# Statistical Methods in AI (CSE/ECE 471)

Representation Learning (Siamese Network, Autoencoders)

Ravi Kiran (ravi.kiran@iiit.ac.in)

Vineet Gandhi (v.gandhi@iiit.ac.in)
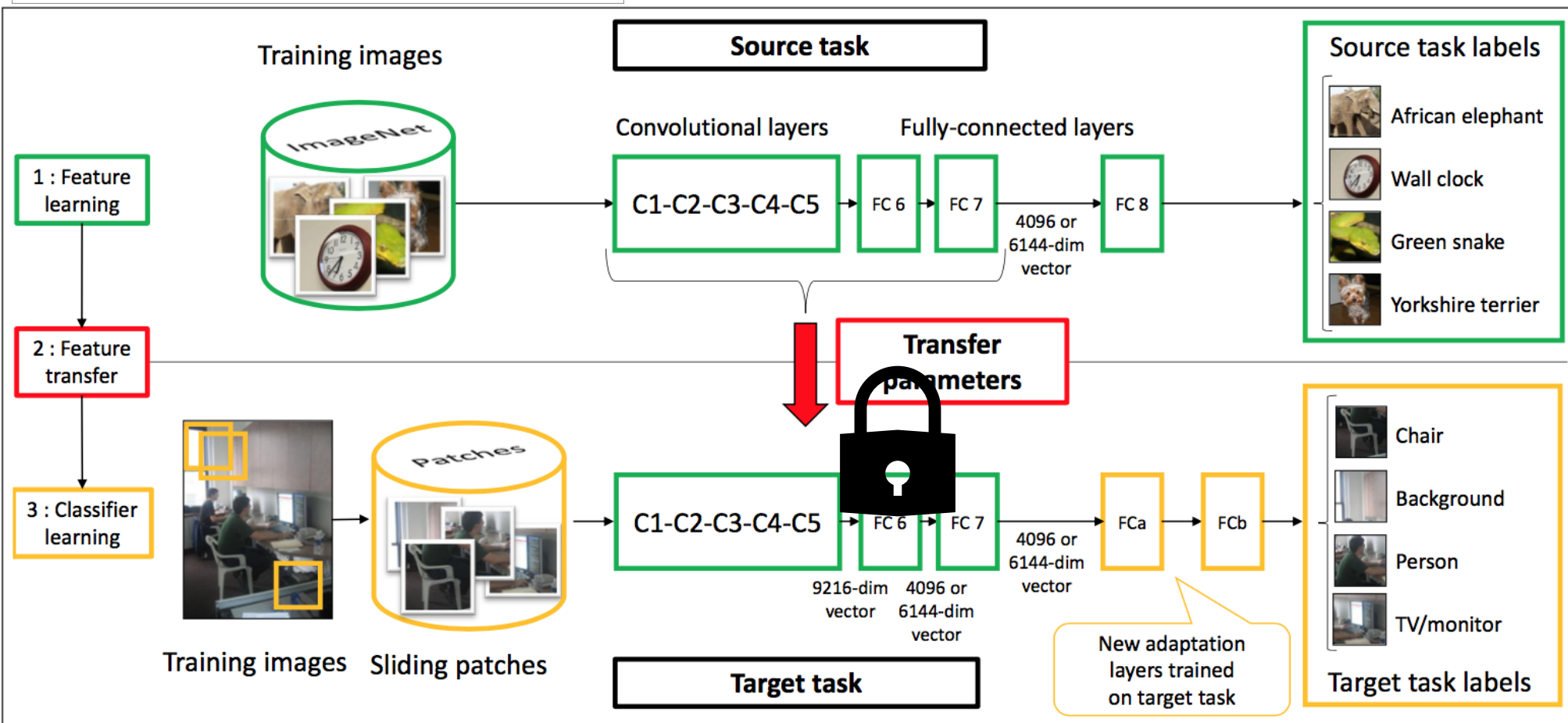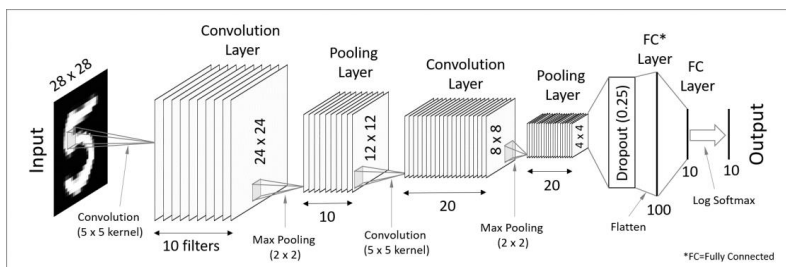
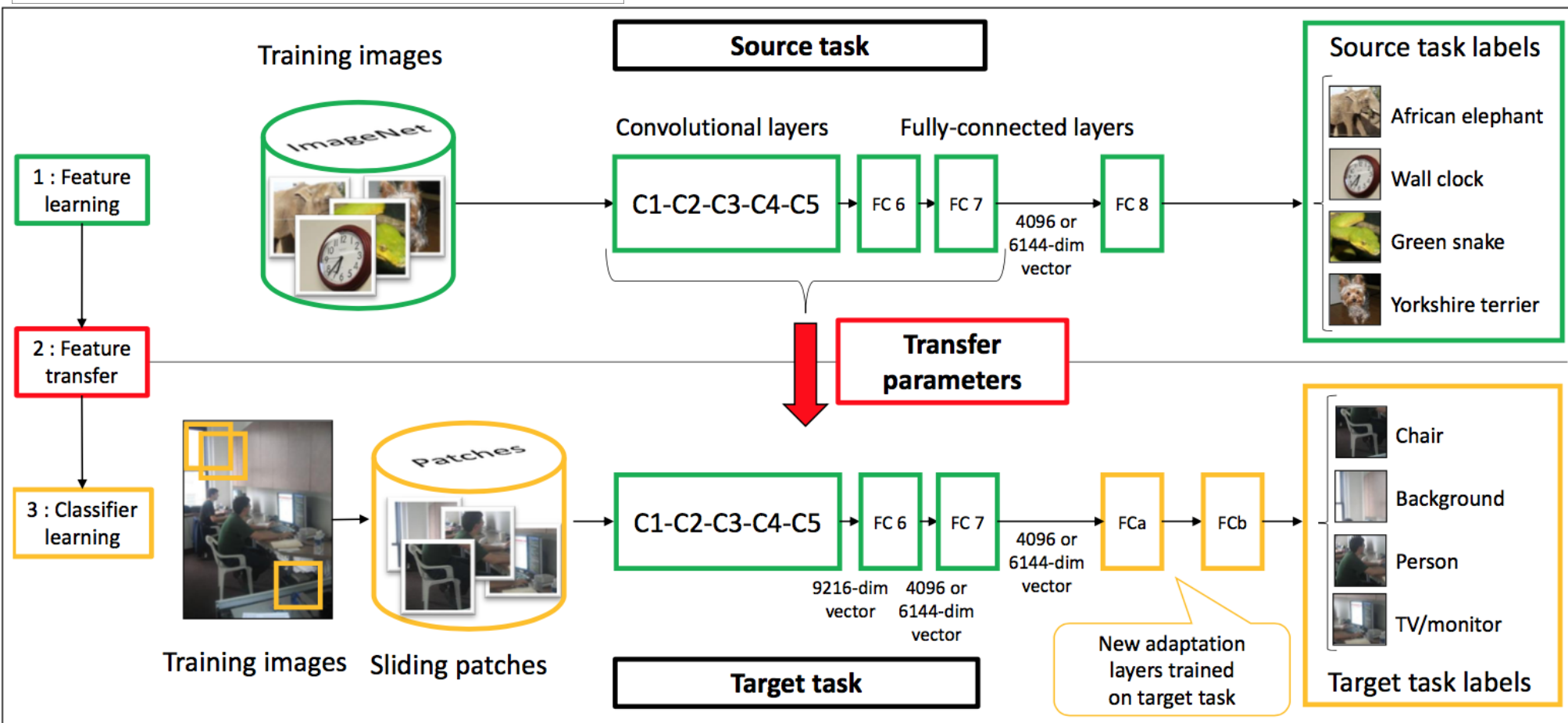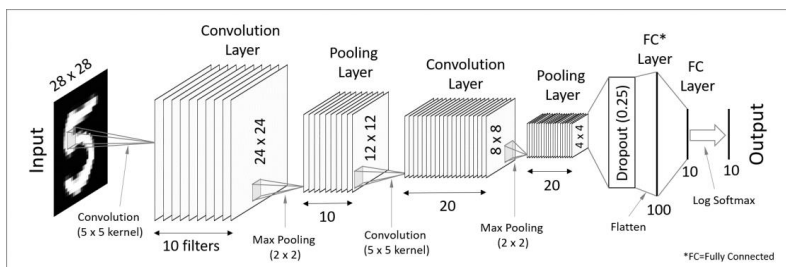Center for Visual Information Technology (CVIT)
IIIT Hyderabad

# Transfer Learning : Approach-1



- Learn only weights for newly added layers.
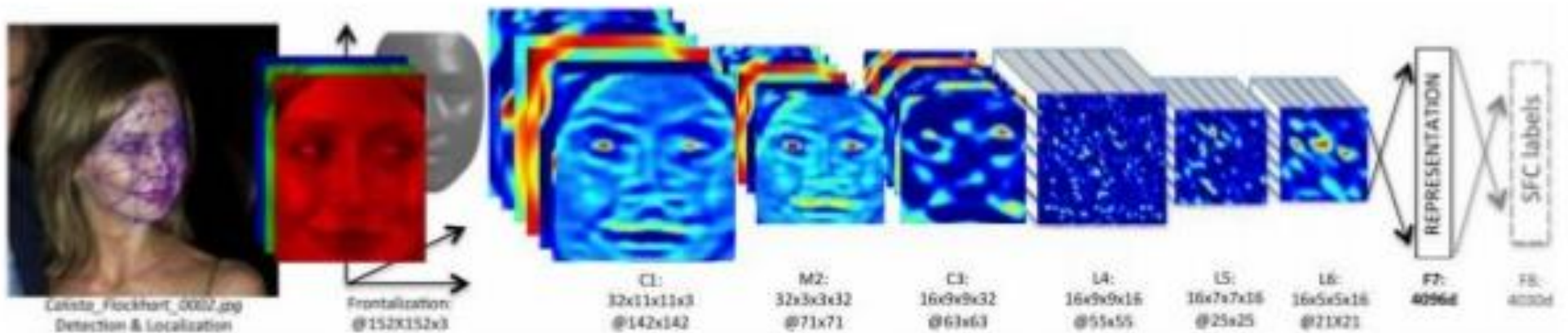- Ideal when 'new domain' data is small in quantity

# Transfer Learning : Approach-2



- LR for new layer weights = 10 * source_lr (for bias, 20 * source_lr)
- Ideal when 'new domain' data is reasonably large or domain shift is significant

