

Deep Learning for Computer Vision (Discriminative way)

Andrii Liubonko
2025

About me

- Theoretical Physics background
 - Ukraine & Germany
- Samsung (7 years)
 - mainly Computer Vision projects
- Grammarly (3 years)
 - NLP related projects
- EPAM (current)
 - mix of LLMs & NLP & Computer Vision



reach me on
Slack if you have
any questions



Course structure

Module 1. Basic image processing - Oles Dobosevych

Module 2. Learning Discriminative Models - Andrii Liubonko

Module 3. Image Segmentation, Image Correspondence -
Maksym Davydov

Module 4. Metric Learning, Representation Learning and Context
Importance, Tracking - Igor Krashenyi

Module 5. Learning Generative Models - Andrii Liubonko

Logistics (this module)

6 lectures

homework:

- assignment 1 : tutorials notebooks 50%
- assignment 2 : mini-project 50%

deadline date: *TBD*

course repo:

<https://github.com/lyubonko/ucu2025cv>

Overview of the module

Lecture I Intro, big picture

Lecture II Essential architectures (CNN)

Lecture III Essential architectures (Transformers)

Lecture IV CV Foundational models in depth (DINO, SAM)

Lecture V Multimodal Foundational models (GEMMA, APIs)

Lecture VI Object Detection

Goals of the Course

- Get big picture of the deep learning for CV
- Working knowledge of essential elements/blocks of **Convolutional Neural Networks [CNNs]**
- Working knowledge of essential elements/blocks of **Transformers** and their use in CV
- Get flavor of CV-related **Foundational models**
- Get deeper with one particular problem (**Object Detection**)

Intro

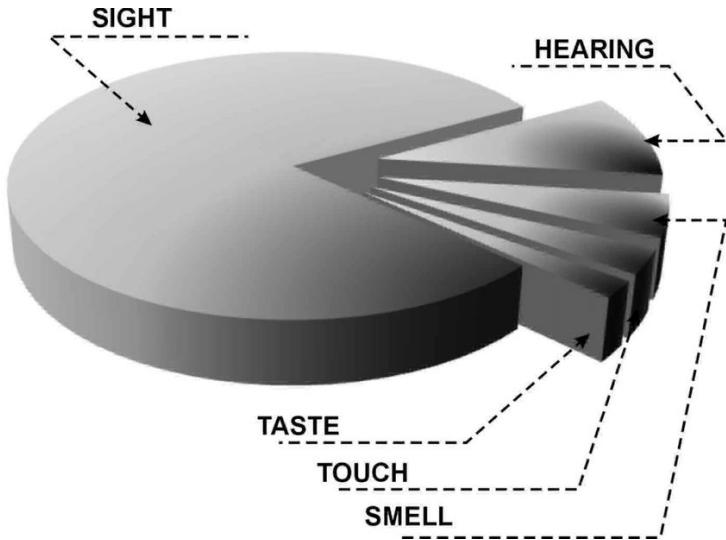
The **State of AI Report** analyses the most interesting developments in AI. We aim to trigger an informed conversation about the state of AI and its implication for the future. The Report is produced by AI investor [Nathan Benaich](#) and [Air Street Capital](#).

The screenshot shows a dark-themed website with a blue header bar. The header contains navigation links: Introduction | Research | Industry | Politics | Safety | Predictions on the left, and #stateofai | 5 on the right. Below the header, the page title is 'Definitions'. The content lists several AI-related terms with their definitions:

- Artificial intelligence (AI):** a broad discipline with the goal of creating intelligent machines, as opposed to the natural intelligence that is demonstrated by humans and animals.
- Artificial general intelligence (AGI):** a term used to describe future machines that could match and then exceed the full range of human cognitive ability across all economically valuable tasks.
- AI Agent:** an AI-powered system that can take actions in an environment. For example, an LLM that has access to a suite of tools and has to decide which one to use in order to accomplish a task that it has been prompted to do.
- AI Safety:** a field that studies and attempts to mitigate the risks (minor to catastrophic) which future AI could pose to humanity.
- Computer vision (CV):** the ability of a program to analyse and understand images and video. This item is circled in red.
- Deep learning (DL):** an approach to AI inspired by how neurons in the brain recognise complex patterns in data. The “deep” refers to the many layers of neurons in today’s models that help to learn rich representations of data to achieve better performance gains.
- Diffusion:** An algorithm that iteratively denoises an artificially corrupted signal in order to generate new, high-quality outputs. In recent years it has been at the forefront of image generation and protein design.
- Generative AI:** A family of AI systems that are capable of generating new content (e.g. text, images, audio, or 3D assets) based on ‘prompts’.
- Graphics Processing Unit (GPU):** a semiconductor processing unit that enables a large number calculations to be computed in parallel. Historically this was required for rendering computer graphics. Since 2012 GPUs have adapted for training DL models, which also require a large number of parallel calculations.

In the bottom right corner of the page, the text 'stateof.ai 2024' is displayed. At the very bottom, there are navigation icons for back, forward, and search, along with a Google Slides icon.

Intro



Computer Vision aims

- to extract useful information from visual input
- to generate novel visual content

(Discriminative Models)
(Generative Models)

→ to extract useful information from visual input **(Discriminative Models)**



- indoor/outdoor? [image classification]
- Where are the objects? [object detection]
- How far is the object ? [depth estimation]
- What people are doing? [activity recognition]
- Is the state of the environment normal? [anomaly detection]
- ...

→ to generate novel visual content

(Generative Models)



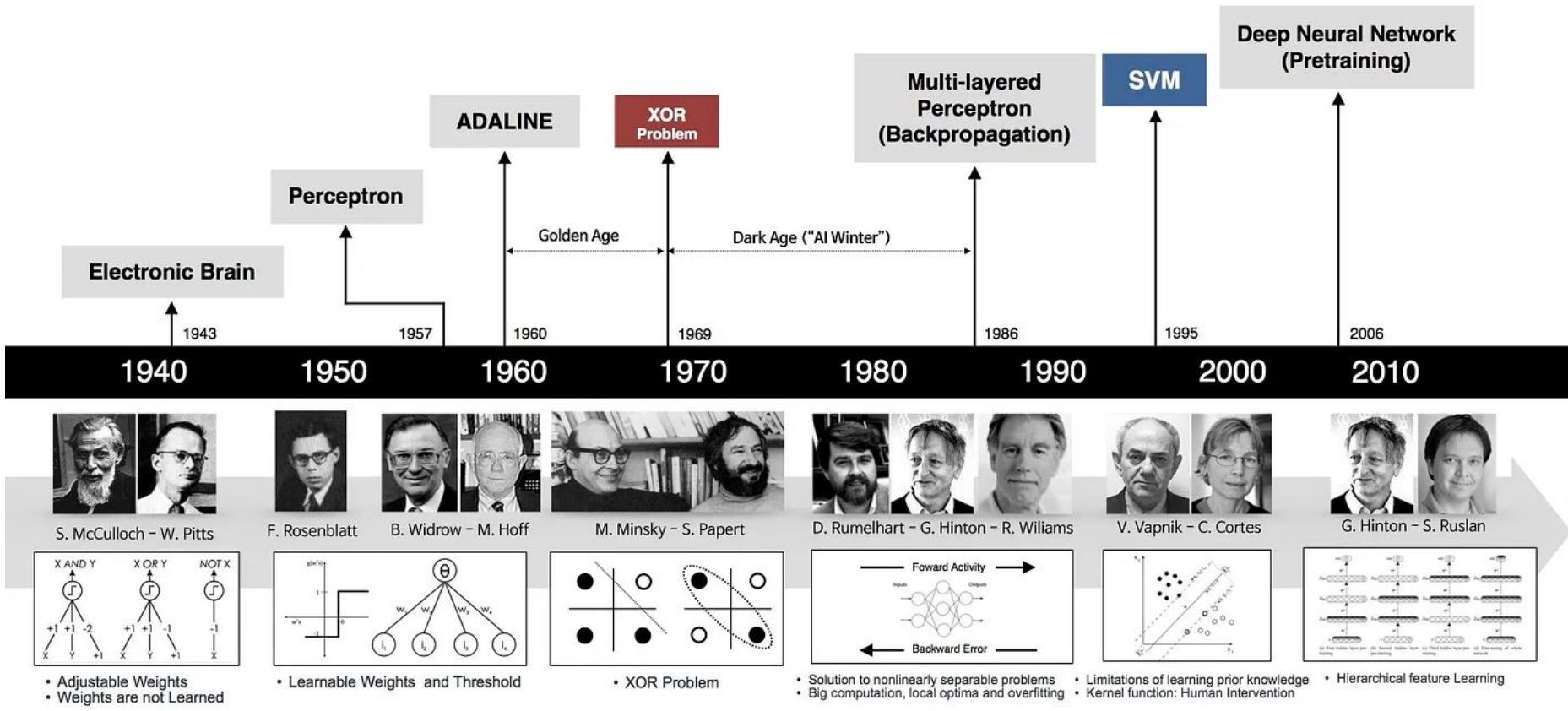
2011.09055



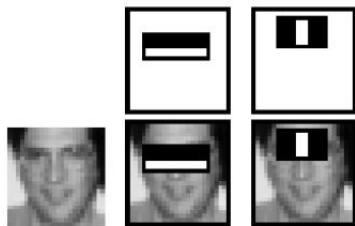
SORA (openAI)
VEO (Google)

