# A Shallow Introduction into the Deep Machine Learning

"A quick tour from old principles to the most recent neural architectures"



Jan Čech

#### Outline of lectures:

- 1. Introduction, basic principles, layers, neural architectures, image recognition
- 2. Object detection, Semantic/Instance segmentation, further insight (Deep fakes, Adversarial examples, Visualization, Style transfer)
- 3. Generative modeling GANs, Diffusion models

#### Deep learning – top awards in science



m p

Deep learning pioneers receivedAlan Touring Prize in 2018







#### Nobel Prize in Physics 2024

 "for foundational discoveries and inventions that enable machine learning with artificial neural networks"



© Nobel Prize Outreach. Photo: Nanaka Adachi John J. Hopfield



© Nobel Prize Outreach. Photo: Clément Morin **Geoffrey Hinton** 

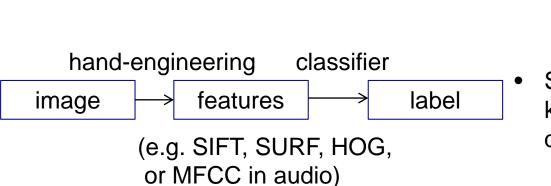
#### What is the "Deep Learning"?



Deep learning (by G. Hinton, DL pioneer, Touring+Nobel prize)

label

= both the classifiers and the features are learned automatically

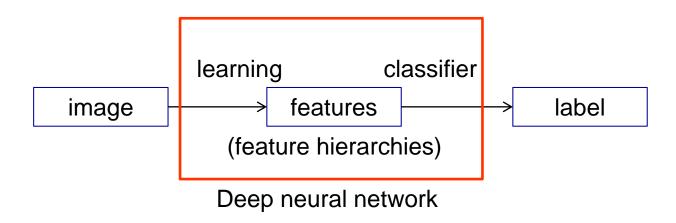


classifier

image

 Typically not feasible, due to high dimensionality

Suboptimal, requires expert knowledge, works in specific domain only



#### What is the "Deep Learning"? Other definitions...



- Andrew Ng (founder of Google Brain, chief of Baidu Al research)
  - "Very large neural networks we can now have and ... huge amounts of data that we have access to."
- Jeff Dean (head of Google AI)
  - "When you hear the term deep learning, just think of a large deep neural net. Deep refers to the number of layers typically and so this kind of the popular term that's been adopted in the press. I think of them as deep neural networks generally."
- Yoshua Bengio (DL pioneer, Turing Award Holder 2018)
  - "Deep learning algorithms seek to exploit the unknown structure in the input distribution in order to discover good representations, often at multiple levels, with higher-level learned features defined in terms of lower-level features."
- Yann LeCun (DL pioneer, Turing Award Holder 2018)
  - "Deep learning [is] ... a pipeline of modules all of which are trainable. ... deep because [has] multiple stages in the process of recognizing an object and all of those stages are part of the training."

#### Deep Learning omnipresent



- Besides the Computer Vision DL is extremely successful in, e.g.
  - Automatic Speech Recognition
    - Speech to text, Speaker recognition
  - Natural Language Processing (LLMs)
    - Machine translation, Question answering, Chatbots (GPT)
  - Robotics / Autonomous driving (e.g., Reinforcement learning)
    - Touring Award 2024 (Adrew G. Barto, R. S. Sutton)
  - Data Science / Bioinformatics (e.g., Alphafold)
    - Nobel Prize in Chemistry 2024 (D. Baker, D. Hassabis, J. Jumper)
- Shift of paradigm started in Computer Vision
  - Large-scale image category recognition (ILSVRC' 2012 challenge)

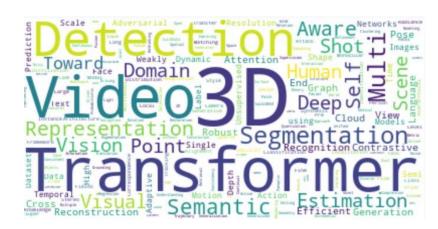
INRIA/Xerox 33%, Uni Amsterdam 30%. Uni Oxford 27%, Uni Tokyo 26%, 16% (deep neural network) [Krizhevsky-NIPS-20 **Uni Toronto** 

# Explosion of interest in "Deep Learning" after 2012

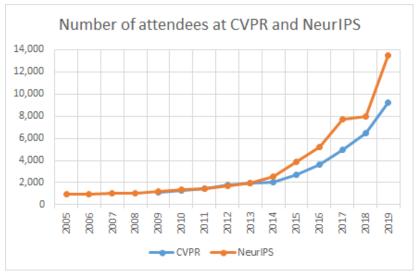


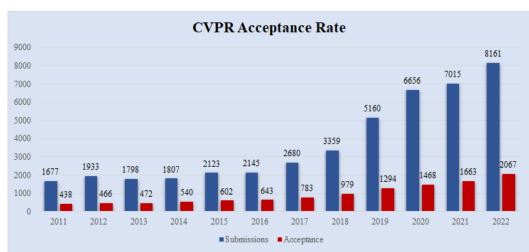
Paper title keywords, CVPR 2019/2022

action adaptation adversarial attention based clouds convolutional data deep depth detection domain efficient estimation face feature generative graph human image instance joint learning local matching model motion network neural object person point pose prediction recognition reconstruction representation robust scene segmentation semantic shape single structure supervised tracking transfer unsupervised video visual



 Number of attendees/submissions in major Computer Vision and Machine Learning grows exponentially





Data Source: <a href="https://hai.stanford.edu/">https://github.com/BIGBALLON/CVPR2022-Paper-Statistics</a>



- Image classification [Krizhevsky-NIPS-2012]
  - Input: RGB-image
  - Output: Single label (Probability Distribution over Classes)



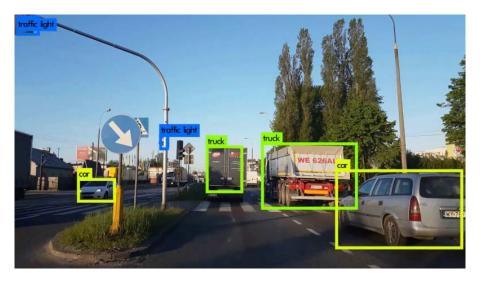


- ImageNet dataset (14M images, 21k classes, Labels by Amazon Mechanical Turk)
- ImageNet Benchmark (1000 classes, 1M training images)





- Object Detection
  - Multiple objects in the image [RCNN, YOLO, ...]



E.g. Face [<u>Hu-Ramanan-2017</u>], Text localization [<u>Busta-2017</u>]







10

- (3D) Pose estimation
  - [<u>Hu-2018</u>], [<u>OpenPose</u>]
  - [Cech-2016]

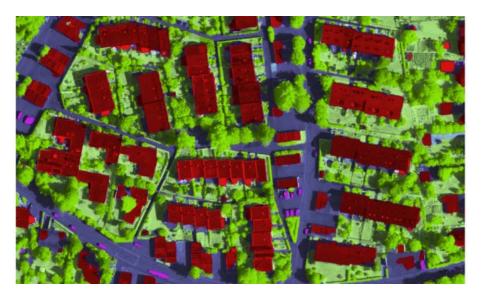




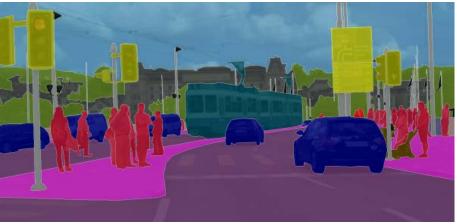


m

- Image Segmentation (Semantic/Instance Segmentation)
  - Each pixel has a label [Long-2015], [Mask-RCNN-2017]





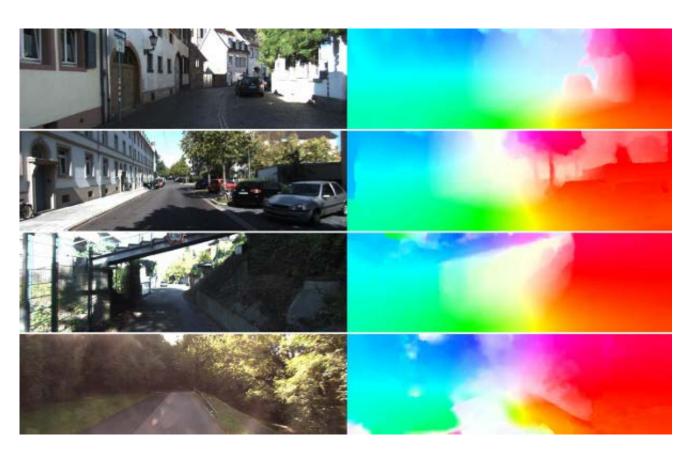


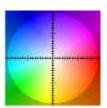
Instance segmentation

Semantic segmentation



- Motion
  - Tracking [Neoral-Serych-2024]
  - Optical Flow [Neoral-2018]
    - Predict pixel level displacements between consecutive frames







- Stereo (depth from two images)
- Depth from a single (monocular) image [Godard-2017]





- Image based novel view synthesis
  - Given: a set of sparse images => arbitrary view (smooth camera path)
  - NeRF (Neural Radiance Field for View Synthesis), [Mildenhall-2020]

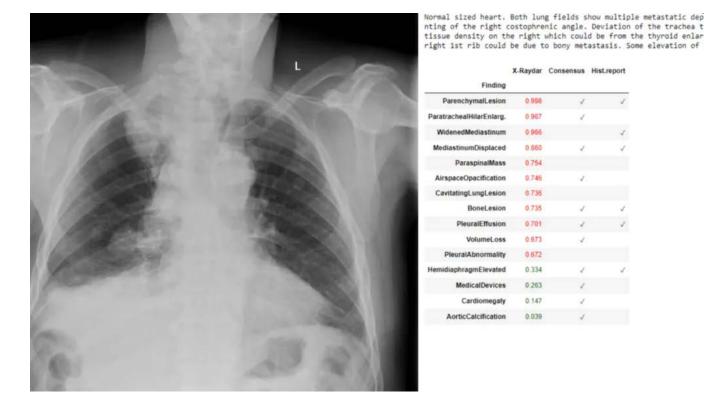






<u>video</u> <u>video</u>] <u>video</u>

- Medical Imaging Computer Aided Diagnosis
  - X-ray, mammography, etc.



