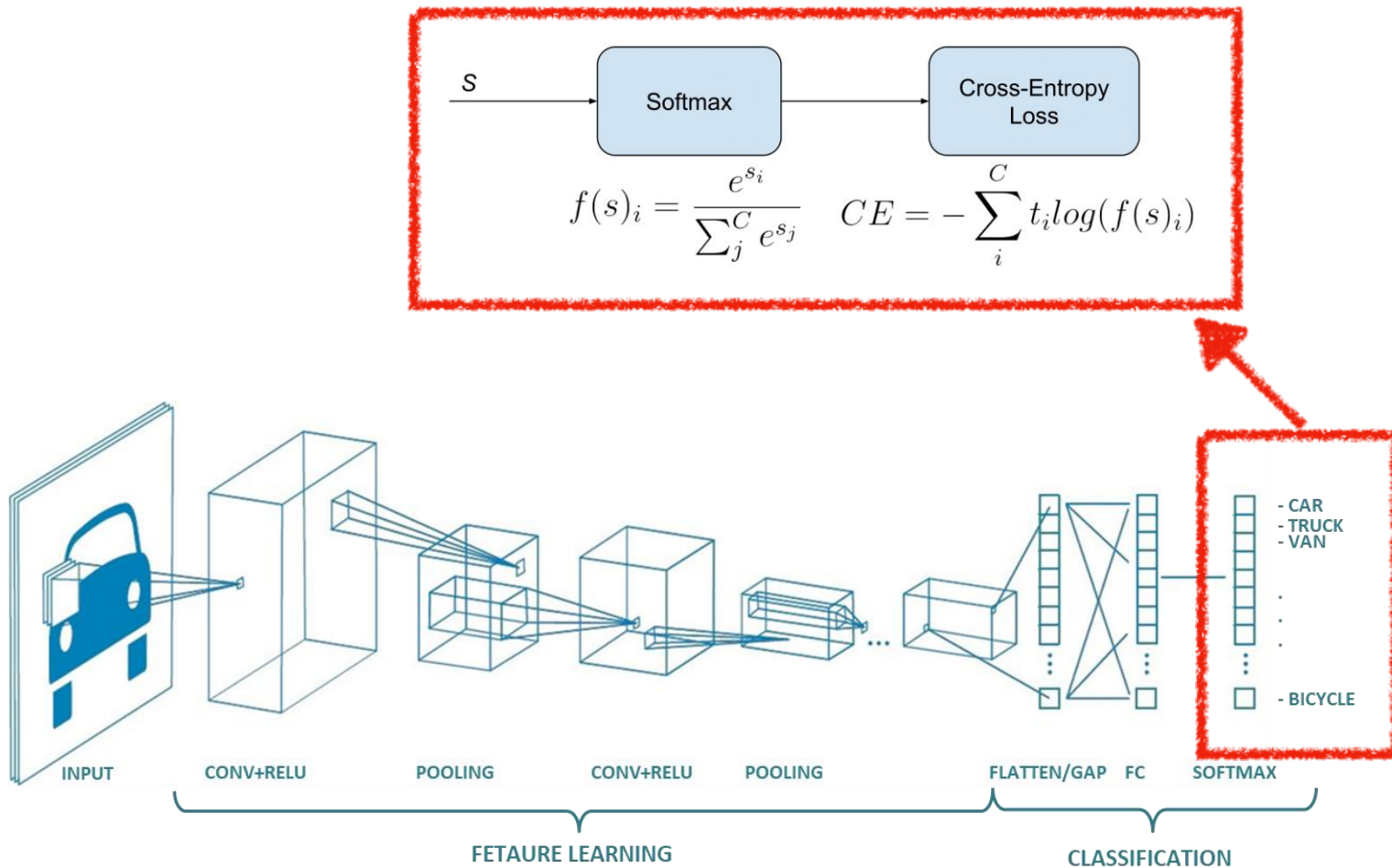


Neural Networks and Deep Learning

Metric Learning

Supervised Learning so far



Few-shot Learning

Supervised ML focused on learning from limited number of examples, inspired by human/animal ability to rapidly generalize from few examples.

Typical scenarios:

- Learning for rare cases – when obtaining labelled data is hard or impossible
- Reducing data gathering effort and computational cost

Variations:

- (<50 , e.g. 5,10)-shot learning: General case, where only few examples are given
- 1-shot learning: Extreme case, where only one data example is given (e.g. 1 image per class)
- 0-shot learning: Instead of image, we have description of new class

e.g.: person re-id hundreds of pedestrians, 2-5 images per pedestrian

Few-shot learning



Ned Stark



Robert Baratheon



Daenerys Targaryen

1-shot learning



Ned Stark



Robert Baratheon



Daenerys Targaryen

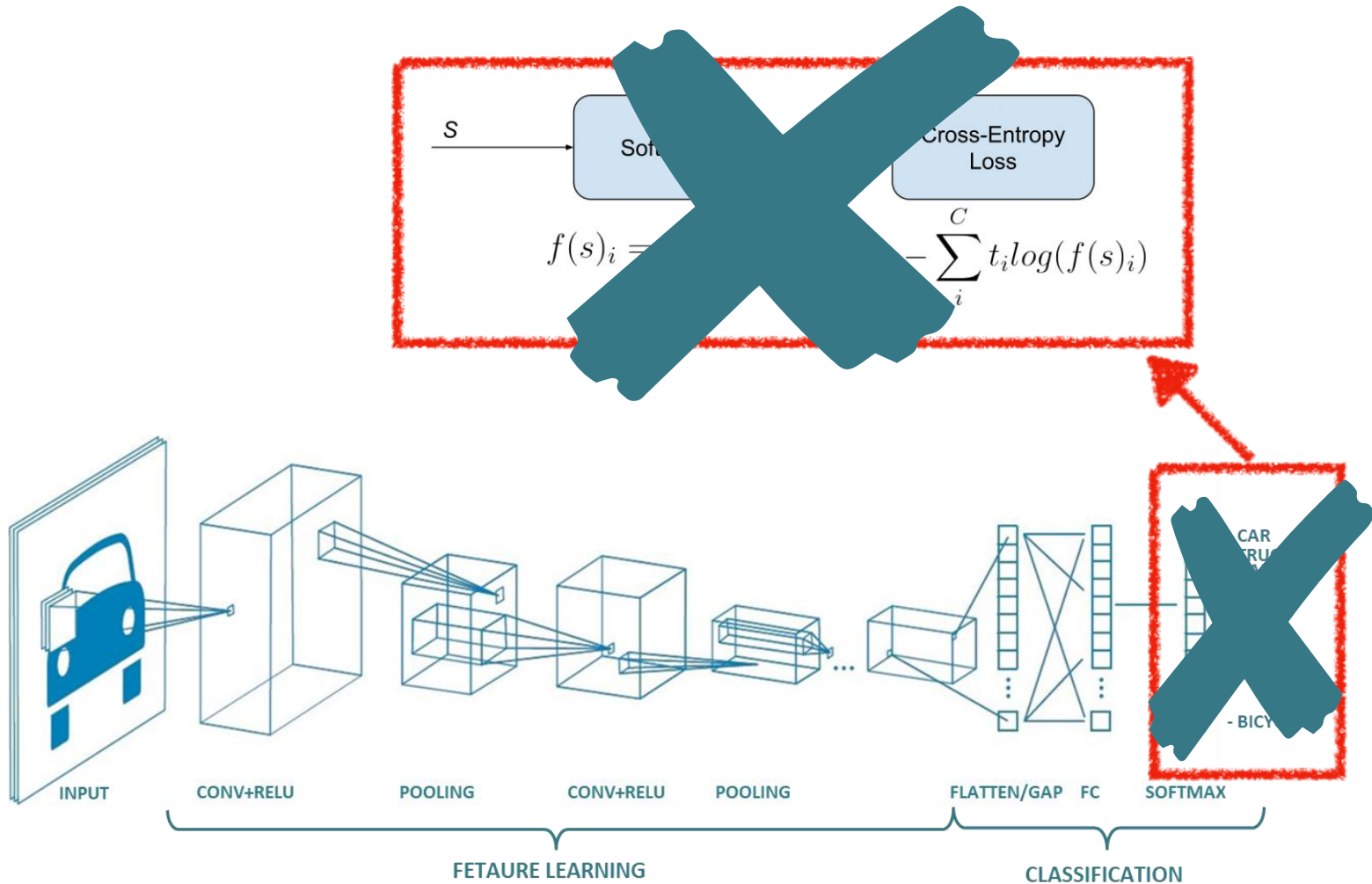
0-shot learning



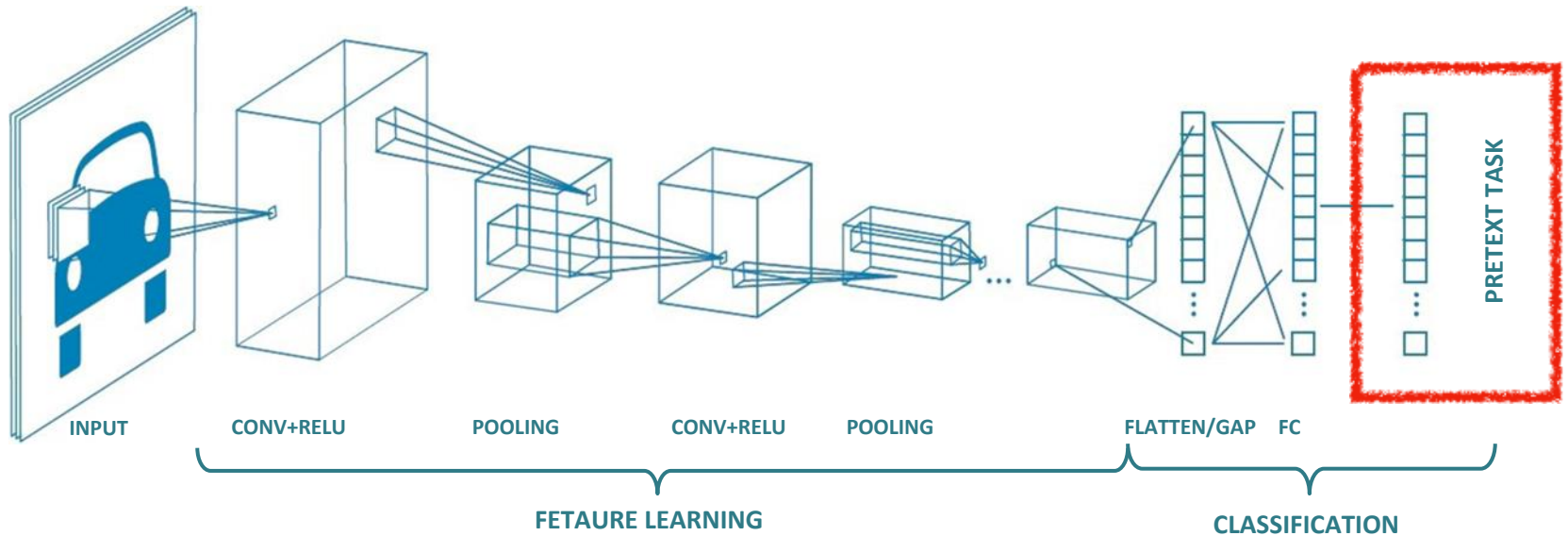
“Blond, long hair, blue eyes”