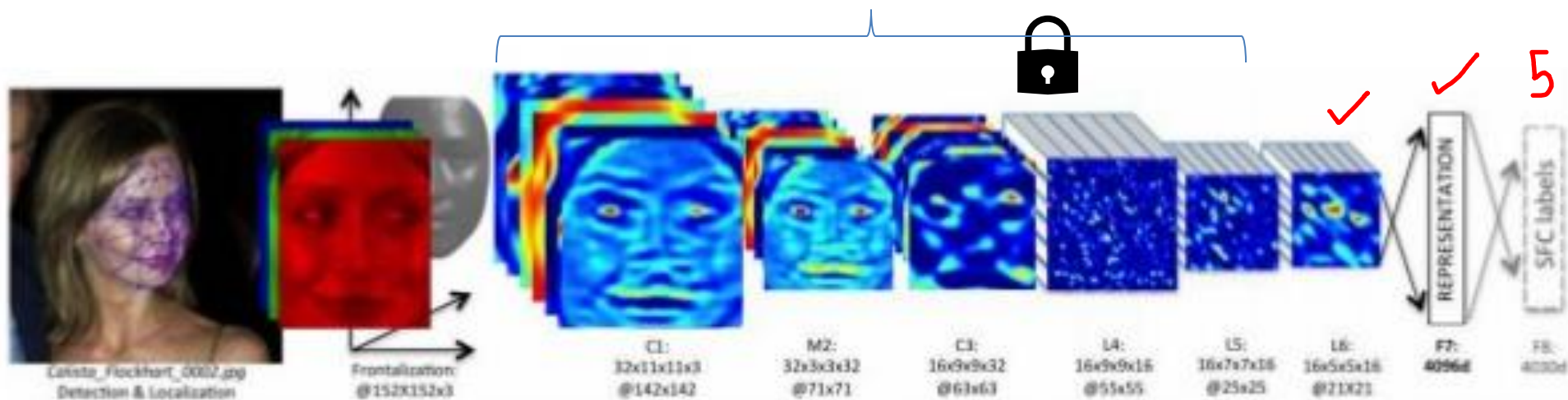# Classification



Face Identification/Recognition (1:N matching)

How to reuse DeepFace (trained on celebrities) for another face dataset ?

# Classification

Classification

Face Identification/Recognition (1:N matching)



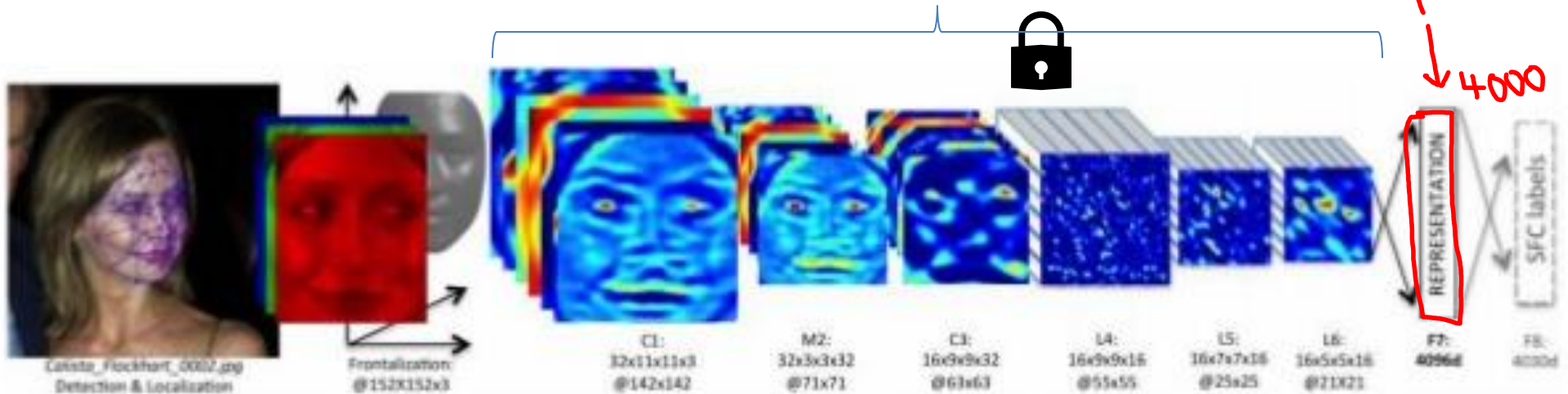How to reuse DeepFace (trained on celebrities) for another face dataset ?
Ans: Fine-tuning

# No-finetuning Classification

Classification



Face Identification/Recognition (1:N matching)

How to reuse DeepFace (trained on celebrities) for another face dataset (without any training) ?
Ans: Use CNN as feature extractor. k-NN on feature representations

# Verification

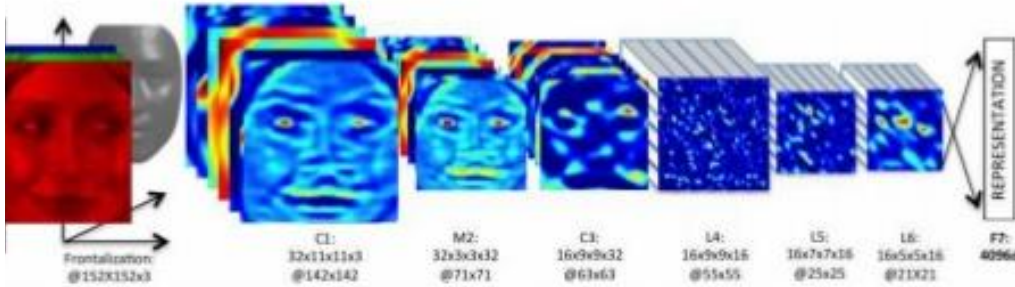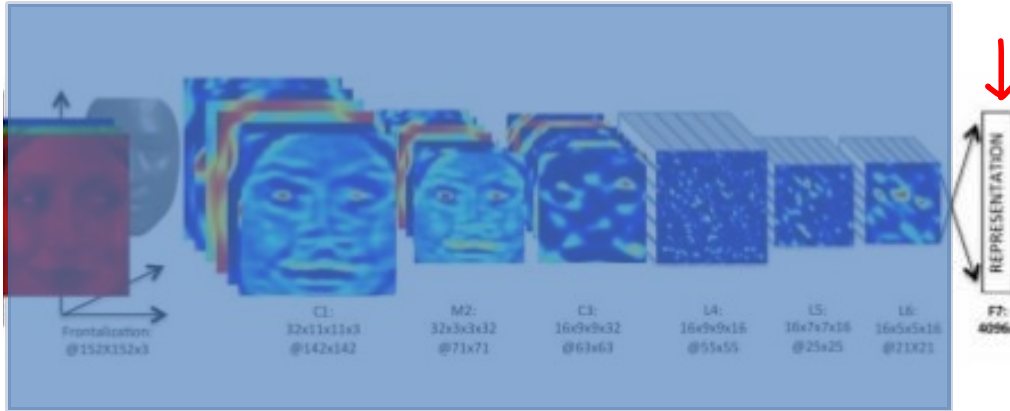Face Authentication/Verification (1:1 matching)
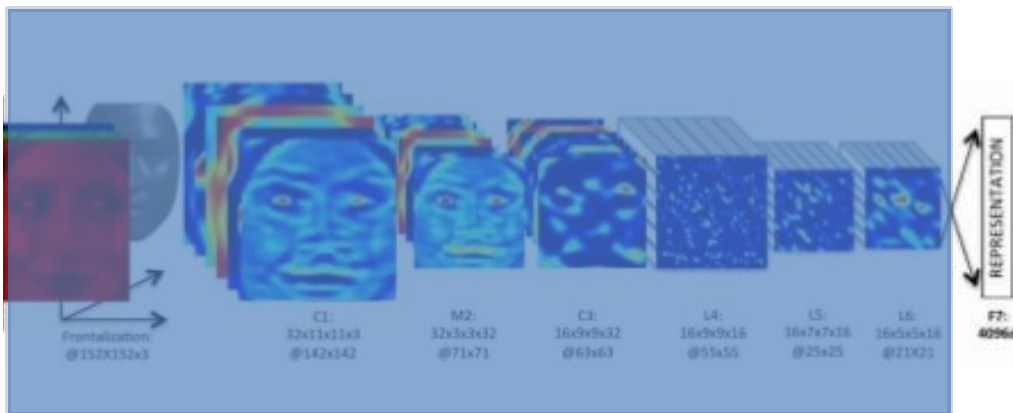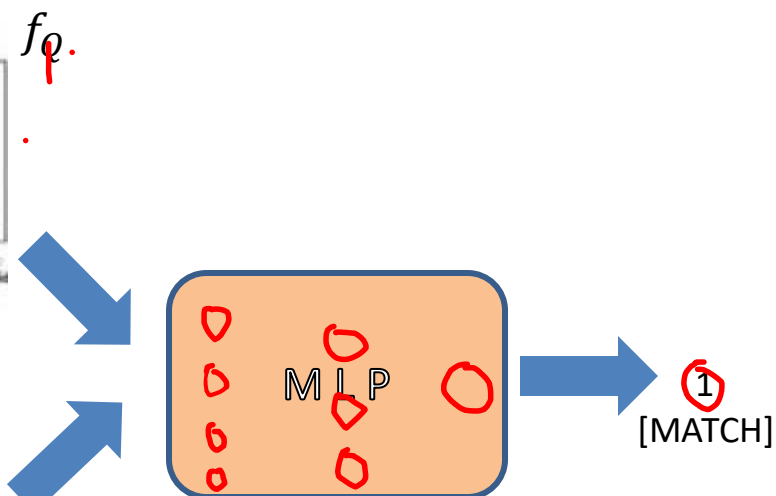
q ✓

