



3000788 Intro to Comp Molec Biol

Lecture 23: Introduction to machine learning and AI

Fall 2025



Sira Sriswasdi, PhD

- Research Affairs
- Center of Excellence in Computational Molecular Biology (CMB)
- Center for Artificial Intelligence in Medicine (CU-AIM)

Today's agenda

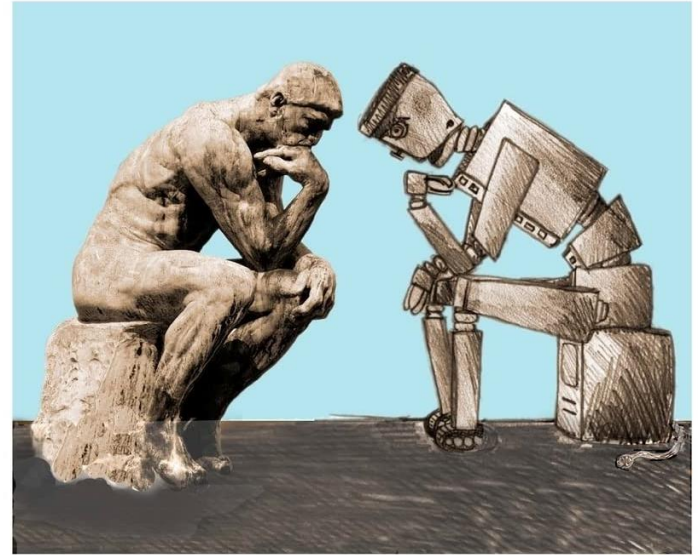


- What is machine learning and AI
- Evolution of AI: from rule-based to generative model

Natural vs artificial intelligence

- AI is a result of **computer algorithms mimicking the natural learning process**
- **Signs of intelligence**
 - Memorization
 - Pattern recognition and generalization
 - Learning from trials and errors
 - Ability to create
 - Reasoning with cause and effect

NATURAL AND ARTIFICIAL
INTELLIGENCE



ROBERT K. LINDSAY

Memorization and information compression



Pattern recognition

Human vs Machine: Pneumonia

Chest X-Rays image the lungs, heart, blood vessels, and bones. AI has been used to read and understand them.

Example:
Pneumonia

Computers:
Score: 0.371

Doctors:
0/15 Detected



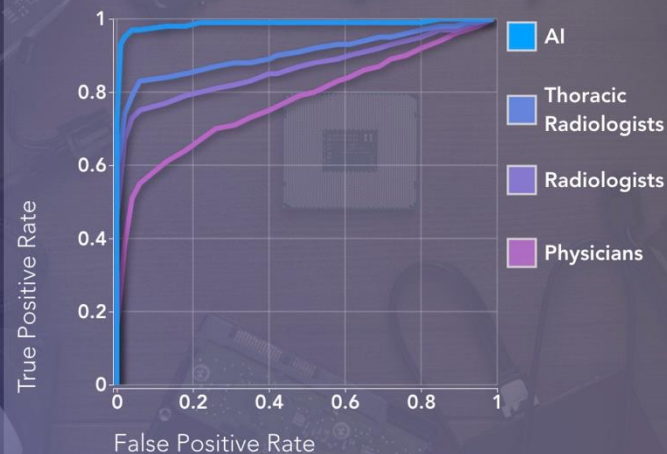
Clearvue Health

Hwang et al

AI vs Doctors: Chest X-Rays

AI was significantly more accurate and precise than radiologists and physicians in diagnosing chest x-rays.

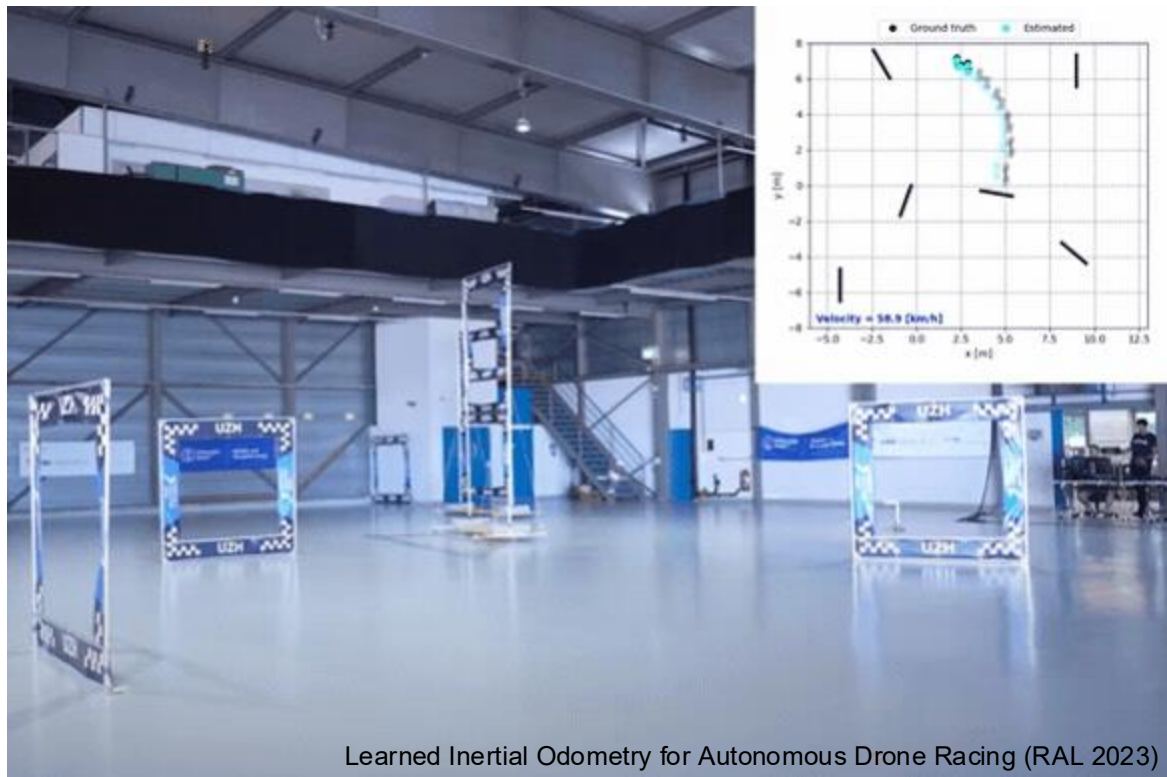
AUC-ROC: Human vs Computer



Clearvue Health

Hwang et al

Trial and error



Google DeepMind's AlphaGo computer beats top player Lee Sedol for third time to sweep competition



Learned Inertial Odometry for Autonomous Drone Racing (RAL 2023)



Machine learning (ML)