A (very) Brief History

History till 2012



Intro



Era of Human-Crafter Features

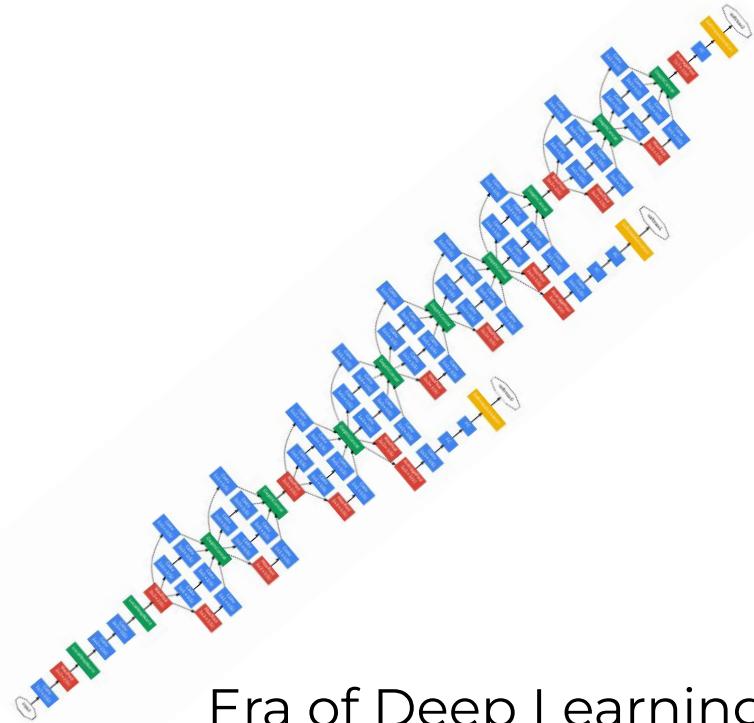
Era of Deep Learning



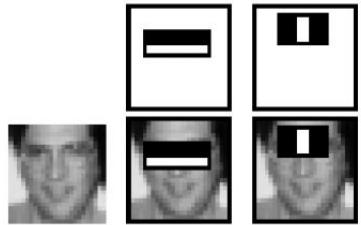
1986
BackProp

1998
LeNet

2012
AlexNet



Intro

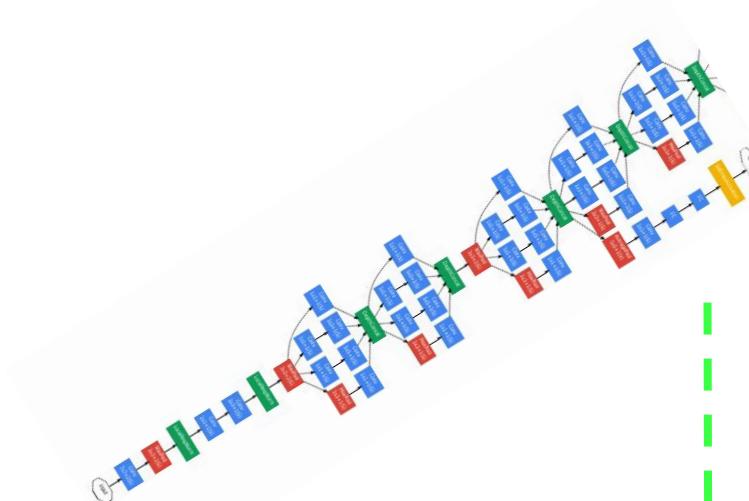


Era of
Human-Crafter
Features

2012
AlexNet

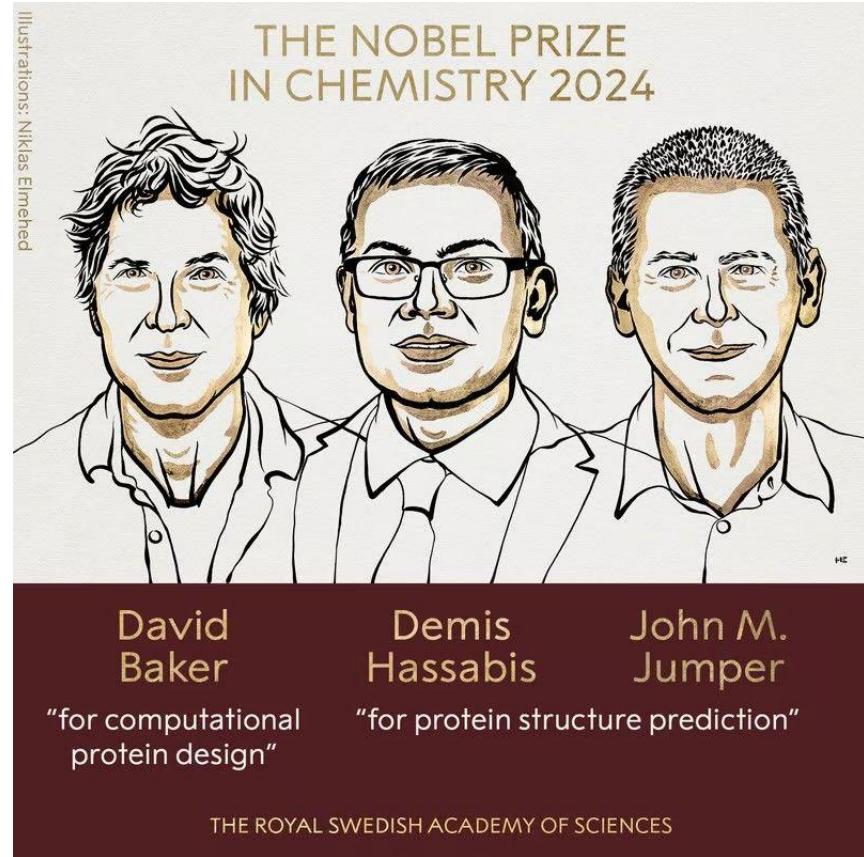
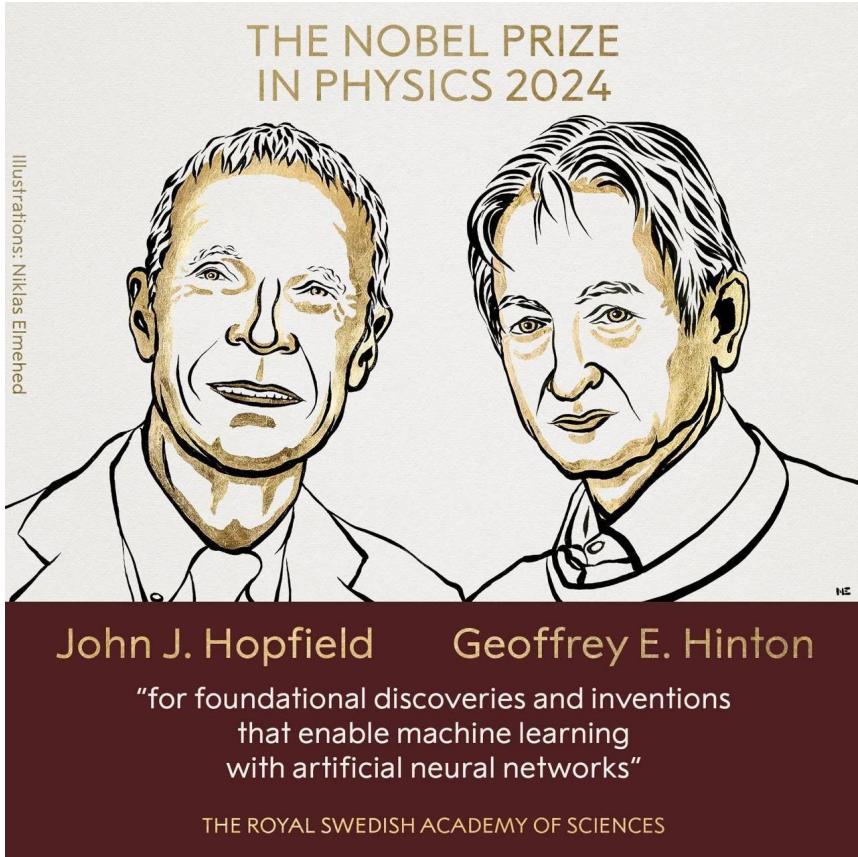
Era of
Deep Learning

