# Transfer Learning

"DON'T BE A HERO"

ANDREJ KARPATHY

# Your first driving lesson

# Your first driving lesson

Imagine learning to drive a car without knowing absolutely nothing about anything

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

# Your first driving lesson

**Randomly initialized Deep Neural Network**

# Your first driving lesson

Any previous learning can be useful

Knowing how to cook is better than knowing nothing at all

# Your first driving lesson

Any previous learning can be useful

Knowing how to cook is better than knowing nothing at all

**We naturally reuse what we previously learnt to be able to solve a new task.**

# How **transfer** learning emerged

# Image classification

• **1998 LeNet-5**

Gradient-based learning applied to document recognition.
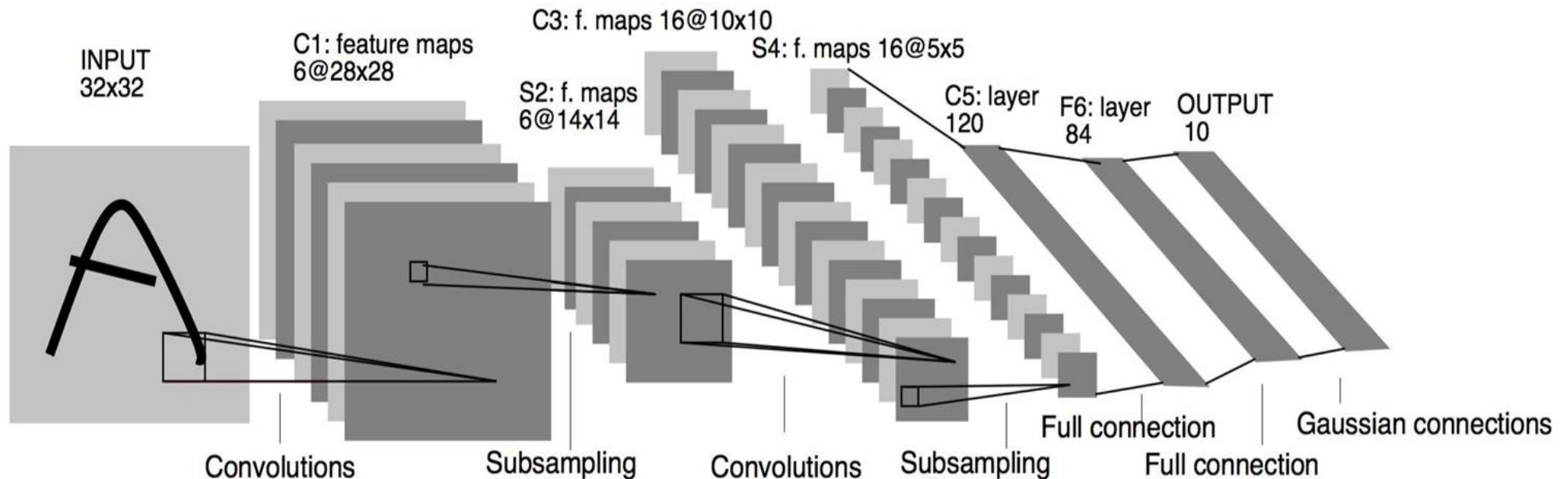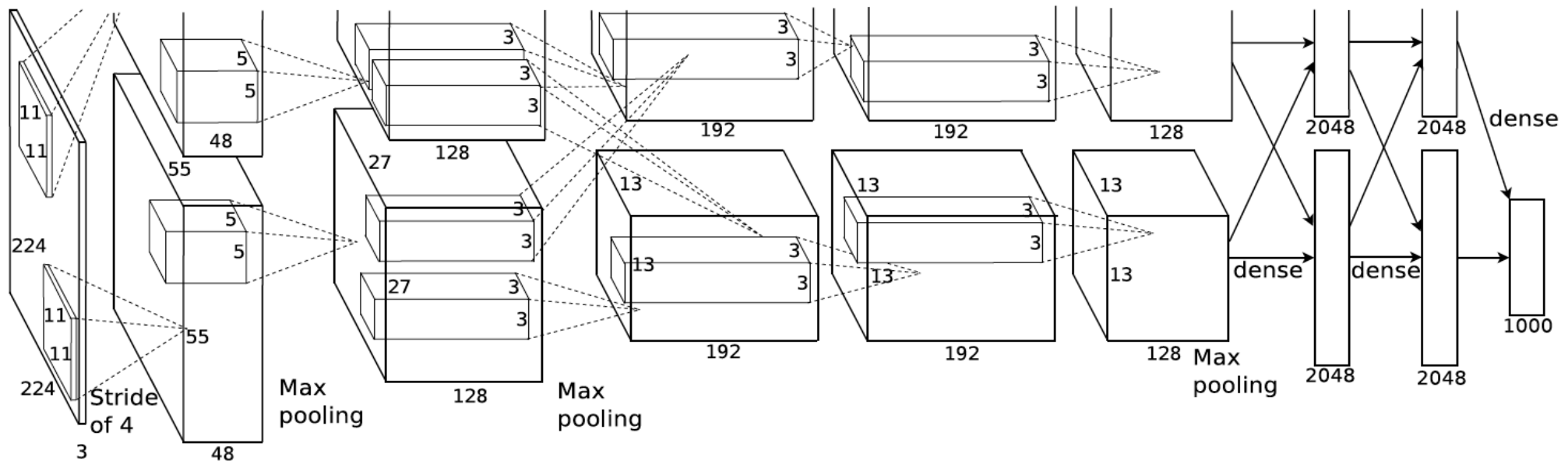Yann LeCun, Léon Bottou, Yoshua Bengio, Patrick Haffner

# Image classification

- **2012 AlexNet**

  ImageNet Classification with Deep Convolutional Neural Networks
  Alex Krizhevsky, Ilya Sutskever and Geoffrey Hinton