- Al as good as doctors at checking X-rays study (<u>BBC news</u>)
- Commercial tools, Startups





- Faces
  - Recognition / Verification
  - Gender/Age
  - Landmarks, pose
  - Expression, emotions

#### ...already in commerce













#### 17

# Examples of Deep learning in Computer Vision



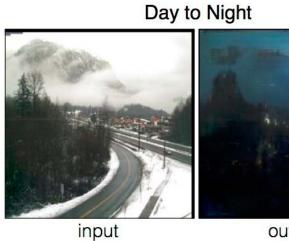
Lip reading [Chung-2017]



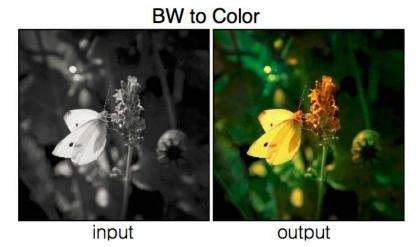




Image-to-Image translation [Isola-2017]







Deblurring, Super-resolution [<u>Šubrtová-2018</u>]



16x16

256x256 (predicted)

256x256 (ground-truth)





- Generative models
  - Generating photo-realistic samples from image distributions
  - Variational Autoencoders, GANs [<u>Nvidia-GAN</u>]







(Images synthetized by a random sampling)





- Generative models (cont.)
  - Large text2image models, 2022+ (DALL-E2, Imagen, Midjourney,
     <u>Stable Diffusion</u> open source, model available)



panda mad scientist mixing sparkling chemicals, artstation



a propaganda poster depicting a cat dressed as french emperor napoleon holding a piece of cheese





- Real image manipulation / editing
  - Instruct Pix2Pix (textual image manipulation) [Brooks-2023]



Hairstyle Transfer [<u>Šubrtová-2021</u>]



<u>video</u>]

- Action/Activity recognition
- Neural Style Transfer
- Image Captioning/Visual **Question Answering**
- and many more...





GPT-4

a brown dog wearing glasses while sitting at a desk



User

What is unusual about this image?



[GPT-4]

Source: https://www.barnorama.com/wp-content/uploads/2016/12/03-Confusing-Pictures.jpg

The unusual thing about this image is that a man is ironing clothes on an ironing board attached to the roof of a moving taxi.

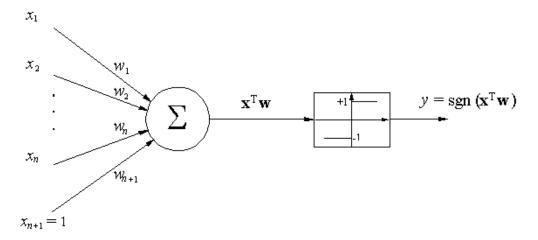
22

#### History: (Artificial) Neural Networks

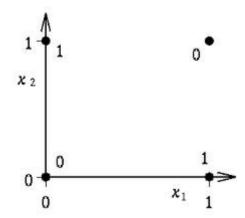


n

- Neural networks are here for almost 70 years
  - Rosenblatt-1956 (perceptron)



Minsky-1969 (xor issue, => skepticism)



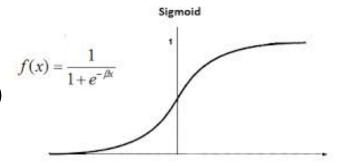
# History: Neural Networks





#### Rumelhart and McClelland – 1986:

- Multi-layer perceptron,
- Back-propagation (supervised training)
  - Differentiable activation function
  - Stochastic gradient descent

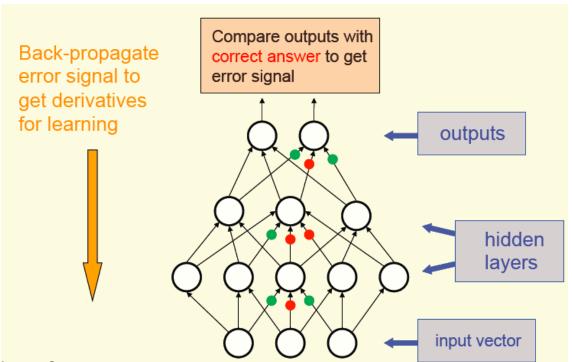


#### Empirical risk

$$Q(w) = \sum_{i=1}^{n} Q_i(w),$$

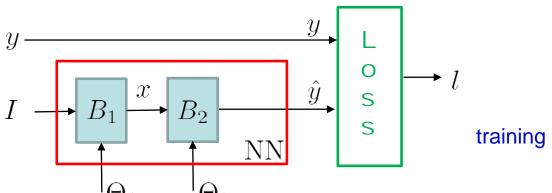
Update weights:

$$w := w - \alpha \nabla Q_i(w).$$



What happens if a network is deep? (it has many layers)





$$NN(I;\Theta) = \hat{y}$$
$$\Theta = (\Theta_1, \Theta_2)$$

training set:  $\{(I_1,y_1),\ldots,(I_n,y_n)\}$ 

total loss: 
$$L(\Theta) = \sum_{i} l_i(y_i, \hat{y}_i) = \sum_{i} l_i\Big(y_i, \mathrm{NN}(I_i; \Theta)\Big)$$

Training:  $\min_{\Theta} L(\Theta)$ 

$$\Theta^0 = \text{init}$$

GD: 
$$\Theta^{t+1} = \Theta^t - \alpha \nabla L(\Theta)$$

SGD: 
$$\Theta^{t+1} = \Theta^t - \alpha \frac{1}{m} \sum_{i=1}^{m} \nabla l_i(y_i, \text{NN}(I_i, \Theta))$$

learning rate  $\alpha$ batch size m << n

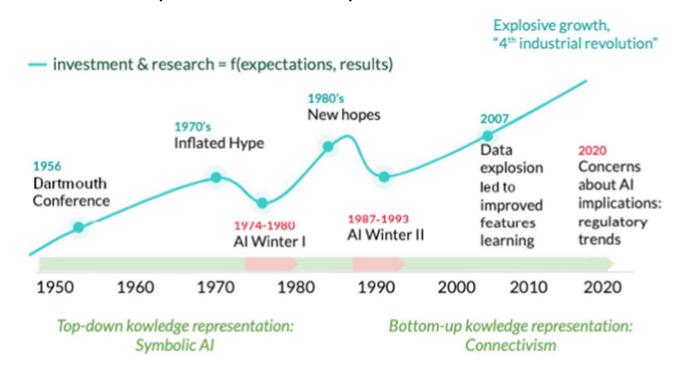
$$\nabla l_i(\Theta)^T = \frac{dl_i(\Theta)}{d\Theta} = \frac{dl_i(\Theta)}{d\hat{y}_i} \frac{dNN(I_i;\Theta)}{d\Theta}$$
$$\frac{dNN(I;\Theta)}{d\Theta_2} = \frac{dB_2(x;\Theta_2)}{d\Theta_2}$$
$$\frac{dNN(I;\Theta)}{d\Theta_1} = \frac{dB_2(x;\Theta_2)}{dx} \frac{dB_1(I;\Theta_1)}{\Theta_1}$$

$$\hat{y}_i = \text{NN}(I_i, \Theta) = B_2(\underline{B_1(I; \Theta_1)}; \Theta_2)$$

#### What was wrong with back propagation?

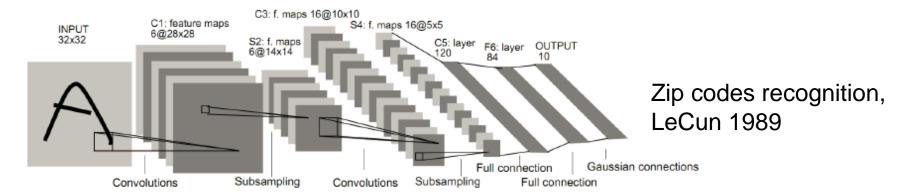


- Local optimization only (needs a good initialization, or re-initialization)
- Prone to over-fitting
  - too many parameters to estimate
  - too few labeled examples
- Computationally intensive
- => Skepticism: A deep network often performed worse than a shallow one



#### Why does it work now?





- However nowadays:
  - Large collections of labeled data available
    - ImageNet (14M images, 21k classes, hand-labeled)
  - Reducing the number of parameters by weight sharing
    - Convolutional layers [<u>LeCun-1989</u>]
  - Novel tricks to prevent overfitting of deep nets
  - Fast enough computers (parallel hardware, GPU)

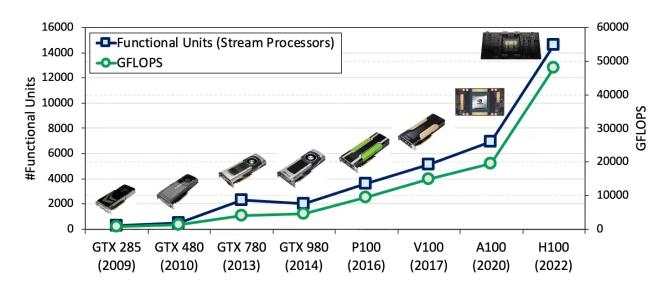
=> Optimism: It works!