The engine behind modern AI

AI 1.0: Hand-crafted algorithms



<https://shapedshed.com/photoshop-101-the-magic-wand-tool/>

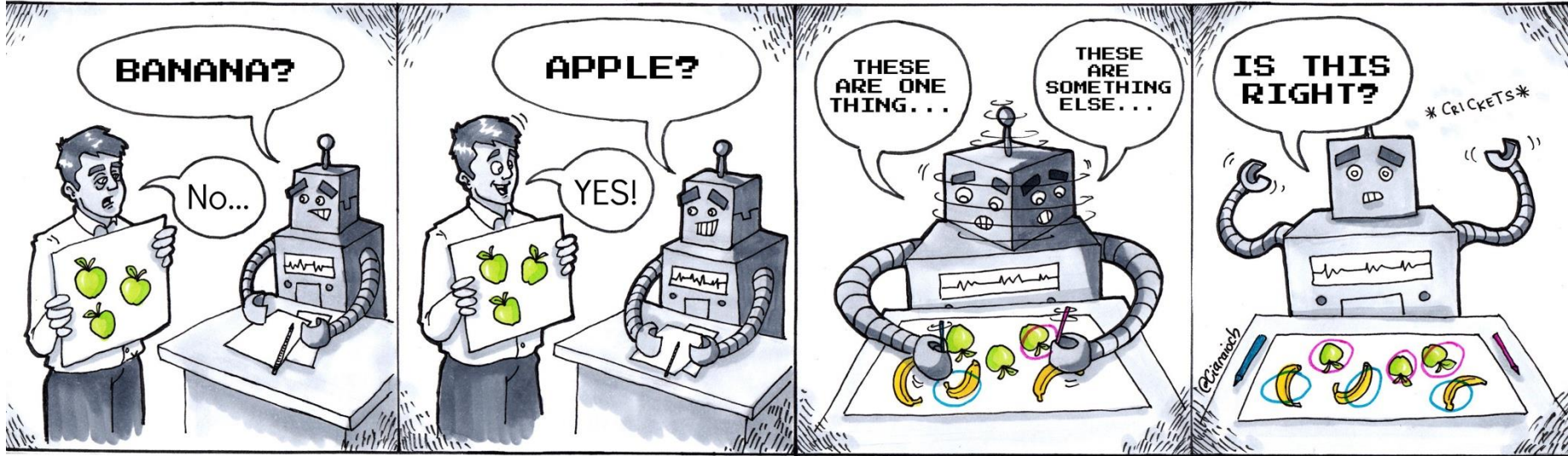
Making computer optimizes itself



Source: Cortes and Sanchez. IEEE Latin America Transaction (2021)

- The relationship is too complex for human to define
- Instead, human provides the data (x, y) and let the computer do the fitting

Machine learning paradigms



Supervised Learning

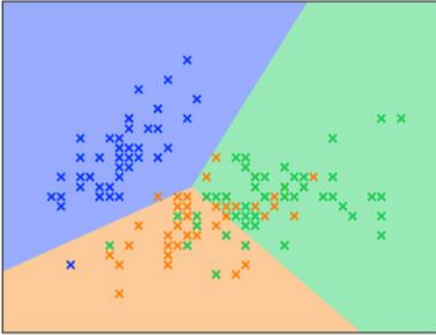
Find accurate decision functions

Unsupervised Learning

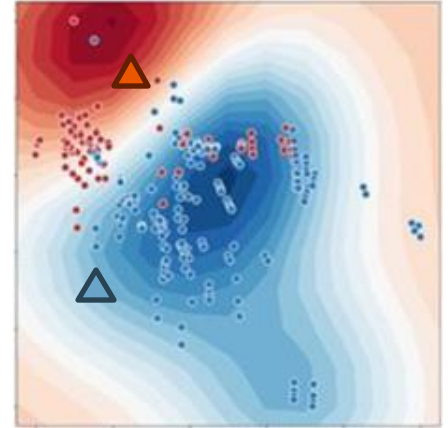
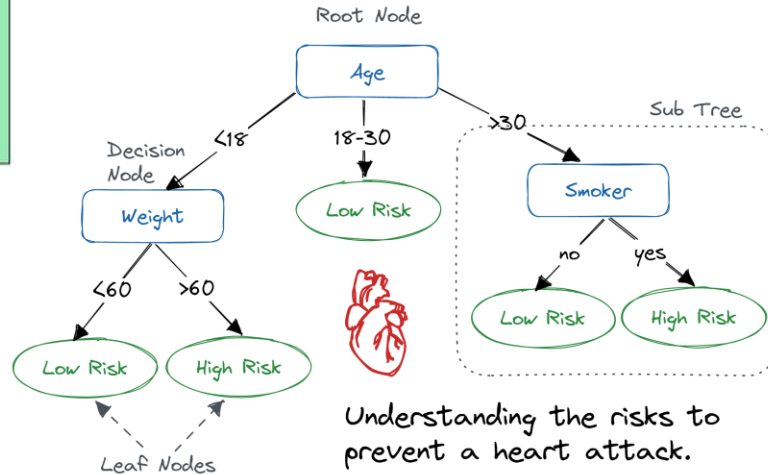
Find similarities among data

AI 2.0: Classical ML models

Linear: Weighted score = $(\text{input}_1 \times w_1) + \dots + (\text{input}_n \times w_n)$

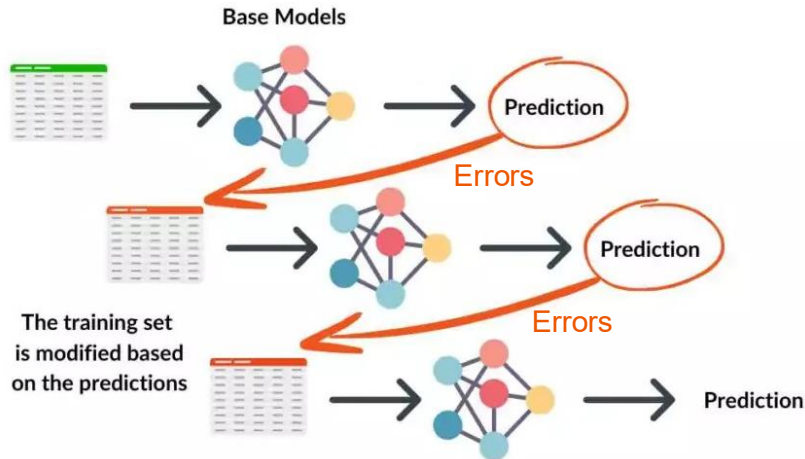


Tree: Collection of decisions



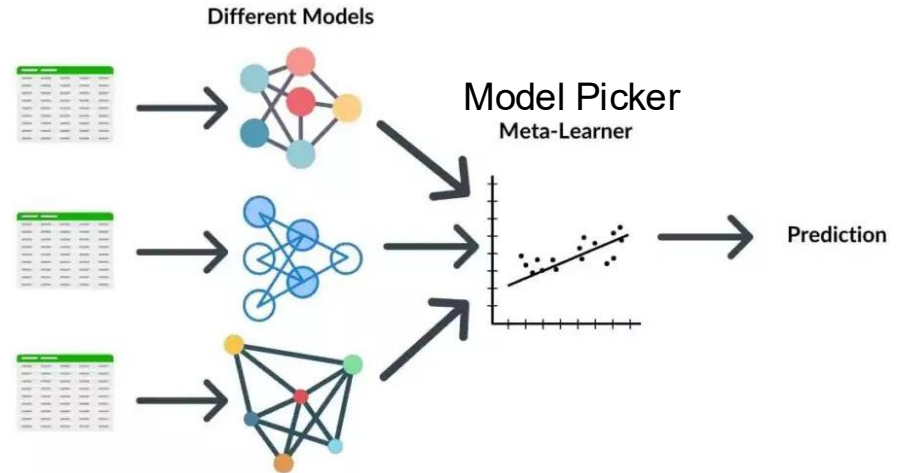
Nearest neighbors:
Predict using similarity to past observations

