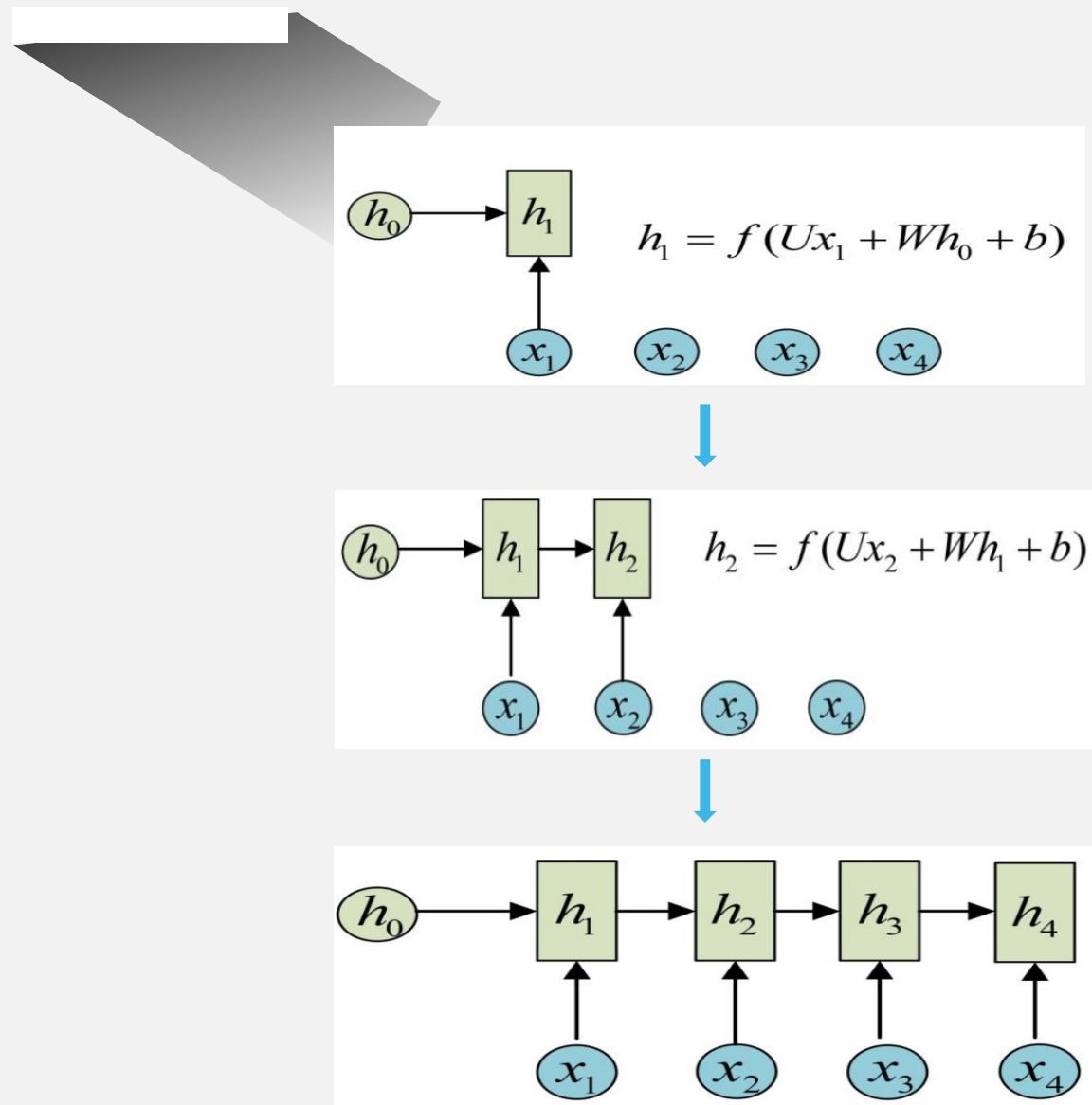
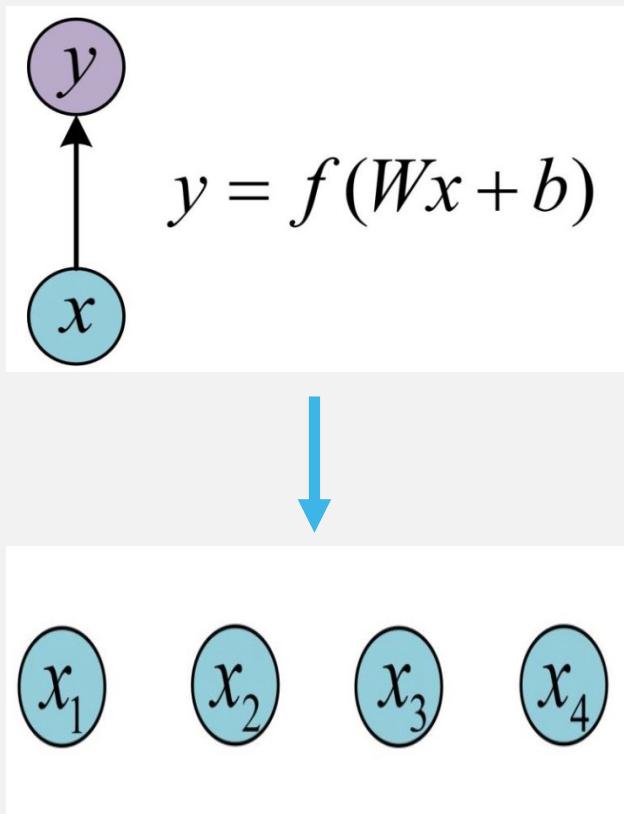


Realtime
navigation

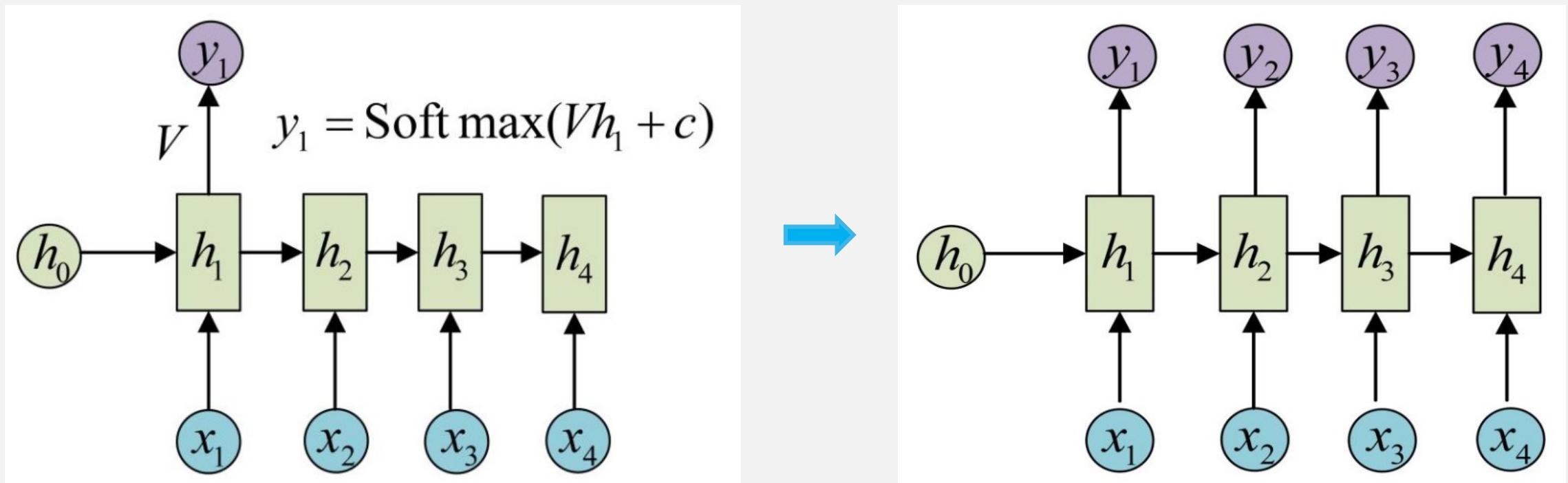
Image Caption



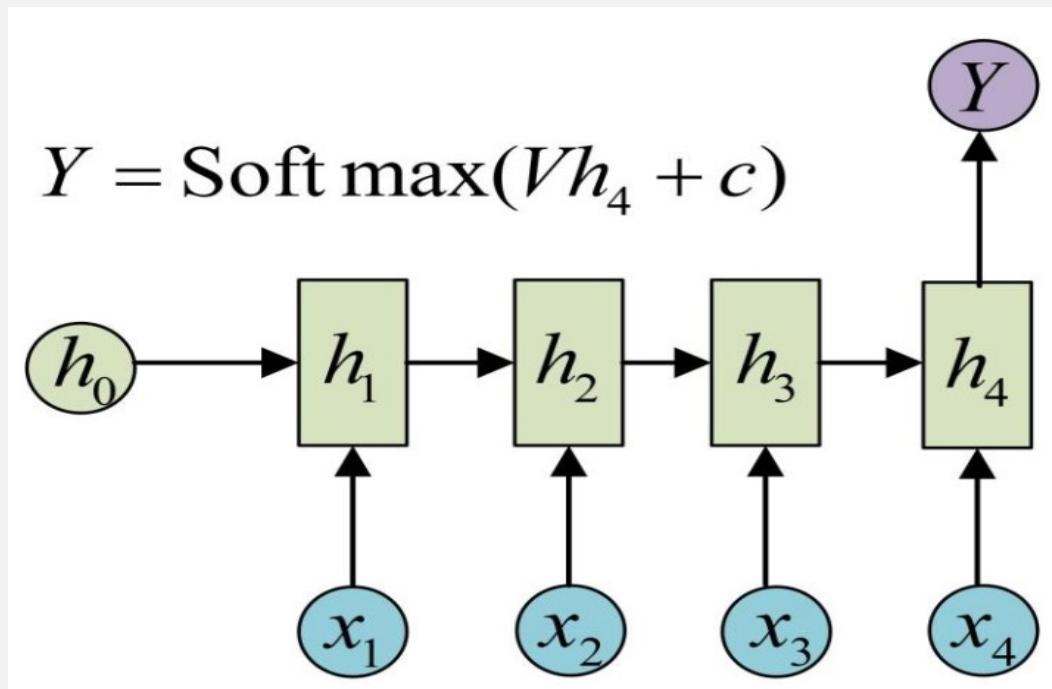
Classical RNN structure



Classical RNN structure



Some variants



N VS 1

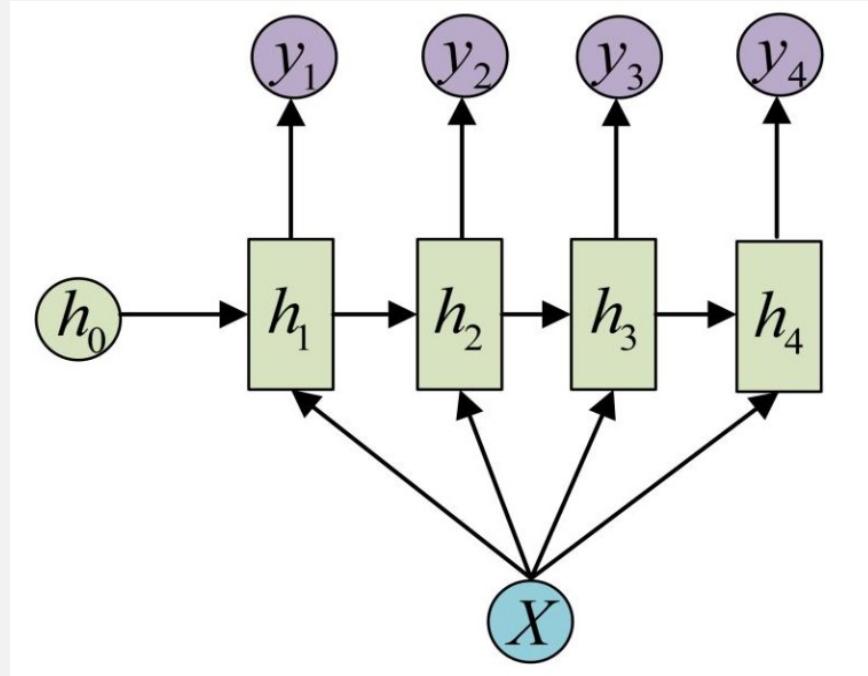
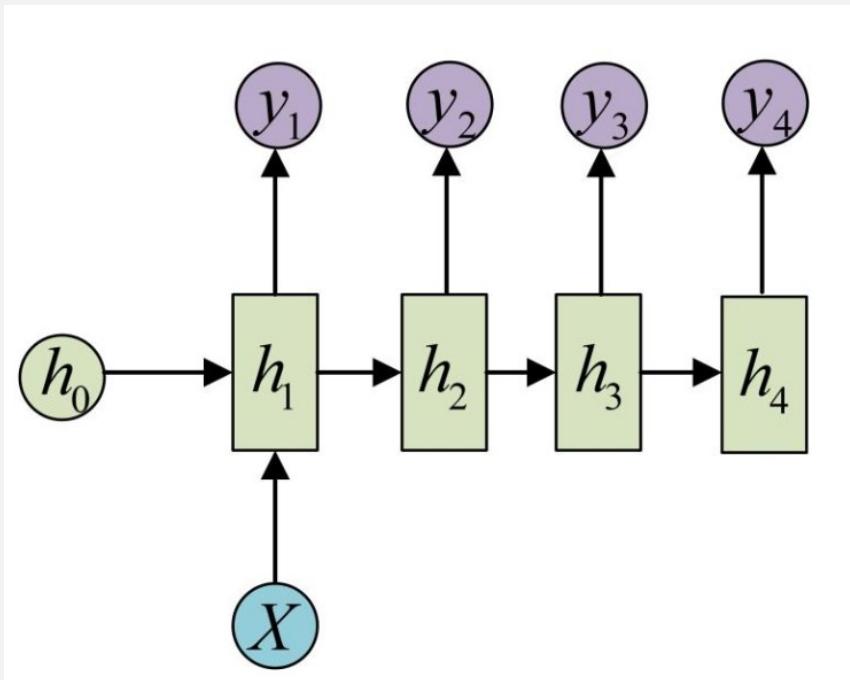


Text classification

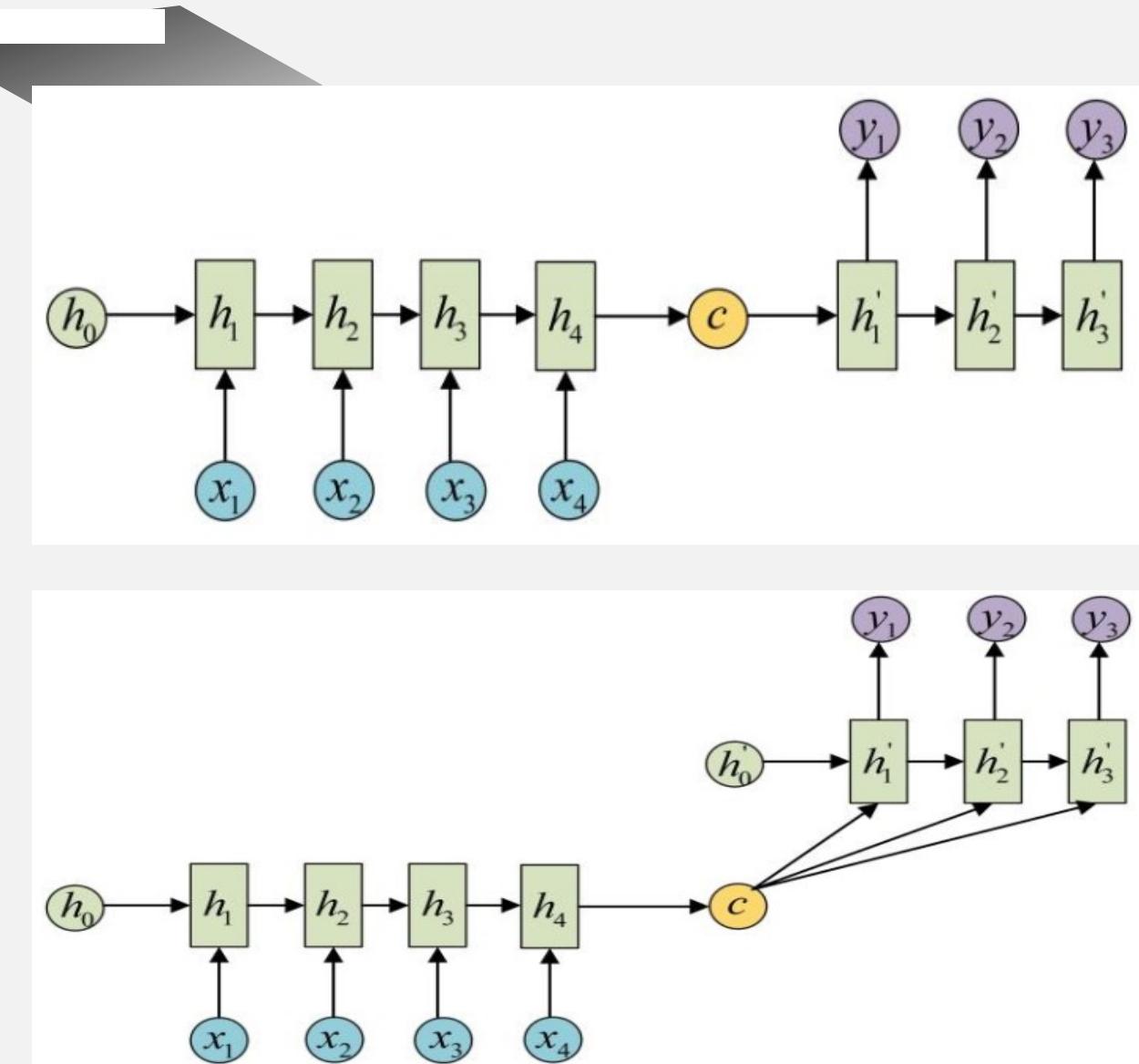
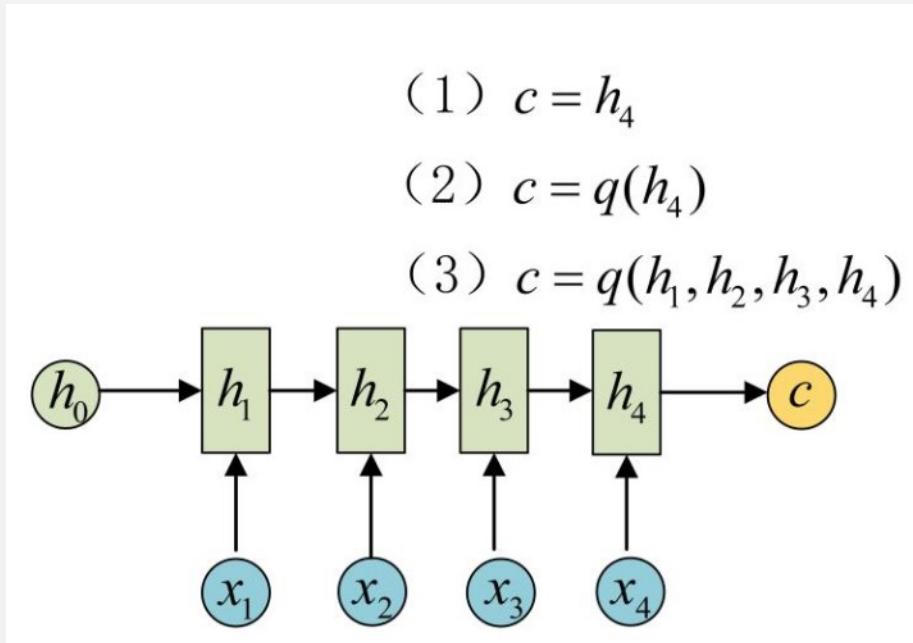
Video classification



1 VS N



Seq2Seq (N VS M)



Machine Comprehension

Sam walks into the kitchen.
Sam picks up an apple.
Sam walks into the bedroom.
Sam drops the apple.

Q: Where is the apple?

A. Bedroom

Speech Recognition



Auto Text Summarization

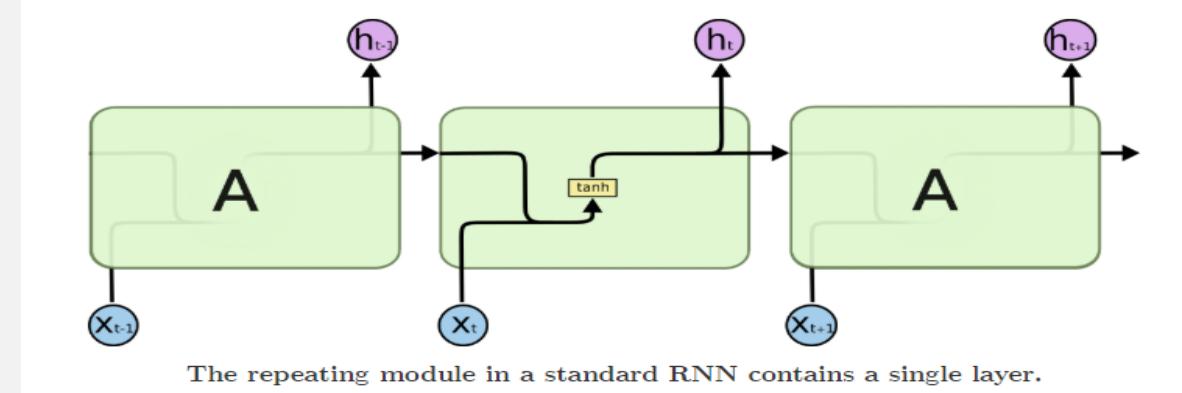
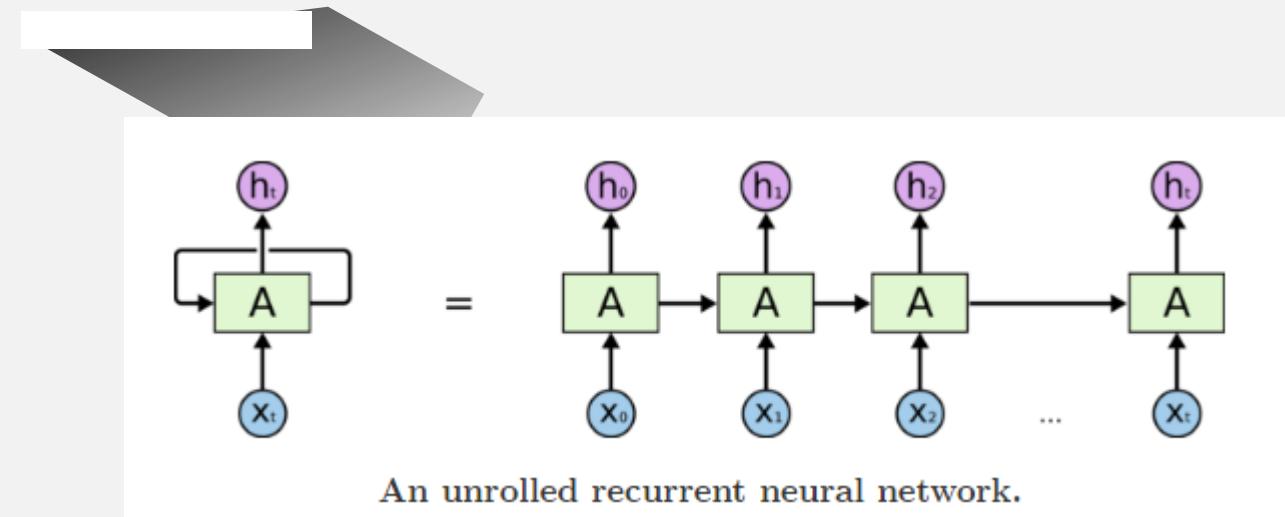
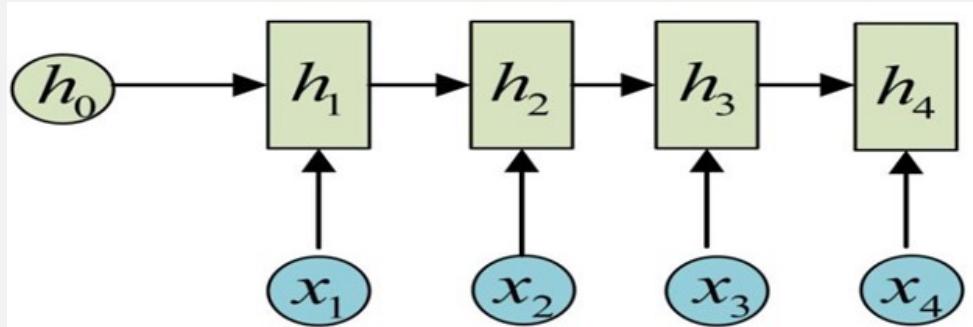
原文：

海湾报刊对美国新当选总统克林顿，能否帮助振兴中东和平进程感到怀疑，但也确实看到了一丝希望。

人类：海湾报刊对克林顿是否会恢复和平进程，持怀疑态度

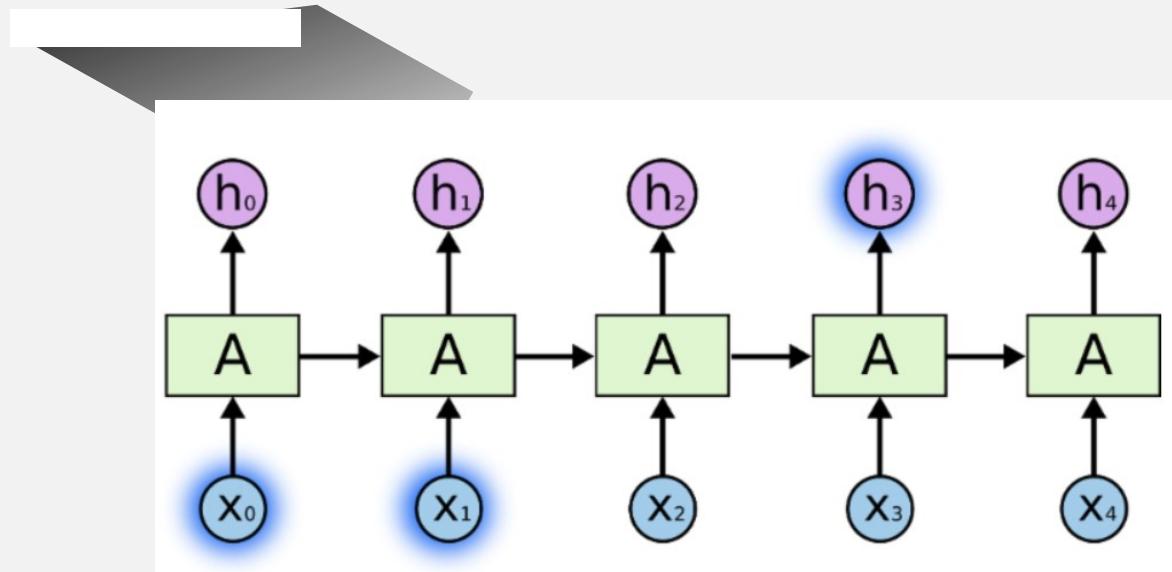
机器：海湾新闻界对克林顿恢复和平进程的前景，持怀疑态度

Understanding LSTM Networks

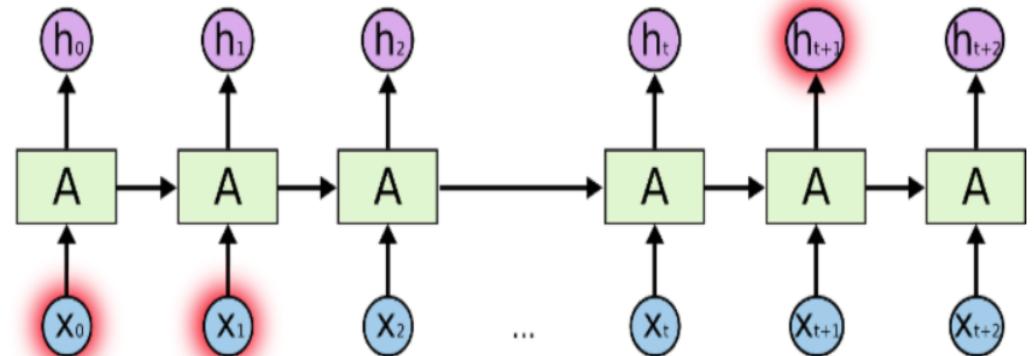


Understanding LSTM Networks

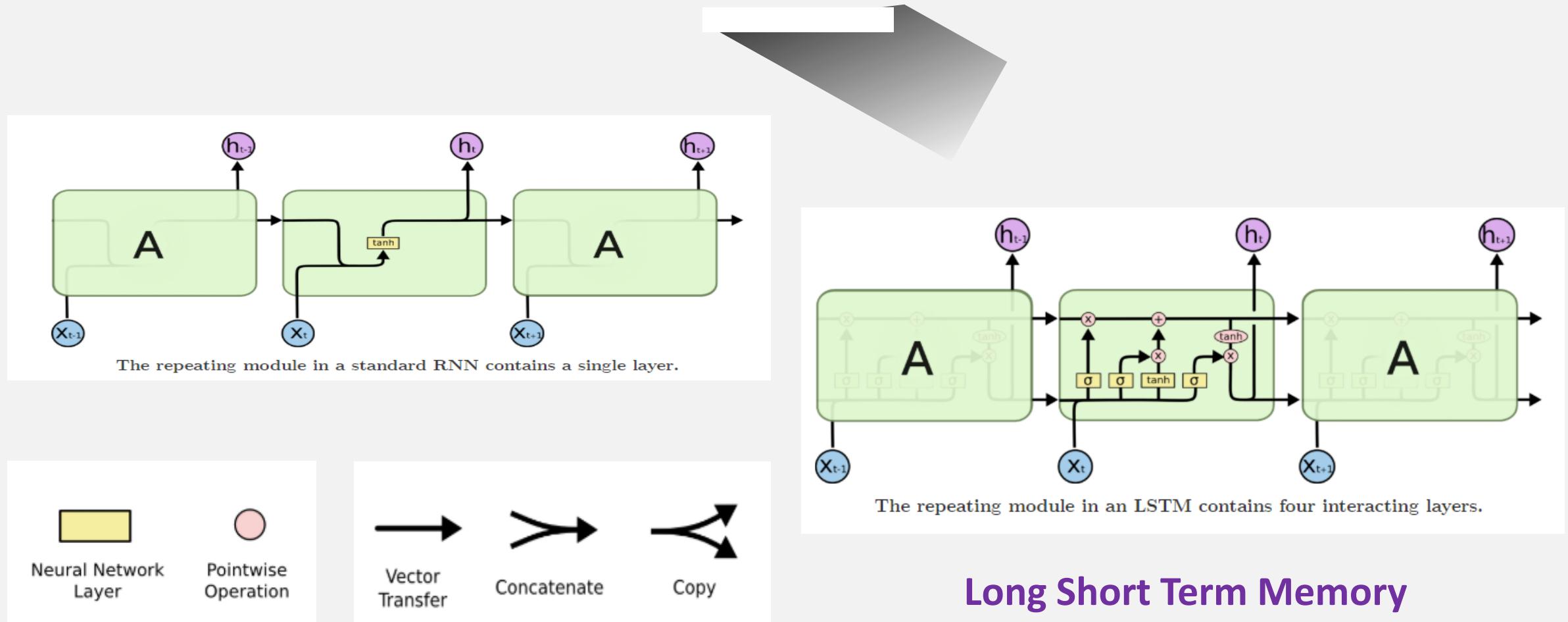
The clouds are in the **sky.**



I grew up in **France**... I speak fluent **French**.