The Self-Supervised Revolution



The Self-Supervised Revolution



Focus on Feature Learning for a 'random' given task (Pretext Task)

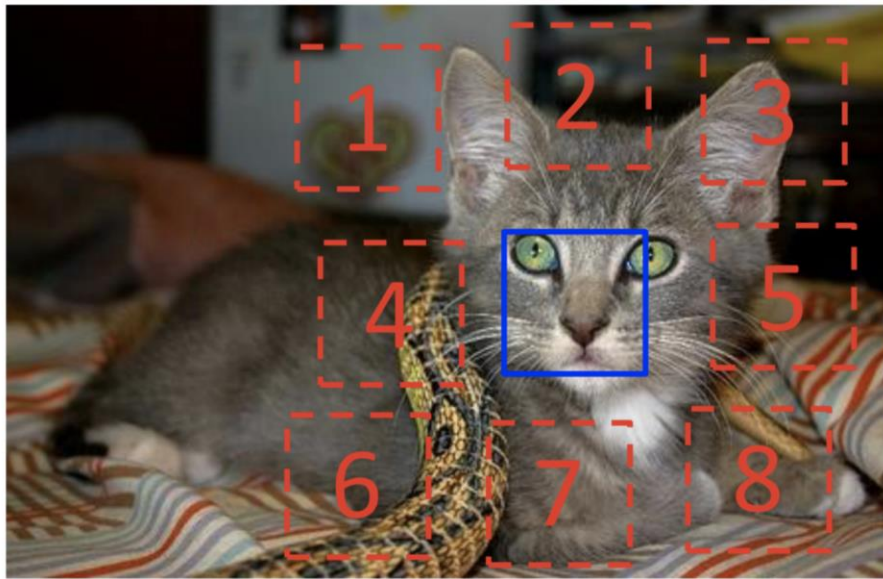
e.g:

- Center word, next sentence, sentence swapping, ...

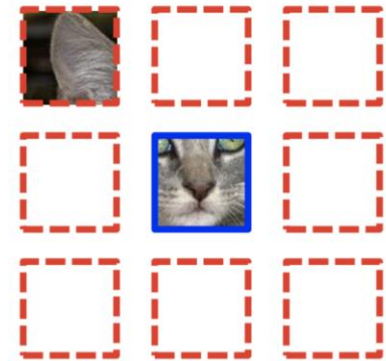
Pretext tasks: Tons of data, easy to extract goal

Learning from pretext tasks

Jigsaw puzzle: what is the relative position



Example:



$$X = (\text{cat face}, \text{snake}); Y = 3$$

Question 1:



?

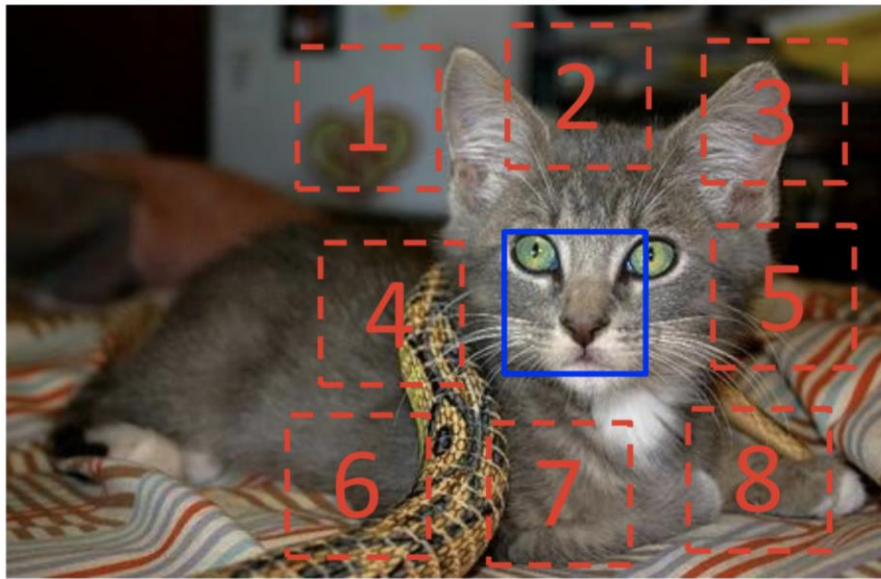
Question 2:



?

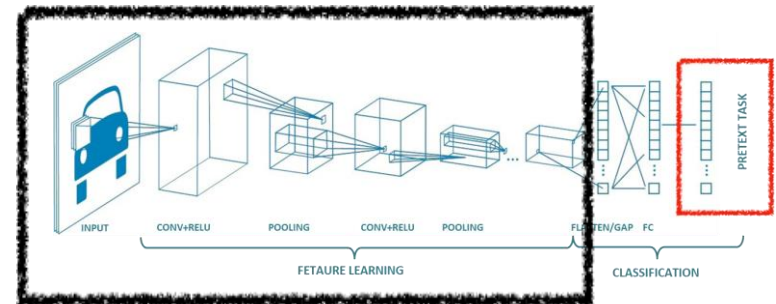
Learning from pretext tasks

Jigsaw puzzle: what is the relative position

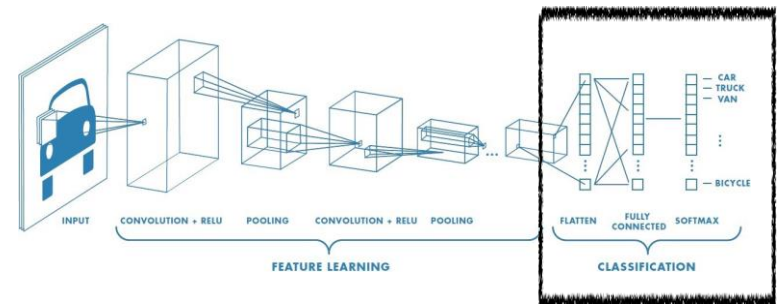


$$X = \left(\begin{array}{c} \text{cat face} \\ \text{cat ear} \end{array} \right); Y = 3$$

1. Train Feature Representation

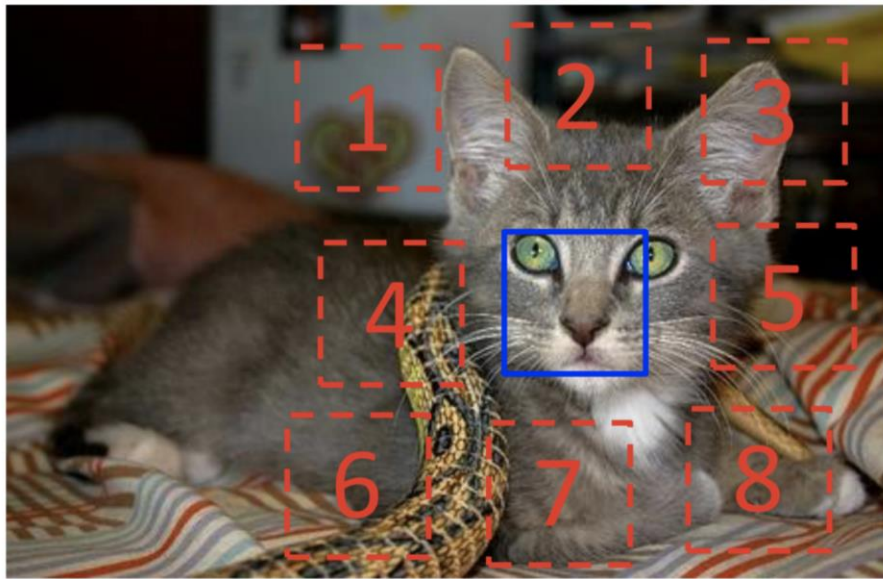


2. Finetune classifier w/ some labeled data



Learning from pretext tasks

Jigsaw puzzle: what is the relative position

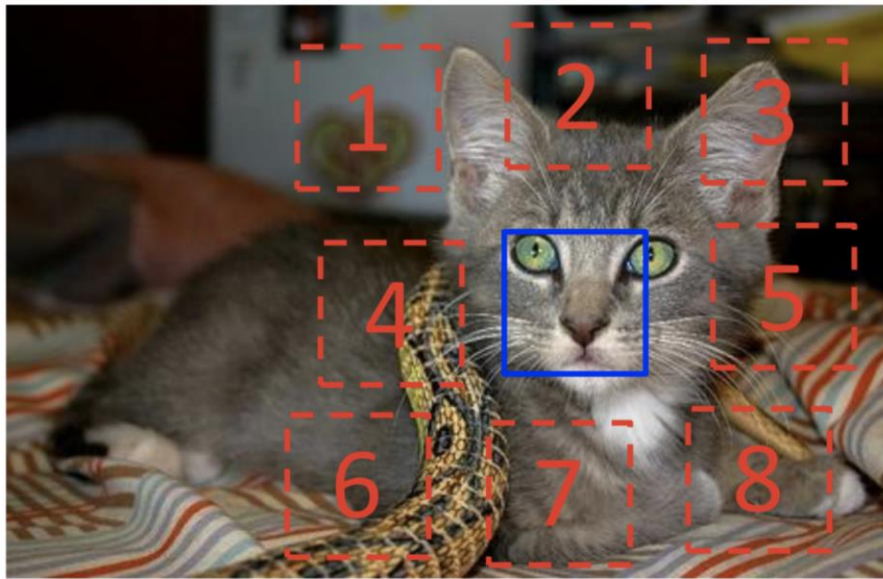


Results are good when a **lot of unlabeled data** is available and **little labeled data** available

$$X = (\text{cat face}, \text{cat ear}); Y = 3$$

Learning from pretext tasks

Jigsaw puzzle: what is the relative position



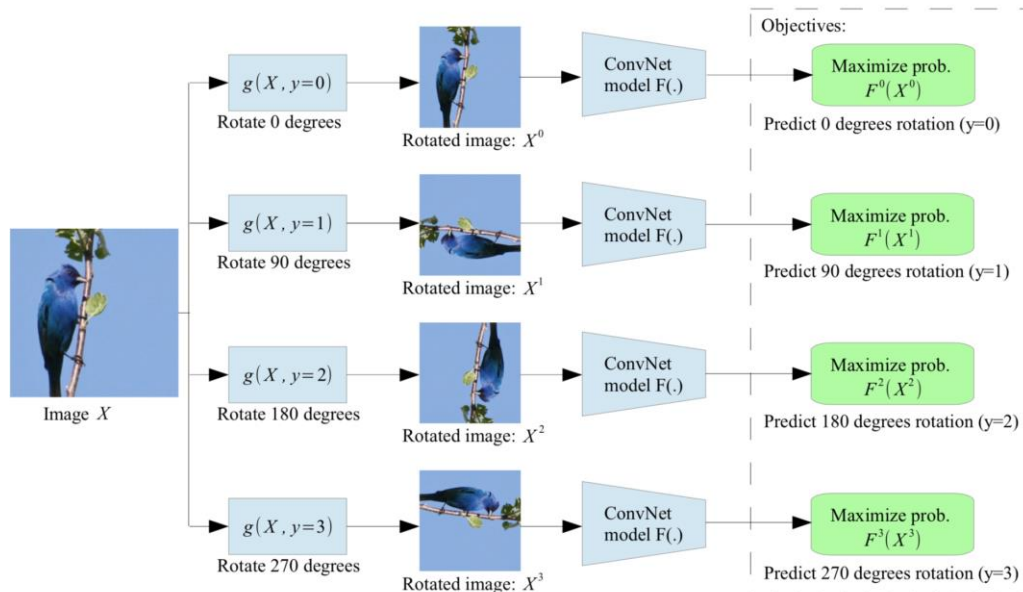
Shortcut problem: The Network always tries to “cheat”.