**1998 LeNet-5**
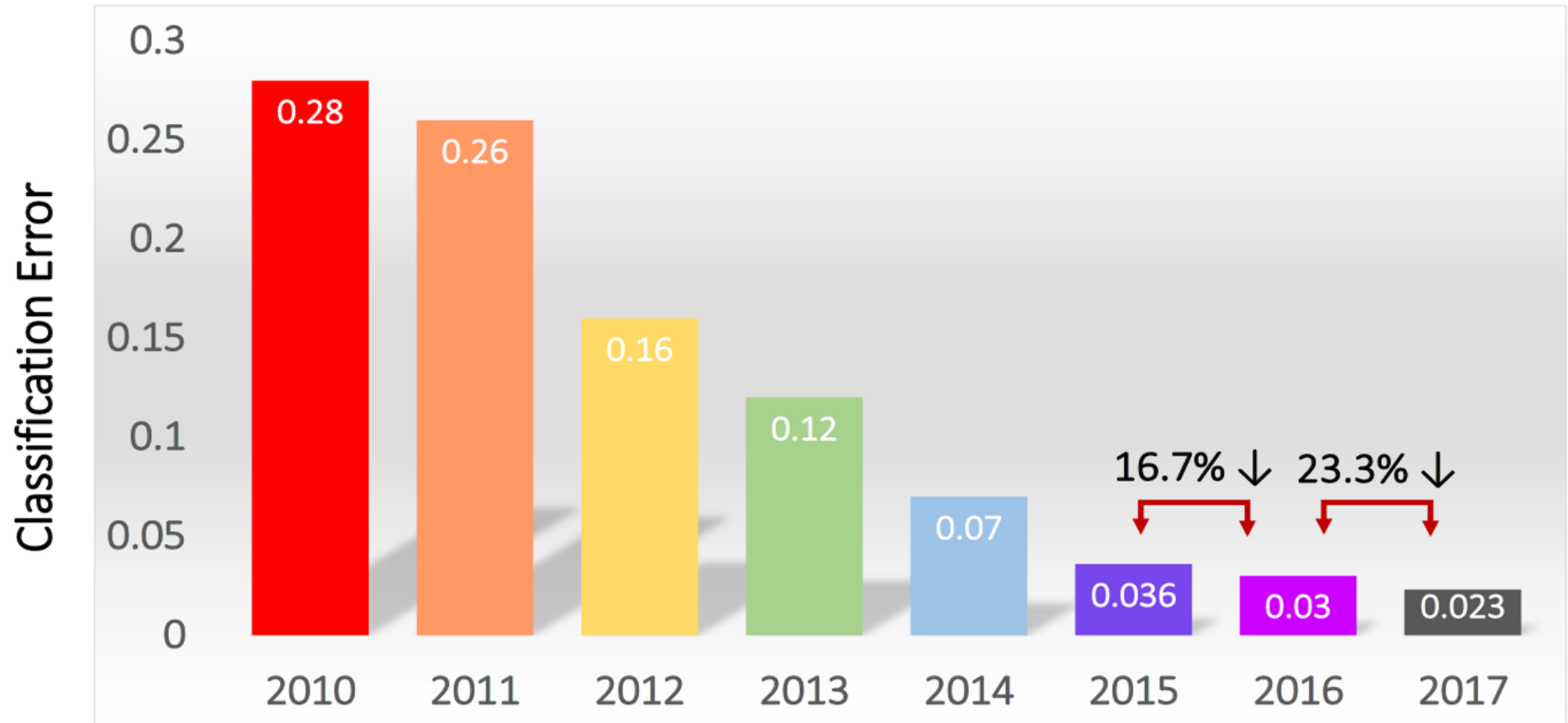
**2012 AlexNet**

2014 VGG19

2014 GoogLeNet

2015 Inception-V3

2015 ResNet-56



THAT'S NOT ENOUGH

WE NEED TO GO DEEPER

memegenerator.net

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

# ImageNet classification results



From image-net.org

# At what price?

- Data available
  - 1,000 images per class
- Computational cost
  - Specific hardware
  - Energy cost
- Human effort
  - Highly skilled professionals
  - Architecture design
  - Hyper-parameter fine tuning

**Barcelona Supercomputing Center**
*Centro Nacional de Supercomputación*

# At what price?

- Data available
  - 1,000 images per class
- Computational cost
  - Specific hardware
  - Energy cost
- Human effort
  - Highly skilled professionals
  - Architecture design
  - Hyper-parameter fine tuning

**We can't do that for every single problem!!**

# At what price?

- Data available
  - 1,000 images per class
- Computational cost
  - Specific hardware
  - Energy cost
- Human effort
  - Highly skilled professionals
  - Architecture design
  - Hyper-parameter fine tuning