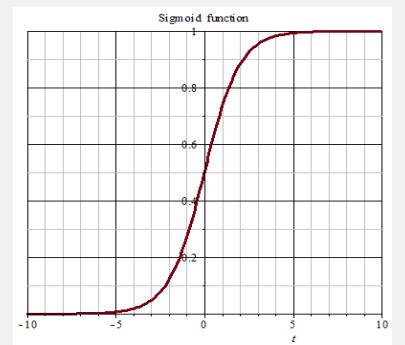
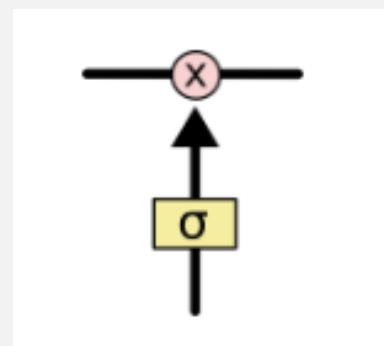
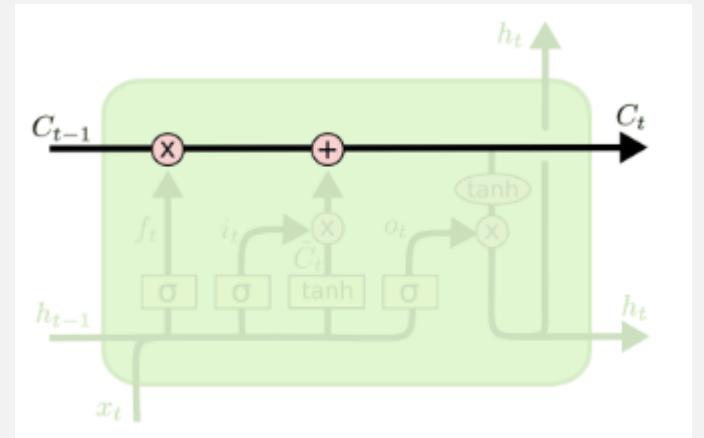
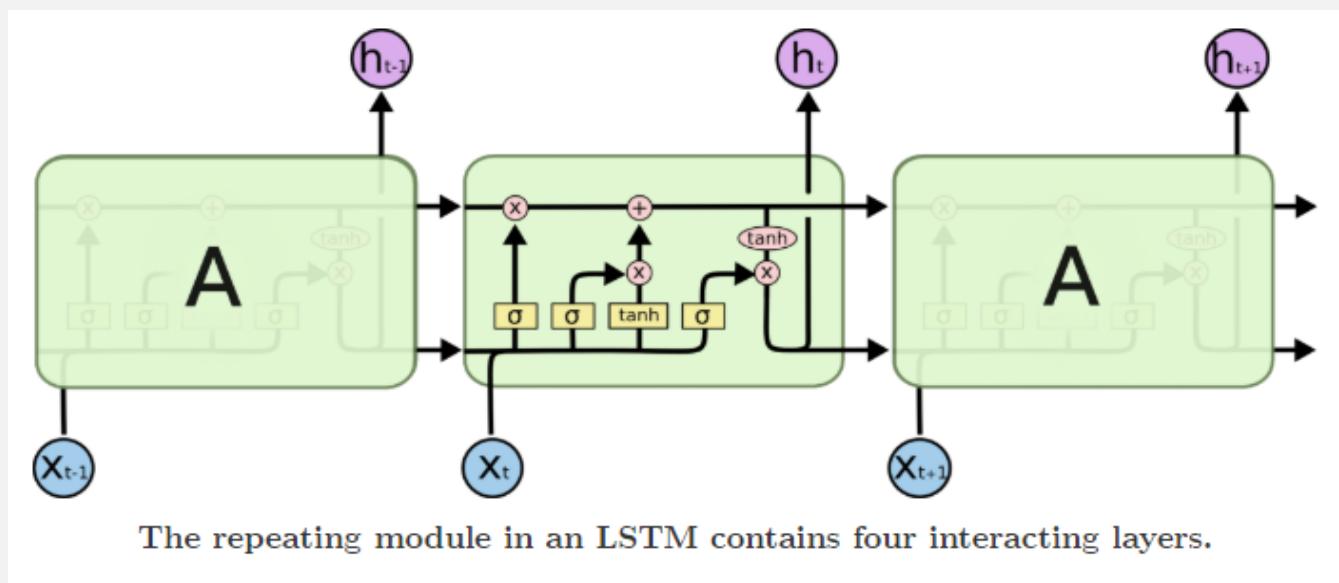


Understanding LSTM Networks

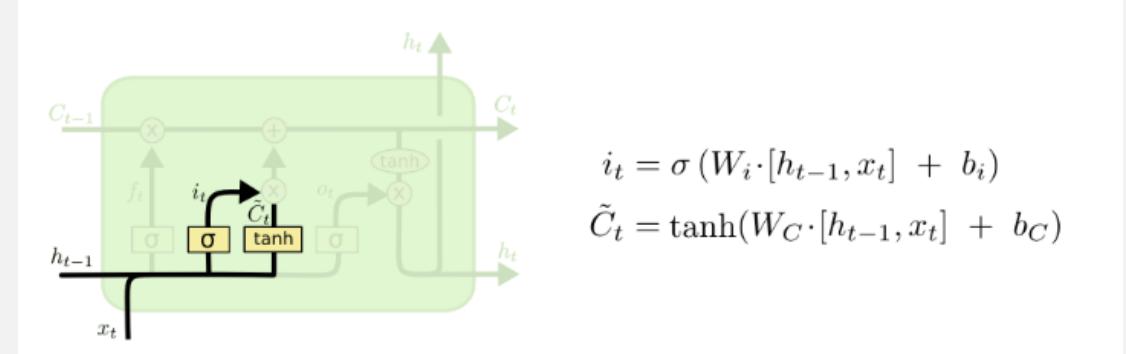
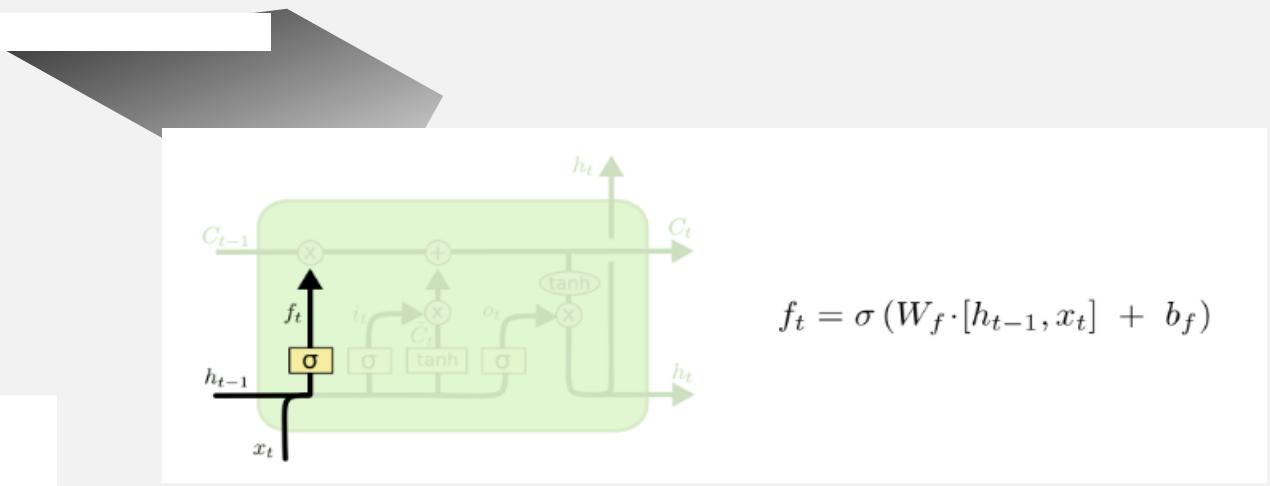
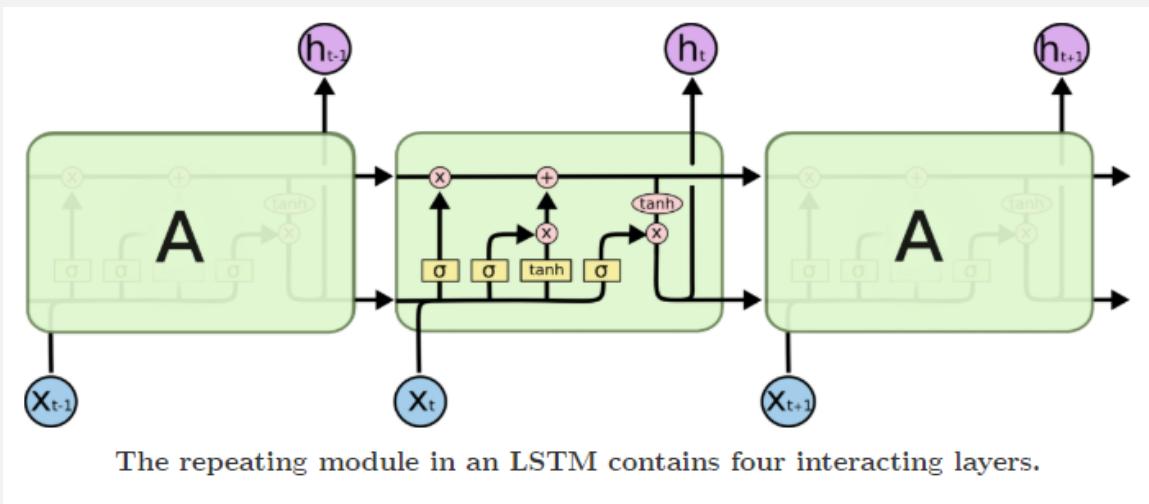
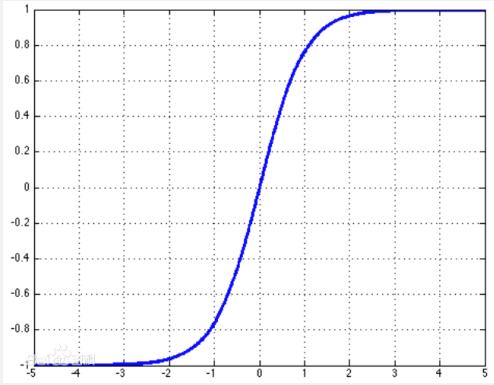


Long Short Term Memory

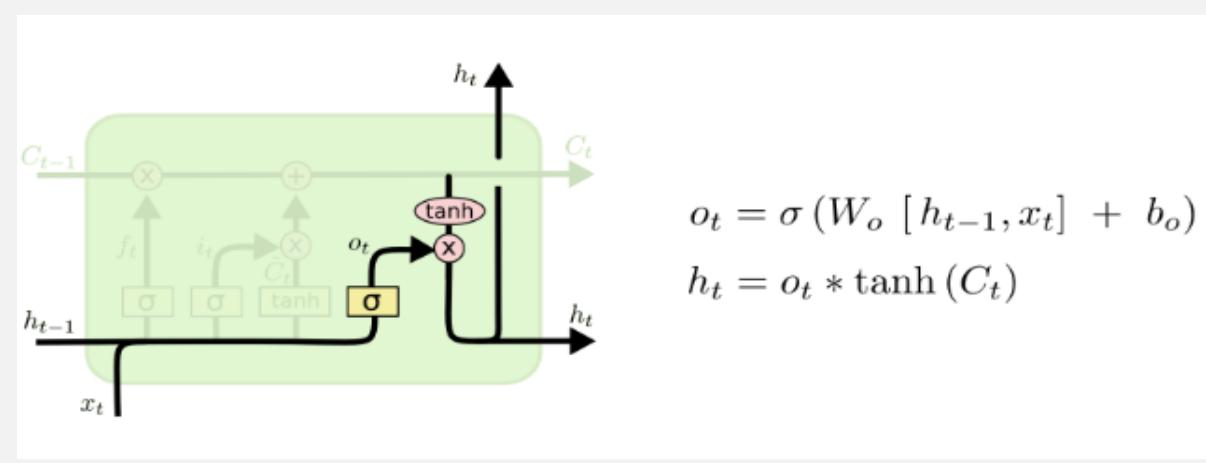
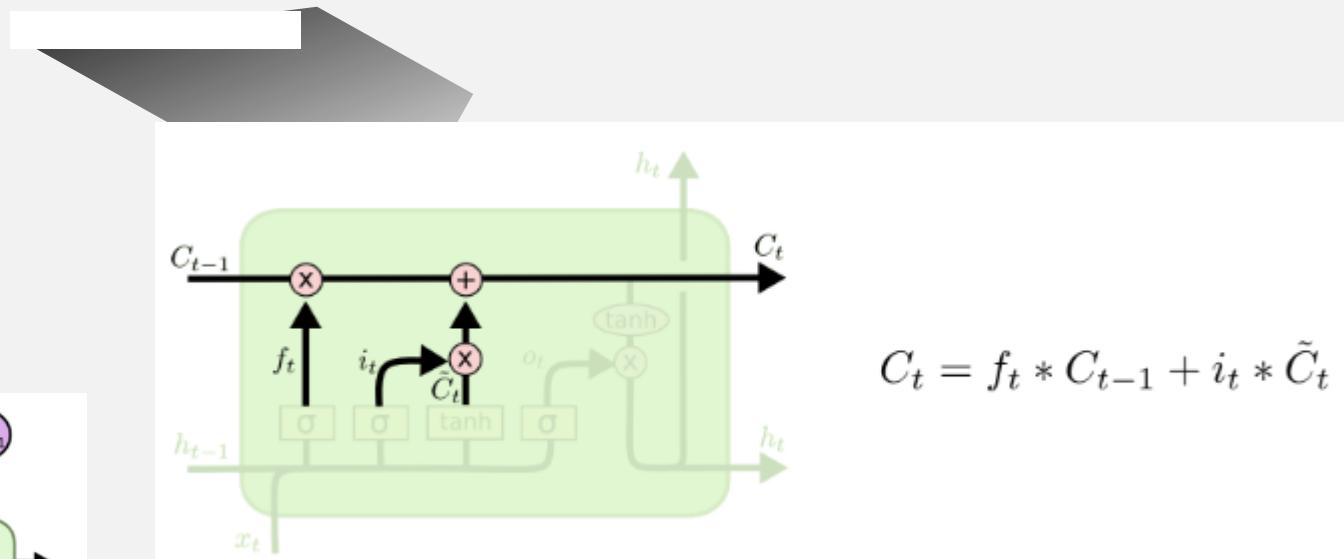
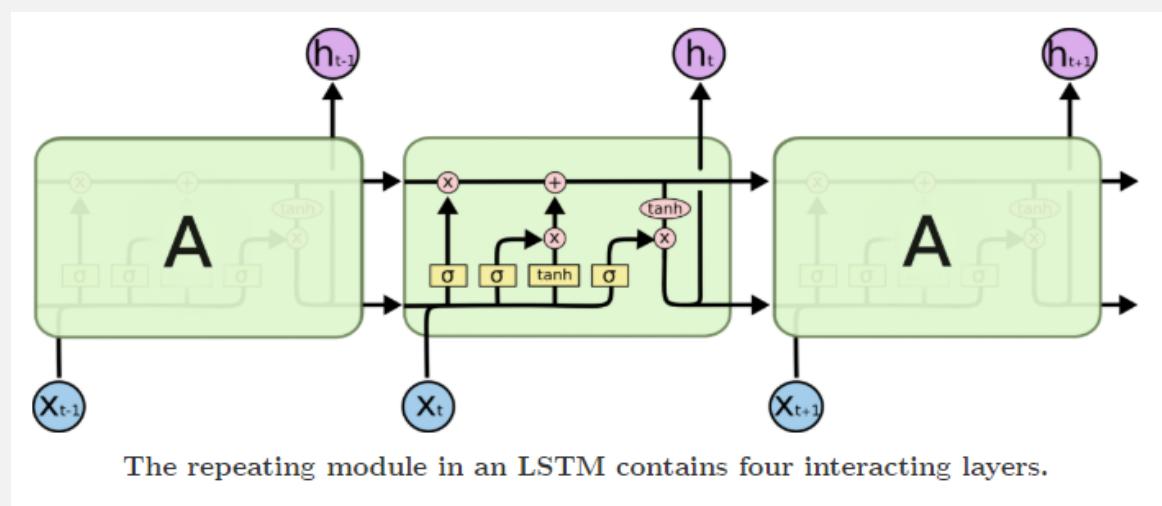
Understanding LSTM Networks



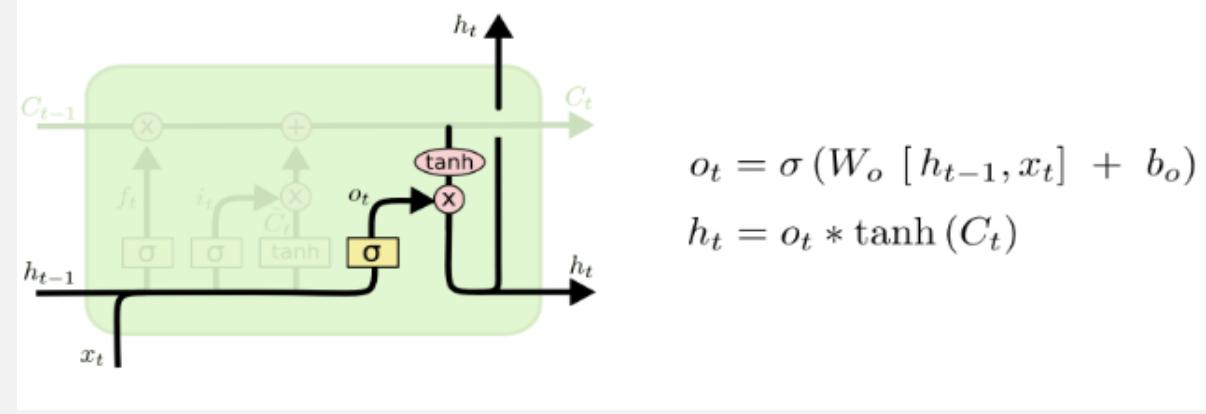
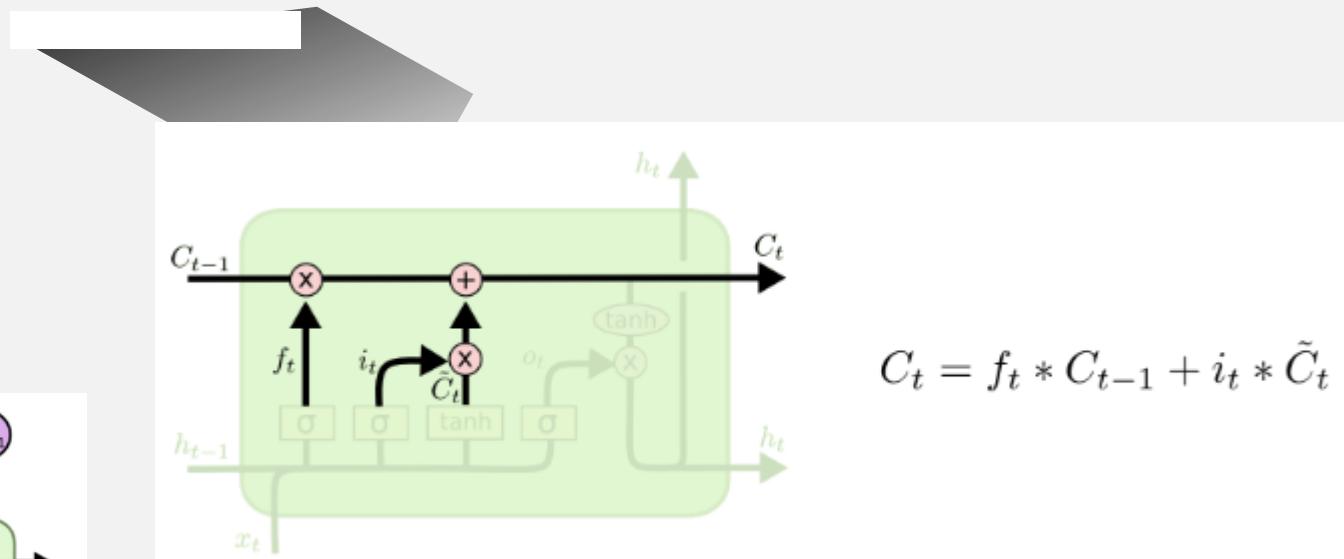
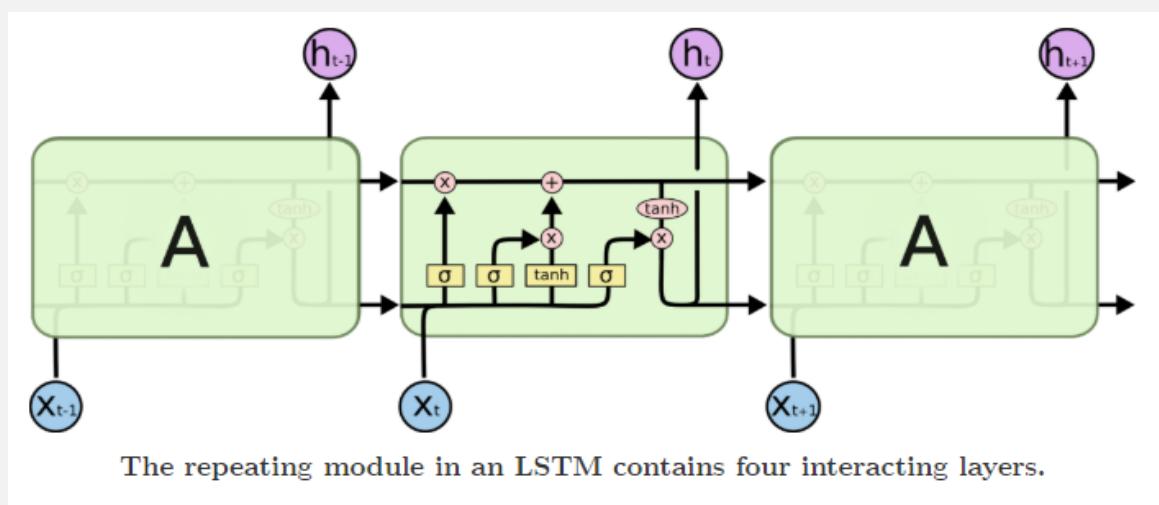
Understanding LSTM Networks