Enhancement with ensemble approaches

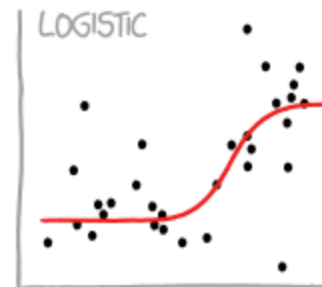
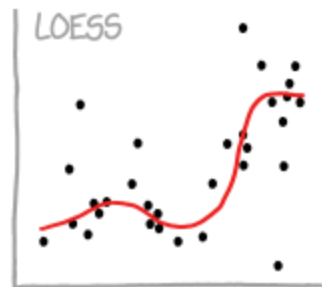
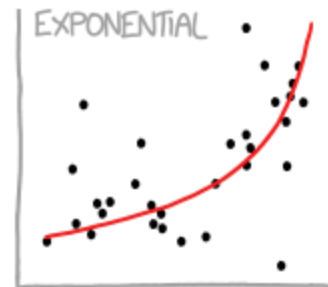
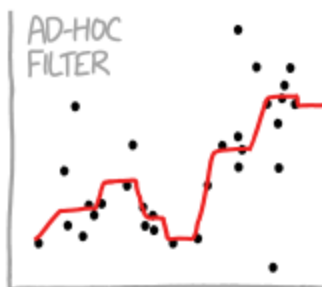
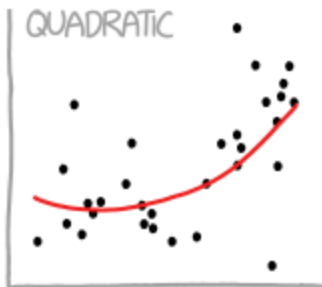
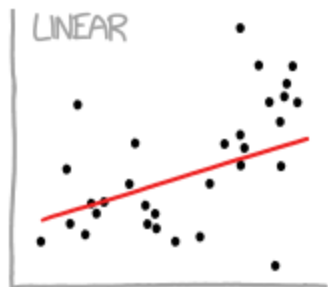


Boosting: Iterative improvement with additional models

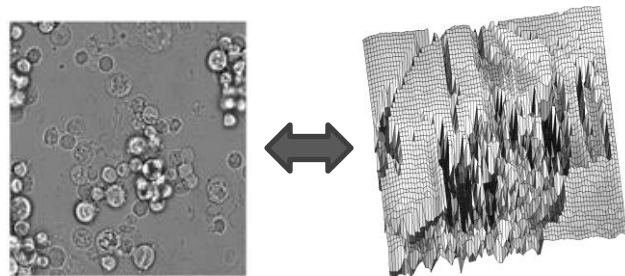
Stacking: Combine multiple models with different capabilities



Limitation of classical ML



- Unable to fit complex relationships
- Unable to handle raw data like text or image





Representation learning

Training computer to parse complex data

Naïve data representation is not useful for prediction

	1	2	3	4	5	6	7	8	9
man	1	0	0	0	man	= 1			
woman	0	1	0	0	woman	= 2			
boy	0	0	1	0	boy	= 3			
girl	0	0	0	1	girl	= 4			
prince	0	0	0	0			
princess	0	0	0	0	1	0	0	0	0
queen	0	0	0	0	0	1	0	0	0
king	0	0	0	0	0	0	1	0	0
monarch	0	0	0	0	0	0	0	0	1

Image from hackermoon.com

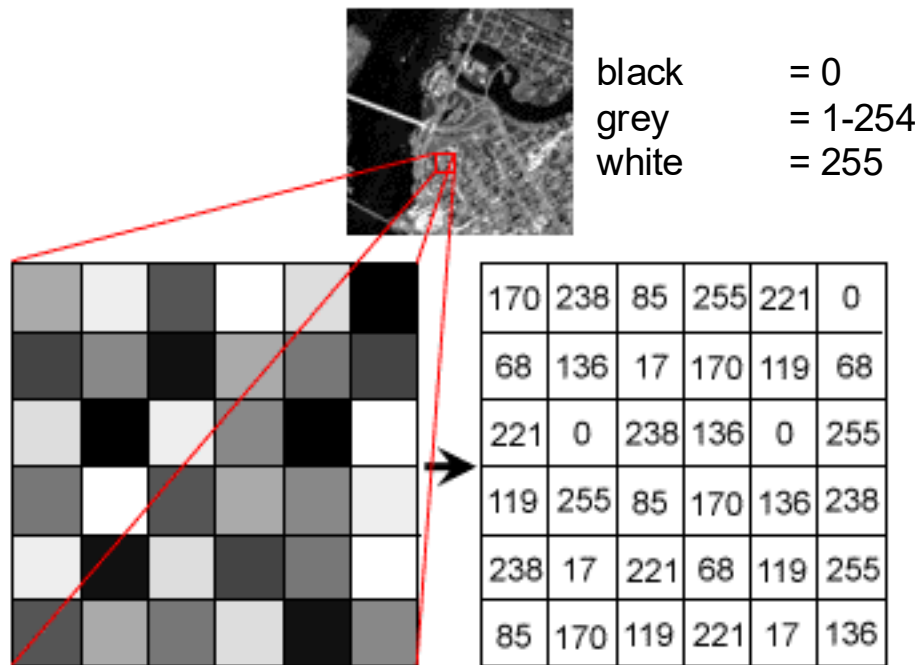
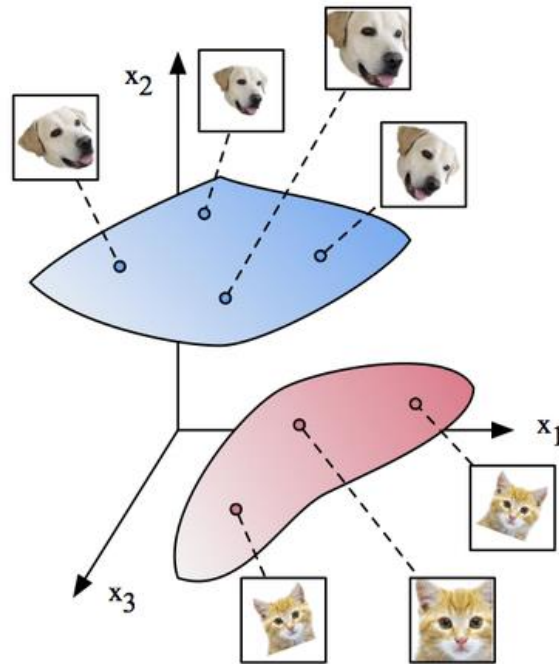


Image from naushardsblog.wordpress.com

What is a (good) representation of the data?