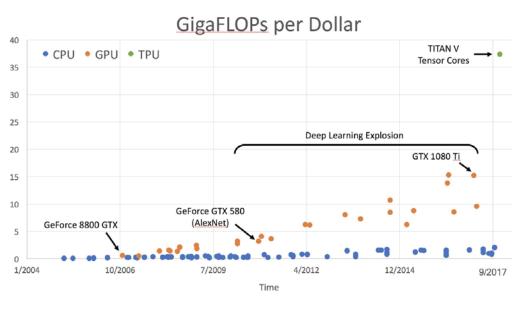
#### Computational power

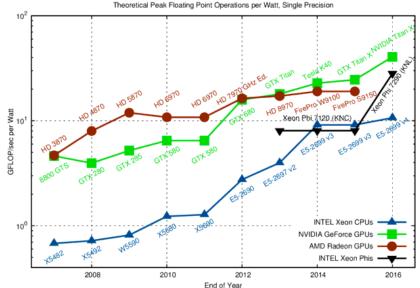








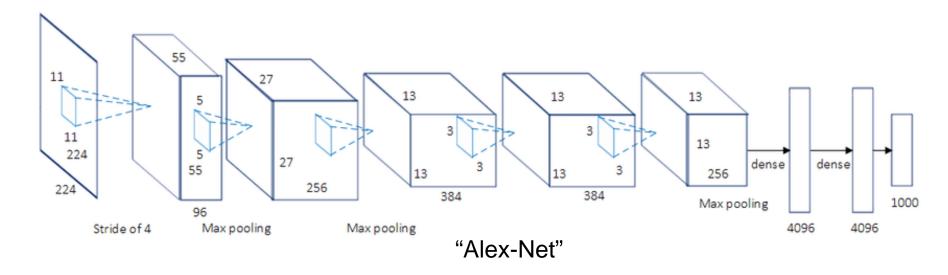




#### Deep convolutional neural networks



- An example for Large Scale Classification Problem:
  - Krizhevsky, Sutskever, Hinton: <u>ImageNet classification with deep</u> convolutional neural networks. NIPS, 2012.
    - Recognizes 1000 categories from ImageNet
    - Outperforms state-of-the-art by significant margin (ILSVRC 2012)

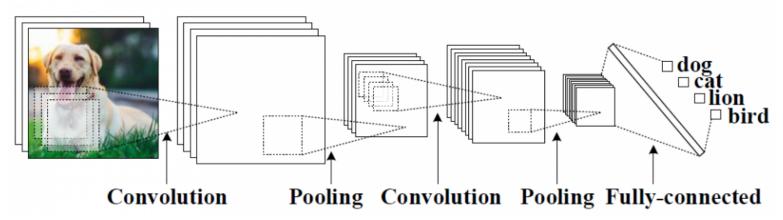


- 5 convolutional layers, 3 fully connected layers
- 60M parameters, trained on 1.2M images (~1000 examples for each category)
- Cross-Entropy loss (softmax log-loss)

#### Deep CNNs – basic building blogs



- A computational graph (chain/directed acyclic graph) connecting layers
  - Each layer has: Forward pass, Backward pass
  - The graph is end-to-end differentiable

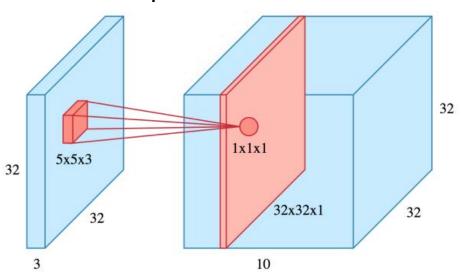


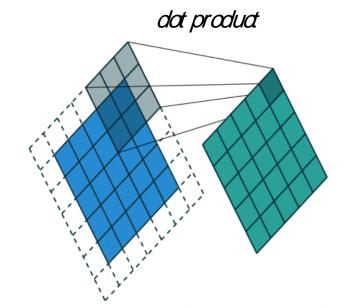
- 1. Input Layer
- 2. Intermediate Layers
  - 1. Convolutions
  - 2. Max-pooling
  - 3. Activations
- Output Layer
- 4. Loss function over the output layer for training

#### Convolutional layer

- **8**
- m p

- Input: tensor (WxHxD)
  - "image" of size WxH with D channels
- Output: tensor (W'xH'xD')
- A bank of D' filters of size (KxKxD) is convolved with the input to produce the output tensor
  - Zero Padding (P), extends the input by zeros
  - Stride (S), mask shifts by more than 1 pixel
  - KxKxDxD' parameters to be learned

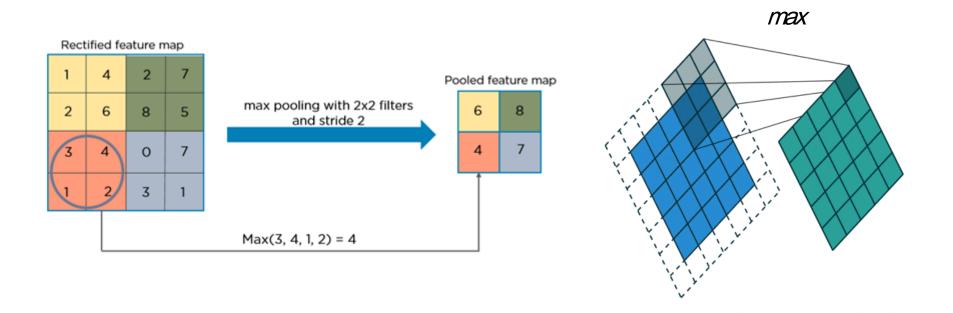




#### Max-pooling layer



- Same inputs  $(W \times H \times D)$  and outputs  $(W' \times H' \times D)$  as convolutional layer
- Same parameters: Mask Size (K), Padding (P), Stride (S)
- Same sliding window as in convolution, but instead of the dot product, pick maximum
- Non-linear operation
- No parameters to be learned

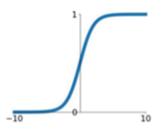


3/

- Non-linearity, applied to every singe cell of the tensor
- Input tensor and output tensor of the same size

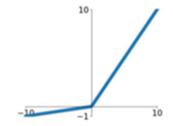
# **Sigmoid**

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$



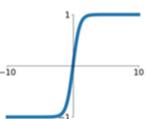
# Leaky ReLU

 $\max(0.1x, x)$ 



#### tanh

tanh(x)

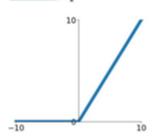


#### **Maxout**

 $\max(w_1^T x + b_1, w_2^T x + b_2)$ 

#### ReLU

 $\max(0, x)$ 



#### **ELU**

$$\begin{cases} x & x \ge 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

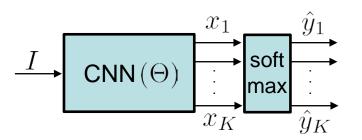
- ReLU is the simplest (used in the AlexNet, good baseline)
- Saturating non-linearity (sigmoid, tanh) causes "vanishing" gradient

#### **Multiclass Classification loss**



n

Cross-Entropy loss (softmax log loss)



$$\hat{y}_i = \frac{e^{x_i}}{\sum_{j=1}^K e^{x_j}}$$

Softmax output as discrete PDF over classes

$$\hat{y}_i \ge 0, \sum_{i=1}^K \hat{y}_i = 1$$

Ground-truth classes "one-hot encoding"

$$y_i = \begin{cases} 1 & i \text{ is the truth class} \\ 0 & \text{otherwise} \end{cases}$$

$$L(\mathbf{y}, \hat{\mathbf{y}}(\Theta)) = -\sum_{i=1}^{K} y_i \log(\hat{y}_i) = -\log(\hat{y}_{i^*})$$

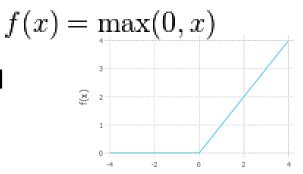
 $i^*$  is index of the truth class

#### Deep convolutional neural networks





- Additional tricks: "Devil is in the details"
  - Rectified linear units instead of standard sigmoid
    - => Mitigate vanishing gradient problem
  - Convolutional layers followed by max-pooling
    - Local maxima selection in overlapping windows (subsampling)
    - => dimensionality reduction, shift insensitivity
  - Dropout
    - 50% of hidden units are randomly omitted during the training, but weights are shared in test time
    - Averaging results of many independent models (similar idea as in Random forests)
    - => Probably very significant to reduce overfitting
  - Data augmentation
    - Images are artificially shifted and mirrored (10 times more images)
    - => transformation invariance, reduce overfitting

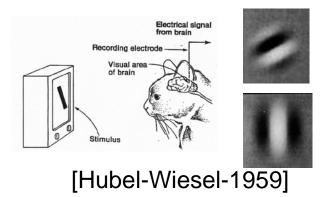


#### Deep convolutional neural networks





- Supervised training
  - The training is done by a standard back-propagation
  - enough labeled data: 1.2M labeled training images for 1k categories
  - Learned filters in the first layer
    - Resemble cells in primary visual cortex





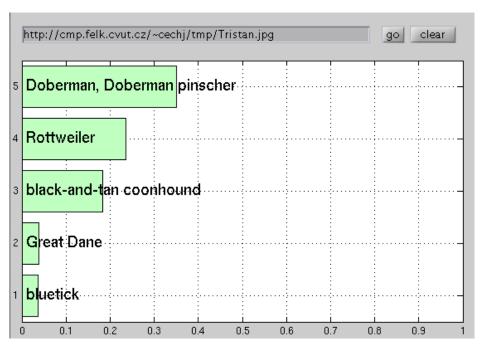
Learned first-layer filters

- Training time:
  - 5 days on NVIDIA GTX 580, 3GB memory (Krizhevsky, today faster)
  - 90 cycles through the training set
- Test time (forward step) on GPU
  - Implementation by Yangqing Jia, <a href="http://caffe.berkeleyvision.org/">http://caffe.berkeleyvision.org/</a>
  - 5 ms/image in a batch mode

# Early experiments 1: Category recognition



- Implementation by Yangqing Jia, 2013, <a href="http://caffe.berkeleyvision.org/">http://caffe.berkeleyvision.org/</a>
  - network pre-trained for 1000 categories provided
- Which categories are pre-trained?
  - 1000 "most popular" (probably mostly populated)
  - Typically very fine categories (dog breeds, plants, vehicles...)
  - Category "person" (or derived) is missing
  - Recognition accuracy subjectively surprisingly good...