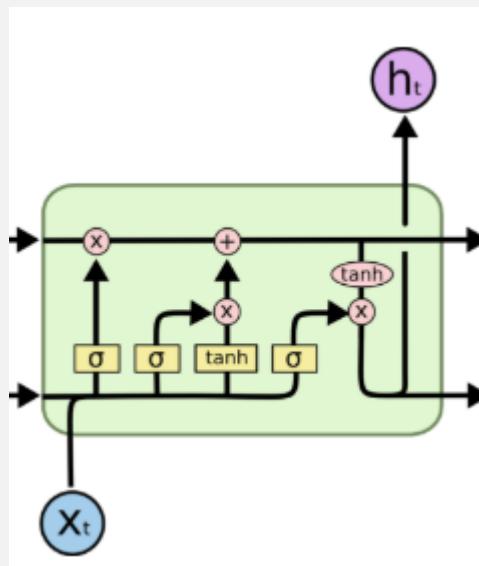
Understanding LSTM Networks



Understanding LSTM Networks

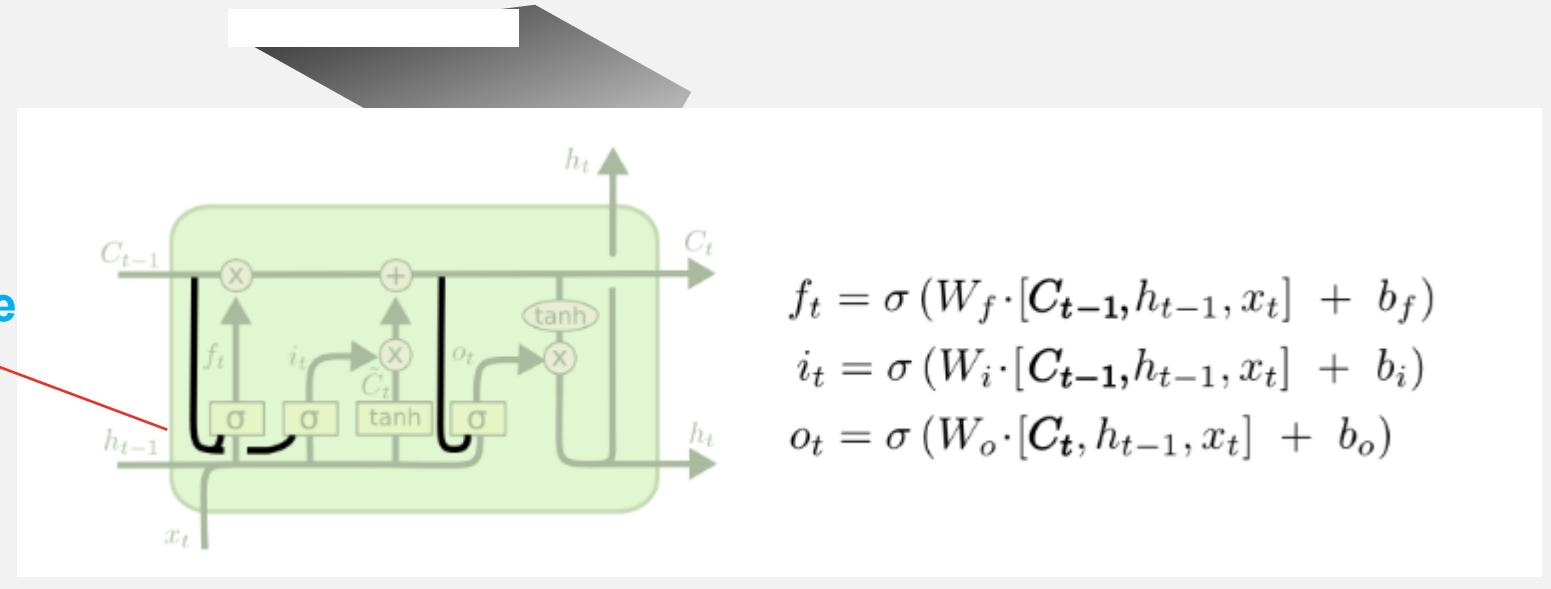


Variants on Long Short Term Memory

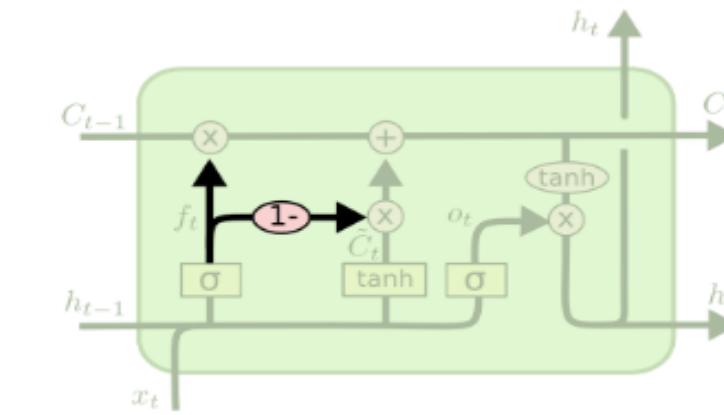


$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

peephole

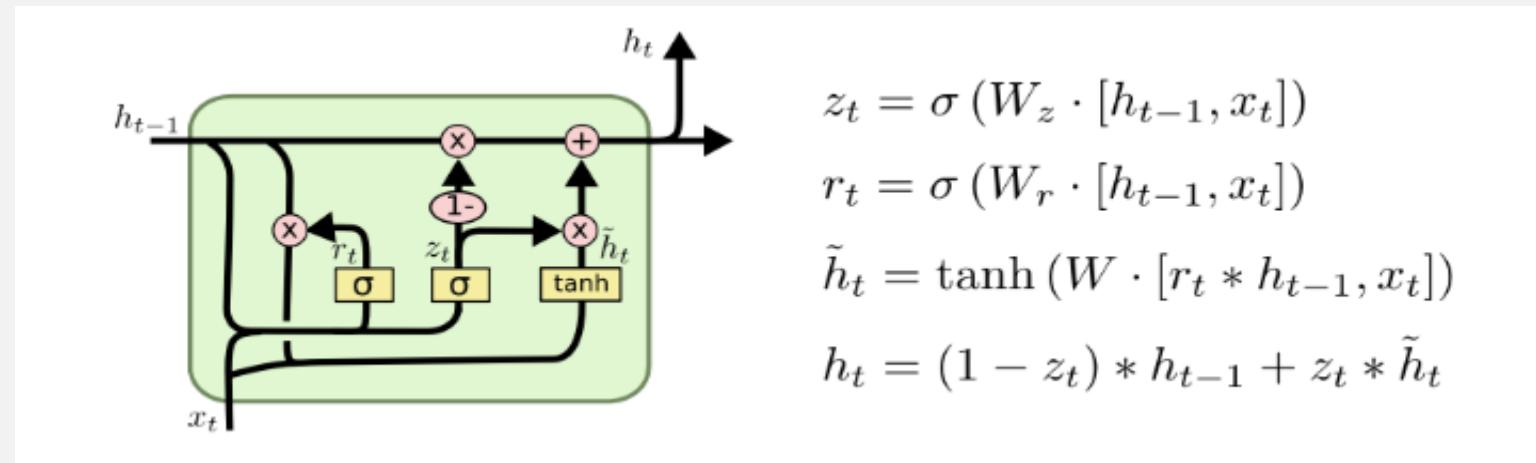
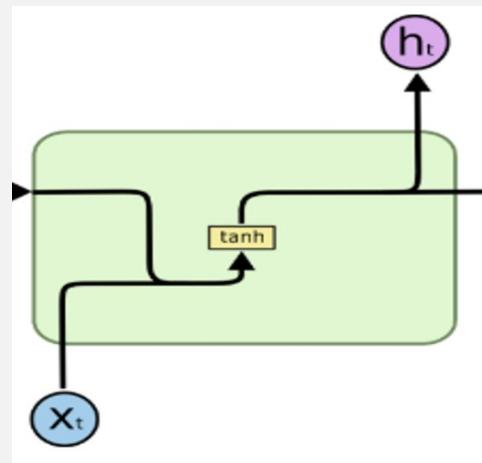
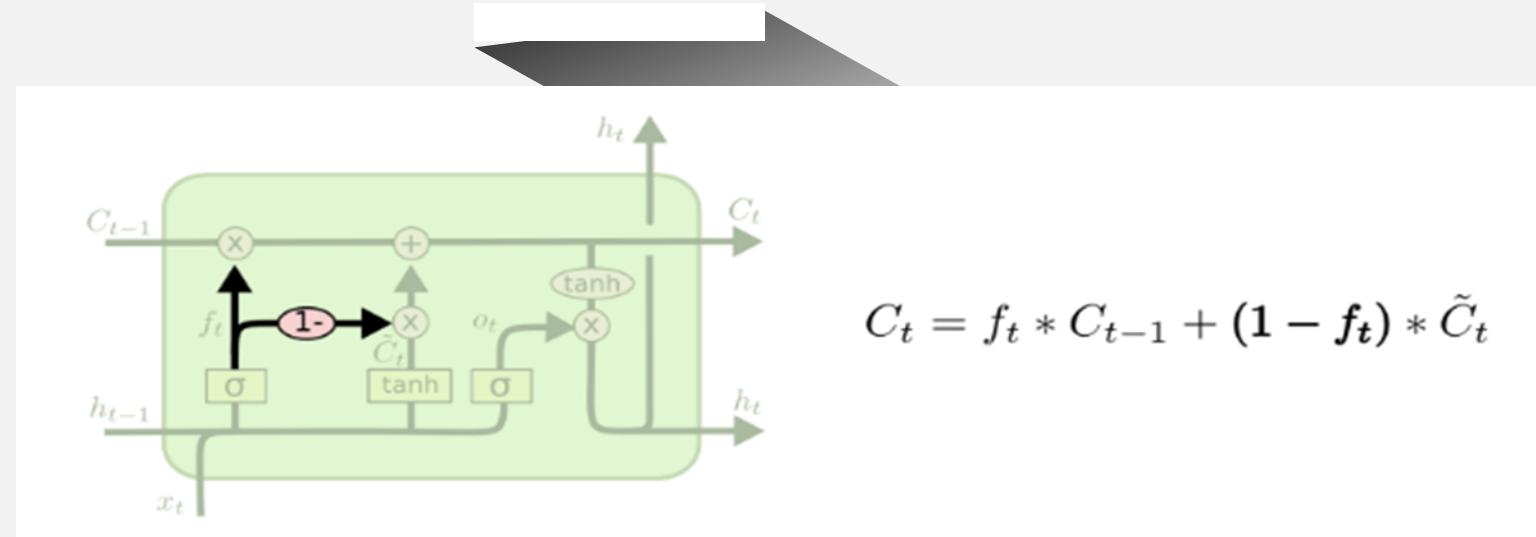
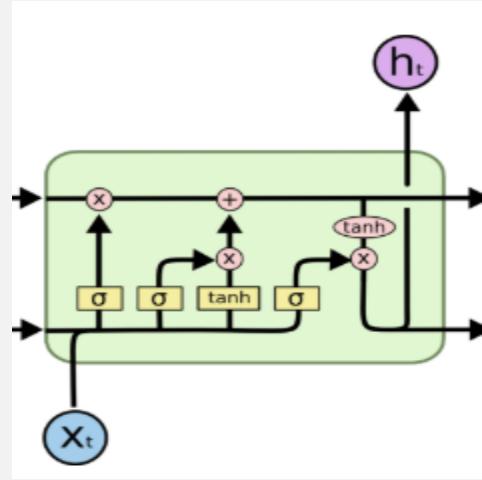


$$f_t = \sigma(W_f \cdot [C_{t-1}, h_{t-1}, x_t] + b_f)$$
$$i_t = \sigma(W_i \cdot [C_{t-1}, h_{t-1}, x_t] + b_i)$$
$$o_t = \sigma(W_o \cdot [C_t, h_{t-1}, x_t] + b_o)$$

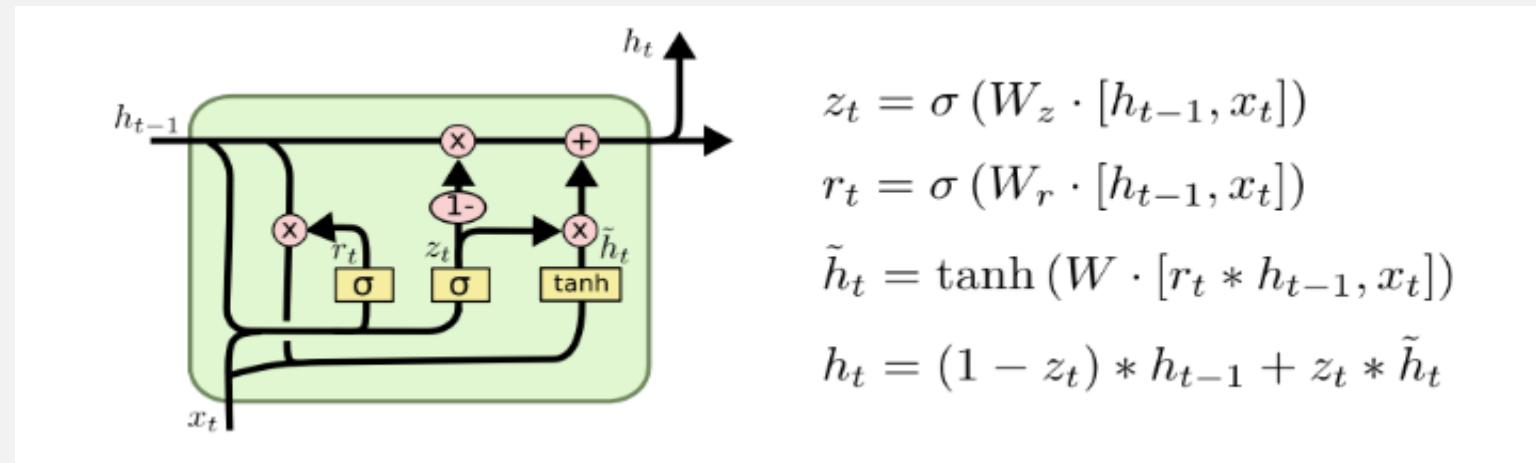
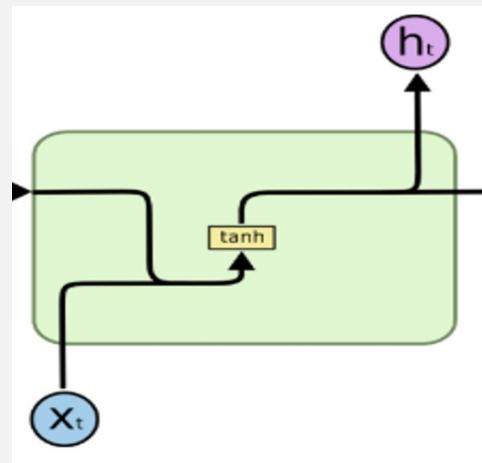
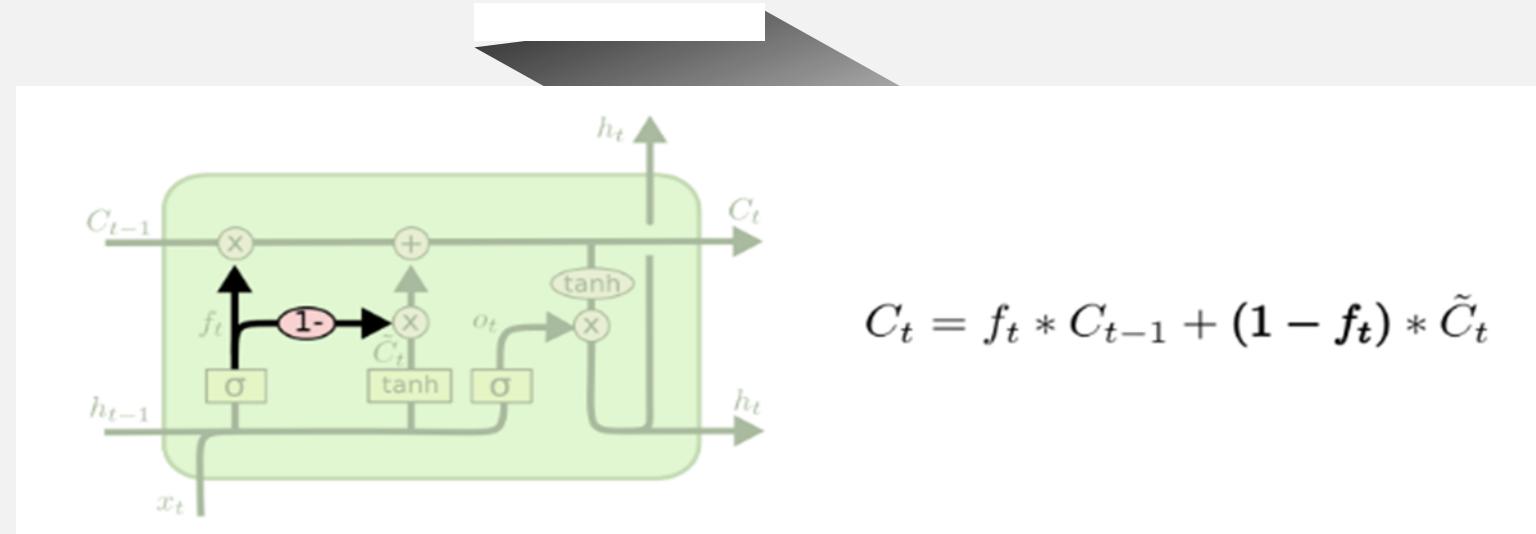
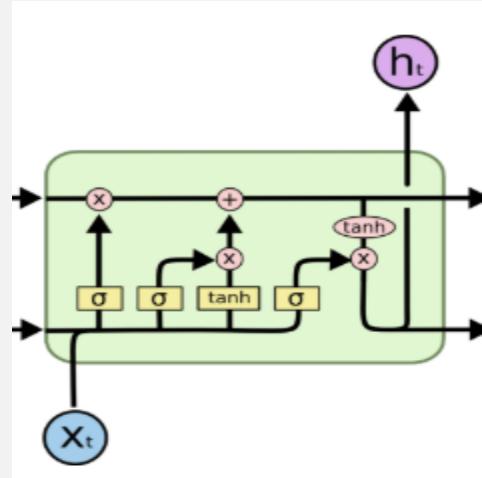


$$C_t = f_t * C_{t-1} + (1 - f_t) * \tilde{C}_t$$

GRU (Gated Recurrent Unit)



GRU (Gated Recurrent Unit)



Representing words

One-hot:

x:

Mike Guo are a good guy.

Vocabulary	Number
A	1
And	2
...	...
Guo	635
Pooter	636
..	...
Zulu	10000

Mike	Guo	are	guy
0	0	0	0
...	0	1	0
0
0	1	0	0
...	0	0	1
1
...	0	0	0

Representing words

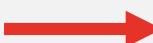
One-hot:

杭州 [0,0,0,0,0,0,1,0,....., 0,0,0,0,0,0]

上海 [0,0,0,1,0,0,0,0,....., 0,0,0,0,0,0]

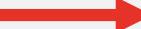
宁波 [0,0,0,1,0,0,0,0,0,....., 0,0,0,0,0,0]

北京 [0,0,0,0,0,0,0,0,....., 1,0,0,0,0,0]



The Curse of Dimensionality in classification



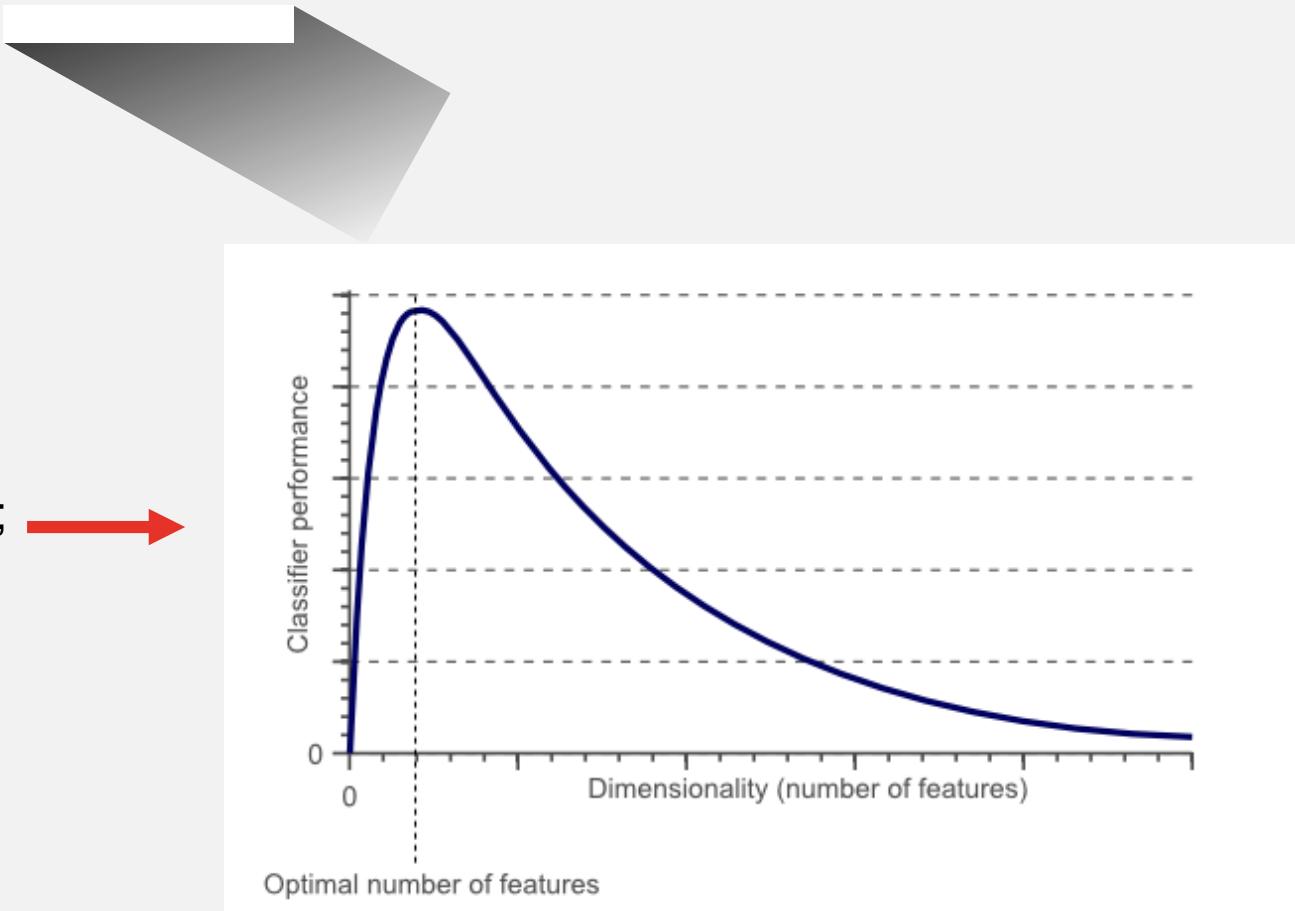
If $0.5*\text{red} + 0.3*\text{green} + 0.2*\text{blue} > 0.6$: return cat; 
else return dog;

More features:

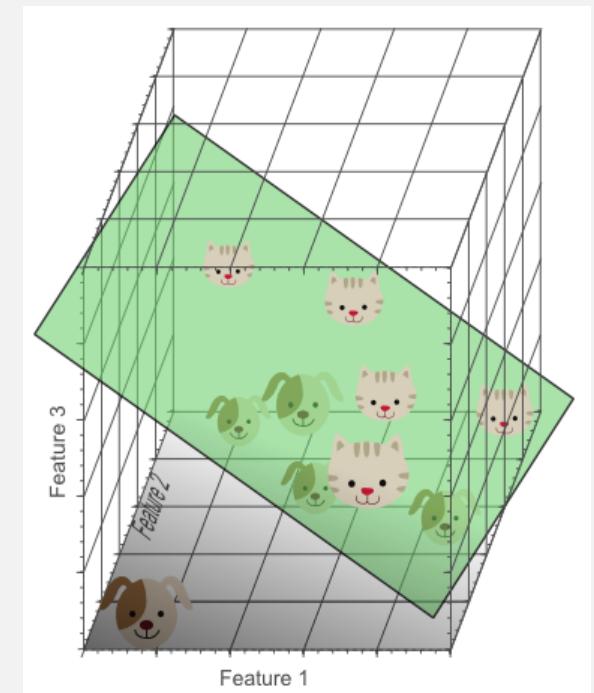
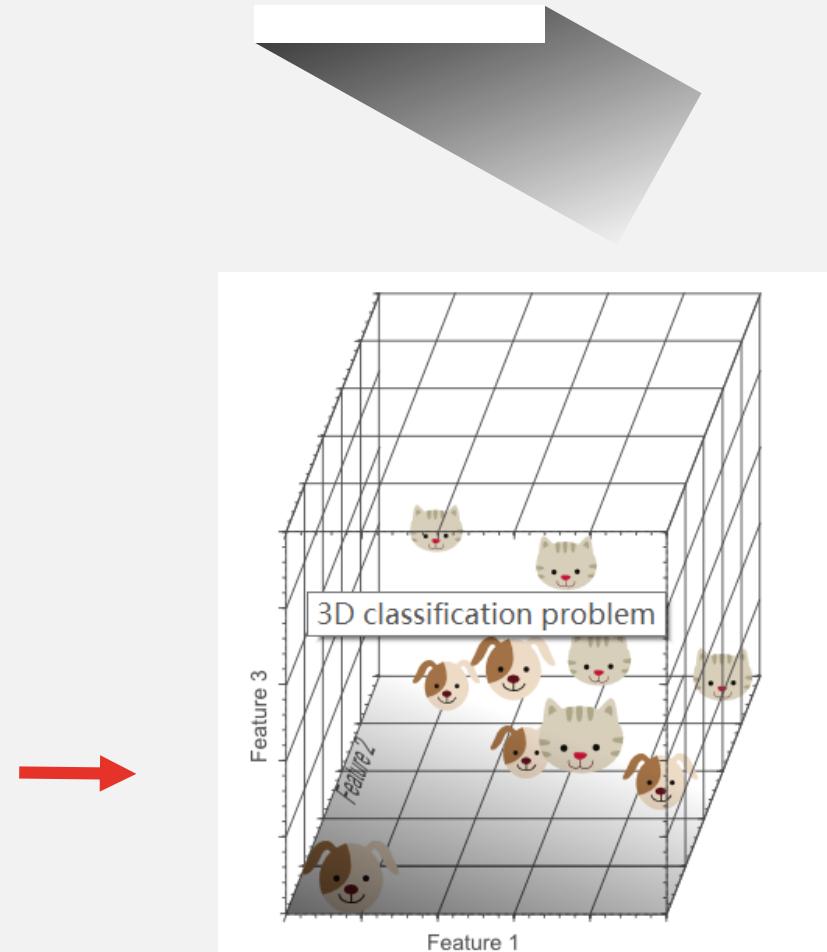
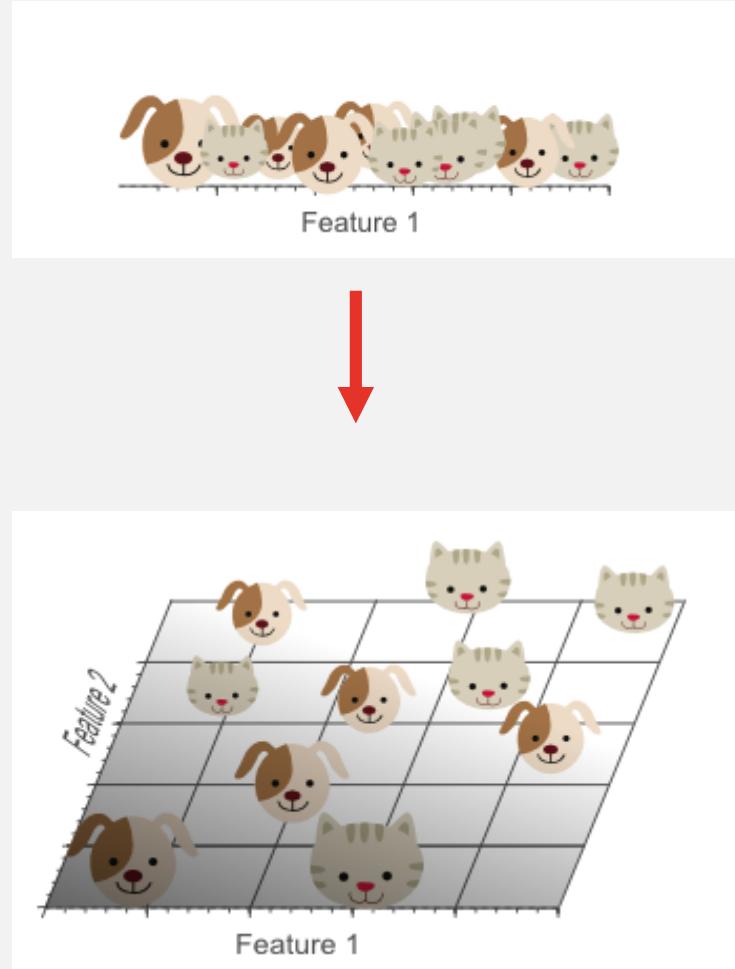
dx

dy

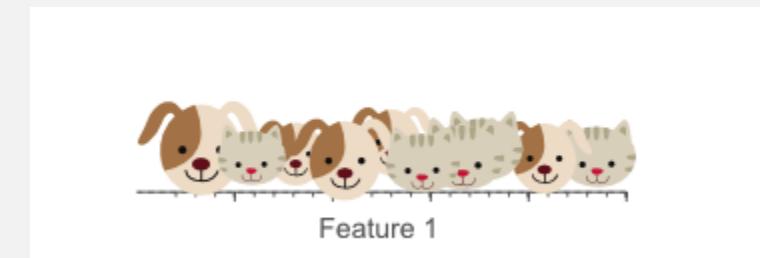
...



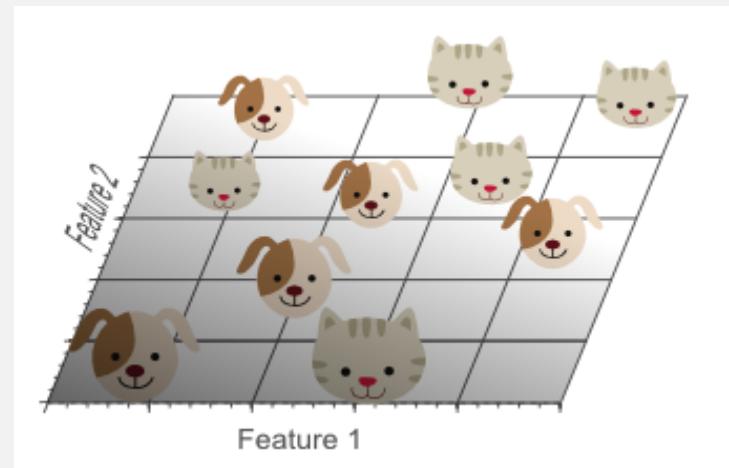
The Curse of Dimensionality in classification



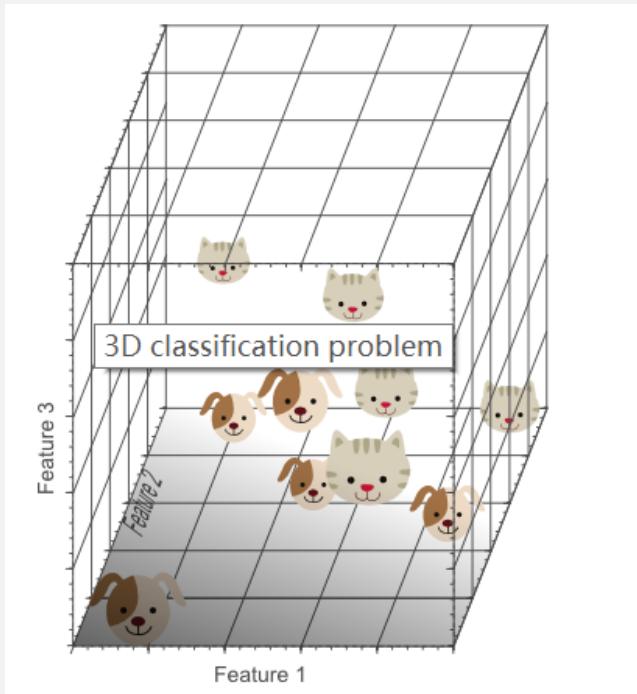
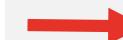
The Curse of Dimensionality in classification



$$10 \div 5 = 2$$

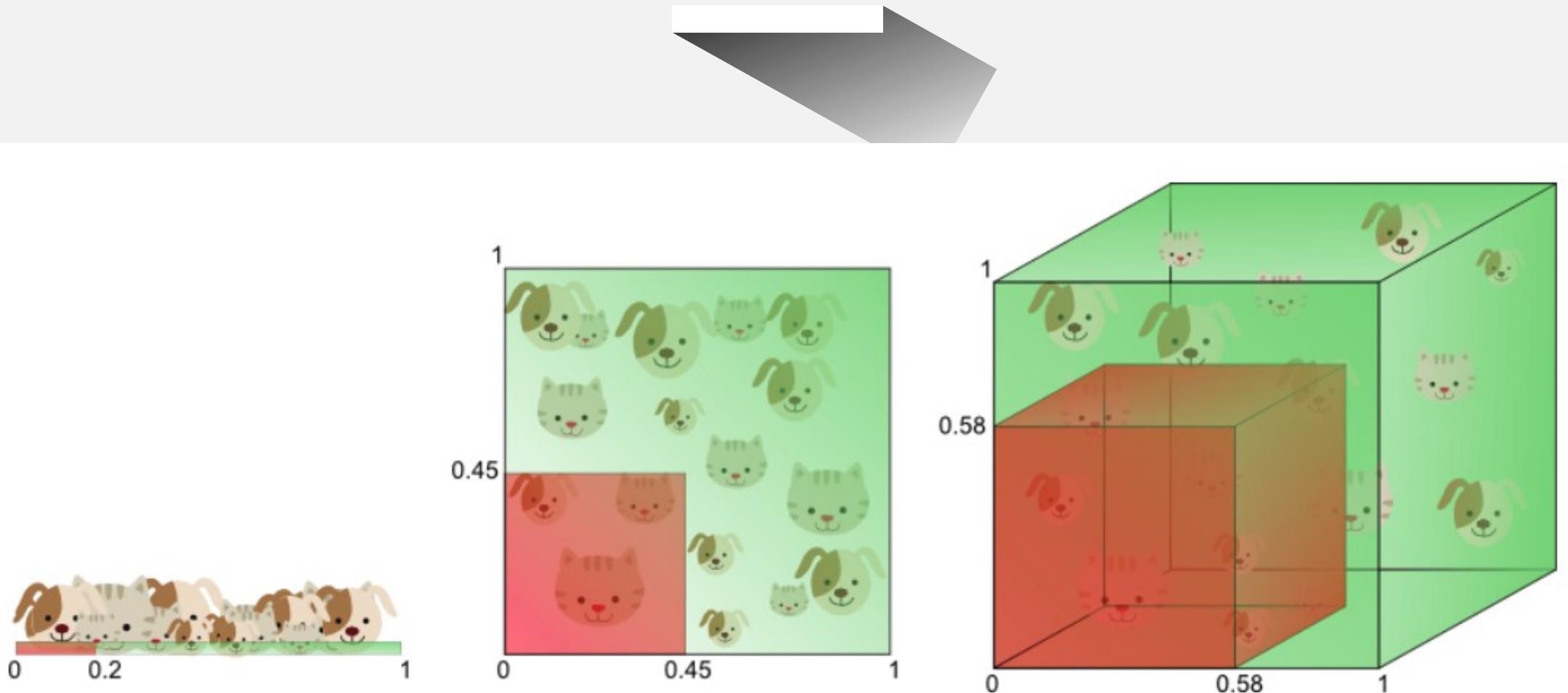


$$10 \div 25 = 0.4$$



$$10 \div 125 = 0.08$$

The Curse of Dimensionality in classification



The amount of training data needed to cover 20% of the feature range grows exponentially with the number of dimensions.

