



深度学习平台与应用

第一讲：绪论

范琦

fanqi@nju.edu.cn

2024年9月4日

大 纲

- 课程概览
- 深度学习历史
- 深度学习应用
- 深度学习平台
- 课程后续安排

- 绪论
- 线性分类器
- 正则化与优化
- 神经网络与反向传播
- 深度学习科研 (Special)
- 图像分类
- 卷积神经网络
- 循环神经网络
- 注意力机制和Transformer
- 目标检测和分割
- 视频理解
- 网络可视化
- 深度自监督学习
- 多模态模型
- 3D视觉模型
- 深度生成模型

■ 前置课程：

- 高等数学，微积分
- 概率论与数理统计，线性代数
- 最优化方法导论，人工智能导论，机器学习导论

■ 后续课程：

- 模式识别与计算机视觉
- 自然语言处理

- 了解深度学习**基础知识**
- 了解深度学习**研究方向和领域**
- 初步锻炼**科研思维**
- 培养论文**写作能力**
- 为后续实习/工作/深造**奠定基础**

授课 & 答疑

- 线下授课

- 助教：

- 段子轩 (2021dwin@sina.com)

- 鲁峰源 (fengyuan_lu@qq.com)

- 答疑：面对面/邮件（请并抄送给助教）

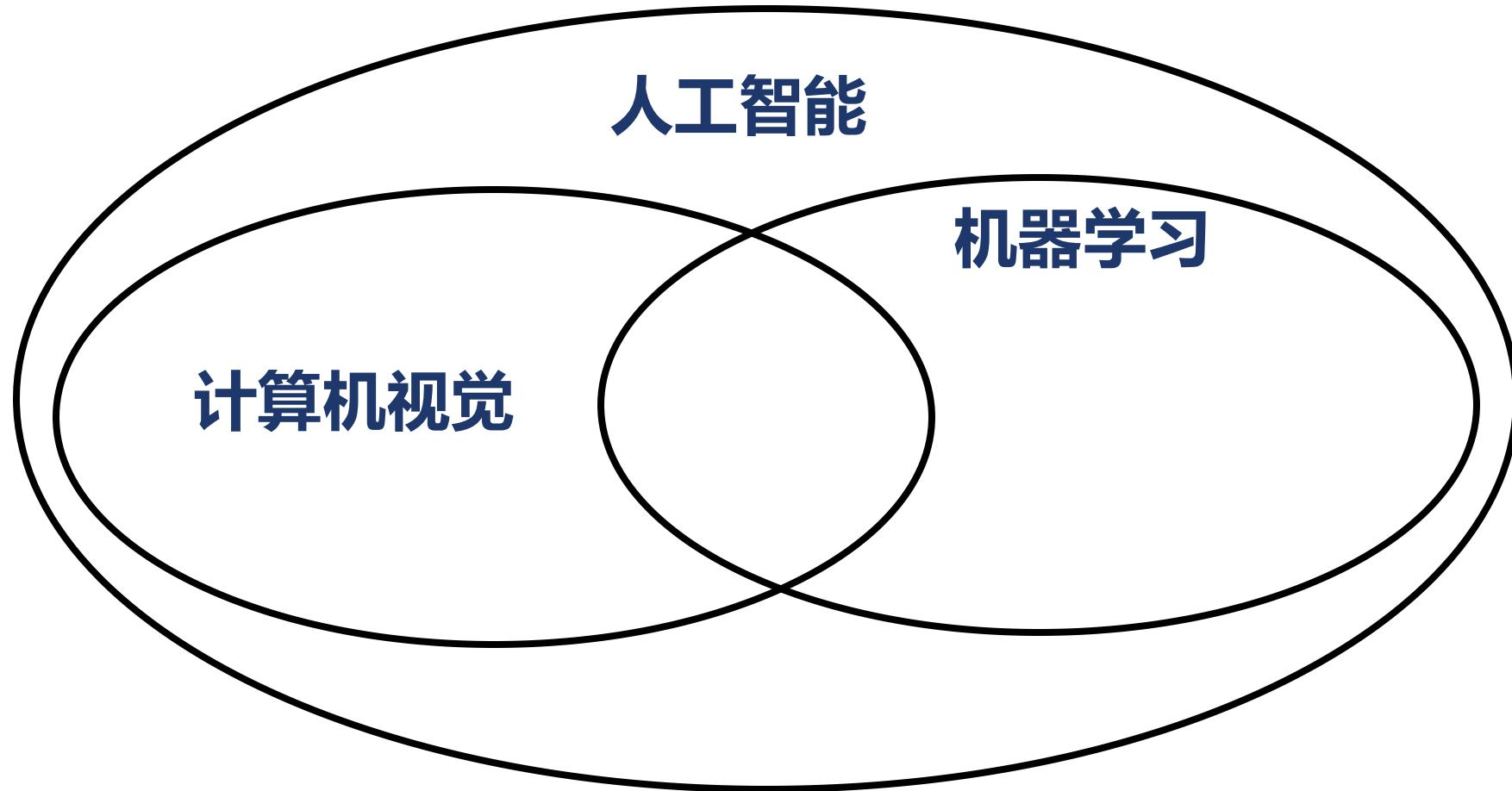
- 课程考核：在第二节课公布

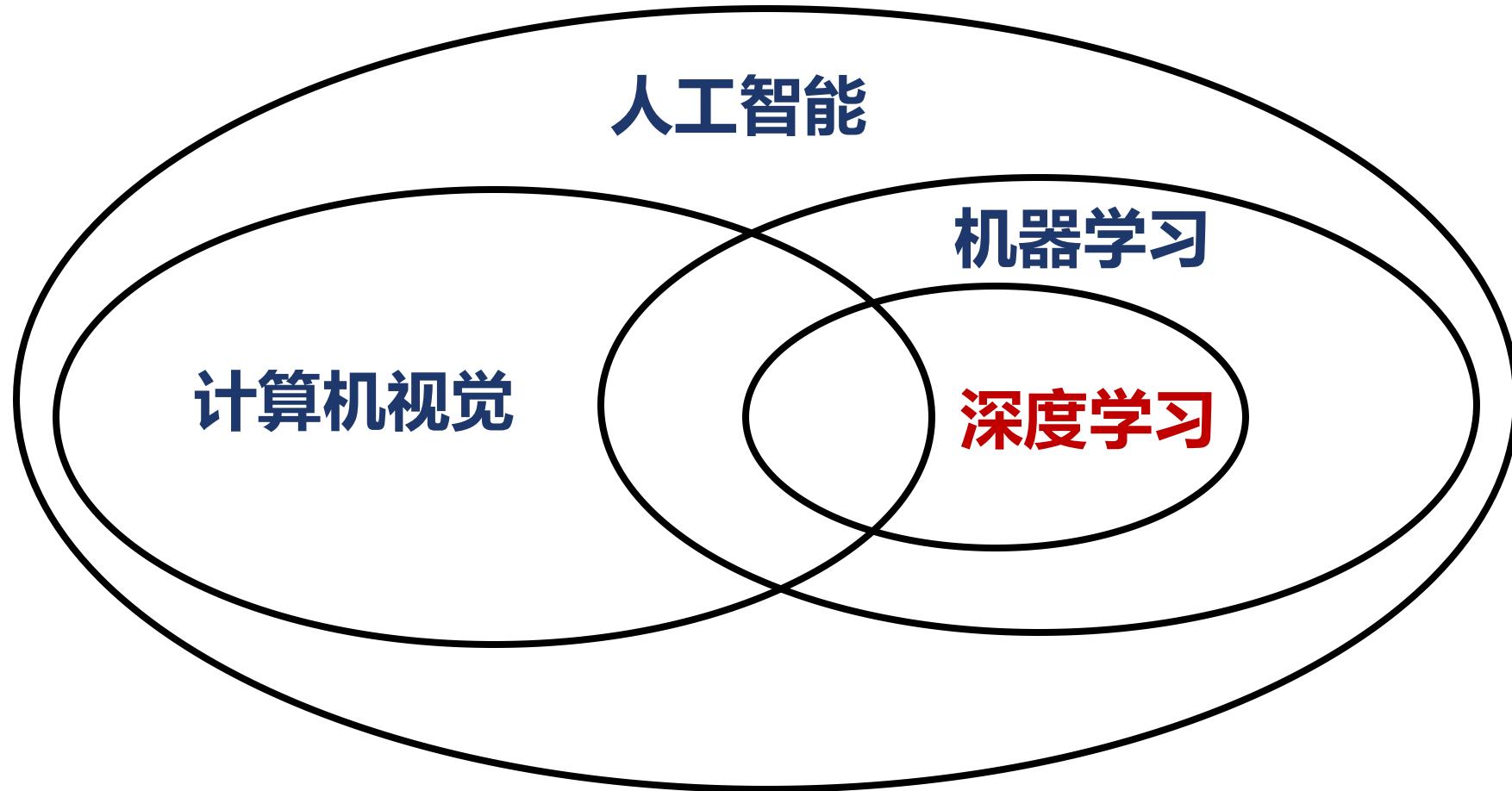


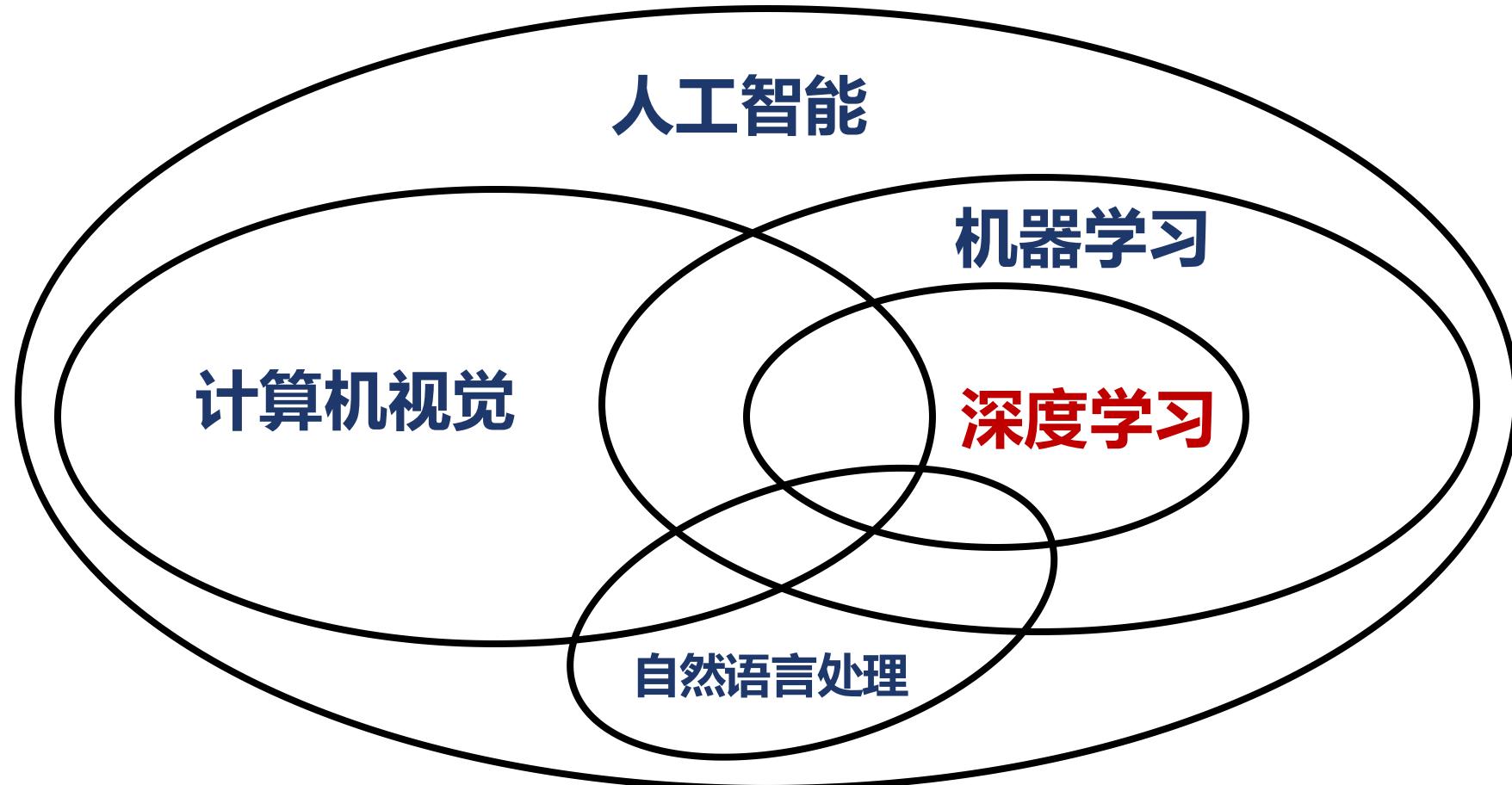
大 纲

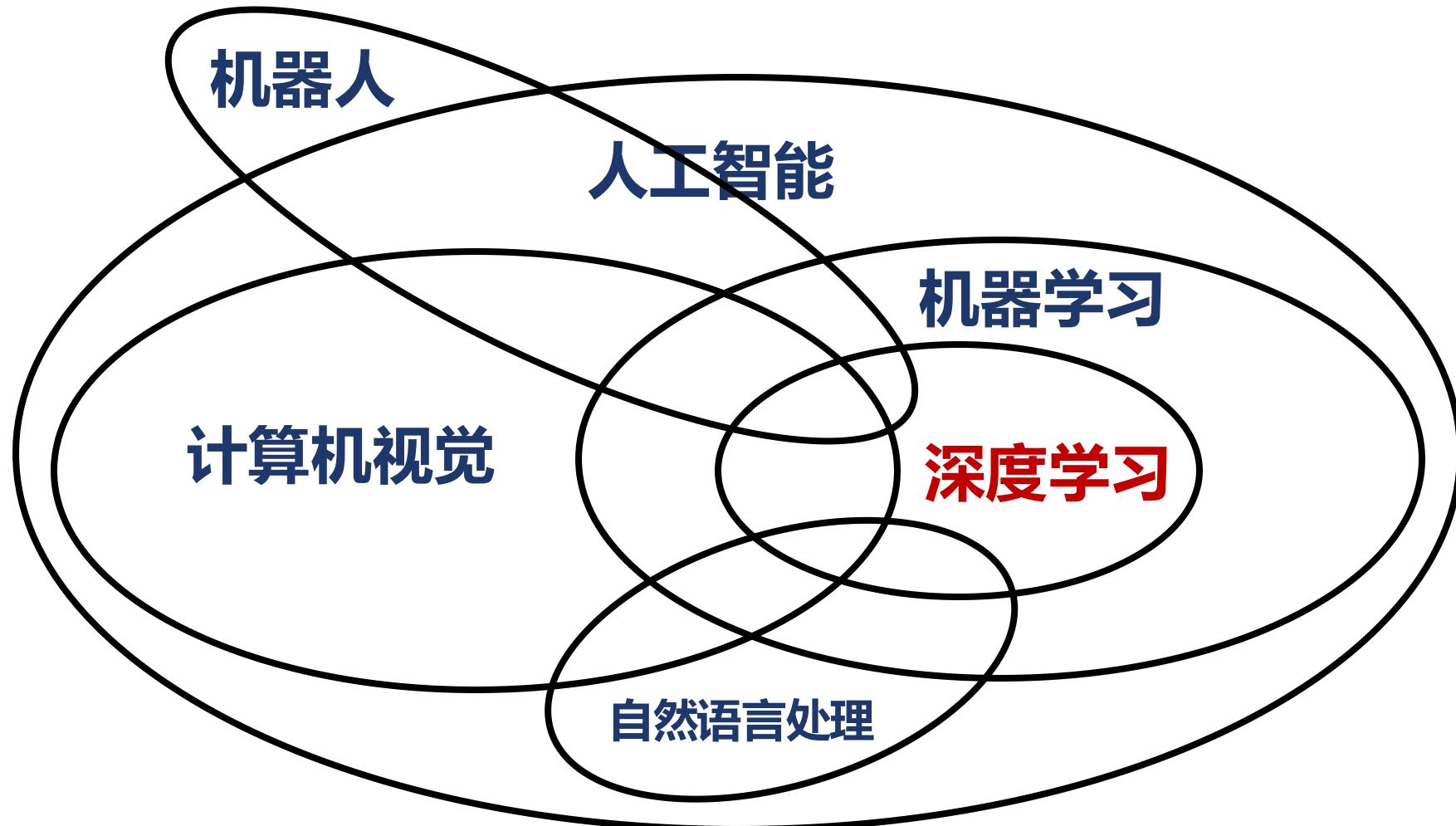
- 课程概览
- 深度学习历史
- 深度学习应用
- 深度学习平台
- 课程后续安排