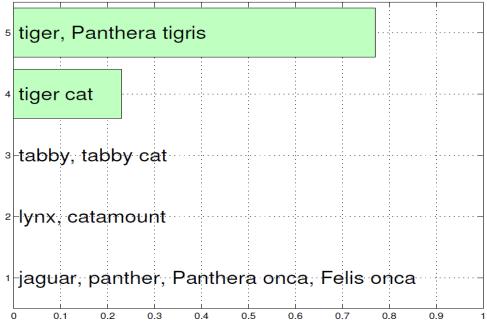




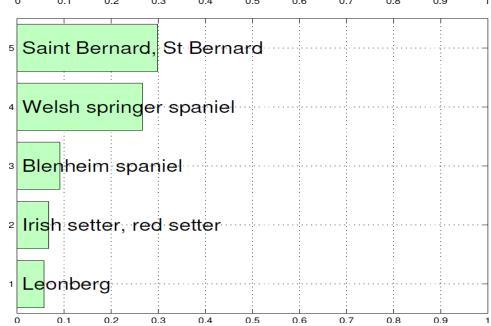
# It is not a texture only...











# Early experiments 2: Category retrieval

•

- 50k randomly selected images from Profimedia dataset
- Category: Restaurant (results out of 50k-random-Profiset)



# Early experiments 2: Category retrieval



m p

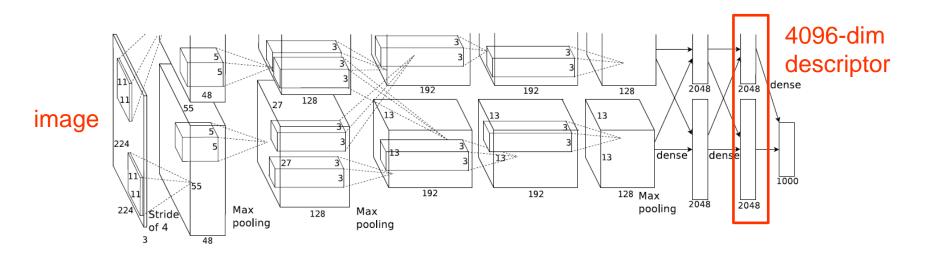
Category: stethoscope (results out of 50k-random-Profiset)



47



- Indications in the literature that the last hidden layer carry semantics
  - Last hidden layer (4096-dim vector), final layer category responses (1000-dim vector)
  - New (unseen) categories can be learned by training (a linear) classifier on top of the last hidden layer
    - Oquab, Bottou, Laptev, Sivic, CVPR, 2014
  - Responses of the last hidden layer can be used as a compact global image descriptor
    - Semantically similar images should have small Euclidean distance
    - Novak, Cech, Zezula, ICSSA, 2015



# Early experiments 3: Similarity search



- Qualitative comparison: (20 most similar images to a query image)
  - MUFIN annotation (web demo), <a href="http://mufin.fi.muni.cz/annotation/">http://mufin.fi.muni.cz/annotation/</a>, [Zezula et al., Similarity Search: The Metric Space Approach. 2005.]
    - Nearest neighbour search in 20M images of Profimedia
    - Standard global image statistics (e.g. color histograms, gradient histograms, etc.)
  - Caffe NN (last hidden layer response + Euclidean distance),
    - Nearest neighbour search in 50k images of Profimedia
    - 400 times smaller dataset!









































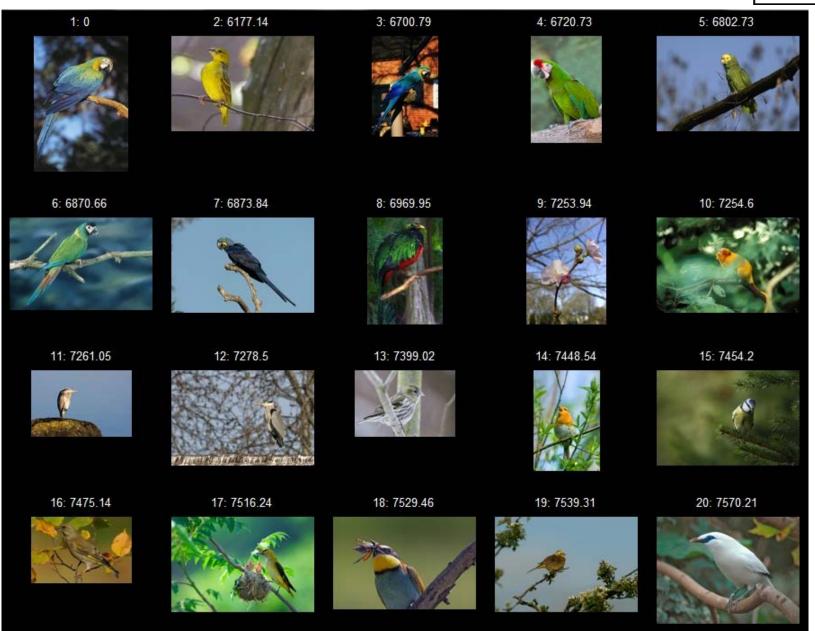
MUFIN results

# Early experiments 3: Similarity search



m p

50















































# Caffe NN results









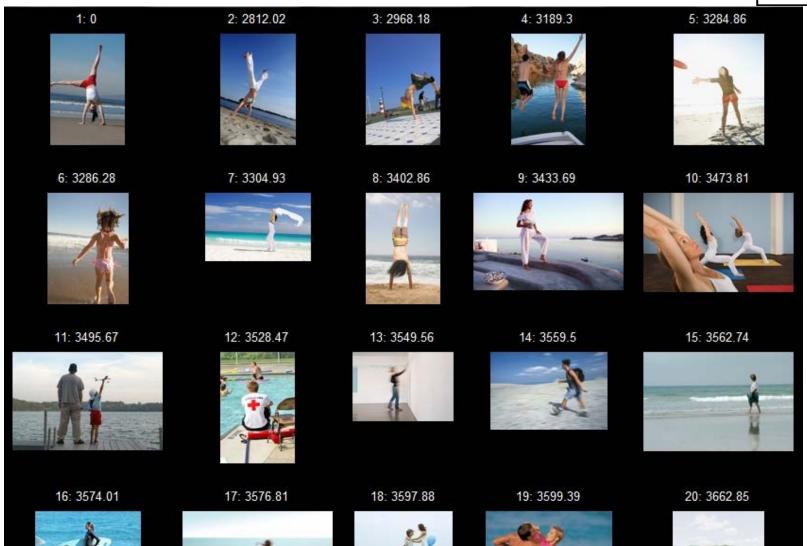
#### MUFIN results



# Early experiments 3: Similarity search



m p



54

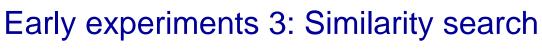
#### **MUFIN** results



# Caffe NN results



















6: 3477.87





























20: 3979.94

#### MUFIN results

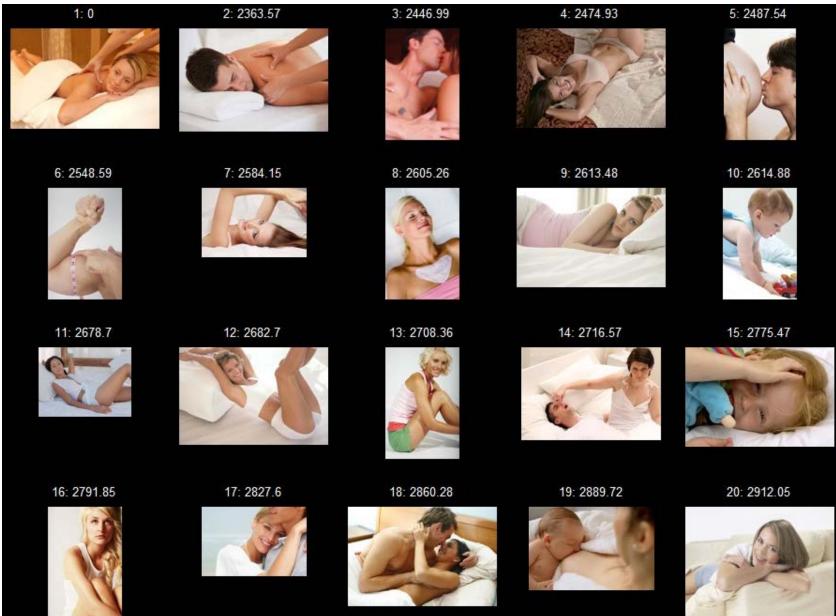


# Early experiments 3: Similarity search



m p

58



#### **Novel tricks**





#### Network initialization

- Mishkin, Matas. <u>All you need is a good init</u>. ICLR 2016
- Weights initialization: zero mean, unit variance, orthogonality
- Batch normalization
  - Iosse, Szegedy. <u>Batch Normalization: Accelerating Deep Network</u>
     <u>Training by Reducing Internal Covariate Shift</u>. NIPS 2015