Representing words

Distributed representation



Representing words



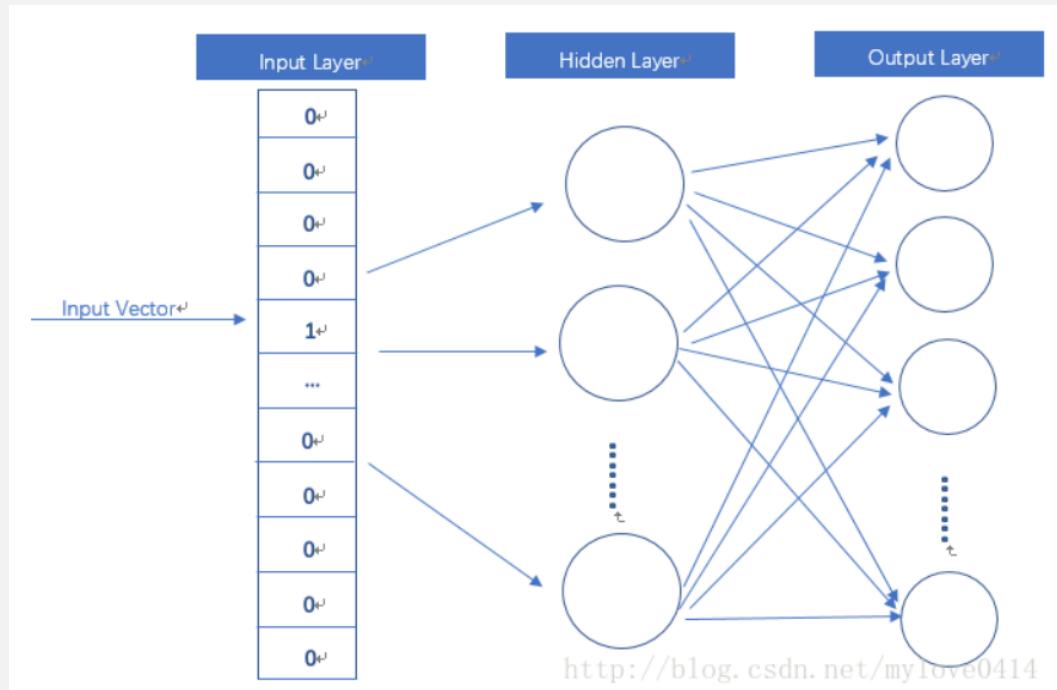
Distributed representation



$$\vec{\text{King}} - \vec{\text{Man}} + \vec{\text{Woman}} = \vec{\text{Queen}}$$

Representing words

Word2Vec, Embeddings



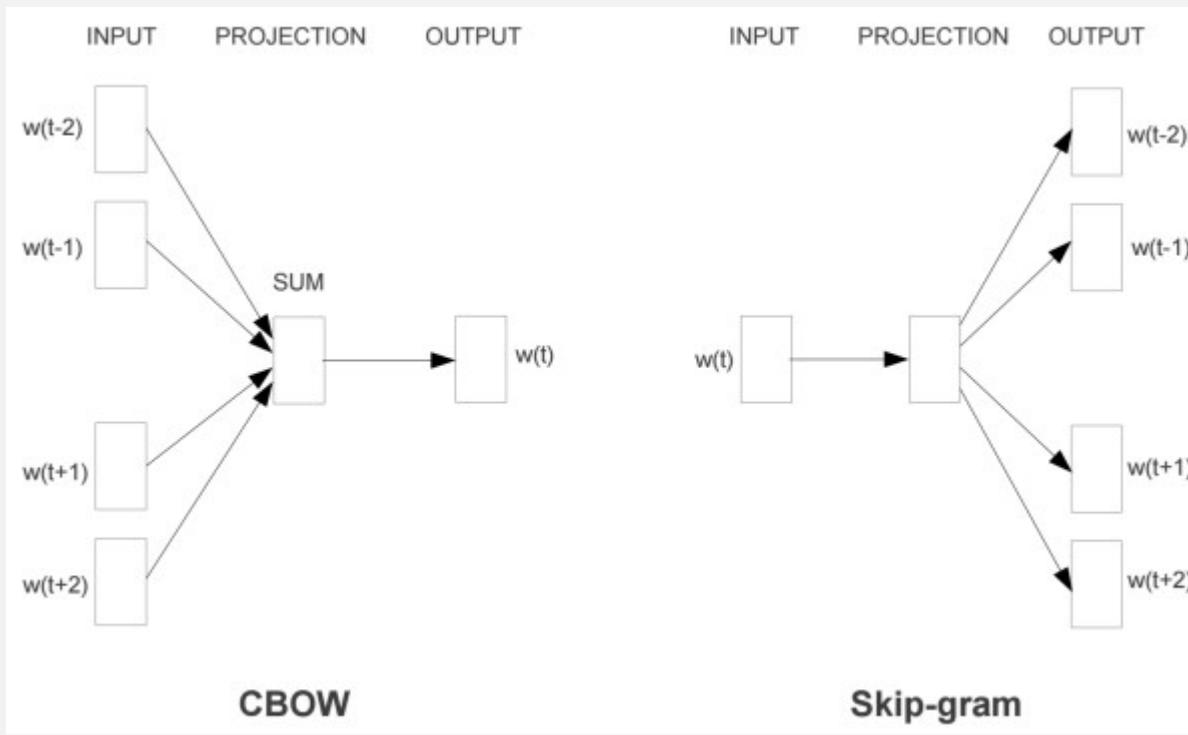
A matrix multiplication diagram illustrating the calculation of a word embedding vector. The input vector [0 0 0 1 0] is multiplied by a weight matrix to produce the output vector [10 12 19].

$$\begin{bmatrix} 0 & 0 & 0 & \textcolor{green}{1} & 0 \end{bmatrix} \times \begin{bmatrix} 17 & 24 & 1 \\ 23 & 5 & 7 \\ 4 & 6 & 13 \\ \textcolor{green}{10} & \textcolor{green}{12} & \textcolor{green}{19} \\ 11 & 18 & 25 \end{bmatrix} = \begin{bmatrix} 10 & 12 & 19 \end{bmatrix}$$

http://blog.csdn.net/mylove0414

Representing words

CBOW & Skip-gram



$$\begin{bmatrix} 0 & 0 & 0 & \textcolor{green}{1} & 0 \end{bmatrix} \times \begin{bmatrix} 17 & 24 & 1 \\ 23 & 5 & 7 \\ 4 & 6 & 13 \\ \textcolor{green}{10} & \textcolor{green}{12} & \textcolor{green}{19} \\ 11 & 18 & 25 \end{bmatrix} = [10 \ 12 \ 19]$$

http://blog.csdn.net/mylove0414

CBOW: Continuous Bag-of-Words Model

Skip-gram: Continuous Skip-gram Model

Representing words



Skip-gram

Input word
The **dog** barked at the mailman

Skip window=2 → ['The', 'dog', 'barked', 'at']

Num skips=2

('dog', 'barked'), ('dog', 'the')

Source Text	Training Samples
The quick brown fox jumps over the lazy dog. →	(the, quick) (the, brown)
The quick brown fox jumps over the lazy dog. →	(quick, the) (quick, brown) (quick, fox)
The quick brown fox jumps over the lazy dog. →	(brown, the) (brown, quick) (brown, fox) (brown, jumps)
The quick brown fox jumps over the lazy dog. →	(fox, quick) (fox, brown) (fox, jumps) (fox, over)

Representing words



Skip-gram

The quick brown fox jumps over lazy dog

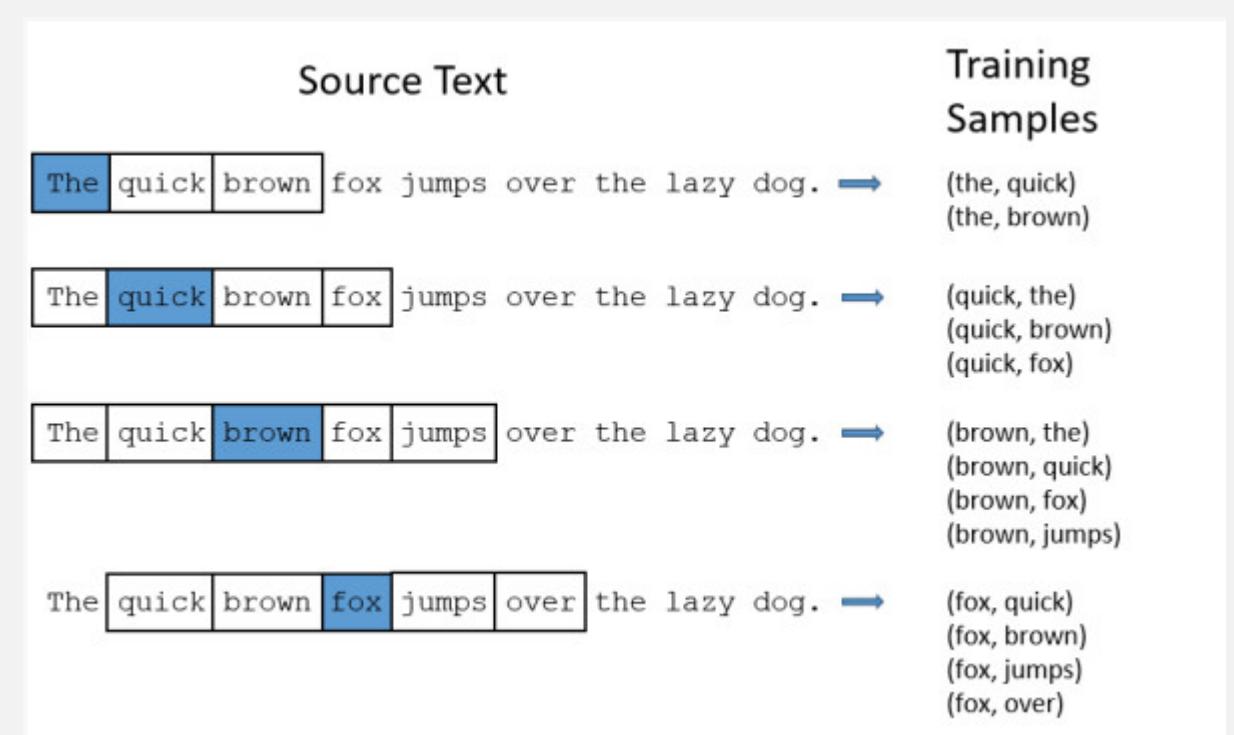
Skip window=2

kangaroo

Soviet
Union

Russia

watermelon

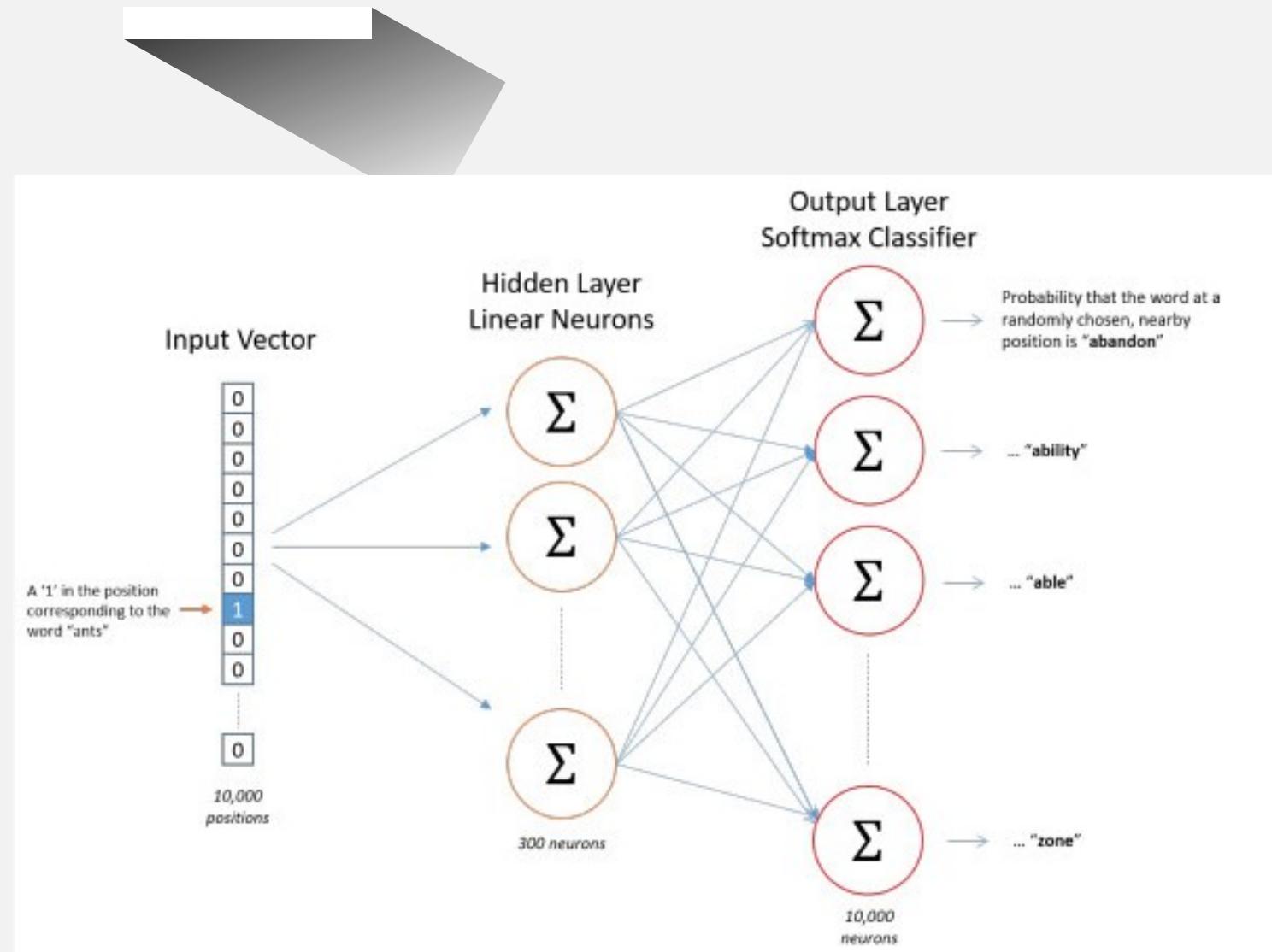


Representing words

Skip-gram

The dog barked at the mailman

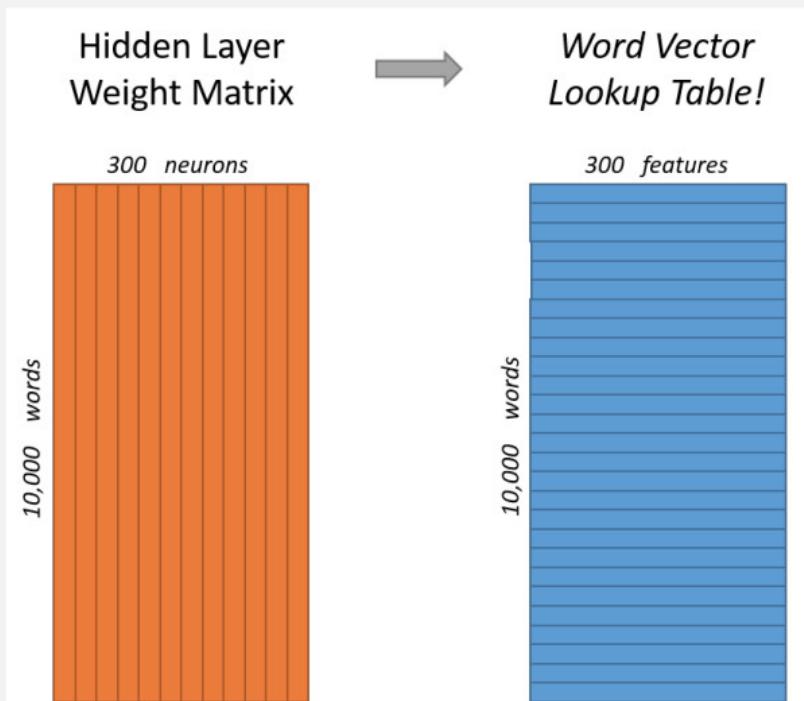
[0,1,0,0,0]



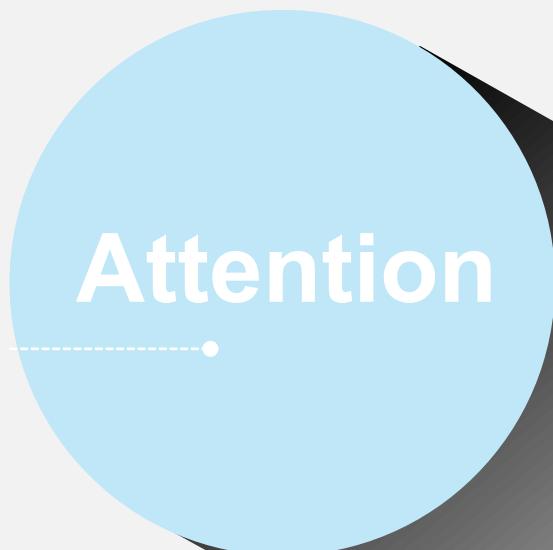
Representing words



Skip-gram

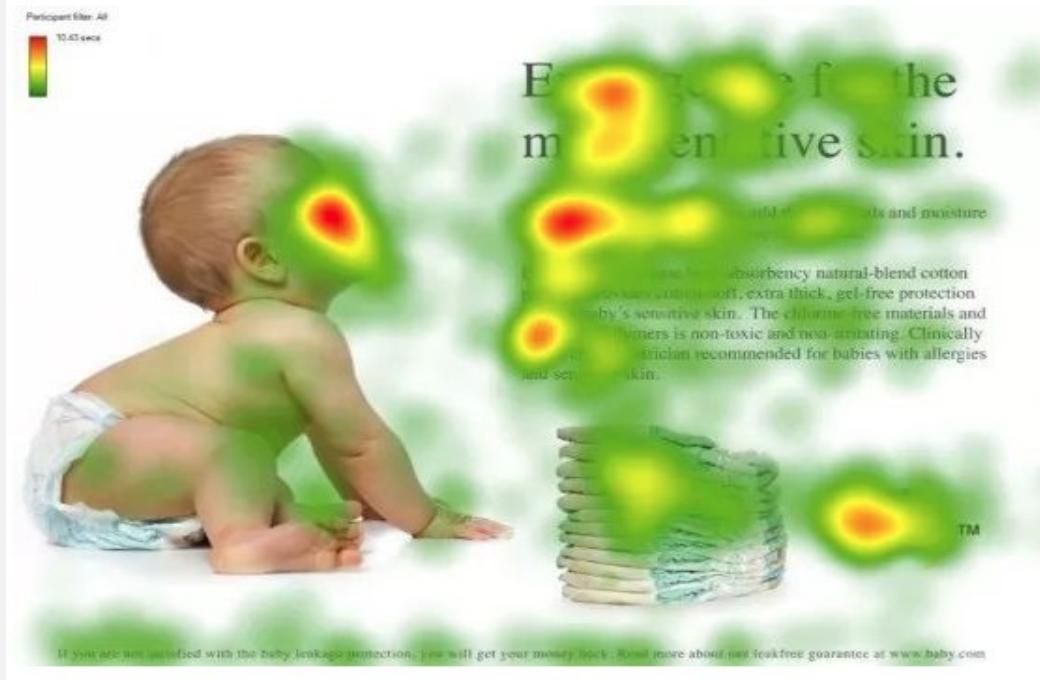


$$\begin{bmatrix} 0 & 0 & 0 & \boxed{1} & 0 \end{bmatrix} \times \begin{bmatrix} 17 & 24 & 1 \\ 23 & 5 & 7 \\ 4 & 6 & 13 \\ \boxed{10} & 12 & 19 \\ 11 & 18 & 25 \end{bmatrix} = [10 \quad 12 \quad 19]$$



- 01 **What are the attention mechanisms ?**
- 02 **Soft attention and hard attention**
- 03 **Self attention**
- 04 **Application of attention mechanism**

What are the attention mechanisms ?



YI REN WEI BEN

Attention Mechanisms in Neural Networks are (very) loosely based on the visual attention mechanism found in humans.

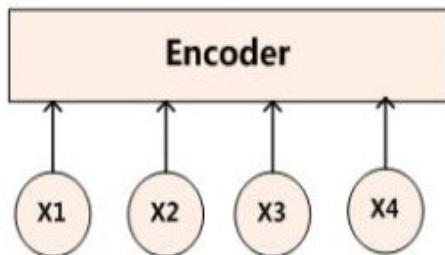


Bottom-up

Top-down



Encoder-Decoder

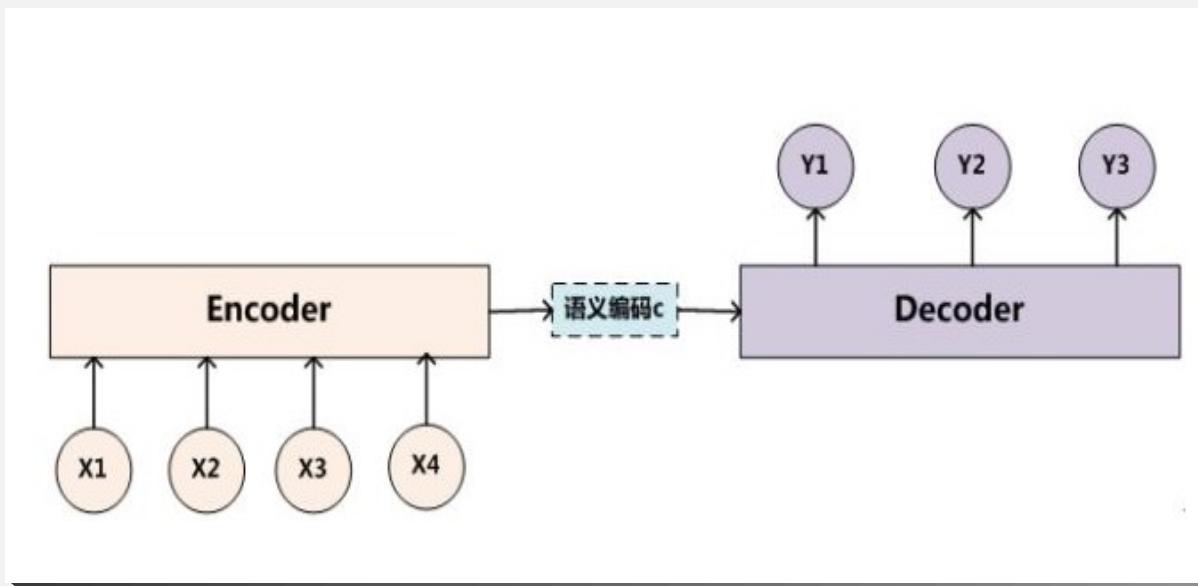


Seq2Seq

Sequence-to-sequence (seq2seq) models have enjoyed great success in a variety of tasks.

- **Source**= $\langle x_1, x_2 \dots x_m \rangle$
- **Target**= $\langle y_1, y_2 \dots y_m \rangle$
- **C=F($x_1, x_2 \dots x_m$)**
- **$y_i=G(C, y_1, y_2 \dots y_{i-1})$**

Attention mechanism



Soft attention

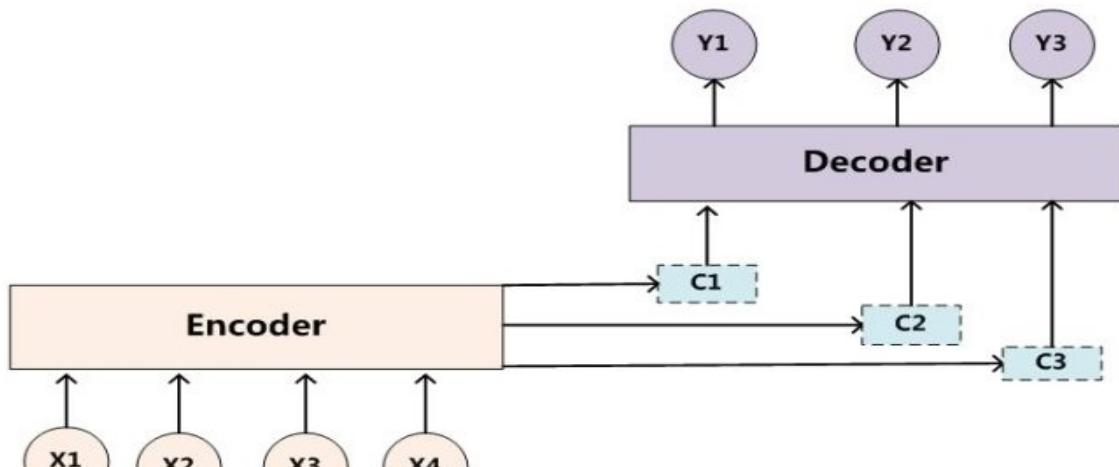
Multiples features with a (soft) mask of values between zero and one.

- $y_1 = f(C)$
- $y_2 = f(C, y_1)$
- $y_3 = f(C, y_1, y_2)$

Tom chase Jerry

汤姆追逐杰瑞

Attention mechanism



Soft attention

Multiples features with a (soft) mask of values between zero and one.