Ex: network can use the camera lens distortion

(Doersch et al. found a solution for this particular problem)

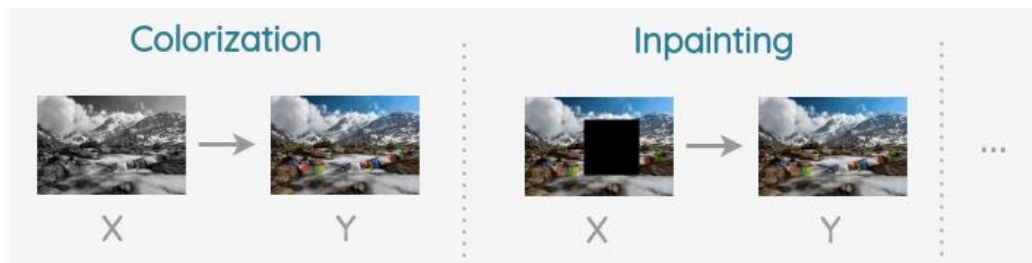
$$X = \left(\begin{array}{c} \text{cat face} \\ \text{cat ear} \end{array} \right); Y = 3$$

Learning from pretext tasks



There are other pretext tasks

- Rotation
- In-painting
- Colorization
- Up to your imagination...



[Gidaris et al., 2018]

Pretext Tasks for NLP

- Center Word Prediction (CBOW)
- Neighbour Word Prediction (skip-gram)
- Neighbour Sentence Prediction (skip-gram, sentence level)
- Auto-regressive Language Modeling (predict next word)
- Masked Language Modeling (e.g. BERT)
- Next Sentence Prediction (classify if sentence comes from same or different document)
- Sentence Order Prediction (classify if correct order or not)
- Sentence Permutation (re-order messed up sentences)
- Emoji Prediction



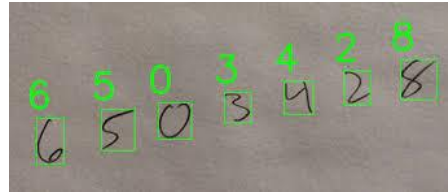
Does any embedded space suffice?

Most ML algorithms are based on comparing samples and to decide if they 'are the same', and to define borders between different concepts.

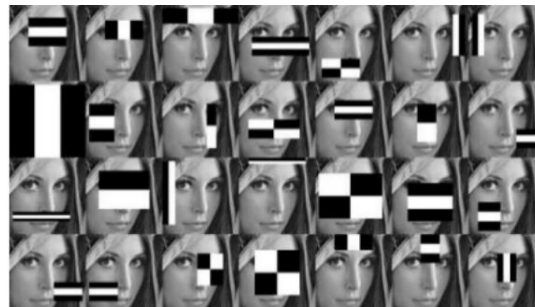
Goal: define the **similarity** between subjects.

Need for Similarity Measures

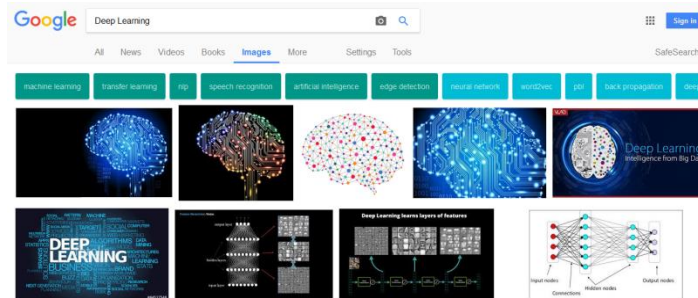
- Recognizing handwriting in checks.



- Automatic detection of faces in a camera image.



- Search Engines, such as Google, matching a **query** (could be text, image, etc.) with a set of **indexed documents** on the web.



Notion of a Metric

- A **Metric** is a function that quantifies a “distance” between every pair of elements in a set, thus inducing a measure of similarity.
- A metric $f(x,y)$ must satisfy the following properties for all x, y, z belonging to the set:
 - *Non-negativity*: $f(x, y) \geq 0$
 - *Identity of Discernible*: $f(x, y) = 0 \iff x = y$
 - *Symmetry*: $f(x, y) = f(y, x)$
 - *Triangle Inequality*: $f(x, z) \leq f(x, y) + f(y, z)$

Types of Metrics

In broad strokes metrics are of two kinds:

- **Pre-defined Metrics:** Metrics which are fully specified without the knowledge of data.

E.g. Euclidian Distance: $f(\mathbf{x}, \mathbf{y}) = (\mathbf{x} - \mathbf{y})^T(\mathbf{x} - \mathbf{y})$

- **Learned Metrics:** Metrics which can only be defined with the **knowledge** of the **data**.

E.g. Mahalanobis Distance: $f(\mathbf{x}, \mathbf{y}) = (\mathbf{x} - \mathbf{y})^T \mathbf{M} (\mathbf{x} - \mathbf{y})$; where **M** is a matrix that is estimated from the data.

Learned Metrics are of two types:

- **Unsupervised** : Use unlabeled data
- **Supervised** : Use labeled data

DEEP LEARNING TECHNIQUES

Deep Learning to the Rescue!

CNNs can **jointly optimize** the representation of the input data conditioned on the “similarity” measure being used, aka end-to-end learning.

State the Problem

Input: Given a pair of input images, we want to know how “similar” they are to each other.

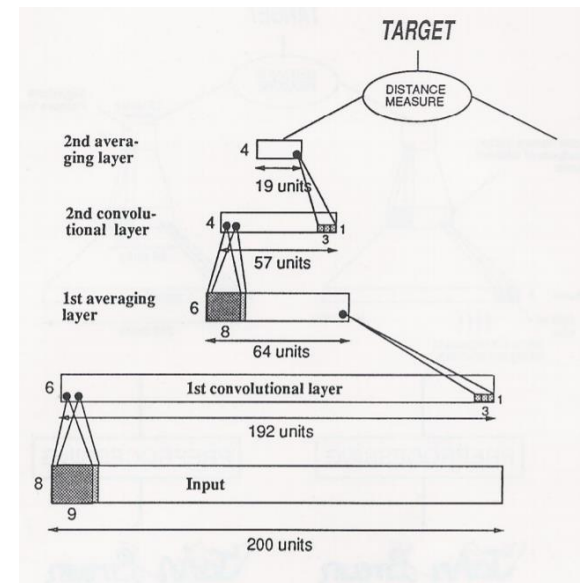
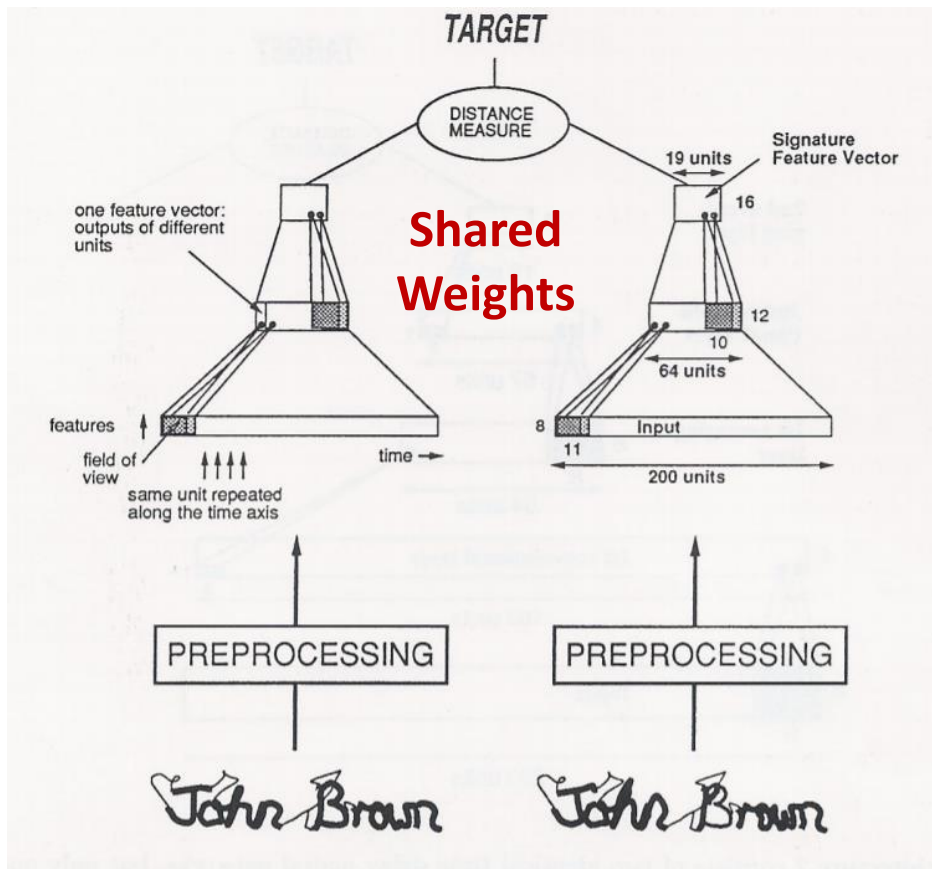
Output: The output can take a variety of forms:

- Either a binary label, i.e. 0 (same) or 1 (different).
- A **Real** number indicating how similar a pair of images are.



Typical Siamese CNN

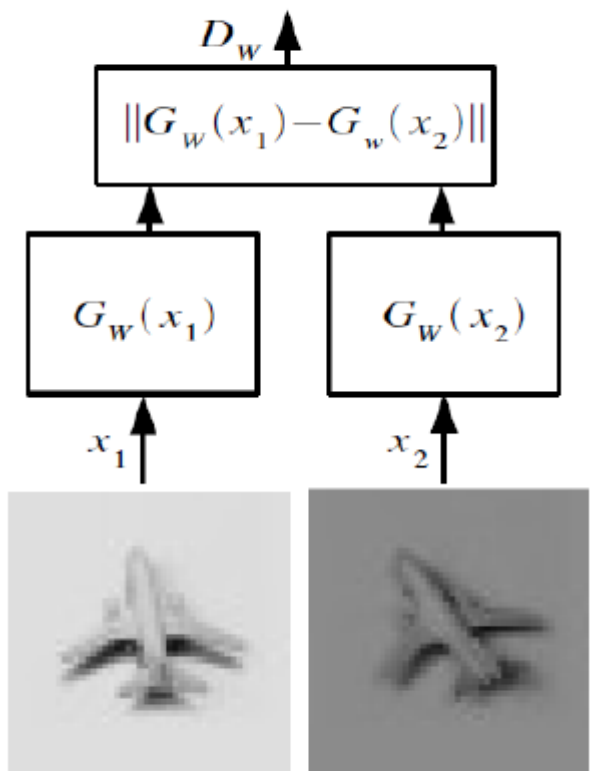
- **Input:** A pair of input signatures.
- **Output (Target):** A label, **0** for similar, **1** else.