2022
ChatGPT

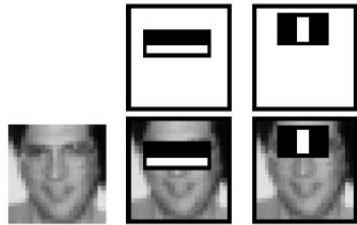


Era of
LLMs

Intro



Intro

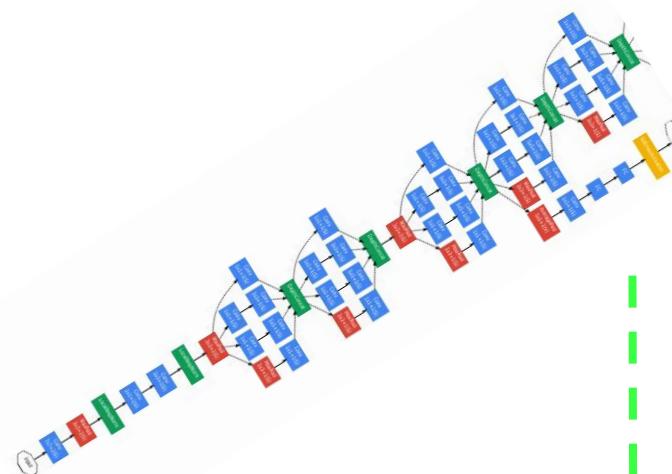


Era of
Human-Crafter
Features

2012
AlexNet

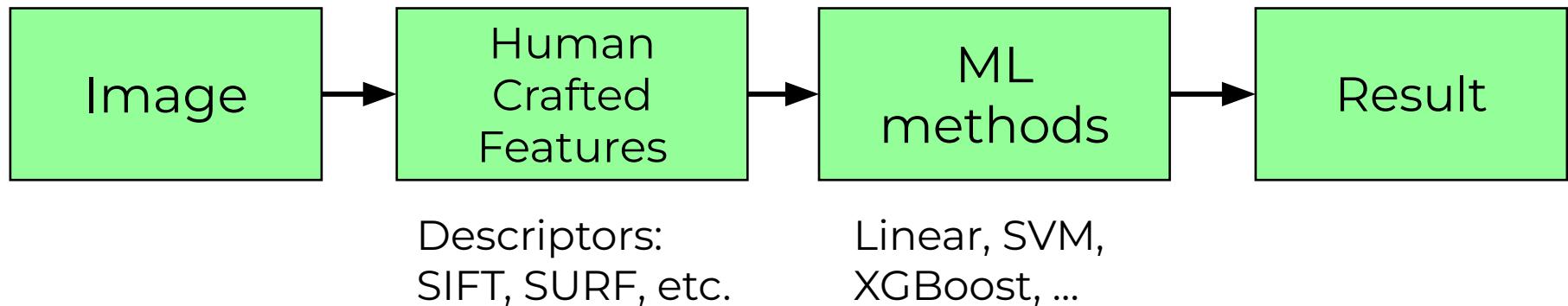
Era of
Deep Learning

2022
ChatGPT

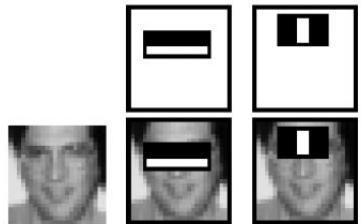


Era of
LLMs

Intro



Intro

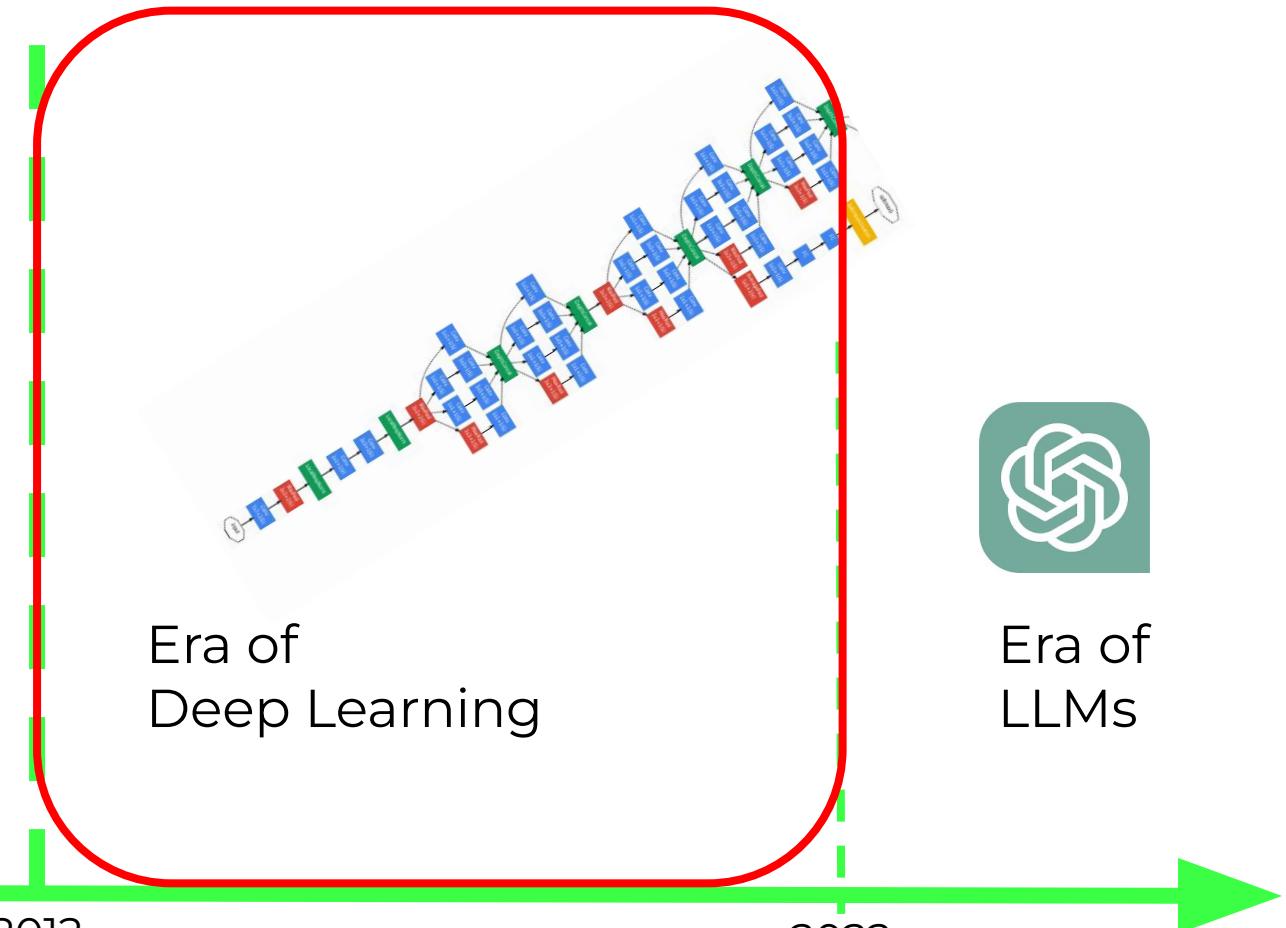


Era of
Human-Crafter
Features

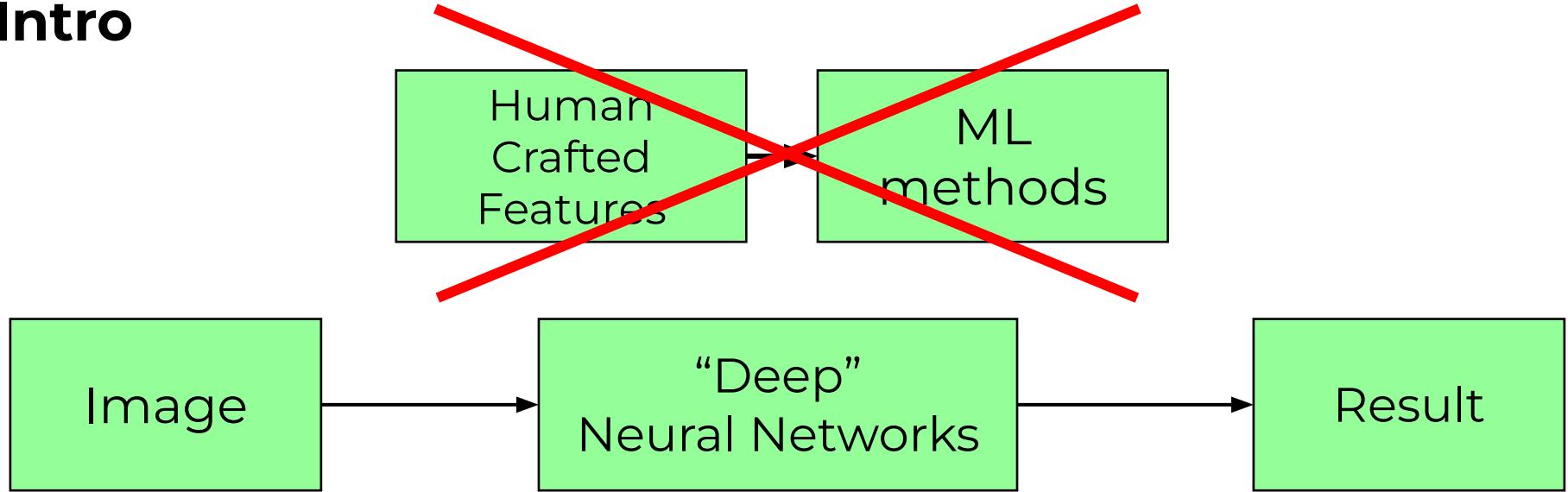
2012
AlexNet

Era of
Deep Learning

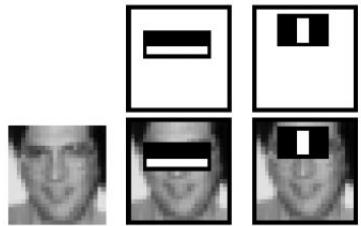
2022
ChatGPT



Intro



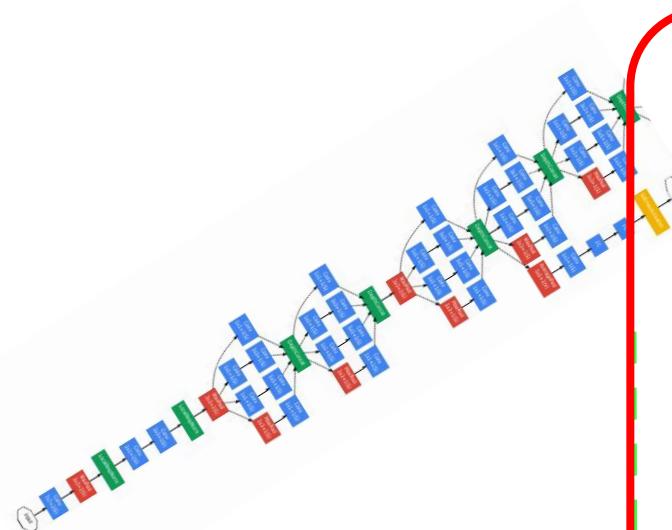
Intro



Era of
Human-Crafter
Features

Era of
Deep Learning

2012
AlexNet

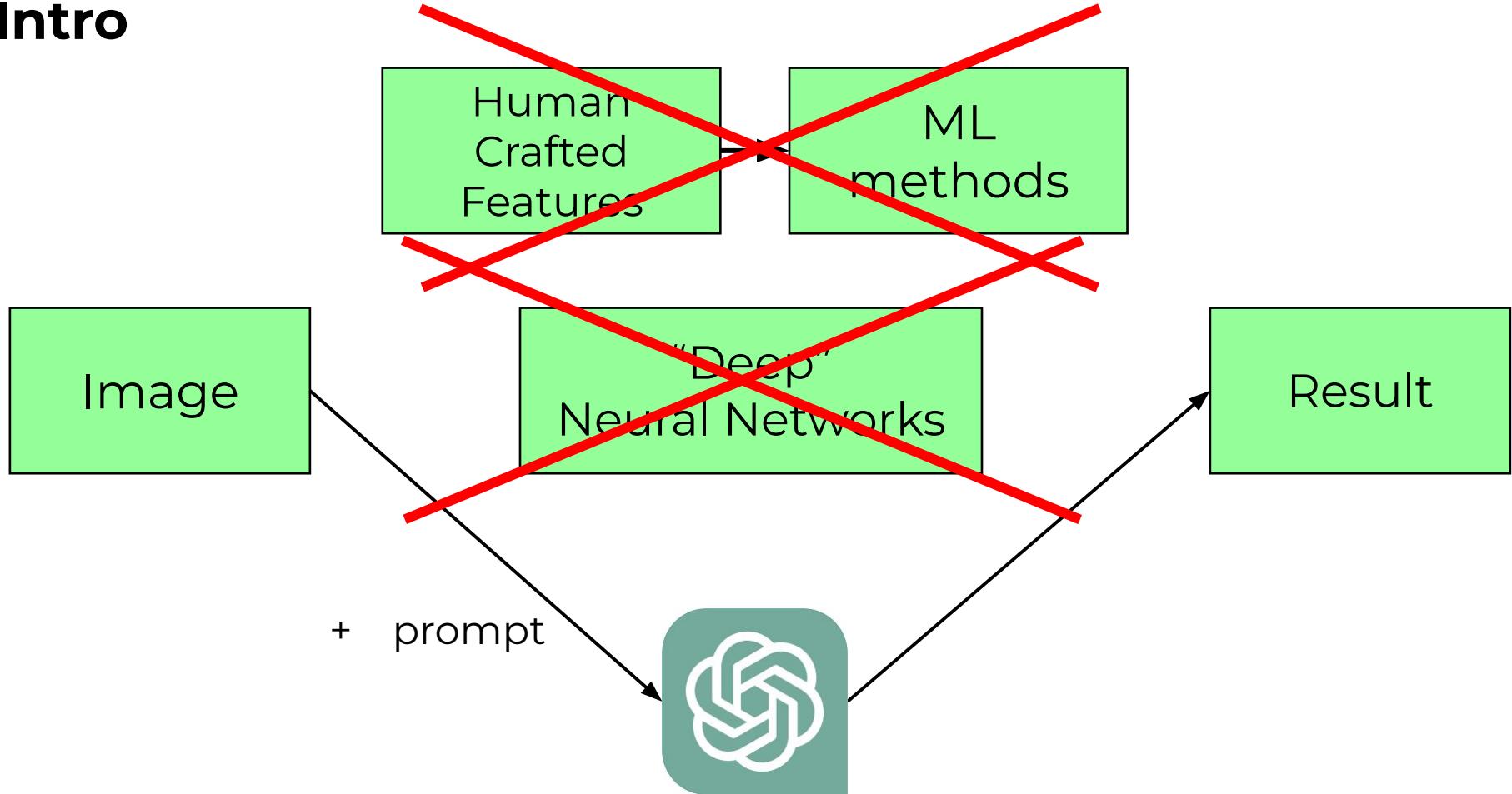


2022
ChatGPT



Era of
LLMs

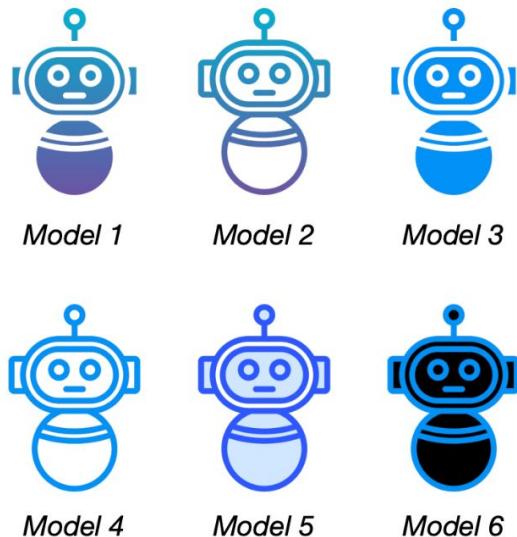
Intro



Intro

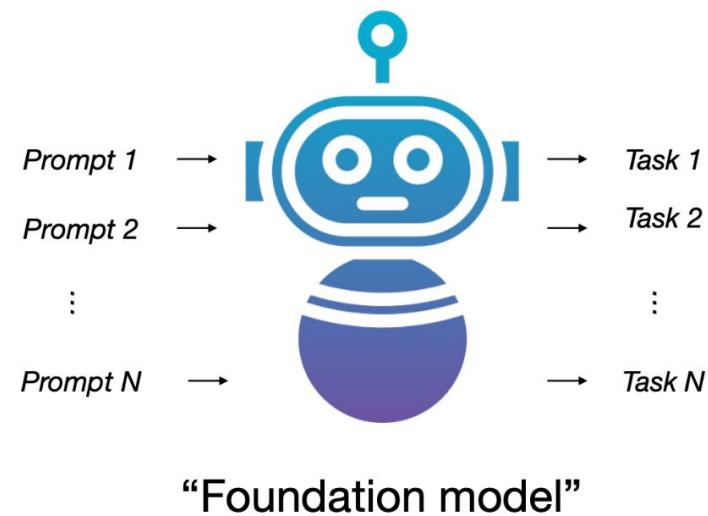
traditional ML

Old days: one model for one purpose



prompt-based ML

Now: one model for multiple purposes



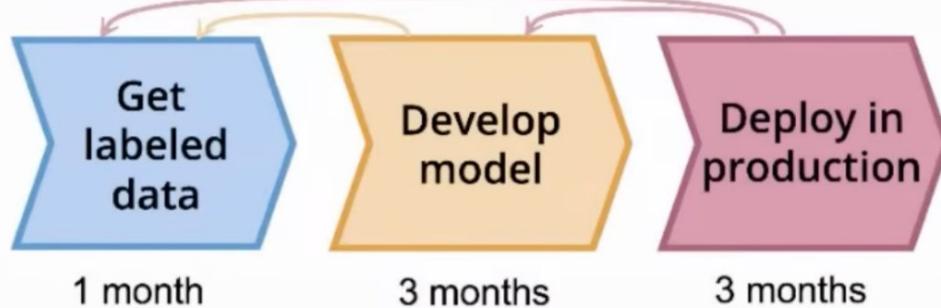
Intro



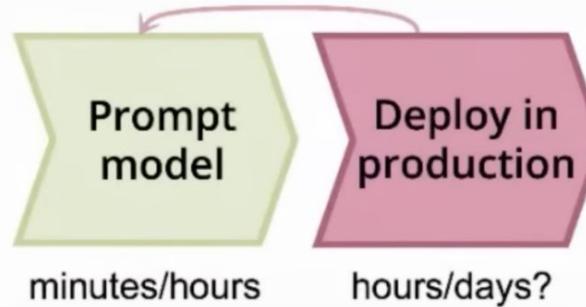
Vision for Future ML Workflow: Iterate Faster

