

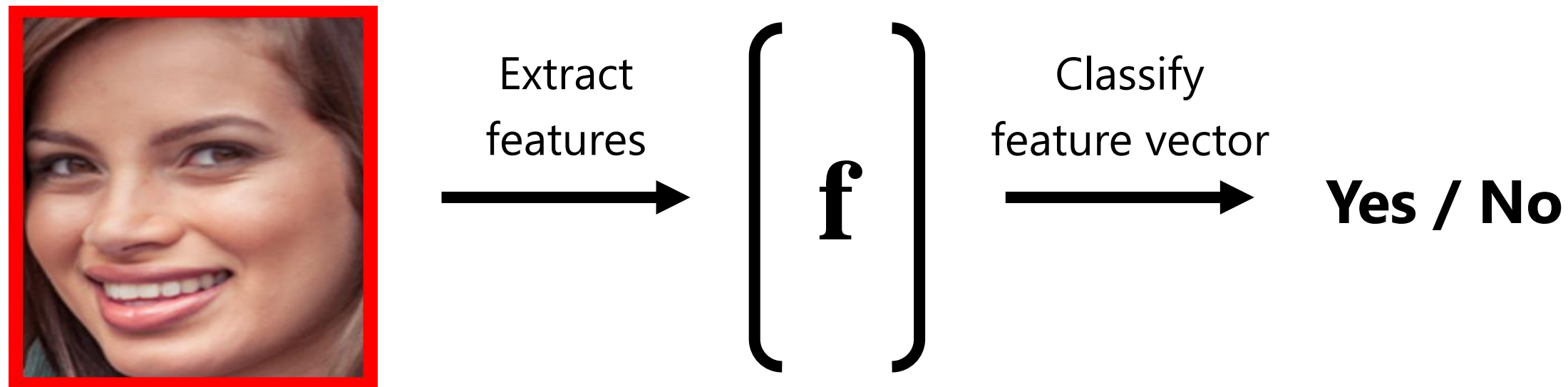
Advanced Computer Vision

Week 11

Nov. 18, 2022
Seokju Lee

Face Detection in the Past: Summary

For each window:



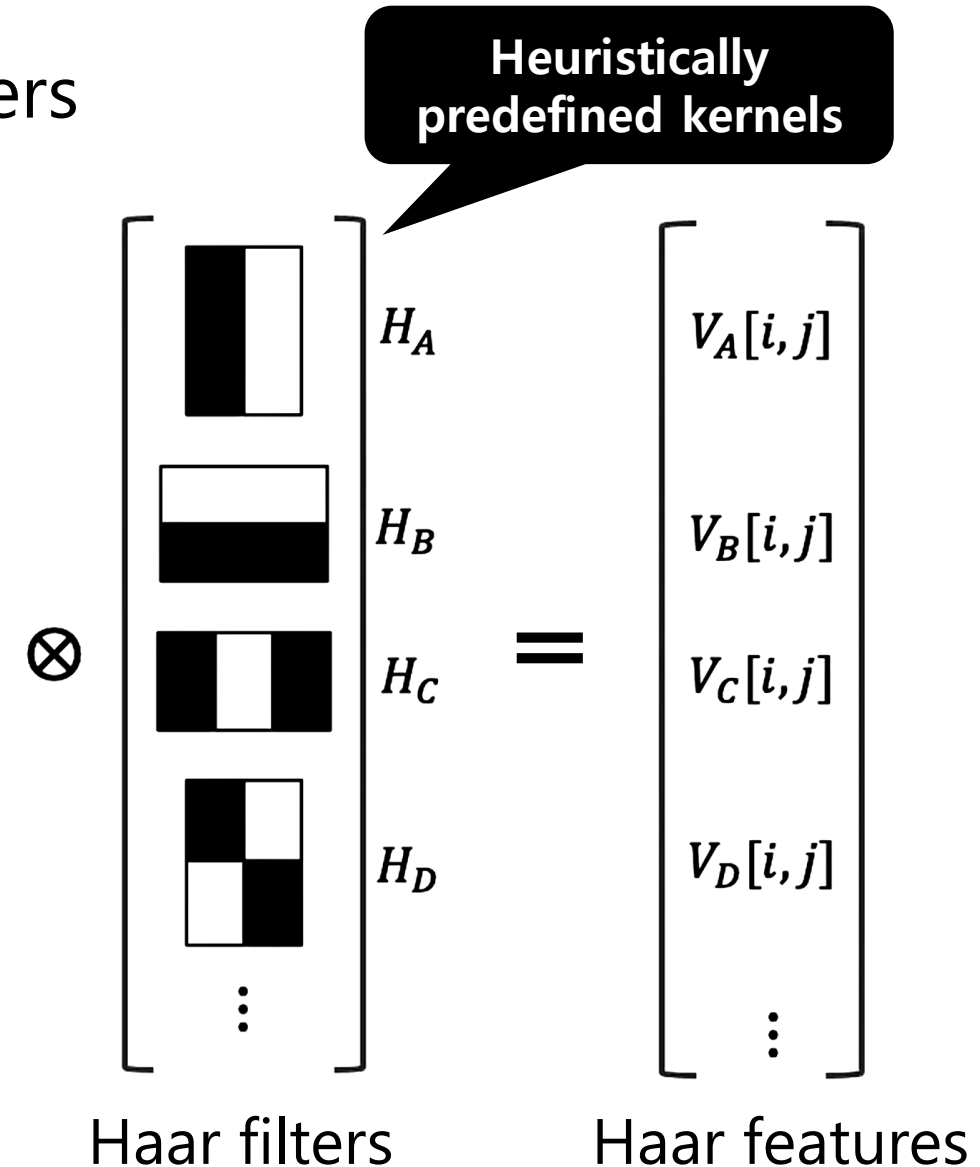
Features: Features to represent faces well.

Classifier: Design a face model and efficiently classify features as face or not.