- Zero mean and unit variance weights are "supported" during training

to avoid vanishing gradient

- ⇒ Small sensitivity to learning rate setting (can be higher, faster training
  - 10 times fewer epochs needed)
- ⇒ Regularizer (dropout can be excluded/smaller) (better optimum found)

```
Input: Values of x over a mini-batch: \mathcal{B} = \{x_{1...m}\};

Parameters to be learned: \gamma, \beta

Output: \{y_i = \mathrm{BN}_{\gamma,\beta}(x_i)\}

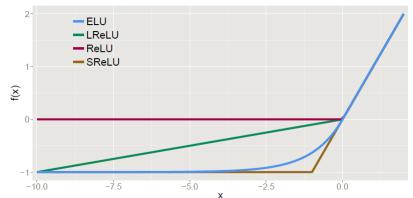
\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \qquad \text{// mini-batch mean}
\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \qquad \text{// mini-batch variance}
\widehat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \qquad \text{// normalize}
y_i \leftarrow \gamma \widehat{x}_i + \beta \equiv \mathrm{BN}_{\gamma,\beta}(x_i) \qquad \text{// scale and shift}
```

Algorithm 1: Batch Normalizing Transform, applied to activation x over a mini-batch.



Exponential Linear Units (ELU) [Clevert et al., ICLR 2016]

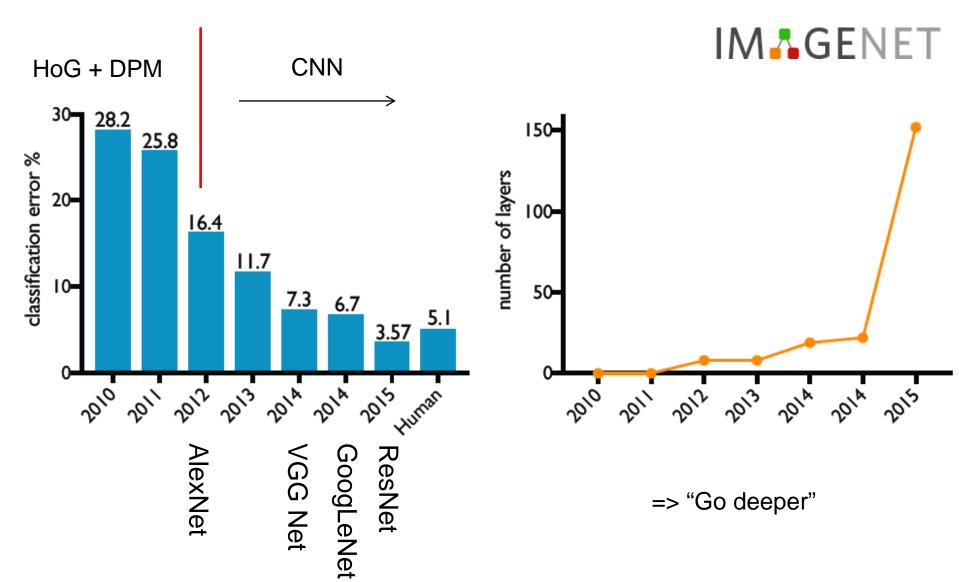
$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha (\exp(x) - 1) & \text{if } x \le 0 \end{cases}$$



- Self normalizing properties, batch normalization unnecessary
- Faster training reported
- ADAM optimizer [Kingma and Ba, ICLR 2015]
  - = (ADAptive Moments)
  - Often improves over SGD (with momentum),
  - Low sensitivity on learning rate setting

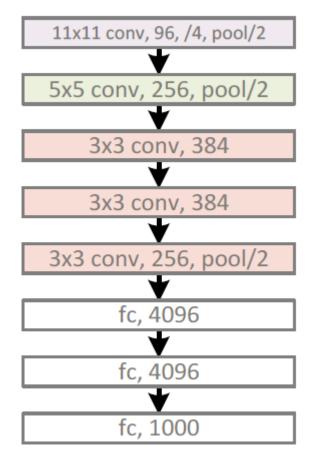
61

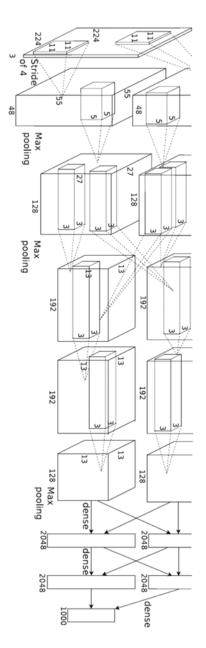
ImageNet Large Scale Visual Recognition Challenge (ILSVRC)



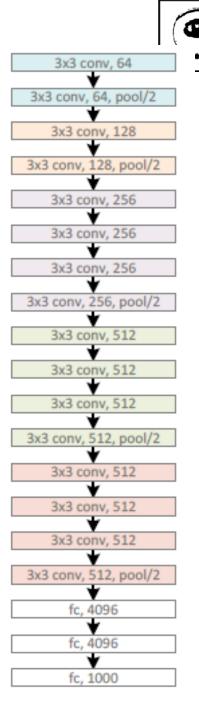
#### AlexNet

Krishevsky et al., NIPS 2012



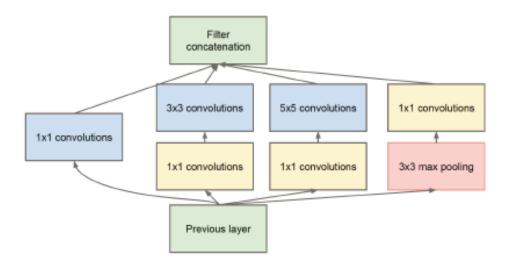


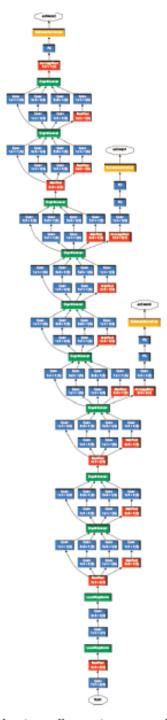
- VGG Net: VGG-16, VGG-19
  - Simonyan and Zisserman, ICLR 2015
  - Deeper than AlexNet
  - Smaller filters (3x3 convolutions), more layers
    - => Same effective receptive field, but more "non-linearity"



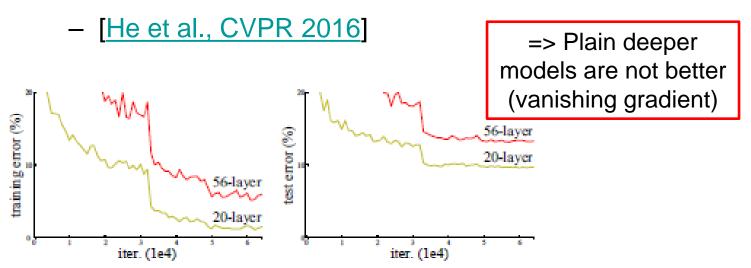
63

- GoogLeNet
  - Szegedy et al., CVPR 2015
  - 22 layers, No Fully-Connected layers
  - Accurate, much less parameters
  - "Inception" module (Net in net)

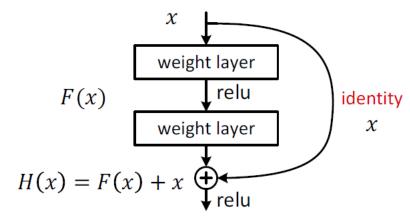


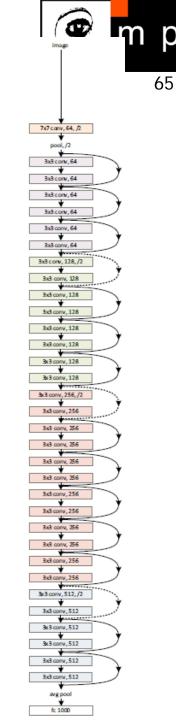






Residual modules, 152 layers





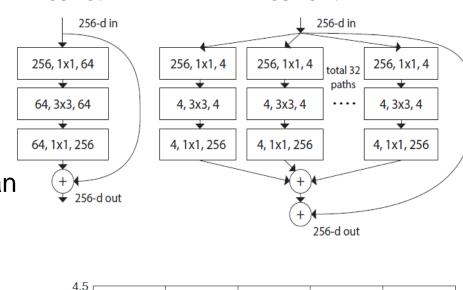
**ResNet** 



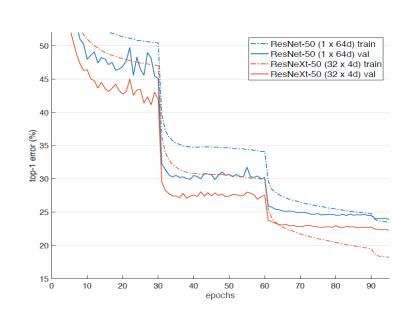
66

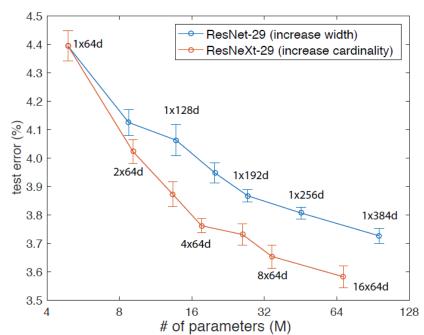
ResNeXt

- Xie-CVPR-2017
- Improvement of ResNet
- Cardinality
  - number of branches in a block
- "Increasing cardinality, better than going wider or deeper"