人工智能



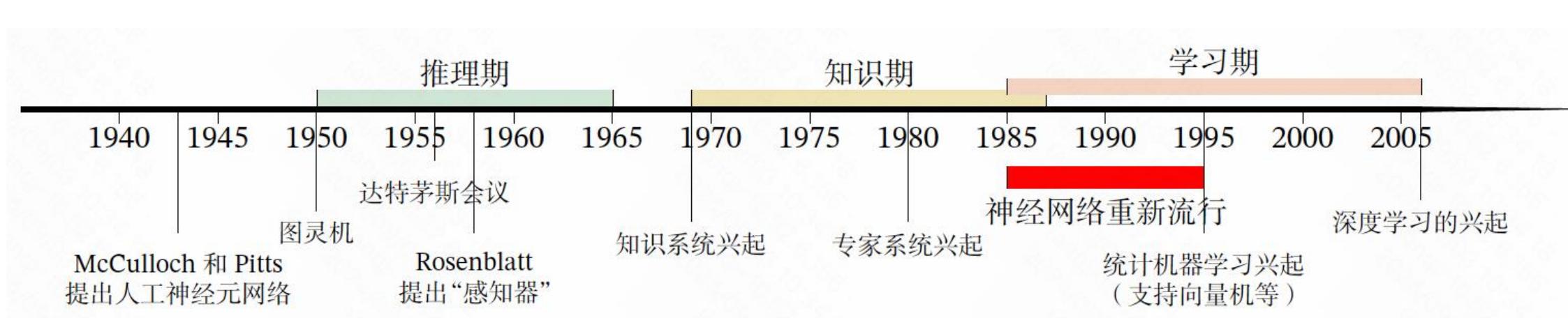






■ 人工智能 (Artificial Intelligence, AI)

- 推理期：基于逻辑或者事实归纳出来一些规则
- 知识期：领域专家构建的知识库+推理机
- 学习期：让计算机从数据中自己学习



- 机器学习 (Machine Learning, ML)
 - 从有限的观测数据中学习出具有一般性的规律，并利用这些规律对未知数据进行预测
 - 浅层学习 (传统机器学习) 关注学习预测模型，**不涉及特征学习**，其特征主要靠**特征工程**来提取



- 表示学习 (Representation Learning)
 - 自动地学习出有效的特征，并提高最终机器学习模型的性能
 - 关键问题是解决语义鸿沟 (Semantic Gap)
 - 语义鸿沟：输入数据的底层特征和高层语义信息之间的不一致性和差异性
 - 核心问题一：什么是一个好的表示
 - 核心问题二：如何学习到好的表示

■ 好的表示

- 具有**很强的表示能力**，即同样大小的向量可以表示更多信息
- 使后续的学习任务变得简单，即需要**包含更高层的语义信息**
- 具有一般性，是任务或领域独立的，可以迁移到其他任务上

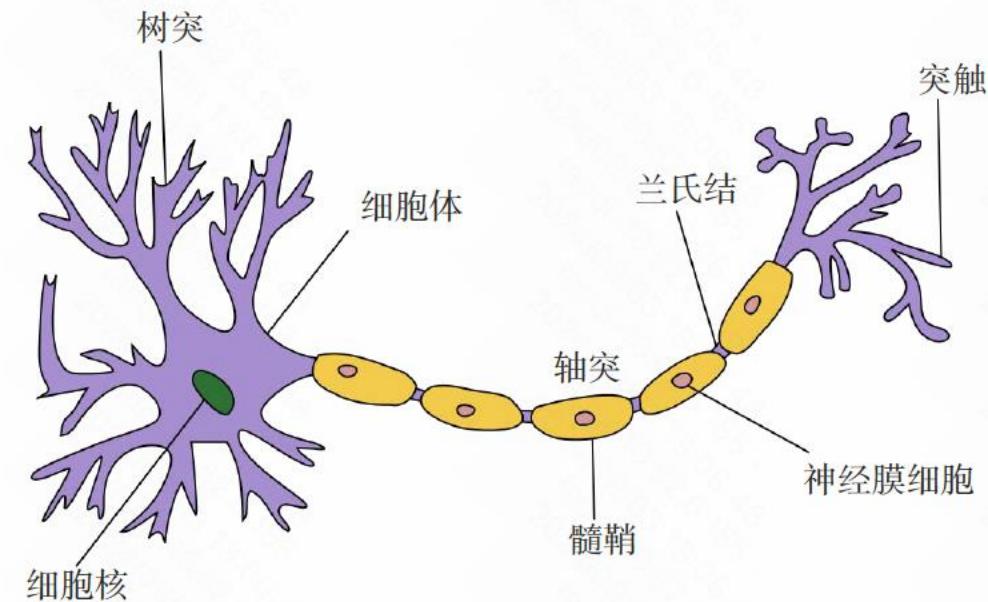
■ 学习好的表示

- 关键是构建**具有一定深度的多层次特征表示**

- 深度学习 (Deep Learning, DL)
 - 是机器学习的一个子问题，特征学习+预测模型
 - 构建并从数据中学习具有一定“深度”的模型
 - 关键问题是贡献度分配问题
 - 神经网络使用误差反向传播算法可以较好解决关键问题

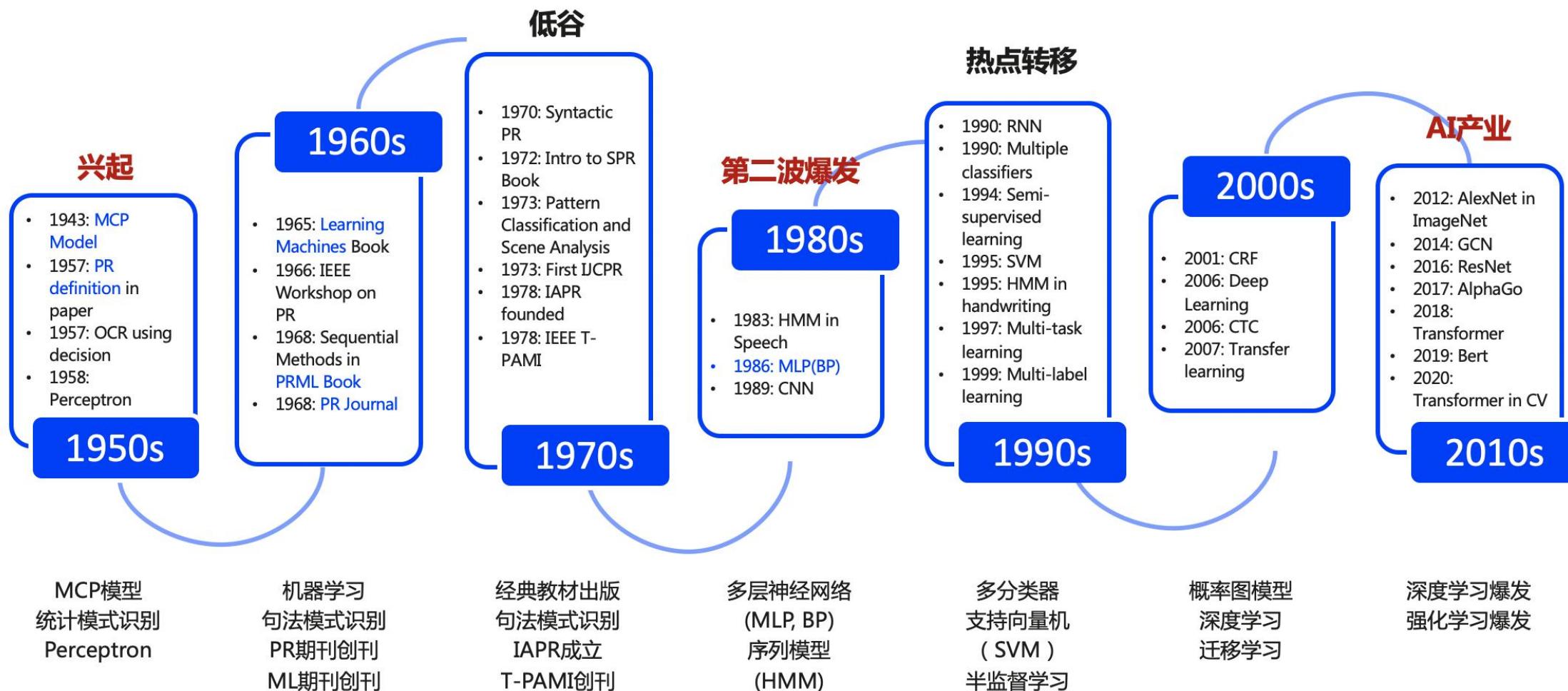


- 人工神经网络
 - 模拟人脑神经网络而设计的一种计算模型
 - 由大量神经元通过极其丰富和完善的连接而构成的自适应非线性动态系统
- 一个通用的函数逼近器
- 反向传播算法



深度学习重要历史节点

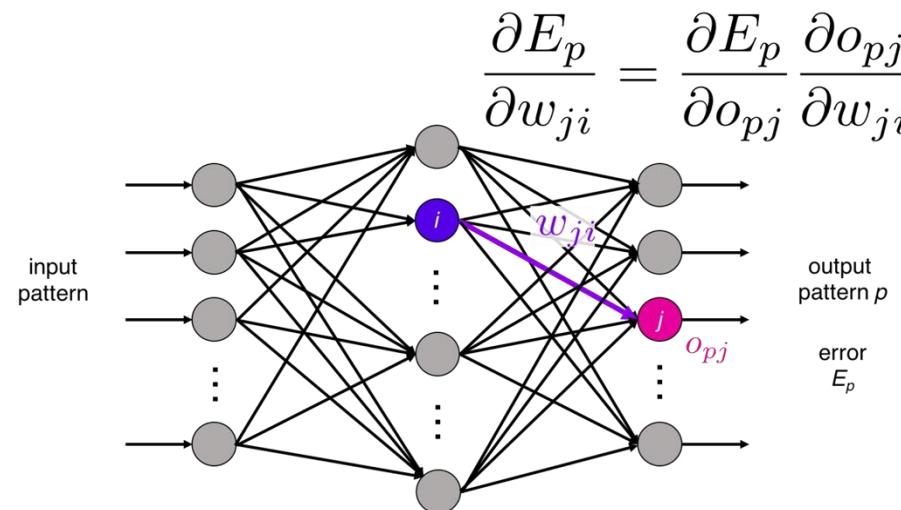
深度学习发展历程



■ 深度学习发展历史阶段

- 1943年 ~ 1969年: 模型提出
- 1969年 ~ 1983年: 冰河期
- 1983年 ~ 1995年: 反向传播算法引起的复兴
- 1995年 ~ 2006年: 流行度降低
- 2006年 ~ 至今: 深度学习的崛起

- 反向传播算法 (Backprop: Rumelhart, Hinton, and Williams, 1986)
- 反向传播计算梯度进行网络训练
- 最成功的神经网络学习算法



第1阶段：激励传播 [编辑]

每次迭代中的传播环节包含两步：

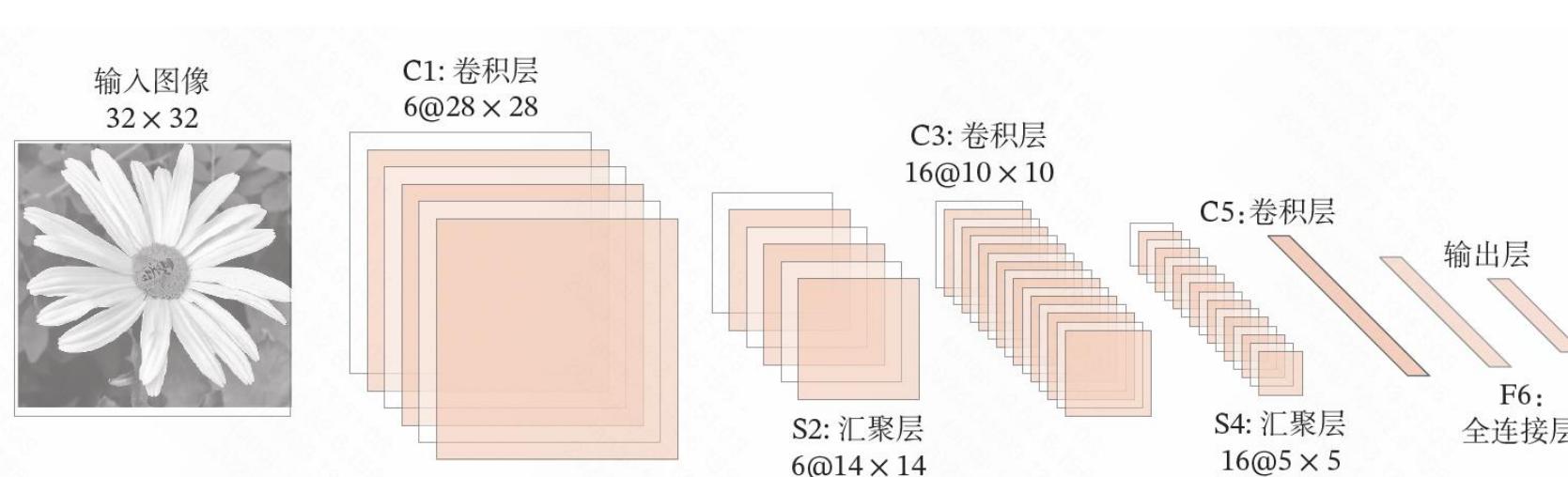
1. (前向传播阶段) 将训练输入送入网络以获得预测结果；
2. (反向传播阶段) 对预测结果同训练目标求差([损失函数](#))。

第2阶段：权重更新 [编辑]

对于每个突触上的权重，按照以下步骤进行更新：

1. 将输入激励和响应误差相乘，从而获得权重的梯度；
2. 将这个梯度乘上一个比例并取反后加到权重上。

- LeNet: Convolutional Networks: LeCun et al, 1998
 - 将反向传播算法引入卷积神经网络
 - 成功应用于手写字体识别系统，被银行系统广泛应用
 - 与现代神经网络非常相似



■ 深度信念网络 (Hinton et al, 2006)