Feature Extractor: Haar Features

Set of correlation responses to **Haar** filters

- Extremely fast by using integral images



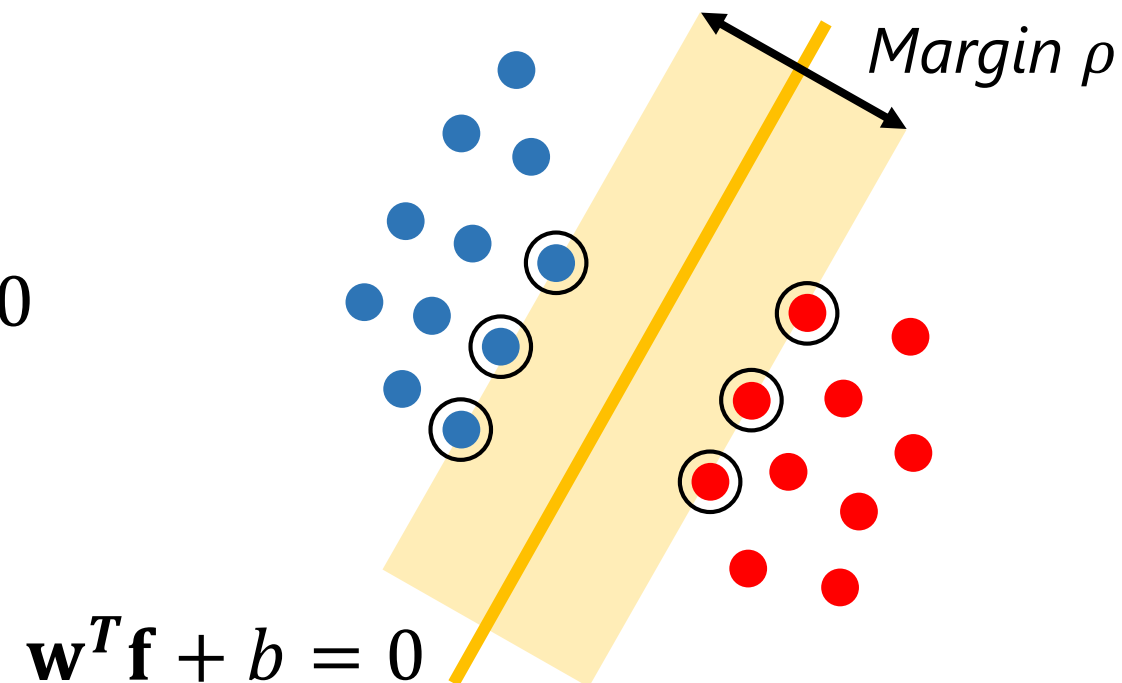
Classifier: Support Vector Machine (SVM)

Given:

- k training images $\{I_1, I_2, \dots, I_k\}$ and their Haar features $\{\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_k\}$.
- k corresponding labels $\{\lambda_1, \lambda_2, \dots, \lambda_k\}$, where $\lambda_j = +1$ if I_j is a face and $\lambda_j = -1$ if I_j is a non-face.

Find:

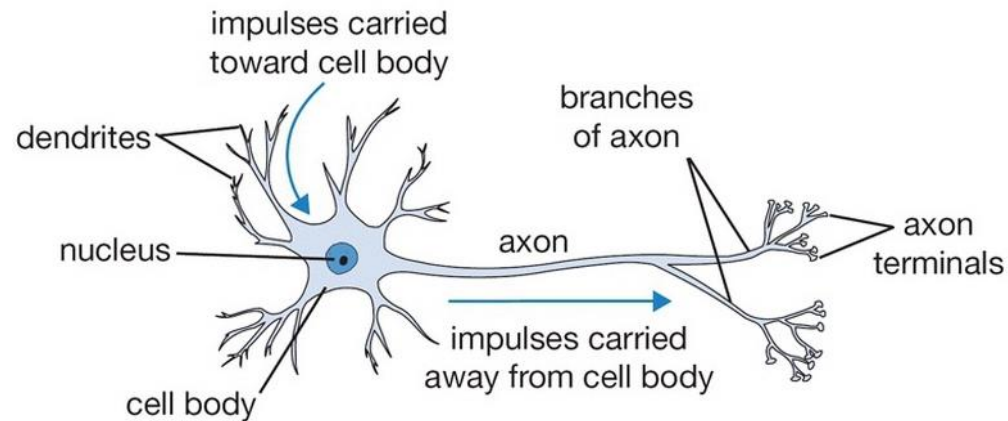
- Decision boundary $\mathbf{w}^T \mathbf{f} + b = 0$
- with maximum margin ρ



Object Recognition in the Present: Deep Learning

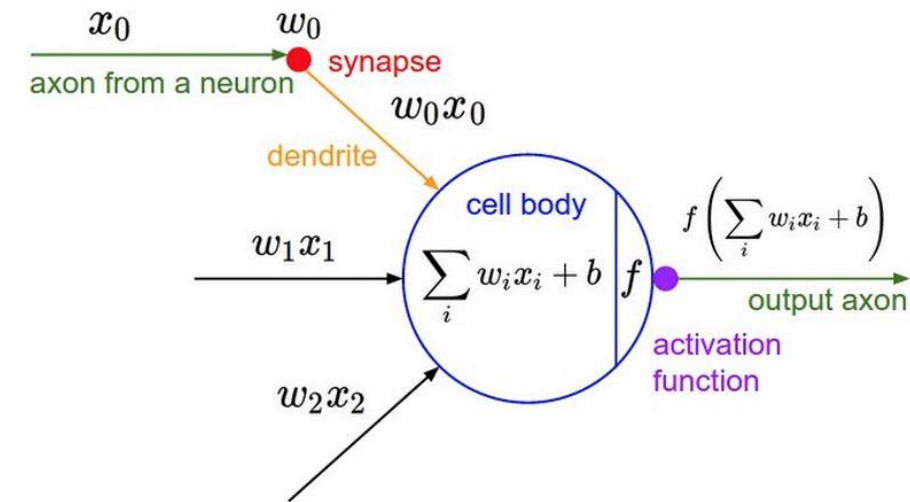
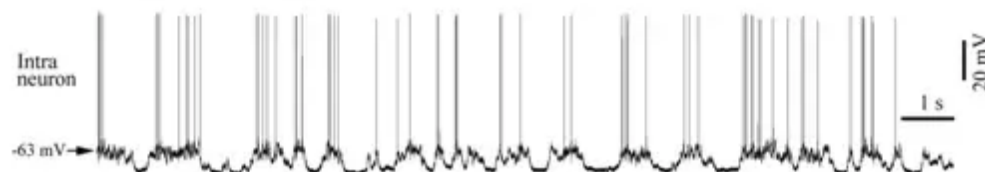
→ Towards Human-like Perception

Neuron: Biological Inspiration for Computation



(Biological) neuron: computational building block for the brain (20W)
→ Analog signals: smooth and continuous
→ Human brain: ~100-1,000 trillion synapses,