**We don't want to** **do that for every single problem!!**
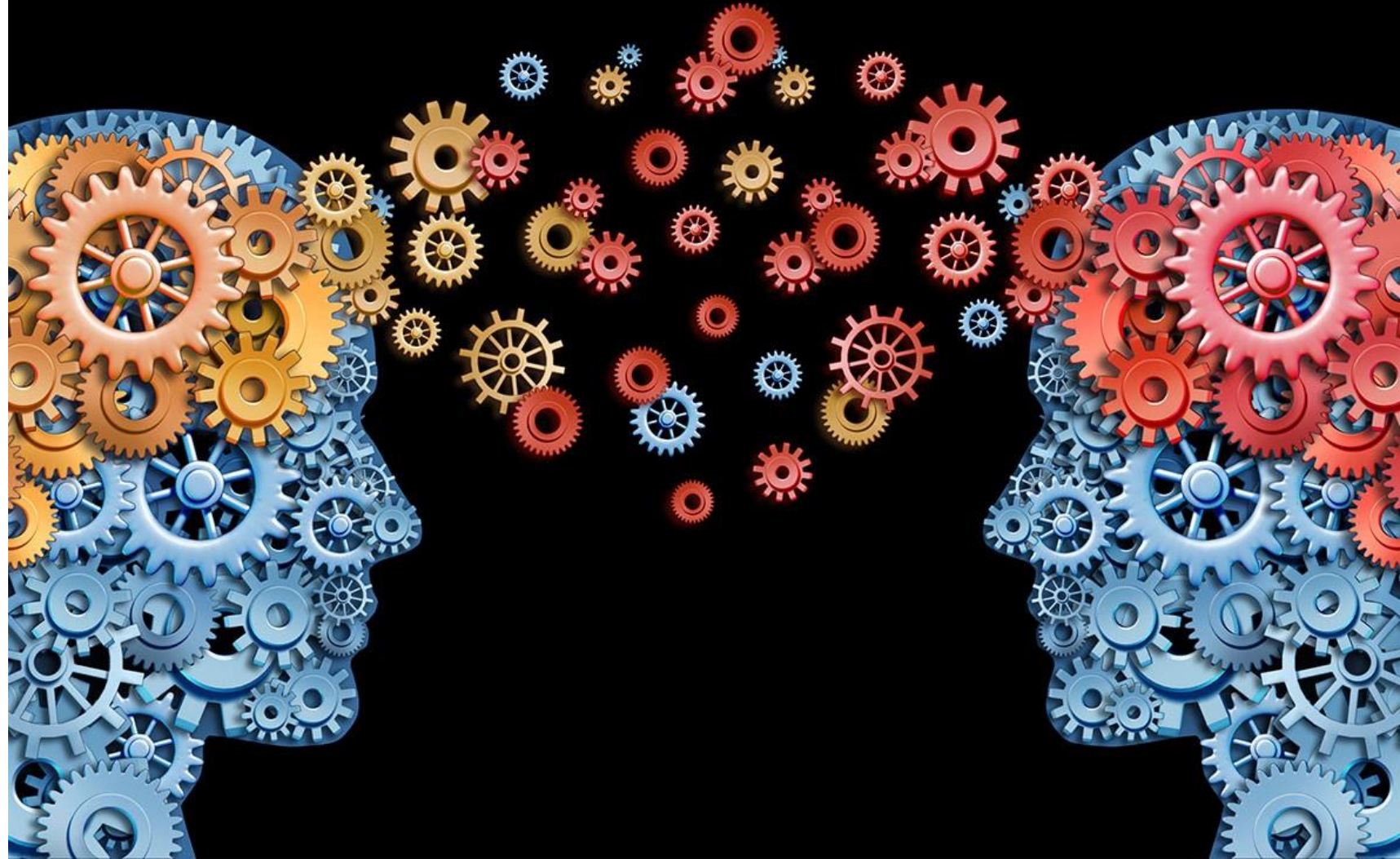
# At what price?

- Data available
  - 1,000 images per class
- Computational cost
  - Specific hardware
  - Energy cost
- Human effort
  - Highly skilled professionals
  - Architecture design
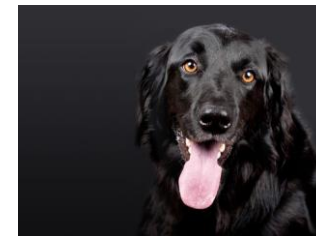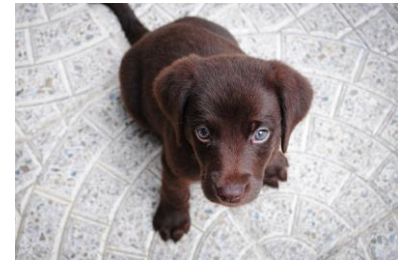  - Hyper-parameter fine tuning

**We don't want to do that for every single problem!!**

→ **Transfer Learning to the rescue**

# What is transfer learning?
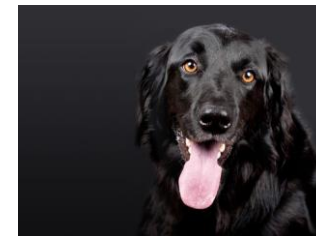
# What is learning about?

# What is learning about?
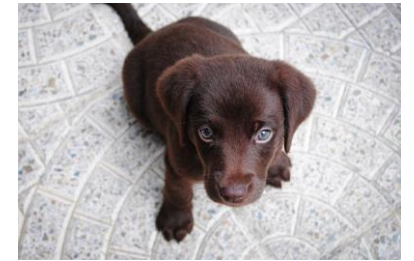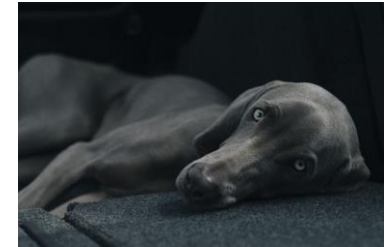
**Train set**

# What is learning about?

**Train set**

**Test set**

# What is learning about?

**MODEL**

TRAIN

PREDICT

**Train set**

**Test set**

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

What is learning about?

MODEL

TRAIN

PREDICT

Test set

Train set

FAIL

# What is learning about?

**MODEL**

**TRAIN**

**PREDICT**

**Test set**

**Train set**

23

# What is learning about?

GREEN BACKGROUND

**Train set**

**Test set**

Barcelona Supercomputing Center
Centro Nacional de Supercomputación

# What is learning about?

GREEN BACKGROUND → DOG

Test set

Train set

Barcelona Supercomputing Center
Centro Nacional de Supercomputación

What is learning about?

GREEN BACKGROUND → DOG

Test set

Train set

Barcelona Supercomputing Center
Centro Nacional de Supercomputación

What is learning about?

GREEN BACKGROUND → DOG

Test set

Train set

Train and test sets have different conditional probability distributions

Barcelona Supercomputing Center
Centro Nacional de Supercomputación

27

What is learning about?

THEY WILL NEVER BE **EXACTLY** EQUAL

We must **seek** to have **similar** conditional probability distributions in train and test sets

Test set

Train set

# What is transfer learning about?

**Train set**

**Test set**

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

# What is transfer learning about?

Train a **Machine Learning Model** on a **train set** with the hope that what has been learnt will be useful to solve a **different task.**

**Train set**



**GENERALIZATION**

**Test set**



**Train** and **test sets** are drawn from a **not so similar** underlying probability distribution.

Barcelona Supercomputing Center
Centro Nacional de Supercomputación

34

# Formalizing transfer learning

Pan, Sinno Jialin, and Qiang Yang.
**"A survey on transfer learning."**
IEEE Transactions on knowledge and data engineering (2010)



$$\mathcal{D_S} \neq \mathcal{D_T} \Rightarrow \mathcal{X_S} \neq \mathcal{X_T} \vee \mathbf{P_S}(X) \neq \mathbf{P_T}(X)$$

$$\mathcal{T} = \{y, f(\cdot)\}$$

$$\mathcal{D_S} \neq \mathcal{D_T}$$

$$\mathcal{D} = \{\mathcal{X}, P(X)\}$$

# Formalizing transfer learning

**Domain:**

**Task:**

# Formalizing **transfer** learning

**Domain:** $\mathcal{D} = \{\mathcal{X}, P(X)\}$

- A feature space $\mathcal{X}$



$\neq$  "The Elgar Concert Hall at the University of Birmingham for our third conference"

➜ Bag of words

➜ Content vector

- A marginal probability distribution $P(X), where\ X = \{x_1, \dots, x_n\} \in \mathcal{X}$



$\neq$



**Task:**

# Formalizing **transfer** learning

**Domain:** $\mathcal{D} = \{\mathcal{X}, P(X)\}$

- A feature space $\mathcal{X}$



"The Elgar Concert Hall at the University of Birmingham for our third conference"  → Bag of words

→ Content vector

- A marginal probability distribution $P(X), where\ X = \{x_1, \dots, x_n\} \in \mathcal{X}$