Bromley, J., Bentz, J.W., Bottou, L., Guyon, I., LeCun, Y., Moore, C., Säckinger, E. and Shah, R., 1993. Signature Verification Using A "Siamese" Time Delay Neural Network. *IJPRAI*, 7(4), pp.669-688. 20

Siamese CNN – Loss Function

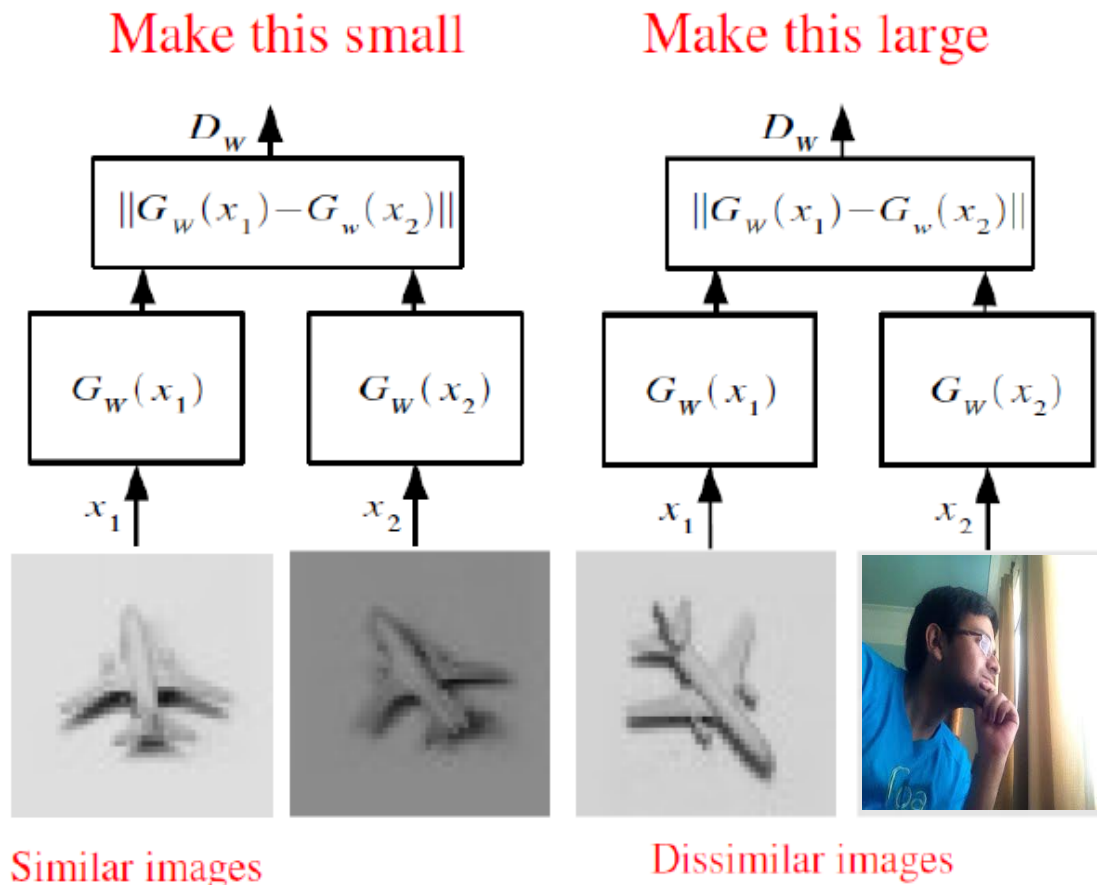
Make this small



Similar images

- Is there a problem with this formulation?
 - The model could learn to embed every input to the same point, i.e. predict a **constant** as **output**.
 - In such a case, every pair of input would be categorized as a positive pair.

Siamese CNN – Loss Function

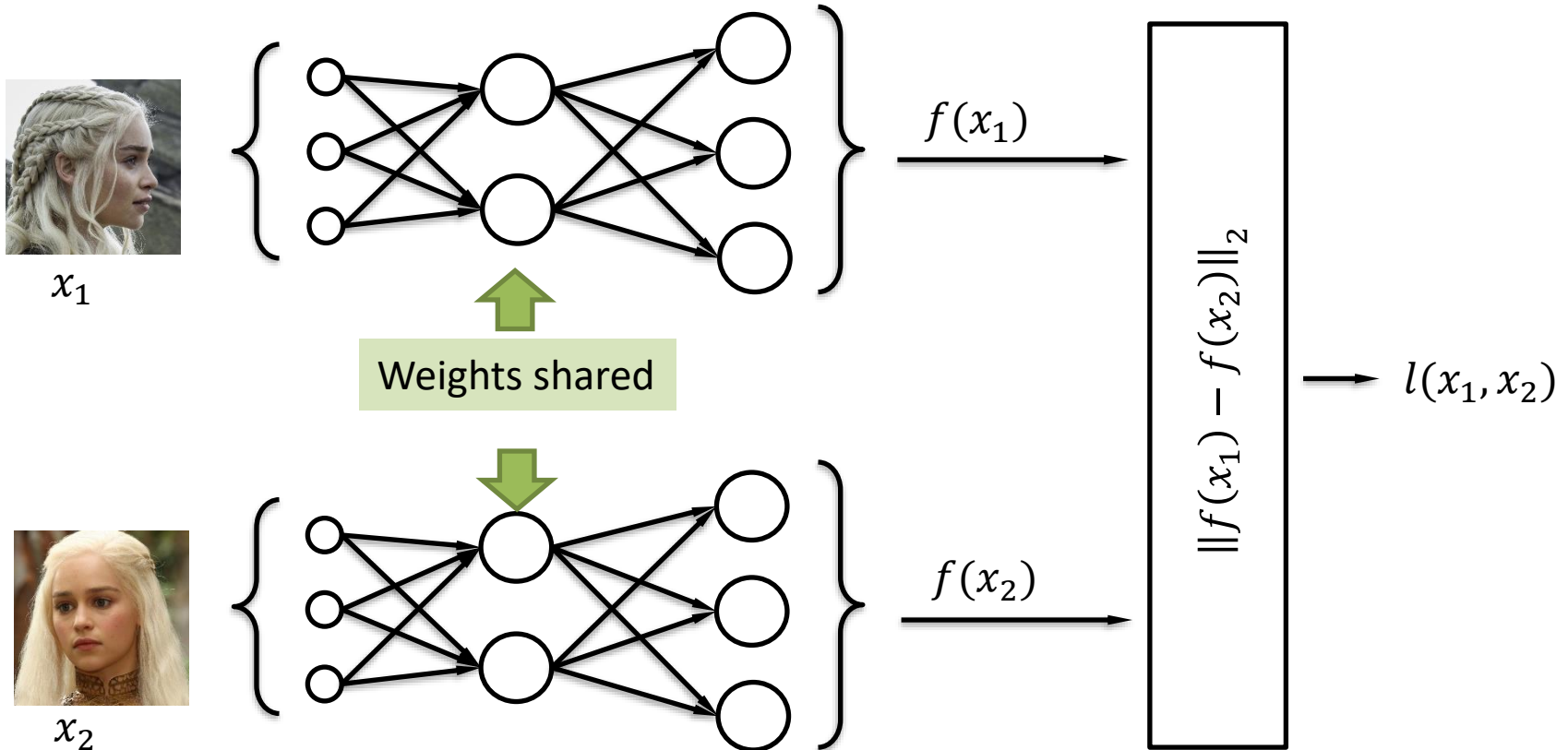


The final loss is defined as :

$$L = \sum \text{loss of positive pairs} + \sum \text{loss of negative pairs}$$

Siamese networks

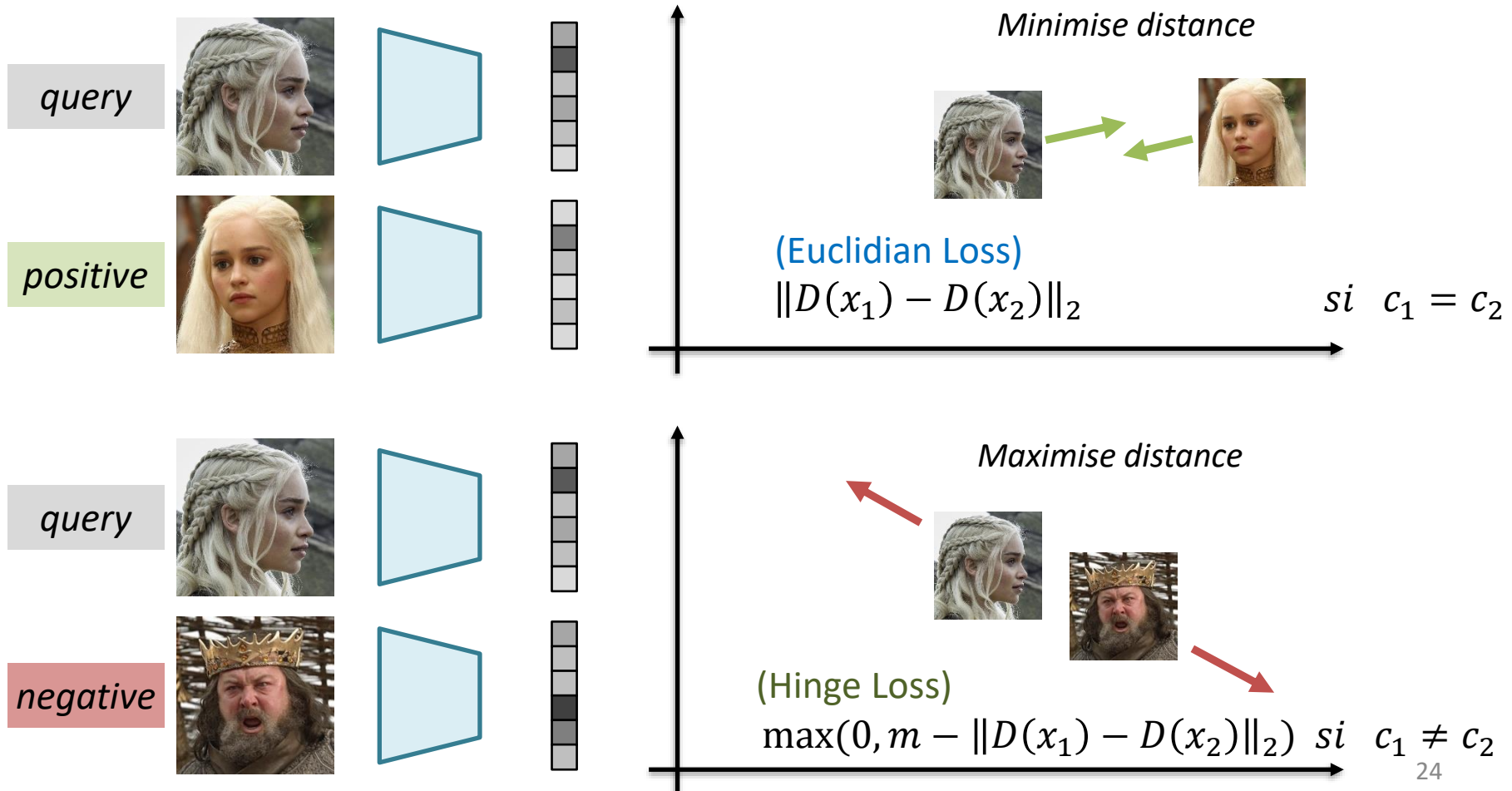
Two parallel feed-forward networks, with **shared weights**. It's the same network actually encoding two different inputs. Shared weights, like in RNNs.