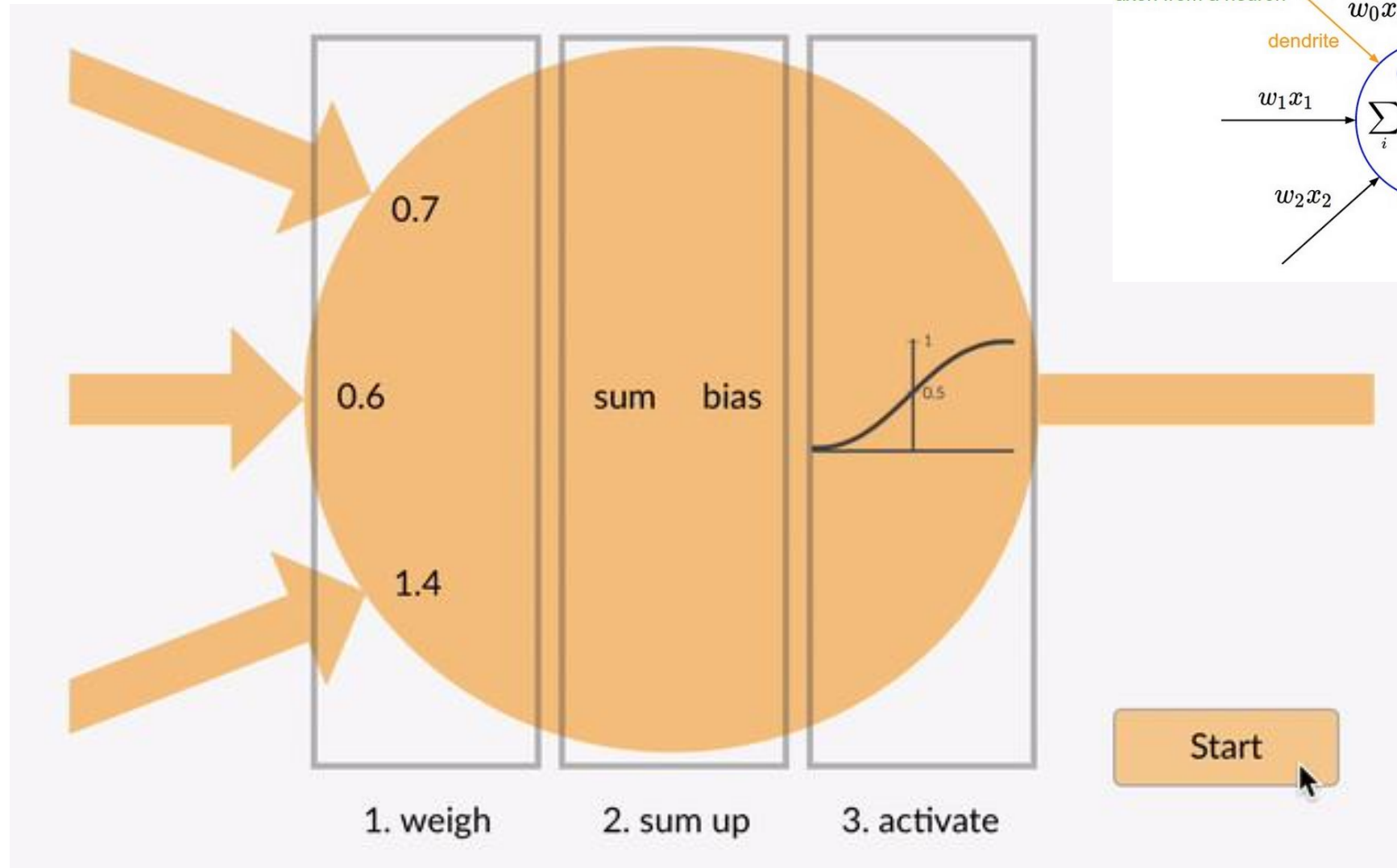
Neuron Spiking Signal



(Artificial) neuron: computational building block for the neural network (20kW)
→ Digital signals: discrete and discontinuous
→ Neural network: ~1-10 billion synapses,

Compared to neural networks, the human brain has **×10k computational power**, and consumes only **0.1% of the power**.

Perceptron: Forward Pass



Output of activation: $f(input \times weight + bias)$

Activation Functions

Why does it need activation?

→ Nonlinearity ↑, complexity ↑ to represent high dimensional information

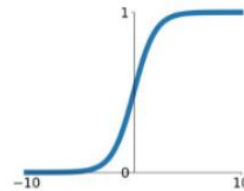
Properties of activation functions

→ Differentiable (for backpropagation)

→ Monotonic (one-to-one correspondence for input & output)

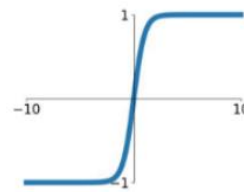
Sigmoid

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



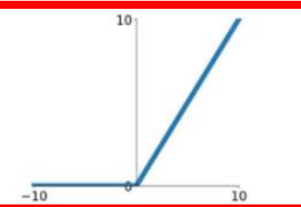
tanh

$$\tanh(x)$$



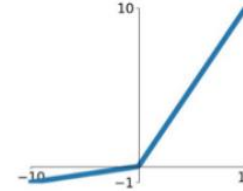
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

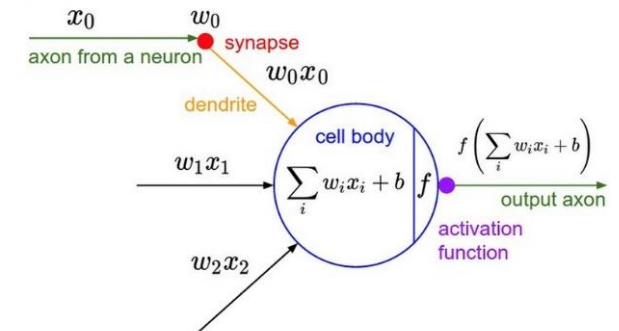
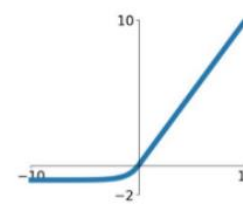


Maxout

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

ELU

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

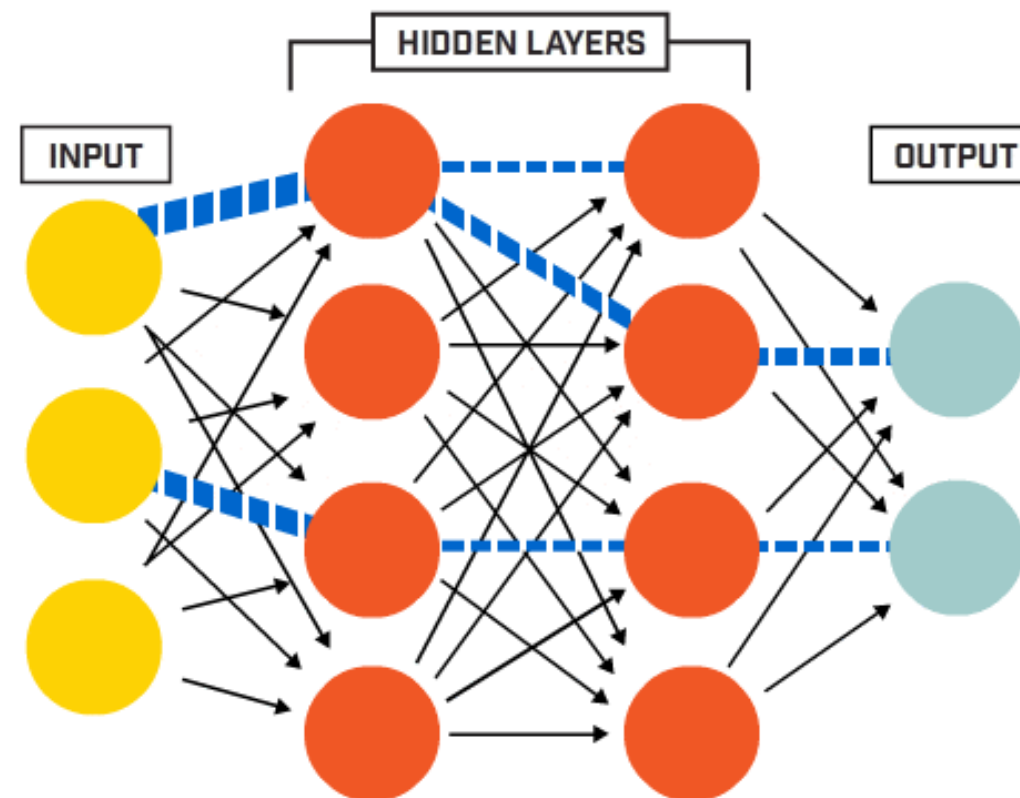


For more details, please check
*vanishing gradient problem

Multi-Layer Perceptron (MLP, or Fully-Connected)