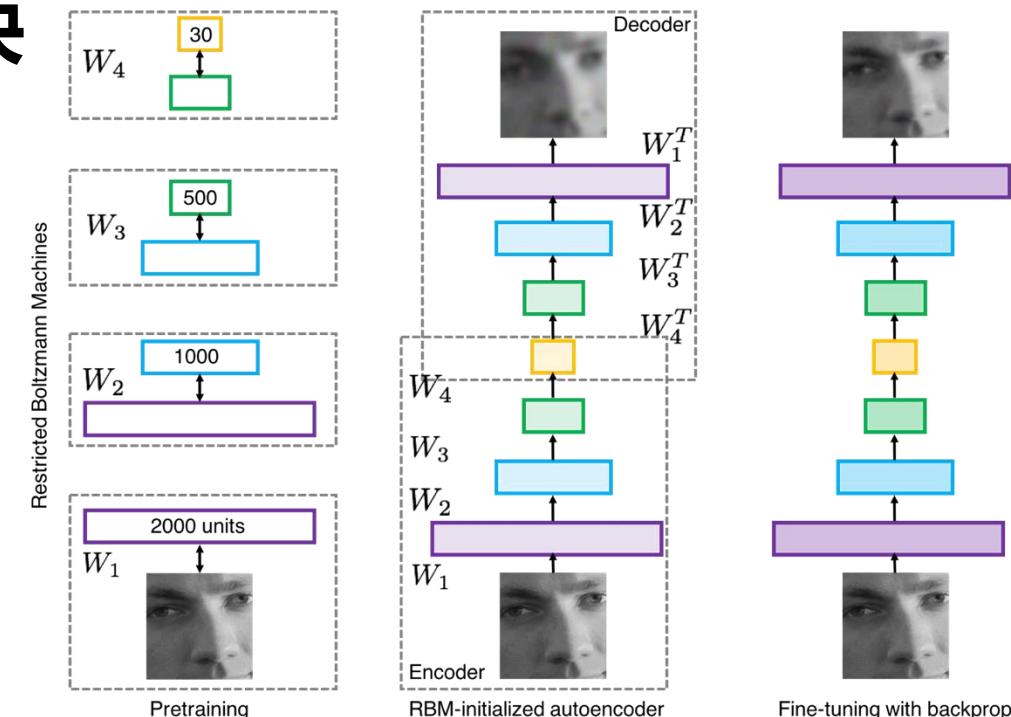
■ 提出“**预训练 + 精调**”方式来解决

深度神经网络难以训练的问题

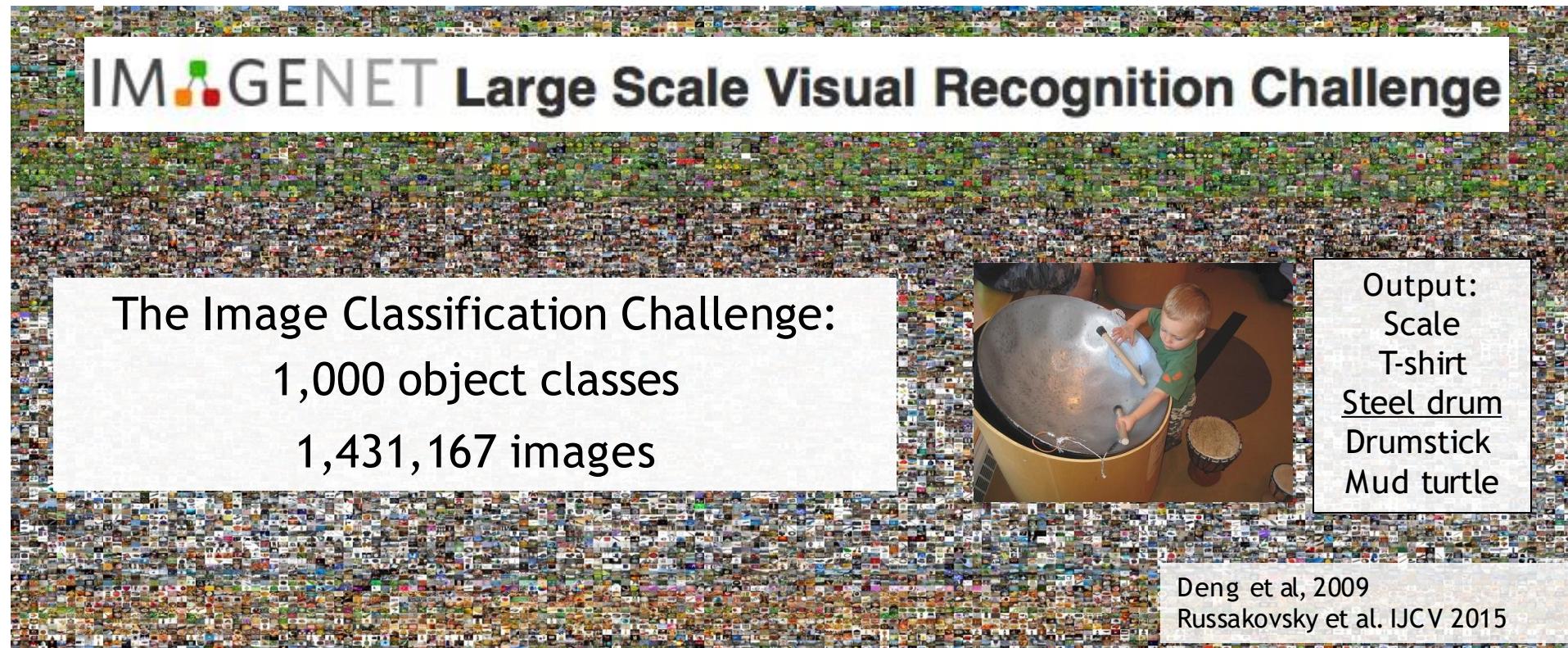
■ 逐层预训练学习一个深度信念网络

■ 将其权重作为一个多层前馈神经网
络的**初始化权重**

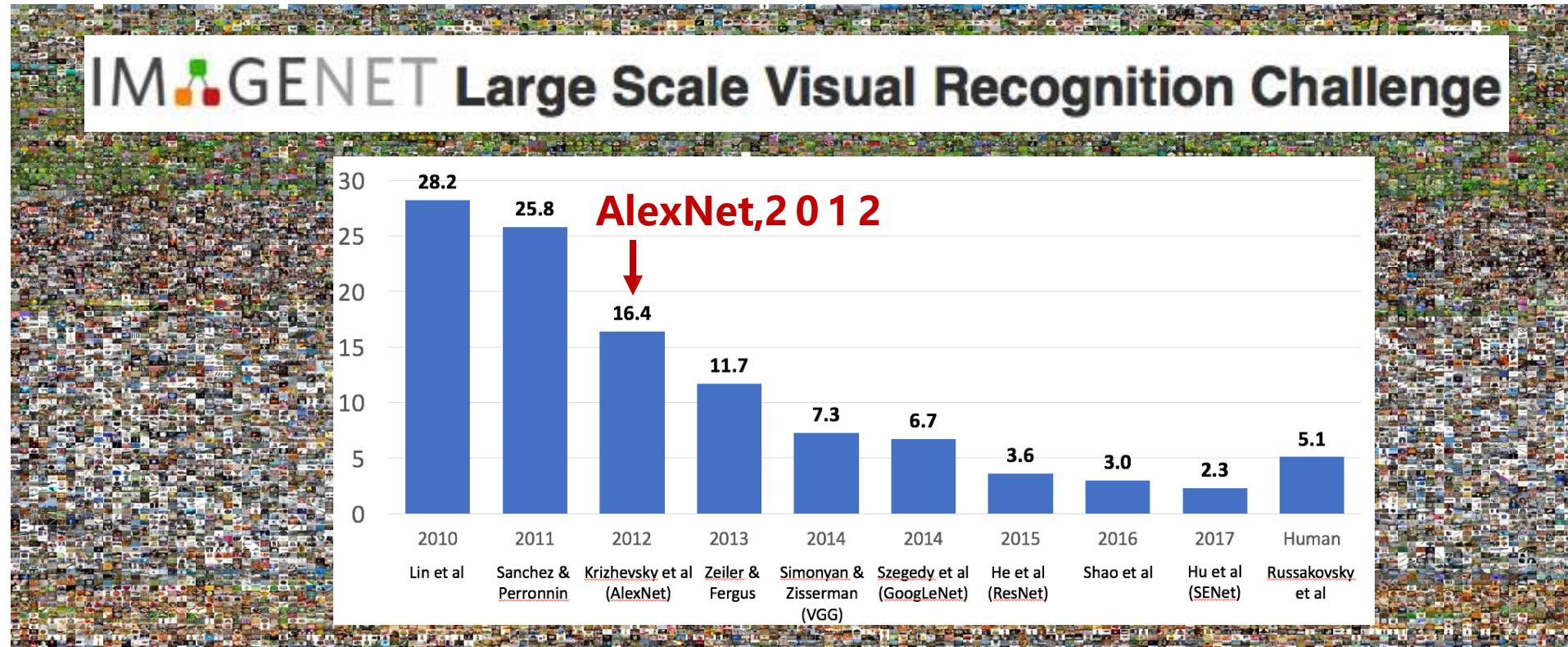
■ 再用**反向传播算法**进行**精调**



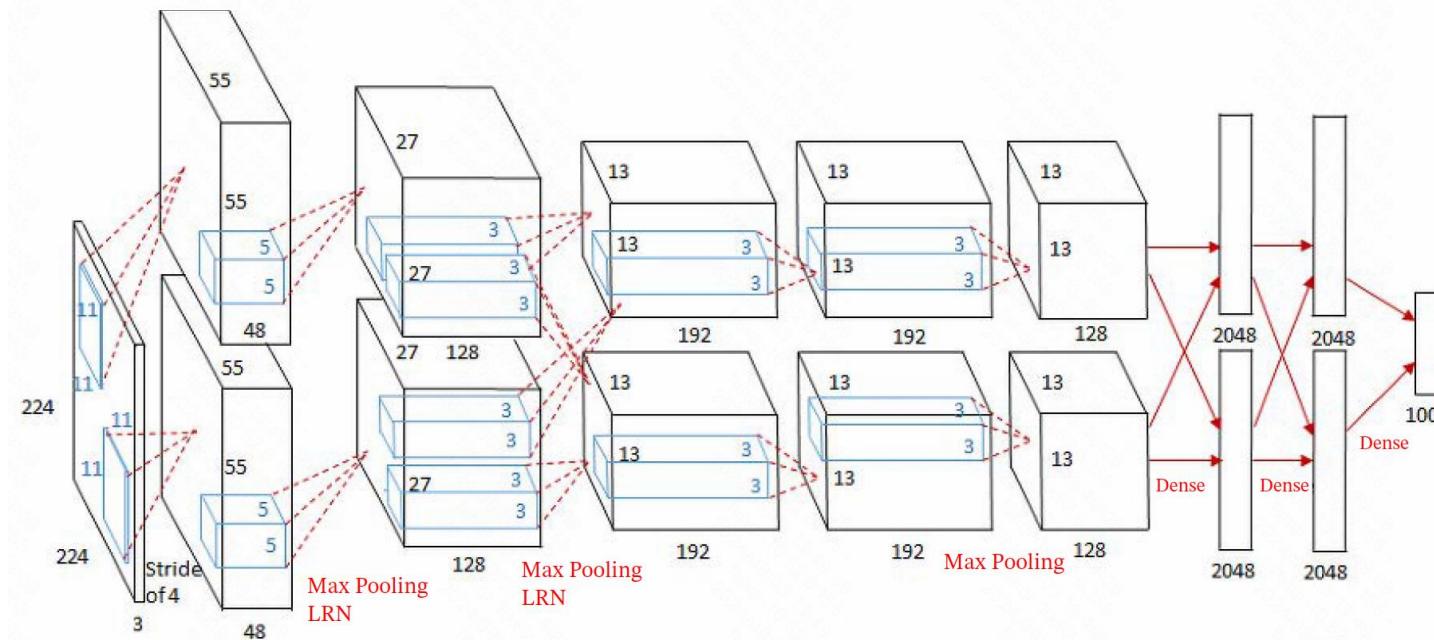
- 大规模数据集 (ImageNet: Deng et al, 2009)
- 1000个类别, 1431167张图片



■ 人们认识到数据对深度学习的重要性

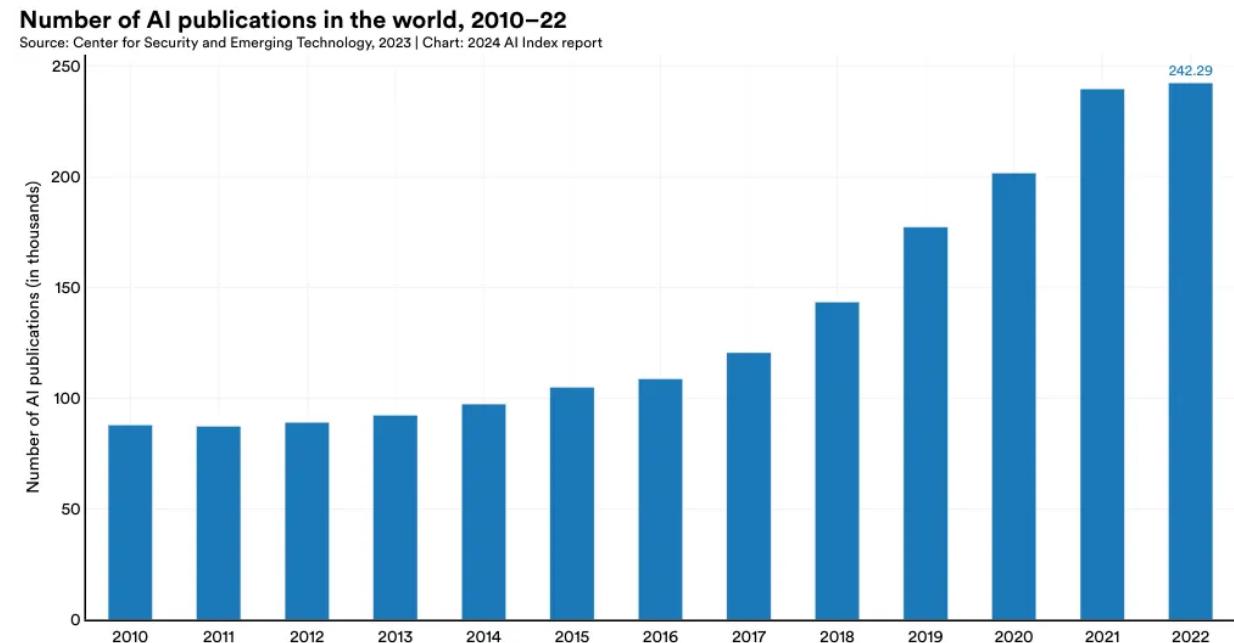
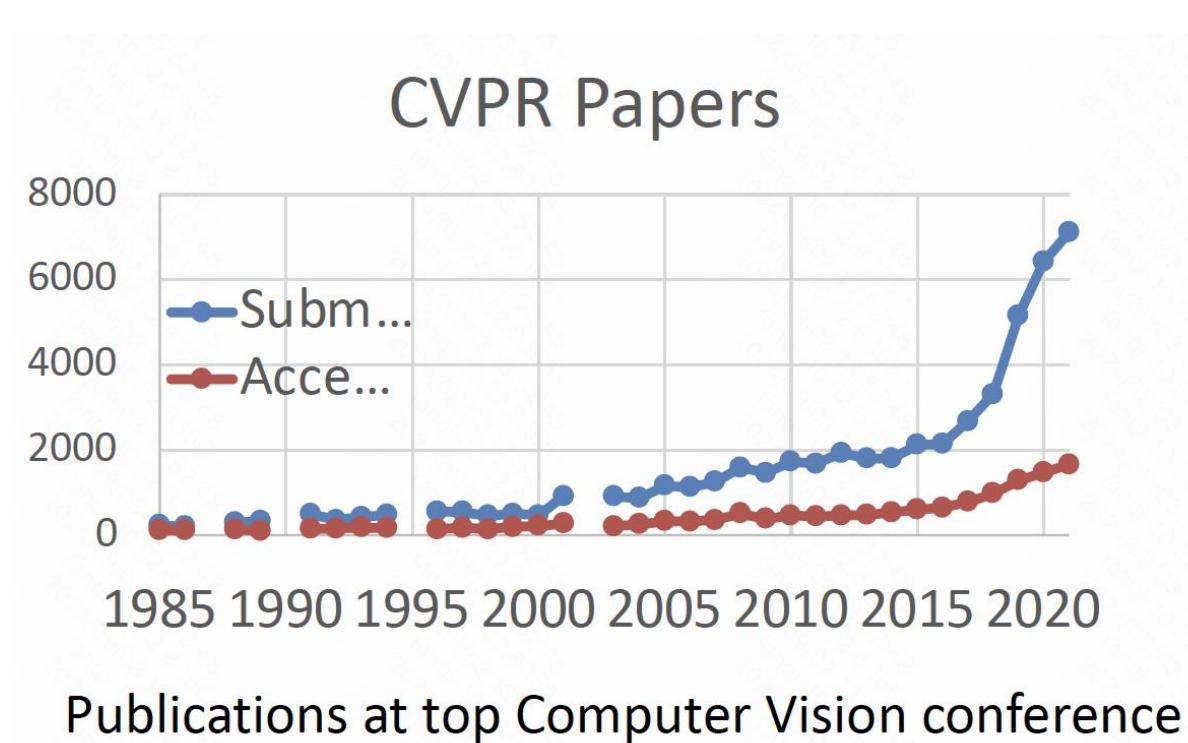


- 第一个现代深度卷积网络 AlexNet: Krizhevsky, et al, 2012
- 端到端训练，使用GPU 进行并行训练
- 使用 ReLU、Dropout、数据增强等技术