$$l(x_1, x_2) = \begin{cases} \|D(x_1) - D(x_2)\|_2 & \text{si } c_1 = c_2 \\ \max(0, m - \|D(x_1) - D(x_2)\|_2) & \text{si } c_1 \neq c_2 \end{cases}$$

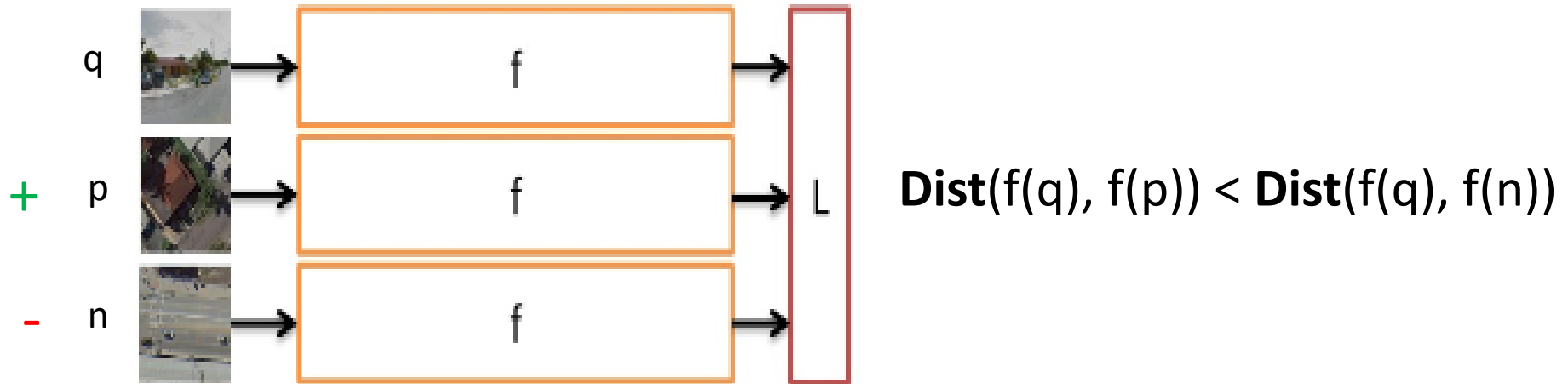
Pairwise comparisons

Given a pair of input images, we want to place them closer together if they are *similar*, and far from each other if they are not



Siamese CNN – Variants

TRIPLET NETWORK



- Compare triplets in one go.
- Check if the sample in the **topmost** channel, is more similar to the one in the **middle** or the one in the **bottom**.
- Allows us to learn ranking between samples.

Siamese CNN – Loss Function

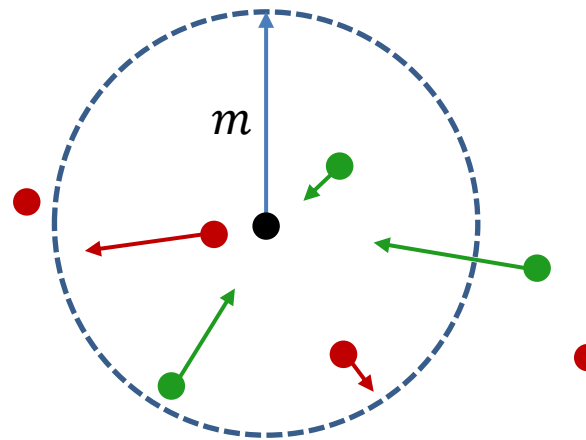
Siamese

$$l(x_q, x_a) = \begin{cases} \|D(x_q) - D(x_a)\|_2 & \text{if } a = p \\ \max(0, m - \|D(x_q) - D(x_a)\|_2) & \text{if } a = n \end{cases}$$

Triplet

Contrastive Loss

$$l(x_q, x_p, x_n) = \|D(x_q) - D(x_p)\|_2 + \max\left(0, m - \|D(x_q) - D(x_n)\|_2\right)$$



- Query sample
- Negative pair
- Positive pair

MINING

Types of negatives

Easy negatives:

$$d(x_q, x_n) > d(x_q, x_p) + m$$

The negative sample is already sufficiently distant to the query (anchor) sample with respect to the positive sample in the embedding space. The loss is 0 and the net parameters are not updated.

Hard Negatives:

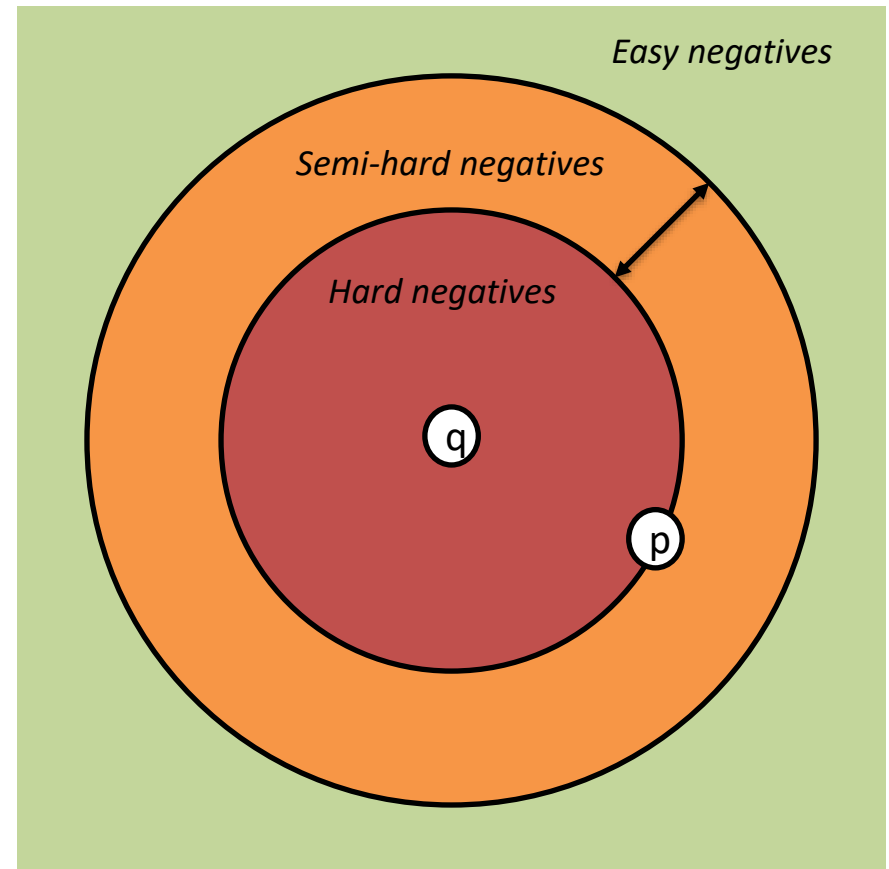
$$d(x_q, x_n) < d(x_q, x_p)$$

The negative sample is closer to the anchor than the positive. The loss is positive and greater than m .

Semi-Hard Negatives:

$$d(x_q, x_p) < d(x_q, x_n) < d(x_q, x_p) + m$$

The negative sample is more distant to the anchor than the positive, but the distance is not greater than the margin, so the loss is still positive (and smaller than m).



Mining training data

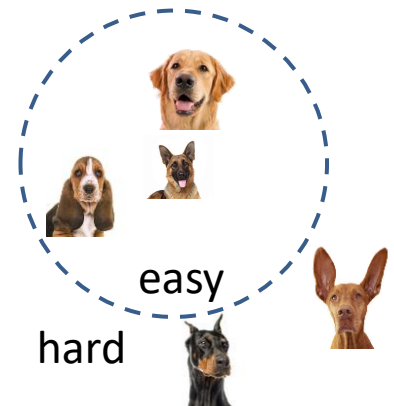
select a random query (e.g. dog):



Positives



Take the hardest positive samples



Mining training data

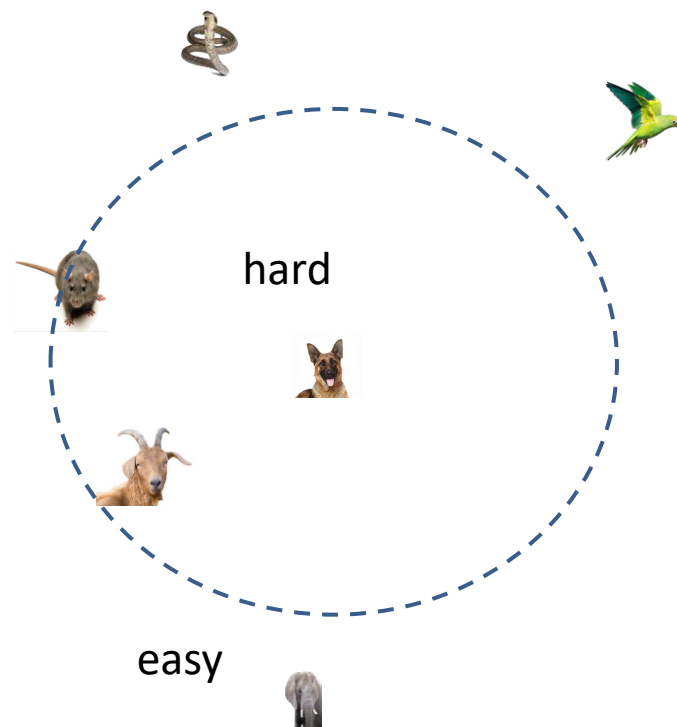
select a random query (dog):