Vectorized feature → Fully connected layers → Non-linear activation

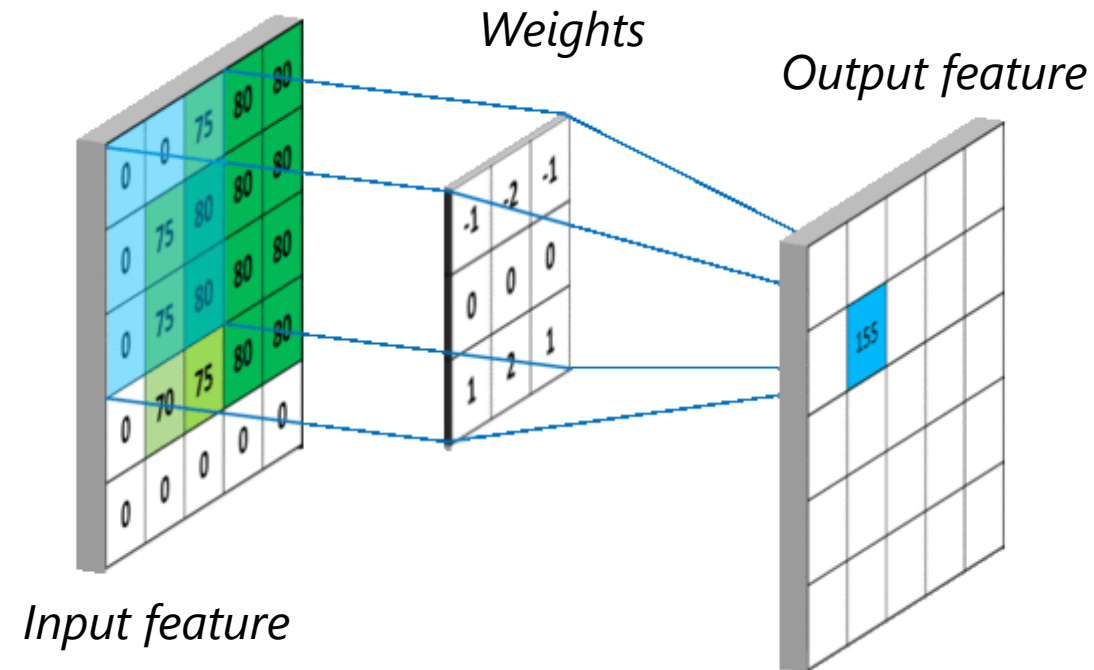
: Number of parameters ↑, hard to optimize, number of hidden layers ↓



Convolutional Neural Network (CNN)

Vectorized feature → Fully connected layers → Non-linear activation

: Number of parameters ↓, easier to optimize, number of hidden layers ↑



→ **Convolution** works on **spatial** & **local** features!

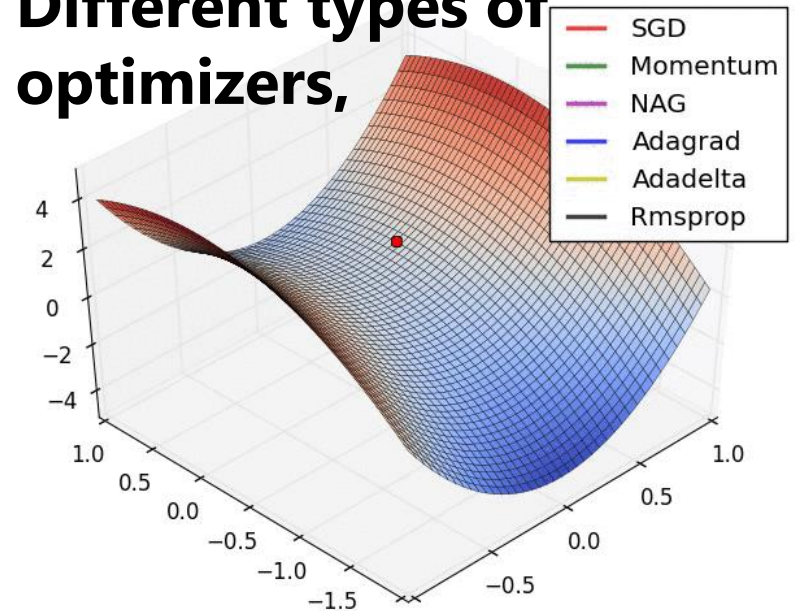
Convolutional Neural Network (CNN)



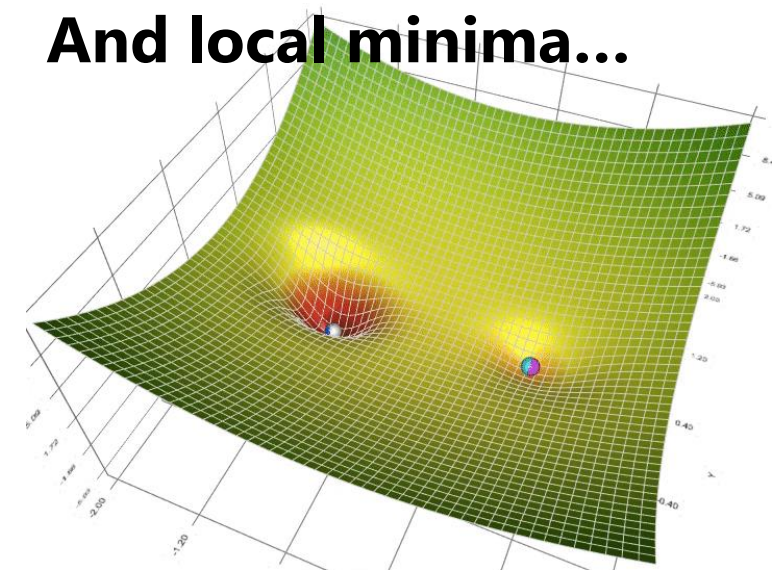
Drawbacks (including but not limited to)

- Works for specifically defined tasks (**weak AI**).
- Requires large amount of **data** (data hungry).
- Requires human **annotation** for real world data.
- **Domain** issues (virtual ↔ real, daytime ↔ night)
- **Manually** select neural architecture.
- Performance is varied over different **loss** functions.
- **Hyperparameter** tuning
 - Learning rate, loss function
 - Batch size, number of iteration, number of kernels
 - Optimizer, momentum

Different types of optimizers,



And local minima...