LANDING AI

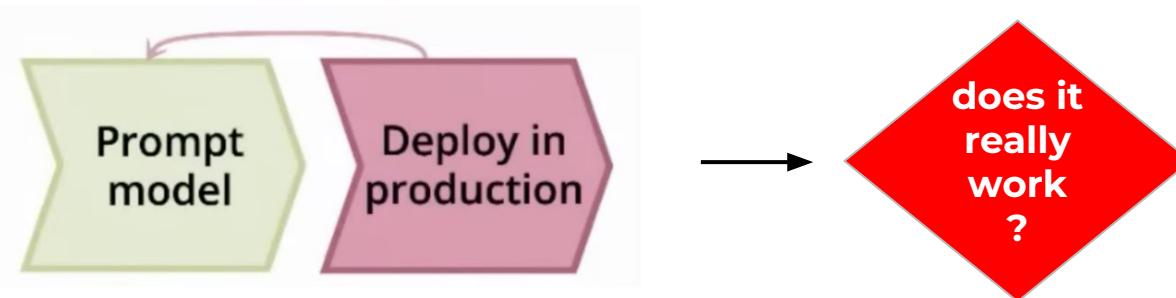
Traditional ML



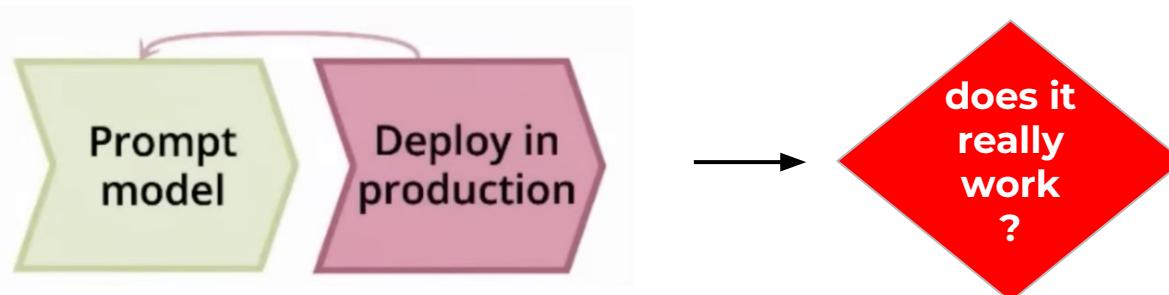
Prompt-based ML



prompt-based ML

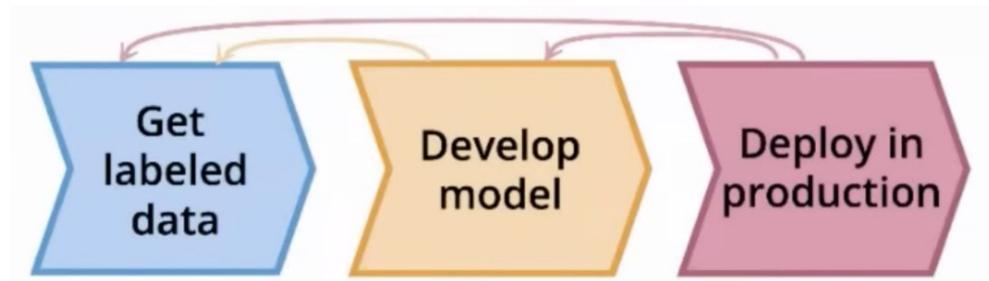


prompt-based ML

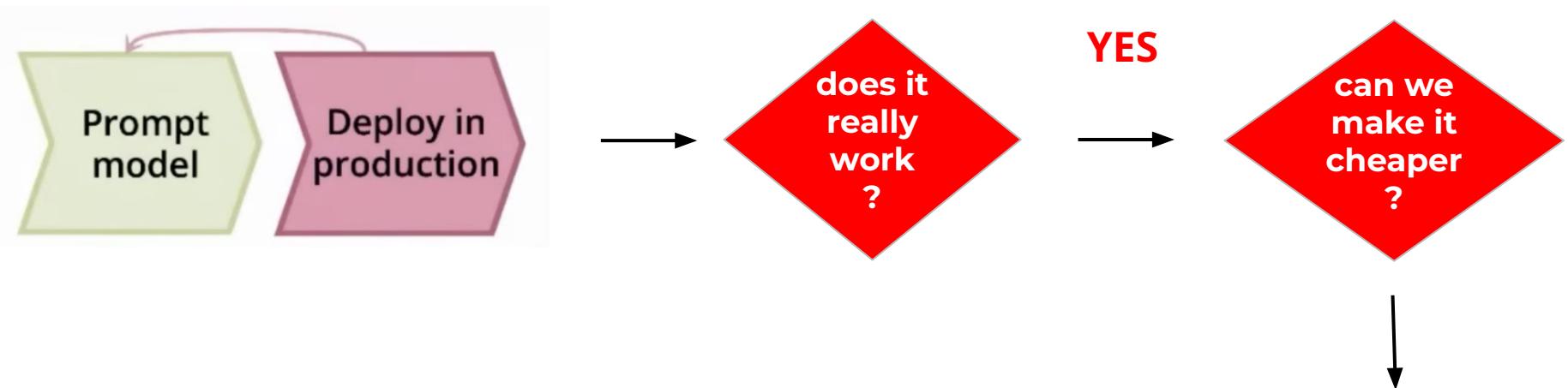


NO

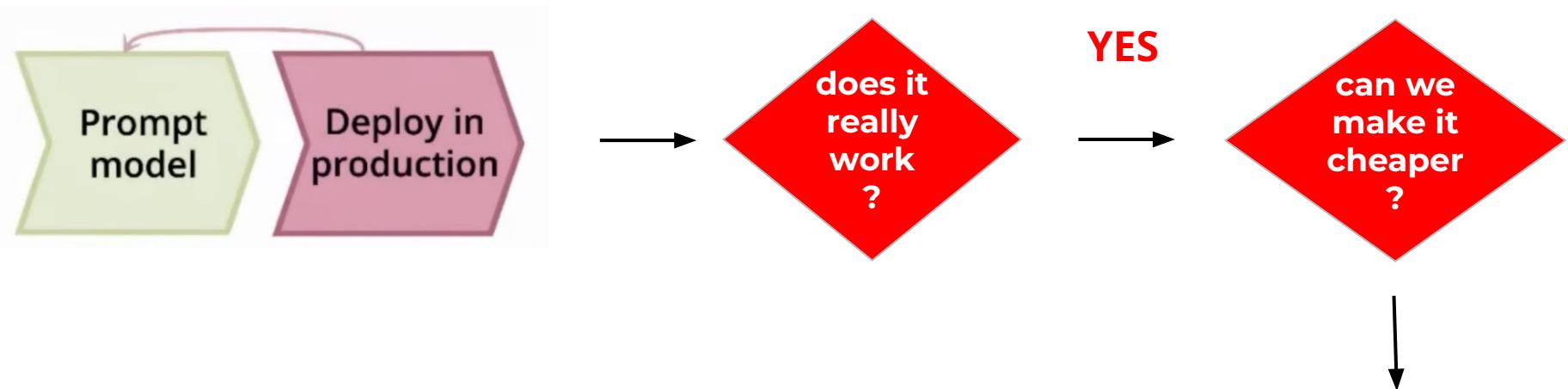
“traditional” ML/DL



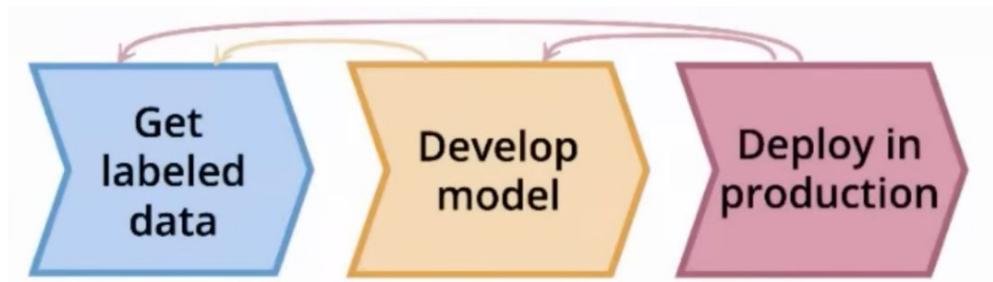
prompt-based ML



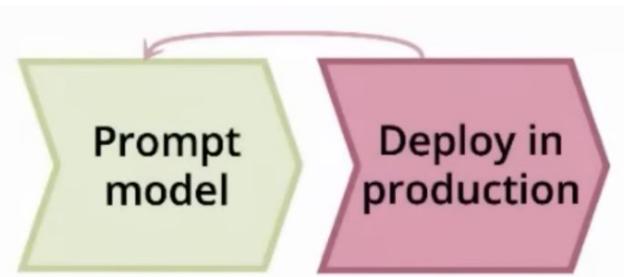
prompt-based ML



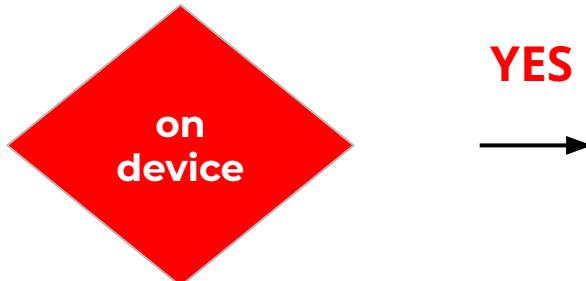
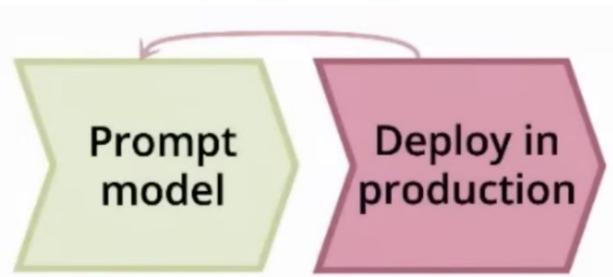
“traditional” ML/DL



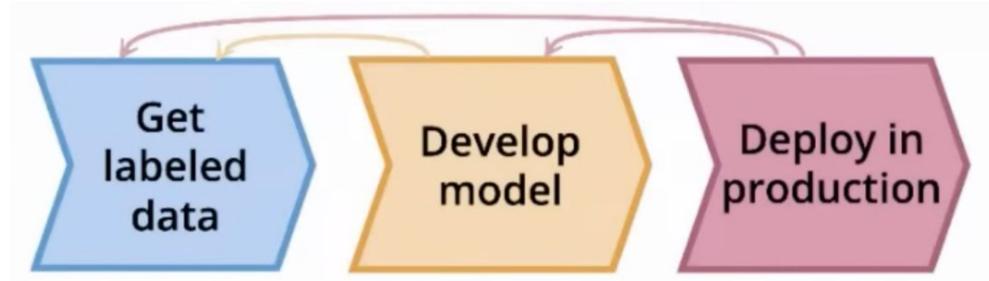
prompt-based ML



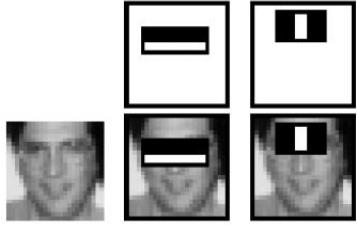
prompt-based ML



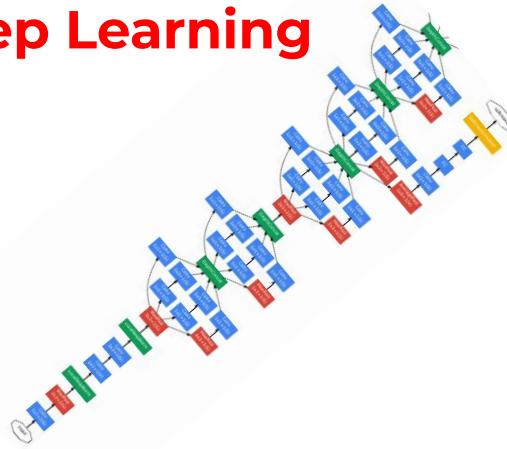
"traditional" ML/DL



Era of Human-Crafter Features



Era of Deep Learning



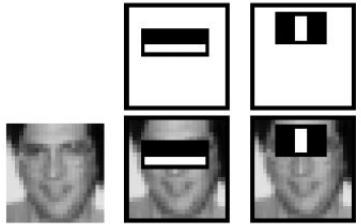
2012
AlexNet

Era of LLMs



2022
ChatGPT

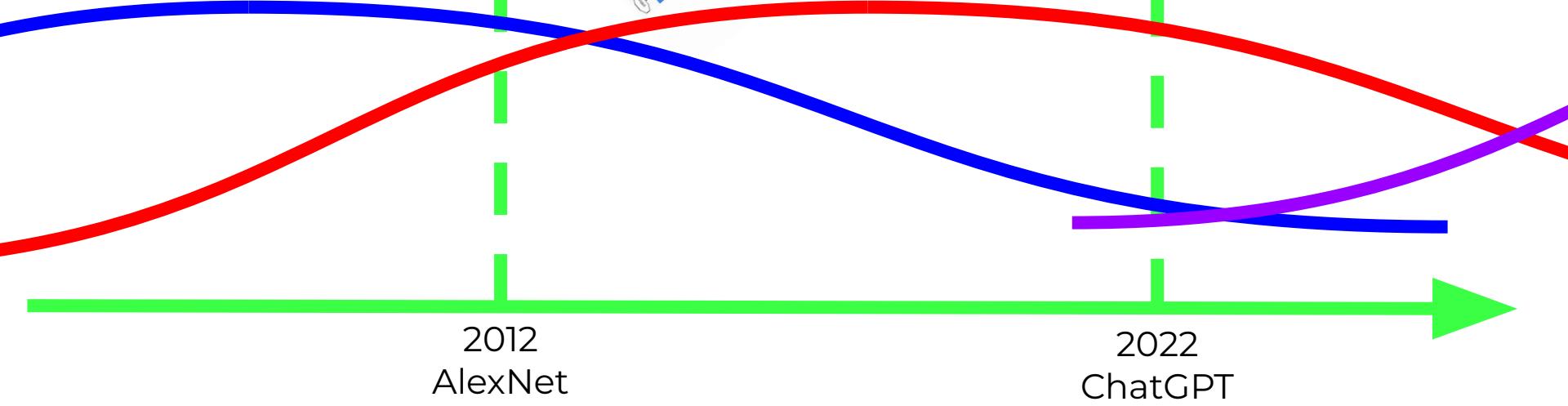
**Era of
Human-Crafter
Features**



**Era of
Deep Learning**



**Era of
LLMs**



APIs multimodal LLMs

Local multimodal LLMs

CV foundational models

“traditional” CV DL

“classical” CV

math

APIs multimodal LLMs

Local multimodal LLMs

CV DL foundational models

CV DL fundamentals

“classical” CV

math

APIs



User

What is funny about this image? Describe it panel by panel.



Source: [hmmm \(Reddit\)](#)

GPT-4 The image shows a package for a "Lightning Cable" adapter with three panels.

Panel 1: A smartphone with a VGA connector (a large, blue, 15-pin connector typically used for computer monitors) plugged into its charging port.

Panel 2: The package for the "Lightning Cable" adapter with a picture of a VGA connector on it.

Panel 3: A close-up of the VGA connector with a small Lightning connector (used for charging iPhones and other Apple devices) at the end.

The humor in this image comes from the absurdity of plugging a large, outdated VGA connector into a small, modern smartphone charging port.

APIs + Fine-tuning



[Image fine-tuning](#)



[Introducing vision to the fine-tuning API | OpenAI](#)

APIs multimodal LLMs

Local multimodal LLMs

DL foundational models

CV DL fundamentals

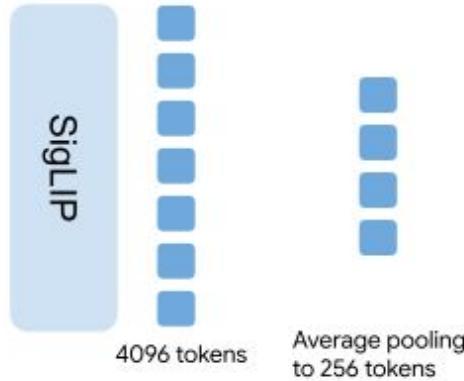
“classical” CV

math

Gemma



896 x 896



“What is this?”

- Image tokens (blue square)
- Input text tokens (orange square)
- Output text tokens (purple square)



[Gemma explained: What’s new in Gemma 3 - Google Developers Blog](#)
[2310.07707](#)

Llama 4: Leading Multimodal Intelligence

Newest model suite offering unrivaled speed and efficiency

Llama 4 Behemoth

288B active parameter, **16** experts

2T total parameters

The most intelligent teacher model for distillation

Preview

Llama 4 Maverick

17B active parameters, **128** experts

400B total parameters

Native multimodal with **1M** context length

Available

Llama 4 Scout

17B active parameters, **16** experts

109B total parameters

Industry leading **10M** context length
Optimized inference

Available

[The Llama 4 herd: The beginning of a new era of natively multimodal AI innovation](#)

APIs multimodal LLMs

Local multimodal LLMs

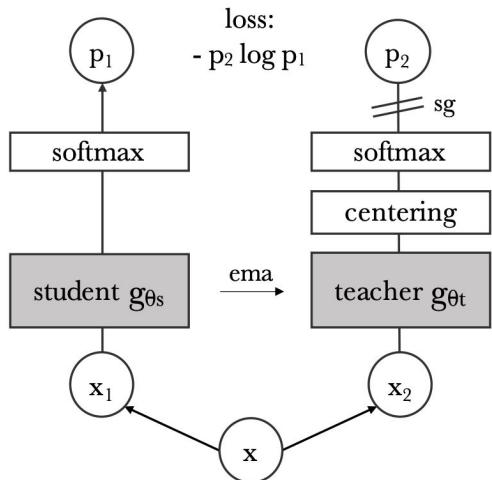
CV foundational models

CV DL fundamentals

“classical” CV

math

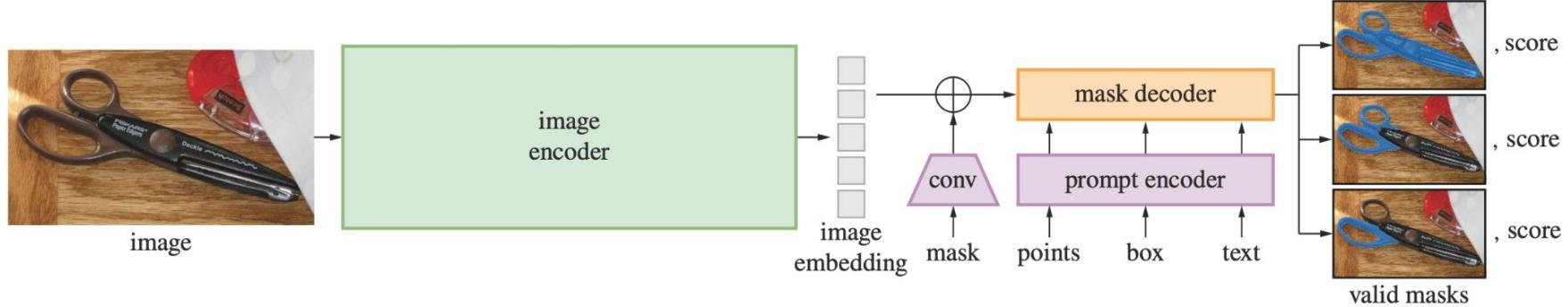
DINO (Self-Distillation with **no** labels)



[2104.14294](#), [2304.07193](#), [2508.10104](#)

[DINOv3](#)

SAM (Segment Anything)



[2304.02643](#); [2408.00714](#)

Introducing Meta Segment
Anything Model 2 (SAM 2)