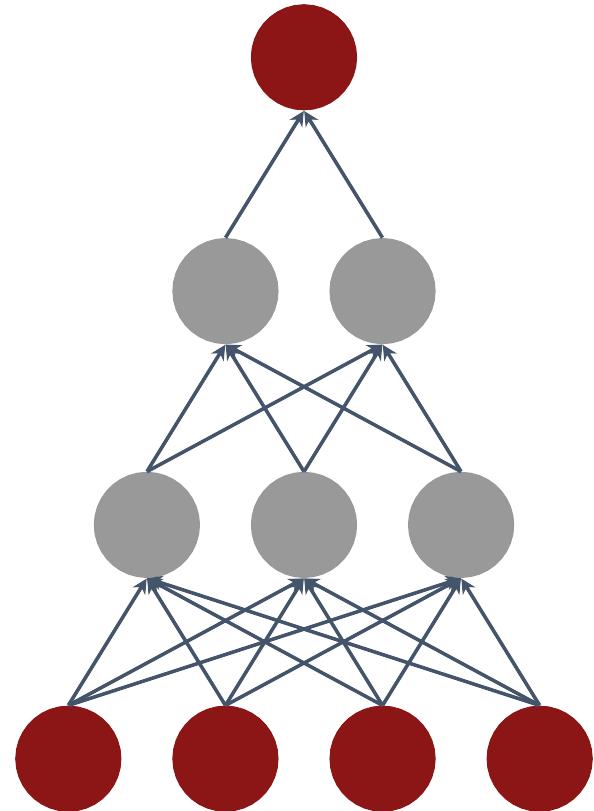
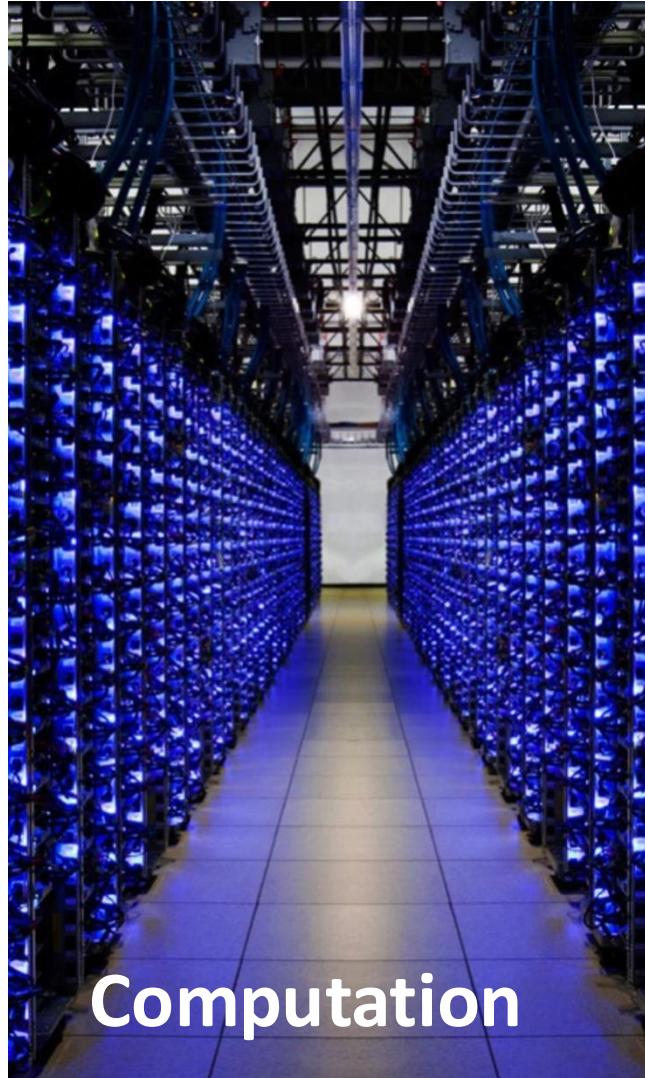


深度学习重要历史节点

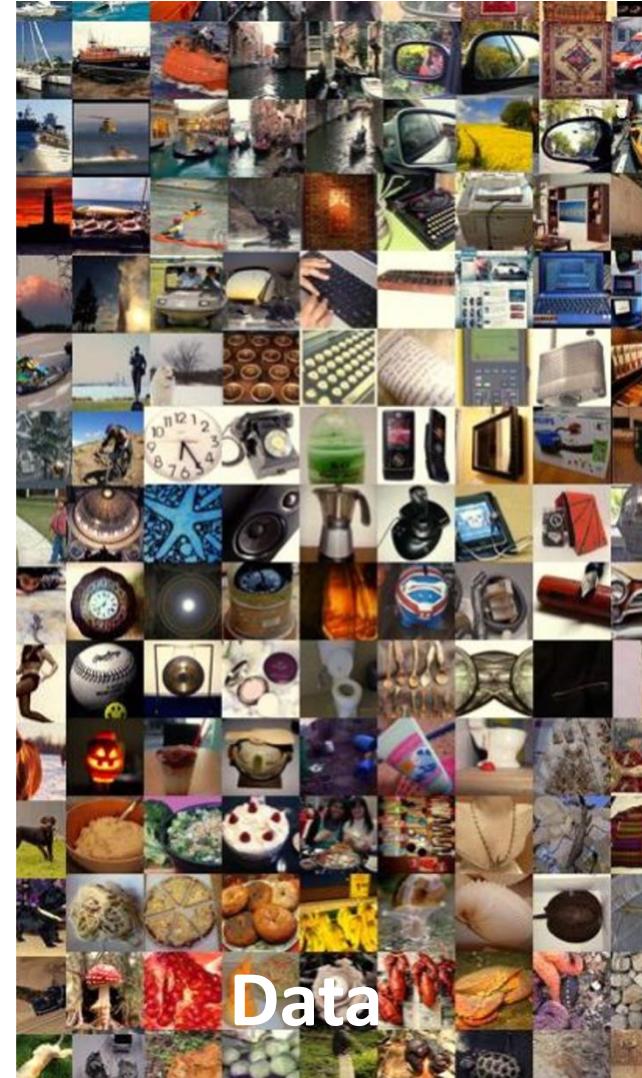
■ 深度学习大爆发：2012年至今



深度学习数据集



Algorithms



- Meta Pointer: A large collection organized by CV Datasets.
- Yet another Meta pointer
- Hugging face datasets: collection of generic datasets available on hugging face
- ImageNet: a large-scale image dataset for visual recognition organized by WordNet hierarchy
- SA-1B: dataset of a large number of images and segmentation masks to segment objects in those images
- COCO: large-scale object detection, segmentation, and captioning dataset
- Open Images: a dataset of ~9M images annotated with image-level labels, object bounding boxes, object segmentation masks, visual relationships, and localized narratives

- **Cityscapes Dataset**: This dataset focuses on semantic understanding of urban street scenes, with pixel-level annotations for various object classes such as cars, pedestrians, and roads
- **DeepFashion**: a large-scale clothes dataset containing over 800,000 diverse fashion images annotated with bounding boxes, clothing categories, and attributes
- **Objaverse**: a large-scale 3D asset database
- **SUN Database**: a benchmark for scene recognition and object detection with annotated scene categories and segmented objects
- **Places Database**: a scene-centric database with 205 scene categories and 2.5 millions of labelled images
- **NYU Depth Dataset v2**: a RGB-D dataset of segmented indoor scenes

- [**Flickr100M**](#): 100 million creative commons Flickr images
- [**Labeled Faces in the Wild**](#): a dataset of 13,000 labeled face photographs
- [**Human Pose Dataset**](#): a benchmark for articulated human pose estimation
- [**YouTube Faces DB**](#): a face video dataset for unconstrained face recognition in videos
- [**UCF101**](#): an action recognition data set of realistic action videos with 101 action categories
- [**HMDB-51**](#): a large human motion dataset of 51 action classes
- [**ActivityNet**](#): A large-scale video dataset for human activity understanding
- [**Moments in Time**](#): A dataset of one million 3-second videos
- Vision-language datasets: [**Visual Genome**](#), [**Flickr30k**](#), [**VQA v2**](#), [**ADE20K**](#), [**LAION**](#)

大 纲

- 课程概览
- 深度学习历史
- 深度学习应用
- 深度学习平台
- 课程后续安排

■ 数值回归

Real world input

6000 square feet,
4 bedrooms,
previously sold for
\$235K in 2005,
1 parking spot.

Model
input

$$\begin{bmatrix} 6000 \\ 4 \\ 235 \\ 2005 \\ 1 \end{bmatrix}$$

Model



Model
output

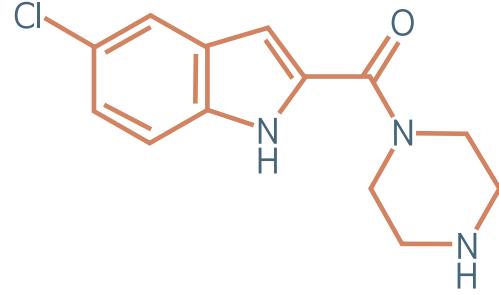
$$[340]$$

Real world output

Predicted price
is \$340k

■ 图回归

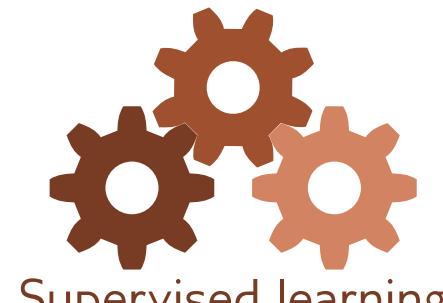
Real world input



Model
input

$$\begin{bmatrix} 1 \\ 0 \\ 1 \\ \vdots \\ 17 \\ 1 \\ 1 \\ \vdots \end{bmatrix}$$

Model



Model
output

$$\begin{bmatrix} -12.9 \\ 56.4 \end{bmatrix}$$

Real world output

Freezing point
is -12.9°C
Boiling point
is 56.4°C

■ 文本分类

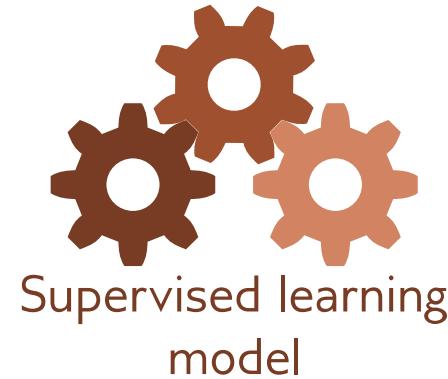
Real world input

“The steak was terrible,
the salad was rotten, and
the soup tasted like socks”

Model
input

$$\begin{bmatrix} 8672 \\ 8194 \\ 9804 \\ 8634 \\ 8672 \\ \vdots \end{bmatrix}$$

Model



Model
output

$$\begin{bmatrix} 0.02 \\ 0.98 \end{bmatrix}$$

Real world output

Positive
Negative

■ 音乐分类

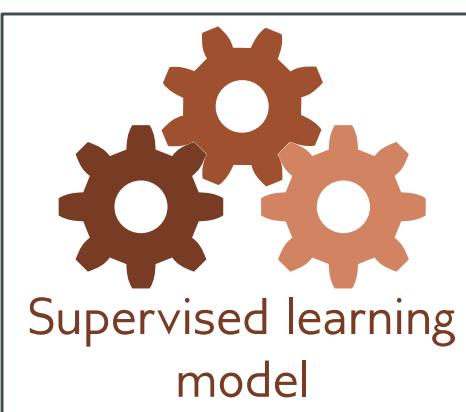
Real world input



Model
input

$$\begin{bmatrix} 125 \\ 12054 \\ 1253 \\ 6178 \\ 24 \\ 4447 \\ \vdots \end{bmatrix}$$

Model



Model
output

$$\begin{bmatrix} 0.03 \\ 0.52 \\ 0.18 \\ 0.07 \\ 0.12 \\ 0.08 \\ \vdots \\ 0.01 \end{bmatrix}$$

Real world output

Classical
Electronica
Hip Hop
Jazz
Pop
Metal
Punk

■ 图像分类

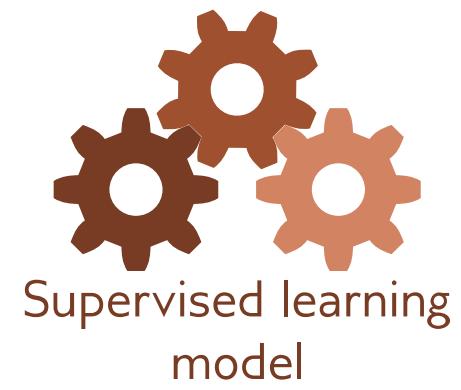
Real world input



Model
input

$$\begin{bmatrix} 124 \\ 140 \\ 156 \\ 128 \\ 142 \\ 157 \\ \vdots \end{bmatrix}$$

Model



Model
output

$$\begin{bmatrix} 0.00 \\ 0.00 \\ 0.01 \\ 0.89 \\ 0.05 \\ 0.00 \\ \vdots \\ 0.01 \end{bmatrix}$$

Real world output

Aardvark
Apple
Bee
Bicycle
Bridge
Clown
⋮

■ 图像分割

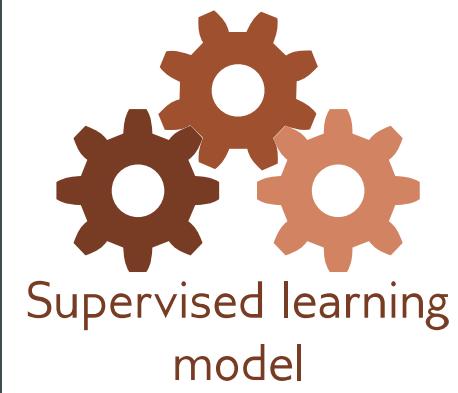
Real world input



Model
input

$$\begin{bmatrix} 183 \\ 204 \\ 231 \\ 185 \\ 204 \\ 232 \\ \vdots \end{bmatrix}$$

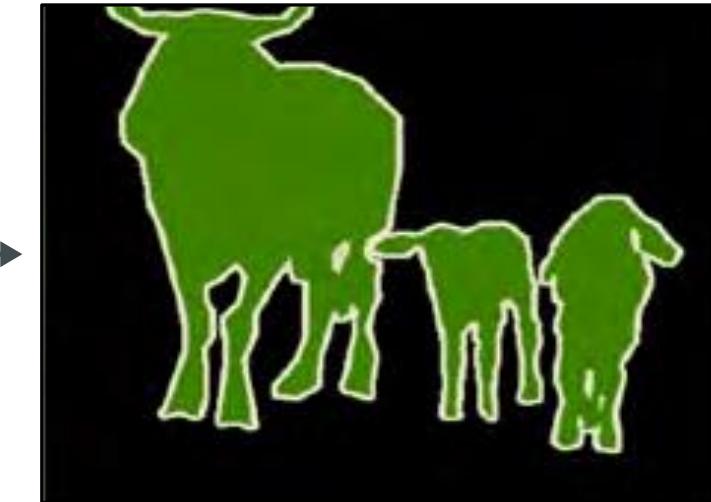
Model



Model
output

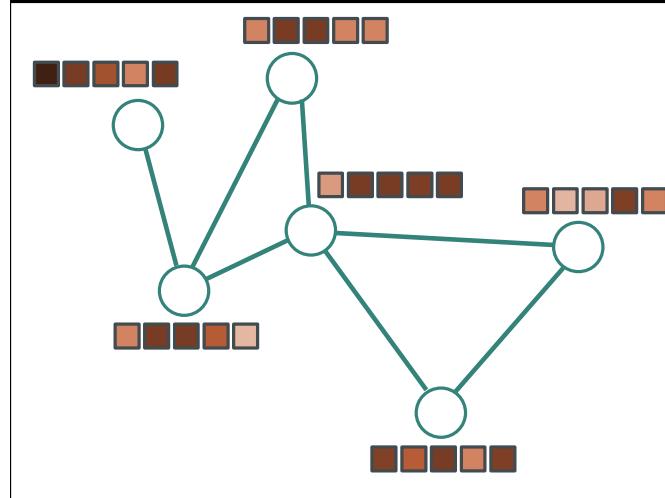
$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \end{bmatrix}$$

Real world output



■ 图节点分类

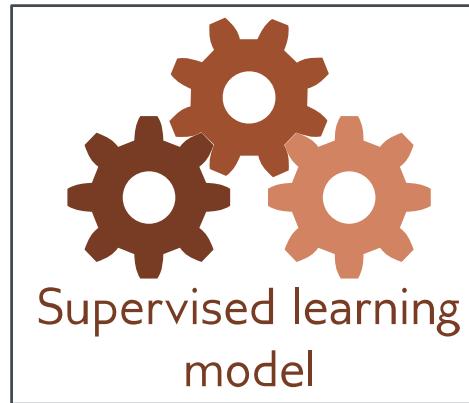
Real world input



Model
input

$$\begin{bmatrix} 1 \\ 0 \\ 1 \\ \vdots \\ 53 \\ 34 \\ 24 \\ \vdots \end{bmatrix}$$