Useful Deep Learning Terms (including but not limited to)

- Deep learning = **Deep neural network (DNN)**
- Deep learning \subset Machine learning \subset AI
- **MLP**: Multi-layer perceptron
- **CNN**: Convolutional neural network
- GAN, GNN, RNN, LSTM, autoencoder
- Spiking neural network
- **Neural network operations**:
 - Convolution, pooling, activation function
 - Feed forward, backpropagation
 - Batch normalization, KL divergence
 - Data augmentation, regularization
- **AlexNet, Inception, VGG, ResNet**
- Others:
 - ViT, Transformer, attention, cost volume
 - PyTorch, Caffe, TensorFlow
- **Supervised** learning, **self-supervised** learning, **unsupervised** learning, **reinforcement** Learning
- Few-shot (one-shot) learning, **adversarial** learning, **domain adaptation**, meta learning, active learning, **multimodal** learning, **contrastive** learning
- **Visual learning tasks**:
 - Image classification, object detection
 - Semantic/instance/panoptic/video segmentation
 - Optical flow, depth, neural rendering
 - Stereo matching, SfM, MVS
 - Image enhancement, super resolution
 - Stylization, image-to-image, VQA, VLN

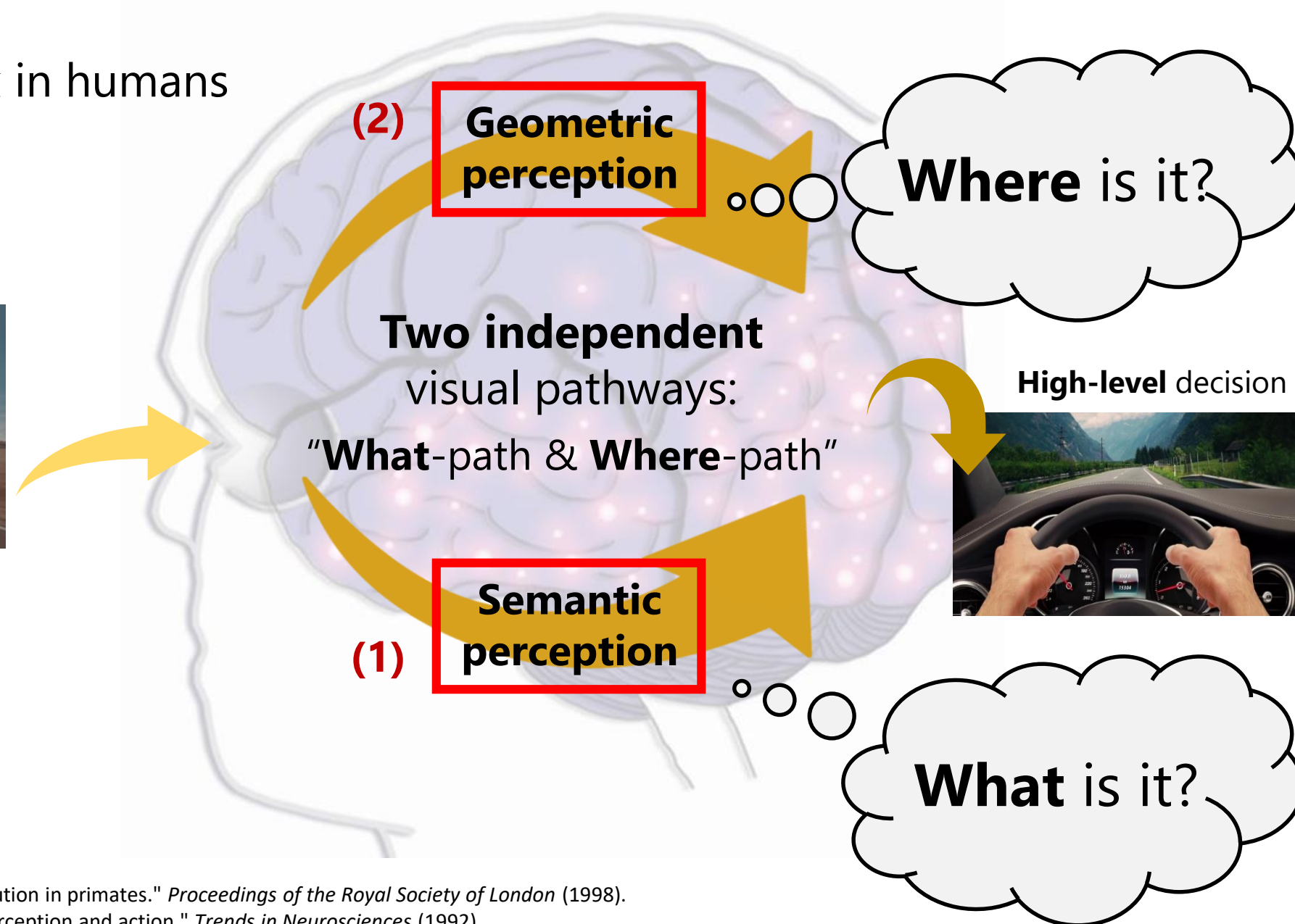
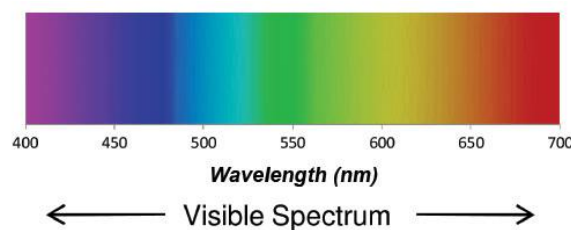
Basic Visual Perception Tasks

Human Visual System: How Do We Perceive the World?

"About **half** of neocortex in humans is devoted to **vision**." [1]



Low-level visual signal



[1] Barton, Robert A. "Visual specialization and brain evolution in primates." *Proceedings of the Royal Society of London* (1998).

[2] M. A. Goodale, et al., "Separate visual pathways for perception and action." *Trends in Neurosciences* (1992).

Visual Perception: Semantics & Geometry

“Semantic” perception

: *Meaning* of an element, *syntax*, *context* of scene, or *relationship* between objects.

Semantic computer vision tasks

- Image classification
- Object detection
- Semantic segmentation
- ...

Video understanding

ex) Video classification

“Geometric” perception

: *Distance*, *shape*, *structure*, *size*, *scale* of an element, *3D space* where we live, *relative position* between objects.

Geometric computer vision tasks

- Depth estimation
- Pose estimation
- 3D reconstruction
- ...