- Numerical mapping of the raw data that are informative of the **key characteristics**
- **Key characteristics** have different meaning based on your purpose for the data
 - Classification / regression
 - Reconstruction
 - Dimensionality reduction



Information extraction from image with kernels



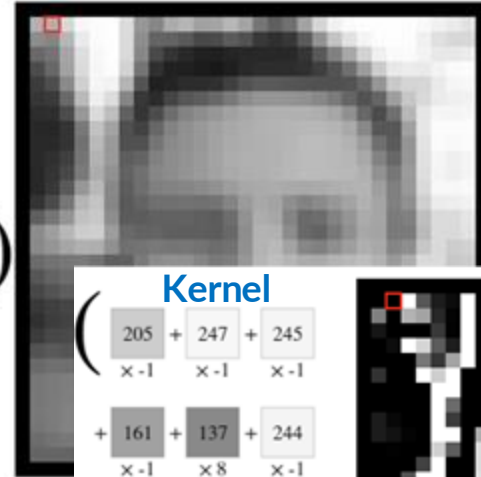
input image

Kernel

$$\begin{pmatrix} 205 & 247 & 245 \\ \times 0.0625 & \times 0.125 & \times 0.0625 \\ + 161 & 137 & 244 \\ \times 0.125 & \times 0.25 & \times 0.125 \\ + 154 & 75 & 200 \\ \times 0.0625 & \times 0.125 & \times 0.0625 \end{pmatrix}$$

= 175

kernel: blur



Kernel

$$\begin{pmatrix} 205 & 247 & 245 \\ \times -1 & \times -1 & \times -1 \\ + 161 & 137 & 244 \\ \times -1 & \times 8 & \times -1 \\ + 154 & 75 & 200 \\ \times -1 & \times -1 & \times -1 \end{pmatrix}$$