d ✓

# Feature Extraction

# Feature Extraction

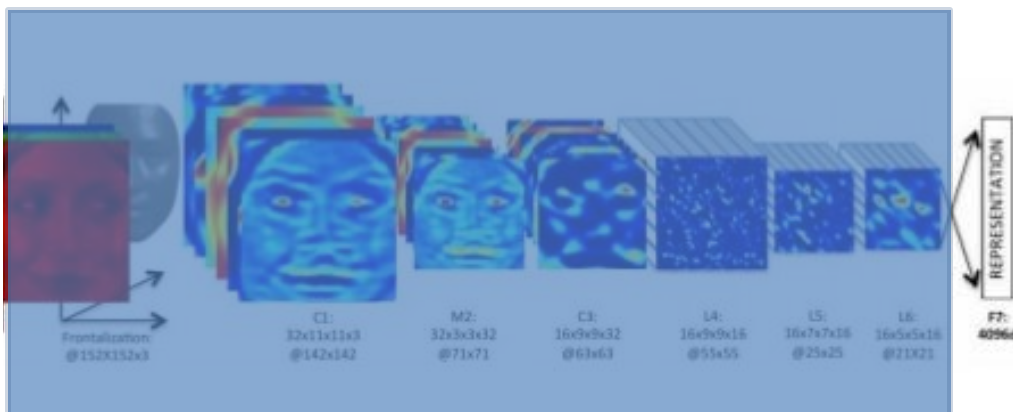

Frontalization
@152x152x3

C1:
32x11x11x3
@142x142

M2:
32x3x3x32
@71x71

C3:
16x9x9x32
@63x63

L4:
16x9x9x16
@55x55

L5:
16x7x7x16
@25x25

L6:
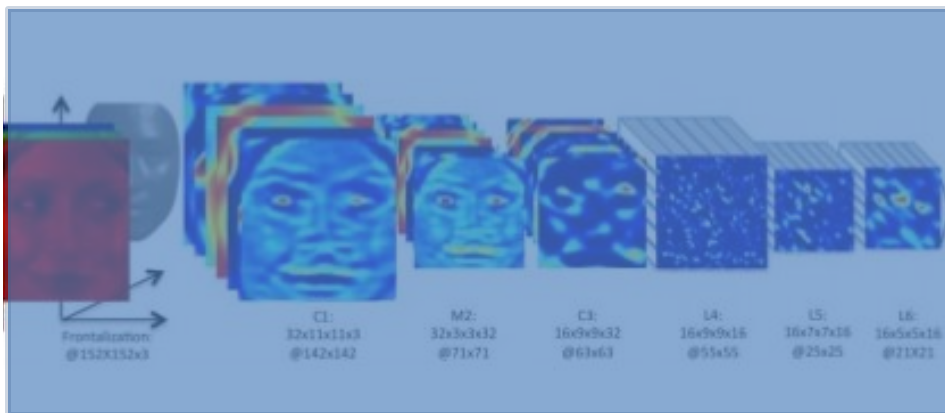16x5x5x16
@21x21

F7:
4096d

REPRESENTATION

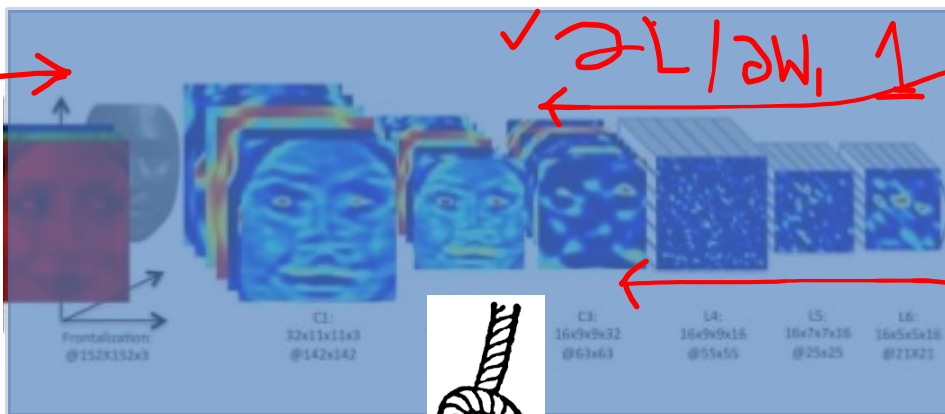# Verification: Approach - 1

# Verification: Approach – 1B

# Verification: Approach – 1C



DB-image

Tied weights

DB image

$f_Q$

5000

$f_B$

400

Contrastive Loss

CNN

Vector Space

CNN

Vector Space

Contrastive Loss:

Learn $f_Q$, $f$ such that:
- dist($f_Q$, $f$ ) is large when ids mismatch
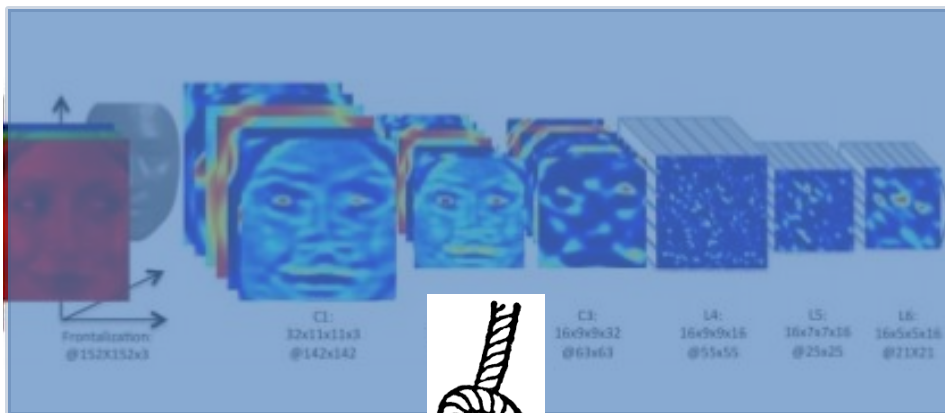- dist($f_Q$, $f$ ) is small when ids match

# Verification: Approach – 1C



Query

DB image

$f_Q$

$f_B$

Tied weights

Contrastive Loss

$f_q$   $f_B$

Contrastive Loss: $yd^2 + (1-y)\max(margin - d, 0)^2$

$>0$

$d < m$

Learn $f_Q$, $f$   such that:
- $d = $dist$(f_Q, f$  $)$ is large when ids mismatch (y=0)
- $d = $dist$(f_Q, f$  $)$ is small when ids match (y=1)

CNN

$d^2$

Vector Space

CNN

$m$

Vector Space

# Verification: Approach – 1C



**Query**

**DB image**

$f_Q$

$f_B$

Tied weights

Learning a similarity function

Contrastive Loss

Contrastive Loss: $yd^2 + (1-y)\max(margin - d, 0)^2$

Learn $f_Q$, $f$  such that:
- $d = \text{dist}(f_Q, f\ )$ is large when ids mismatch (y=0)
- $d = \text{dist}(f_Q, f\ )$ is small when ids match (y=1)
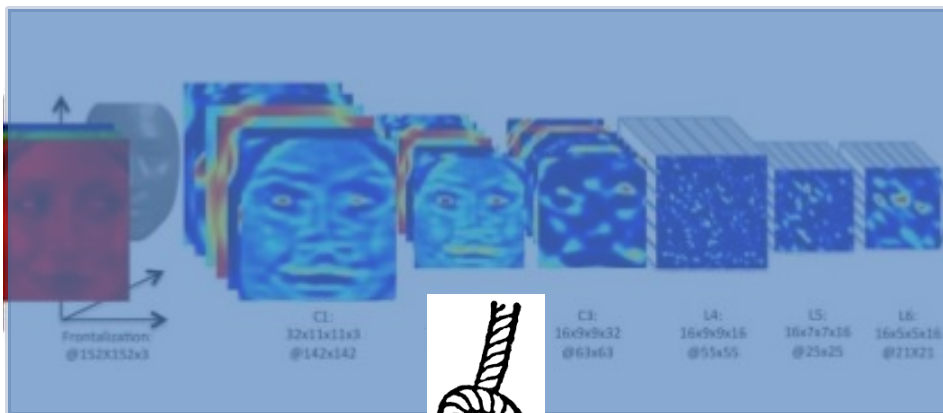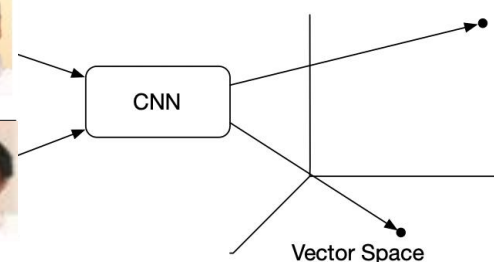
Vector Space

Vector Space
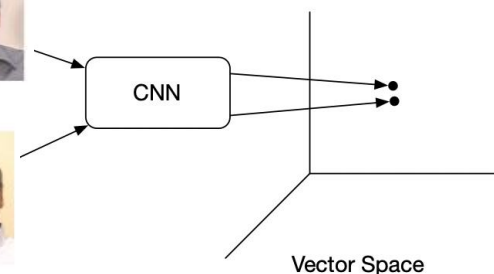
CNN

CNN

# Verification Approach 2

# Popular Architecture Varieties

- No one "architecture" fits all!
- Design largely governed by what performs well empirically on the task at hand.



Inputs are merged right at the onset

Inputs are first embedded independently, then merged.

Zagoruyko, S. and Komodakis, N., 2015. Learning to compare image patches via convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 4353-4361).

# Siamese Network

## Application in Signature Verification

- The input is 8(feature) x 200(time) units.