Segment Anything 3

Coming 2025

APIs multimodal LLMs

Local multimodal LLMs

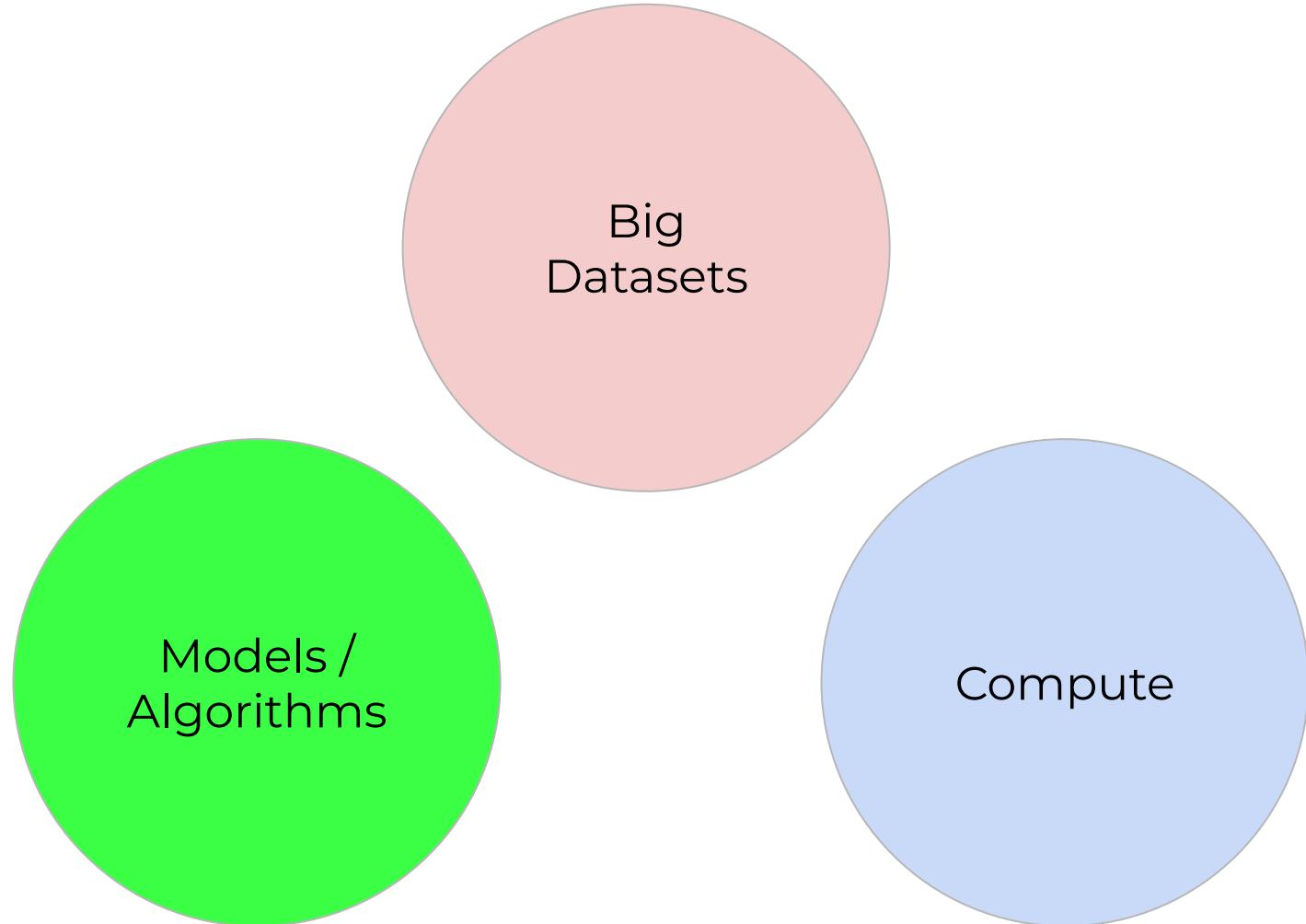
CV DL foundational models

CV DL fundamentals

“classical” CV

math

Intro



Intro



V. Vapnik
(creator of SVM)

L. Jackel
(proponent of NN)



Y. LeCun

1. Jackel bets (one fancy dinner) that by March 14, 2000, people will understand quantitatively why big neural nets working on large databases are not so bad. (Understanding means that there will be clear conditions and bounds)

Vapnik bets (one fancy dinner) that Jackel is wrong.

But .. If Vapnik figures out the bounds and conditions, Vapnik still wins the bet.

2. Vapnik bets (one fancy dinner) that by March 14, 2005, no one in his right mind will use neural nets that are essentially like those used in 1995.

Jackel bets (one fancy dinner) that Vapnik is wrong

V. Vapnik

3/14/95

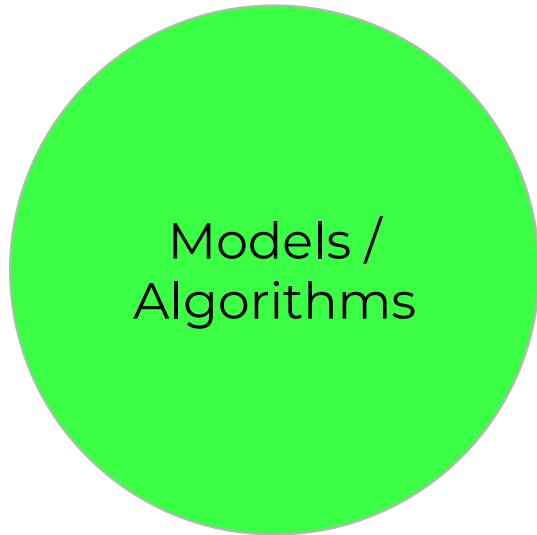
L. Jackel

3/14/95

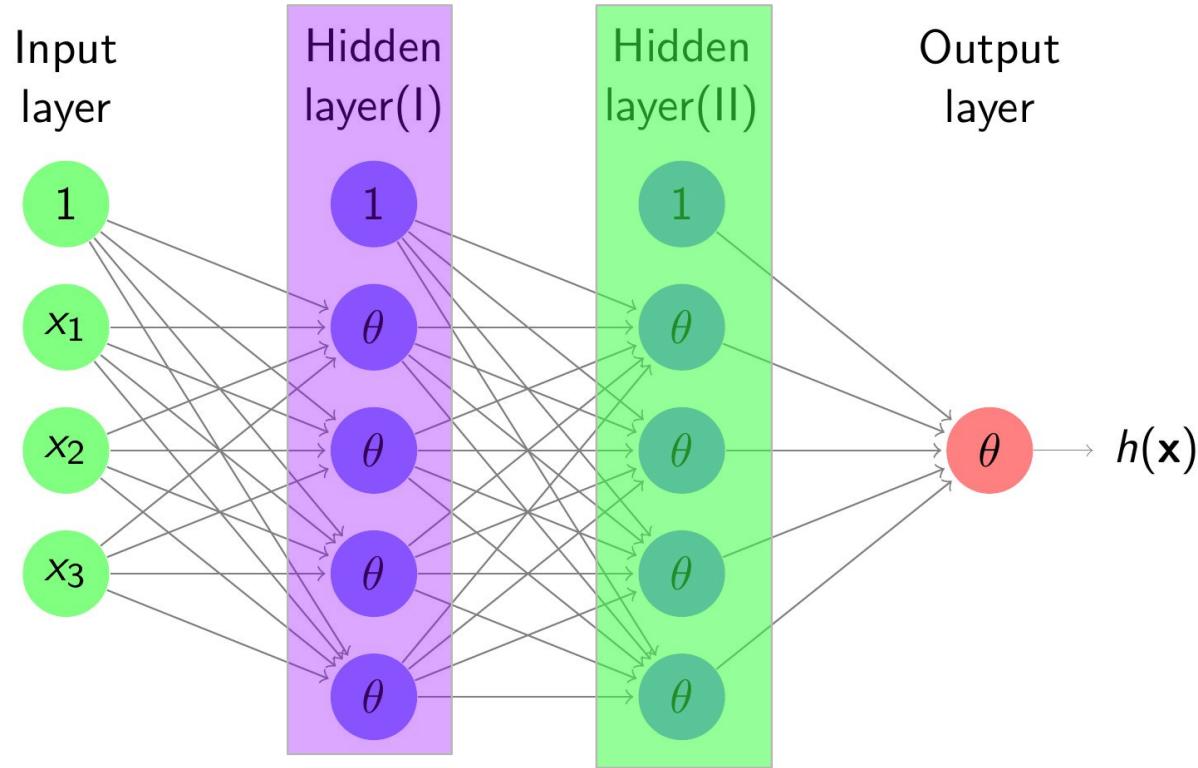
Witnessed by Y. LeCun

3/14/95

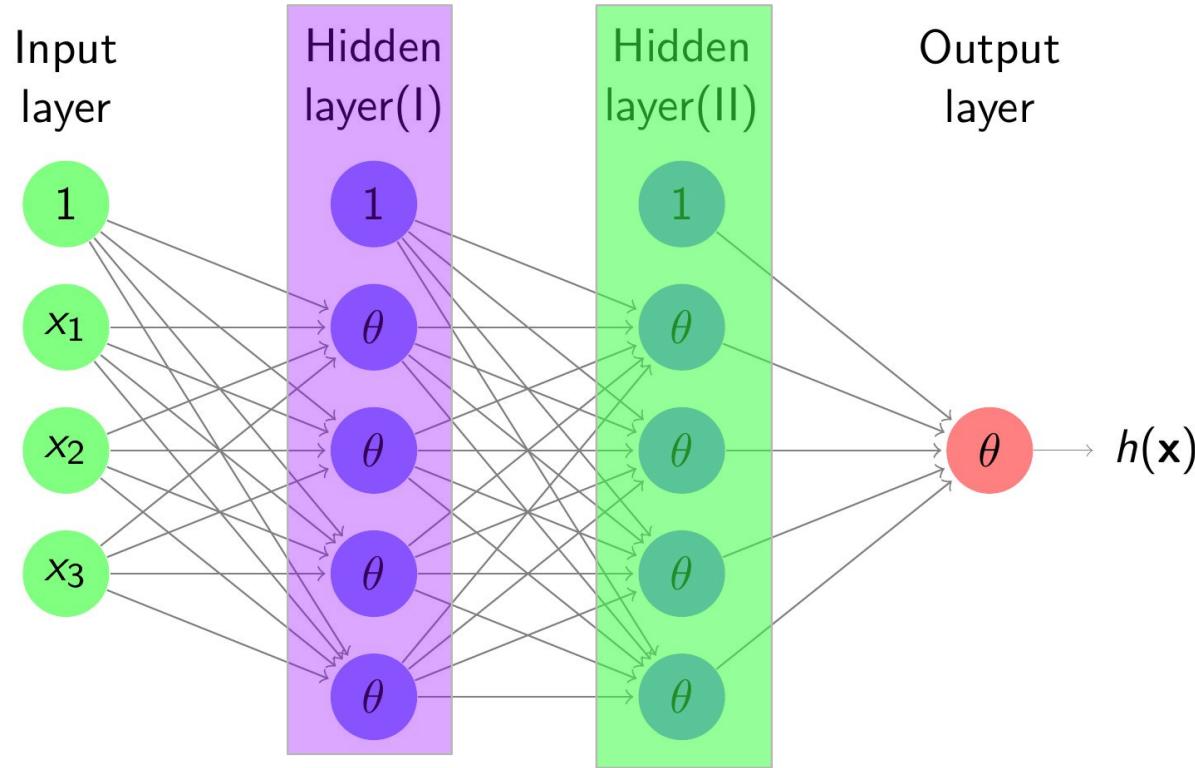
Intro



- Neural Network
- Number of layers => Deep
- Important specific components
(convolutions, attention, ...)

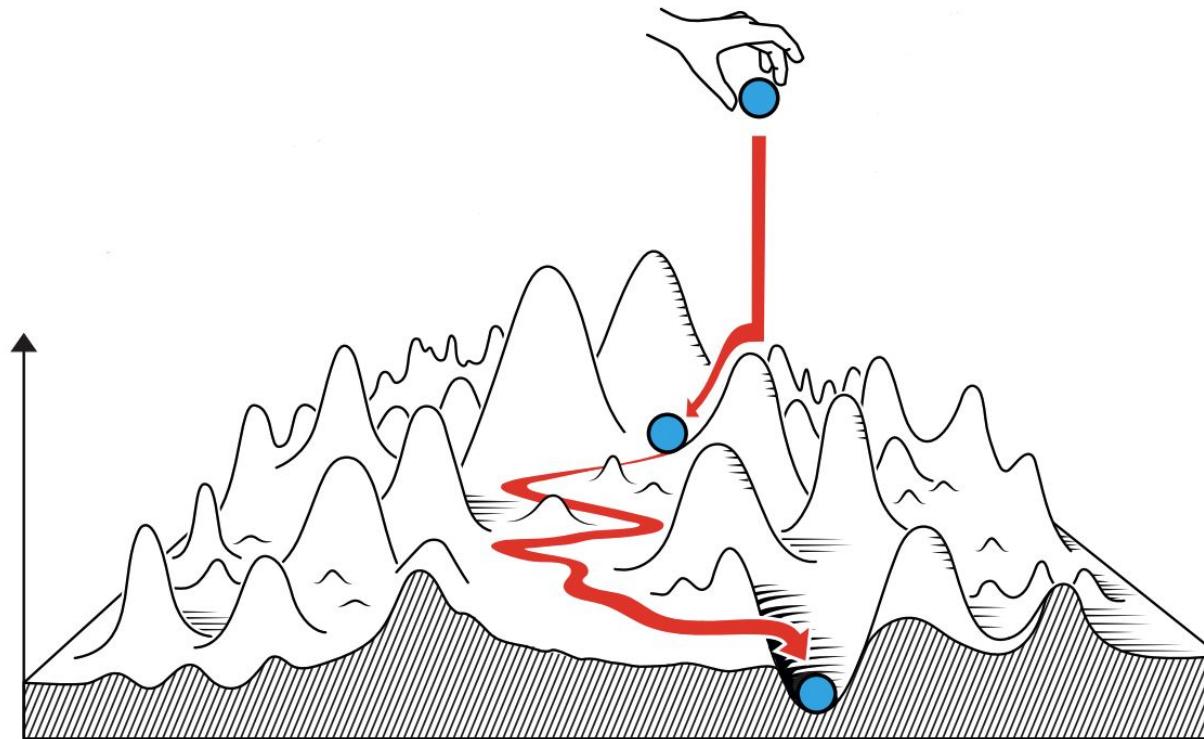


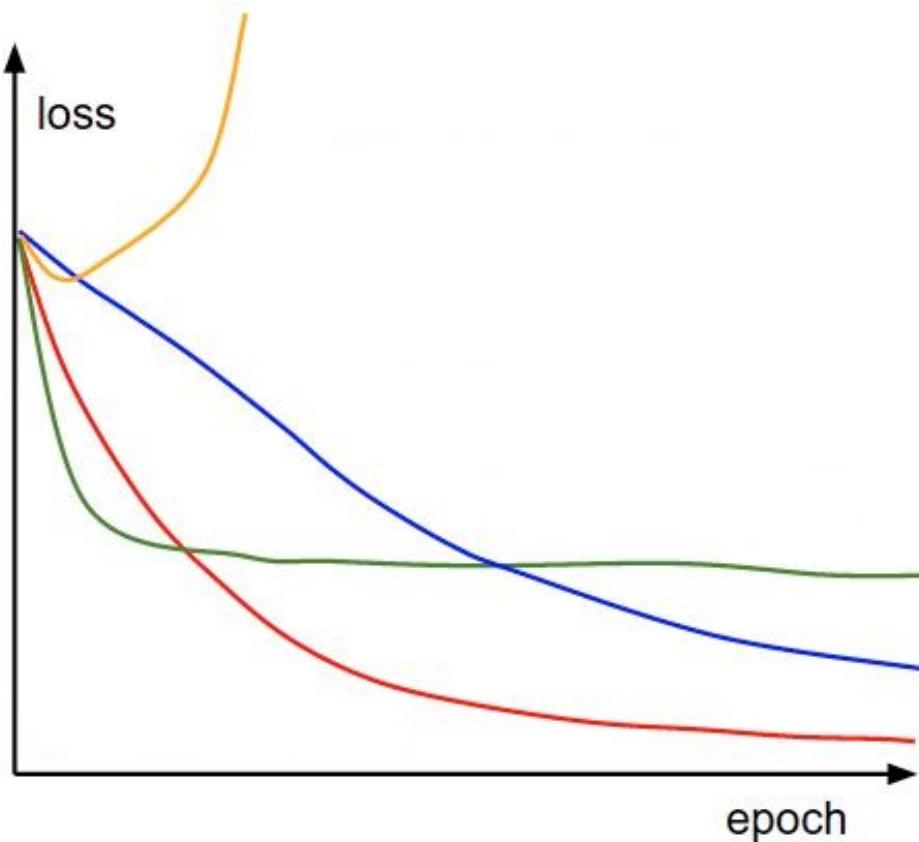
$$\{f(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{W}_L \sigma_L(\mathbf{W}_{L-1} \cdots \sigma_2(\mathbf{W}_2 \sigma_1(\mathbf{W}_1 \mathbf{x}))) \mid \boldsymbol{\theta} = \{\mathbf{W}_1, \dots, \mathbf{W}_L\}\}$$