Motion understanding

ex) 3D motion estimation

+ “Temporal”



Computer Vision Tasks

*Slide by Kim, et al., "Video Panoptic Segmentation" (CVPR 2020)

Model Complexity ↑ Output dimension ↑

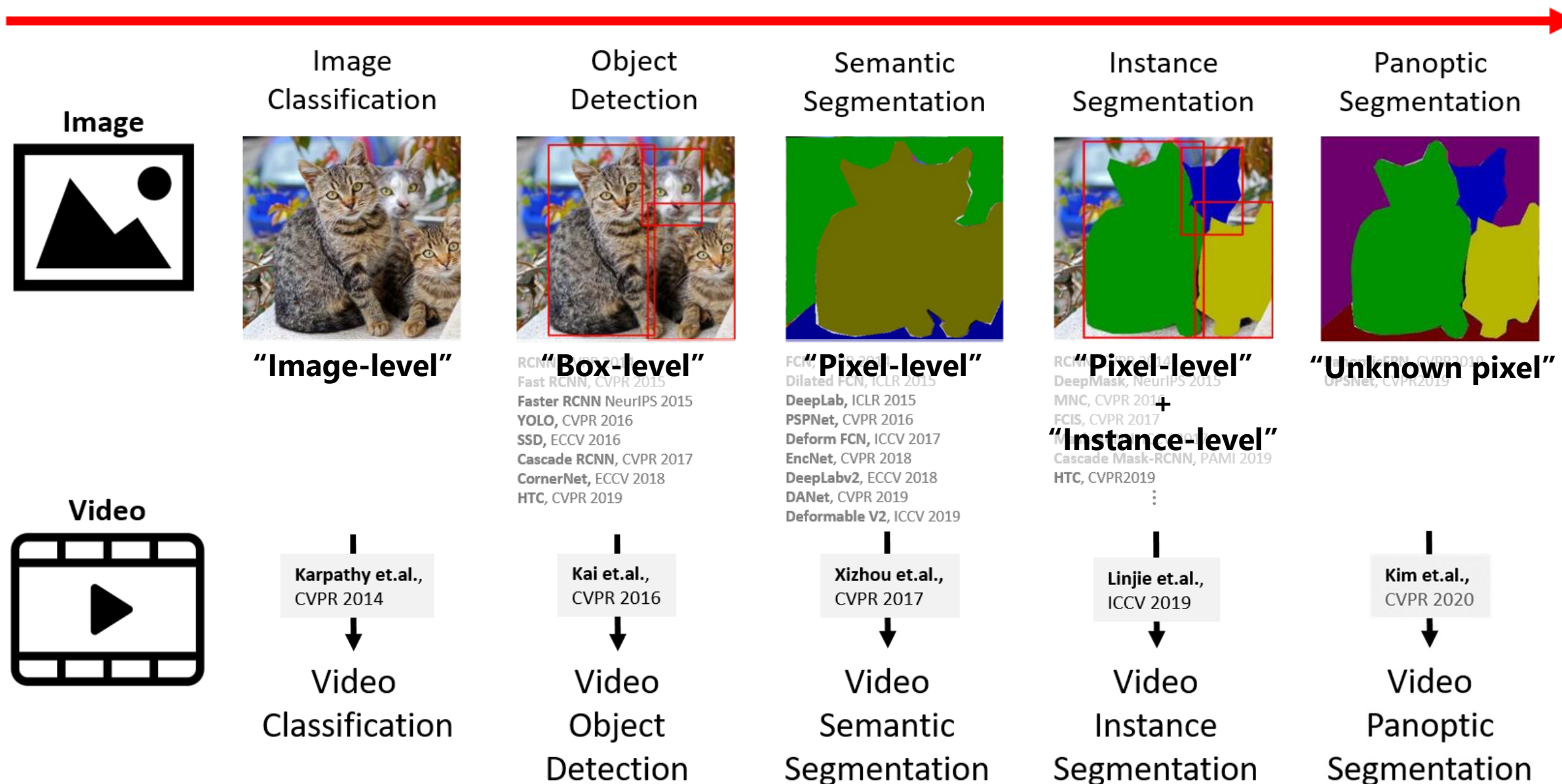


Image Classification

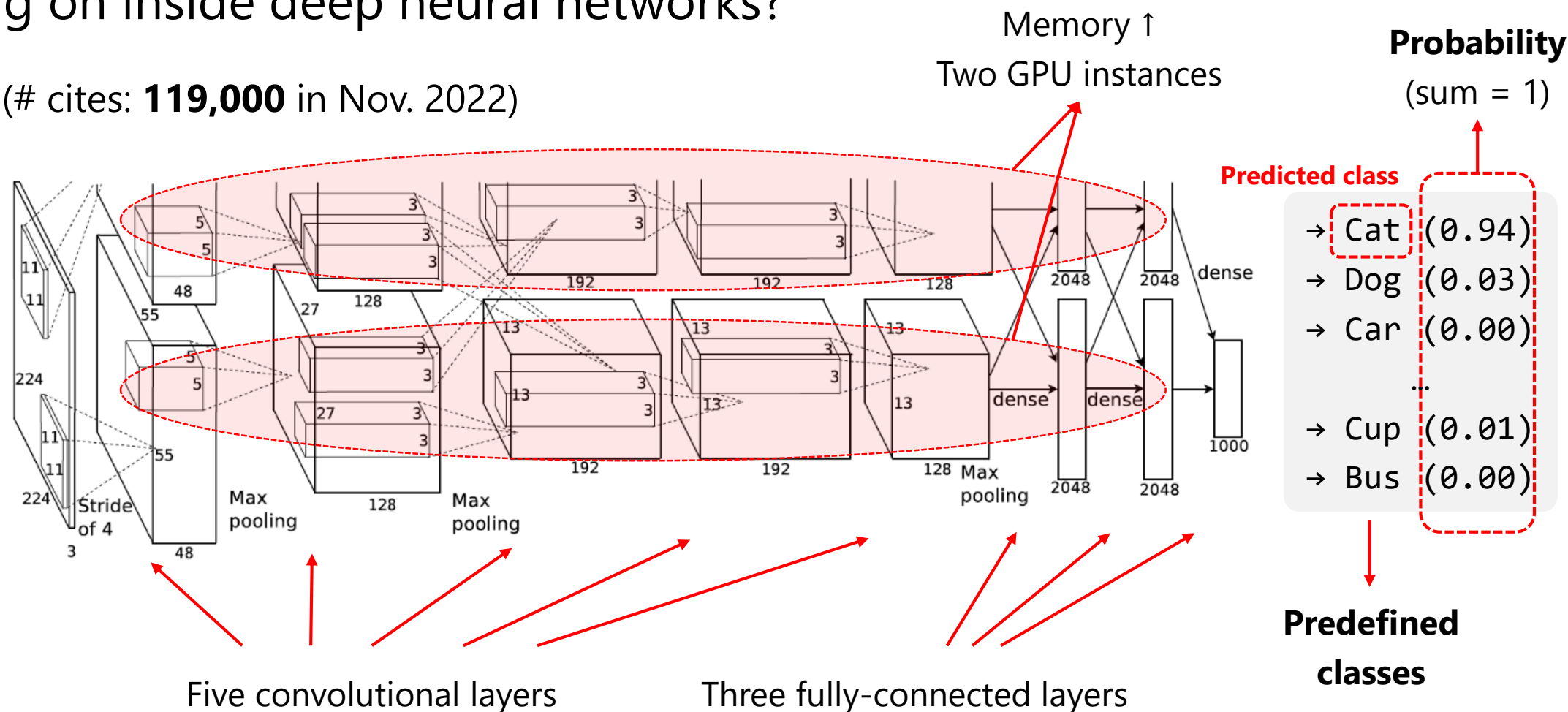
- AlexNet, VGG, GoogleNet, ResNet

Image Classification

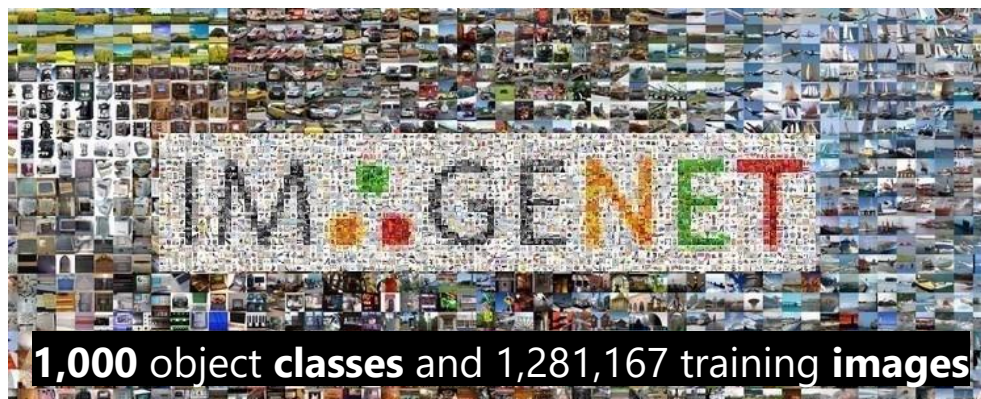
The most fundamental task using deep learning!

→ What's going on inside deep neural networks?

→ AlexNet [1] (# cites: **119,000** in Nov. 2022)



AlexNet: Breakthrough in 2012

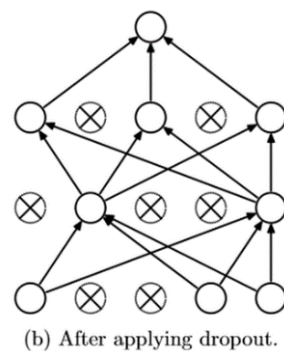