- The cosine distance was used, (1 for genuine pairs, -1 for forgery pairs )
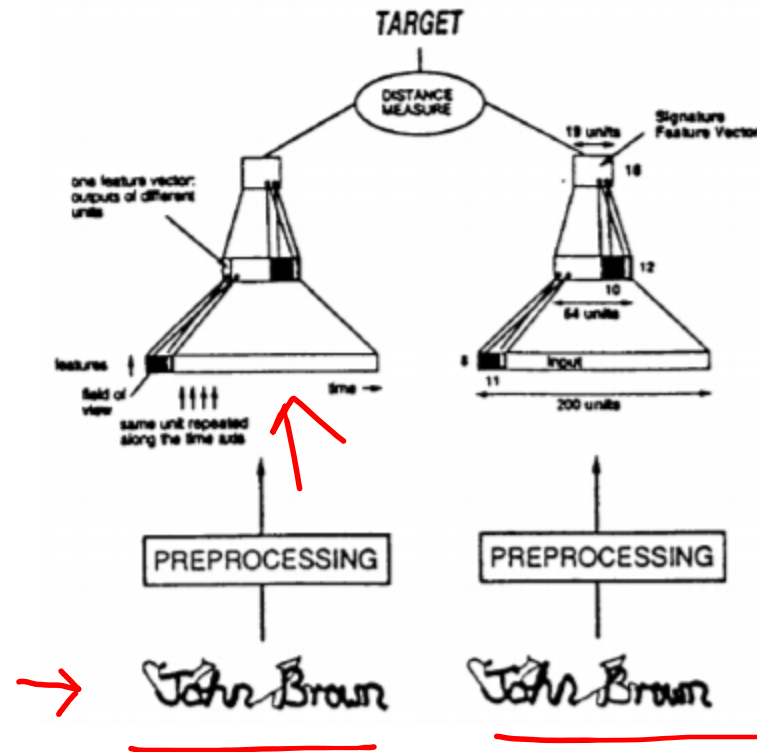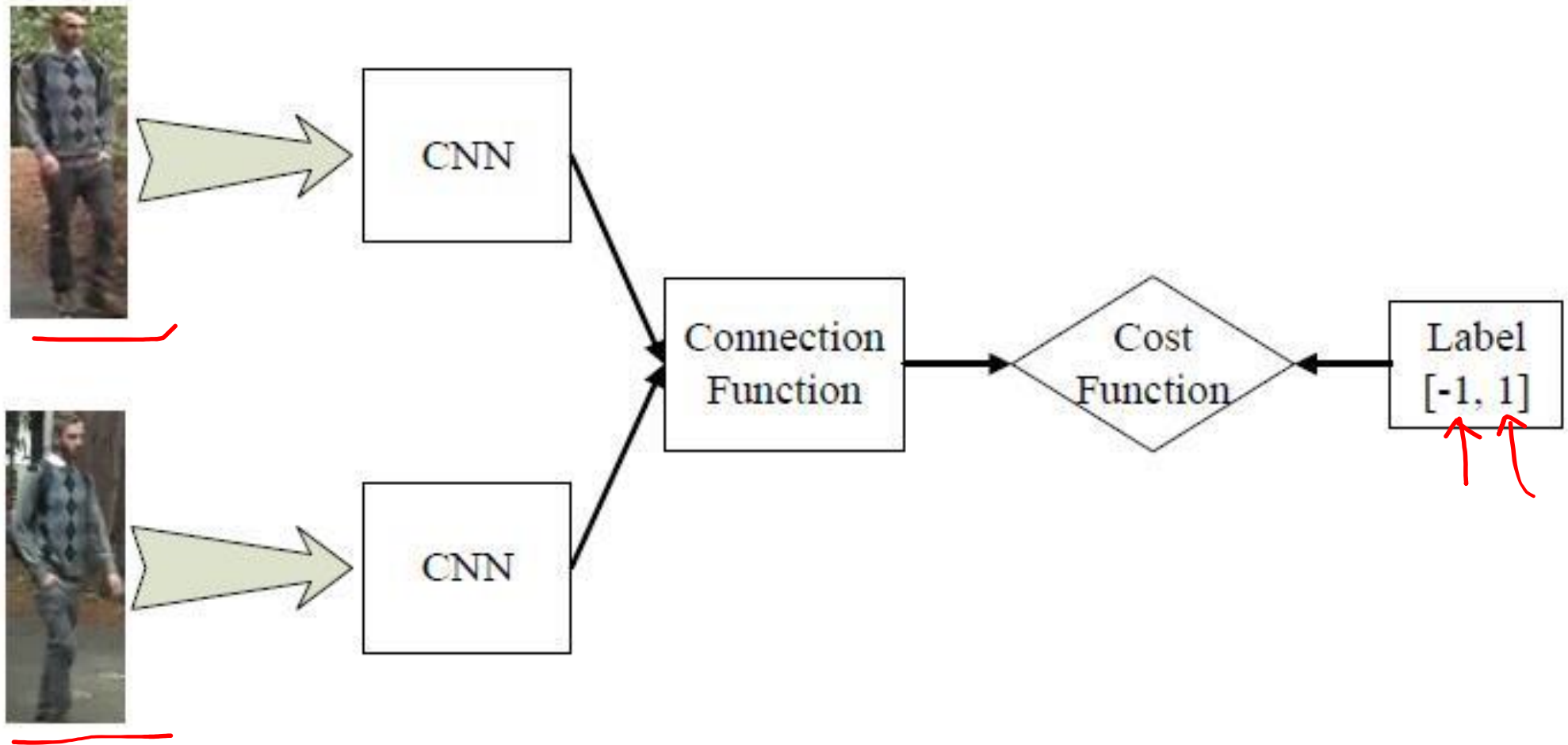


*Bromley J, Guyon I, Lecun Y, et al. Signature Verification using a" Siamese" Time Delay Neural Network, NIPS Proc. 1994*
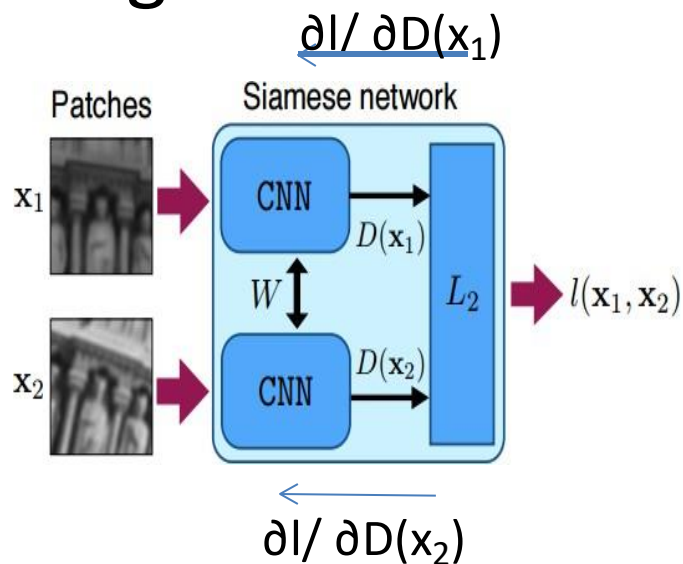
# Siamese Network (Person re-id)

# Siamese CNN – Training

- Update each of the two streams independently and then average the weights.

$\partial l / \partial D(x_1)$



Patches    Siamese network

$x_1$   CNN   $D(\mathbf{x}_1)$

$W$   $L_2$   $l(\mathbf{x}_1, \mathbf{x}_2)$
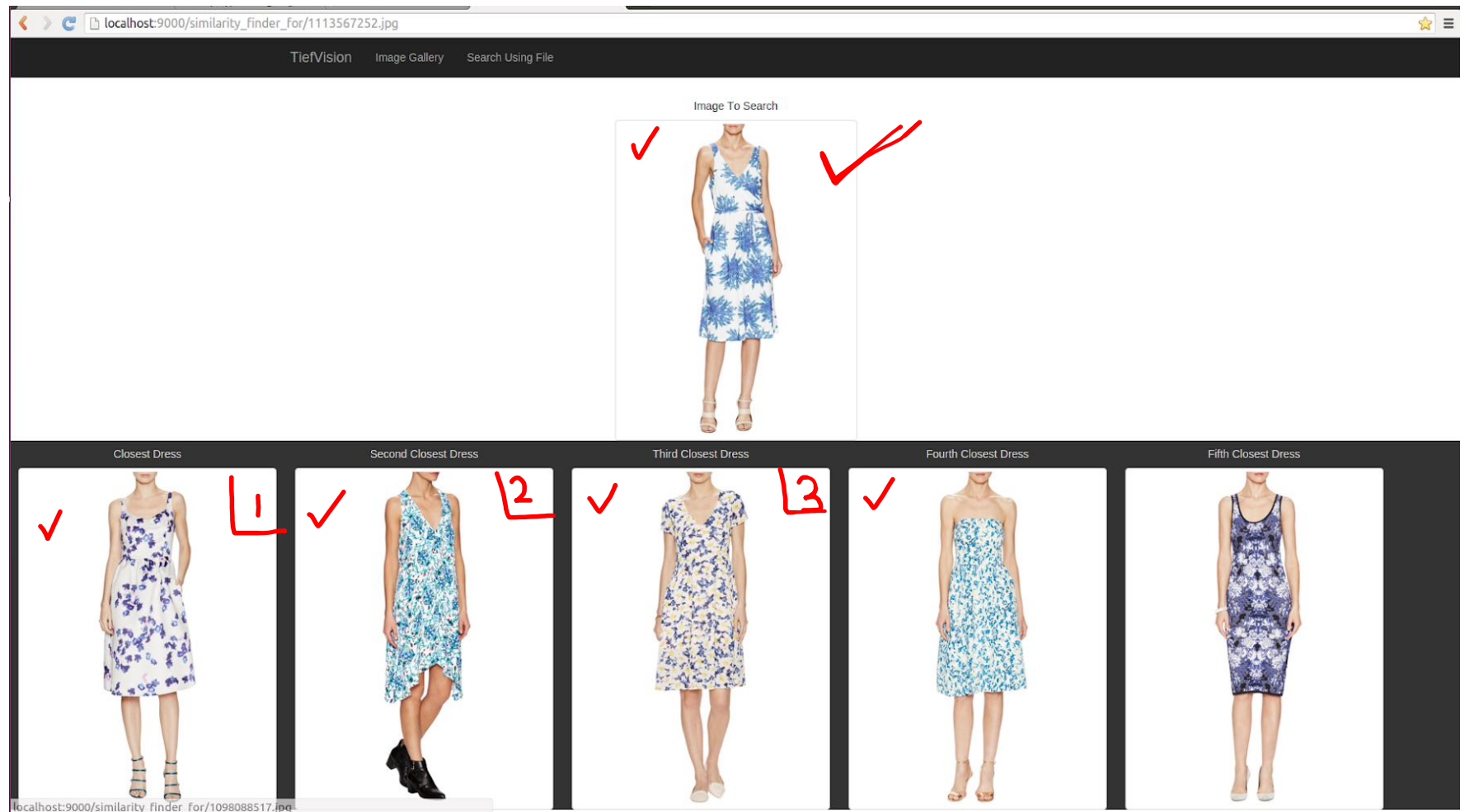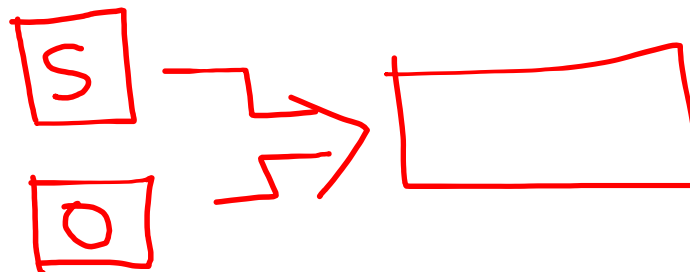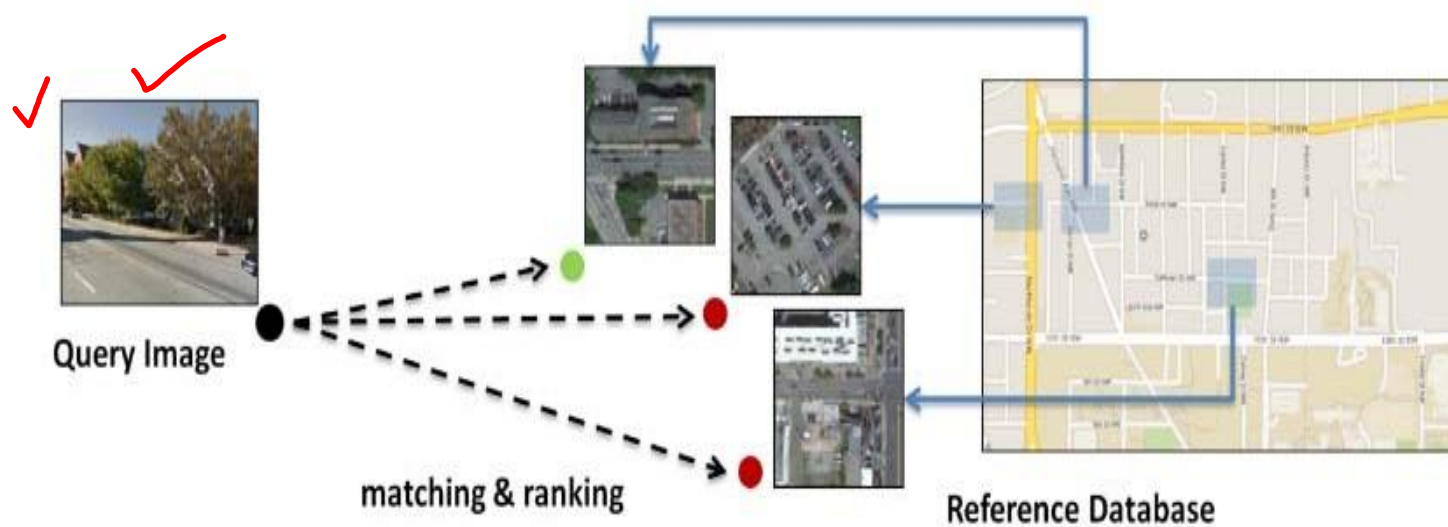
$x_2$   CNN   $D(\mathbf{x}_2)$