**Task:** $\mathcal{T} = \{\mathcal{Y}, f(\cdot)\}$

- A label space $\mathcal{Y}$

   **CAT, DOG ≠ LION, WOLF**

- An objective predictive function $f(\cdot) \Leftrightarrow P(y|x)$



Barcelona Supercomputing Center
Centro Nacional de Supercomputación

38

# Formalizing transfer learning

**Domain:** $\mathcal{D} = \{\mathcal{X}, P(X)\}$

**Task:** $\mathcal{T} = \{\mathcal{Y}, f(\cdot)\}$

# Formalizing transfer learning

**Source**  **Target**

**Domain:** $\mathcal{D} = \{\mathcal{X}, P(X)\}$

- A feature space $\mathcal{X}$
  - **The Same (different)**
- A marginal probability distribution $P(X)$
  - **Different**
  - **Similar**

**Task:** $\mathcal{T} = \{\mathcal{Y}, f(\cdot)\}$

# Formalizing transfer learning

**Source**          **Target**

**Domain:** $\mathcal{D} = \{\mathcal{X}, P(X)\}$

- A feature space $\mathcal{X}$
  - **The Same (different)**



- A marginal probability distribution $P(X)$
  - **Different**
  - **Similar**



**Task:** $\mathcal{T} = \{\mathcal{Y}, f(\cdot)\}$

- A label space $\mathcal{Y}$
  - **Different**
  - **The same**

**{CAT, DOG}**          **{LION, WOLF}**
**{FELINE, CANINE}**      **{FELINE, CANINE}**

- An objective predictive function          $f_S(\cdot)$          $f_T(\cdot)$
  - **Different (but similar?)**

# What is transfer learning about?

**Train set**

**Test set**

Barcelona Supercomputing Center
Centro Nacional de Supercomputación

# What is **transfer** learning about?

**Source domain**

**Target domain**

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

# What is transfer learning about?

**Source domain**

**Source Task** $\mathcal{Y}=\{CAT, DOG\}, f_S(\cdot)$

**Target domain**

**Target Task** $\mathcal{Y}=\{LION, WOLF\}, f_T(\cdot)$

# Formalizing **transfer** learning

**Domain:** $\mathcal{D} = \{\mathcal{X}, P(X)\}$

- A feature space $\mathcal{X}$
  - **The Same (different)**



- A marginal probability distribution $P(X)$
  - **Different**
  - **Similar**



**Task:** $\mathcal{T} = \{\mathcal{Y}, f(\cdot)\}$

- A label space $\mathcal{Y}$
  - **Different**
  - **The same**

**{CAT, DOG}**          **{LION, WOLF}**
**{FELINE, CANINE}**          **{FELINE, CANINE}**

- An objective predictive function          $f_S(\cdot)$          $f_T(\cdot)$
  - **Different (but similar?)**

45

# Formalizing transfer learning

**Source**    **Target**

**Domain:** $\mathcal{D} = \{\mathcal{X}, P(X)\}$

- A feature space $\mathcal{X}$
  - **The Same (different)**

- A marginal probability distribution $P(X)$
  - **Different**
  - **Similar**

**Task:** $\mathcal{T} = \{\mathcal{Y}, f(\cdot)\}$

- A label space $\mathcal{Y}$
  - **Different**
  - **The same**

**{CAT, DOG}**    **{LION, WOLF}**
**{FELINE, CANINE}**    **{FELINE, CANINE}**

- An objective predictive function    $f_S(\cdot)$    $f_T(\cdot)$
  - **Different (but similar?)**

46

# Formalizing transfer learning

$\hat{f}_S(\cdot) =$ 

$\hat{f}_T(\cdot) =$ 

- Are they similar?
- Can we just use $\hat{f}_S(\cdot)$ to approximate $f_T(\cdot)$?
- Can we reuse part of it?

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

# Representation learning

# Deep Neural Networks are **representation** learning techniques

- **Support Vector Machine** (SVM) is just a **classifier** (a very good one).

- SVM find the best boundary separating the data instances into different classes in a **given** feature space.

# Deep Neural Networks are **representation** learning techniques

- **Support Vector Machine** (SVM) is just a **classifier** (a very good one).

- SVM find the best boundary separating the data instances into different classes in a **given** feature space.

# Deep Neural Networks are **representation** learning techniques

- SVMs using the **kernel trick** can overcome the linear limitation through an **implicit** mapping to a higher dimensional feature space

# Deep Neural Networks are representation learning techniques

# Deep Neural Networks are **representation** learning techniques



Input

Linear Classifier

Output

maxpool

maxpool

depth=64
3x3 conv
conv1_1
conv1_2

depth=128
3x3 conv
conv2_1
conv2_2

maxpool
depth=256
3x3 conv
conv3_1
conv3_2
conv3_3
conv3_4

maxpool
depth=512
3x3 conv
conv4_1
conv4_2
conv4_3
conv4_4

maxpool
depth=512
3x3 conv
conv5_1
conv5_2
conv5_3
conv5_4

size=4096
FC1
FC2
size=1000
softmax

# Deep Neural Networks are **representation** learning techniques



**Intermediate representations**

**Linear Classifier**

Input

Output

maxpool

maxpool

depth=256
3x3 conv
conv3_1
conv3_2
conv3_3
conv3_4

maxpool

depth=512
3x3 conv
conv4_1
conv4_2
conv4_3
conv4_4

maxpool

depth=512
3x3 conv
conv5_1
conv5_2
conv5_3
conv5_4

maxpool

depth=64
3x3 conv
conv1_1
conv1_2

depth=128
3x3 conv
conv2_1
conv2_2

size=4096
FC1
FC2
size=1000
softmax

# Reusing DNNs knowledge

SAVE THE EARTH

REUSE DNNs

# Reusing DNNs knowledge

- Feature Extraction

- Fine-tuning

# Reusing DNNs knowledge

- **Feature Extraction**

- Fine-tuning

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

# Feature extraction



Intermediate representations

Linear Classifier

Input

Output

maxpool

depth=64
3x3 conv
conv1_1
conv1_2

depth=128
3x3 conv
conv2_1
conv2_2

maxpool

depth=256
3x3 conv
conv3_1
conv3_2
conv3_3
conv3_4

maxpool

depth=512
3x3 conv
conv4_1
conv4_2
conv4_3
conv4_4

maxpool

depth=512
3x3 conv
conv5_1
conv5_2
conv5_3
conv5_4

maxpool

size=4096
FC1
FC2
size=1000
softmax

# Feature extraction

# Feature extraction



Intermediate representations

Linear Classifier

Input

Output

maxpool

maxpool

maxpool

maxpool

maxpool

depth=64
3x3 conv
conv1_1
conv1_2

depth=128
3x3 conv
conv2_1
conv2_2

depth=256
3x3 conv
conv3_1
conv3_2
conv3_3
conv3_4

depth=512
3x3 conv
conv4_1
conv4_2
conv4_3
conv4_4

depth=512
3x3 conv
conv5_1
conv5_2
conv5_3
conv5_4

size=4096
FC1
FC2
size=1000
softmax

# Feature extraction

# Feature extraction

# Simple solutions

- DNN last layer features + SVM
  (Feature extraction)