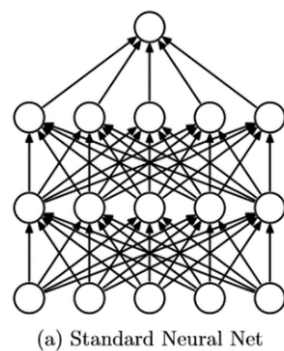
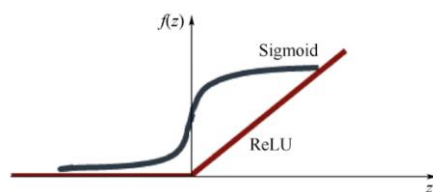
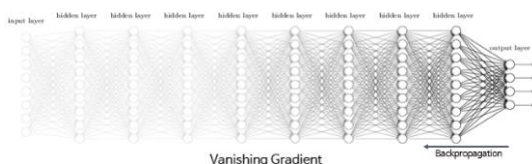
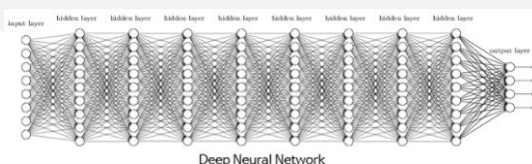


Troublesome of previous neural networks

- Local minimum or slow learning
- Overfitting
- Small data
- Time complexity
- Vanishing gradients

AlexNet

- Big data: ImageNet Challenge
- GPUs
- ReLU (Rectified Linear Unit)
- Dropout
- Deeper network



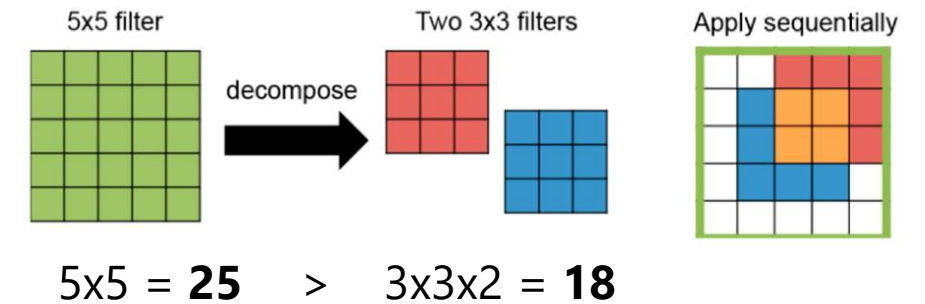
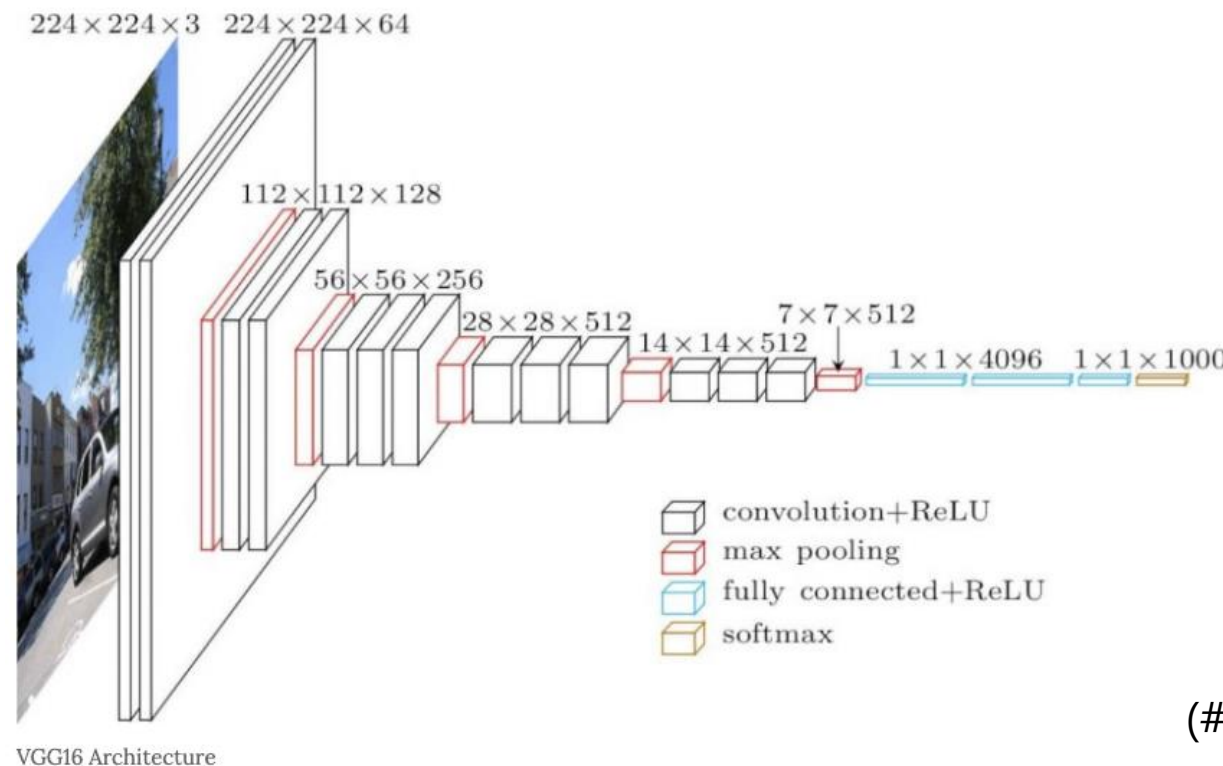
nVIDIA®

VGG: ImageNet Challenge (2014) 2nd Place

Small filters + Deeper networks + Beautifully uniform design

→ Why use smaller filters?