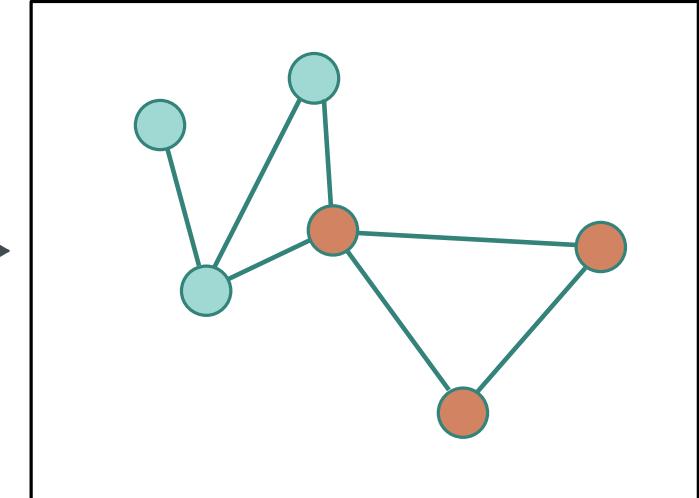
Model



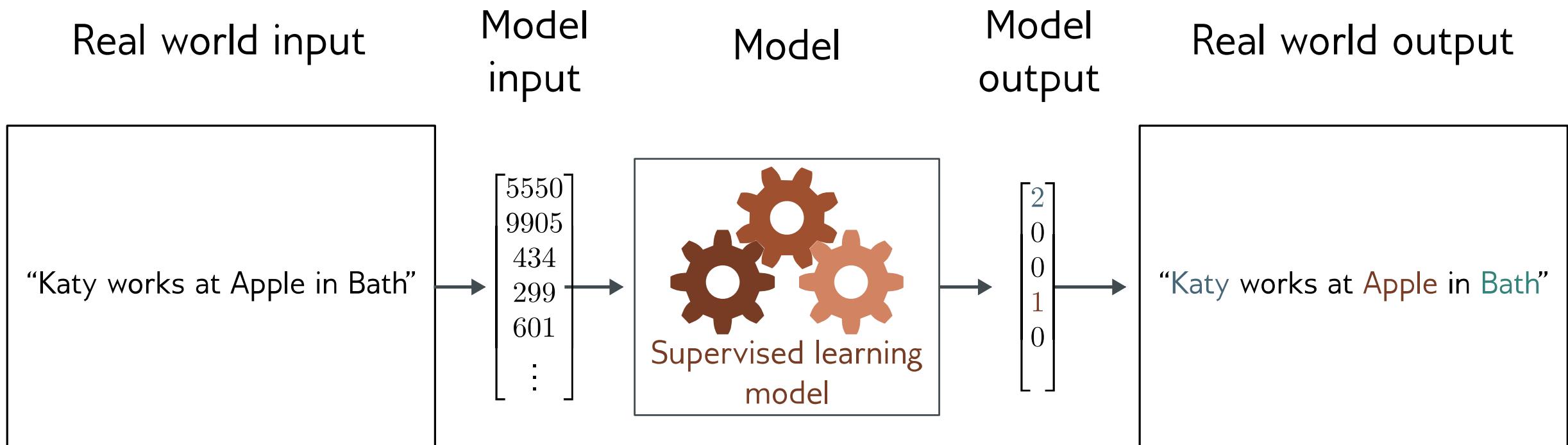
Model
output

$$\begin{bmatrix} 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ \vdots \end{bmatrix}$$

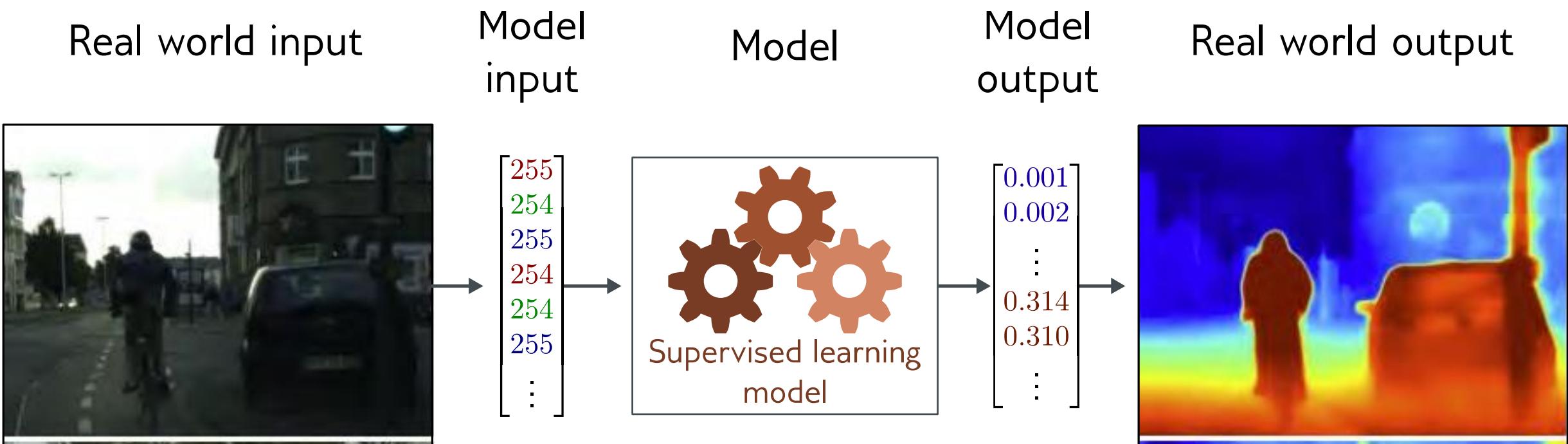
Real world output



■ 实体识别



■ 深度估计



■ 姿态估计

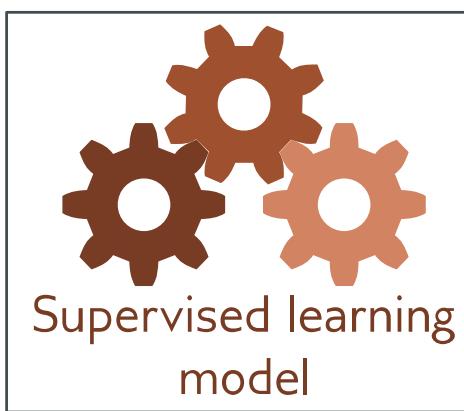
Real world input



Model
input

$$\begin{bmatrix} 3 \\ 5 \\ 4 \\ 3 \\ 5 \\ 5 \\ \vdots \end{bmatrix}$$

Model



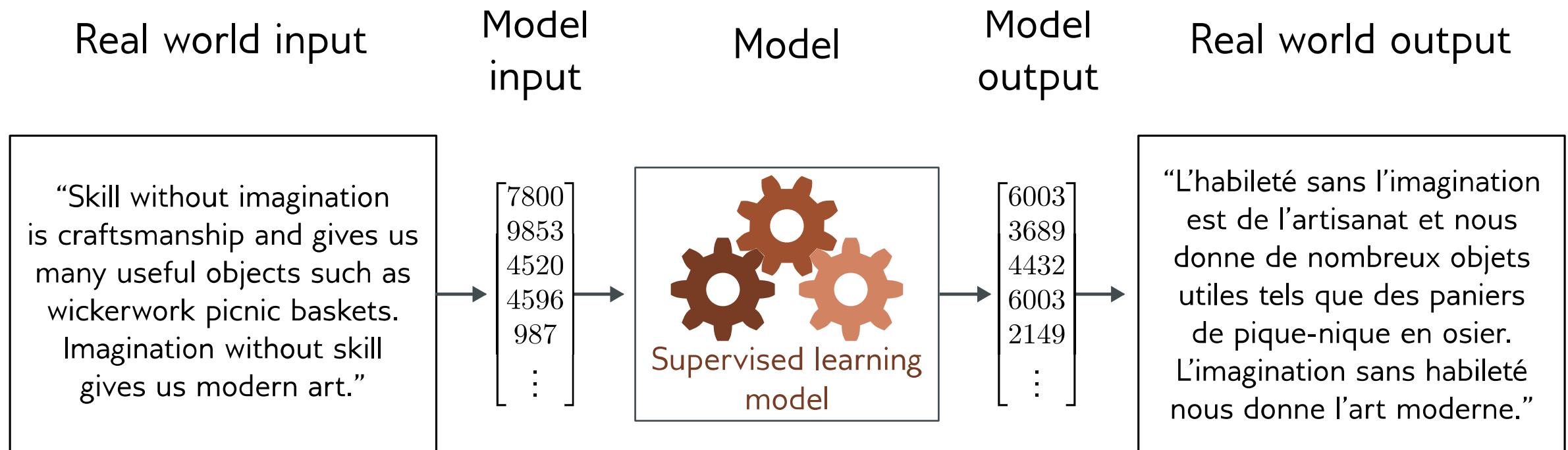
Model
output

$$\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 3 \\ \vdots \end{bmatrix}$$

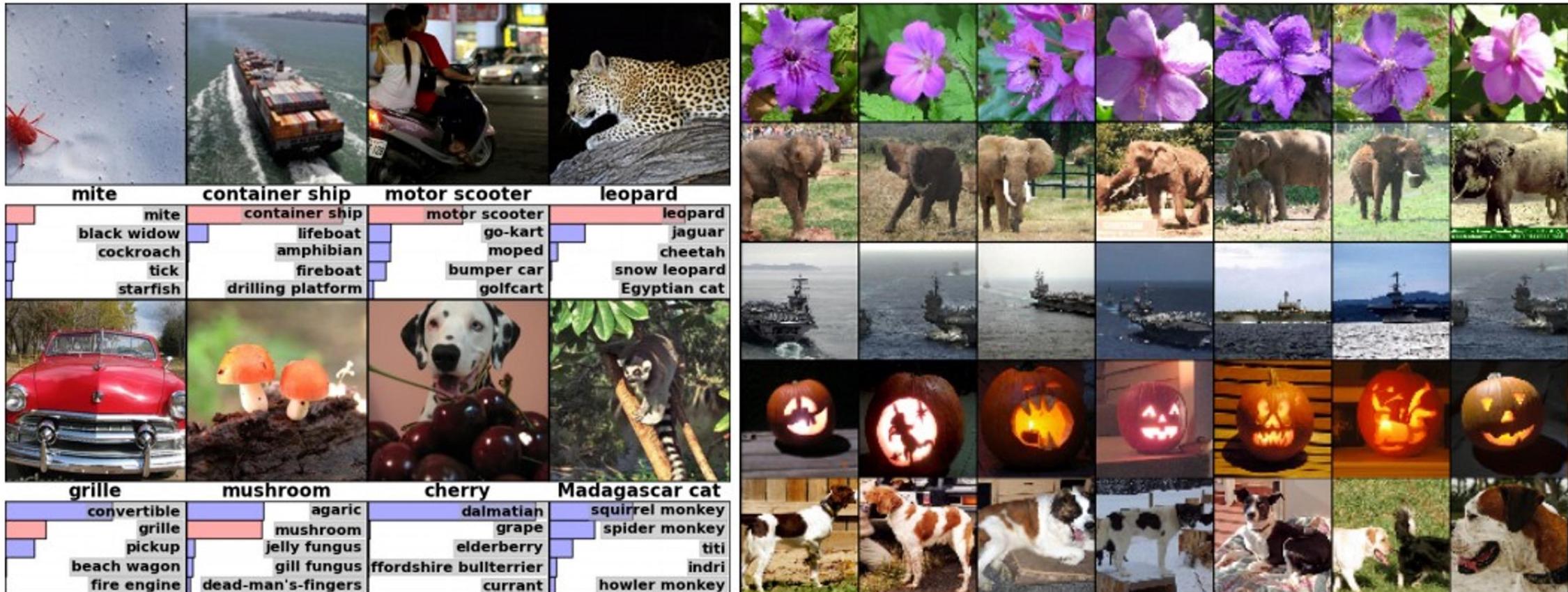
Real world output



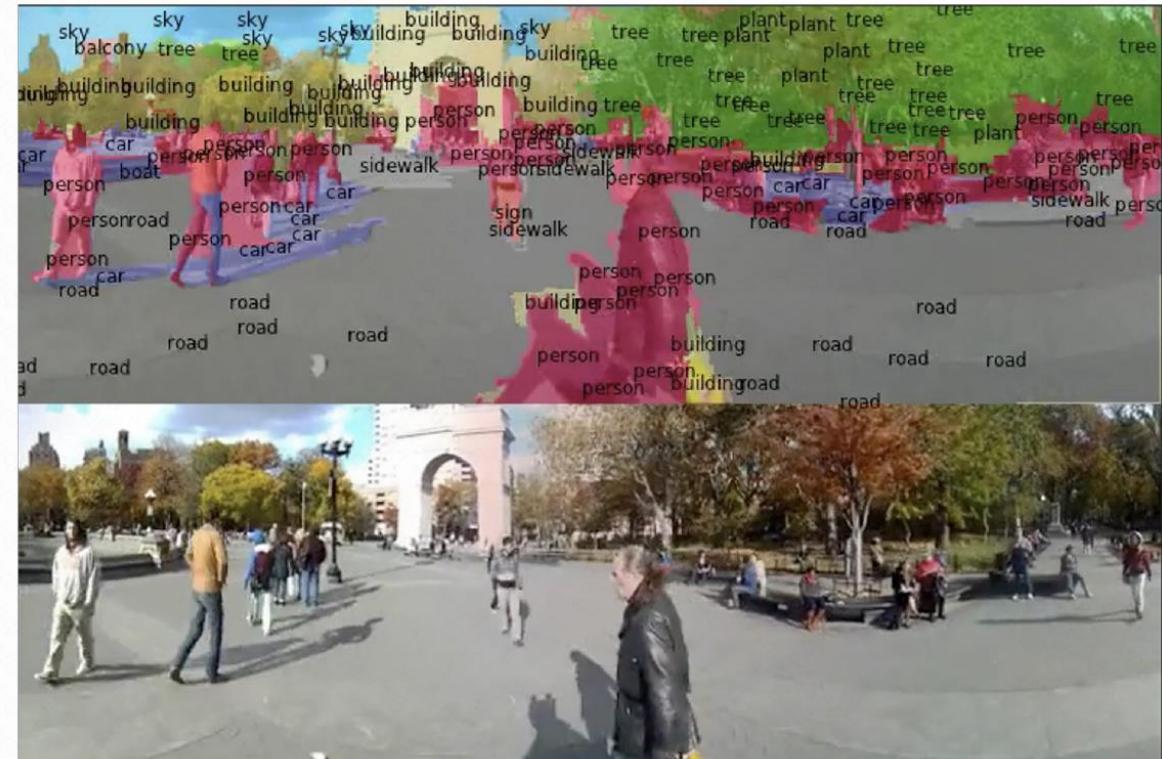
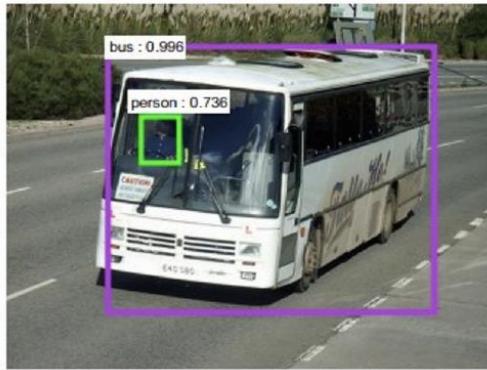
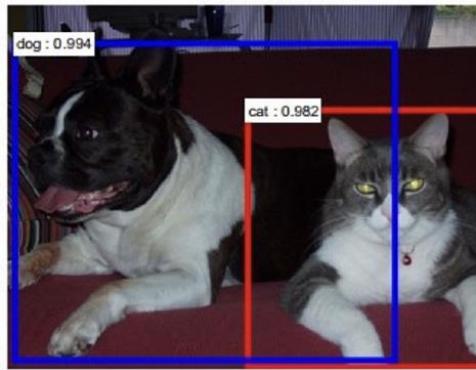
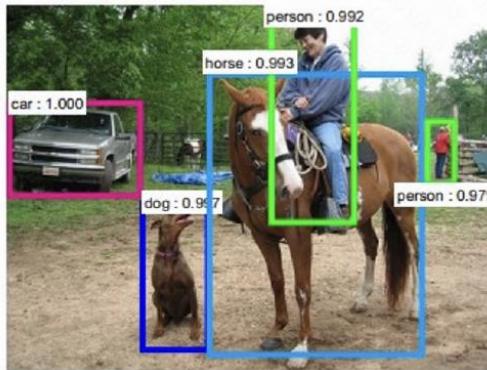
■ 翻译