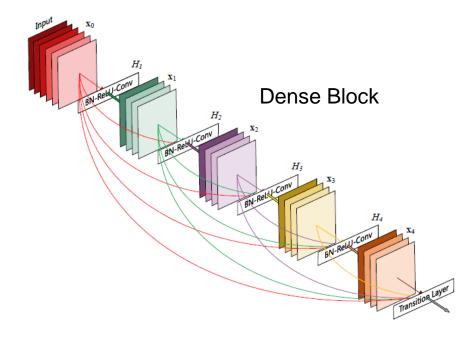
ResNeXt

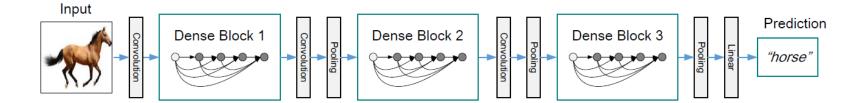






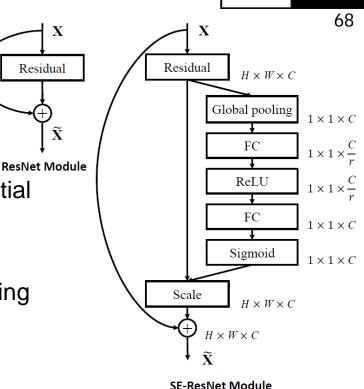
- Huang-CVPR-2017
- Densifying Skip connections
- Chain of several "dense blocks"
- Argument: Features are reused
- Higher accuracy with fewer parameters over ResNet reported
- Best paper award @ CVPR

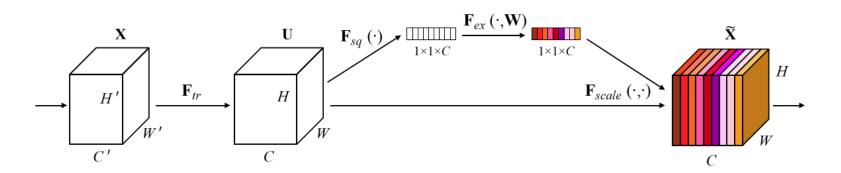




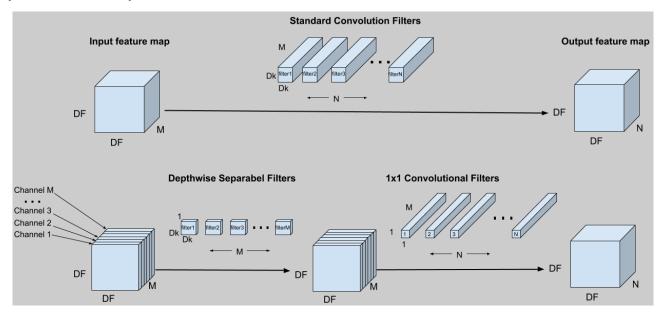
- **(B)**
- m k

- Squeeze-and-Excitation Network (SE-Net)
  - <u>Hu-CVPR-2018</u>, <u>Hu-TPAMI-2019</u>]
  - Chain of SE-blocks
  - Squeeze:
    - Channel descriptor by aggregating over spatial dimension
  - Excitation
    - Small bottleneck fully connected net producing scale of each channel
  - Capture channel interdependences
  - Winner of ILSVRC 2017 (Top-5 err 2.25%)
  - Negligible extra computational cost

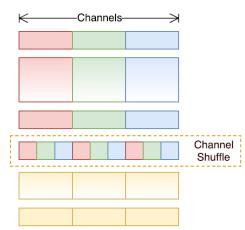




- Computationally efficient architectures
  - MobileNet [Howard-2017, Google Inc.]
    - depth wise separable convolutions



- ShuffleNet [Zhang-CVPR-2018, Face++]
  - Comparable accuracy with AlexNet, 13x speed up





Taken from Natural Language Processing



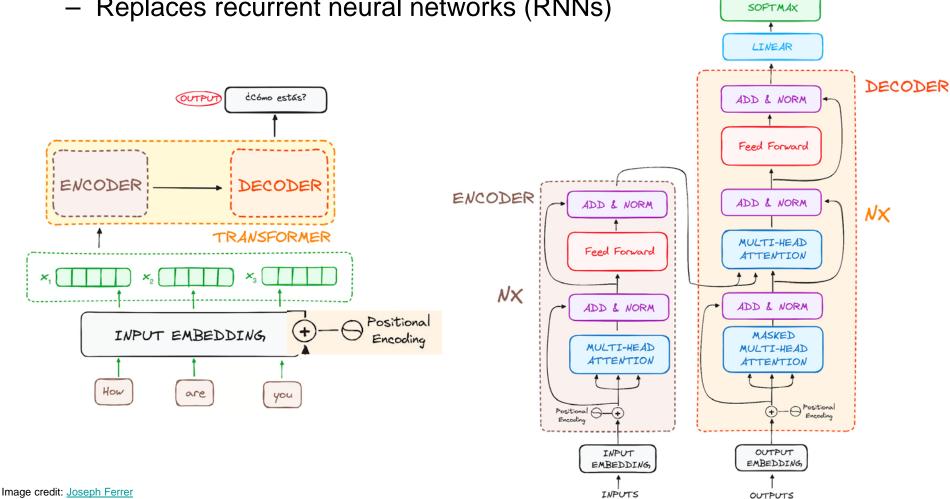


PROBABILITIES

"Attention is all you need" [Wasvani-2017]

Originally for machine translation (seqence2sequence)

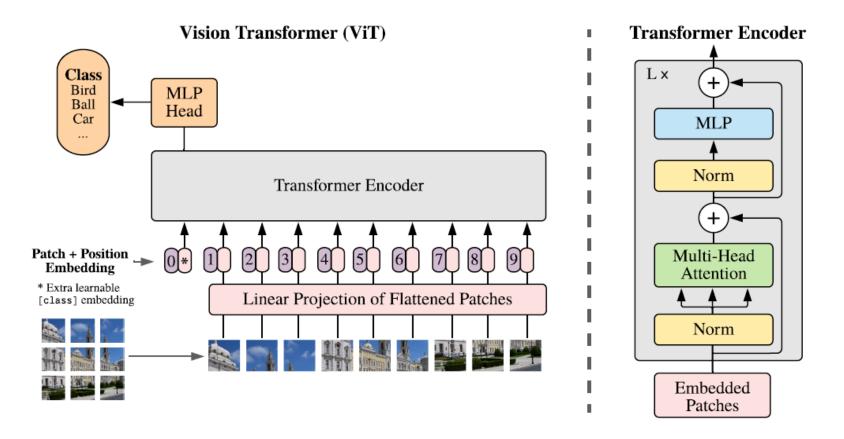
Replaces recurrent neural networks (RNNs)



70



- Vision Transformers [<u>Dosovitskiy-2021</u>]
  - No Convolutions, Encoder only transformer, Parallel processing
  - Image is cut into fixed-size patches and the sequence of vectorized patches (tokens/words) is fed into the transformer



Outperforms ResNET on ImageNet, but needs 100M image pretraining

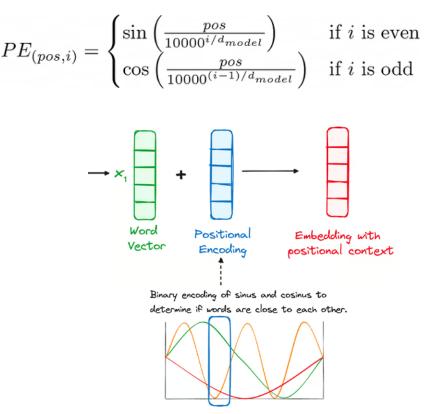


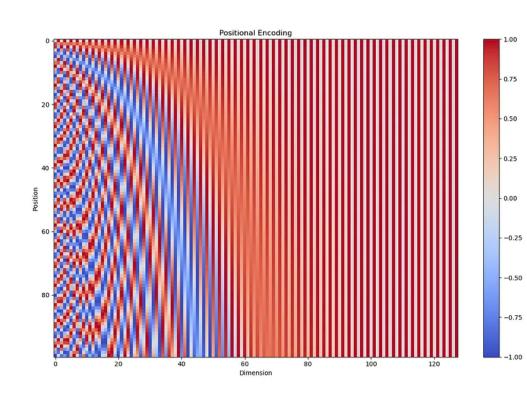


- (Vision) Transformer
  - Input tokens treated equally, but order of the sequence is important
    - "Dog bites man" vs. "Man bites dog"

## ⇒ Positional Encoding

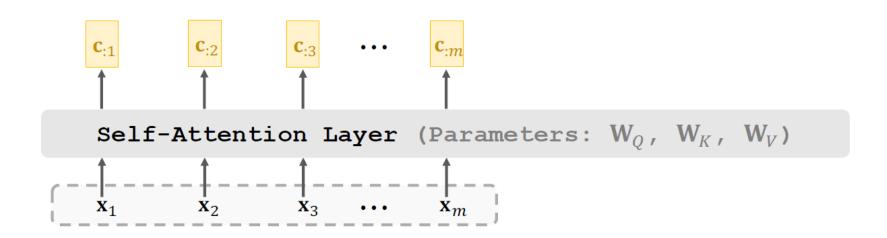
- Encodes absolute position of each token
- Using smooth functions (sin, cos) each token's position gives a vector