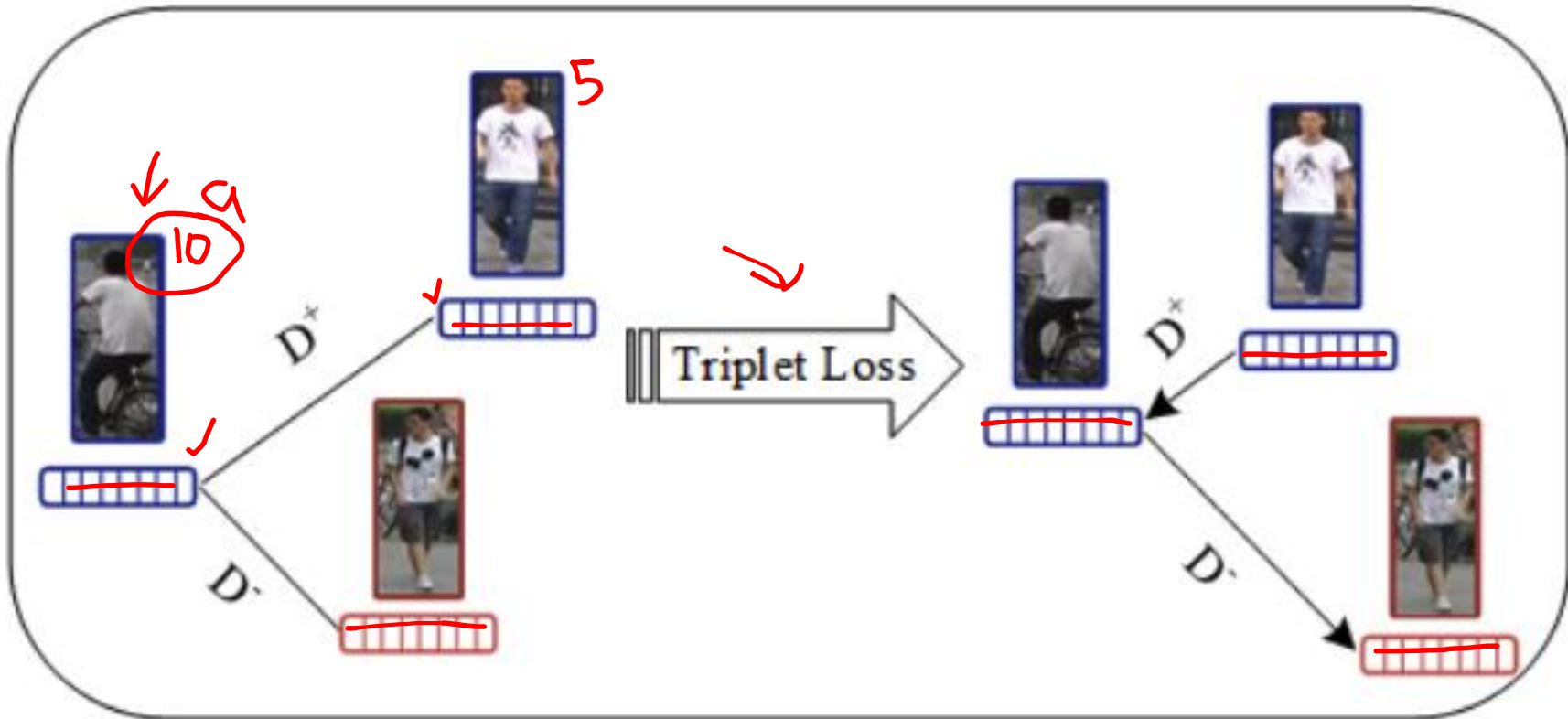$\partial l / \partial D(x_2)$

# Applications

Retrieval

Ranking



https://github.com/paucarre/tiefvision

# Street-View to Overhead-View Image Matching



Query Image

matching & ranking

Reference Database

Vo, N.N. and Hays, J., 2016, October. Localizing and orienting street views using overhead imagery. In European Conference on Computer Vision (pp. 494-509).

# Many variants exist

- Popular Loss Function – Triplet Loss

# Code Reference

https://medium.com/@prabhnoor0212/siamese-network-keras-31a3a8f37d04
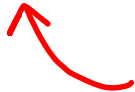
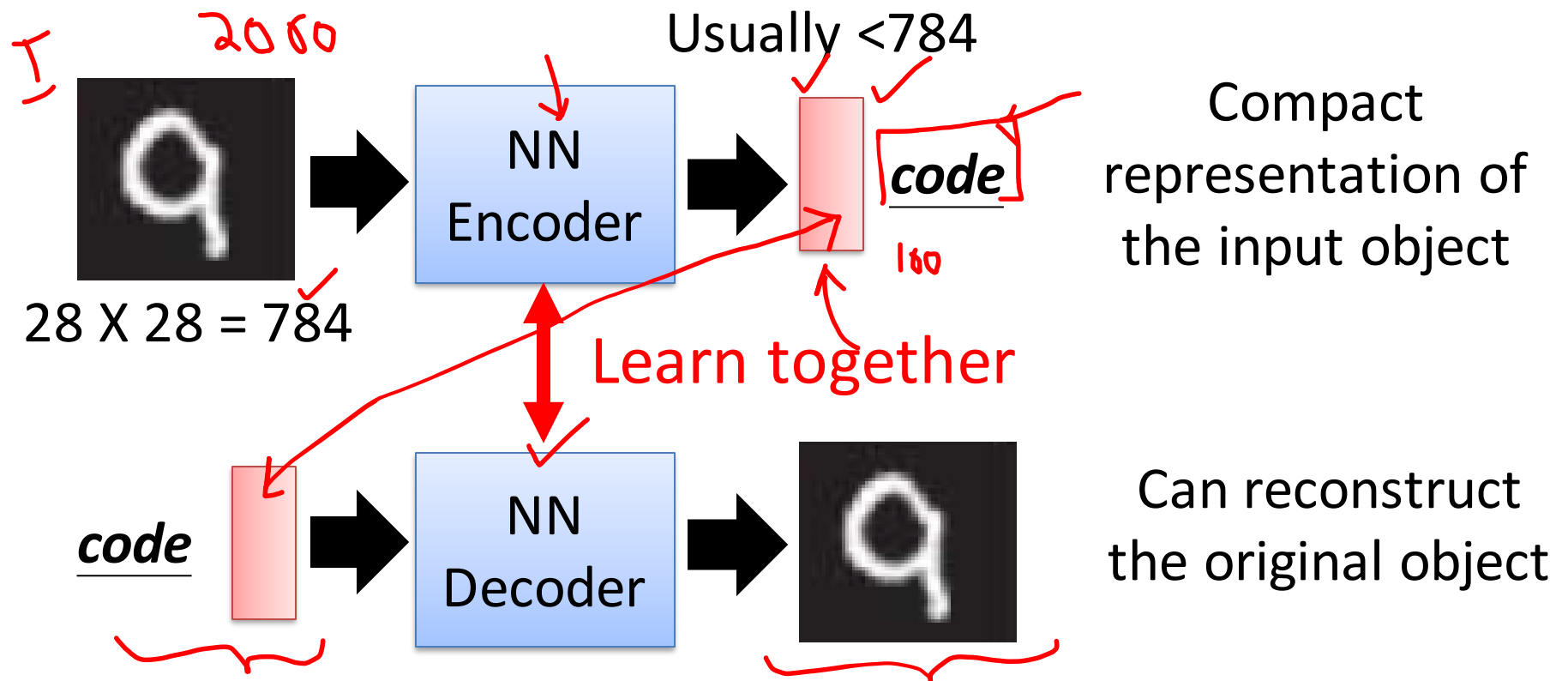# Unsupervised Learning:
## Deep Auto-encoder

# Unsupervised Learning

"We expect unsupervised learning to become far more important in the longer term. Human and animal learning is largely unsupervised: we discover the structure of the world by observing it, not by being told the name of every object."
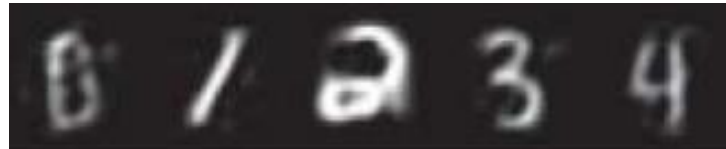– LeCun, Bengio, Hinton, Nature 2015