# Simple solutions

- DNN last layer features + SVM
    (Feature extraction)
    We need: **Similar task and domain**

# Reusing DNNs knowledge

- Feature Extraction

- **Fine-tuning**

# Fine tuning

# Fine tuning

**What if the tasks are quite different?**

# Fine tuning

**What if the tasks are quite different?**



Source Task

Intermediate representations

Linear Classifier

Input

Output

maxpool

depth=64
3x3 conv
conv1_1
conv1_2

depth=128
3x3 conv
conv2_1
conv2_2

maxpool

depth=256
3x3 conv
conv3_1
conv3_2
conv3_3
conv3_4

maxpool

depth=512
3x3 conv
conv4_1
conv4_2
conv4_3
conv4_4

maxpool

depth=512
3x3 conv
conv5_1
conv5_2
conv5_3
conv5_4

maxpool

size=4096
FC1
FC2
size=1000
softmax

# Fine tuning



**Intermediate representations**

**Input**

Features learned for the **Source Task**

maxpool

depth=64
3x3 conv
conv1_1
conv1_2

depth=128
3x3 conv
conv2_1
conv2_2

maxpool
depth=256
3x3 conv
conv3_1
conv3_2
conv3_3
conv3_4

maxpool
depth=512
3x3 conv
conv4_1
conv4_2
conv4_3
conv4_4

maxpool
depth=512
3x3 conv
conv5_1
conv5_2
conv5_3
conv5_4

maxpool
size=4096
FC1
FC2
size=1000
softmax

# Fine tuning



**Intermediate representations**

**Input**

maxpool

depth=64
3x3 conv
conv1_1
conv1_2

depth=128
3x3 conv
conv2_1
conv2_2

maxpool

maxpool

depth=256
3x3 conv
conv3_1
conv3_2
conv3_3
conv3_4

depth=512
3x3 conv
conv4_1
conv4_2
conv4_3
conv4_4

maxpool

depth=512
3x3 conv
conv5_1
conv5_2
conv5_3
conv5_4

maxpool

size=4096
FC1
FC2
size=1000
softmax

Features learned for the **Source Task**

**Can we make them better?**

# Fine tuning

# Fine tuning

# Fine tuning

# Fine tuning

# Fine tuning

# Fine tuning

# Fine tuning



**Error back-propagation**

**Intermediate representations**

**Input**

D N N

D N N

**Target Task Labels**

- Difficult to control how much do we re-train.
  - → Reduced learning rate (1/10)
  - → Early stopping
  - → Alternate source/target sampling

maxpool   maxpool   maxpool   maxpool   maxpool

# Simple solutions

- DNN last layer features + SVM
  (Feature extraction)
  We need: **Similar task and domain**



- Add one or several NN layers +
  Fine-tuning pre-trained layers

# Simple solutions

- DNN last layer features + SVM
  (Feature extraction)
  We need: **Similar task and domain**


- Add one or several NN layers +
  Fine-tuning pre-trained layers
  We need: **Enough data**

# Beyond the last layer

# Knowledge inside DNN

# Knowledge inside DNN

# Knowledge inside DNN

# Knowledge inside DNN



# Filters
Depth

Image height

Image width

maxpool    maxpool    maxpool

depth=256    depth=512    depth=512
3x3 conv    3x3 conv    3x3 conv
conv3_1    conv4_1    conv5_1
conv3_2    conv4_2    conv5_2
conv3_3    conv4_3    conv5_3
conv3_4    conv4_4    conv5_4

conv1_2    conv2_2

size=4096
FC1
FC2
size=1000
softmax

# Knowledge inside DNN

# Knowledge inside DNN

# Knowledge inside DNN

# Knowledge inside DNN

# Knowledge inside DNN



Conv3_3

FC7

# Knowledge **inside** DNN



Conv3_3

FC7

**Input Domain**

More influenced by **domain**

More influenced by **task**

**Task Labels**

# Knowledge inside DNN



Conv3_3

FC7

**Input Domain**

**Task Labels**

More influenced by **domain**

More influenced by **task**

**Simple features**

**Complex features**

Visualizations from: Yosinski, Jason, et al. "Understanding neural networks through deep visualization." *arXiv preprint arXiv:1506.06579* (2015).

# Fine tuning inside DNN?



Conv3_3

FC7

**Input Domain**

**Task Labels**

Fine-tuning

Intermediate representations

Error back-propagation

Input

Target Task Labels

# Fine tuning inside DNN?



Conv3_3    FC7

**Input Domain**

**Task Labels**

Fine-tuning

Intermediate representations

Error back propagation

Input

maxpool    maxpool    maxpool    maxpool    maxpool

depth=64    depth=128    depth=256    depth=512    depth=512    size=4096
3x3 conv    3x3 conv    3x3 conv    3x3 conv    3x3 conv    FC1
conv1_1    conv2_1    conv3_1    conv4_1    conv5_1    FC2
conv1_2    conv2_2    conv3_2    conv4_2    conv5_2    size=1000
                       conv3_3    conv4_3    conv5_3    softmax
                       conv3_4    conv4_4    conv5_4