- $(\text{Tom}, 0.3)$ $(\text{Chase}, 0.2)$ $(\text{Jerry}, 0.5)$
- $y_1 = f(C_1)$
- $y_2 = f(C_2, y_1)$
- $y_3 = f(C_3, y_1, y_2)$

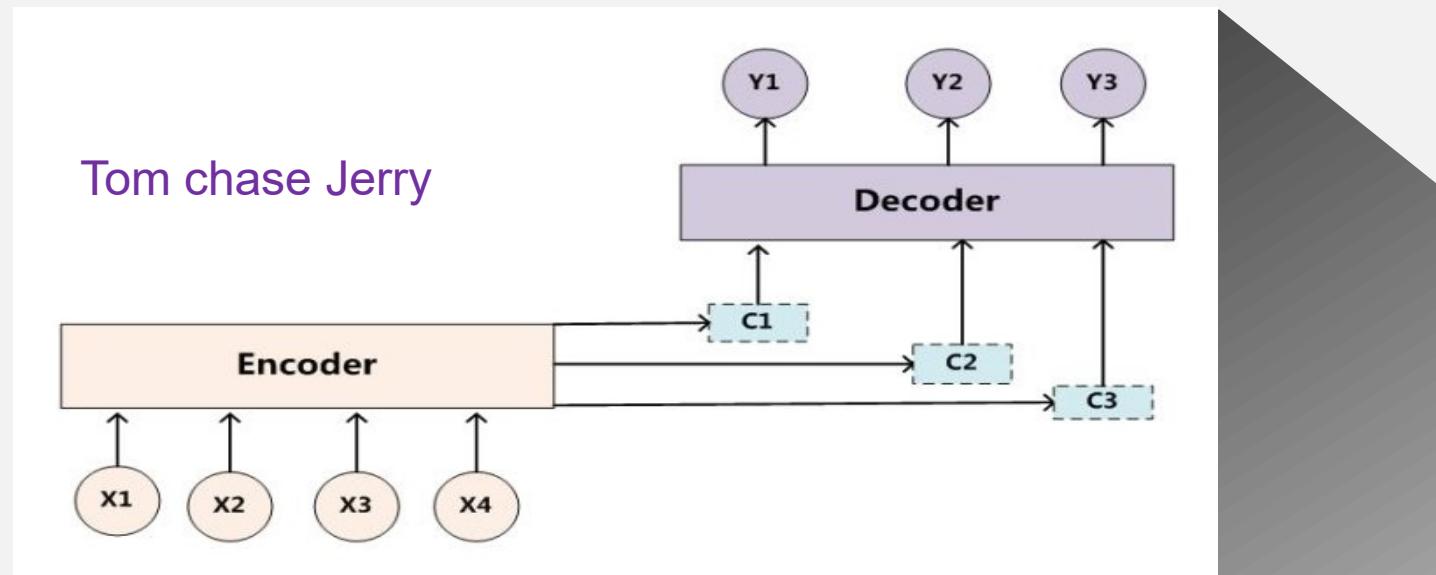
$$c_t = \sum_{j=1}^{Lx} a_{ij} h_j$$

$0.2 * f2(\text{Chase}), 0.2 * f2("Jerry")$

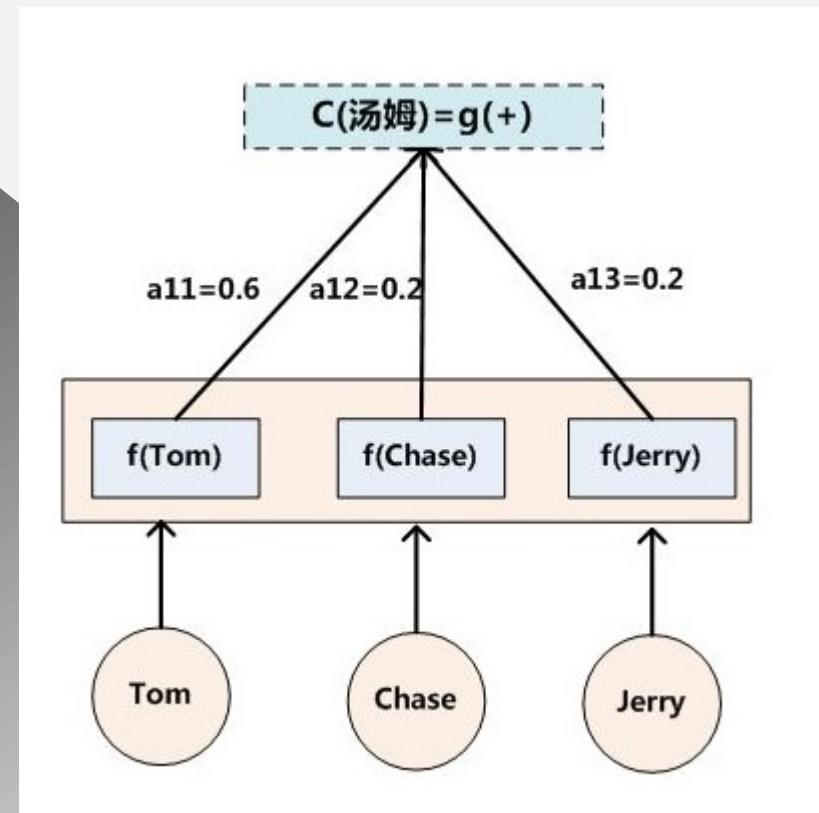
$0.3 * f2(\text{Tom}), 0.7 * f2(\text{Chase}), 0.1 * f2("Jerry")$

$0.2 * f2("Tom"), 0.2 * f2(\text{Chase}), 0.5 * f2("Jerry")$

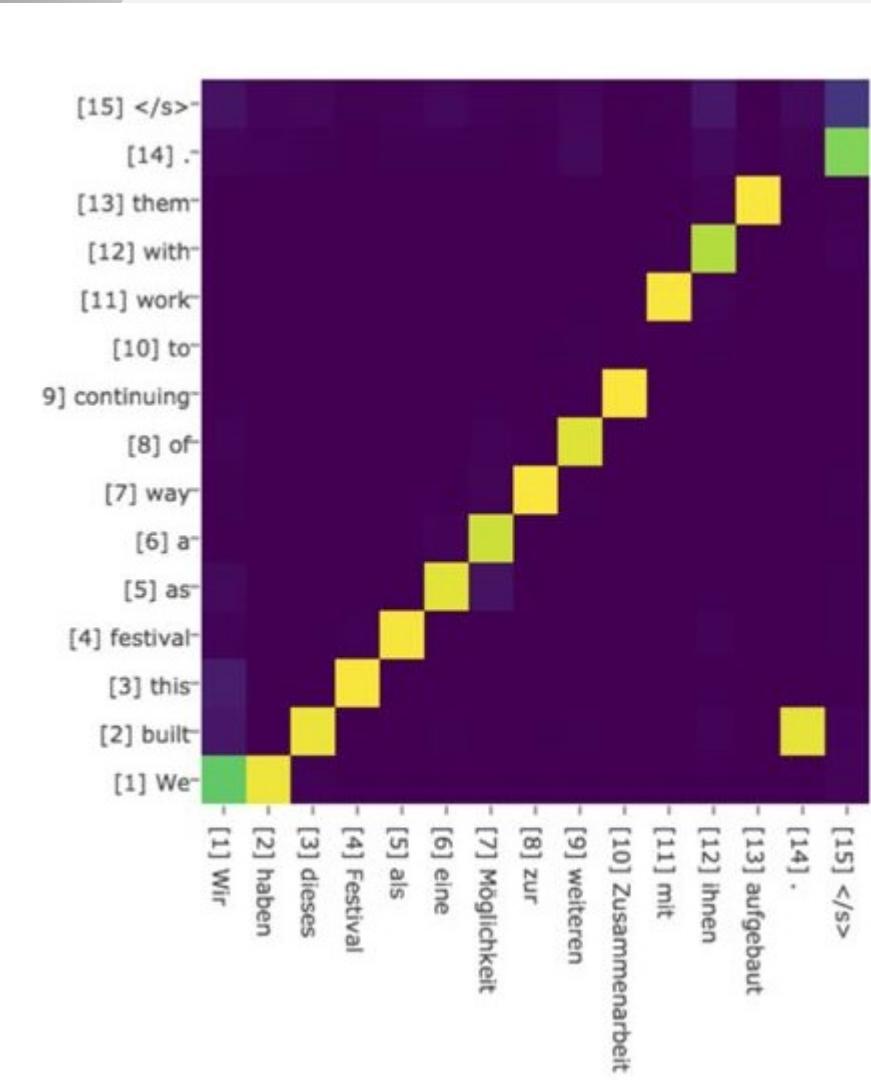
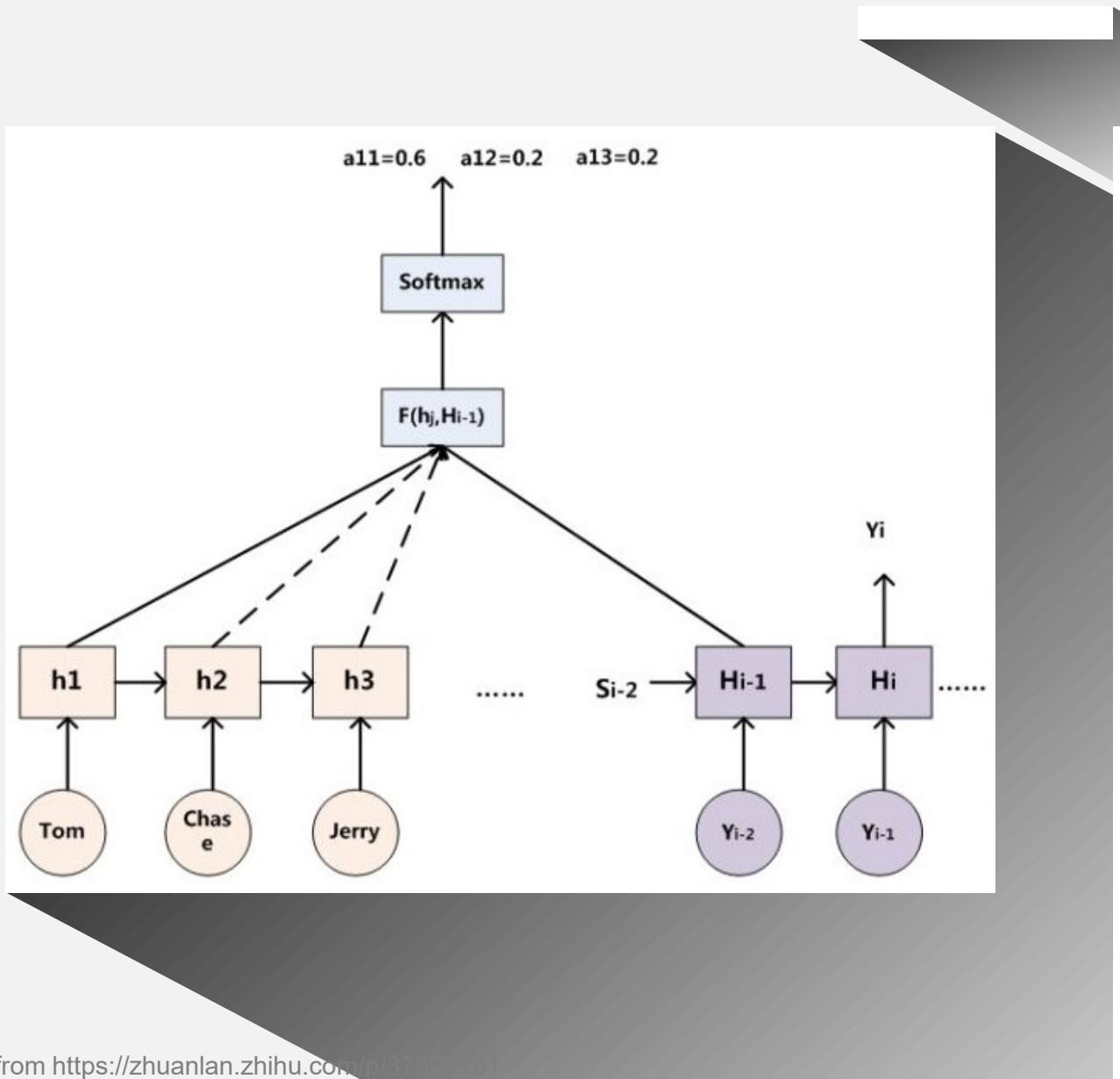
Attention mechanism



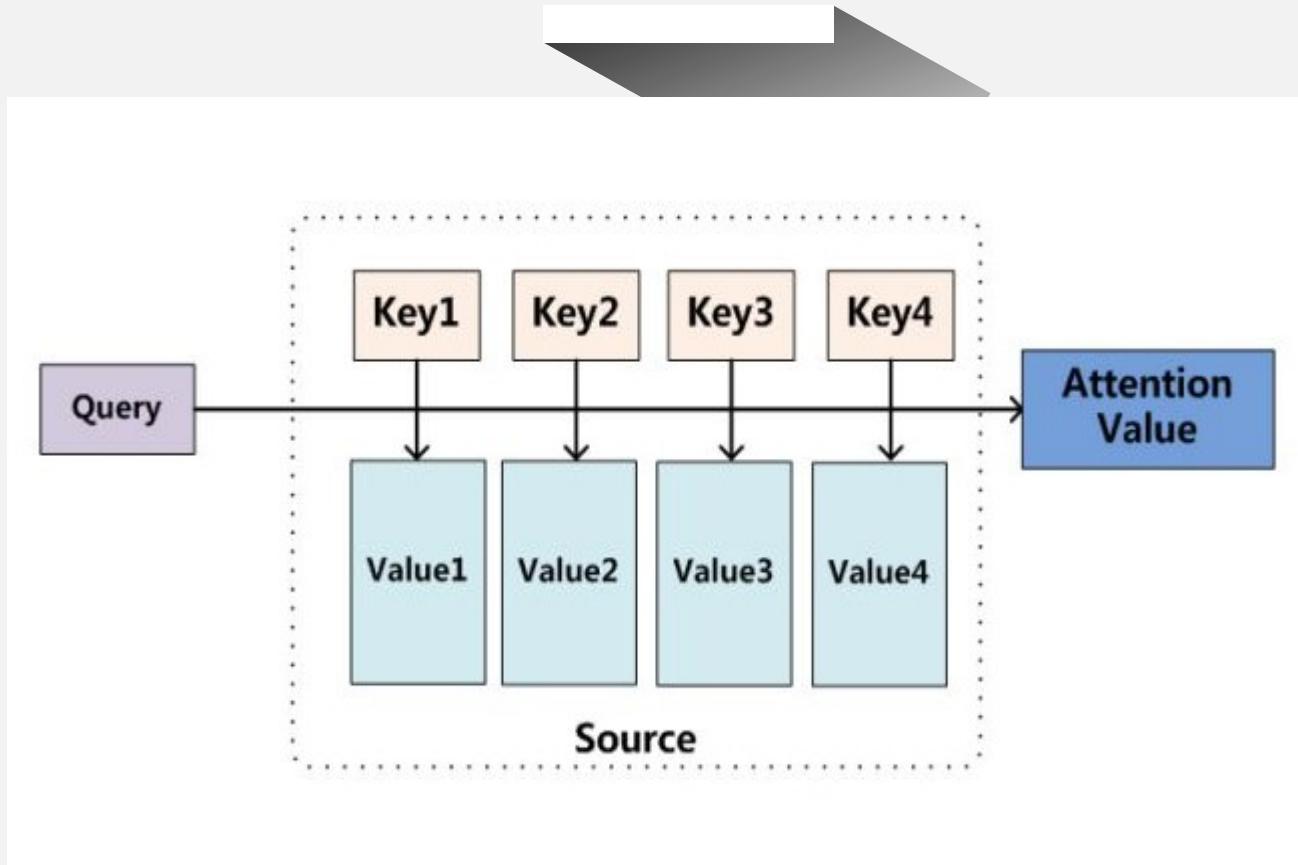
$$c_t = \sum_{j=1}^{Lx} a_{ij} h_j$$



Attention mechanism

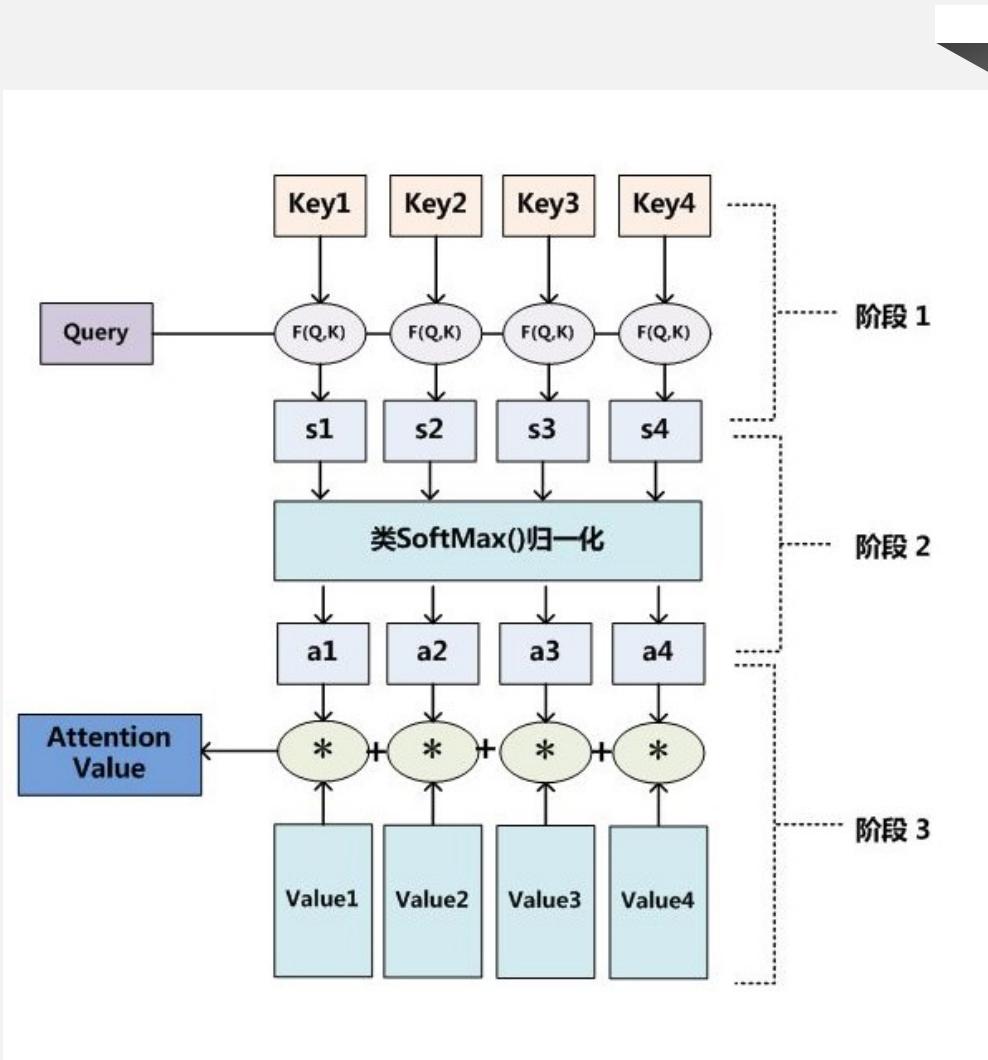


Attention mechanism



$$\text{Attention}(\text{Query}, \text{Source}) = \sum_{i=1}^{L_x} \text{Similarity}(\text{Query}, \text{Key}_i) * \text{Value}_i$$

Attention mechanism



$$\text{Similarity}(\text{Query}, \text{Key}_i) = \text{Query} * \text{Key}_i$$

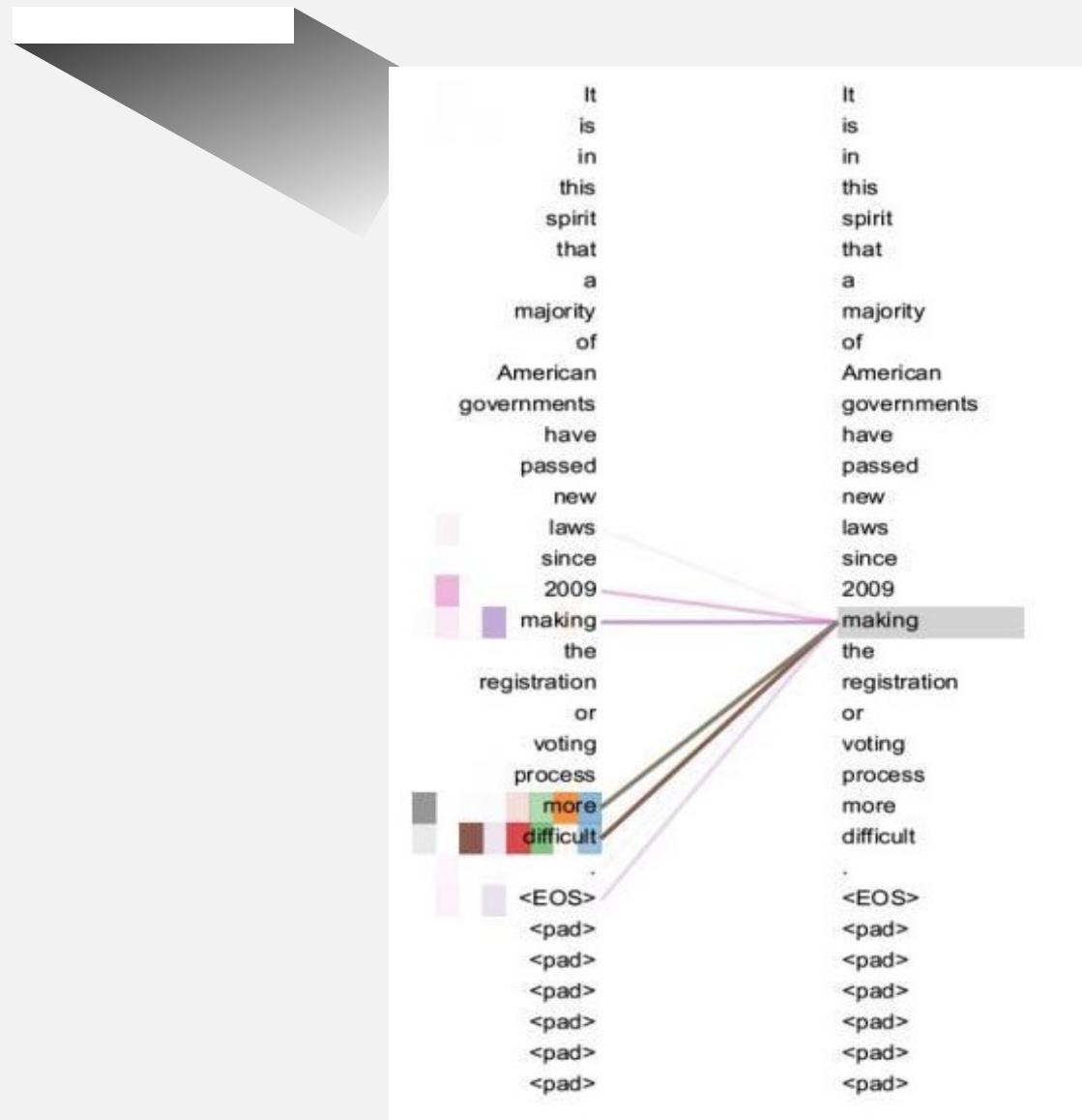
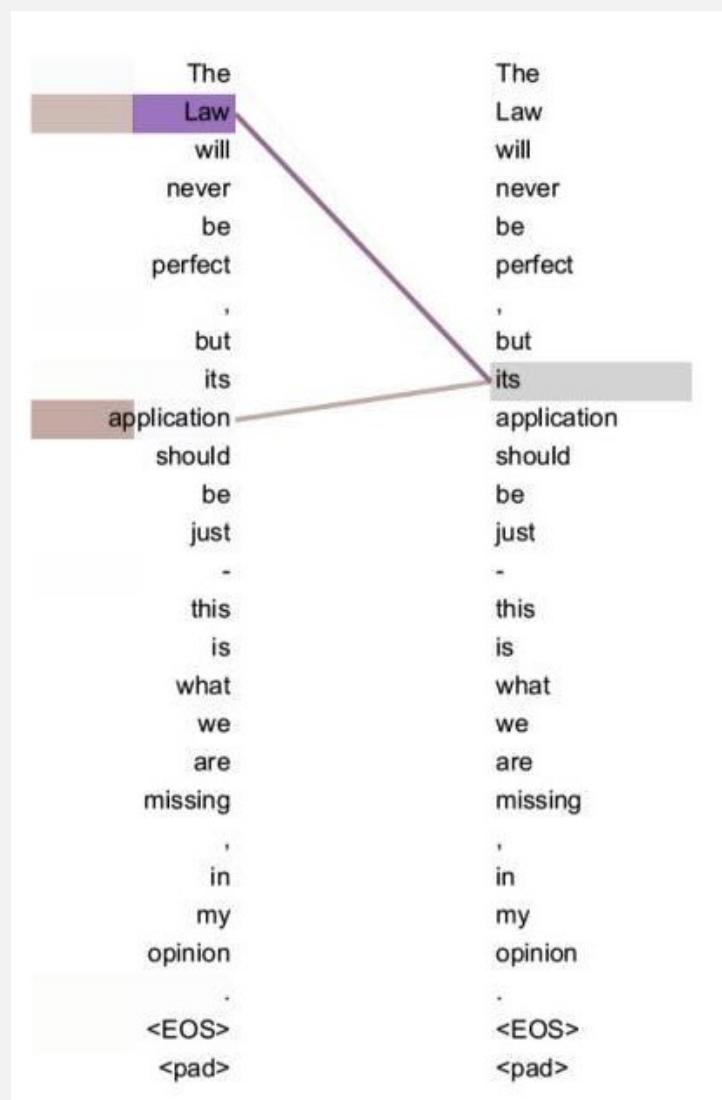
$$\text{Similarity}(\text{Query}, \text{Key}_i) = \frac{\text{Query} * \text{Key}_i}{\|\text{Query}\| * \|\text{Key}_i\|}$$

$$\text{Similarity}(\text{Query}, \text{Key}_i) = \text{MLP}(\text{Query}, \text{Key}_i)$$

$$a_i = \text{softmax}(\text{Sim}_i) = \frac{e^{\text{Sim}_i}}{\sum_{j=1}^{L_x} e^{\text{Sim}_i}}$$

$$\text{Attention}(\text{Query}, \text{Source}) = \sum_{j=1}^{L_x} a_i * \text{Value}_i$$

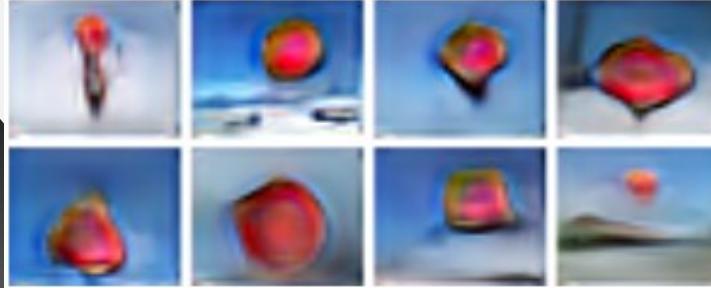
Self Attention



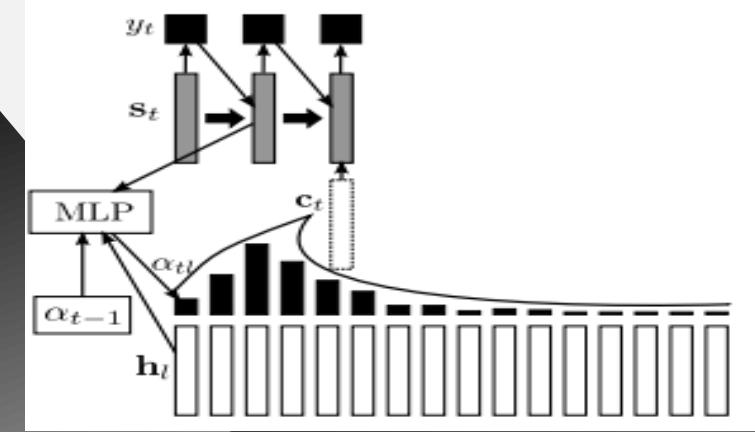
Application of attention mechanism



Image caption



Text2image



Speech recognition

Image caption



A dog is standing on a hardwood floor.