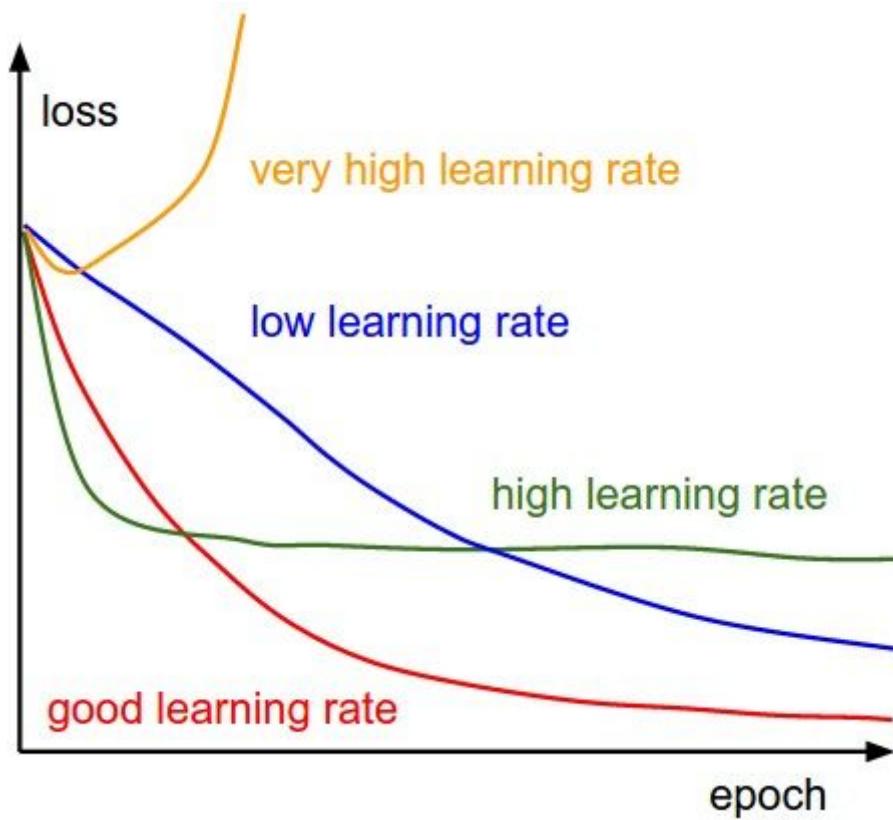


← Back-propagate error signal to get derivatives for learning

Optimization

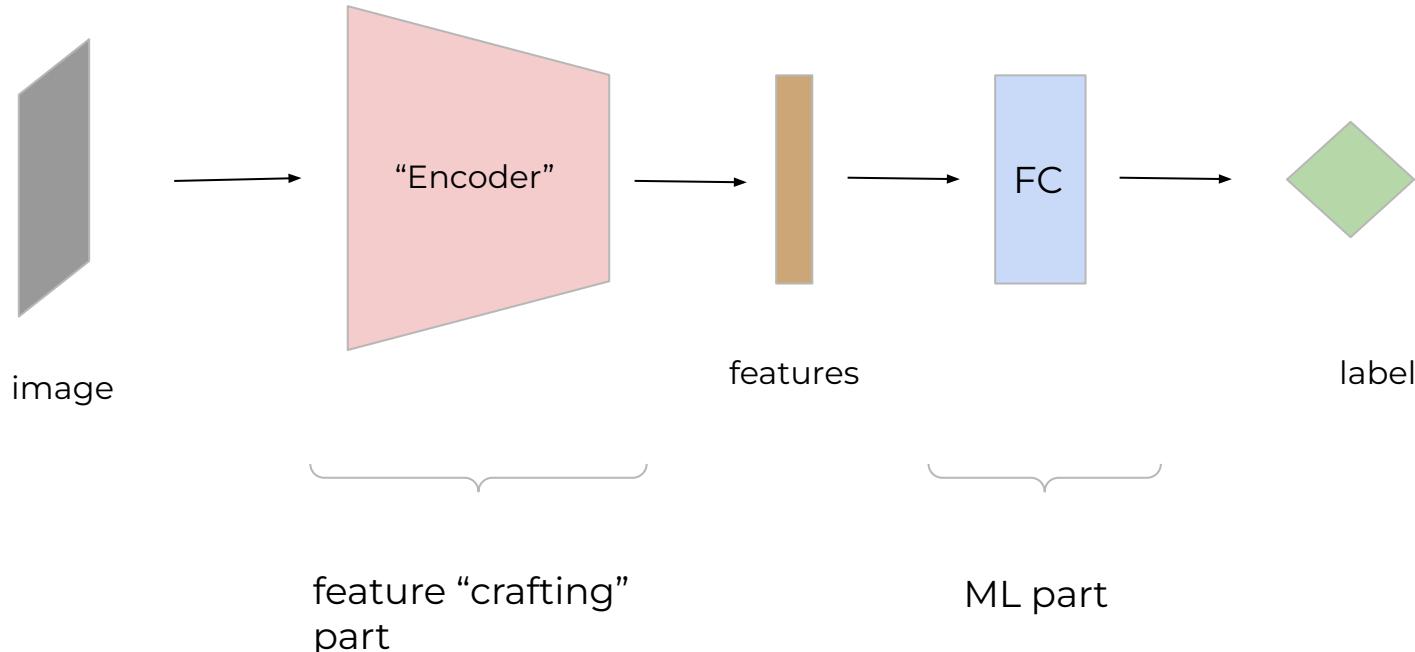




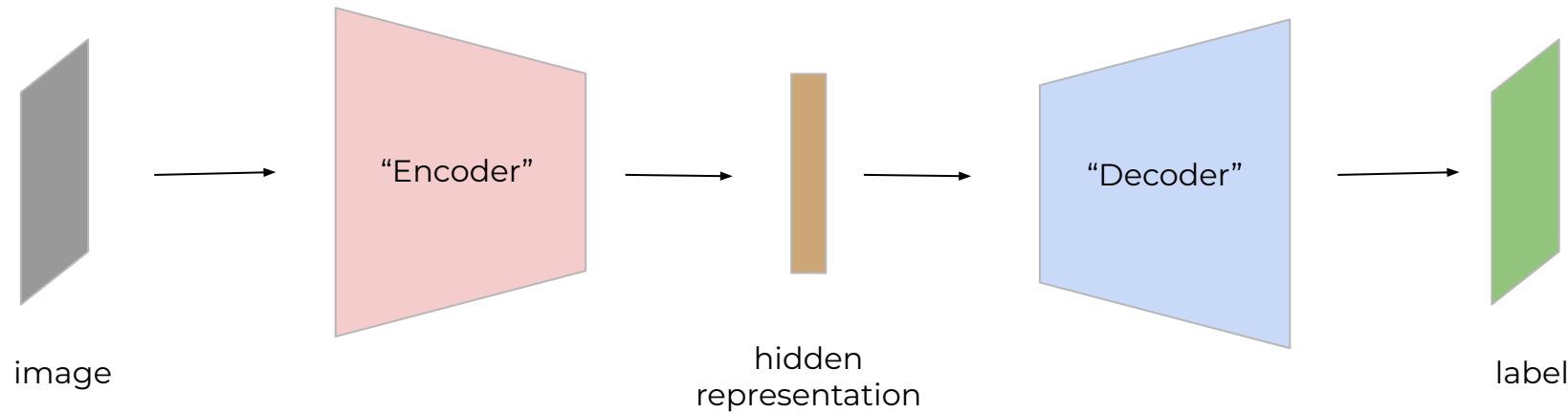


Typical CV architecture

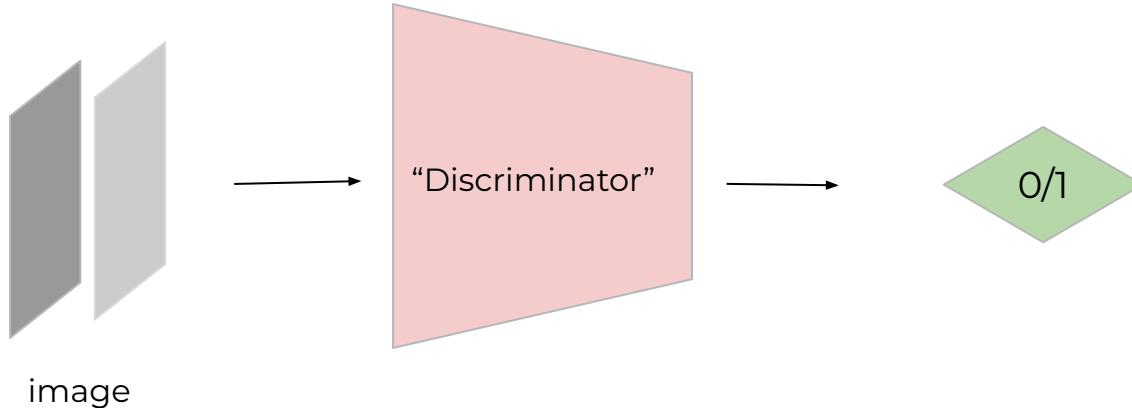
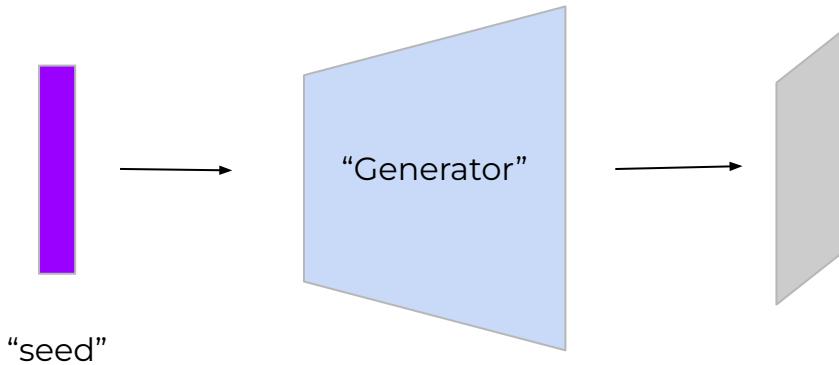
Typical CV architecture [classification]



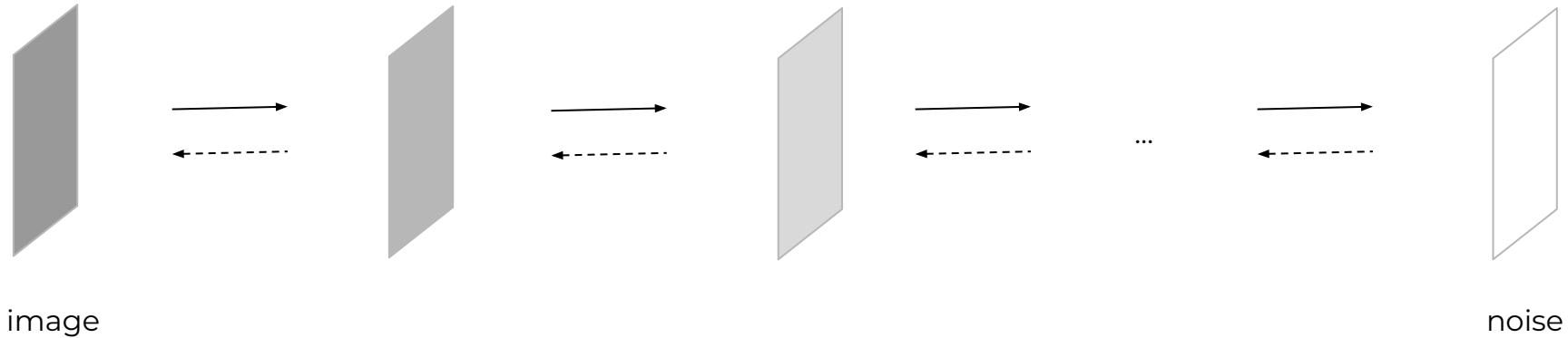
Typical CV architecture [segmentation]



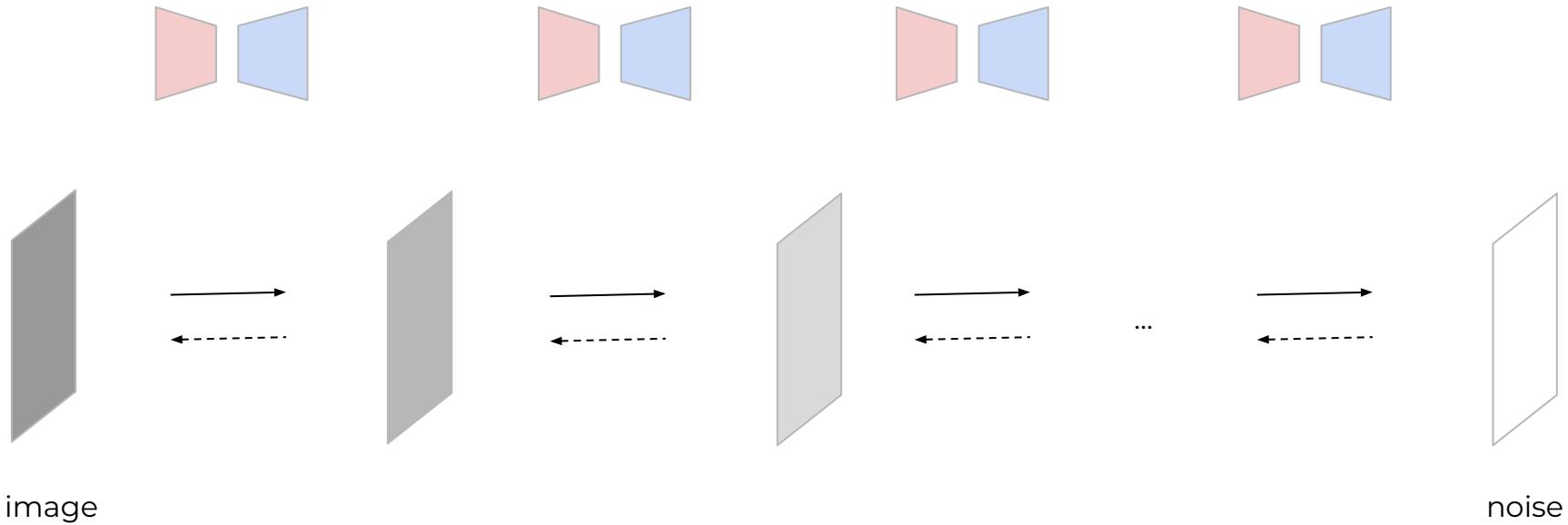
Typical CV architecture [Generation, GAN]