Positives



Negatives



Mining training data

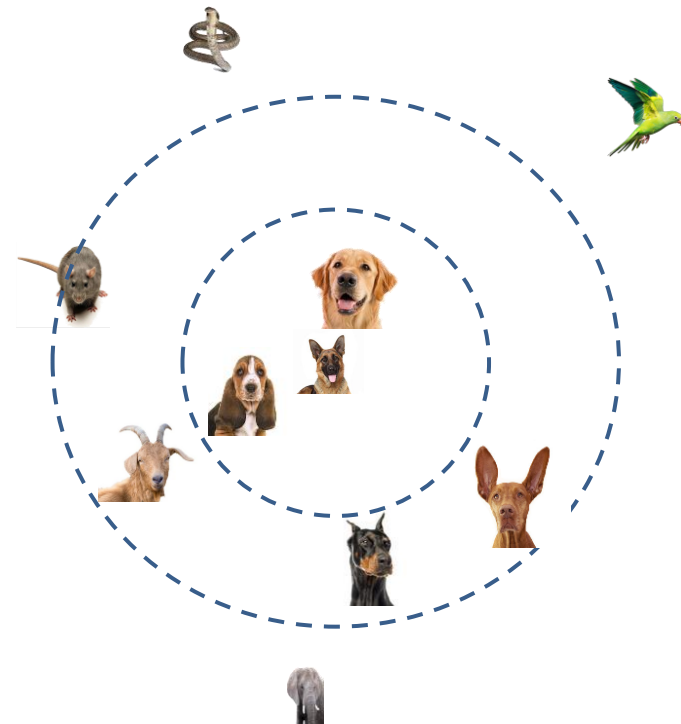
select a random query (dog):



Positives



Negatives



Mining training data

select a random query (dog):



Fix a positive



DOG

Negatives

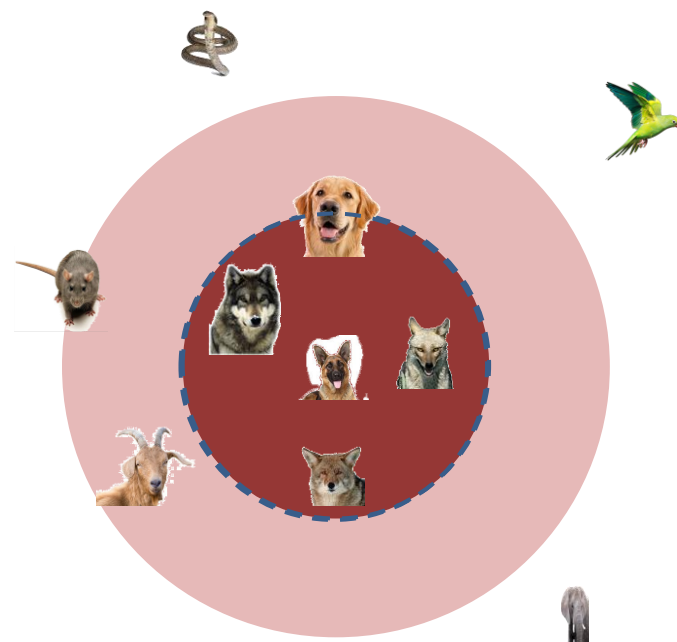


WOLF

COYOTE

JACKAL

FOX



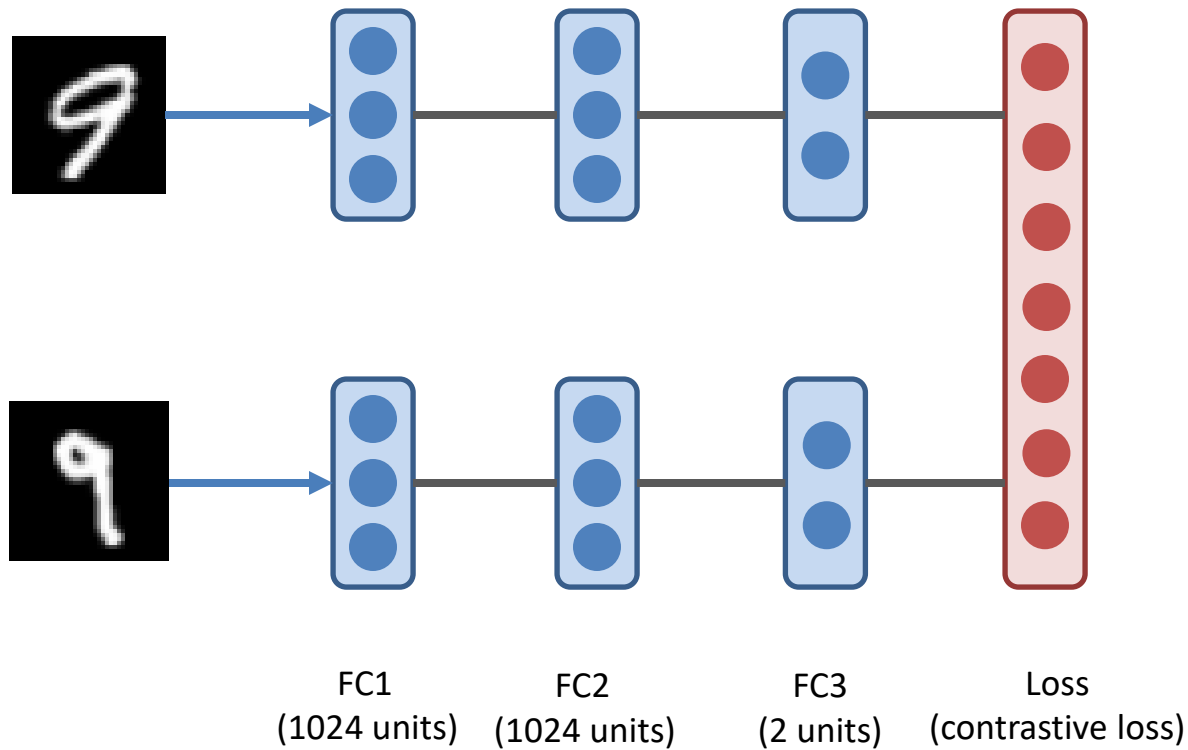
Hard negatives

Semi-hard negatives

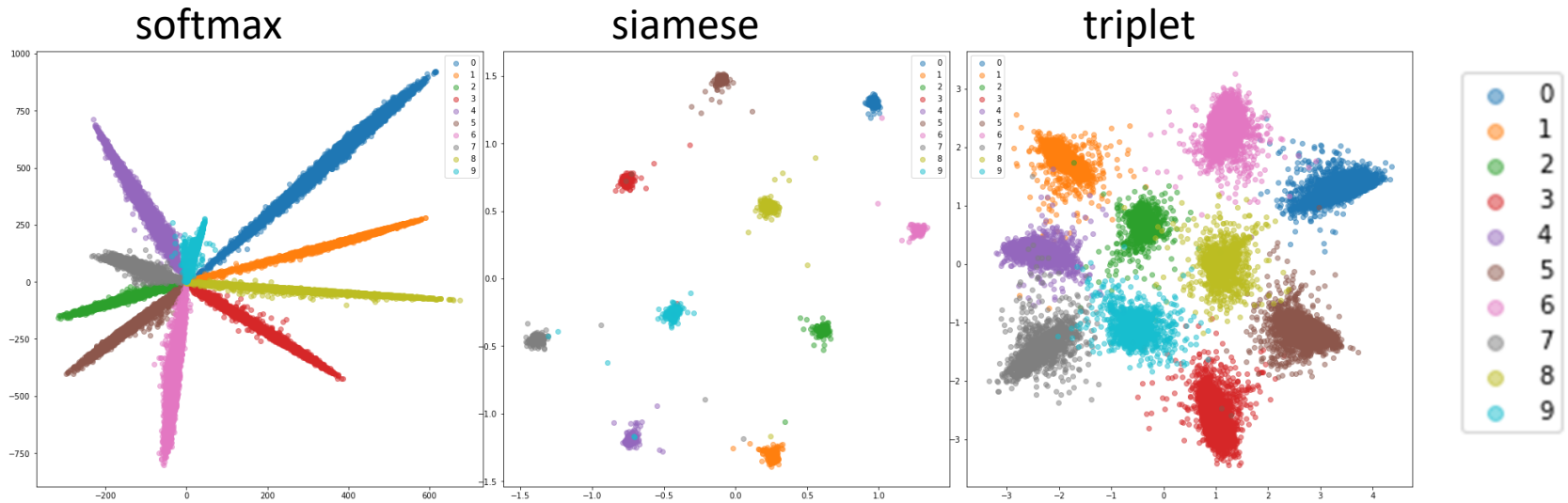
ARE EMBEDDINGS MEANINGFUL?

Metric

MNIST Digit Similarity Assessment



Mnist test



Network architecture

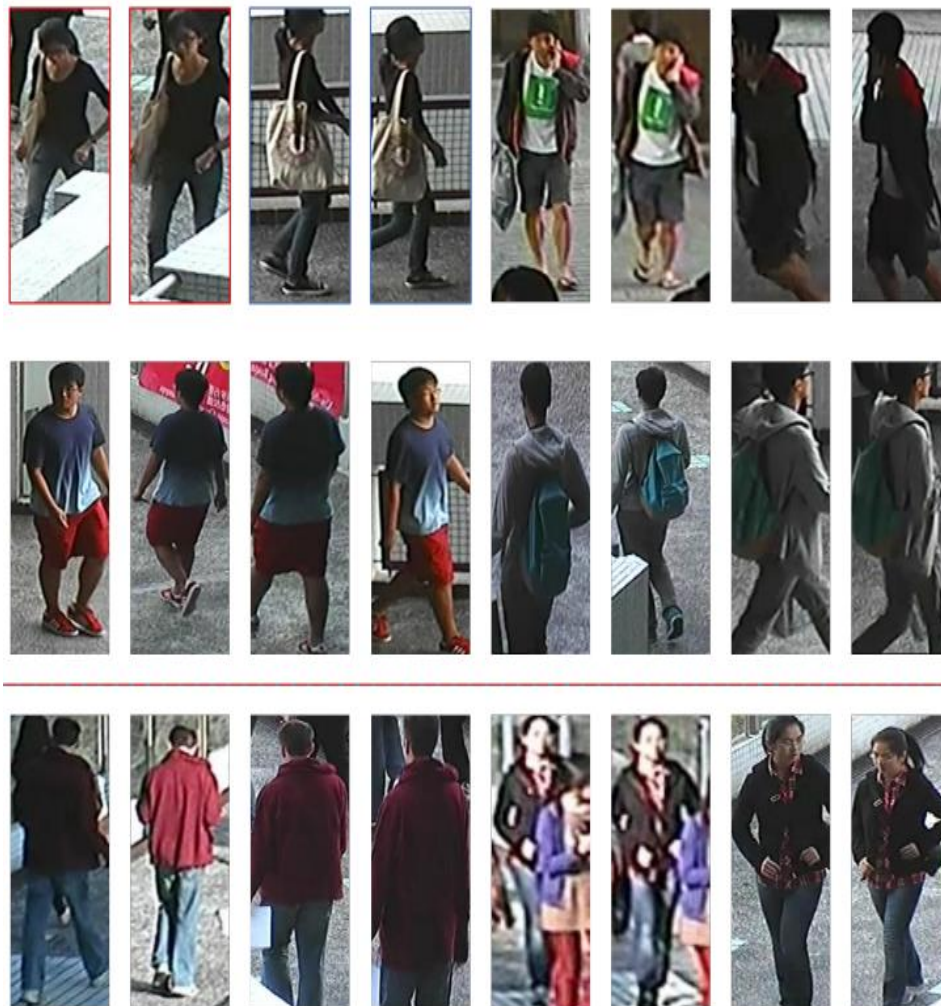
conv 32 5x5 -> PReLU -> MaxPool 2x2 ->

conv 64 5x5 -> PReLU -> MaxPool 2x2 ->

Dense 256 -> PReLU -> Dense 256 -> PReLU -> Dense 2

APPLICATIONS

Person Re-Identification



Person Re-Identification

**True
positive**



**True
negative**



MULTIMODAL LEARNING



Predicting the “topic” of an image

VISUAL INFORMATION



TEXTUAL INFORMATION

Scene Text Annotations

TELEPHONE

Image Captions

A big red telephone booth that a man is standing in.