D N N    D N N

Target Task Labels

Barcelona Supercomputing Center
Centro Nacional de Supercomputación

# Fine tuning inside DNN?



Conv3_3

FC7

**Input Domain**

**Task Labels**

Fine-tuning

# Feature extraction inside DNN?



Conv3_3

FC7

Input Domain

Task Labels

Feature extraction

Intermediate representations

Input

S V M

depth=64
3x3 conv
conv1_1
conv1_2

maxpool
depth=128
3x3 conv
conv2_1
conv2_2

maxpool
depth=256
3x3 conv
conv3_1
conv3_2
conv3_3
conv3_4

depth=512
3x3 conv
conv4_1
conv4_2
conv4_3
conv4_4

maxpool
depth=512
3x3 conv
conv5_1
conv5_2
conv5_3
conv5_4

maxpool
size=4096
FC1
FC2
size=1000
softmax

Target Task Labels

# Feature extraction inside DNN?



Conv3_3

FC7

Input Domain

Task Labels

Feature extraction

Intermediate representations

Input

maxpool
depth=64
3x3 conv
conv1_1
conv1_2

maxpool
depth=128
3x3 conv
conv2_1
conv2_2

maxpool
depth=256
3x3 conv
conv3_1
conv3_2
conv3_3
conv3_4

depth=512
3x3 conv
conv4_1
conv4_2
conv4_3
conv4_4

maxpool
depth=512
3x3 conv
conv5_1
conv5_2
conv5_3
conv5_4

size=4096
FC1
FC2
size=1000
softmax

S V M

Target Task Labels

Barcelona Supercomputing Center
Centro Nacional de Supercomputación

96

# Feature extraction inside DNN?



Conv3_3

FC7

**Input Domain**

**Task Labels**

Feature extraction

Intermediate representations

Input

maxpool
depth=64
3x3 conv
conv1_1
conv1_2

depth=128
3x3 conv
conv2_1
conv2_2

maxpool
depth=256
3x3 conv
conv3_1
conv3_2
conv3_3
conv3_4

maxpool
depth=512
3x3 conv
conv4_1
conv4_2
conv4_3
conv4_4

maxpool
depth=512
3x3 conv
conv5_1
conv5_2
conv5_3
conv5_4

size=4096
FC1
FC2
size=1000
softmax

S V M

Target Task Labels

Barcelona Supercomputing Center
Centro Nacional de Supercomputación

# Feature extraction inside DNN?



Conv3_3

FC7

Input Domain

Task Labels

Intermediate representations

Input

Feature extraction

maxpool
maxpool maxpool maxpool maxpool
depth=64 depth=128 depth=256 depth=512 depth=512
3x3 conv 3x3 conv 3x3 conv 3x3 conv 3x3 conv
conv1_1 conv2_1 conv3_1 conv4_1 conv5_1
conv1_2 conv2_2 conv3_2 conv4_2 conv5_2
conv3_3 conv4_3 conv5_3
conv3_4 conv4_4 conv5_4

size=4096
FC1
FC2
size=1000
softmax

S V M

Target Task Labels

# Feature extraction inside DNN?

**VGG16 dim: 12,416**

Conv3_3

FC7

**Input Domain**

**Task Labels**

Feature extraction

Barcelona Supercomputing Center
Centro Nacional de Supercomputación

# Feature extraction inside DNN?

VGG16 dim: 12,416

Conv3_3

FC7



Input Domain

Task Labels

Feature extraction

Different range of values

Intermediate representations

Input

maxpool
maxpool
maxpool
maxpool
maxpool

depth=64
3x3 conv
conv1_1
conv1_2

depth=128
3x3 conv
conv2_1
conv2_2

depth=256
3x3 conv
conv3_1
conv3_2
conv3_3
conv3_4

depth=512
3x3 conv
conv4_1
conv4_2
conv4_3
conv4_4

depth=512
3x3 conv
conv5_1
conv5_2
conv5_3
conv5_4

size=4096
FC1
FC2
size=1000
softmax

SVM

Target Task Labels

Barcelona Supercomputing Center
Centro Nacional de Supercomputación

# Feature behavior in Transfer Learning



Conv3_3

FC7

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

# Standardization: features in same range

Typically, it is used a L2-normalization

# Standardization: features in same range

**Typically, it is used a L2-normalization**

L2 norm

0,00  0,13  0,27  0,40  0,54  0,67  0,80  0,94

0   50   100   150   200   250   300   350

Feature Standardisation

-5  -4  -3  -2  -1  0  1  2  3  4  5

**Each feature independently**

Barcelona Supercomputing Center
Centro Nacional de Supercomputación

# Standardization: features in same range

**Typically, it is used a L2-normalization**



**L2 norm**

**Instance based**

0,00    0,13    0,27    0,40    0,54    0,67    0,80    0,94

**Feature Standardisation**

**Feature based**

0    50    100    150    200    250    300    350

-5  -4  -3  -2  -1   0   1   2   3   4   5

**Each feature independently**

# Standardization: features in context

- Distinct feature behavior for a subset of data wrt rest

Conv3_3

FC7

Barcelona Supercomputing Center
Centro Nacional de Supercomputación

# Stardardization: features in context



fc7 n1946

Garcia-Gasulla et al.
*On the behavior of convolutional nets for feature extraction.* 2017

Gadwall
Brown Pelican
White Pelican
Heermann Gull

CUB-200 - **birds**

# Curse of dimensionality



VGG16 dim: 12,416

Conv3_3

FC7

Input Domain

Task Labels

# Curse of dimensionality

# Curse of dimensionality



VGG16 dim: 12,416

Conv3_3

FC7

Input Domain

Task Labels

FAIL

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

# Curse of dimensionality

**VGG16 dim: 12,416**



Conv3_3

FC7

**Input Domain**

**Task Labels**

**FAIL**

- Unsupervised: Only takes into account the **domain**

# Curse of dimensionality

VGG16 dim: 12,416

Conv3_3

FC7

Input Domain

Task Labels

Supervised Feature Selection

?

Barcelona Supercomputing Center
Centro Nacional de Supercomputación

# Curse of dimensionality

VGG16 dim: 12,416

Conv3_3

FC7

Input Domain

Task Labels

Supervised Feature Selection

?

- High computational **cost!**

# Curse of dimensionality

# Curse of dimensionality



VGG16 dim: 12,416

Conv3_3

FC7

Input Domain

Task Labels

- Which is the real problem?

Barcelona Supercomputing Center
Centro Nacional de Supercomputación

# Curse of dimensionality



**VGG16 dim: 12,416**

Conv3_3

FC7

**Input Domain**

**Task Labels**

- Which is the real problem?
  - Too many features?
  - Too few images?

# Curse of dimensionality



**VGG16 dim: 12,416**

Conv3_3

FC7

**Input Domain**

**Task Labels**

- Which is the real problem?
  - Too many features?
  - ~~Too few images?~~  **A requirement**

# Curse of dimensionality

VGG16 dim: 12,416



**Input Domain**

**Task Labels**

Conv3_3

FC7

- Which is the real problem?
  - ~~Too many features?~~
  - **Too much information!**
  - ~~Too few images?~~

# Meaningful discretization

Real values

Conv3_3    FC7

Quantitization to {-1,0,1} based on feature values

Discrete values

Conv3_3    FC7

1 -1 0 1 1 0 0 0 1 -1 0 1 1 0 0 0 1 -1 0 1 1 0 0 0 1 -1 0 1 1 0 0 0 1 0 0 0 1 -1 0 1 1 0 0 0 1 -1 0 1 1

# Full-Network Embedding **Recipe**

Garcia-Gasulla, et al. *An out-of-the-box full-network embedding for convolutional neural networks.* 2018.