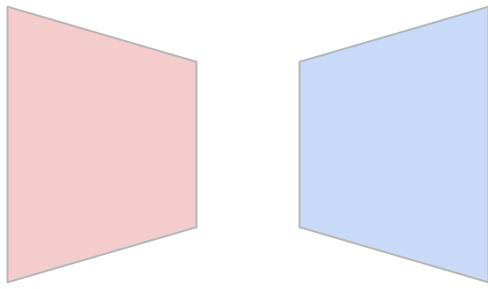


Typical CV architecture [Generation, Diffusion]



Typical CV architecture

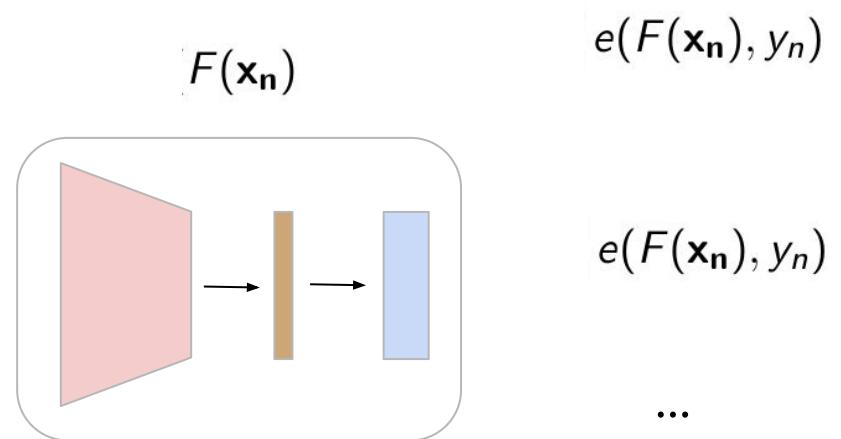
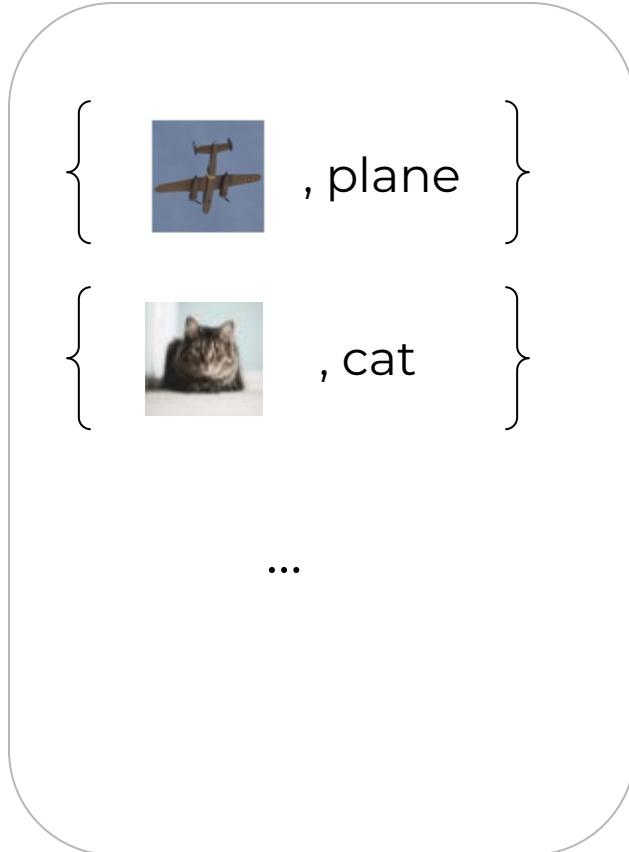




- convolutions
 - ...
- attention
 - ...
- ...

Training strategies

Supervised way



$$L_{train}(\omega) = \frac{1}{N} \sum_{n=1}^N e(F(\mathbf{x}_n), y_n)$$

require images paired with high-quality metadata

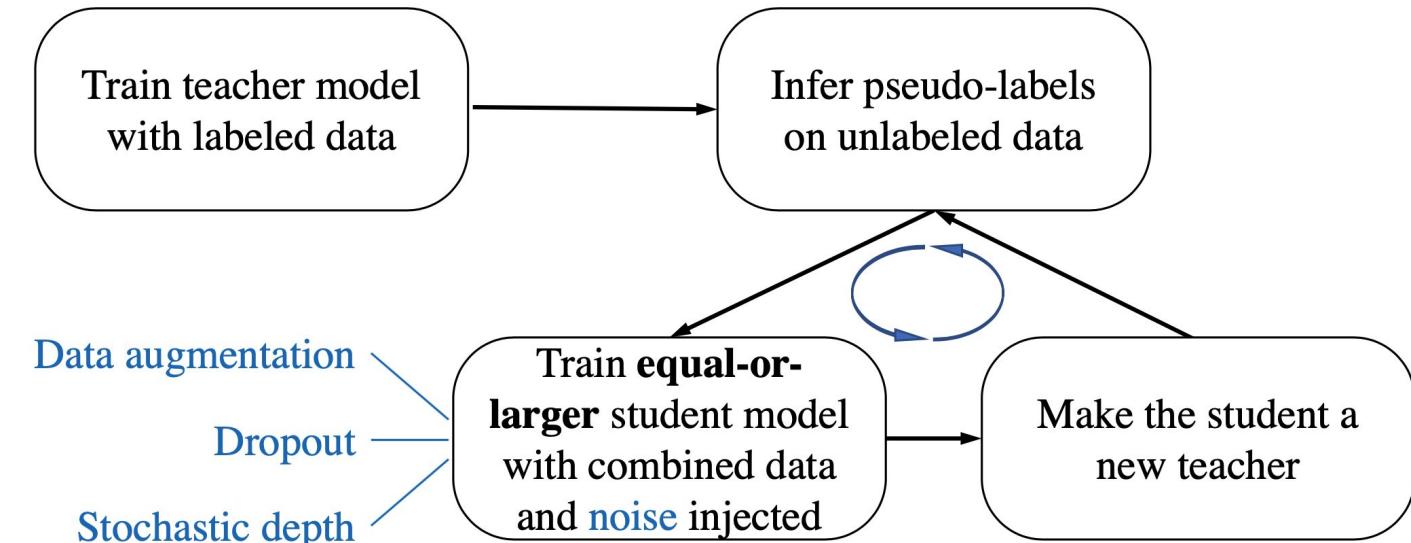
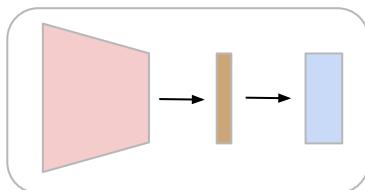
Weakly-supervised way

{
  , plane
}

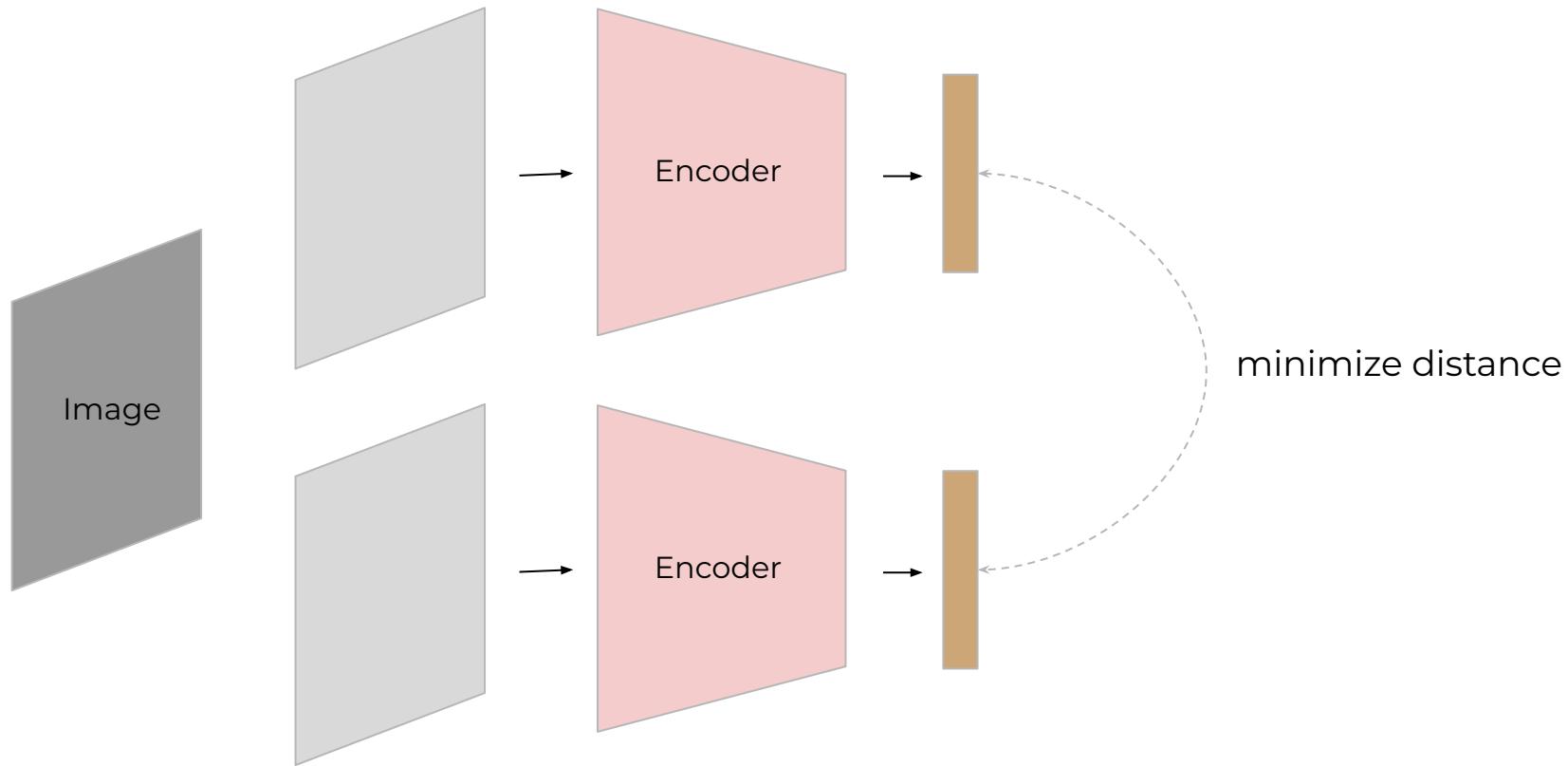
{
  , cat
}



...