Image caption

Automatically generating captions of an image is a task very close to the heart of scene understanding

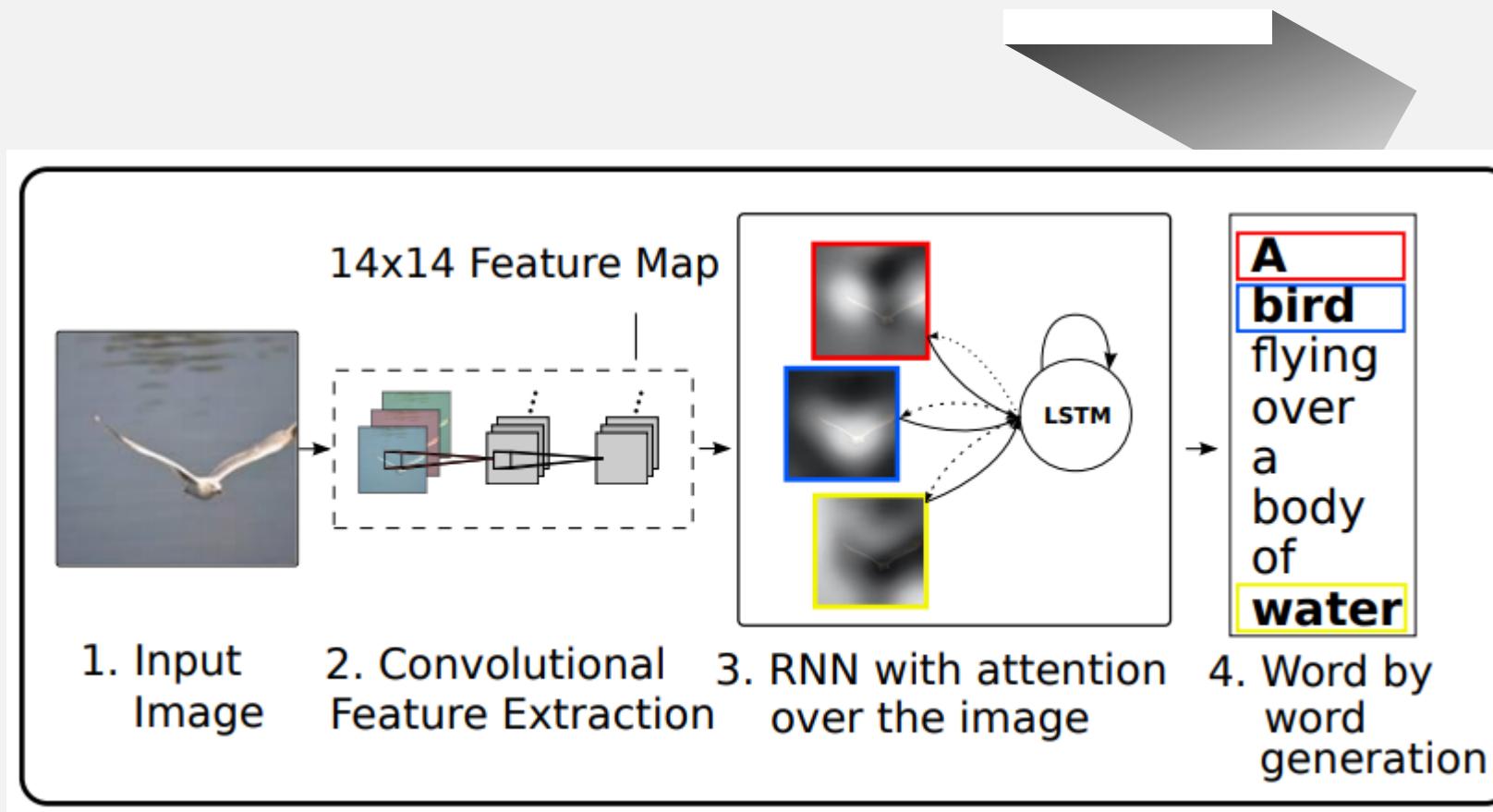
Encoder-decoder



Attention mechanism



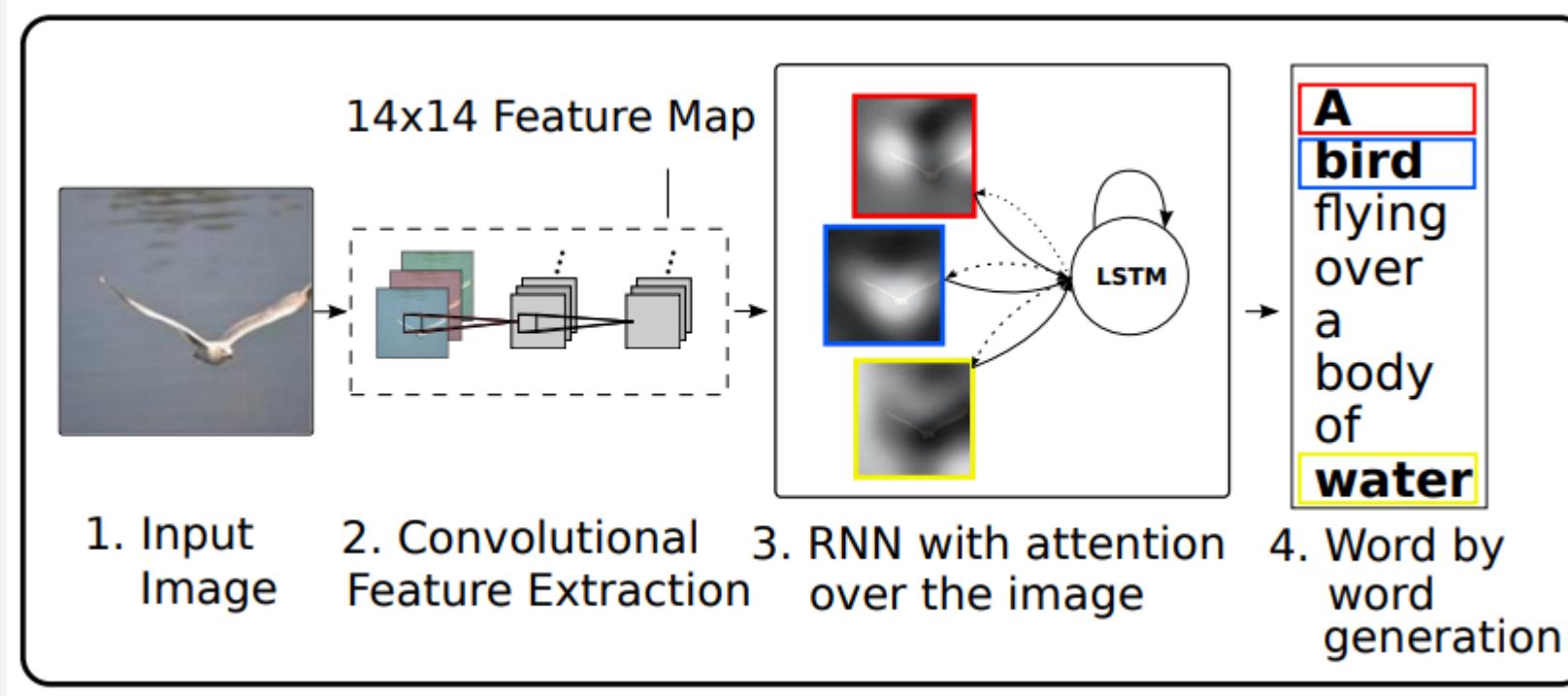
Show, attend and tell



$$\mathbf{a} = \mathbf{a}_1, \dots, \mathbf{a}_L, \quad \mathbf{a}_i \in \mathbb{R}^D$$

VGG19		
序号	层结构	
1	conv1-1	1
2	relu1-1	
3	conv1-2	2
4	relu1-2	
5	pool1	
6	conv2-1	3
7	relu2-1	
8	conv2-2	4
9	relu2-2	
10	pool2	
11	conv3-1	5
12	relu3-1	
13	conv3-2	6
14	relu3-2	
15	conv3-3	7
16	relu3-3	
17	conv3-4	8
18	relu3-4	
19	pool3	
20	conv4-1	9
21	relu4-1	
22	conv4-2	10
23	relu4-2	
24	conv4-3	11
25	relu4-3	
26	conv4-4	12
27	relu4-4	
28	pool4	
29	conv5-1	13
30	relu5-1	
31	conv5-2	14
32	relu5-2	
33	conv5-3	15
34	relu5-3	
35	conv5-4	16
36	relu5-4	
37	pool5	
38	fc6(4096)	17
39	relu6	
40	fc7(4096)	18
41	relu7	
42	fc8(1000)	19
43	prob(softmax)	

Show, attend and tell



$$\mathbf{a} = \mathbf{a}_1, \dots, \mathbf{a}_L, \quad \mathbf{a}_i \in \mathbb{R}^D$$

Hard attention

$$\hat{\mathbf{z}}_t = \sum_i s_{t,i} \mathbf{a}_i$$

Soft attention

$$L_{p(s_t|a)}[\hat{\mathbf{z}}_t] = \sum_{i=1}^L \alpha_{t,i} \mathbf{a}_i$$

Datasets

IMAGE 6620094641
The entity refers to other



SENTENCES

1. Two combatants in a rugby stadium making a move on each other , one wearing green with ball and defended by # 5 in the maroon colors .
2. Two muddy athletes from opposing teams going head to head in a game of rugby while fans watch expectantly .
3. Two soccer players on opposing teams trying to keep the soccer ball from one another .
4. A rugby player is running with the ball as a defender moves to intervene .
5. Two men from opposing teams with dirty uniforms in a match of rugby .

Flickr30k Entities

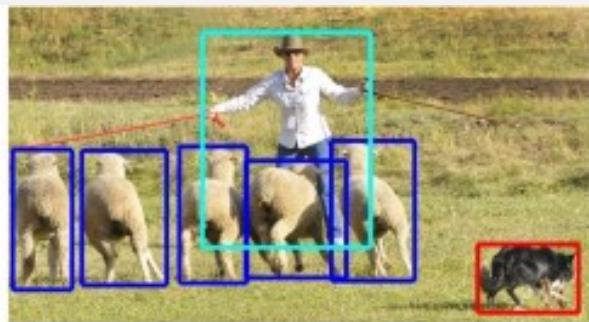
The Flickr30k dataset has become a standard benchmark for sentence-based image description.

- 276k manually annotated bounding boxes
- 158k captions
- 244k coreference chains

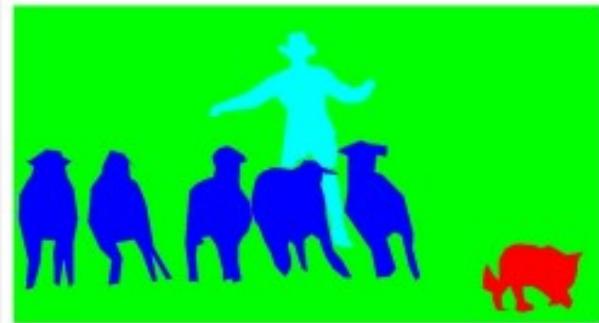
Datasets



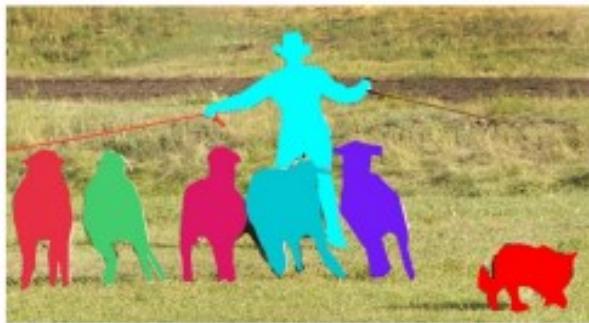
(a) Image classification



(b) Object localization



(c) Semantic segmentation



(d) This work

Microsoft COCO

COCO is a large-scale object detection, segmentation, and captioning dataset.

- Object segmentation
- Recognition in context
- Superpixel stuff segmentation
- 330K images (>200K labeled)
- 1.5 million object instances
- 80 object categories
- 91 stuff categories
- 5 captions per image
- 250,000 people with keypoints

Show, attend and tell



Table 1. BLEU-1,2,3,4/METEOR metrics compared to other methods, † indicates a different split, (—) indicates an unknown metric, ◊ indicates the authors kindly provided missing metrics by personal communication, Σ indicates an ensemble, ^a indicates using AlexNet

Dataset	Model	BLEU				METEOR
		B-1	B-2	B-3	B-4	
Flickr8k	Google NIC(Vinyals et al., 2014) ^{†Σ}	63	41	27	—	—
	Log Bilinear (Kiros et al., 2014a) [◊]	65.6	42.4	27.7	17.7	17.31
	Soft-Attention	67	44.8	29.9	19.5	18.93
	Hard-Attention	67	45.7	31.4	21.3	20.30
Flickr30k	Google NIC ^{†◊Σ}	66.3	42.3	27.7	18.3	—
	Log Bilinear	60.0	38	25.4	17.1	16.88
	Soft-Attention	66.7	43.4	28.8	19.1	18.49
	Hard-Attention	66.9	43.9	29.6	19.9	18.46
COCO	CMU/MS Research (Chen & Zitnick, 2014) ^a	—	—	—	—	20.41
	MS Research (Fang et al., 2014) ^{†a}	—	—	—	—	20.71
	BRNN (Karpathy & Li, 2014) [◊]	64.2	45.1	30.4	20.3	—
	Google NIC ^{†◊Σ}	66.6	46.1	32.9	24.6	—
	Log Bilinear [◊]	70.8	48.9	34.4	24.3	20.03
	Soft-Attention	70.7	49.2	34.4	24.3	23.90
	Hard-Attention	71.8	50.4	35.7	25.0	23.04

A(0.99)



Show, attend and tell



Example 1.

Candidate 1: It is a guide to action which ensures that the military always obeys the commands of the party.

Candidate 2: It is to insure the troops forever hearing the activity guidebook that party direct.

Reference 1: It is a guide to action that ensures that the military will forever heed Party commands.

modified n-gram precision

$$p_n = \frac{\sum_{C \in \{Candidate\}} \sum_{n\text{-gram} \in C} Count_{clip}(n\text{-gram})}{\sum_{C' \in \{Candidate\}} \sum_{n\text{-gram}' \in C'} Count(n\text{-gram})}$$

Example 2.

Candidate: the the the the the the.

Reference 1: The cat is on the mat.

Reference 2: There is a cat on the mat.

$$Bleu = BP \cdot \exp\left(\sum_{n=1}^N w_n \log p_n\right)$$

Show, attend and tell

Figure 3. Examples of attending to the correct object (*white* indicates the attended regions, *underlines* indicated the corresponding word)



A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



A little girl sitting on a bed with a teddy bear.

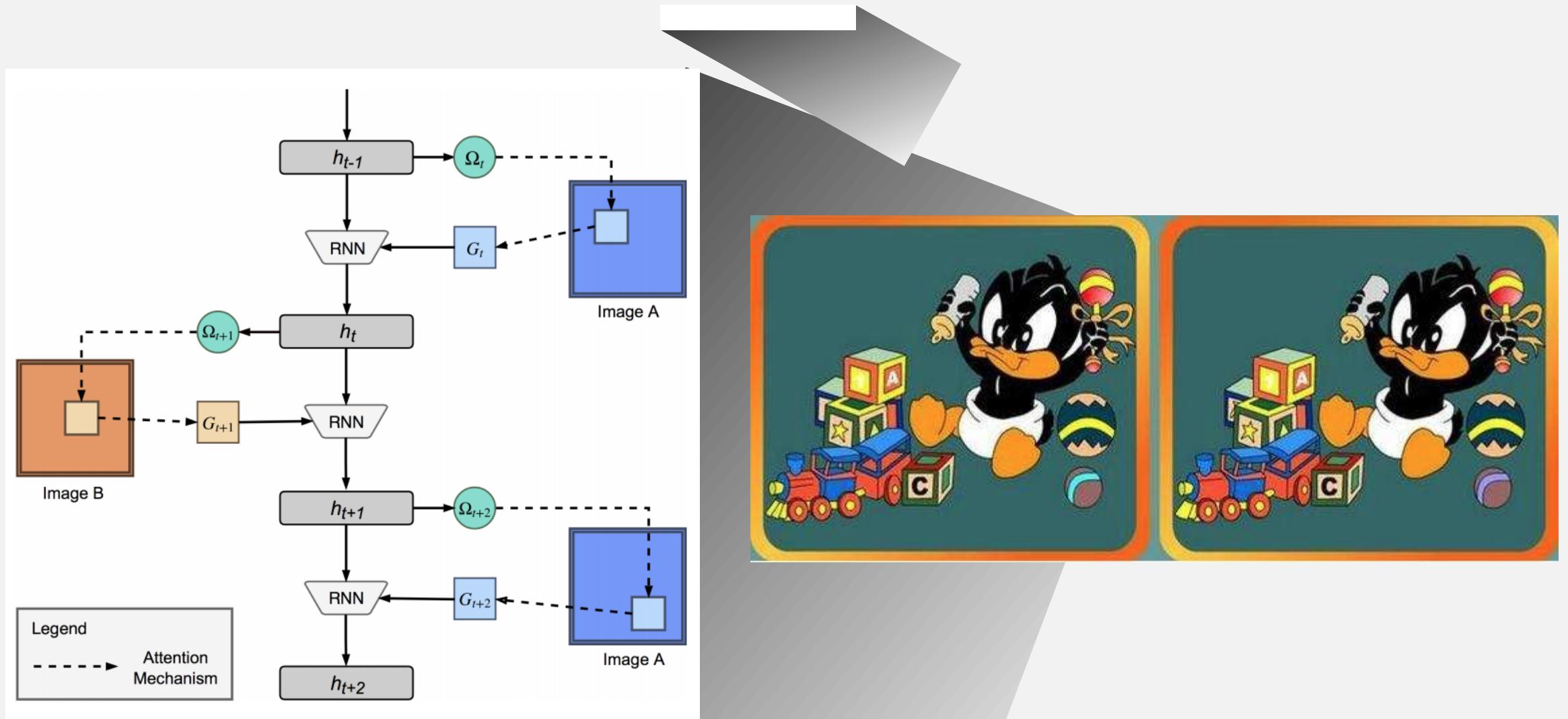


A group of people sitting on a boat in the water.

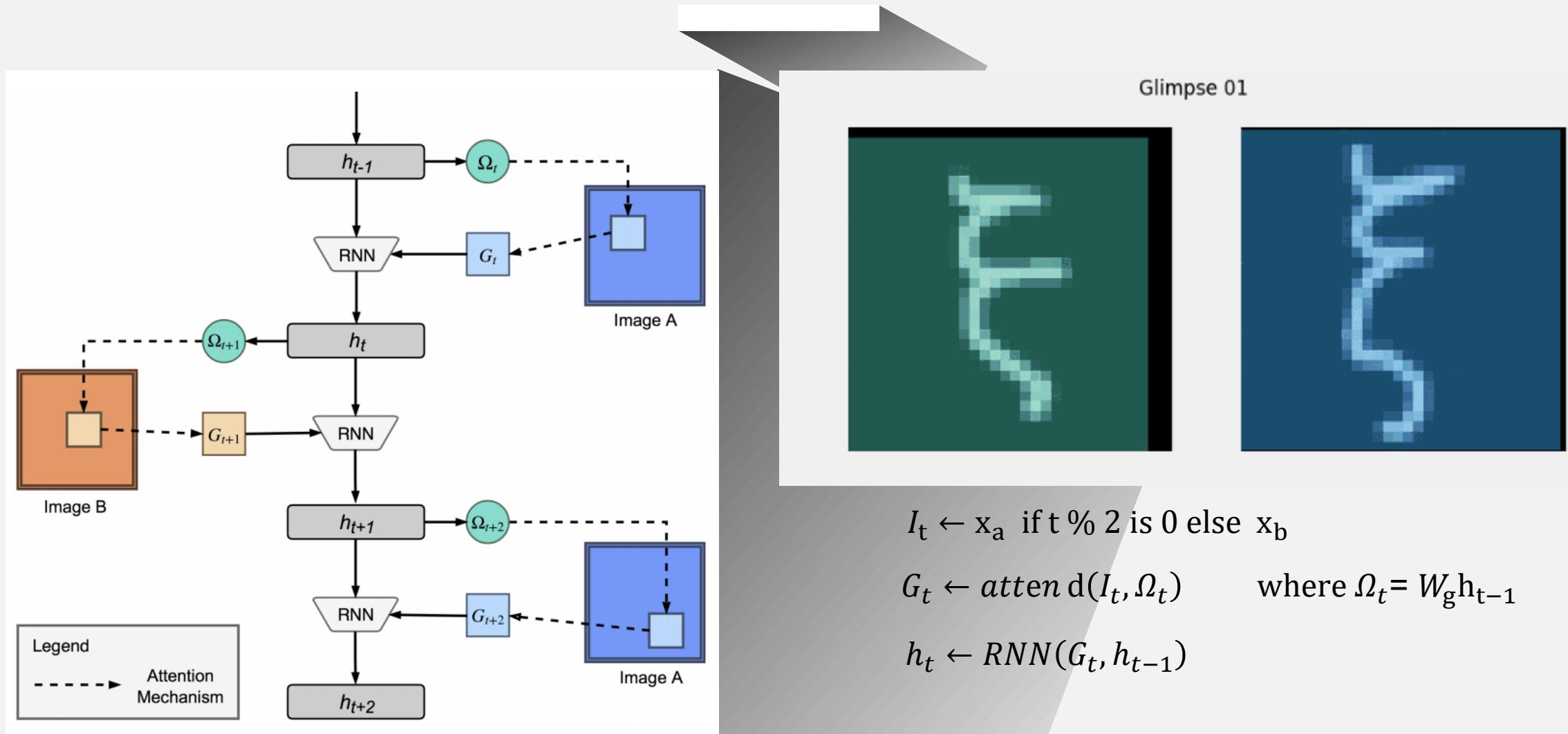


A giraffe standing in a forest with trees in the background.

Attentive Recurrent Comparators

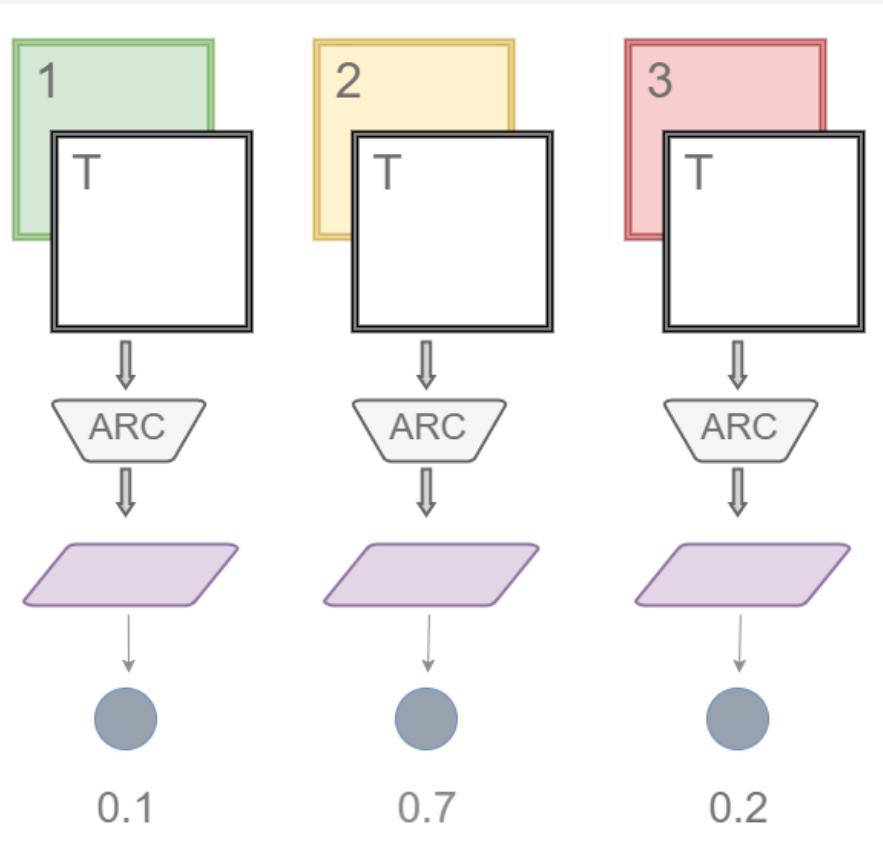


Attentive Recurrent Comparators

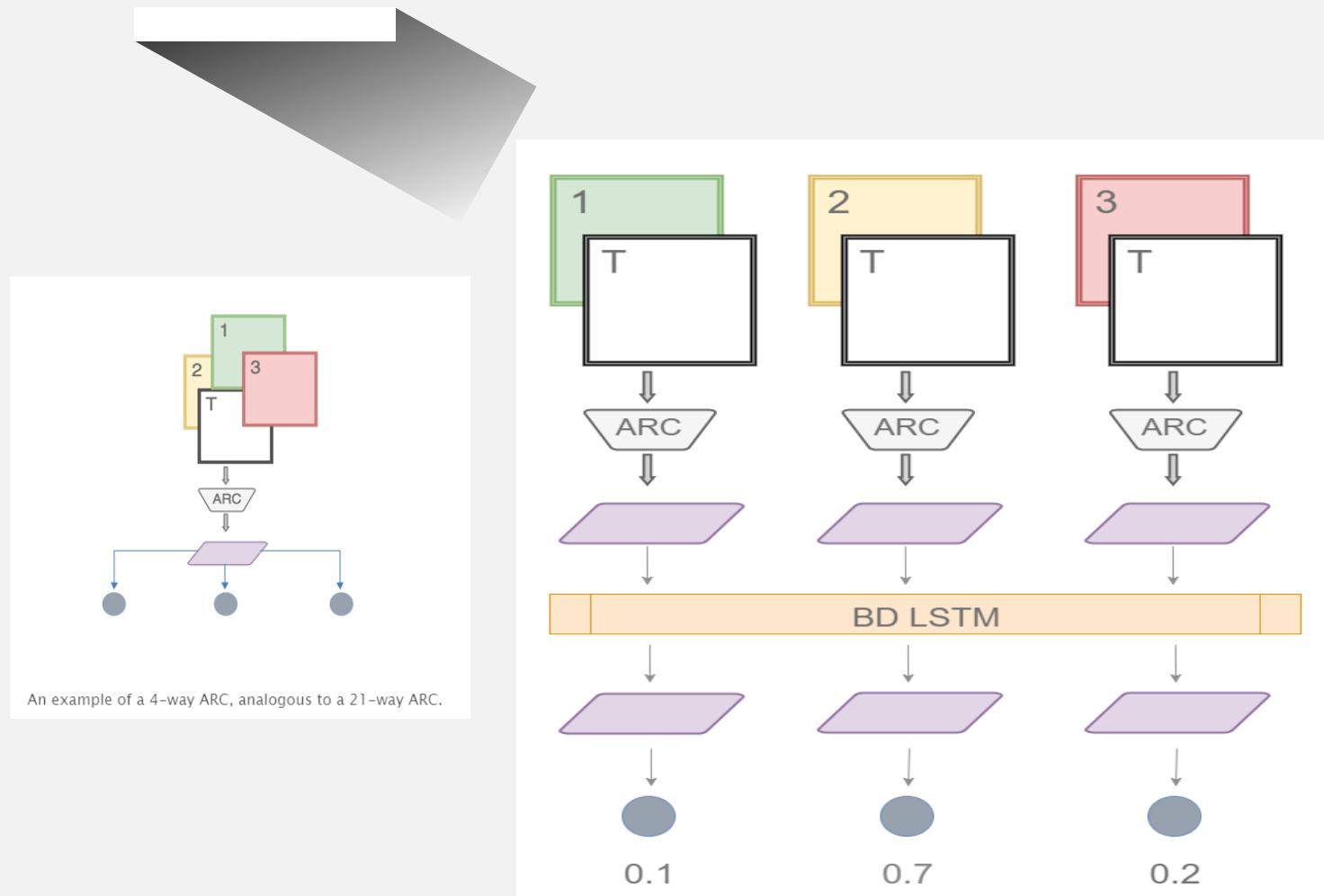


Attentive Recurrent Comparators

20Way one-shot



Naive ARC



Full Context ARC

Attentive Recurrent Comparators



(a) It can be seen that the two characters look very similar in their stroke pattern and differ only in their looping structure. ARC has learnt to focus on these crucial aspects.

Table 1. Glimpses per image versus classification accuracy of ARC. Out of the 50 alphabets provided in the Omniglot dataset, 30 were used for training and validation and the last 20 for testing

GLIMPSSES	ACCURACY (TEST SET)
1	58.2%
2	65.0%
4	80.8%
6	89.25%
8	92.08%

Attentive Recurrent Comparators

Table 5. One-shot classification accuracies of various methods and our ARC models on Omniglot dataset - Within Alphabets

MODEL	ACCURACY
KNN	21.7%
SIAMESE NETWORK	58.3%
DEEP SIAMESE NETWORK (KOCH ET AL.)	92.0%
HUMANS	95.5%
BPL	96.7%
NAIVE ARC	91.75%
NAIVE CONVARC	97.75%
FULL CONTEXT CONVARC	98.5%

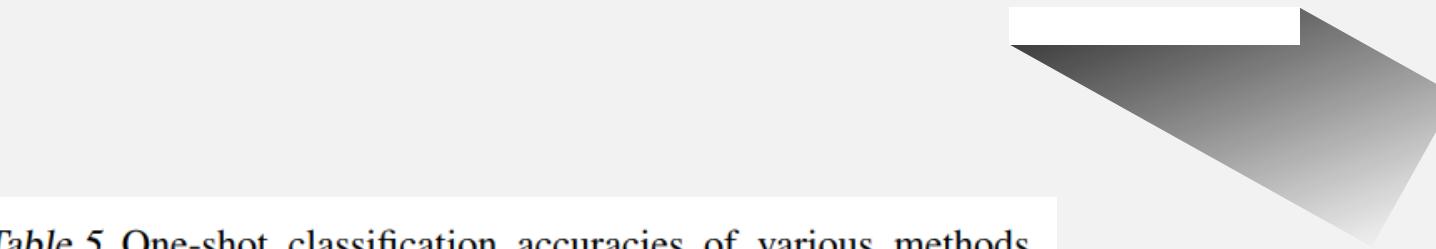


Table 6. 5 way one-shot Classification accuracies of various methods and our ARC models - miniImageNet

MODEL	ACCURACY
RAW PIXELS w/ COSINE SIMILARITY	23.0%
BASELINE CLASSIFIER	38.4%
MATCHING NETWORKS	46.6%
NAIVE CONVARC	49.14%

Visual Question Answering



DAQUAR



Q: How many red chairs are there?

H: ()

M: 6

C: blinds

Q: How many chairs are at the table?

H: wall

M: 4

C: chair

Q: What is on the right side of cabinet?