A person standing inside of a telephone booth.

This is an image of a man in a phone booth.

A man is standing in a red phone booth.

A man is using a phone to in a phone booth.

TRAINING SAMPLE

Predicting the “topic” of an image

VISUAL INFORMATION



TEXTUAL INFORMATION

Scene Text Annotations

TELEPHONE

Image Captions

A big red telephone booth that a man is standing in.

A person standing inside of a telephone booth.

This is an image of a man in a phone booth.

A man is standing in a red phone booth.

A man is using a phone to in a phone booth.

Topic Modelling (LDA)

LDA Discovered Topics

university	0.04
Student	0.02
President	0.01
...	

Album	0.02
Band	0.01
Music	0.01
...	

⋮

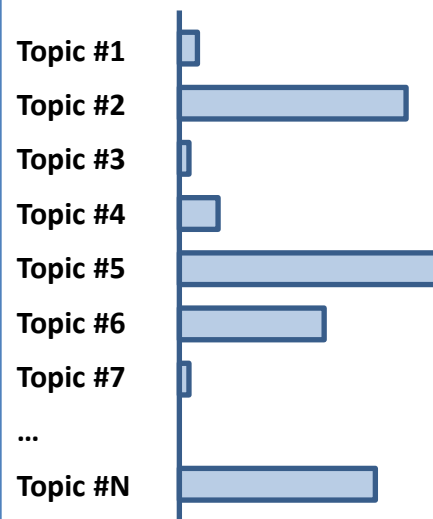
Car	0.04
Model	0.02
Engine	0.01
...	

Predicting the “topic” of an image

VISUAL INFORMATION



P (TOPIC | TEXT)



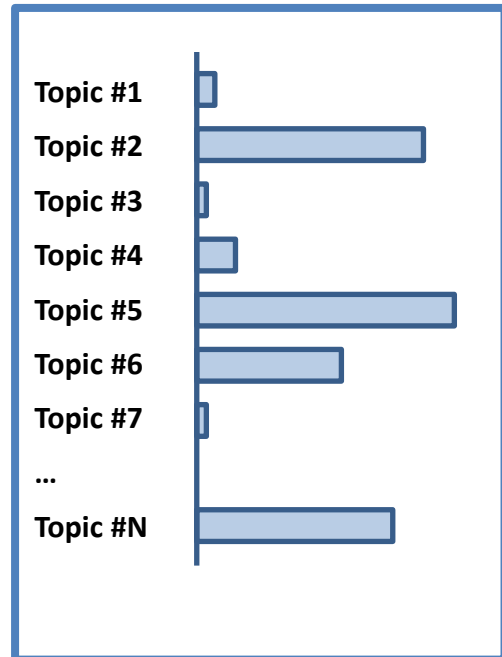
Predicting the “topic” of an image

VISUAL INFORMATION



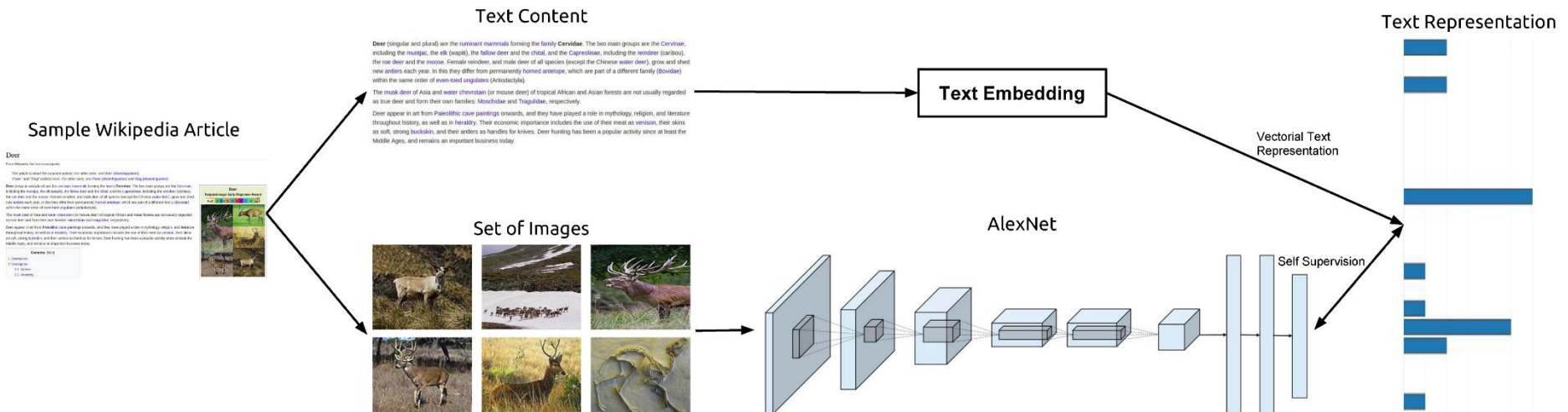
DNN

$P(\text{TOPIC} \mid \text{TEXT})$



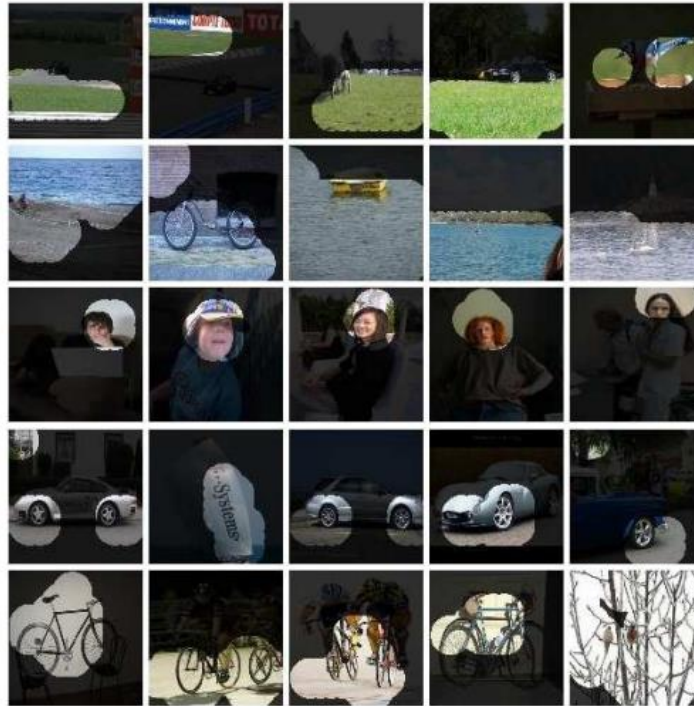
Learning to understand images by reading the Wikipedia

Task: Look at the image and predict what kind of article (topic) it illustrates



Learning to understand images by reading the Wikipedia

Features learnt are meaningful. Units are selective to generic textures (grass, water) or shapes, objects and object-parts



Self-supervised learning from Web Data



Wikipedia:

1.7M articles in English with
4.2M associated illustrative
images.



WIKIPEDIA
The Free Encyclopedia

WebVision:

2.4M Flickr and Google
images associated to
ImageNet classes.

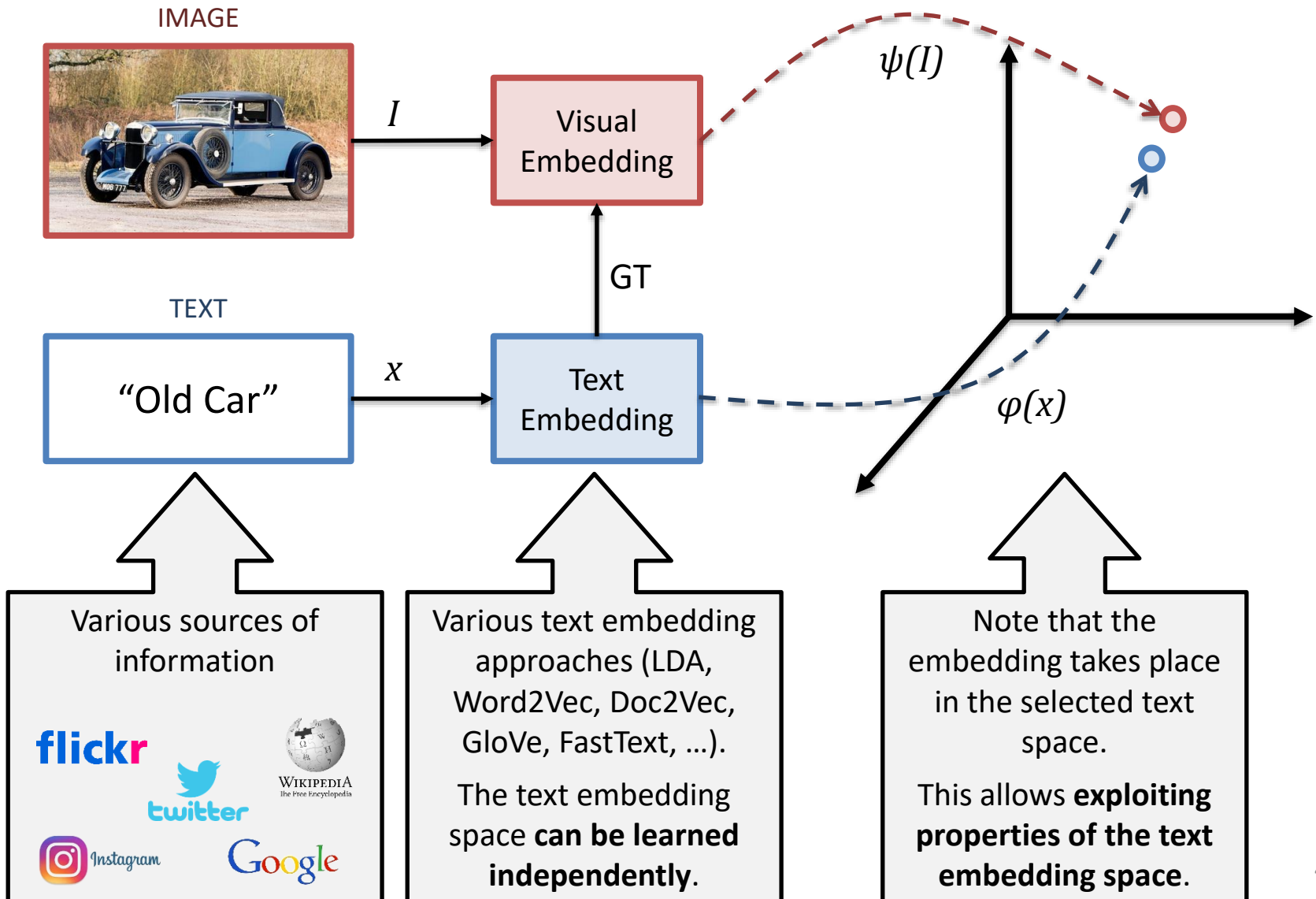


InstaCities1M:

1M Instagram images
associated with one of the 10
most populated English
speaking cities.