# Auto-encoder



**2060**

28 X 28 = 784

NN Encoder

Usually <784

*code*

**Ibo**

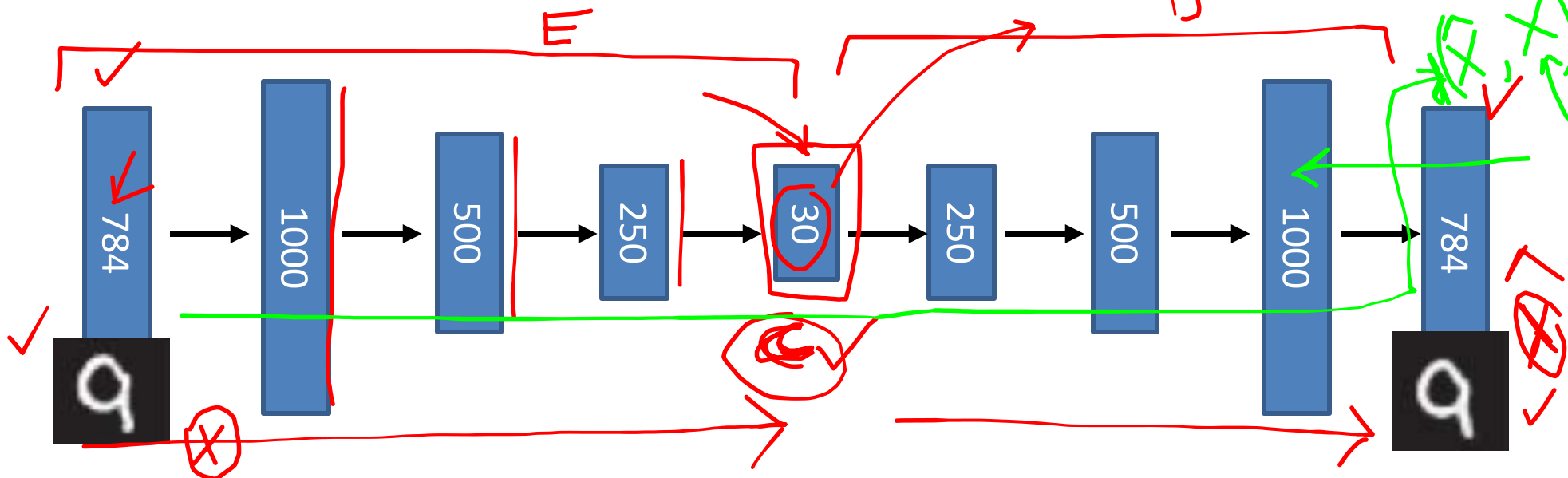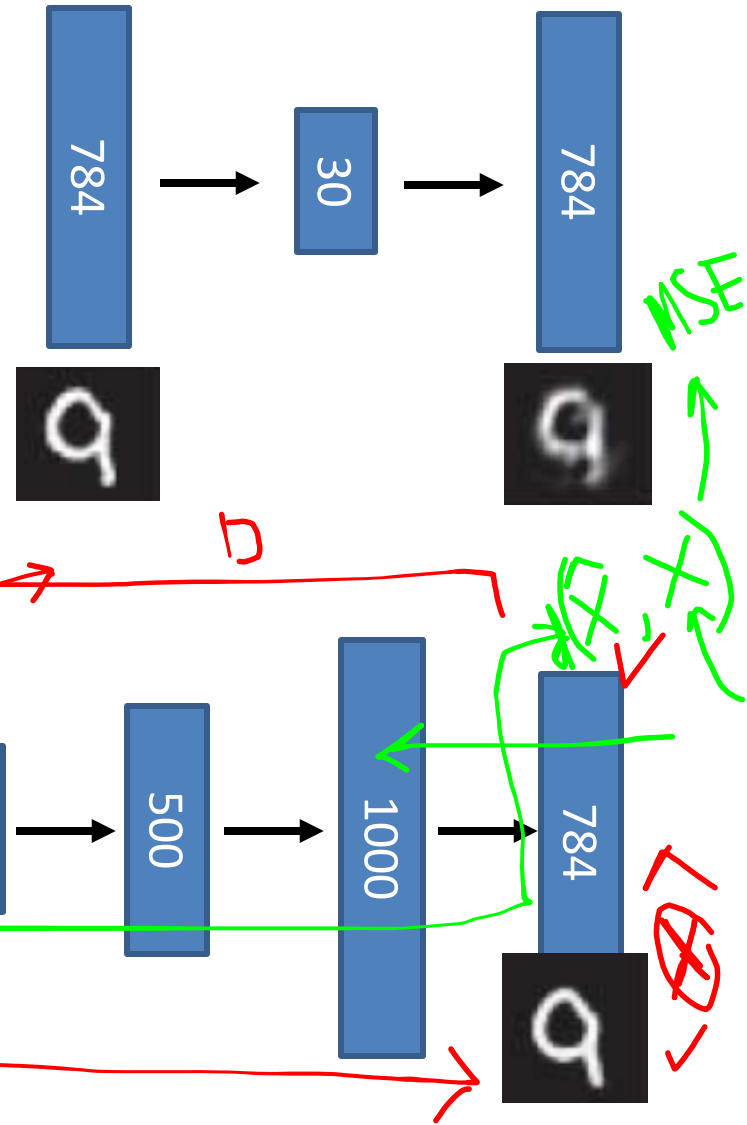Compact representation of the input object
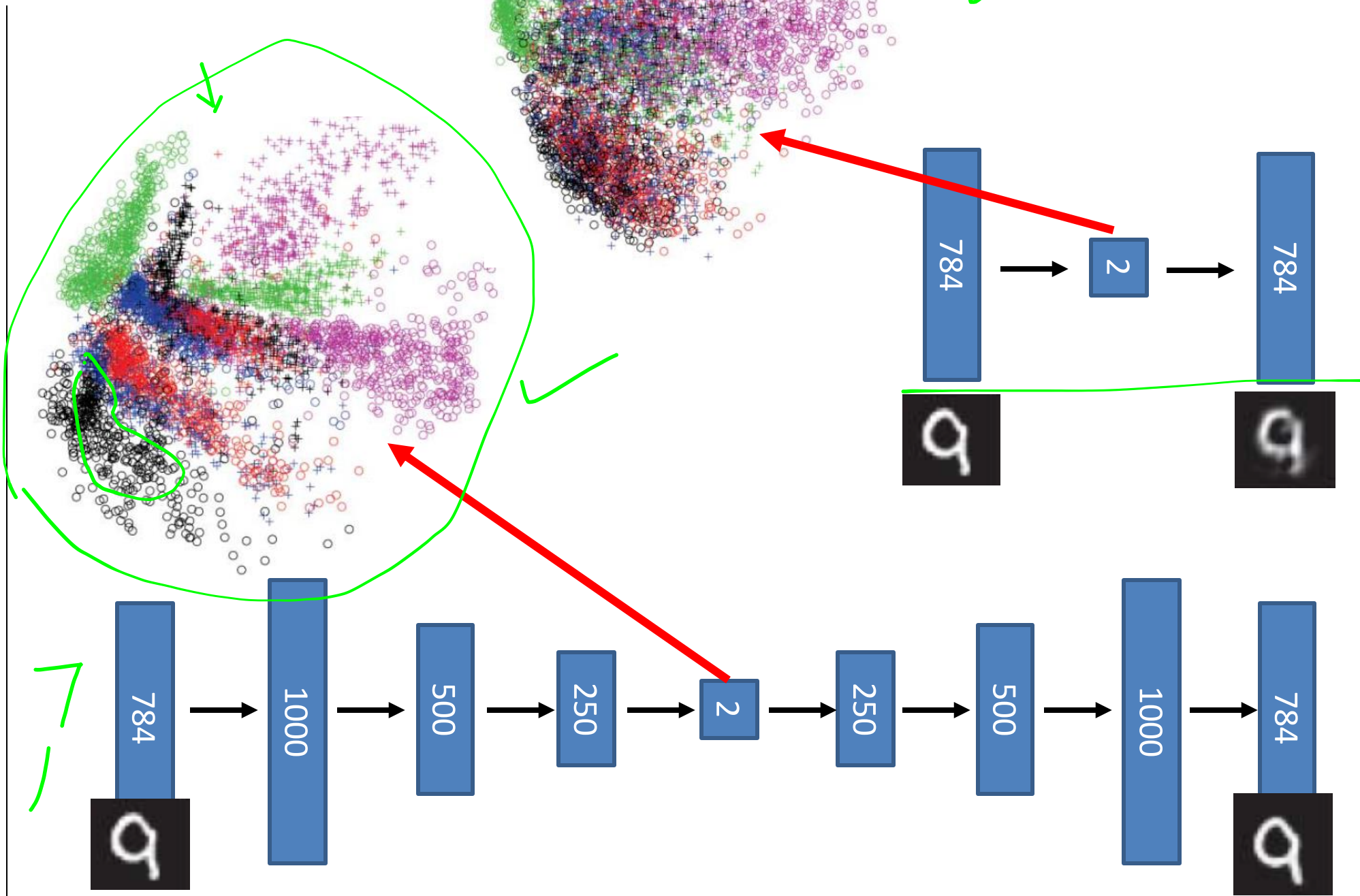
*code*

NN Decoder

Learn together

Can reconstruct the original object

# Deep Auto-encoder

# Auto-encoder

De-noising auto-encoder

As close as possible

encode                    decode

$x$      Add
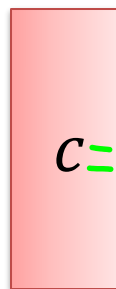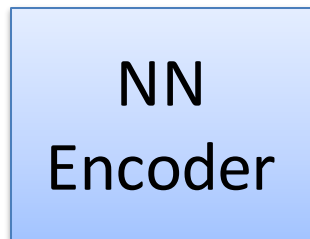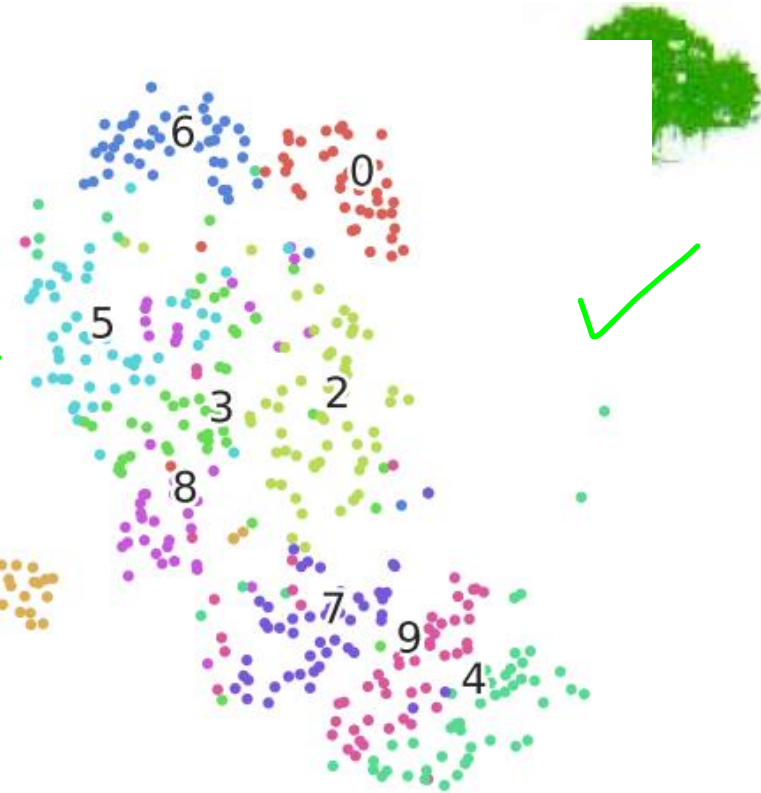         noise            $x'$                    $c$                    $\hat{x}$

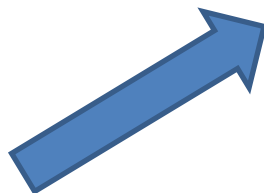Vincent, Pascal, et al. "Extracting and composing robust features
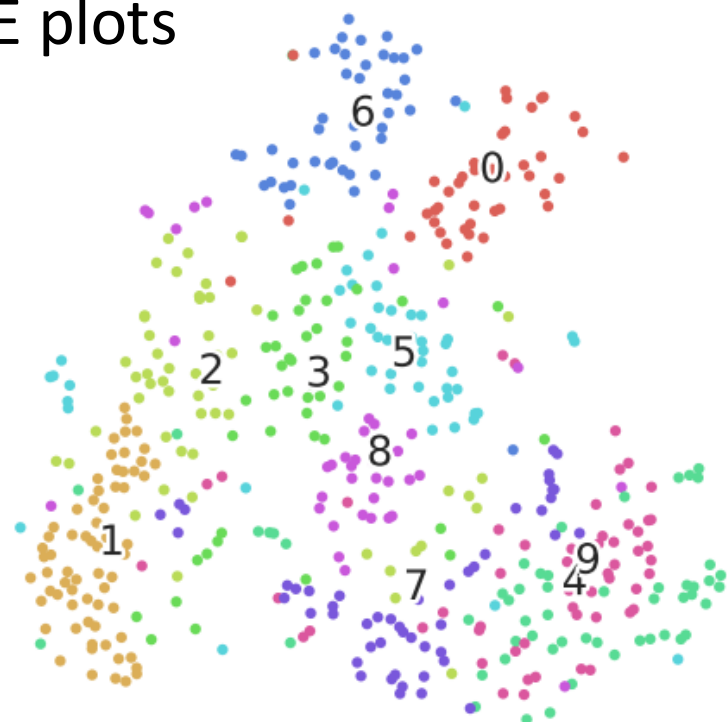with denoising autoencoders." *ICML,* 2008.