H: picture

M: bed

C: bed

Q: What is on the wall?

H: mirror

M: bed

C: picture

DAQUAR

The DAset for QUestion Answering on Real-world images



12,000 question-answer pairs
on RGBD images

Based on NYU-Depth V2



COCO-QA



COCO-QA: What does an intersection show on one side and two double-decker buses and a third vehicle,?
Ground Truth: Building

COCO-QA

Exploring Models and Data for Image Question Answering

- Automatically generated from image captions.
- 123287 images
- 78736 train questions
- 38948 test questions
- 4 types of questions: object, number, color, location
- Answers are all one-word.

VQA

Who is wearing glasses?
man woman



Is the umbrella upside down?
yes no



Where is the child sitting?
fridge arms



How many children are in the bed?
2 1

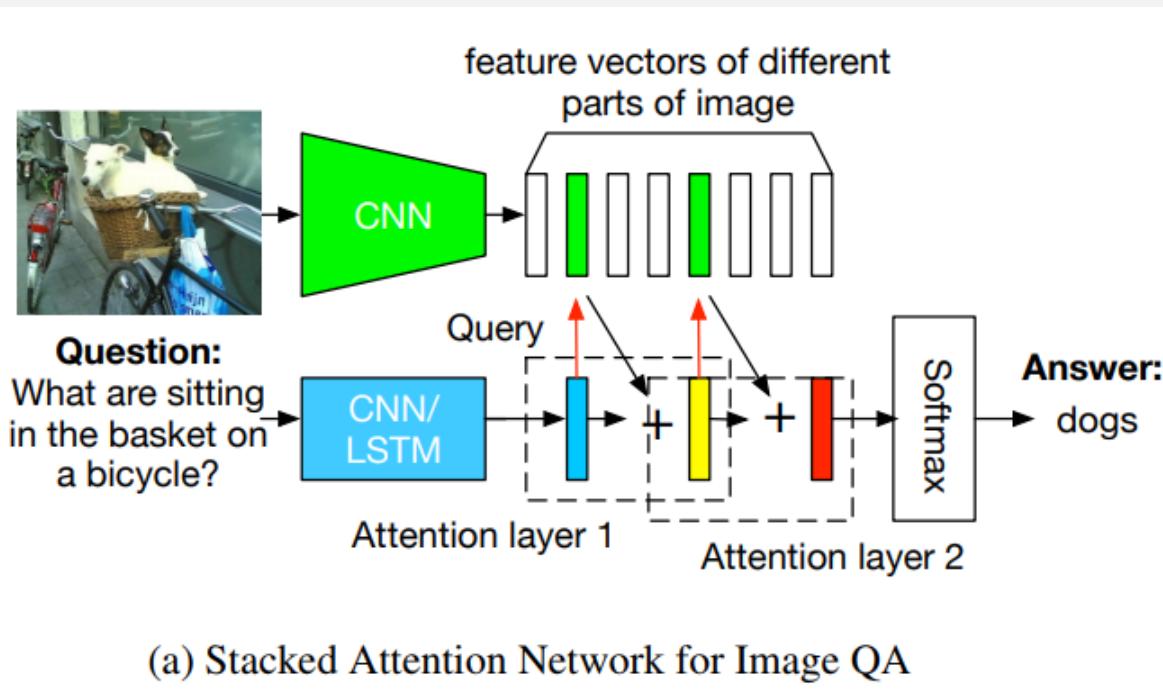


VQA

VQA is a new dataset containing open-ended questions about images. These questions require an understanding of vision language and commonsense knowledge to answer.

- 265,016 images (COCO and abstract scenes)
- At least 3 questions (5.4 questions on average) per image
- 10 ground truth answers per question
- 3 plausible (but likely incorrect) answers per question
- Automatic evaluation metric

Stacked Attention Networks for Image Question Answering



SANs

This paper presents stacked attention networks (SANs) that learn to answer natural language questions from images.

- Image model
- Question model
- Stacked attention networks

Image model

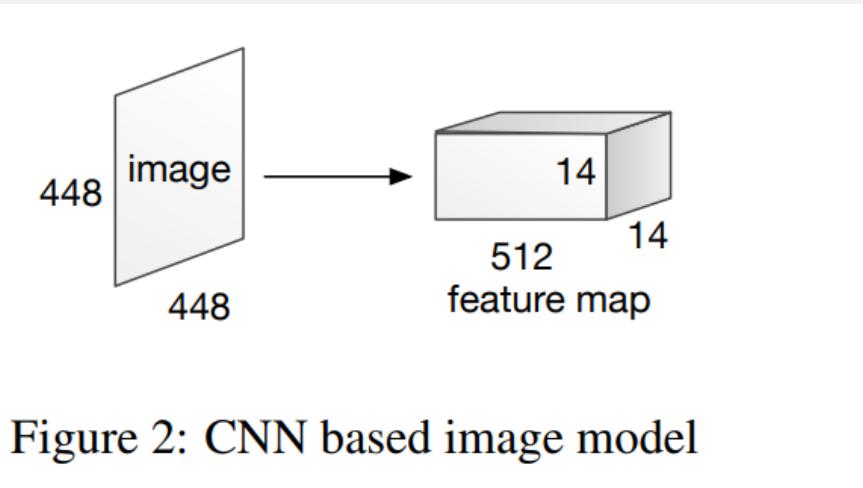
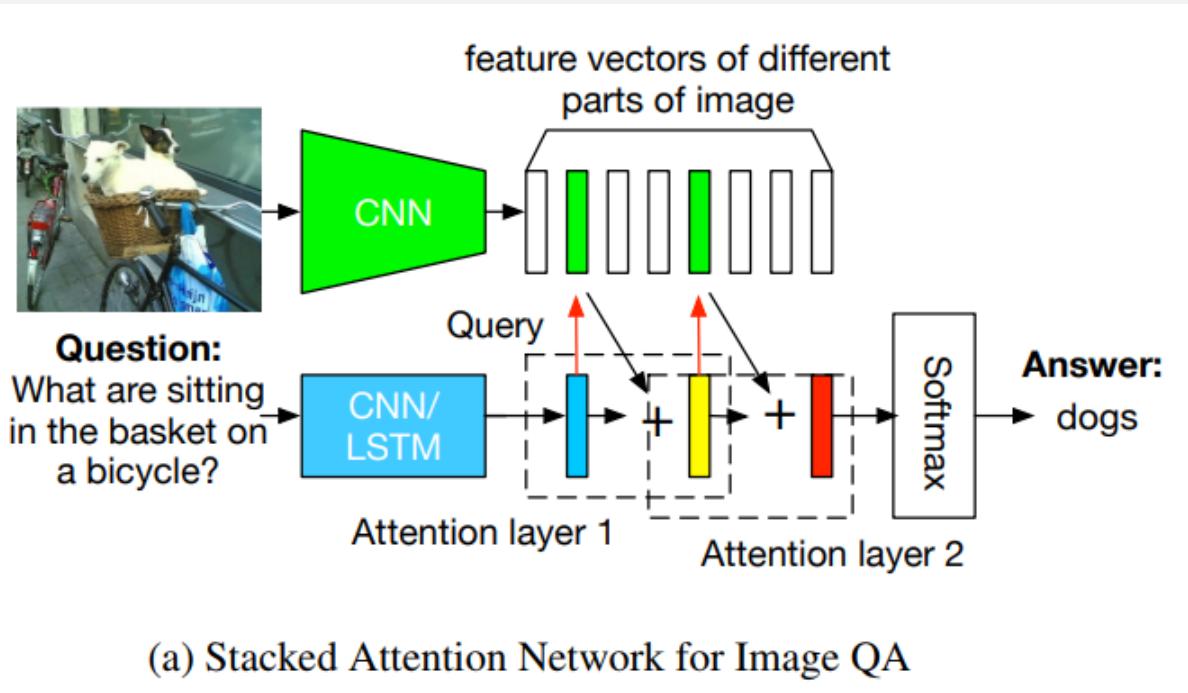
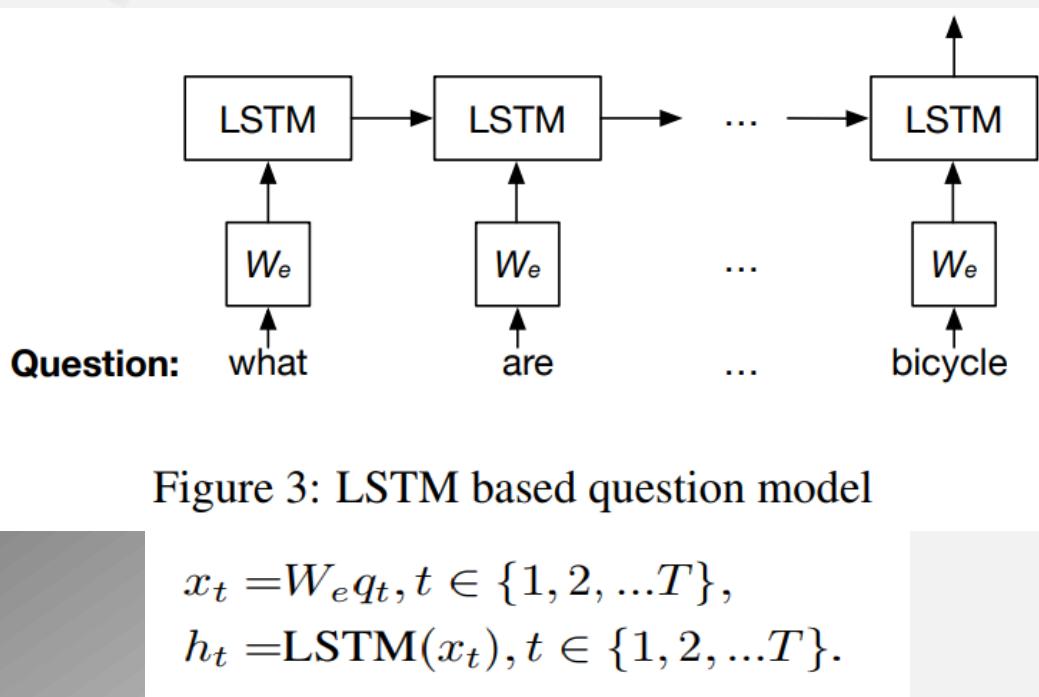
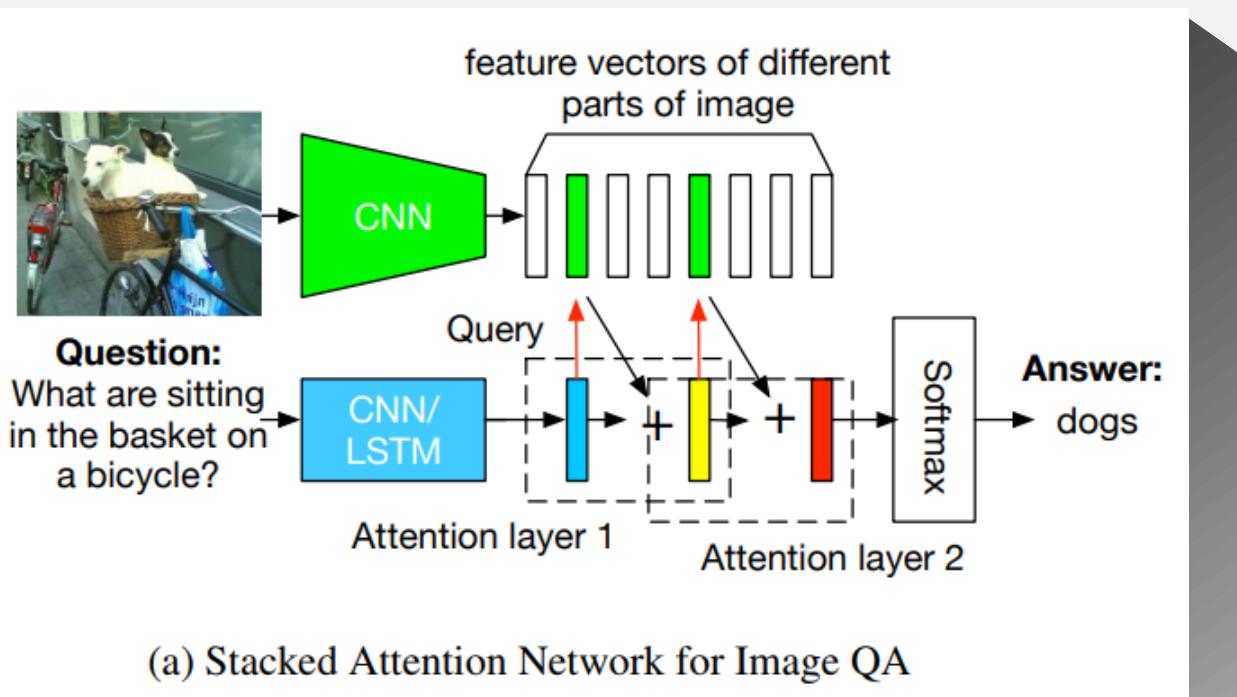


Figure 2: CNN based image model

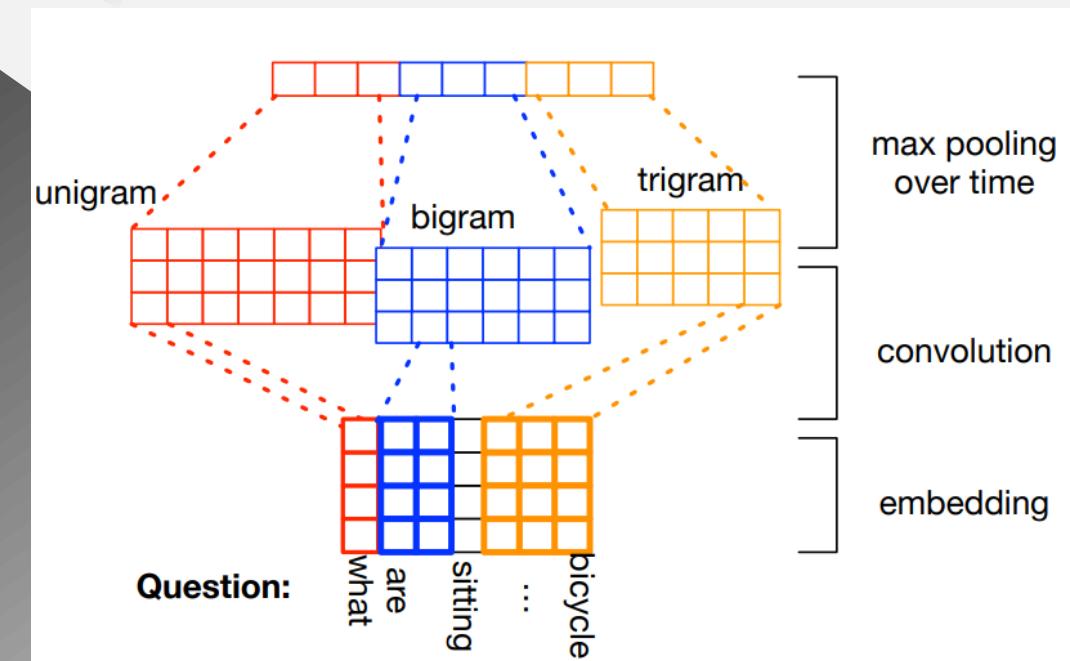
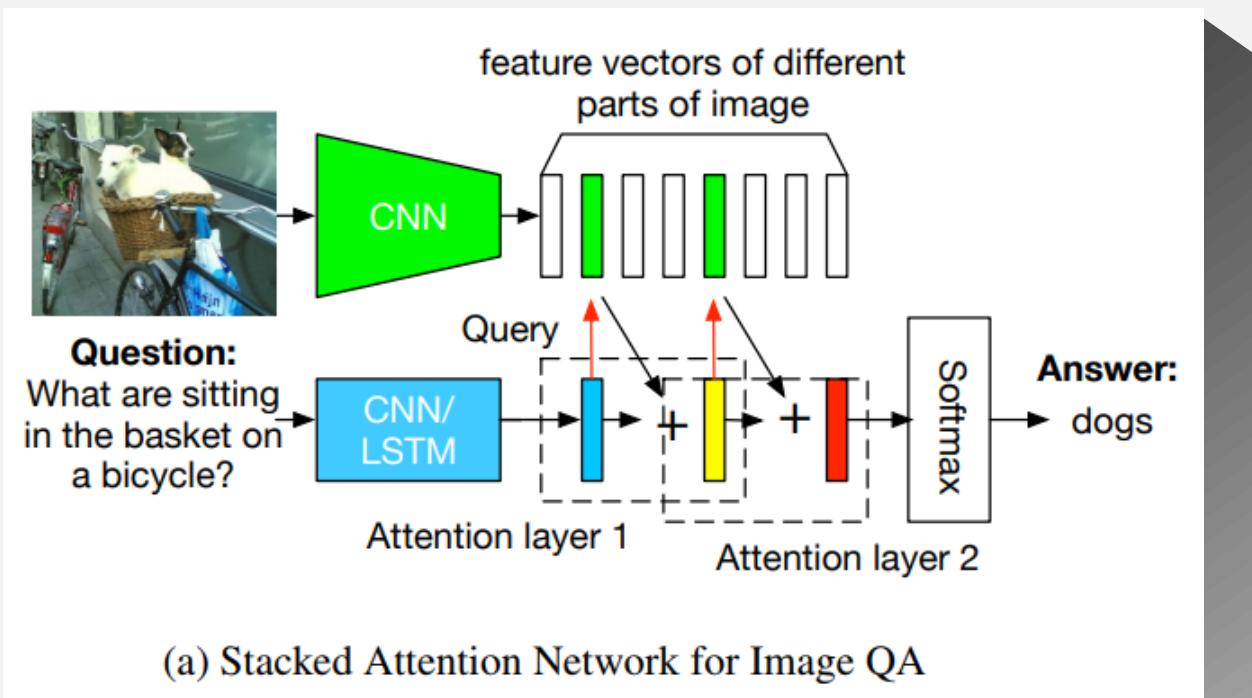
$$f_I = \text{CNN}_{vgg}(I).$$

$$v_I = \tanh(W_I f_I + b_I),$$

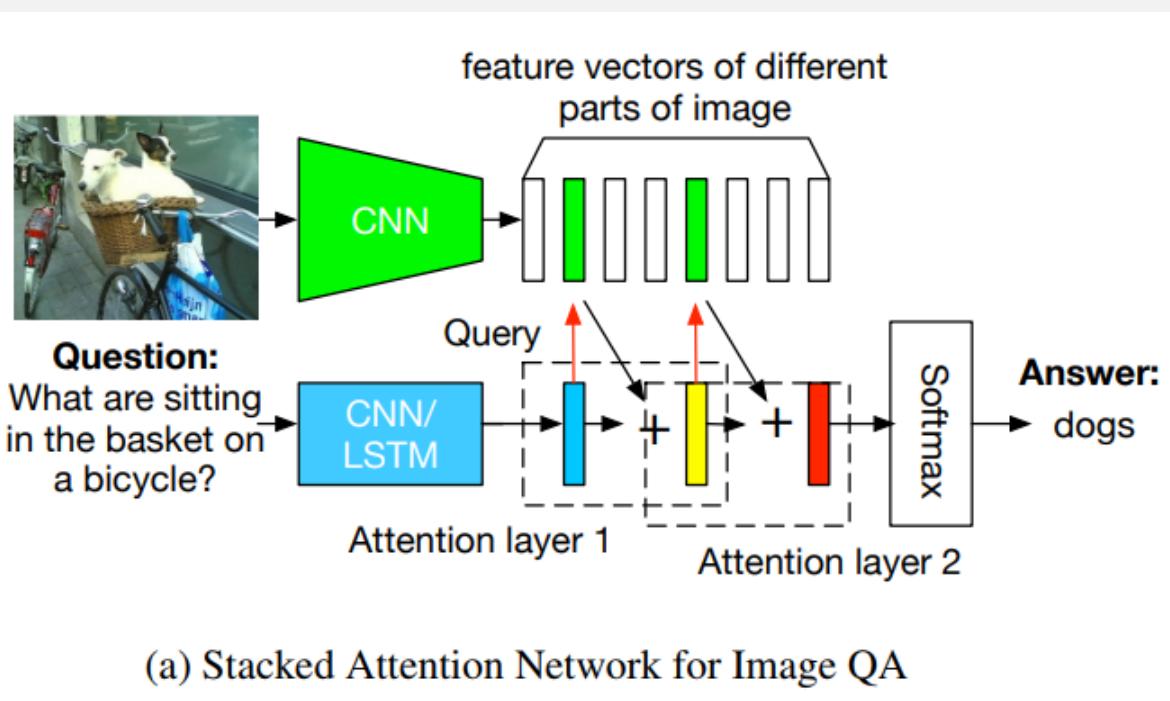
Question model



Question model



Stacked Attention Networks



$$h_A = \tanh(W_{I,A}v_I \oplus (W_{Q,A}v_Q + b_A))$$

$$p_I = \text{softmax}(W_P h_A + b_P)$$

$$\tilde{v}_I = \sum_i p_i v_i$$

$$\mu = \tilde{v}_I + v_Q$$

$$h_A^k = \tanh(W_{I,A}^k v_I \oplus (W_{Q,A}^k \mu^{k-1} + b_A^k))$$

$$p_I^k = \text{softmax}(W_P^k h_A^k + b_P^k)$$

$$\tilde{v}_I^k = \sum_i p_i^k v_i$$

$$\mu^k = \tilde{v}_I^k + \mu^{k-1}$$

$$p_{\text{ans}} = \text{softmax}(W_u u^K + b_u)$$

Stacked Attention Networks

Methods	Accuracy	WUPSO.9	WUPSO.0
VSE: [21]			
GUESS	6.7	17.4	73.4
BOW	37.5	48.5	82.8
LSTM	36.8	47.6	82.3
IMG	43.0	58.6	85.9
IMG+BOW	55.9	66.8	89.0
VIS+LSTM	53.3	63.9	88.3
2-VIS+BLSTM	55.1	65.3	88.6
CNN: [17]			
IMG-CNN	55.0	65.4	88.6
CNN	32.7	44.3	80.9
Ours:			
SAN(1, LSTM)	59.6	69.6	90.1
SAN(1, CNN)	60.7	70.6	90.5
SAN(2, LSTM)	61.0	71.0	90.7
SAN(2, CNN)	61.6	71.6	90.9

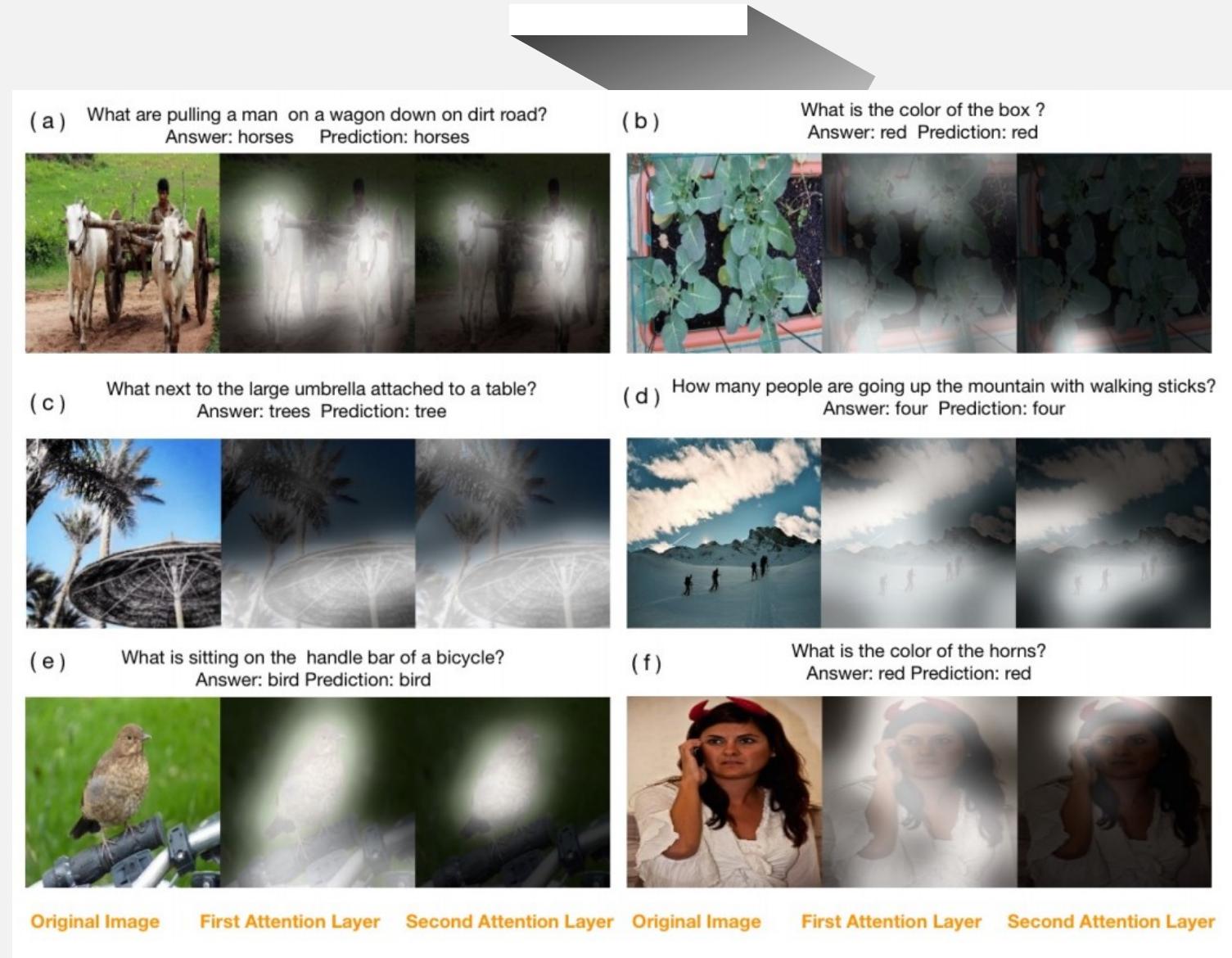
Table 3: COCO-QA results, in percentage



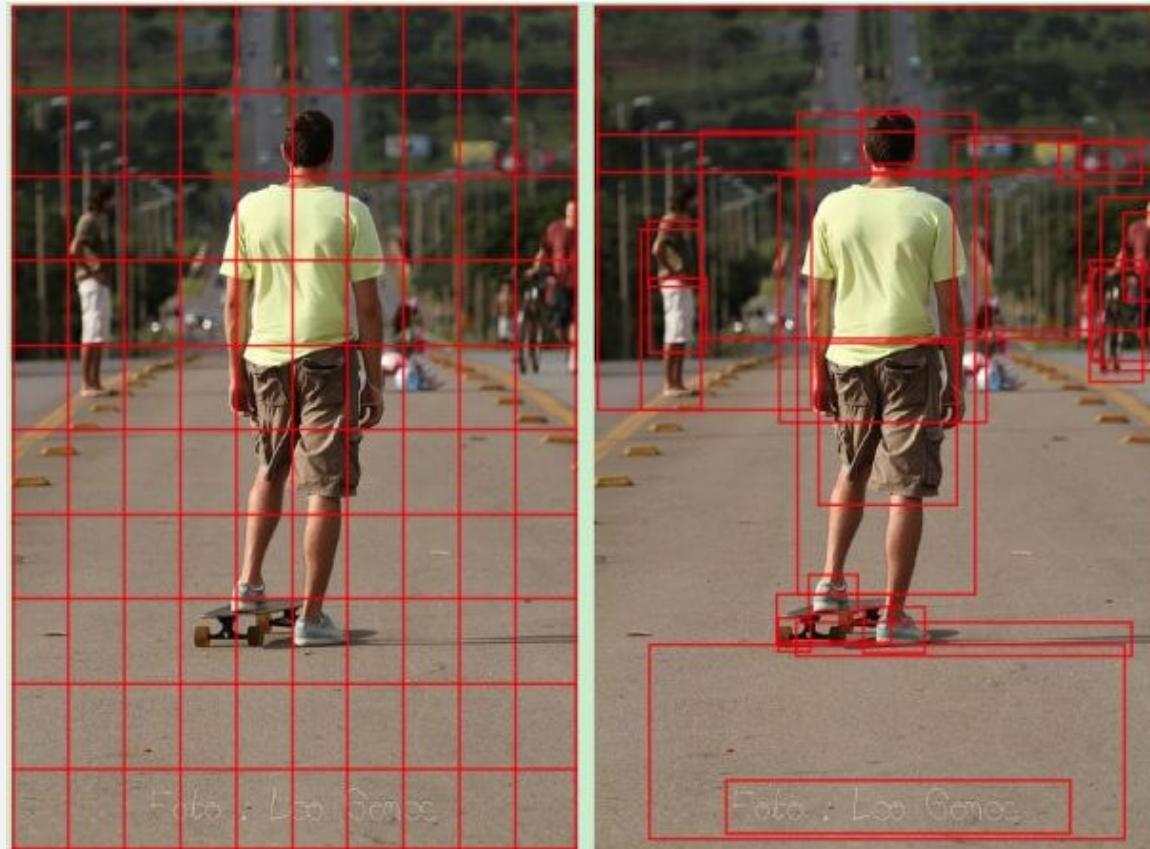
Methods	All	Yes/No 36%	Number 10%	Other 54%
SAN(1, LSTM)	56.6	78.1	41.6	44.8
SAN(1, CNN)	56.9	78.8	42.0	45.0
SAN(2, LSTM)	57.3	78.3	42.2	45.9
SAN(2, CNN)	57.6	78.6	41.8	46.4