1. Spatial Average Pooling

2. Standardisation

3. Discretization

Barcelona Supercomputing Center
Centro Nacional de Supercomputación

# Full Network embedding

# FNE - Results: Similar source task

Network pre-trained on **Places2** for mit67 and on **ImageNet** for the rest.

| Dataset | mit67 | cub200 | flowers102 | cats-dogs | sdogs | caltech101 | food101 | textures | wood |
|---------|-------|--------|------------|-----------|-------|------------|---------|----------|------|
| Baseline fc6 | 80.0 | 65.8 | 89.5 | 89.3 | 78.0 | $91.4_{\pm0.6}$ | $61.4_{\pm0.2}$ | 69.6 | $70.8_{\pm6.6}$ |
| Baseline fc7 | 81.7 | 63.2 | 87.0 | 89.6 | 79.3 | $89.7_{\pm0.3}$ | $59.1_{\pm0.6}$ | 69.0 | $68.9_{\pm6.8}$ |
| Full-network | 83.6 | 65.5 | 93.3 | 89.2 | 78.8 | $91.4_{\pm0.6}$ | $67.0_{\pm0.7}$ | 73.0 | $74.1_{\pm6.9}$ |
| SotA | 86.9 [5] | 92.3 [10] | 97.0 [5] | 91.6 [6] | 90.3 [5] | 93.4 [31] | 77.4 [4] | 75.5 [17] | - - |
| ED | ✓ | ✓ | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | - |
| FT | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | - |

# FNE - Results: Similar source task

Network pre-trained on **Places2** for mit67 and on **Imager...**

**Best case scenario!**

| Dataset | mit67 | cub200 | flowers102 | cats-dogs | sdogs | caltech101 | food101 | textures | wood |
|---|---|---|---|---|---|---|---|---|---|
| Baseline `fc6` | 80.0 | 65.8 | 89.5 | 89.3 | 78.0 | $91.4_{\pm0.6}$ | $61.4_{\pm0.2}$ | 69.6 | $70.8_{\pm6.6}$ |
| Baseline `fc7` | 81.7 | 63.2 | 87.0 | 89.6 | 79.3 | $89.7_{\pm0.3}$ | $59.1_{\pm0.6}$ | 69.0 | $68.9_{\pm6.8}$ |
| Full-network | 83.6 | 65.5 | 93.3 | 89.2 | 78.8 | $91.4_{\pm0.6}$ | $67.0_{\pm0.7}$ | 73.0 | $74.1_{\pm6.9}$ |
| SotA | 86.9 [5] | 92.3 [10] | 97.0 [5] | 91.6 [6] | 90.3 [5] | 93.4 [31] | 77.4 [4] | 75.5 [17] | - - |
| ED | ✓ | ✓ | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | - |
| FT | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | - |

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

# FNE - Results: Similar source task

Network pre-trained on **Places2** for mit67 and on **ImageNet** for the rest.

| Dataset | mit67 | cub200 | flowers102 | cats-dogs | sdogs | caltech101 | food101 | textures | wood | |
|---------|-------|--------|------------|-----------|-------|------------|---------|----------|------|---|
| Baseline fc6 | 80.0 | 65.8 | 89.5 | 89.3 | 78.0 | $91.4 \pm 0.6$ | $61.4 \pm 0.2$ | 69.6 | $70.8 \pm 6.6$ | +2.9 |
| Baseline fc7 | 81.7 | 63.2 | 87.0 | 89.6 | 79.3 | $89.7 \pm 0.3$ | $59.1 \pm 0.6$ | 69.0 | $68.9 \pm 6.8$ | +4.2 |
| Full-network | 83.6 | 65.5 | 93.3 | 89.2 | 78.8 | $91.4 \pm 0.6$ | $67.0 \pm 0.7$ | 73.0 | $74.1 \pm 6.9$ | |
| SotA | 86.9 [5] | 92.3 [10] | 97.0 [5] | 91.6 [6] | 90.3 [5] | 93.4 [31] | 77.4 [4] | 75.5 [17] | - | |
| ED | ✓ | ✓ | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | - | |
| FT | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | - | |

123

# FNE - Results: Similar source task

Network pre-trained on **Places2** for mit67 and on **ImageNet** for the rest.

| Dataset | mit67 | cub200 | flowers102 | cats-dogs | sdogs | caltech101 | food101 | textures | wood | |
|---|---|---|---|---|---|---|---|---|---|---|
| Baseline `fc6` | 80.0 | 65.8 | 89.5 | 89.3 | 78.0 | $91.4_{\pm0.6}$ | $61.4_{\pm0.2}$ | 69.6 | $70.8_{\pm6.6}$ | +2.9 |
| Baseline `fc7` | 81.7 | 63.2 | 87.0 | 89.6 | 79.3 | $89.7_{\pm0.3}$ | $59.1_{\pm0.6}$ | 69.0 | $68.9_{\pm6.8}$ | +4.2 |
| Full-network | 83.6 | -0.3 | 93.3 | -0.4 | -0.5 | $91.4_{\pm0.6}$ | $67.0_{\pm0.7}$ | 73.0 | $74.1_{\pm6.9}$ | |
| SotA | 86.9 [5] | 92.3 [10] | 97.0 [5] | 91.6 [6] | 90.3 [5] | 93.4 [31] | 77.4 [4] | 75.5 [17] | - - | |
| ED | ✓ | ✓ | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | - | |
| FT | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | - | |

Barcelona Supercomputing Center
Centro Nacional de Supercomputación

# FNE - Results: Similar source task

Network pre-trained on **Places2** for mit67 and on **ImageNet** for the rest.

| Dataset | mit67 | cub200 | flowers102 | cats-dogs | sdogs | caltech101 | food101 | textures | wood | |
|---|---|---|---|---|---|---|---|---|---|---|
| Baseline fc6 | 80.0 | 65.8 | 89.5 | 89.3 | 78.0 | 91.4±0.6 | 61.4±0.2 | 69.6 | 70.8±6.6 | +2.9 |
| Baseline fc7 | 81.7 | 63.2 | 87.0 | 89.6 | 79.3 | 89.7±0.3 | 59.1±0.6 | 69.0 | 68.9±6.8 | +4.2 |
| Full-network | 83.6 | -0.3 | 93.3 | -0.4 | -0.5 | 91.4±0.6 | 67.0±0.7 | 73.0 | 74.1±6.9 | |
| SotA | 86.9 [5] | 92.3 [10] | 97.0 [5] | 91.6 [6] | 90.3 [5] | 93.4 [31] | 77.4 [4] | 75.5 [17] | - - | |
| ED | ✓ | ✓ | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | - | |
| FT | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | - | |

# FNE - Results: Dissimilar source task

Network pre-trained on **ImageNet** for mit67 and on **Places2** for the rest.