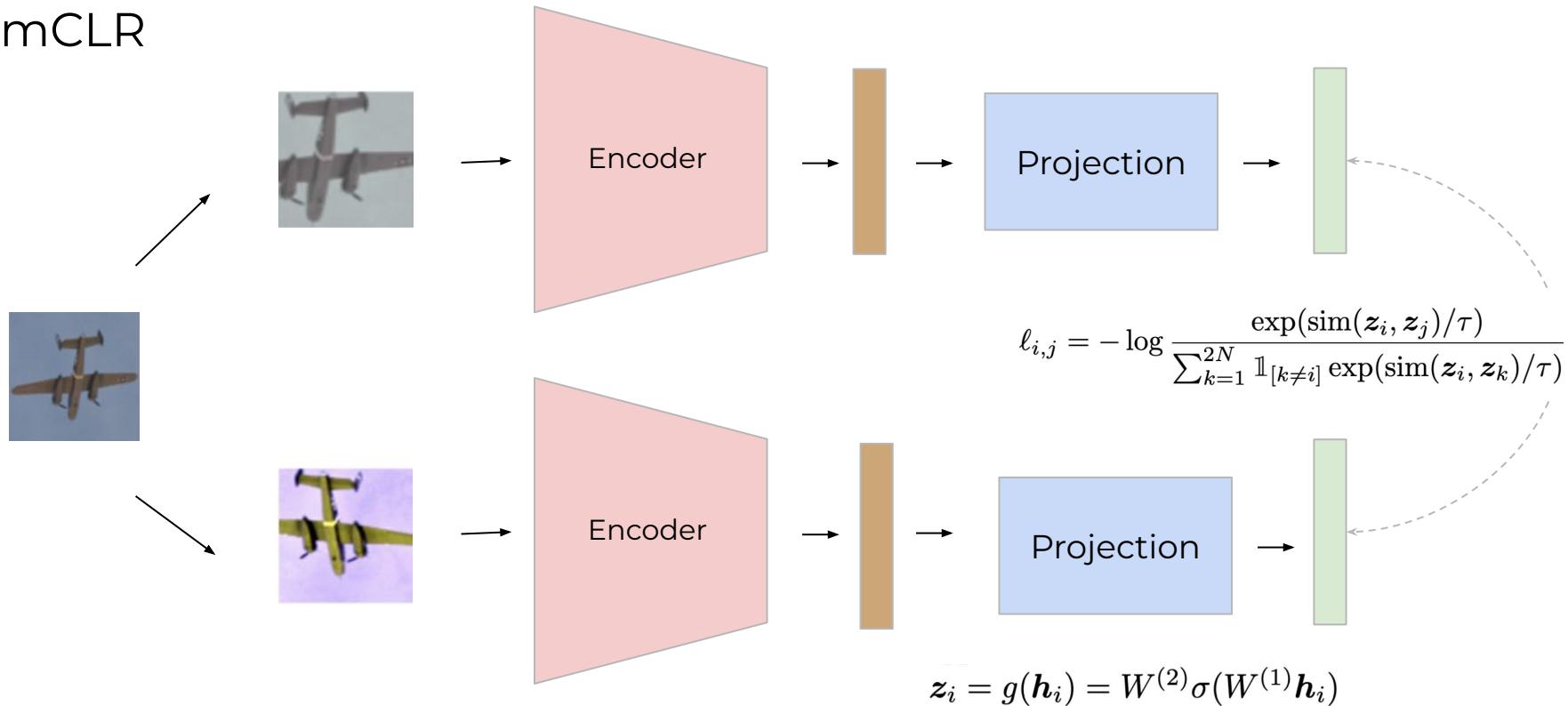


Self-supervised way



Self-supervised way

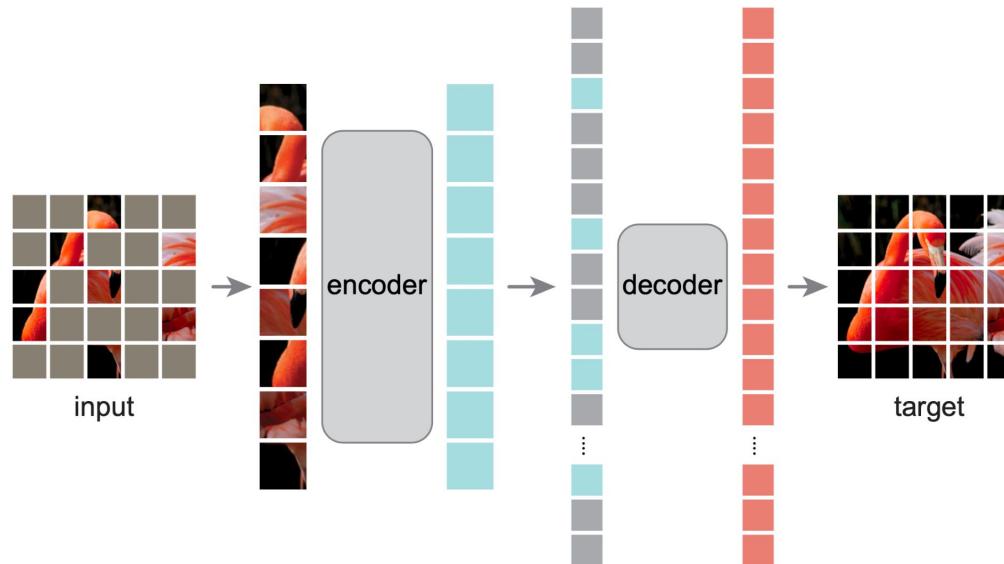
SimCLR



2002.05709

Self-supervised way

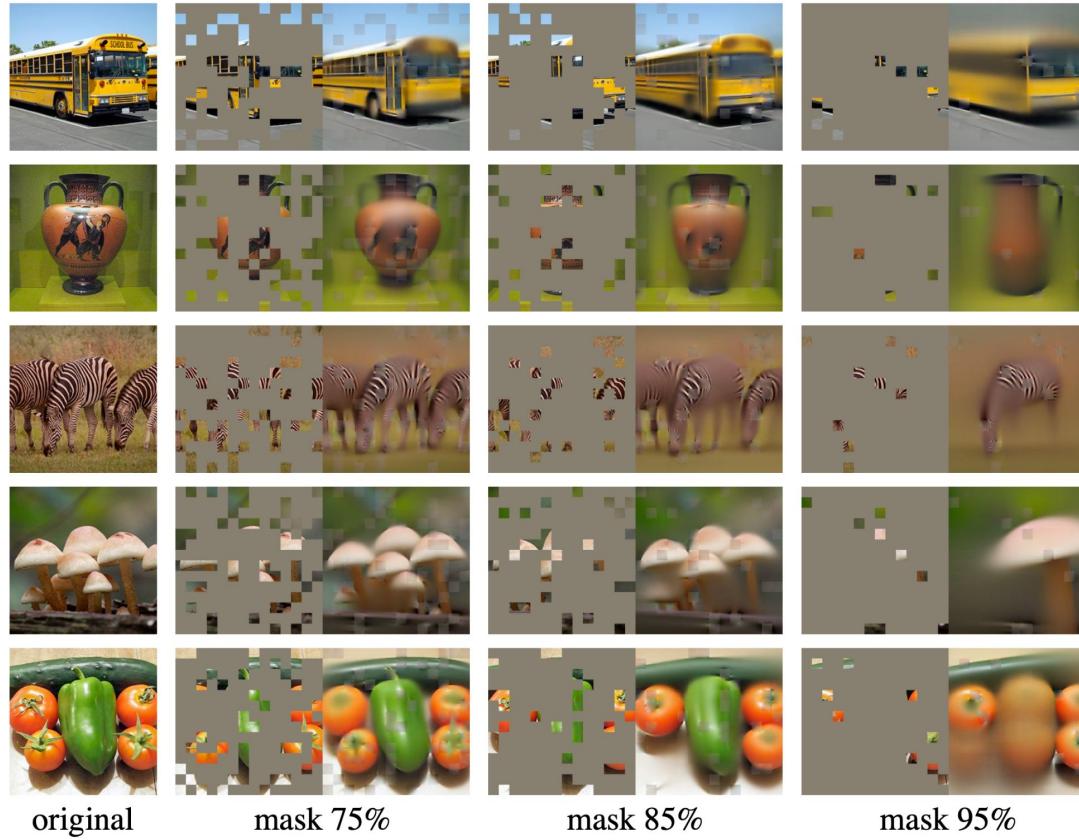
MAE



2111.06377

Self-supervised way

MAE



2111.06377

Self-supervised way

- do not require annotated data
- robust to input distribution shifts, provide strong global and local features, and generate rich embeddings that facilitate physical scene understanding.
- produce versatile and robust generalist features since SSL models are not trained for any specific downstream task
- requiring no human intervention, is well-suited for lifelong learning amid the growing volume of web data

Self-supervised way

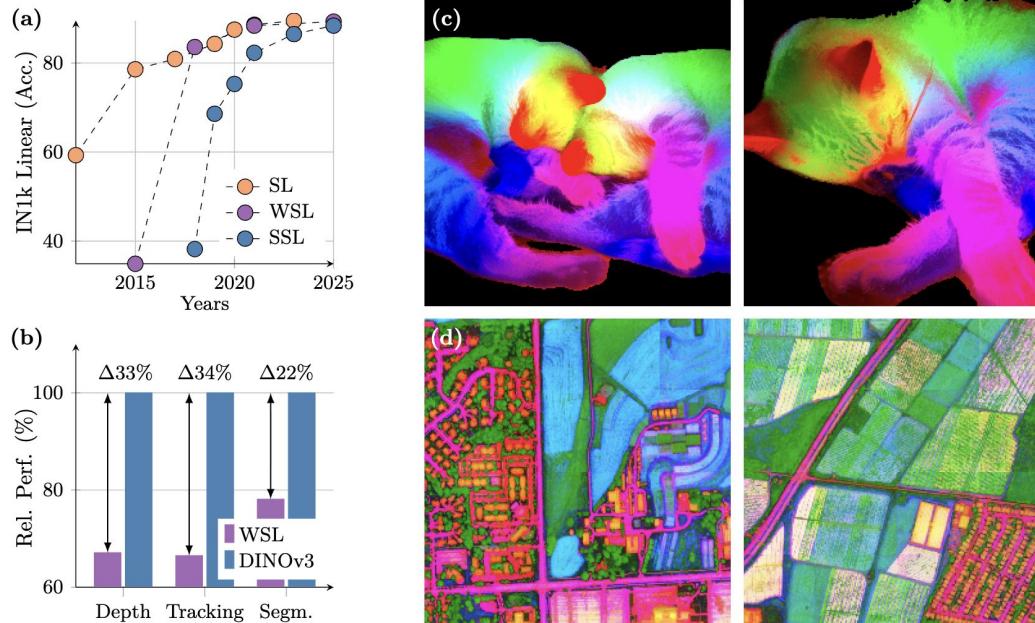


Figure 1: (a) Evolution of linear probing results on ImageNet1k (IN1k) over the years, comparing fully-supervised (SL), weakly-supervised (WSL) and self-supervised learning (SSL) methods. Despite coming into the picture later, SSL has quickly progressed and now reached the Imagenet accuracy plateau of recent years. On the other hand, we demonstrate that SSL offers the unique promise of high-quality dense features. With DINOv3, we markedly improve over weakly-supervised models on dense tasks, as shown by the relative performance of the best-in-class WSL models to DINOv3 (b). We also produce PCA maps of features obtained from high resolution images with DINOv3 trained on natural (c) and aerial images (d).

APIs multimodal LLMs

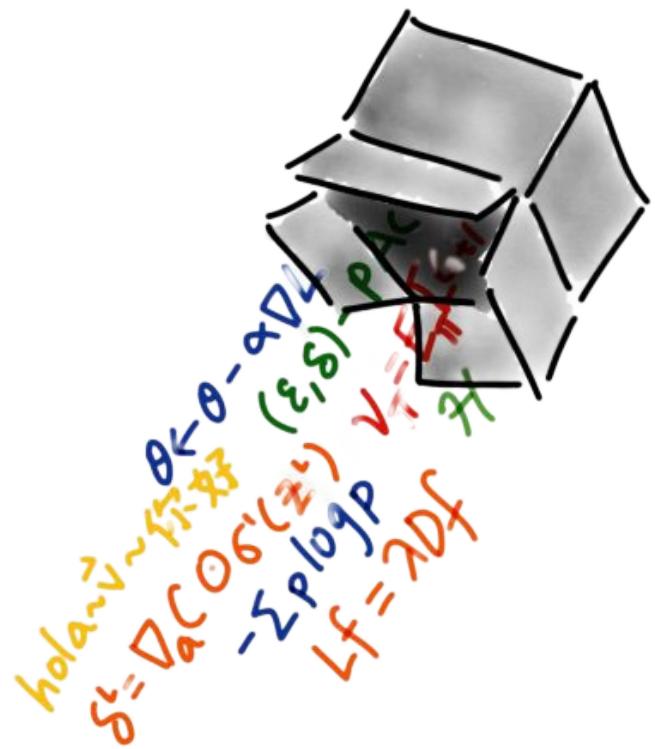
Local multimodal LLMs

CV DL foundational models

CV DL fundamentals

“classical” CV

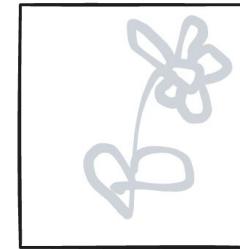
math





w

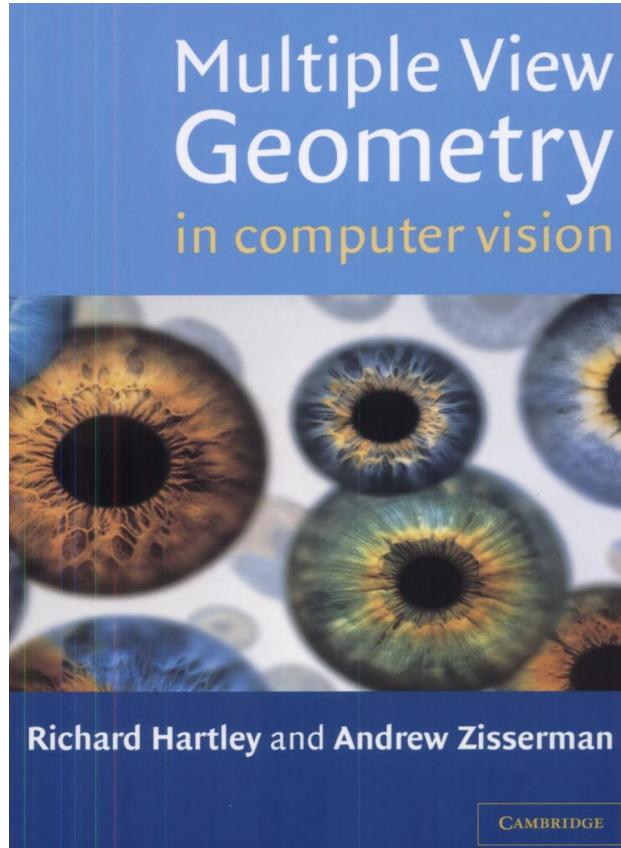
state of the world



x

measurements

- **discriminative** model $\Pr(w|x)$
- **generative** model $\Pr(x|w)$
- $\Pr(w|x) = \Pr(x|w) * \Pr(w) / \int [\Pr(x | w) * \Pr(w)] dw$



Richard Hartley and Andrew Zisserman
"Multiple View Geometry in Computer Vision"

[CS231A: Computer Vision, From 3D Perception to 3D Reconstruction and beyond](#)

https://www.google.com.ua/books/editio_n/Multiple_View_Geometry_in_Computer_Visio/

Summary

APIs multimodal LLMs

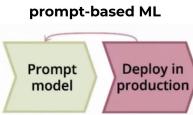
Local multimodal LLMs

CV foundational models

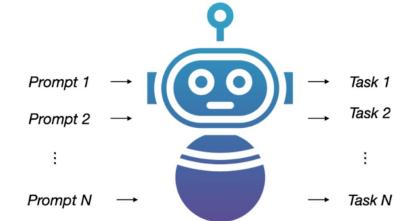
CV DL fundamentals

“classical” CV

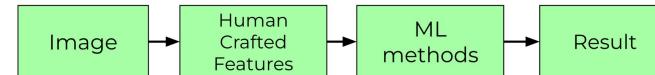
math



GEMMA, LLAMA, ...



DINO, SAM, ...



Descriptors:
SIFT, SURF, etc.
Linear, SVM,
XGBoost, ...

A 3D coordinate system with vectors labeled \vec{v}_1 , \vec{v}_2 , and \vec{v}_3 . A vector \vec{f} is shown as the sum of components along these axes: $\vec{f} = f_1 \vec{v}_1 + f_2 \vec{v}_2 + f_3 \vec{v}_3$.

Overview of the module

Lecture I Intro, big picture

Lecture II Essential architectures (CNN)

Lecture III Essential architectures (Transformers)

Lecture IV CV Foundational models in depth (DINO, SAM)

Lecture V Multimodal Foundational models (GEMMA, APIs)

Lecture VI Object Detection

Time to practice!



notebook

CIFAR10



DTD

freckled



banded



bubbly



swirly



bubbly



striped



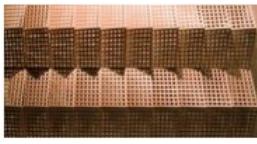
marbled



interlaced



meshed



swirly



frilly



marbled



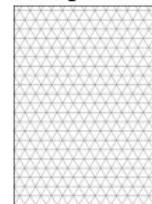
smeared



stained



grid



meshed



COCO-Q

painting



weather



cartoon



cartoon



painting



painting



painting



handmake



tattoo



weather



painting



painting



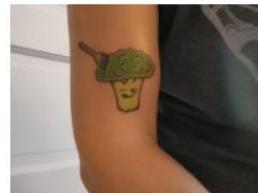
painting



sketch



tattoo



weather