- (Vision) Transformer
  - Main idea: Self-Attention Mechanism
    - Inputs (vectors x<sub>1</sub>, ..., x<sub>m</sub>)
    - Parameters (matrices W<sub>Q</sub>, W<sub>K</sub>, W<sub>V</sub>)

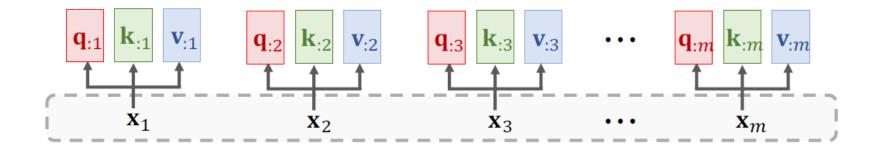


```
Query: \mathbf{q}_{:i} = \mathbf{W}_Q \mathbf{x}_i, Key: \mathbf{k}_{:i} = \mathbf{W}_K \mathbf{x}_i, Value: \mathbf{v}_{:i} = \mathbf{W}_V \mathbf{x}_i. \mathbf{c}_{:j} = \mathbf{V} \cdot \text{Softmax}(\mathbf{K}^T \mathbf{q}_{:j}).
```



- (Vision) Transformer
  - Main idea: Self-Attention Mechanism
    - Inputs (vectors x<sub>1</sub>, ..., x<sub>m</sub>)
    - Parameters (matrices W<sub>Q</sub>, W<sub>K</sub>, W<sub>V</sub>)

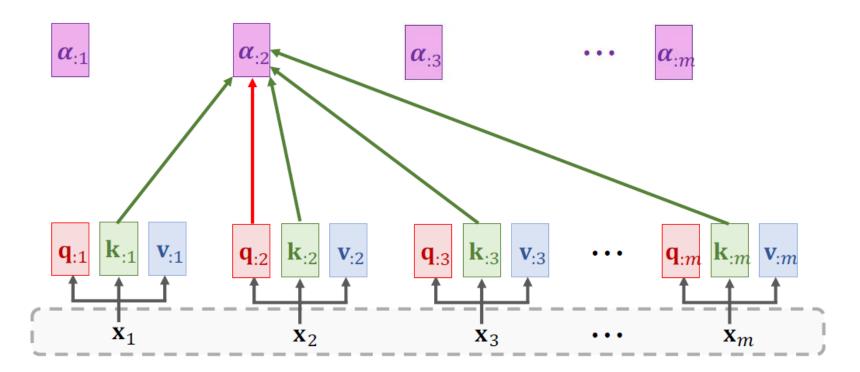
Query:  $\mathbf{q}_{:i} = \mathbf{W}_Q \mathbf{x}_i$ , Key:  $\mathbf{k}_{:i} = \mathbf{W}_K \mathbf{x}_i$ , Value:  $\mathbf{v}_{:i} = \mathbf{W}_V \mathbf{x}_i$ .





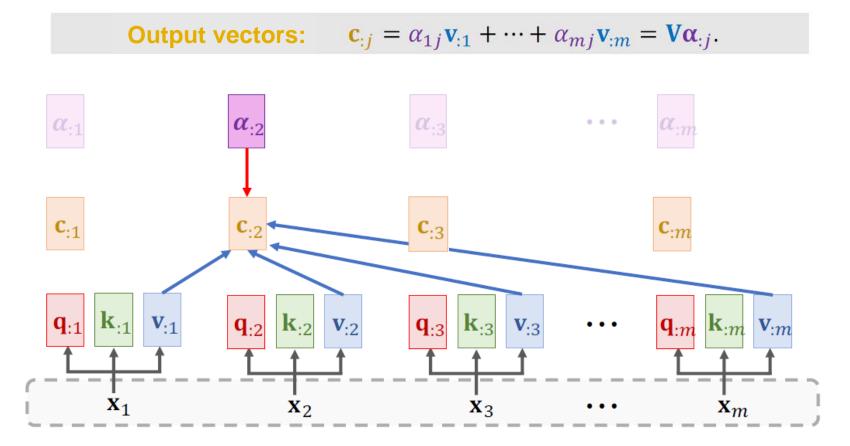
- (Vision) Transformer
  - Main idea: Self-Attention Mechanism
    - Inputs (vectors x<sub>1</sub>, ..., x<sub>m</sub>)
    - Parameters (matrices W<sub>Q</sub>, W<sub>K</sub>, W<sub>V</sub>)

Weights:  $\alpha_{:j} = \operatorname{Softmax}(\mathbf{K}^T \mathbf{q}_{:j}) \in \mathbb{R}^m$ .



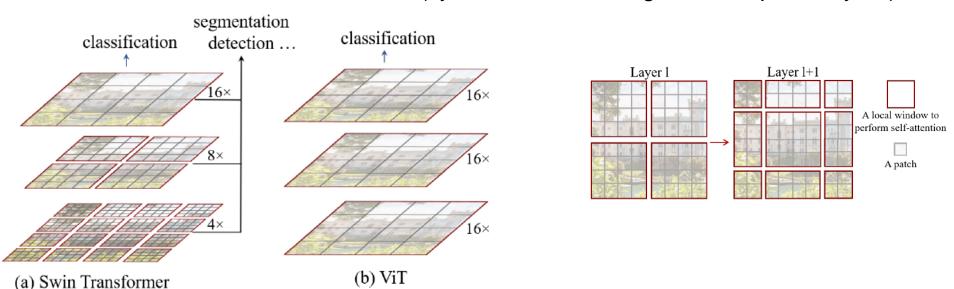


- (Vision) Transformer
  - Main idea: Self-Attention Mechanism
    - Inputs (vectors x<sub>1</sub>, ..., x<sub>m</sub>)
    - Parameters (matrices W<sub>Q</sub>, W<sub>K</sub>, W<sub>V</sub>)





- SWIN Transformer [<u>Liu-2021</u>] ("Shifted Windows")
  - Improvement of ViT transformer
    - data hungry (needs large set pretraining)
    - Image tokens too large unsuitable for object detection, semantic segmentation
  - Hierarchical features
    - Self attention within windows (linear complexity w.r.t. image size)
    - Cross-window connection (cyclic window shifting in subsequent layers)



 State-of-the-art general purpose backbone (recognition, detection, segmentation, ....)

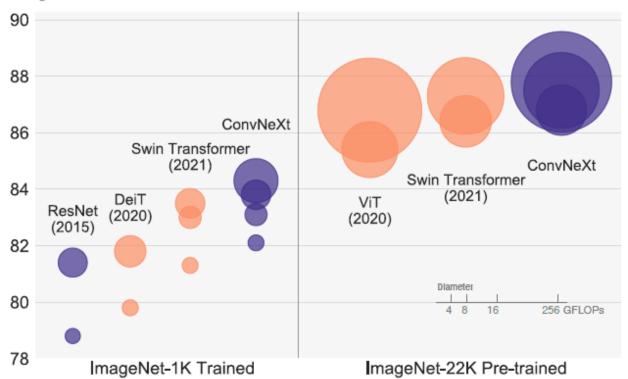
#### **DNN Architectures – ConvNext**

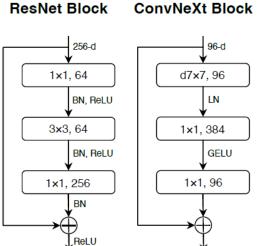


ConvNeXt [<u>Liu-2022</u>]

- Pure Convolutional Neural Network (again)
- Similar to ResNet, but tweaked
- Larger kernel size, BatchNorm -> <u>LayerNorm</u>
- ReLU -> GeLU (smoother)

#### ImageNet-1K Acc.

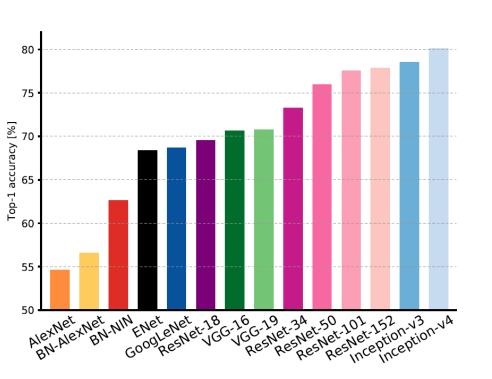


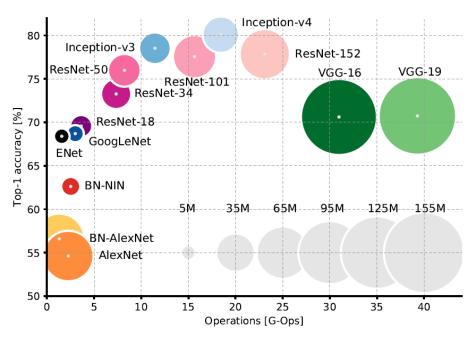


78

# **CNN** models (comparison)



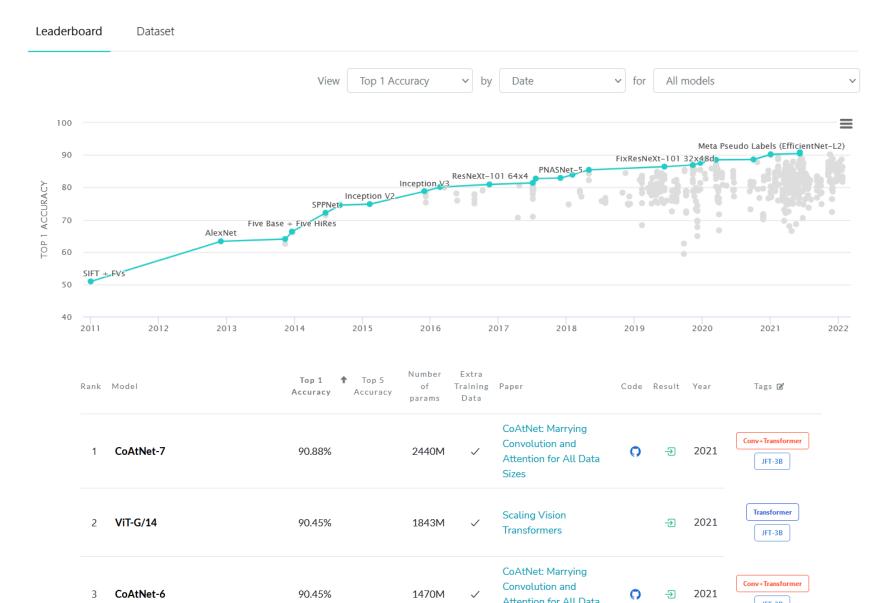




[Canziani et al., <u>An Analysis of Deep Neural Network Models for Practical Applications</u>, 2017. arXiv:1605.07678v4]