# Deep Auto-encoder - Example



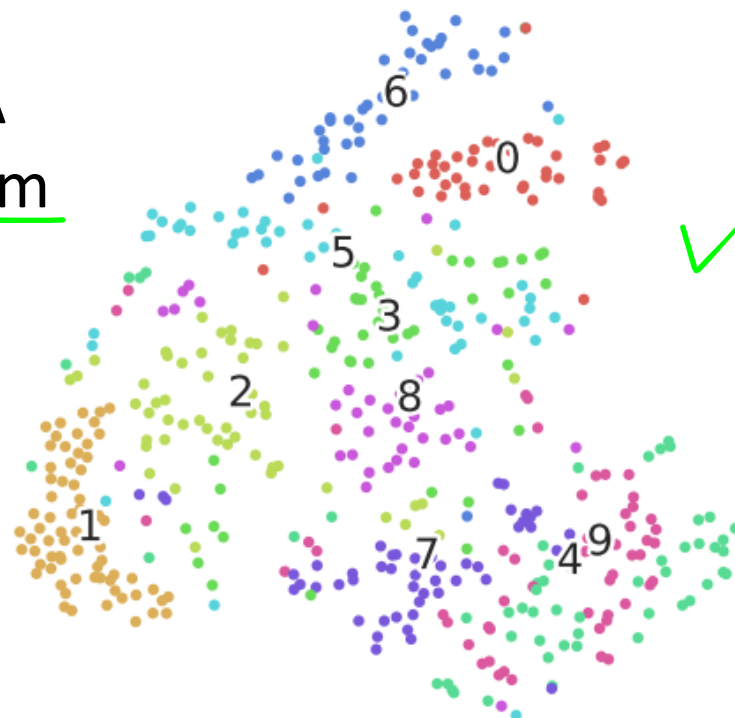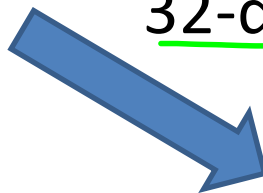tSNE plots

NN Encoder

$c = 32$

PCA 32-dim

# Auto-encoder – Text Retrieval

## *Vector Space Model*

query

document

## *Bag-of-word*

word string:
"This is an apple"

| | |
|---|---|
| this | 1 |
| is | 1 |
| a | 0 |
| an | 1 |
| apple | 1 |
| pen | 0 |
| ⋮ | |

Semantics are not considered.

# Auto-encoder – Text Retrieval

The documents talking about the same thing will have close code.

2 → 125 → 250 → 500 → 2000

Bag-of-word
(document or query)



Interbank markets

European Community monetary/economic

Energy markets

query

Disasters and accidents

Leading economic indicators

Legal/judicial

Accounts/ earnings

Government borrowings

LSA: project documents to 2 latent topics

# Auto-encoder – Similar Image Search

Retrieved using Euclidean distance in pixel intensity space



(Images from Hinton's slides on Coursera)

Reference: Krizhevsky, Alex, and Geoffrey E. Hinton. "Using very deep autoencoders for content-based image retrieval." *ESANN*. 2011.

# Auto-encoder – Similar Image Search

32x32

8192

4096

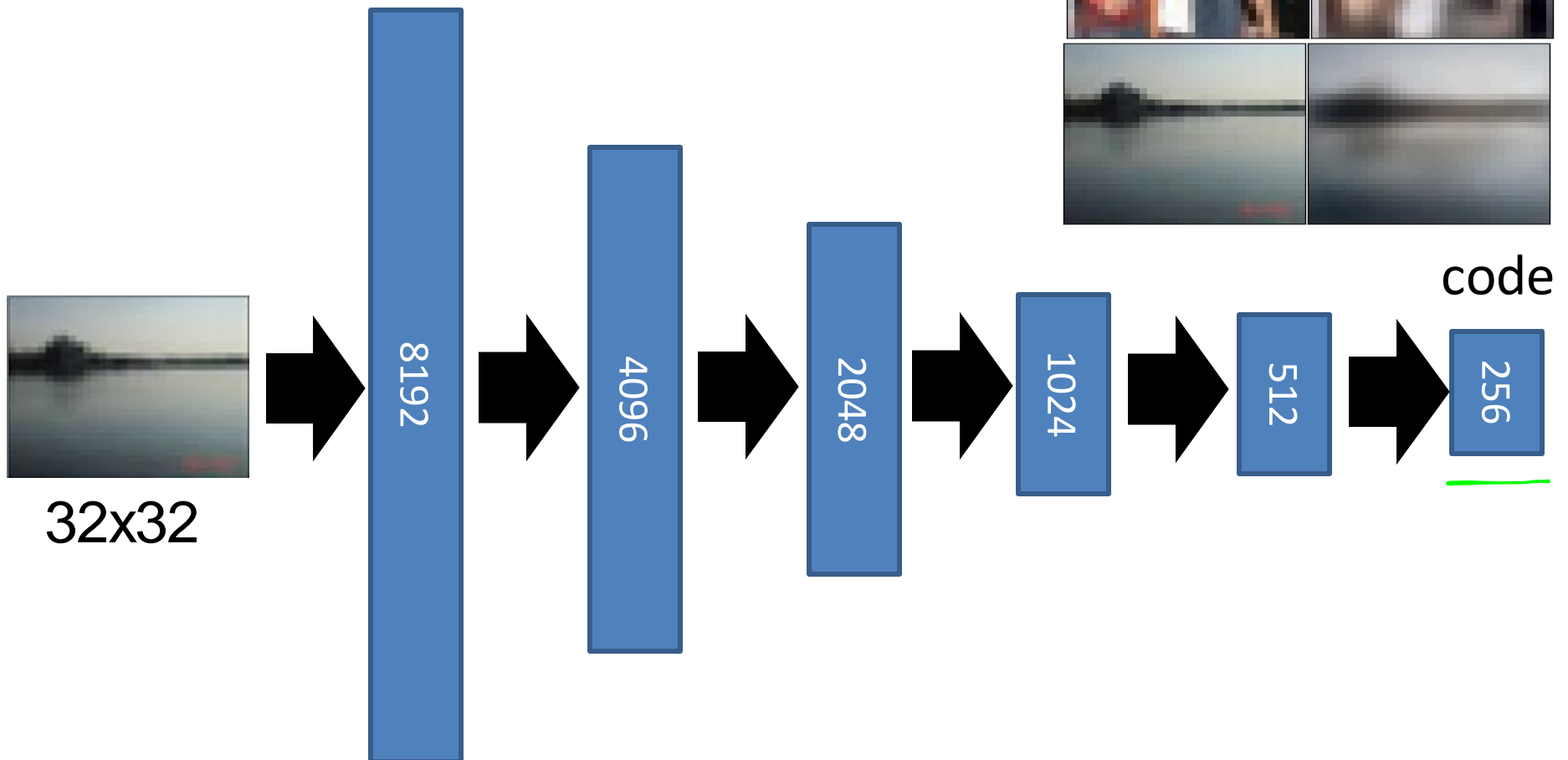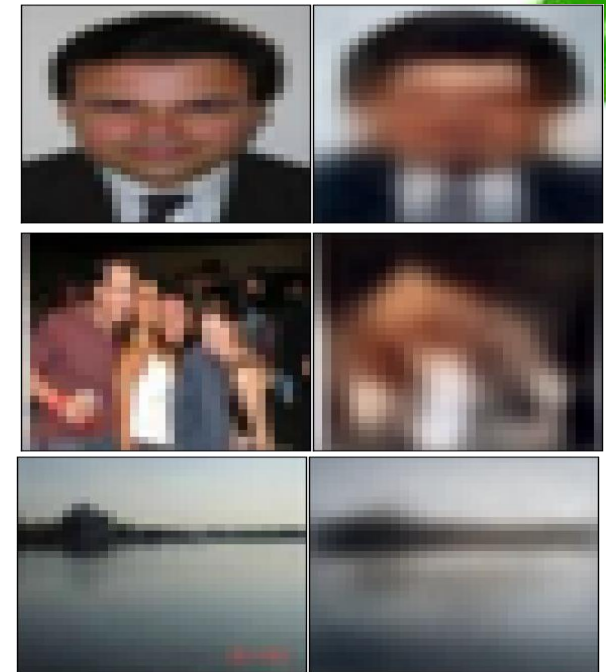2048

1024

512

256

code

(crawl millions of images from the Internet)

# Retrieved using Euclidean distance in pixel intensity space



dist: 0.0 — dist: 3064.2 — dist: 3094.1 — dist: 3132.4

dist: 3139.2 — dist: 3147.0 — dist: 3150.9 — dist: 3154.8

## retrieved using 256 codes



dist: 0

# Auto-encoder for CNN



As close as possible

Deconvolution

Unpooling

Deconvolution

Unpooling

Deconvolution

code

Convolution

Pooling

Convolution

Pooling

# CNN -Unpooling



14 x 14                28 x 28



Max Locations "Switches"