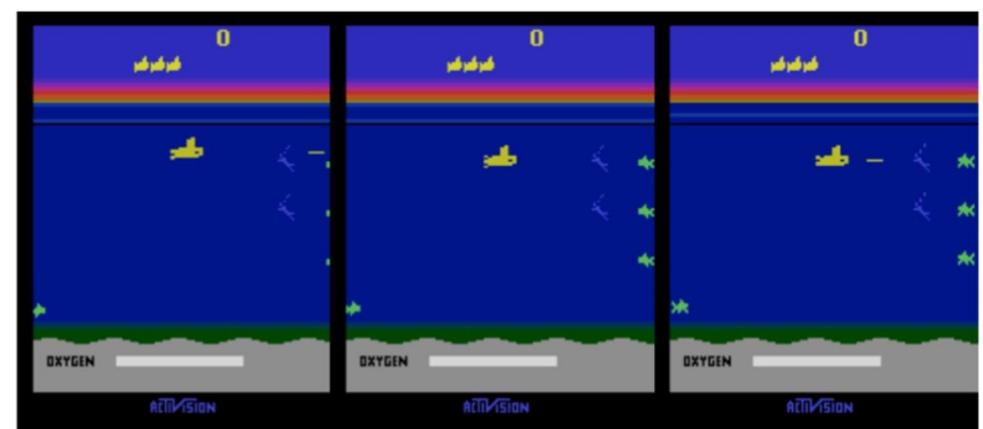
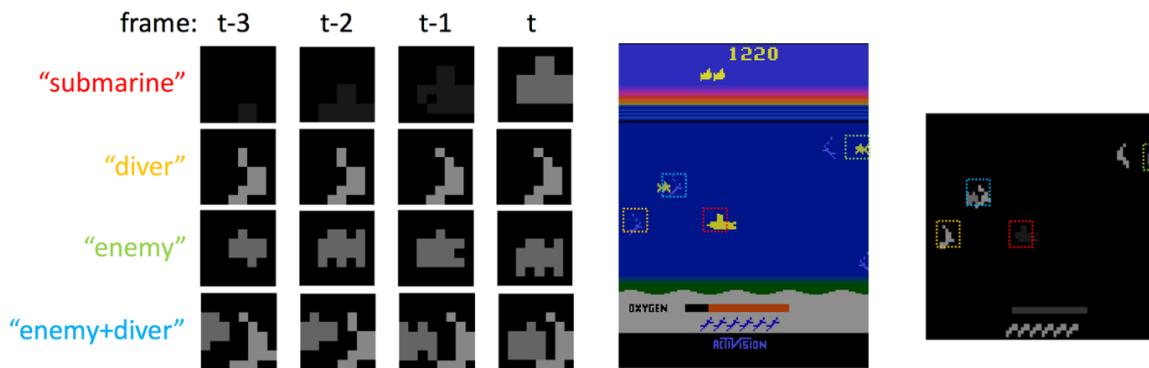
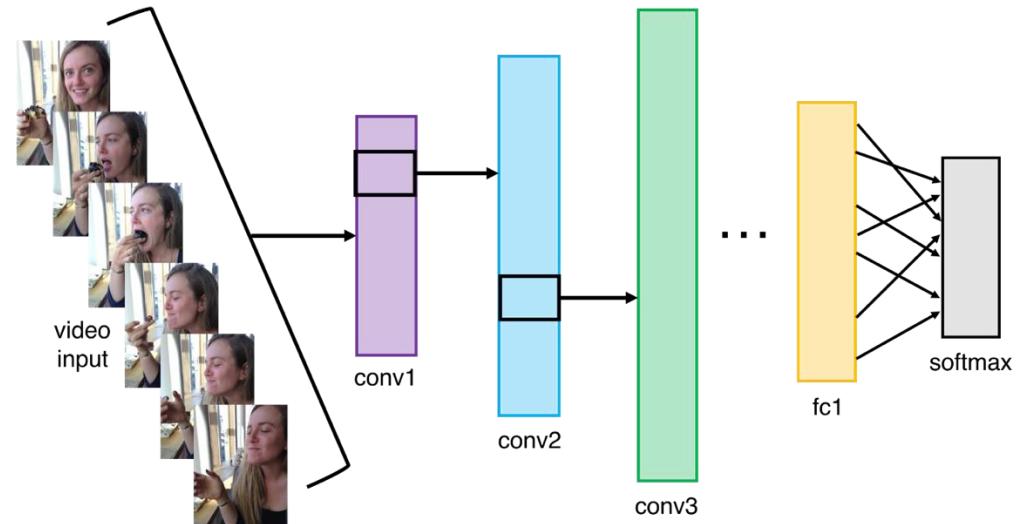
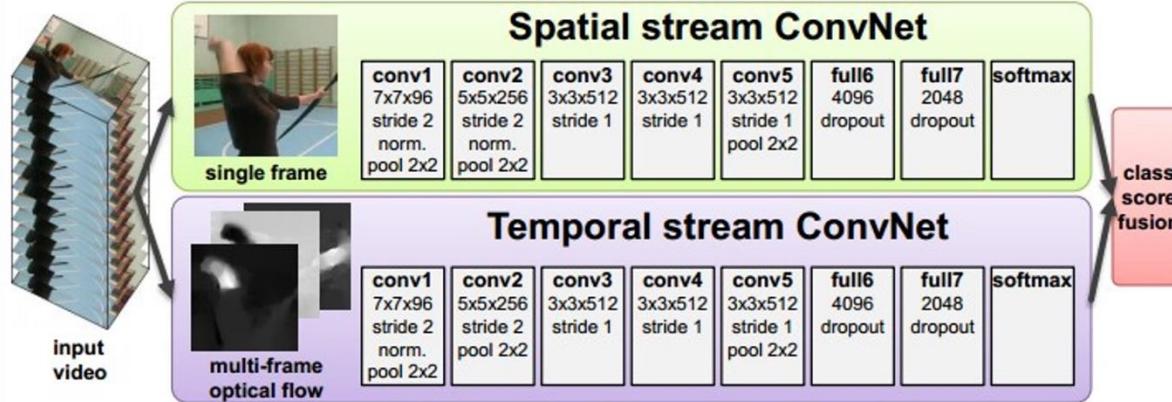
■ 图像分类 & 图像检索



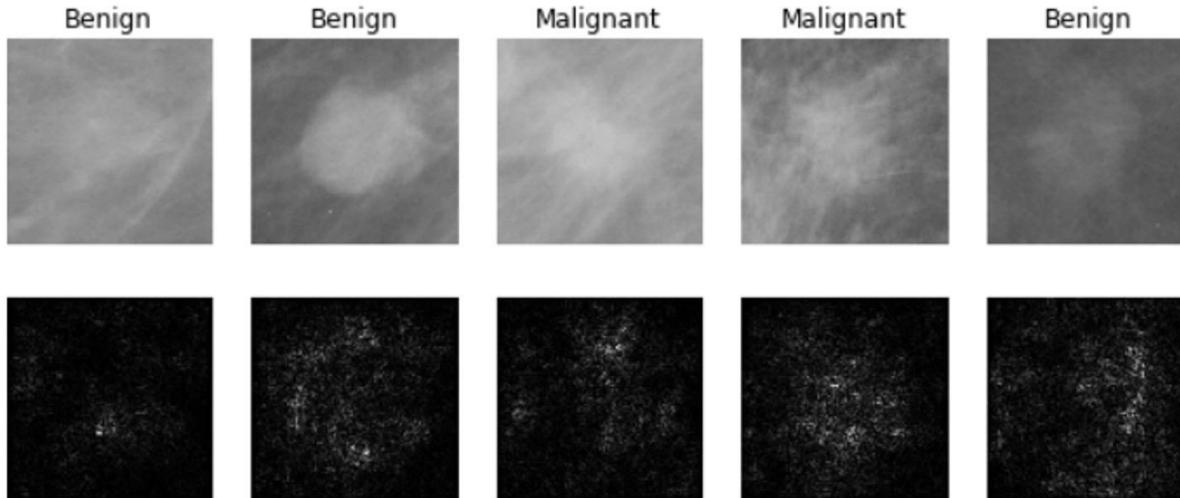
■ 目标检测 & 图像分割



■ 视频分类 & 行为识别 & 强化学习



■ 医疗图像识别 & 星系识别 & 生物识别



■ 图像描述



*A white teddy bear
sitting in the grass*



*A man in a baseball
uniform throwing a ball*



*A woman is holding
a cat in her hand*



*A man riding a wave
on top of a surfboard*

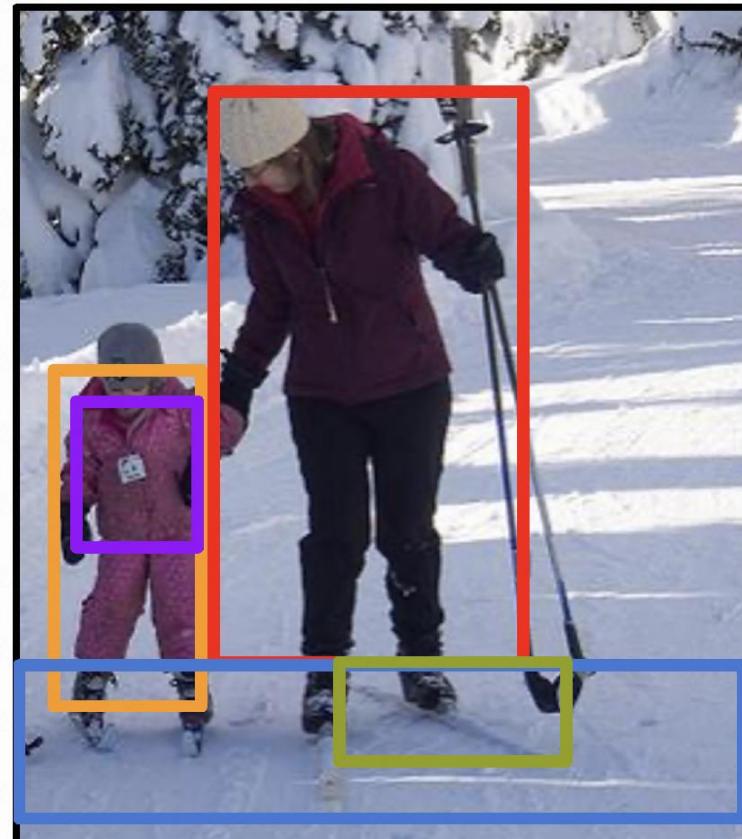


*A cat sitting on a
suitcase on the floor*

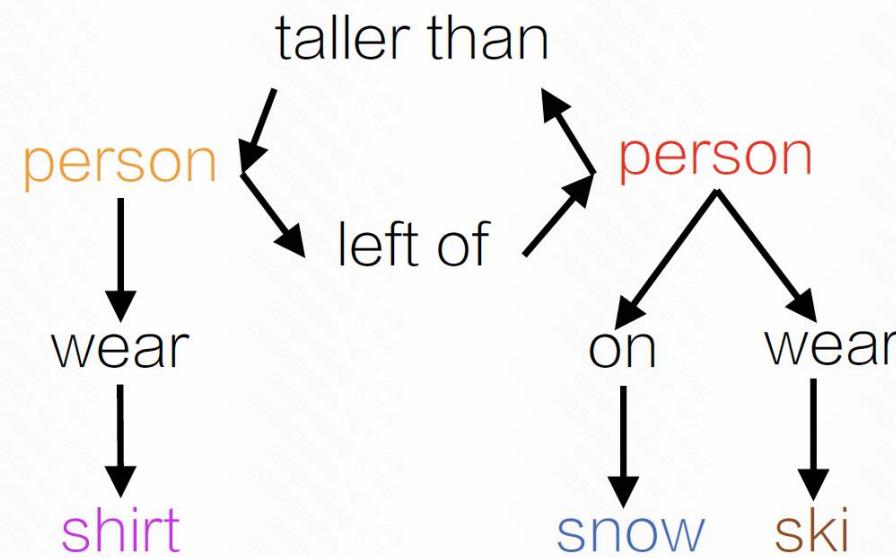


*A woman standing on a
beach holding a surfboard*

■ 视觉关系检测



Results:
spatial, comparative, asymmetrical,
verb, prepositional



■ 风格迁移



■ 图像生成

TEXT PROMPT

an armchair in the shape of an avocado. an armchair imitating an avocado.

AI-GENERATED IMAGES

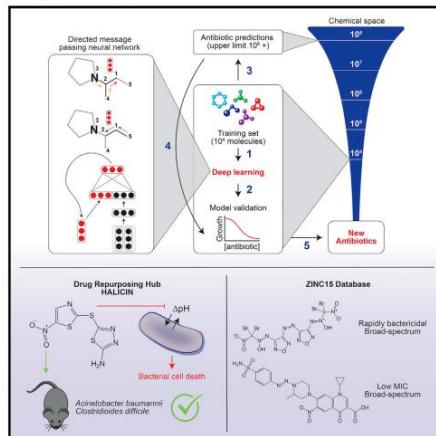


■ AI4Science

Cell

A Deep Learning Approach to Antibiotic Discovery

Graphical Abstract



Authors

Jonathan M. Stokes, Kevin Yang, Kyle Swanson, ..., Tommi S. Jaakkola, Regina Barzilay, James J. Collins

Correspondence

regina@csail.mit.edu (R.B.), jmjc@mit.edu (J.J.C.)

In Brief

A trained deep neural network predicts antibiotic activity in molecules that are structurally different from known antibiotics, among which Halicin exhibits efficacy against broad-spectrum bacterial infections in mice.

nature
biomedical engineering

ARTICLES
<https://doi.org/10.1038/s41551-018-0195-0>

Prediction of cardiovascular risk factors from retinal fundus photographs via deep learning

Ryan Poplin^{1,4}, Avinash V. Varadarajan^{1,4}, Katy Blumer¹, Yun Liu¹, Michael V. McConnell^{2,3}, Greg S. Corrado¹, Lily Peng^{1,4*} and Dale R. Webster^{1,4}

Traditionally, medical discoveries are made by observing associations, making hypotheses from them and then designing and running experiments to test the hypotheses. However, with medical images, observing and quantifying associations can often be difficult because of the wide variety of features, patterns, colours, values and shapes that are present in real data. Here, we show that deep learning can extract new knowledge from retinal fundus images. Using deep-learning models trained on data from 284,335 patients and validated on two independent datasets of 12,026 and 999 patients, we predicted cardiovascular risk factors not previously thought to be present or quantifiable in retinal images, such as age (mean absolute error within 3.26 years), gender (area under the receiver operating characteristic curve (AUC)=0.97), smoking status (AUC=0.71), systolic blood pressure (mean absolute error within 11.23 mmHg) and major adverse cardiac events (AUC=0.70). We also show that the trained deep-learning models used anatomical features, such as the optic disc or blood vessels, to generate each prediction.

nature

Explore content ▾ About the journal ▾ Publish with us ▾

nature > articles > article

Article | Open access | Published: 15 July 2021

Highly accurate protein structure prediction with AlphaFold

John Jumper  Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, Alex Bridgland, Clemens Meyer, Simon A. Kohl, Andrew J. Ballard, Andrew Cowie, Bernardino Romera-Paredes, Stanislav Nikolov, Rishabh Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman, Ellen Clancy, Michal Zielinski, ... Demis Hassabis 

+ Show authors

Nature 596, 583–589 (2021) | Cite this article

1.42m Accesses | 12k Citations | 3493 Altmetric | Metrics



Can you write me a report analyzing this chest X-ray?



Findings:

- Devices: None.
- Lungs: No pneumothorax. No substantial pleural effusion. Lungs appear clear.
- Cardiomediastinal: Normal heart size. Mediastinal contours within normal limits.
- Other: No acute skeletal abnormality.

Impression:

No active disease seen in chest.