# CNN models (comparison)

ImageNet <u>leaderboard</u> (Top-1 accuracy)



## Face interpretation problems



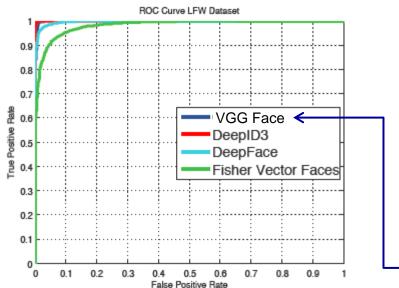
- Face recognition, face verification
  - Architecture similar to AlexNet deep CNN (softmax at the last layer)

[<u>Taigman-ECVV-2014</u>] DeepFace: Closing the Gap to Human-Level Performance in Face Verification (authors from Facebook)

[Parkhi-BMVC-2015] Deep Face recognition (authors from Oxford Uni)

- 2.6M images of 2.6k celebrities, trained net available

[Deng-CVPR-2019] ArcFace (faces mapped onto a unit sphere)



No.	Method	# Training Images	# Networks	Accuracy
1	Fisher Vector Faces	-	-	93.10
2	DeepFace (Facebook)	4 M	3	97.35
3	DeepFace Fusion (Facebook)	500 M	5	98.37
4	DeepID-2,3	Full	200	99.47
5	FaceNet (Google)	200 M	1	98.87
6	FaceNet+ Alignment (Google)	200 M	1	99.63
7	(VGG Face)	2.6 M	1	98.78

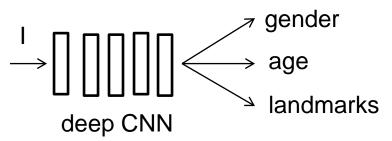
 Face represented by penultimate layer response, similarity search, large scale indexing

# Face interpretation problems





- Facial landmarks, Age / Gender estimation
  - Multitask network
    - Shared representation
    - Combination of both classification and regression problems





# Age estimation – How good the network is?



### Our survey

~20 human subjects, ~100 images of 2 datasets

#### MORPH dataset

True: 22, MAE: 18.8



True: 36, MAE: 17.8



True: 33, MAE: 16.3



True: 22, MAE: 16.1



True: 25, MAE: 16.0



**IMDB** dataset

True: 25, MAE: 0.5



True: 66, MAE: 1.0



True: 29, MAE: 1.0



True: 19, MAE: 1.0



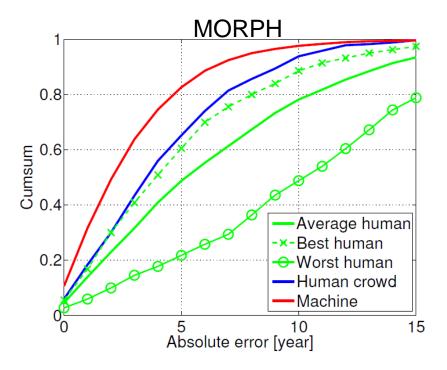
True: 43, MAE: 1.0



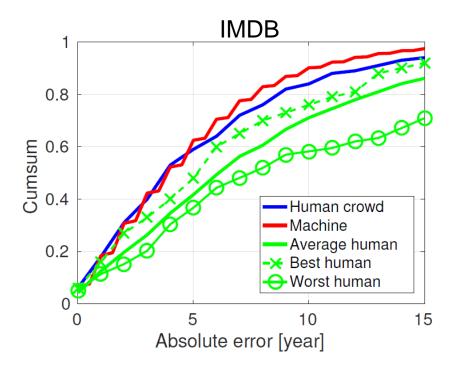
# Age estimation – How good the network is?



Better than average human...



		MAE	CS5	MaxAL
Average human	:	6.8	48.6	24.1
Human crowd	:	4.7	65.1	19.0
Machine		3.2	82.6	26.0



		MAŁ	CS5	MaxAL
Average human	:	8.2	41.7	31.5
Human crowd	:	5.7	59.0	21.0
Machine	:	5.1	62.5	42.7

- Franc-Cech-IVC-2018
- Network runs real-time on CPU

# Predicting Decision Uncertainty from Faces



- [Jahoda, Vobecky, Cech, Matas. <u>Detecting Decision Ambiguity from</u> Facial Images. In Face and Gestures, 2018]
- Can we train a classifier to detect uncertainty?



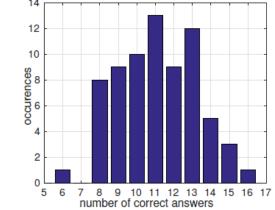




Training set: 1,628 sequences

Test set: 90 sequences

- => YES, we can...
  - CNN 25% error rate, while human volunteers 45%

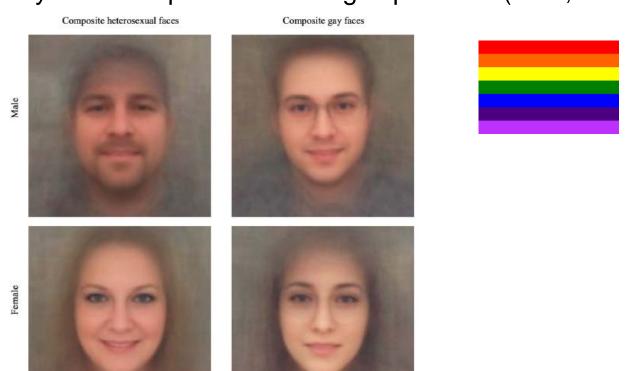


## Sexual Orientation from Face Images



- [Wang and Kosinki. Deep neural networks are more accurate than humans at detecting sexual orientation from facial images. Journal of Personality and Social Psychology, 2018]
- Better accuracy than human in (gay vs. heterosexual)

  - 81% accuracy (for men), average human accuracy (61%)
  - 71% accuracy (for women)
- average human accuracy (54%)
  - Accuracy further improved if 5 images provided (91%, 83%)



# General recipe to use deep neural networks



- m p
- Recipe to use deep neural network to "solve any problem" (G. Hinton 2013) 87
- Have a deep net
  - If you do not have enough labeled data, pre-train it by unlabeled data; otherwise do not bother with pre-initialization
  - Use rectified linear units instead of standard neurons (sigmoid)
  - Use dropout to regularize it (you can have many more parameters than training data)
  - If there is a spatial structure in your data, use convolutional layers
- Novel:
  - Use Batch Normalization [loffe-Szegedy-NIPS-2015]
  - ReLU => ELU, GELU
  - Adaptive Optimizers (ADAM)
  - Various architectures (AlexNet, VGG, GoogLeNet, ResNet, ResNeXt, DenseNet, SE-Net, MobileNet, ShuffleNet, Transformers, Swin, ConvNext)
- Experience:
  - Data matters (the more data the better), transfer learning, data augmentation

#### Conclusions





- DNNs efficiently learns the abstract representation
- Low computational demands for running, Training needs GPU
- Many "deep" toolboxes: Caffe (Berkeley), MatconvNet (Oxford),
   TensorFlow (Google), Theano (Montreal), PyTorch (Facebook), ...
- NNs are (again) in the "Golden Age" (or witnessing a bubble), as many practical problems seem solvable in near future
- Explosion of interest of DNN in literature, graduates get incredible offers, start-ups appear all the time

Do we understand enough what is going on? <a href="http://www.youtube.com/watch?v=LVLoc6FrLi0">http://www.youtube.com/watch?v=LVLoc6FrLi0</a>





#### **Further Resources**



- Deep Learning Textbook
  - Ian Goodfellow and Yoshua Bengio and Aaron Courville, Deep Learning, MIT Press, 2016
  - Available <u>on-line</u> for free.
- Lectures / video-lectures
  - Stanford University course on Deep Learning (cs231n)
  - MIT lectures on Introduction in Deep Learning (MIT 6.S191)
  - Various blogs and on-line journals
    - Google Al blog (<a href="https://ai.googleblog.com/">https://ai.googleblog.com/</a>)
    - OpenAl blog (<a href="https://openai.com/blog">https://openai.com/blog</a>)
    - MetaAl blog (<a href="https://ai.facebook.com/blog/">https://ai.facebook.com/blog/</a>)
    - Andrej Karpathy (blog)
    - ...