Number of parameters ↓ (efficiency ↑) + Deeper layer (nonlinearity ↑)

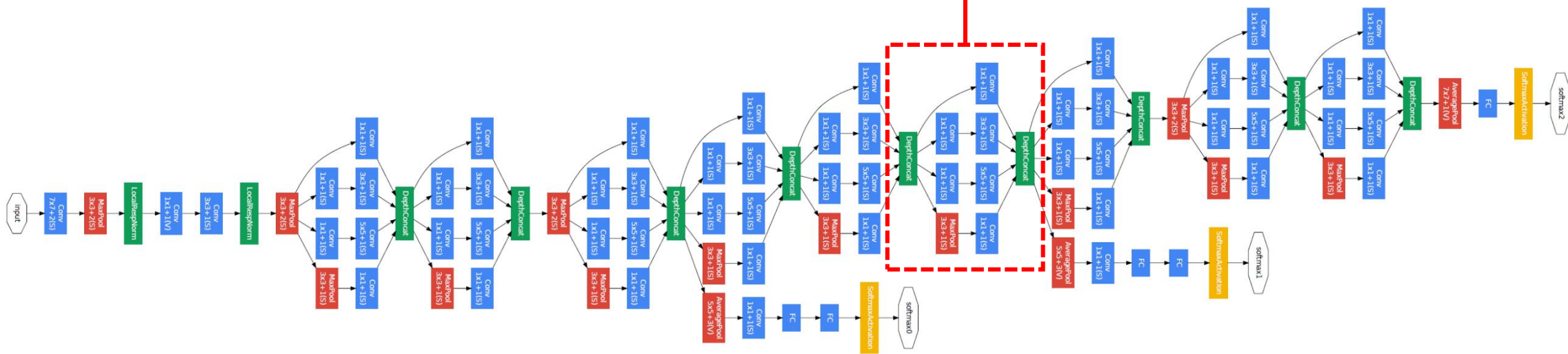


(# cites: **89,600** in Nov. 2022)

GoogleNet: ImageNet Challenge (2014) Winner

Deeper networks with a computational efficiency

- **Inception** module: Local network topology (network within a network)
- 5M params. ($\times 12$ less than AlexNet)
- Fully convolutional networks

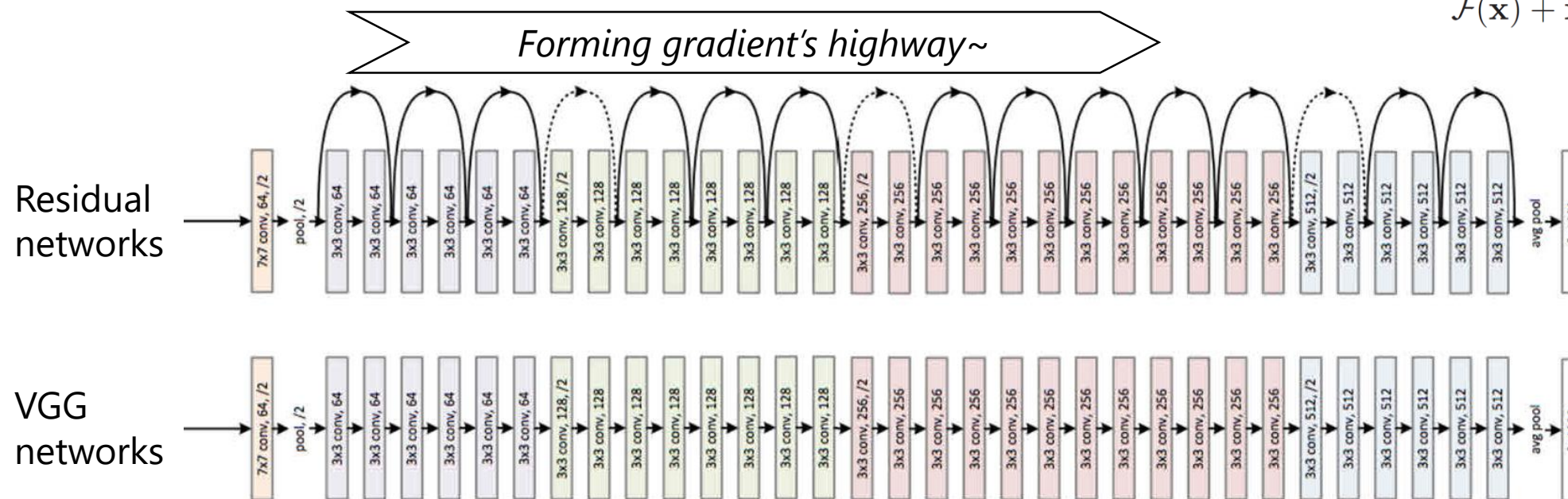
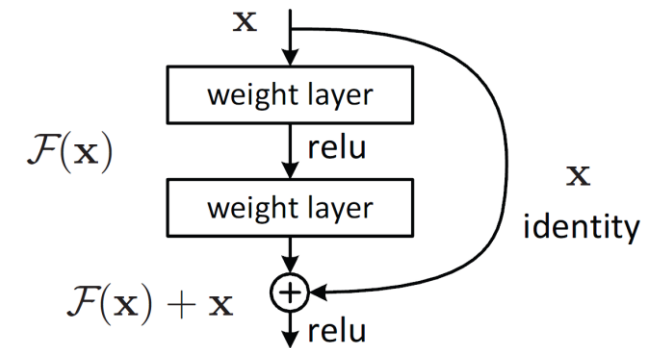


(# cites: **43,600** in Nov. 2022)

ResNet: ImageNet Challenge (2015) Winner

Major breakthrough in the network architecture

- Better than "human performance" in ImageNet Challenge
- **Residual units** make models easy to train!



(# cites: **136,800** in Nov. 2022)

Summary

Basic computer vision tasks inspired by human visual system

- Semantic and geometric scene understanding
- Basic visual perception tasks: image classification, detection, segmentation, ...

Basic deep neural networks

- AlexNet, VGG, GoogleNet, ResNet, ...
- Take-home message:

“Not all **complex** and **deep** networks are good, but how well you **regularize (represent)** **multi-dimensional features** is the key to improve the performance.”

Experiments

Image classification

Code is available in <https://view.kentech.ac.kr/f088fa7f-874e-44bc-bd6d-6084b42dfdf7>

```
$ python alexnet.py
```