Enter a question here

■ 不同任务间的共性

- 输入和输出之间具有非常复杂的关系（浅层学习难以处理）
- 但输入和输出都遵循一定的规则（可以通过深度学习处理）

“A Kazakh man on a horse holding a bird of prey”

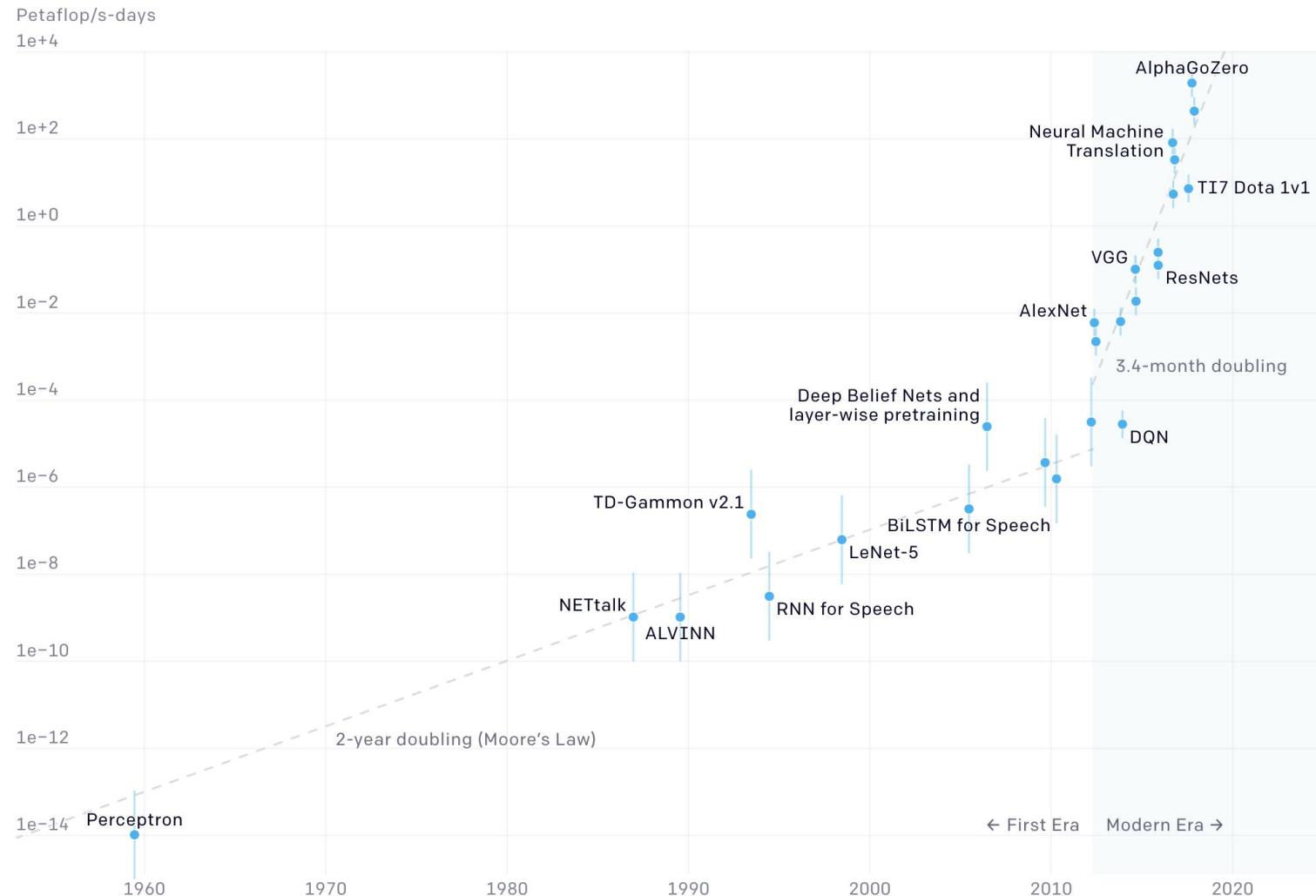
语言遵循语法规则



自然图像也有规则

深度学习算力需求

Two Distinct Eras of Compute Usage in Training AI Systems



深度学习算力需求

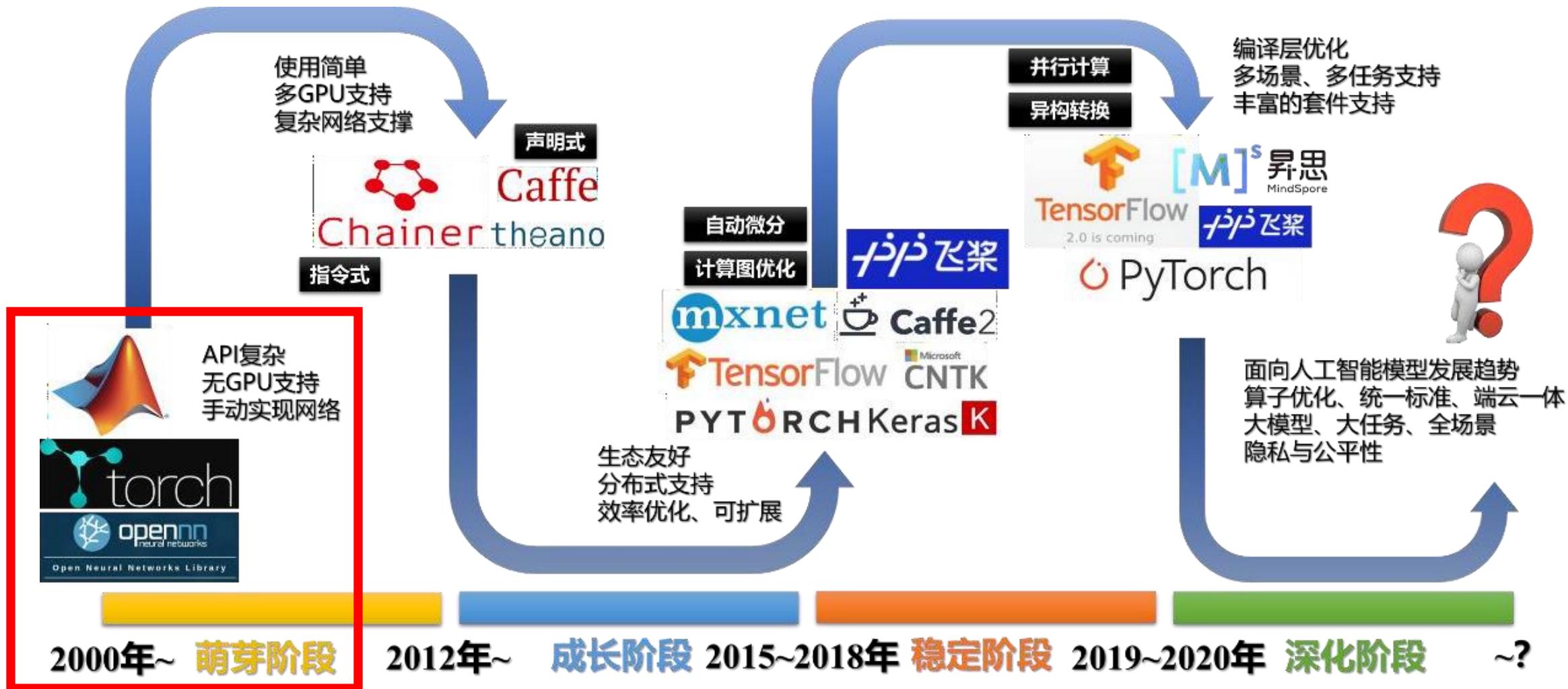
AlexNet to AlphaGo Zero: A 300,000x Increase in Compute (Log Scale)