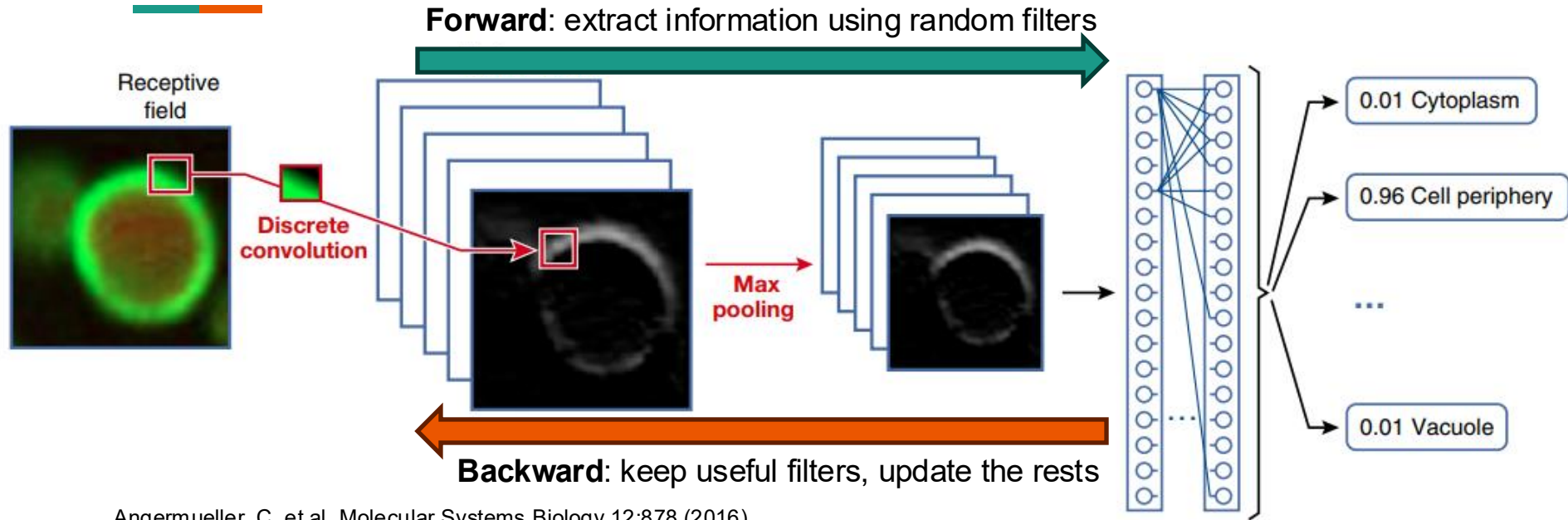
= -435

kernel: outline



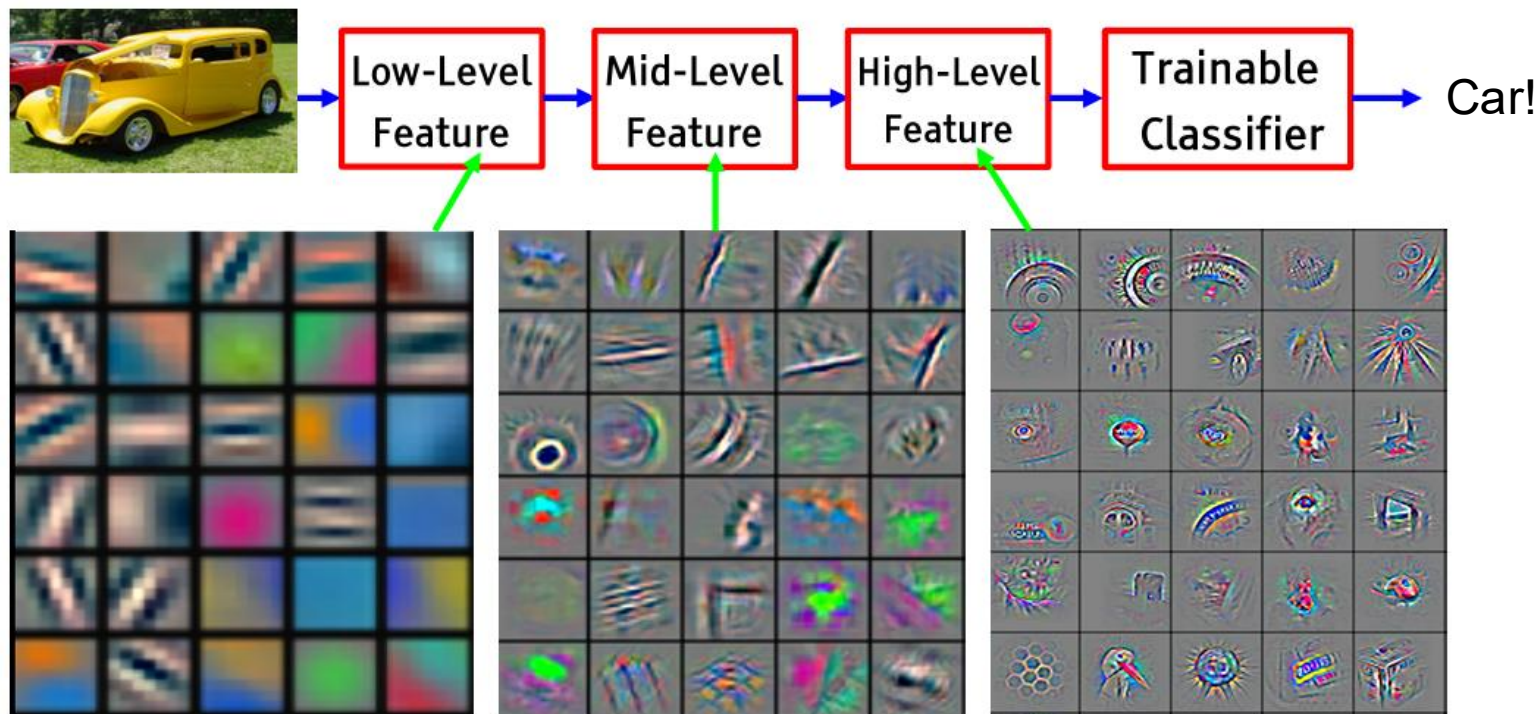
output image

AI 3.0: Data-driven representation learning

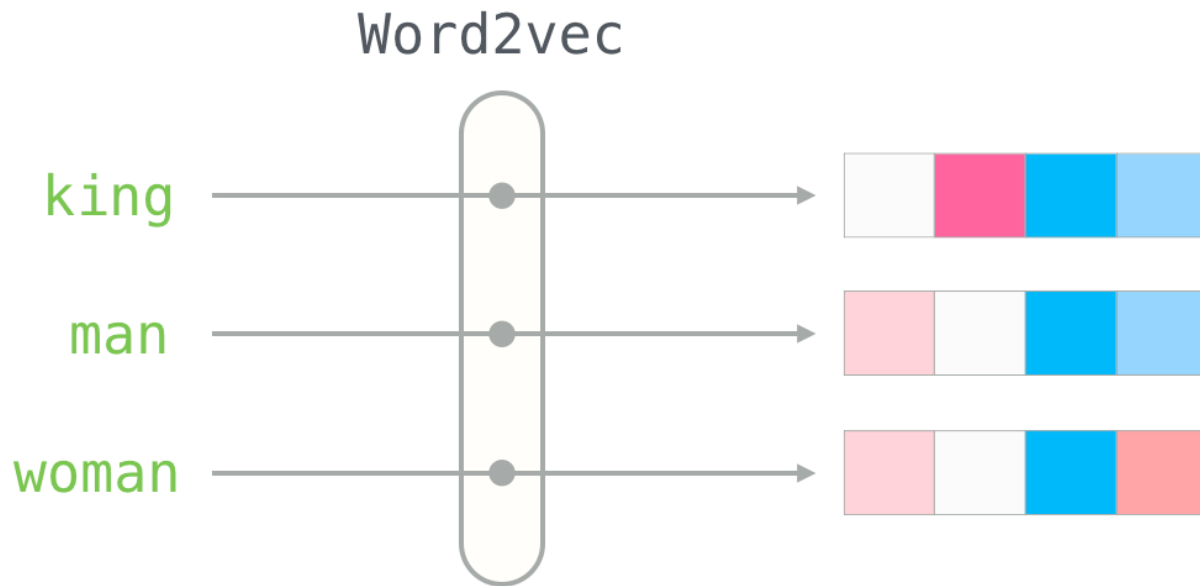


- Finding the right kernels is too difficult for human
- Instead, human provides the data (x, y) and let the computer do the fitting

Convolution kernels create image representation

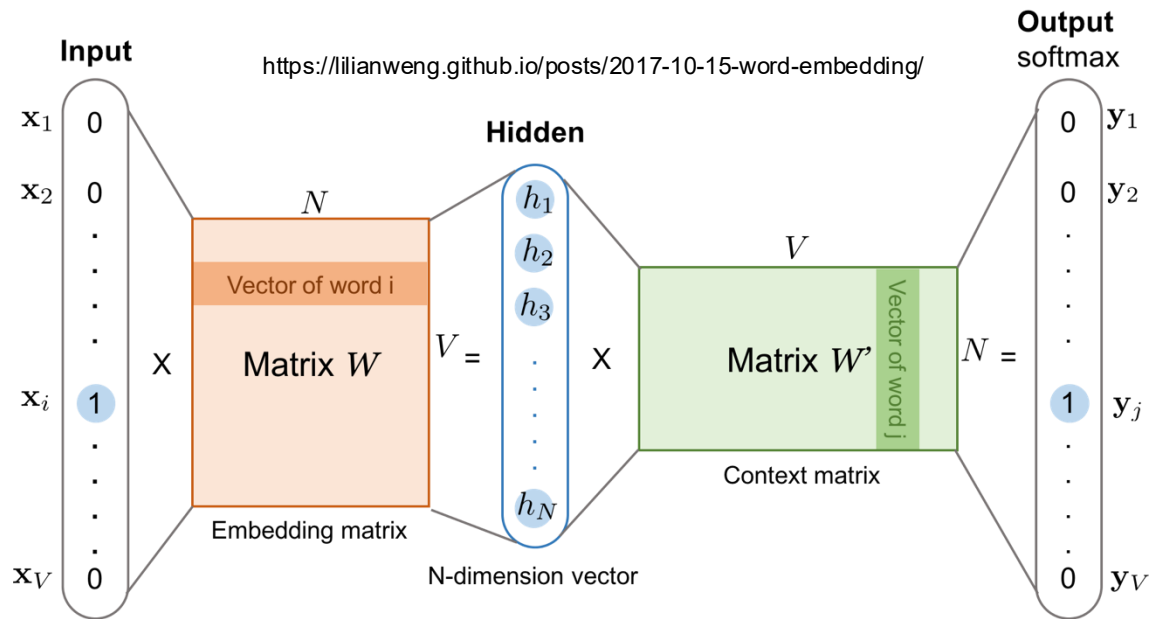


Word embedding



- Transformation of a word into a useful vector that captures its meaning and other characteristics. But how?

How to learn a good word embedding?



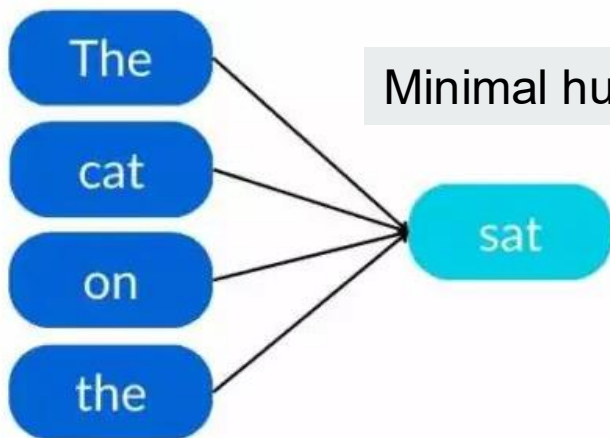
- Map each word to an n -dimensional vector (representation)
- Use these representations to make prediction, but what to predict?

Predict the central word or surrounding words

Example Sentence: The cat sat on the mat.

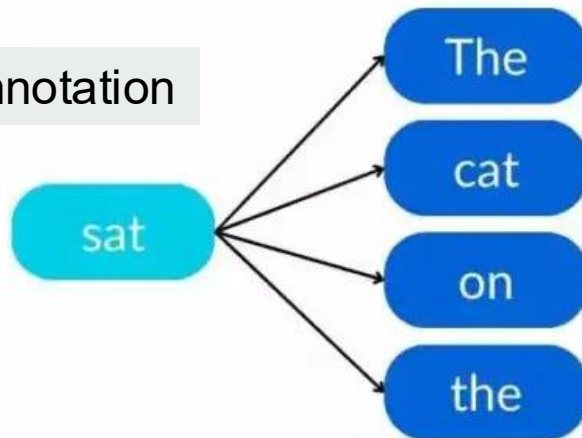
Continuous Bag-of-Words (CBOW)

Goal: Given context words,
predict the target word.



Skip-gram Model

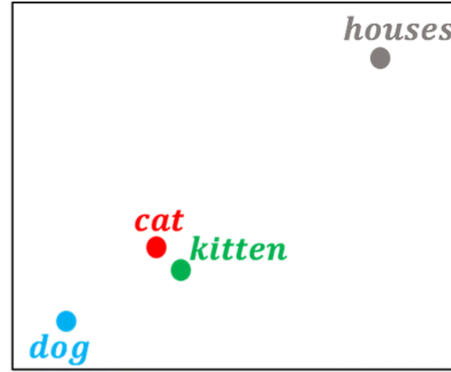
Goal: Given a word,
predict the surrounding context words.