Table 6: VQA results on our partition, in percentage

Stacked Attention Networks



Bottom-Up and Top-Down Attention



Bottom-Up and Top-Down Attention

Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering

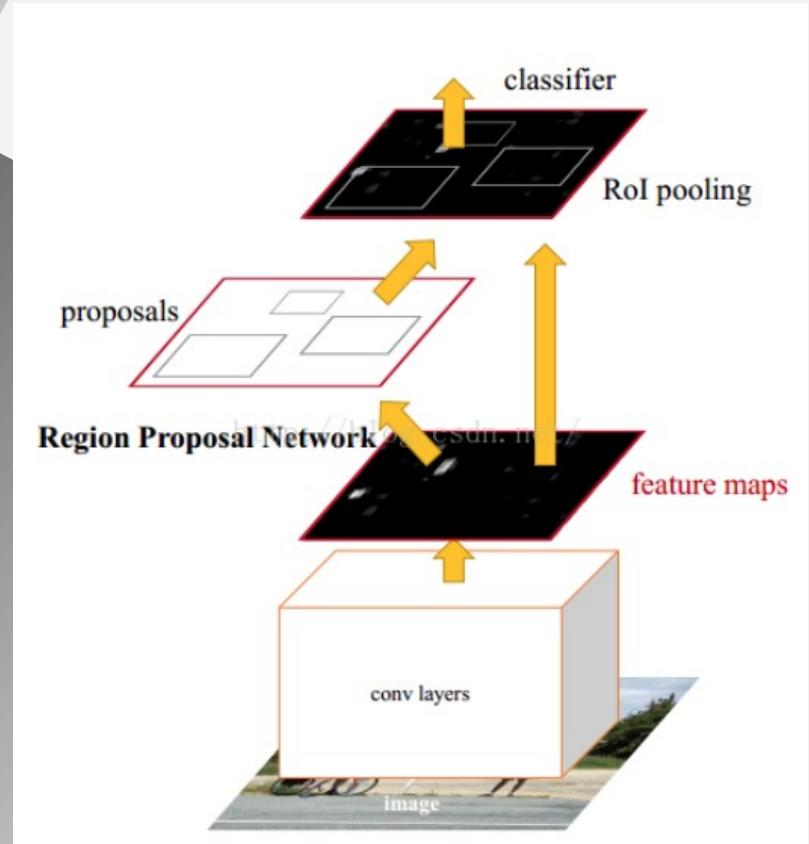
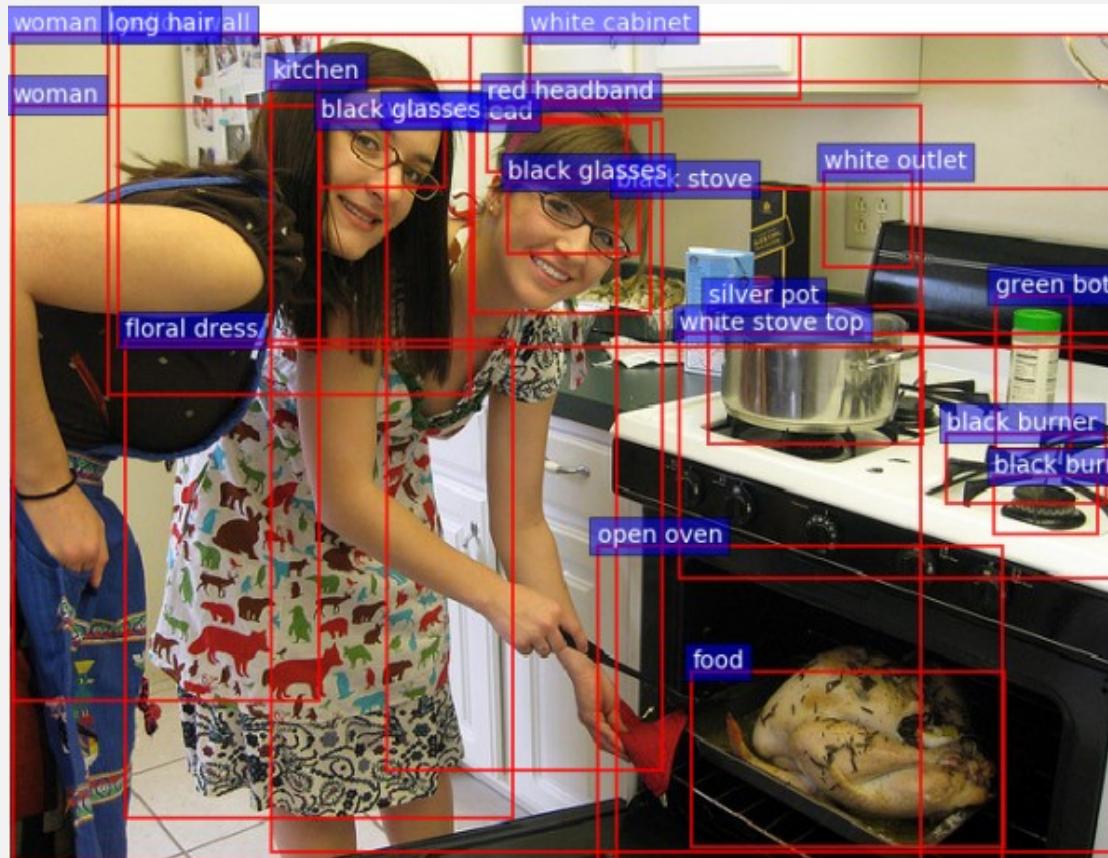


Bottom-Up



Top-Down

Bottom-Up attention model



Bottom-Up attention model

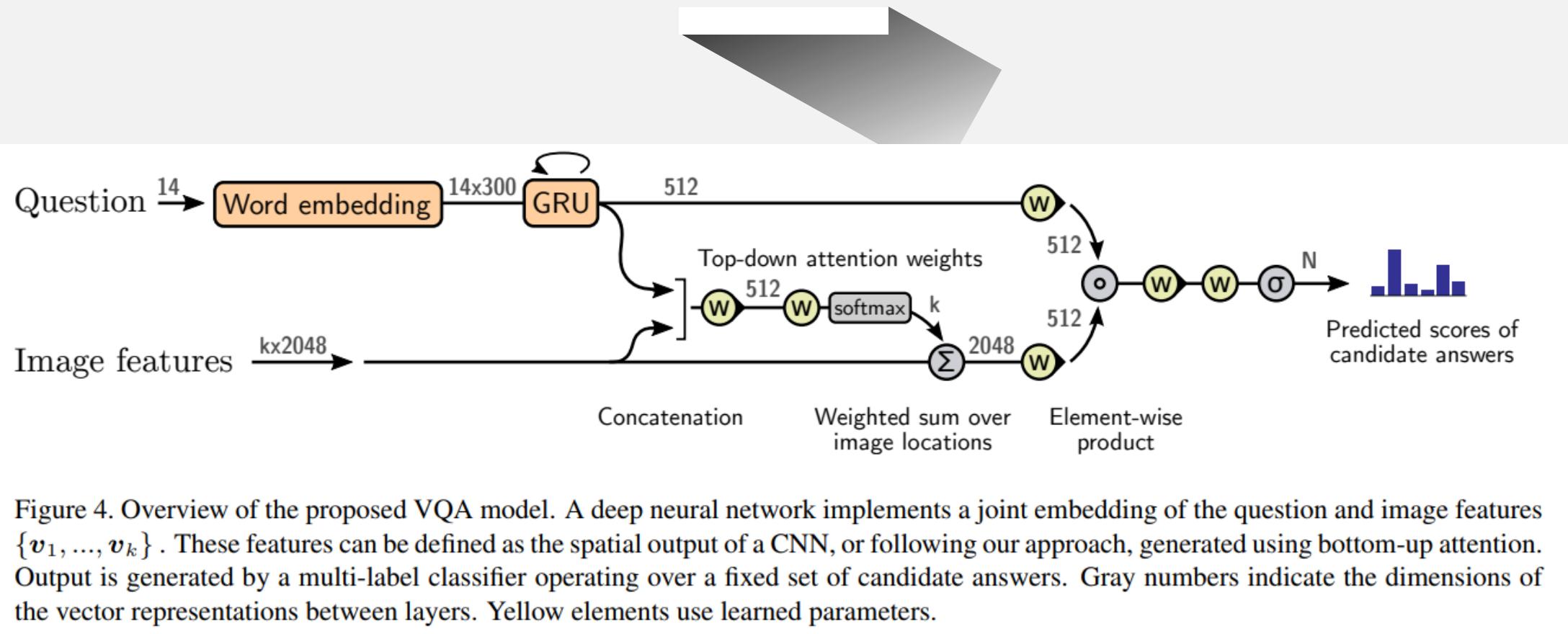


Figure 4. Overview of the proposed VQA model. A deep neural network implements a joint embedding of the question and image features $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$. These features can be defined as the spatial output of a CNN, or following our approach, generated using bottom-up attention. Output is generated by a multi-label classifier operating over a fixed set of candidate answers. Gray numbers indicate the dimensions of the vector representations between layers. Yellow elements use learned parameters.

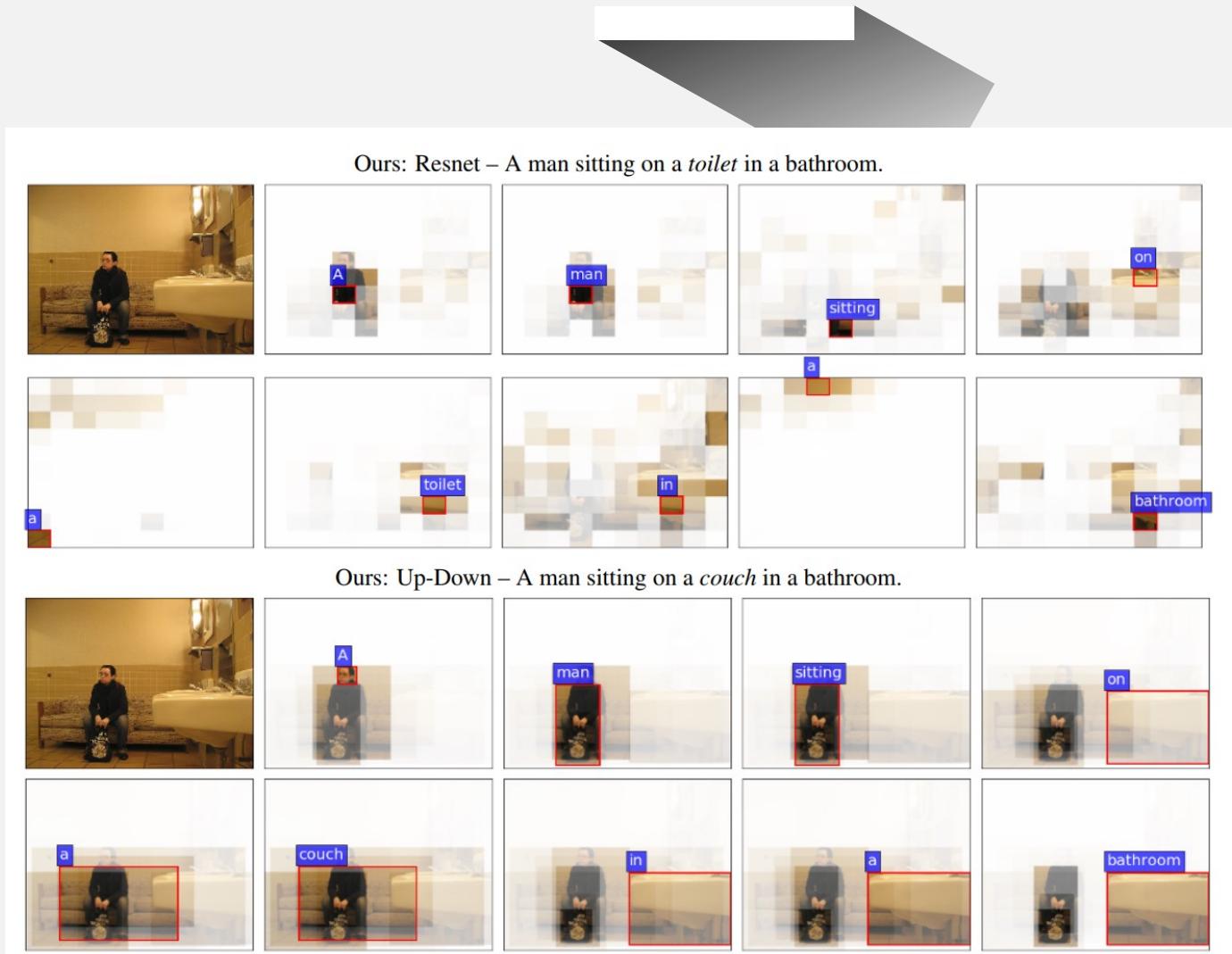
Bottom-Up attention model



	Yes/No	Number	Other	Overall
Prior [13]	61.20	0.36	1.17	25.98
Language-only [13]	67.01	31.55	27.37	44.26
d-LSTM+n-I [28, 13]	73.46	35.18	41.83	54.22
MCB [12, 13]	78.82	38.28	53.36	62.27
UPMC-LIP6	82.07	41.06	57.12	65.71
Athena	82.50	44.19	59.97	67.59
HDU-USYD-UNCC	84.50	45.39	59.01	68.09
Ours: Up-Down	86.60	48.64	61.15	70.34

Table 5. VQA v2.0 test-standard server accuracy, ranking our submission against recent published and unpublished work for each question type. Our approach, an ensemble of 30 models trained with different random seeds, outperforms all other leaderboard entries.

Bottom-Up attention model



Multi-attention network for one shot learning

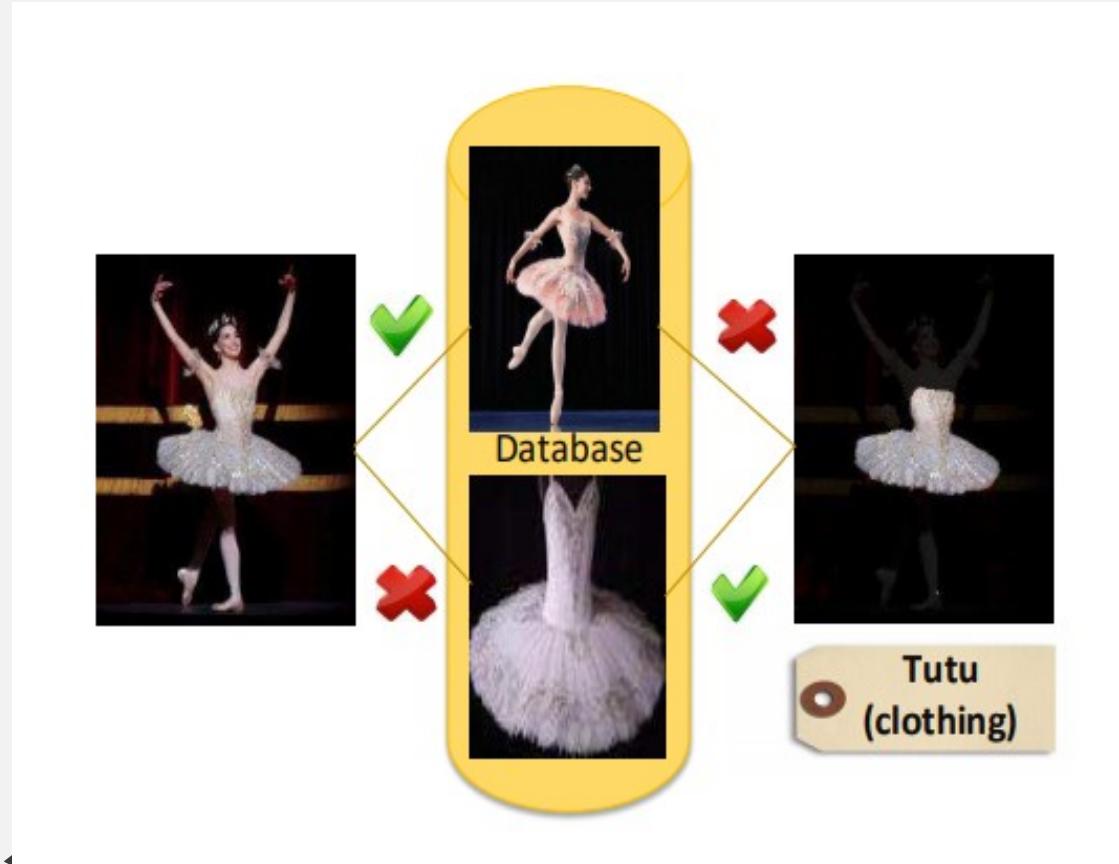


Car=class 1



Motorcycle=class2

Multi-attention network for one shot learning



Multi-attention network

We design a neural network architecture which takes the semantic embedding of the class tag to generate attention maps.

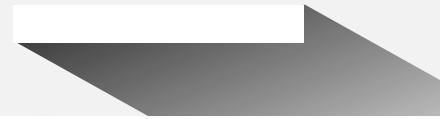


The use of the class tag

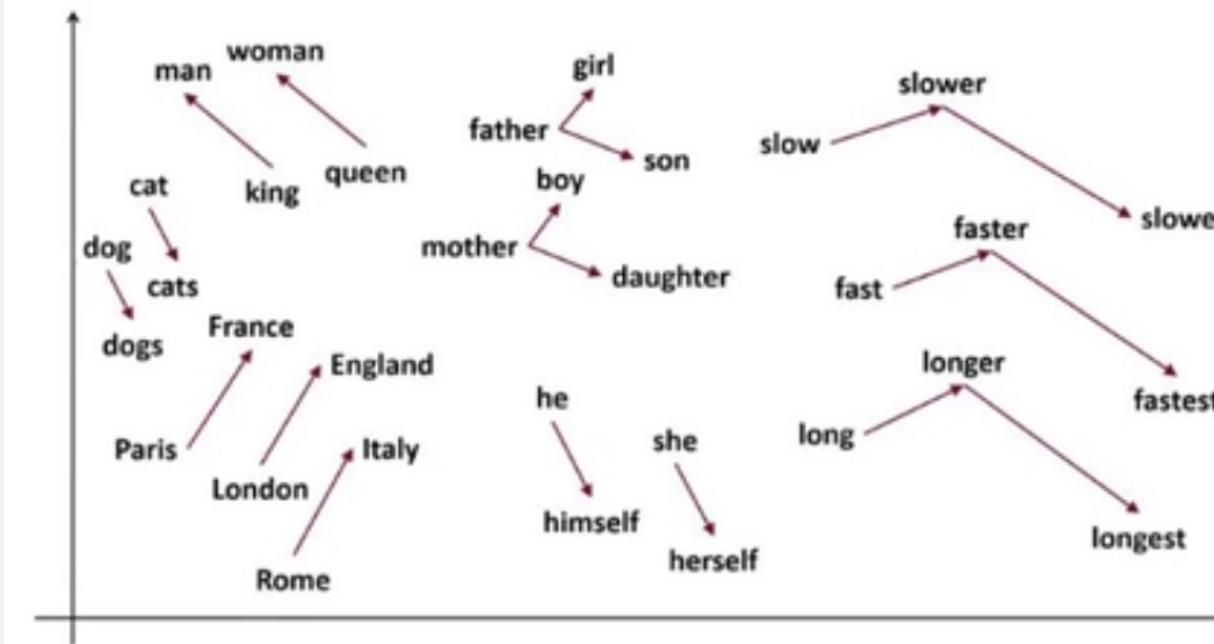
Visual attention



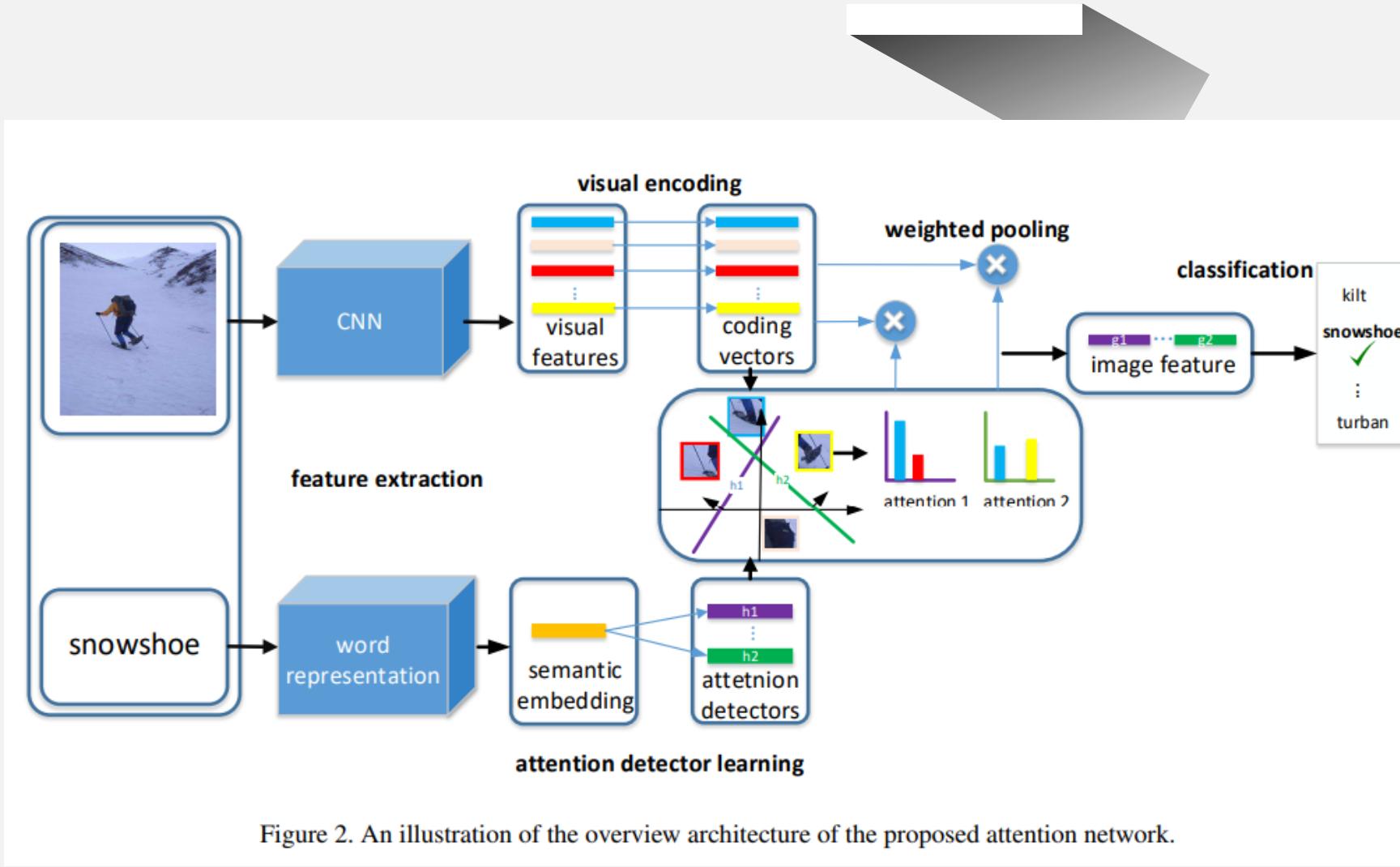
Multi-attention network for one shot learning



Distributed Word Embedding



Multi-attention network for one shot learning



$$\mathbf{v}_i = f(\mathbf{W}_v \mathbf{x}_i + \mathbf{b}_v),$$

$$\mathbf{h} = \mathbf{W}_s \mathbf{c} + \mathbf{b}_s,$$

$$a'_i = \mathbf{h}^\top \mathbf{v}_i.$$

$$a_i = \frac{b(a'_i)}{\sum_i b(a'_i)},$$

Figure 2. An illustration of the overview architecture of the proposed attention network.

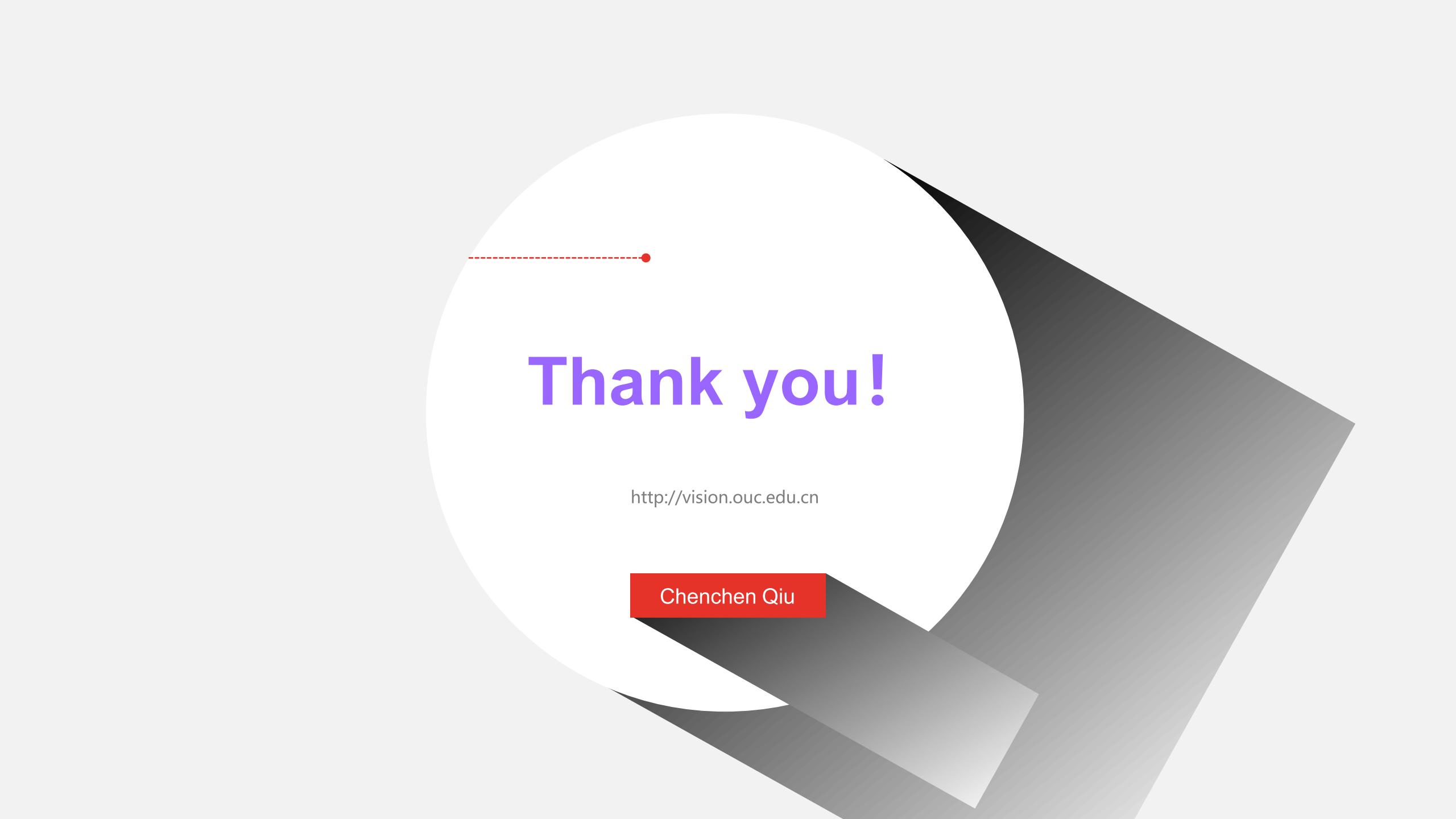
Multi-attention network for one shot learning

Table 1. Comparison of the attention network to alternative solutions on the Animal Dataset.

close-world	Global FC	68.9%
	SE	67.7%
	Zero shot learning	55.2%
	Attention (attribute)	72.4%
	Attention (word vector)	74.0%
	Multi-Attention (2)	82.4%
	Multi-Attention (5)	77.7%
open-world	Global FC	42.8%
	SE	42.1%
	Zero shot learning	15.3%
	Attention (attribute)	47.4%
	Attention (word vector)	49.0%
	Multi-Attention (2)	55.7%
	Multi-Attention (5)	56.8%

Table 2. Comparison of the attention network to alternative solutions on the Artifact Dataset.

close-world	SE (256D)	27.8%
	SE (512D)	28.2%
	SE + Joint Bayesian	31.5%
	Attention	34.5%
	Attention + Joint Bayesian	36.8%
open-world	Multi-Attention (2)	50.5%
	SE (256D)	11.6%
	SE (512D)	12.2%
	SE + Joint Bayesian	11.4%
	Attention	15.2%
	Multi-Attention (2)	28.2%



Thank you!

<http://vision.ouc.edu.cn>

Chenchen Qiu