| Dataset | mit67 | cub200 | flowers102 | cats-dogs | caltech101 | textures | wood |
|---|---|---|---|---|---|---|---|
| Baseline `fc7` | 72.2 | 23.6 | 73.3 | 38.7 | 72.0 | 55.8 | 65.3 |
| Full-network | 75.5 | 35.5 | 88.7 | 56.2 | 80.0 | 65.1 | 74.0 |

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

# FNE - Results: Dissimilar source task

**Most frequent real-world scenario!**

Network pre-trained on **ImageNet** for mit67 and on **Places2** for the rest.

| Dataset | mit67 | cub200 | flowers102 | cats-dogs | caltech101 | textures | wood |
|---|---|---|---|---|---|---|---|
| Baseline `fc7` | 72.2 | 23.6 | 73.3 | 38.7 | 72.0 | 55.8 | 65.3 |
| Full-network | 75.5 | 35.5 | 88.7 | 56.2 | 80.0 | 65.1 | 74.0 |

# FNE - Results: Dissimilar source task

Network pre-trained on **ImageNet** for mit67 and on **Places2** for the rest.

| Dataset | mit67 | cub200 | flowers102 | cats-dogs | caltech101 | textures | wood |
|---|---|---|---|---|---|---|---|
| Baseline fc7 | 72.2 | 23.6 | 73.3 | 38.7 | 72.0 | 55.8 | 65.3 |
| Full-network | 75.5 | 35.5 | 88.7 | 56.2 | 80.0 | 65.1 | 74.0 |
| | +3.3 | +11.9 | +15.4 | +17.5 | +8.0 | +9.3 | +10.6 |

# Simple solutions

- DNN last layer features + SVM
  (Feature extraction)
  We need: **Similar task and domain**



- Add one or several NN layers +
  Fine-tuning pre-trained layers
  We need: **Enough data**



- Full Network Embedding
  – **Robust to different task and domain**
  – **Works with little data**

# Practical tips

# Practical tips

**Fine-tuning**

- Whenever possible **don't start from scratch.**

# Practical tips

**Fine-tuning**

- Whenever possible **don't start from scratch.**
- **External data** can help prevent overfitting (even from a different problem).

# Practical tips

**Fine-tuning**

- Whenever possible **don't start from scratch.**
- **External data** can help prevent overfitting (even from a different problem).
- Begin **freezing** as much as possible and proceed with caution (particularly for large models)

# Practical tips

**Feature extraction**

- Easy **baseline** for every problem.

# Practical tips

**Feature extraction**

- Easy **baseline** for every problem.
- **ImageNet,** a model to pre-train them all. (but not always)

# Practical tips

**Feature extraction**

- Easy **baseline** for every problem.
- **ImageNet,** a model to pre-train them all. (but not always)
- Always **normalize** features.

# Practical tips

**Feature extraction**

- Easy **baseline** for every problem.
- **ImageNet,** a model to pre-train them all. (but not always)
- Always **normalize** features.
- If source and target task are closely related:
  - → **Last two layers** are your best chance.

Barcelona Supercomputing Center
Centro Nacional de Supercomputación

# Practical tips

**Feature extraction**

- Easy **baseline** for every problem.
- **ImageNet,** a model to pre-train them all. (but not always)
- Always **normalize** features.
- If source and target task are closely related:
  - → **Last two layers** are your best chance.
- If source and target task are quite different:
  - → **Try everything**
  - → **Use FNE**

# Useful References

**Don't start training from scratch**

- Zhe Xu, Shaoli Huang, Ya Zhang, and Dacheng Tao. *Augmenting strong supervision using webdata for fine-grained categorization.* 2015.

- Steve Branson, Grant Van Horn, Serge Belongie, and Pietro Perona. *Bird species categorizationusing pose normalized deep convolutional nets*. 2014.

- Chang Liu, Yu Cao, Yan Luo, Guanling Chen, Vinod Vokkarane, and Yunsheng Ma. *Deep-food: Deep learning-based food image recognition for computer-aided dietary assessment.* 2016.

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

# Useful References

**Transfer learning 101**

- Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. *How transferable are features in deep neural networks?* 2014.

- Dario Garcia-Gasulla, Ferran Parés, Armand Vilalta, Jonatan Moreno, Eduard Ayguadé, Jesús Labarta, Ulises Cortés, and Toyotaro Suzumura. *On the behavior of convolutional nets for feature extraction.* 2017

# Useful References

**Basic feature extraction**

- Hossein Azizpour, Ali Sharif Razavian, Josephine Sullivan, Atsuto Maki, and Stefan Carlsson. *Factors of transferability for a generic convnet.* 2016.

- Ali Sharif Razavian, Hossein Azizpour, Josephine Sullivan, and Stefan Carlsson. *CNN features off-the-shelf: an astounding baseline for recognition.* 2014.

- Yunchao Gong, Liwei Wang, Ruiqi Guo, and Svetlana Lazebnik. *Multi-scale orderless pooling of deep convolutional activation features.* 2014.

- Jeff Donahue, Yangqing Jia, Oriol Vinyals, Judy Hoffman, Ning Zhang, Eric Tzeng, and Trevor Darrell. *Decaf: A deep convolutional activation feature for generic visual recognition.* 2014.

- Arsalan Mousavian and Jana Kosecka. *Deep convolutional features for image based retrieval and scene categorization.* 2015.

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

# Useful References

**Multi-layer feature extraction**

- Garcia-Gasulla, Dario, Armand Vilalta, Ferran Parés, Eduard Ayguadé, Jesus Labarta, Ulises Cortés, and Toyotaro Suzumura. *An out-of-the-box full-network embedding for convolutional neural networks.* 2018.
- Vilalta, Armand, Dario Garcia-Gasulla, Ferran Parés, Jonathan Moreno, Eduard Ayguadé Jesús Labarta, Ulises Cortés, and Toyotaro Suzumura. *Full-Network Embedding in a Multimodal Embedding Pipeline.* 2017.

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

# Useful References

**Advanced fine tuning to solve small tasks**

- Weifeng Ge and Yizhou Yu. *Borrowing treasures from the wealthy: Deep transfer learning through selective joint fine-tuning.* 2017.
- Marcel Simon and Erik Rodner. *Neural activation constellations: Unsupervised part model discovery with convolutional networks.* 2015.

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

# CTE-Power9

52 computing nodes. Each one:
- 2 x IBM Power9 8335-GTH @ 2.4GHz(total 160 threads)
- 512GB of main memory
- 4 x GPU NVIDIA V100 (Volta) with 16GB each

Available through PRACE and RES

**Hands on session 16:00**

**2 GPUs per account**



![Barcelona Supercomputing Center - Centro Nacional de Supercomputación](logo)

# thanks.

dario.garcia@bsc.es
armand.vilalta@bsc.es