Meaningful word embedding

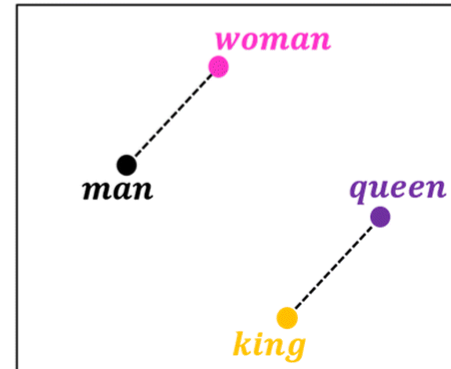
<i>cat</i> →	0.6	0.9	0.1	0.4	-0.7	-0.3	-0.2
<i>kitten</i> →	0.5	0.8	-0.1	0.2	-0.6	-0.5	-0.1
<i>dog</i> →	0.7	-0.1	0.4	0.3	-0.4	-0.1	-0.3
<i>houses</i> →	-0.8	-0.4	-0.5	0.1	-0.9	0.3	0.8

Dimensionality
reduction of
word
embeddings
from 7D to 2D

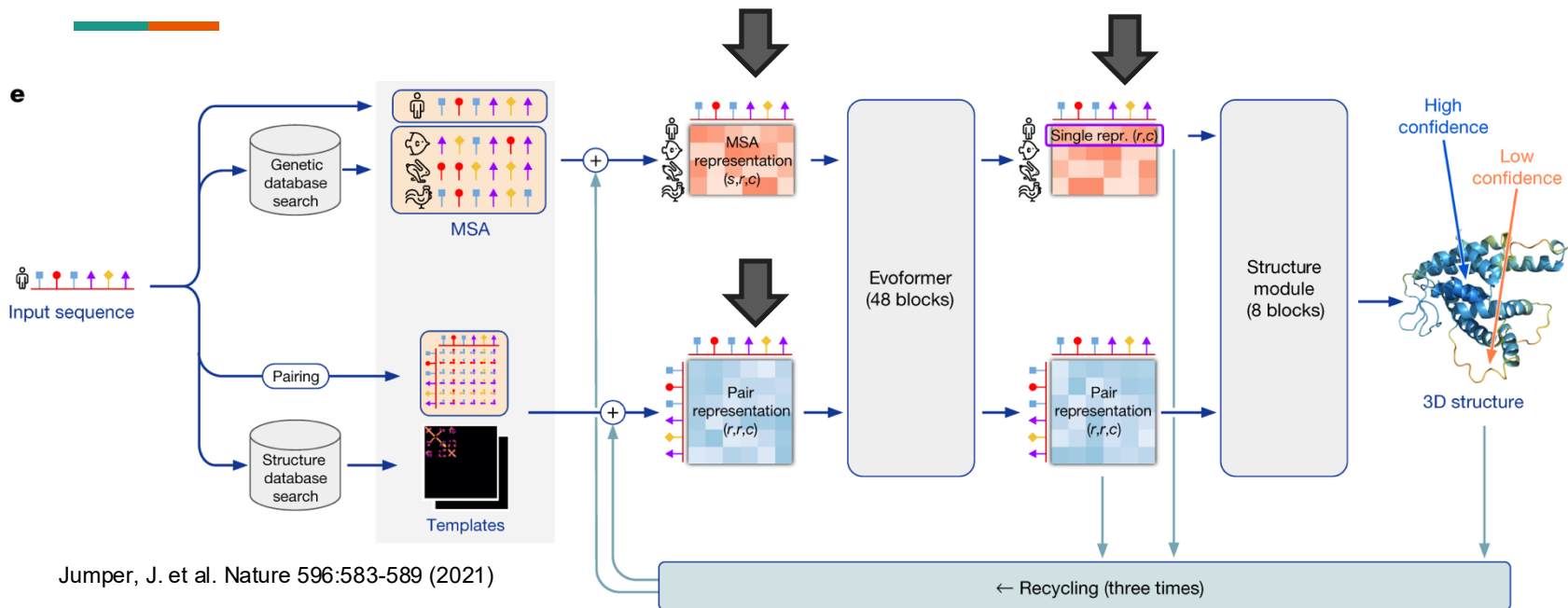


<i>man</i> →	0.6	-0.2	0.8	0.9	-0.1	-0.9	-0.7
<i>woman</i> →	0.7	0.3	0.9	-0.7	0.1	-0.5	-0.4
<i>king</i> →	0.5	-0.4	0.7	0.8	0.9	-0.7	-0.6
<i>queen</i> →	0.8	-0.1	0.8	-0.9	0.8	-0.5	-0.9

Dimensionality
reduction of
word
embeddings
from 7D to 2D



AlphaFold v2 architecture



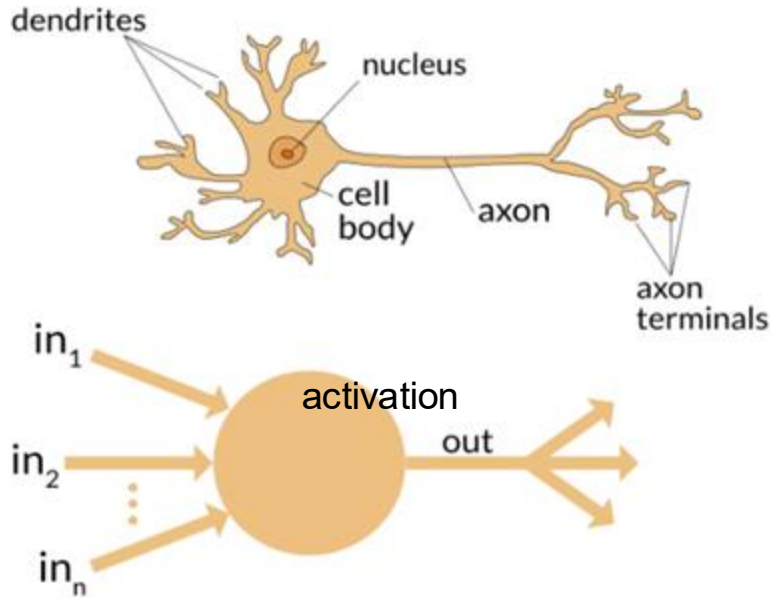
- Multiple representations of different biological concepts are trained to be computationally compatible inside the model



The rise of artificial neural network

An 80-year journey

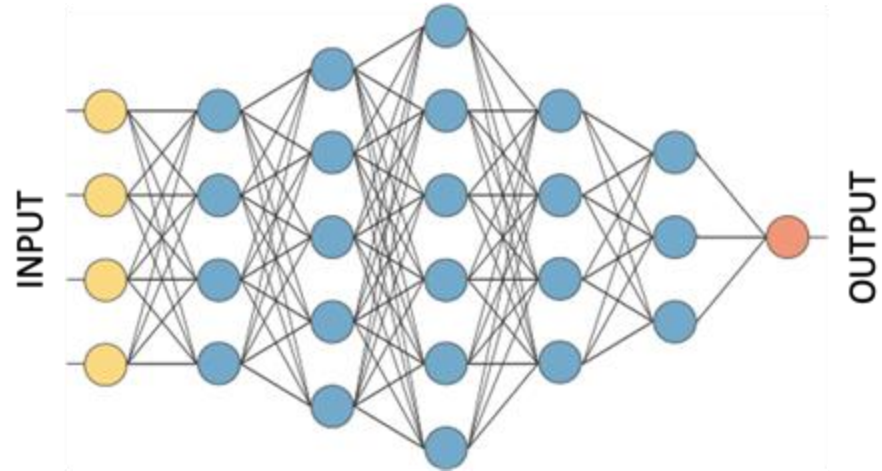
Inspired by biology



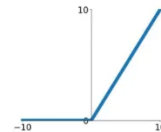
$$out = f(w_1 in_1 + w_2 in_2 + \dots + w_n in_n)$$