Self-supervised learning from Web Data



Text-based semantic retrieval

“wild”



“happy”



“monday”



Model trained with Word2Vec on InstaCites1M

Text-based semantic retrieval

“haircut”



Text-based semantic retrieval

“haircut”



“haircut”
+
“man”



“haircut”
+
“woman”



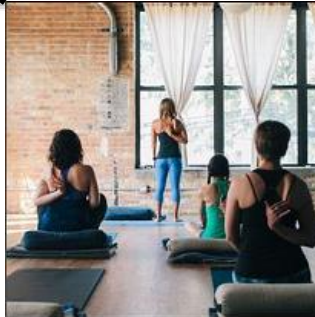
“haircut”
-
“man”



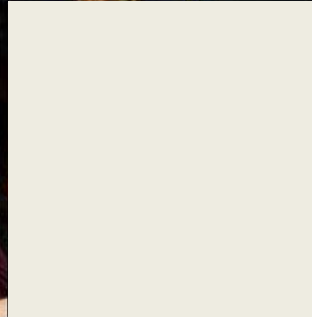
“people”

“people” + “art”

“art”



“people”
+
“dog”



“art”
+
“city”



“dog”

“dog” + “city”

“city”





-wedding



+animal



-sea

-sunset



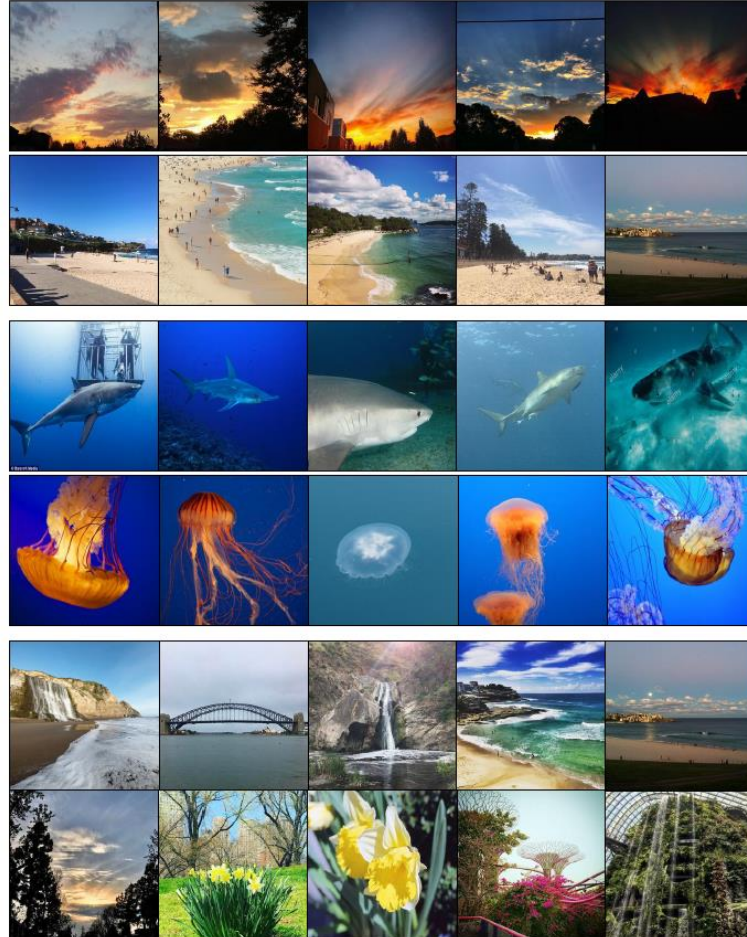
+tooth

+dangerous



-forest

-river



Scene Text Visual Question Answering



Question

What is written on the blue shirt the boy is wearing?

Answer

I Think Somebody Needs A Tickle

Reasoning Capacities

Object recognition,
Action recognition

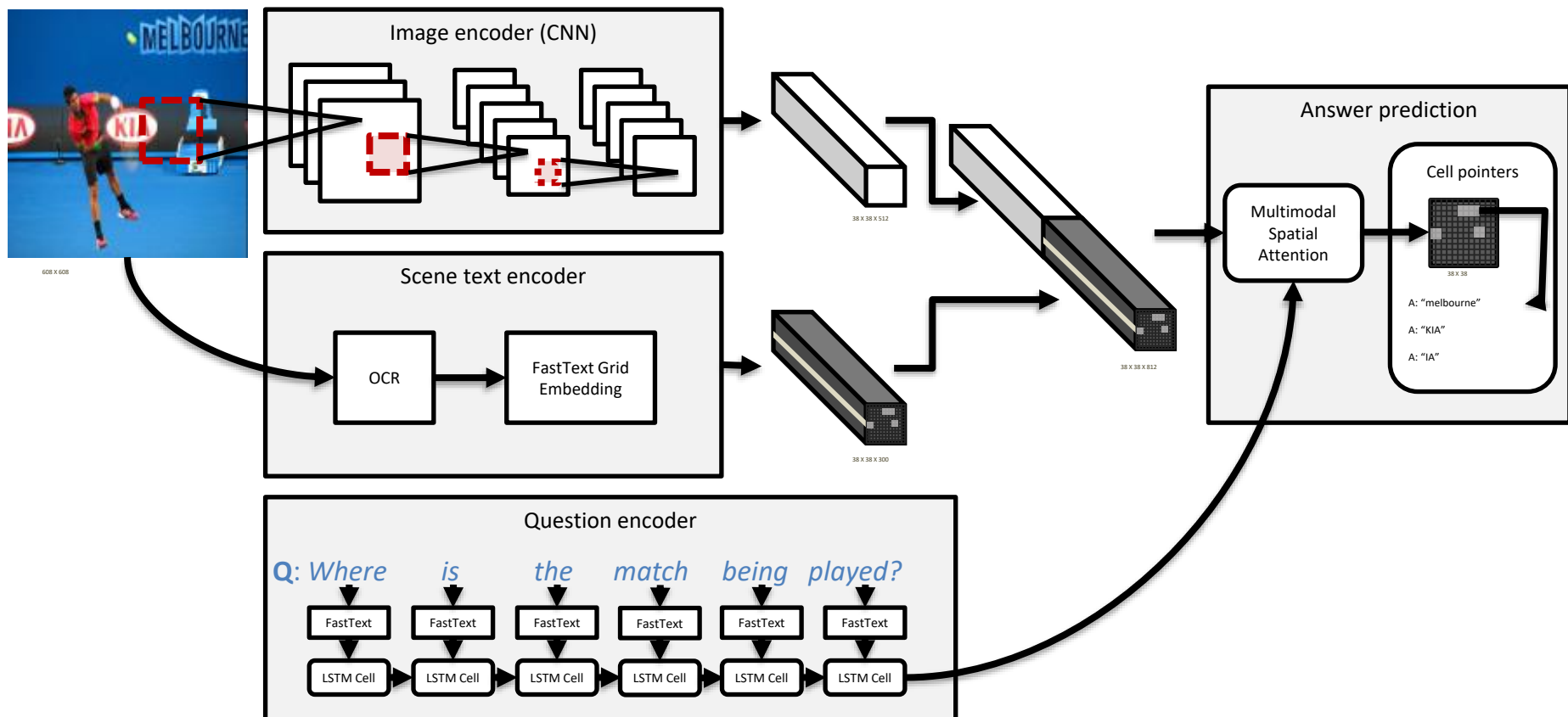


What temperature does the air conditioner show?

24C

Prior world knowledge

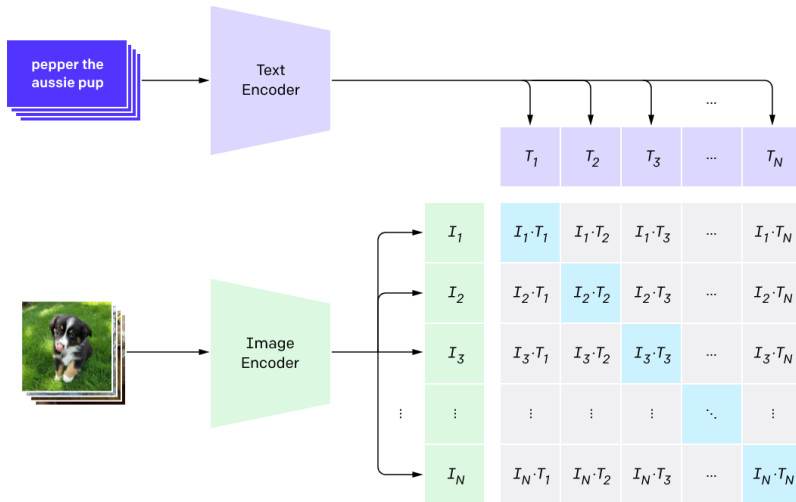
Visual Question Answering



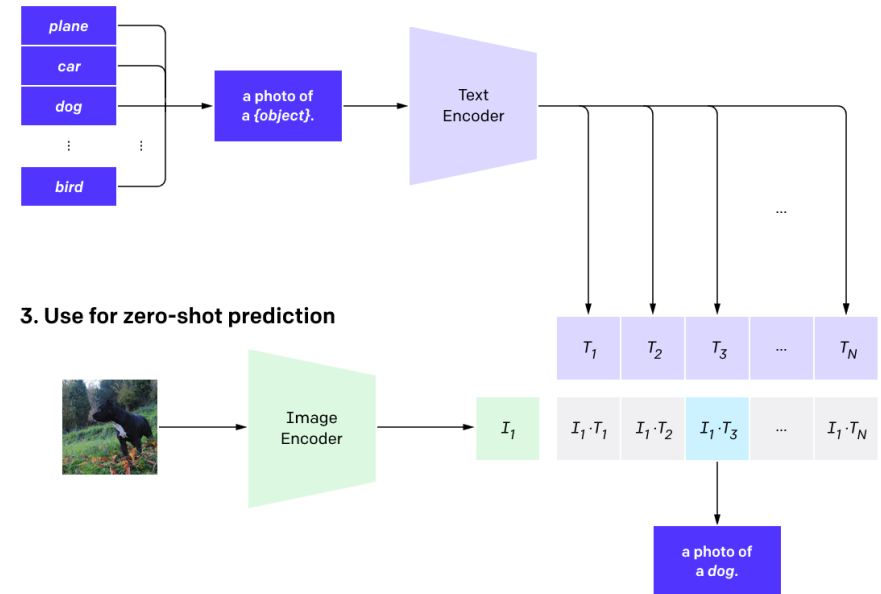
CLIP: Contrastive Language-Image Pre-training

“Our method uses an abundantly available source of supervision: the text paired with images found across the internet”

1. Contrastive pre-training



2. Create dataset classifier from label text



“our best performing CLIP model trains on 256 GPUs for 2 weeks”