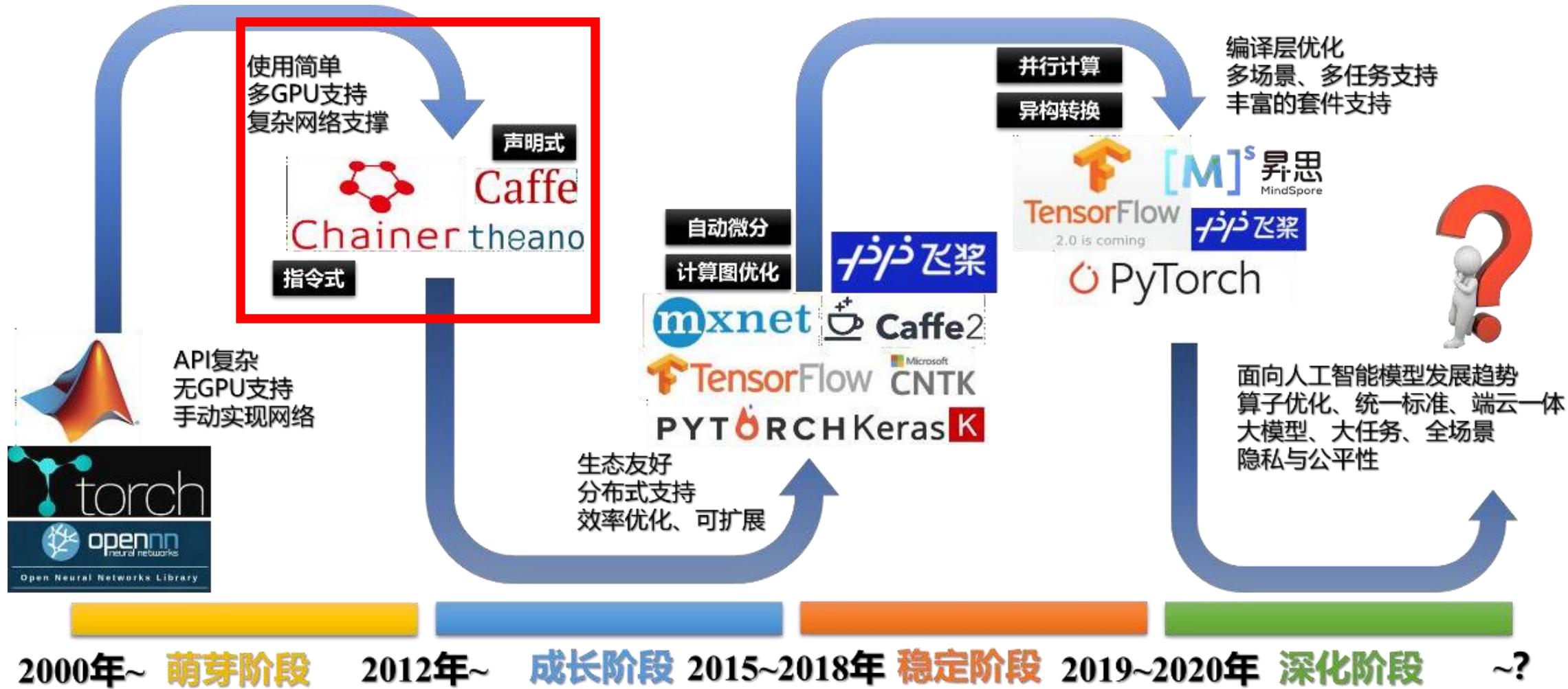
大 纲

- 课程概览
- 深度学习历史
- 深度学习应用
- 深度学习平台
- 课程后续安排

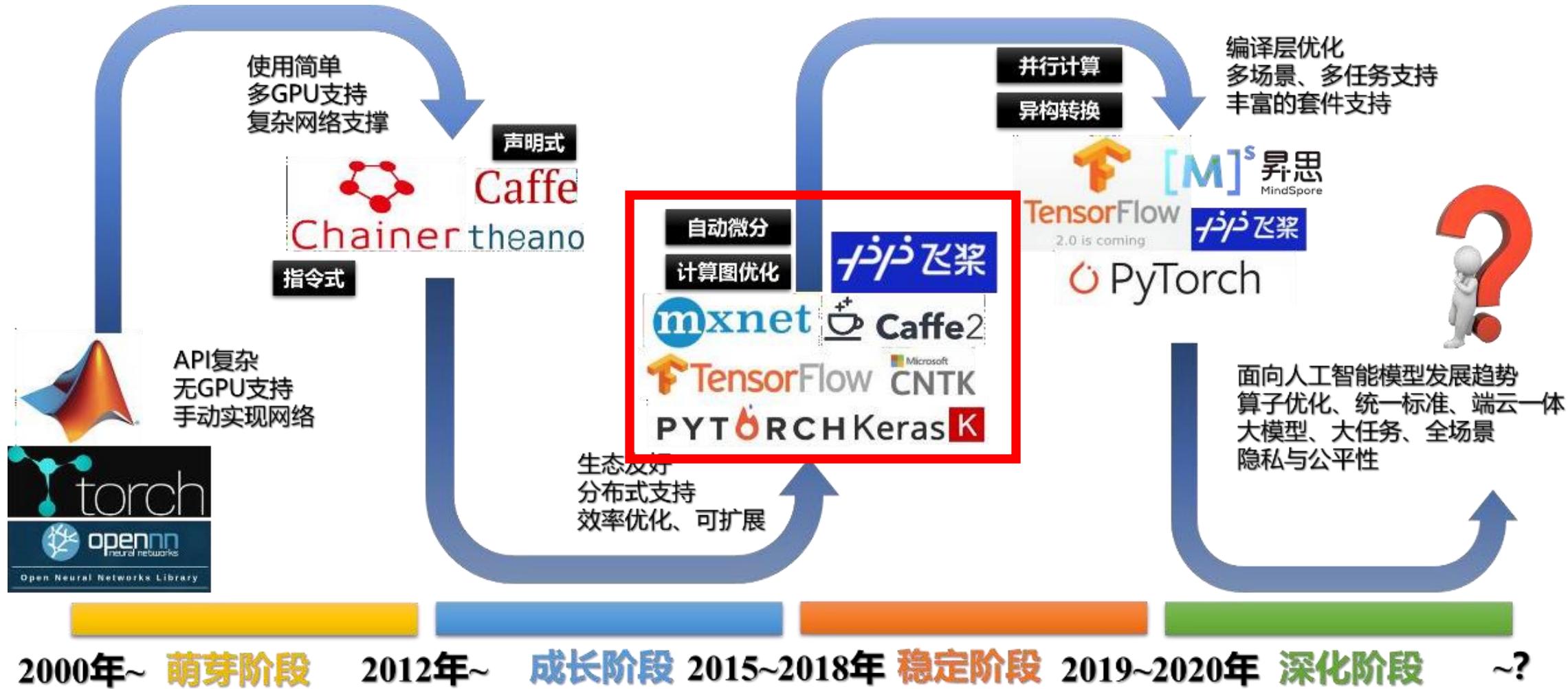
深度学习平台



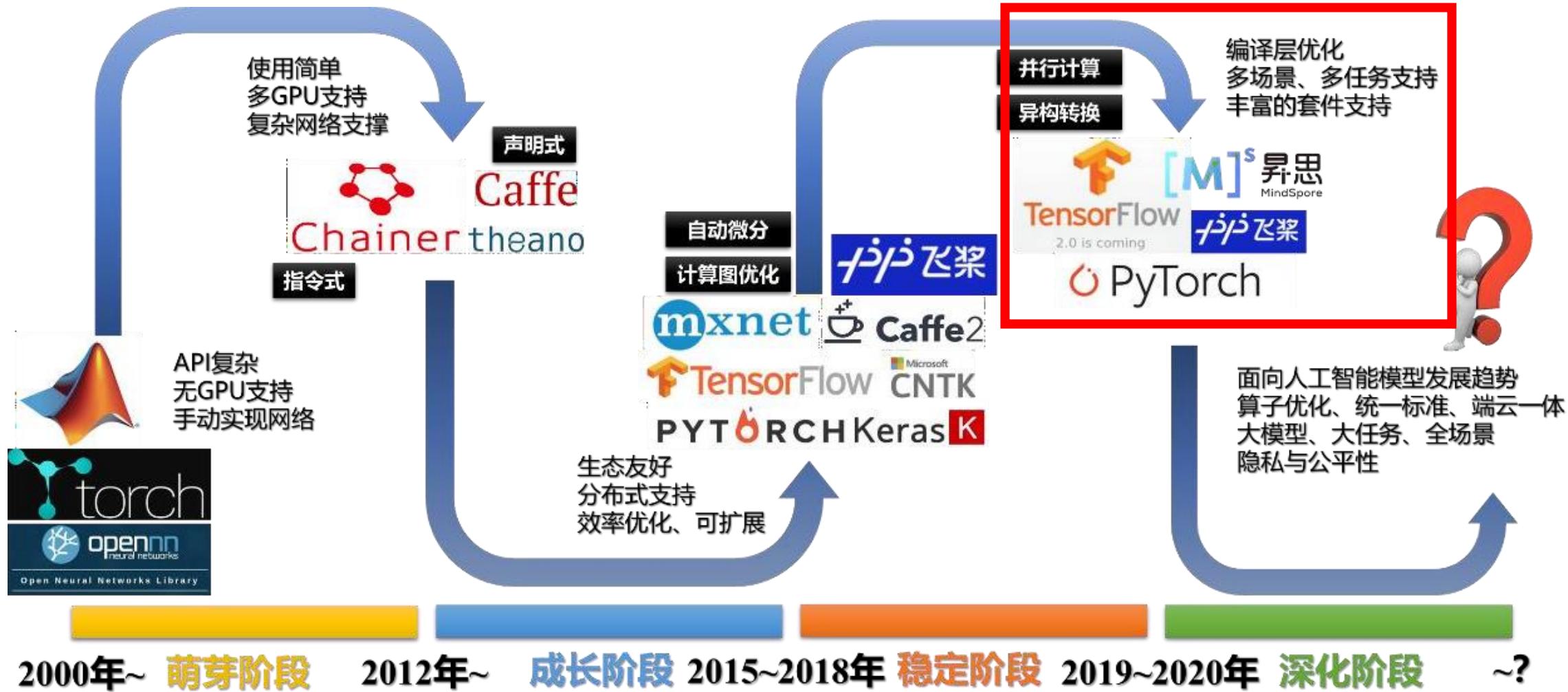
深度学习平台



深度学习平台



深度学习平台

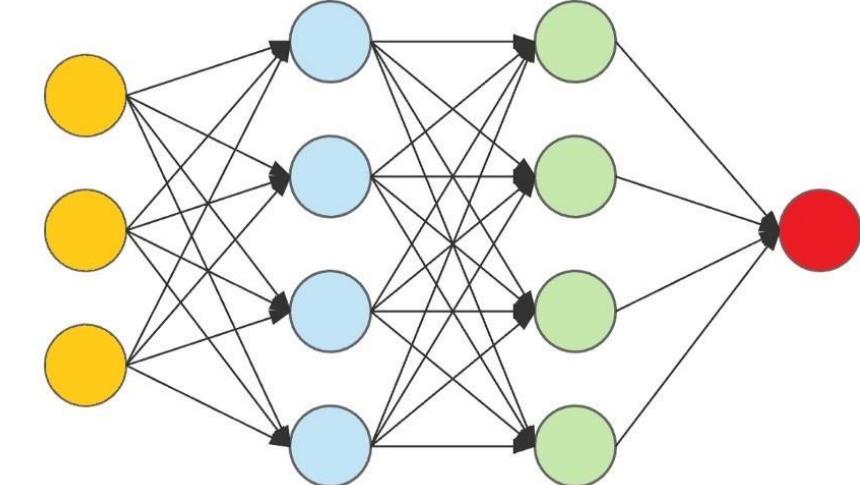
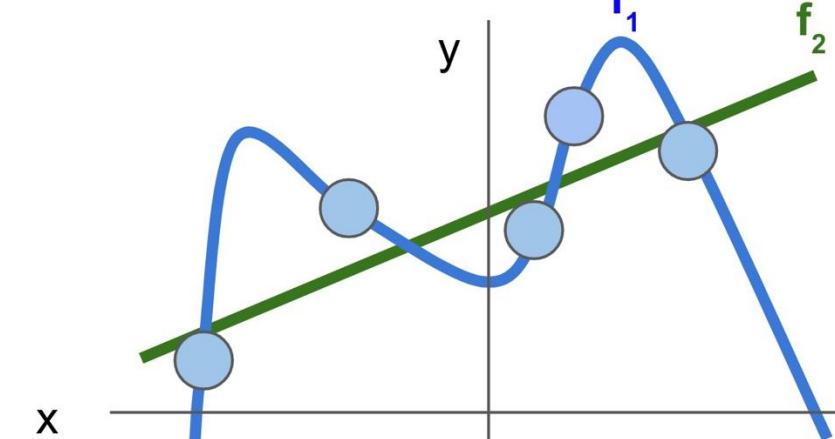
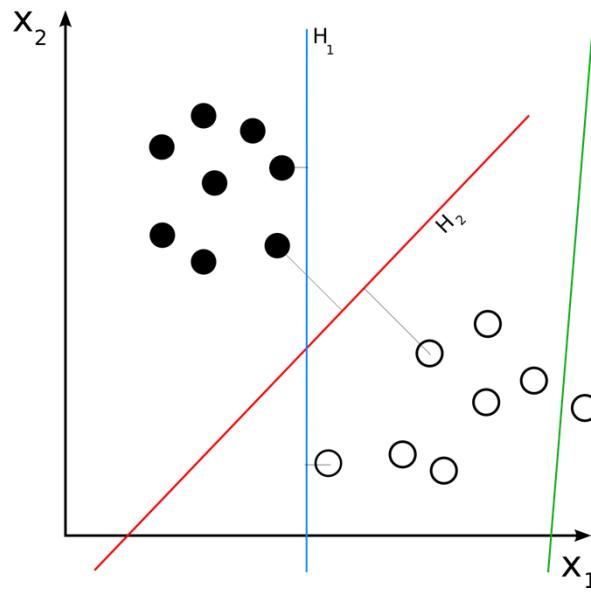


大 纲

- 课程概览
- 深度学习历史
- 深度学习应用
- 深度学习平台
- 课程后续安排

■ 深度学习基础

■ 线性分类器、正则化与优化、神经网络



■ 基础视觉任务

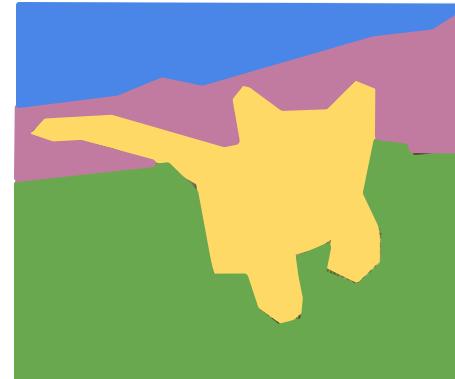
■ 图像理解（分类、检测、分割）、视频理解，网络可视化

分类



CAT

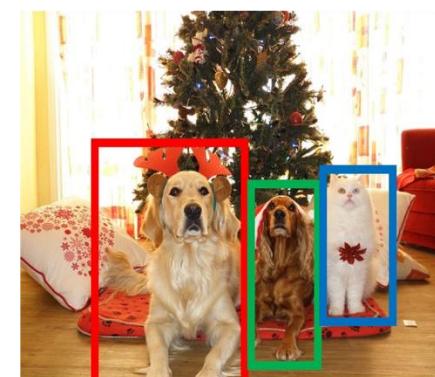
语义分割



GRASS, CAT, TREE,
SKY

No spatial extent

目标检测



DOG, DOG, CAT

实例分割

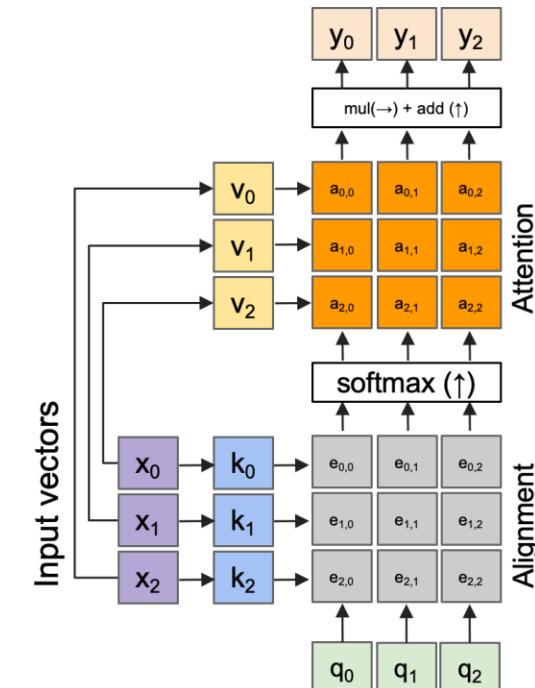
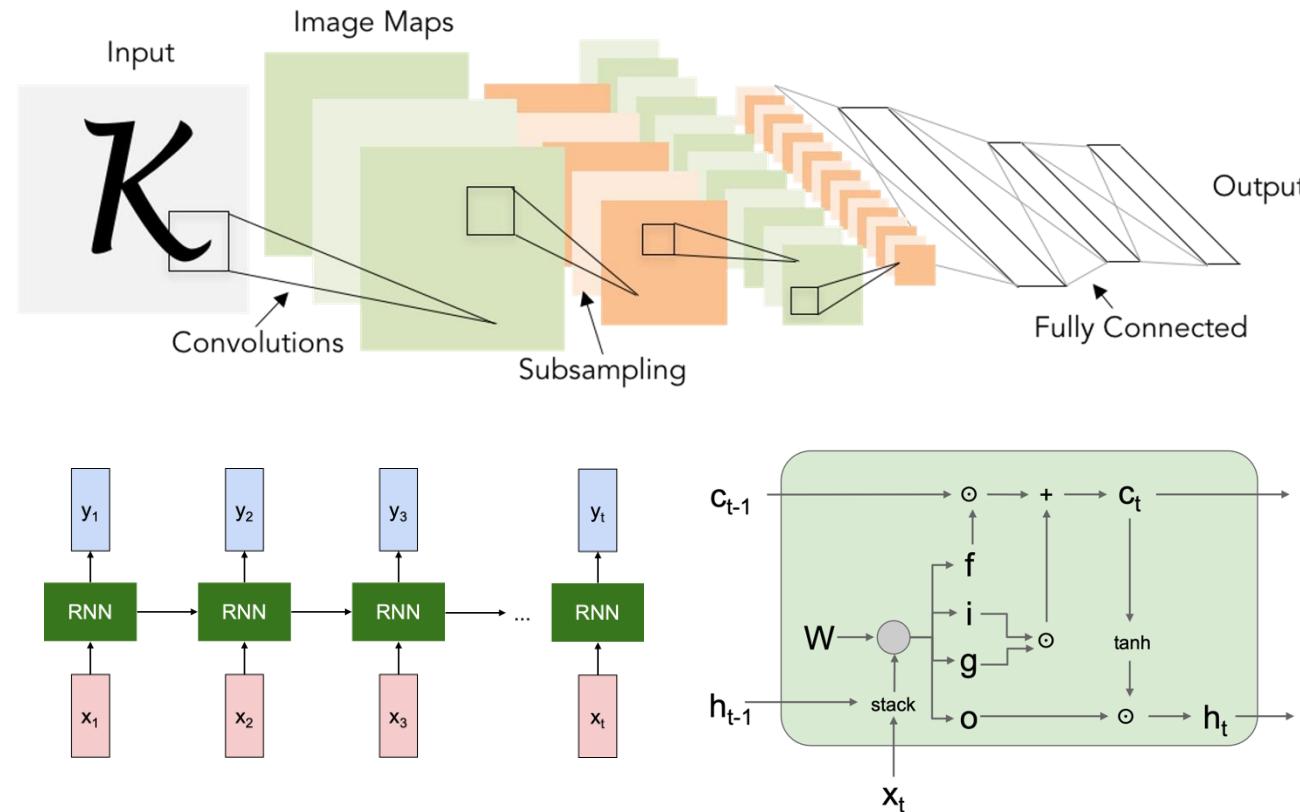


DOG, DOG, CAT

Multiple Object

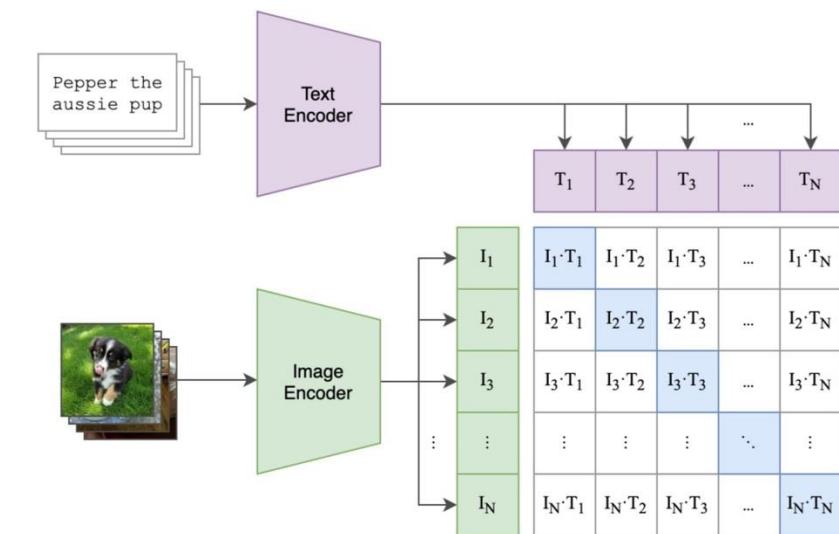
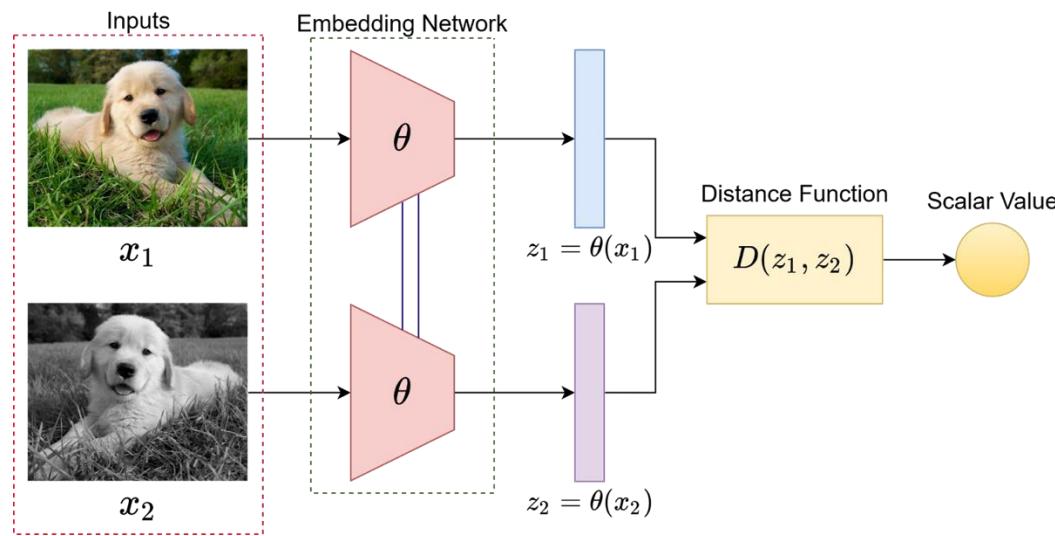
■ 基础神经网络

■ 卷积神经网络、循环神经网络、Transformer



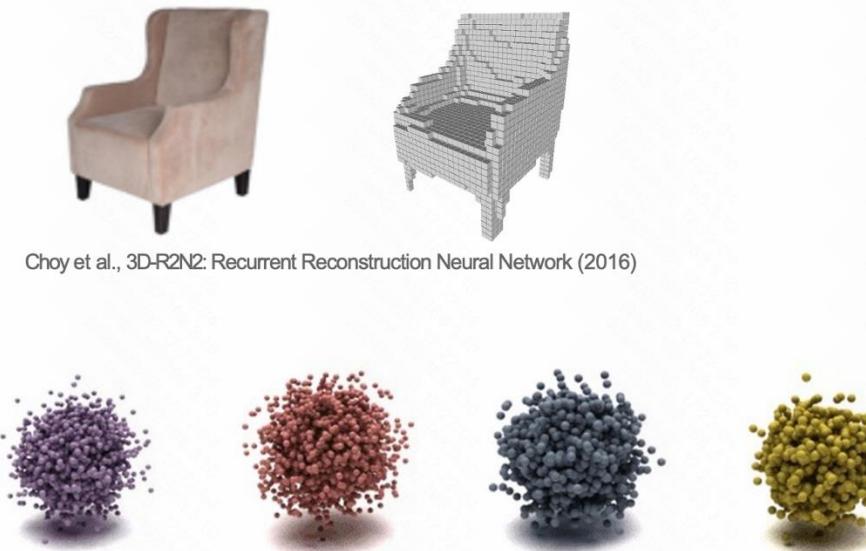
■ 更多深度学习任务

■ 自监督学习、多模态模型



■ 更多深度学习任务

■ 生成模型、3D模型



Zhou et al., 3D Shape Generation and Completion through Point-Voxel Diffusion (2021)