$f()$ is an activation function

Artificial Neural Network

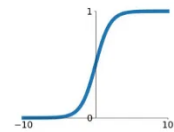


ReLU
 $\max(0, x)$



Sigmoid

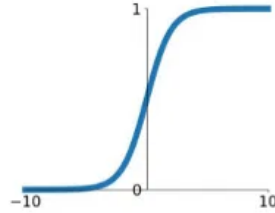
$$\sigma(x) = \frac{1}{1+e^{-x}}$$



Activation function

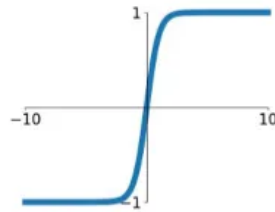
Sigmoid

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



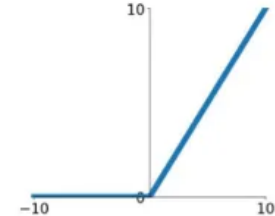
tanh

$$\tanh(x)$$



ReLU

$$\max(0, x)$$



- Simple, non-linear function
- Mimic the activation of a biological neuron
- Without non-linear activation function, ANN is just a linear regression

Universal approximation theorem (Cybenko 1989)

Universal Approximation Theorem: Fix a continuous function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ (activation function) and positive integers d, D . The function σ is not a polynomial if and only if, for every continuous function $f : \mathbb{R}^d \rightarrow \mathbb{R}^D$ (target function), every compact subset K of \mathbb{R}^d , and every $\epsilon > 0$ there exists a continuous function $f_\epsilon : \mathbb{R}^d \rightarrow \mathbb{R}^D$ (the layer output) with representation

$$f_\epsilon = W_2 \circ \sigma \circ W_1,$$

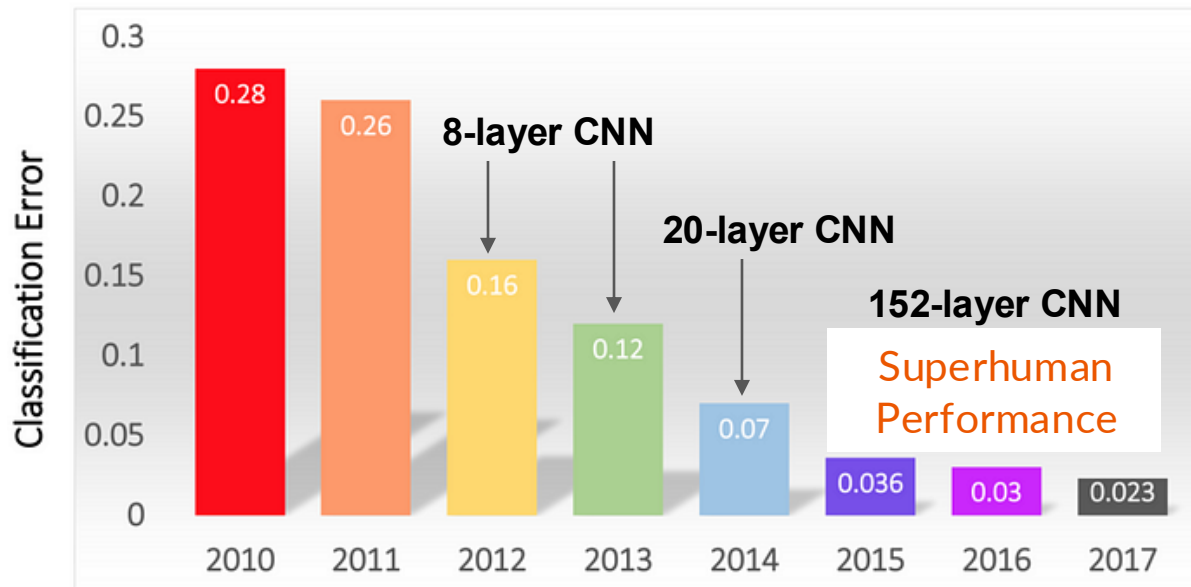
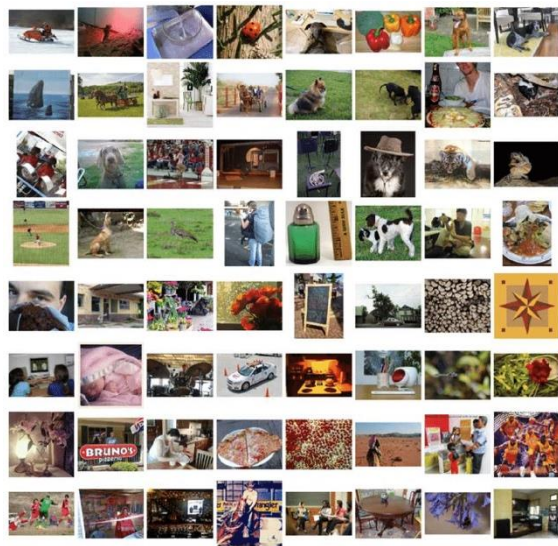
where W_2, W_1 are composable affine maps and \circ denotes component-wise composition, such that the approximation bound

$$\sup_{x \in K} \|f(x) - f_\epsilon(x)\| < \epsilon$$

holds for any ϵ arbitrarily small (distance from f to f_ϵ can be infinitely small).

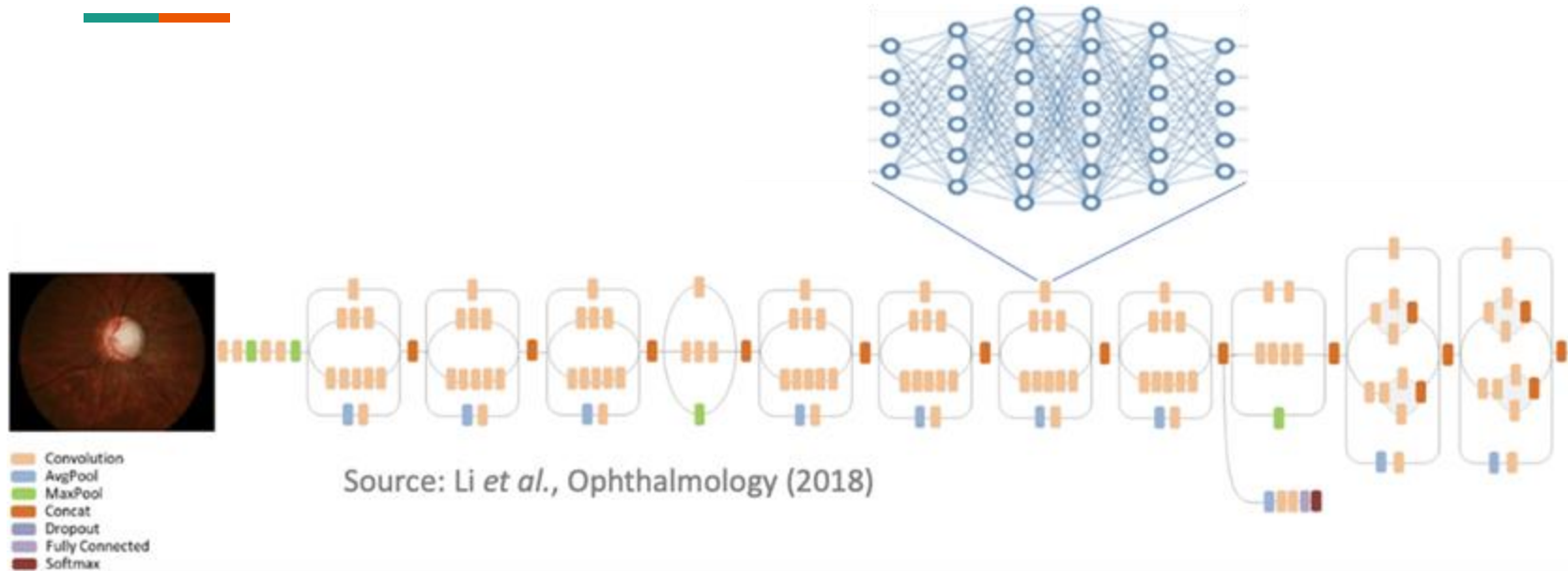
- ANN with just one layer of neurons and a non-polynomial activation function can capture any continuous mathematical relationship