Pooled Maps

Pooling

Rectified Feature Maps

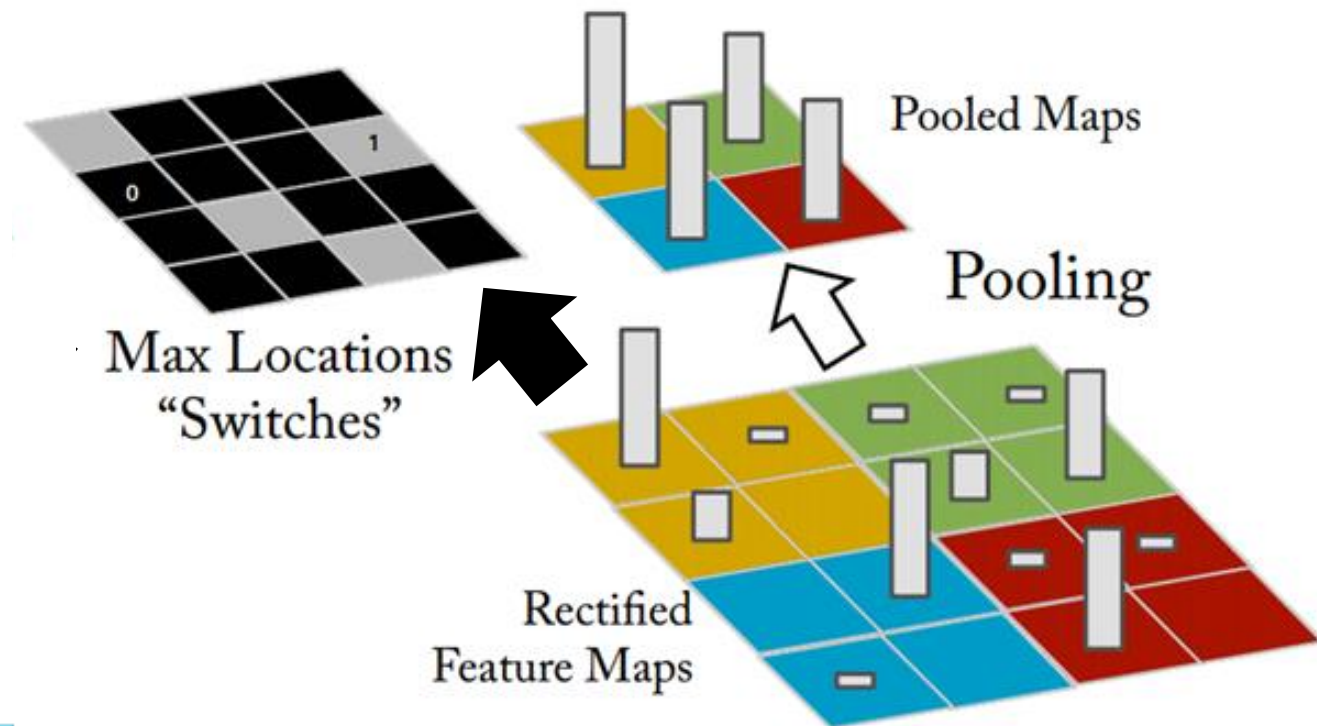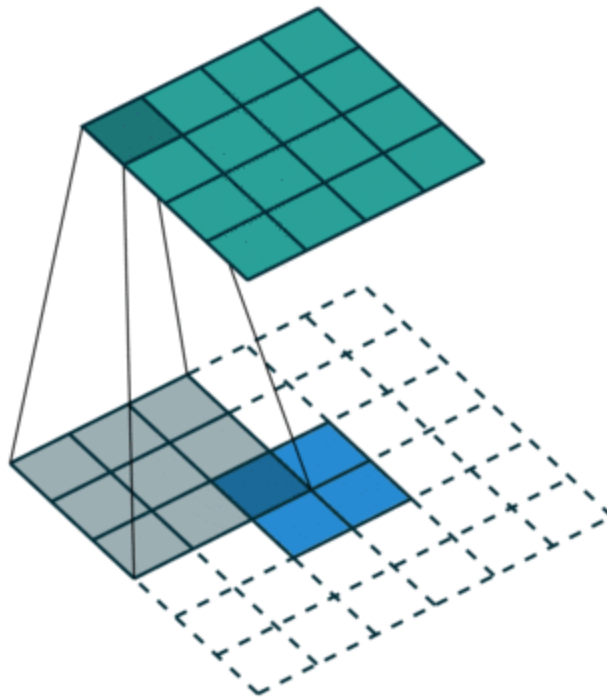Alternative: simply repeat the values

Source of image :
https://leonardoaraujosantos.gitbooks.io/artificial-inteligence/content/image_segmentation.html
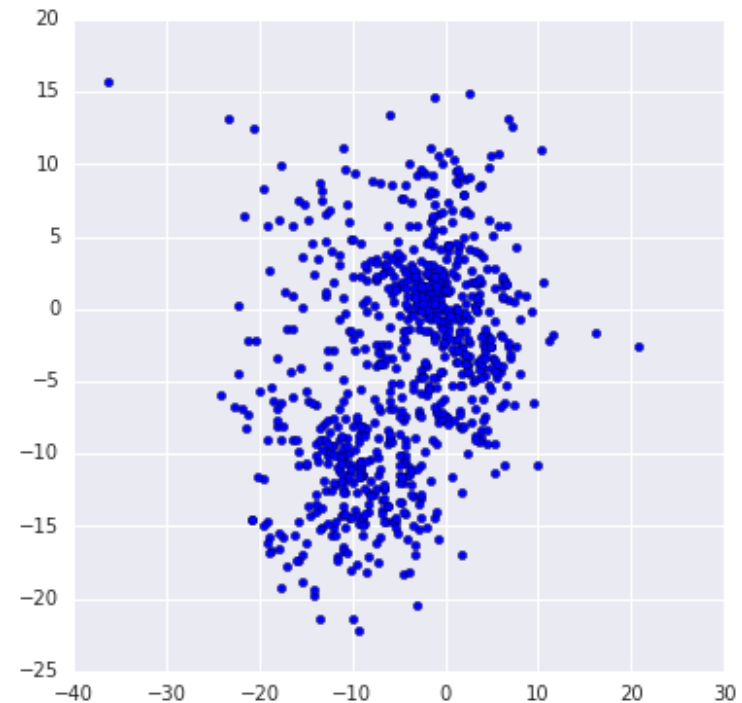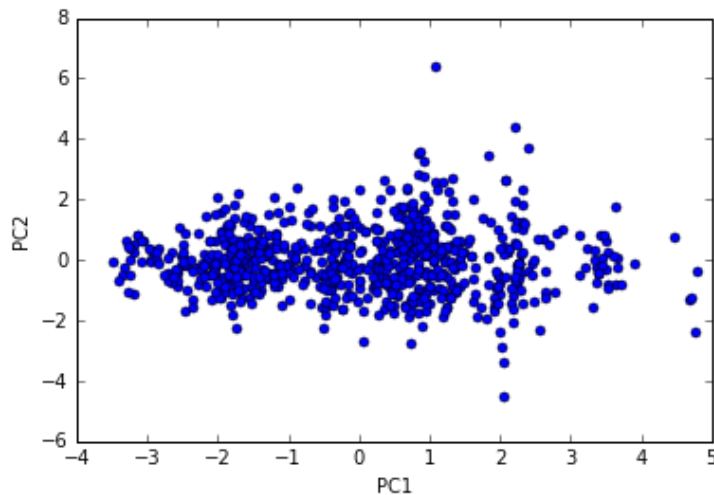
# Deconvolution
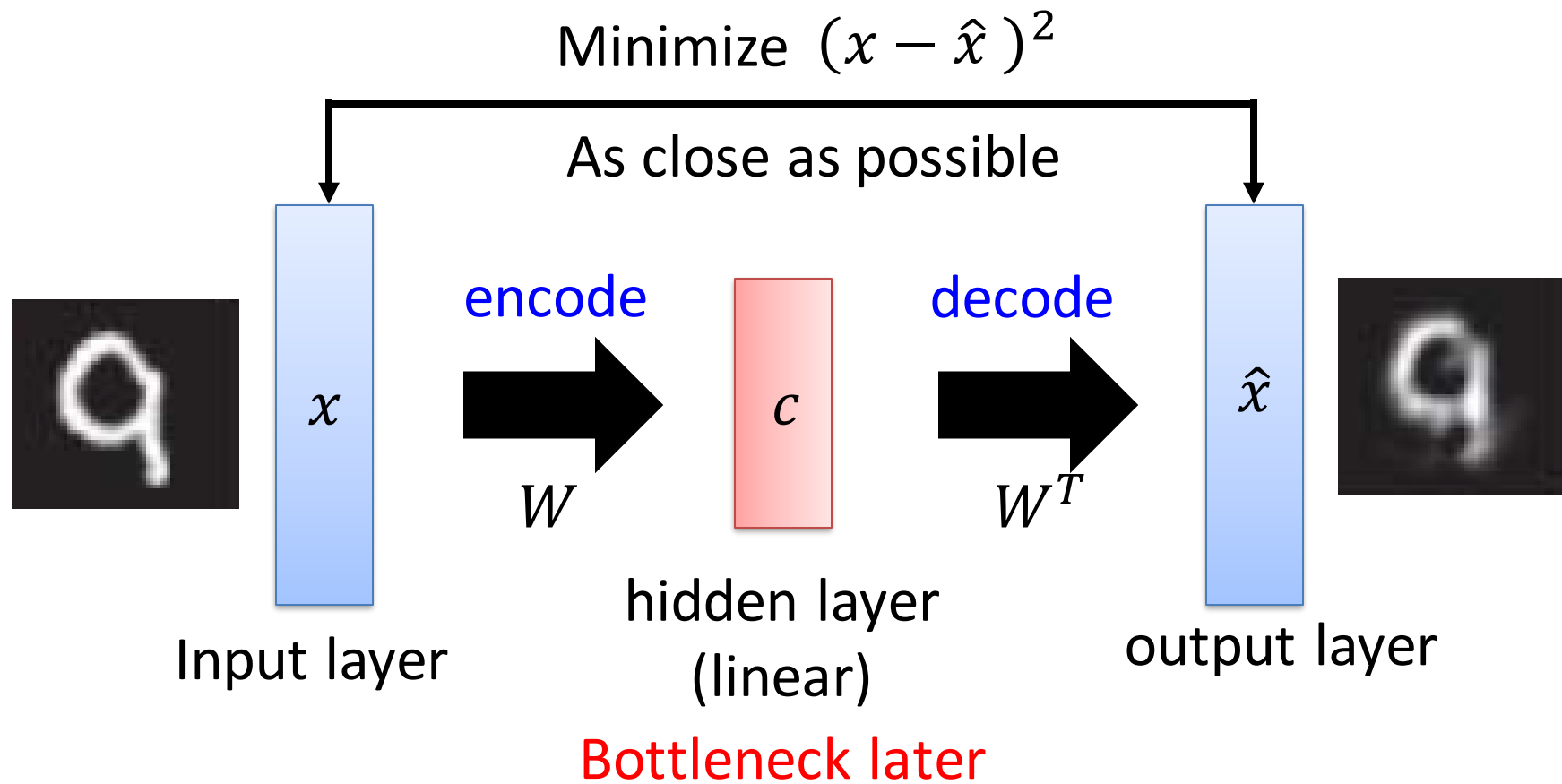
Actually, deconvolution is convolution.

# Pokémon

http://140.112.21.35:2880/~tlkagk/pokemon/pca.html
http://140.112.21.35:2880/~tlkagk/pokemon/auto.html

The code is modified from

http://jkunst.com/r/pokemon-visualize-em-all/

# PCA ~ Autoencoder with linear layers

Minimize $(x - \hat{x})^2$

As close as possible

encode $\qquad$ decode

$x$ $\qquad$ $c$ $\qquad$ $\hat{x}$

$W$ $\qquad\qquad$ $W^T$

Input layer $\qquad$ hidden layer (linear) $\qquad$ output layer

Bottleneck later

Output of the hidden layer is the code

# Code Reference

https://blog.keras.io/building-autoencoders-in-keras.html