ImageNet: a showcase of ANN on real-world data



- Emergence of neural network capability due to **data** (internet and digital technology) and **computing resources** (GPU)

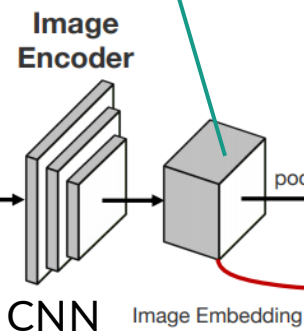
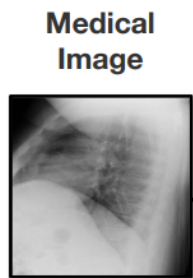
Deep learning is ML for deep neural network



- Explosion in dataset and model sizes
- Techniques for guiding the model to learn embedding

AI 3.0 learns and utilizes embedding

Capture key characteristics about the image – lesions?



Sentence Decoder

Word Decoder

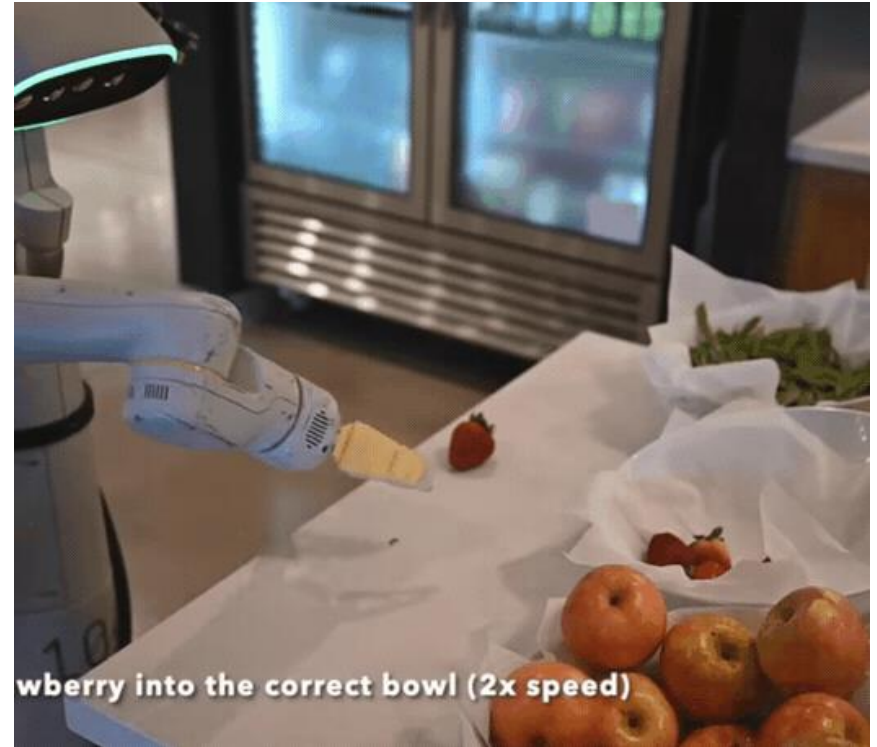
Generated Report


heart size is normal.
there is no focal consolidation,
effusion or pneumothorax.
the lungs are clear.
there is no acute osseous
abnormalities.

Co-utilization of
image and word
embeddings to
provide explanation

Recurrent Neural Network (RNN)

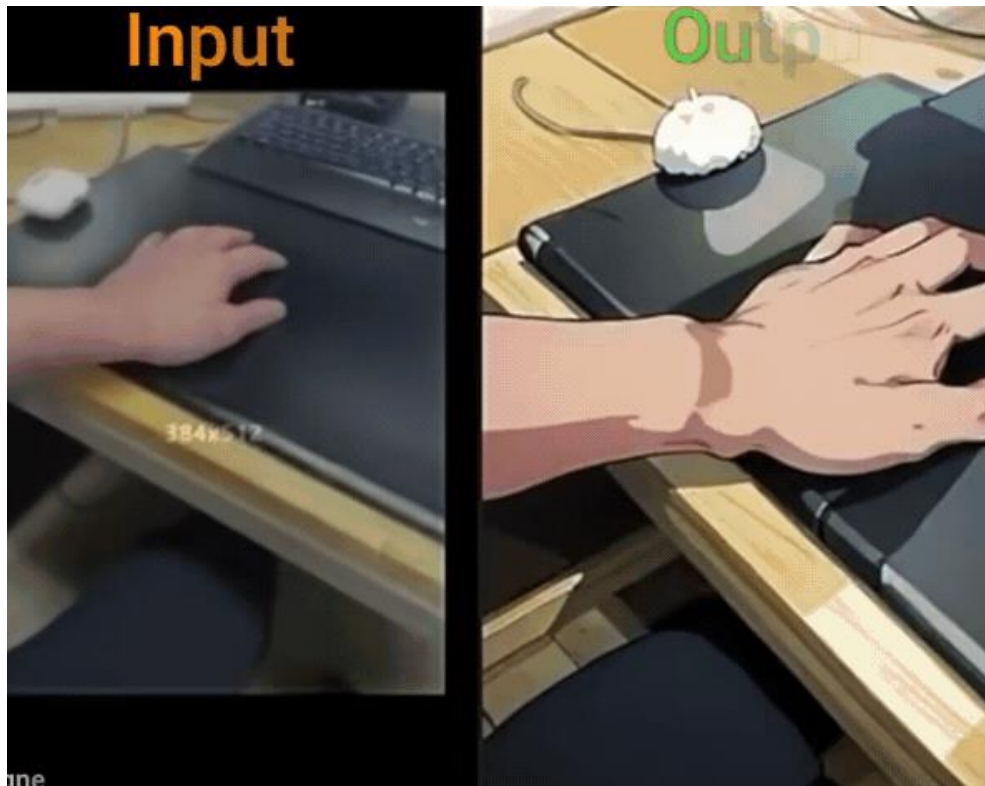
Linking visual and text embeddings with robot control





AI 4.0: Generative and foundational models

Going beyond prediction



Talking head anime: <https://github.com/pkhungurn>



Stable Diffusion: 8

The rise of large language model



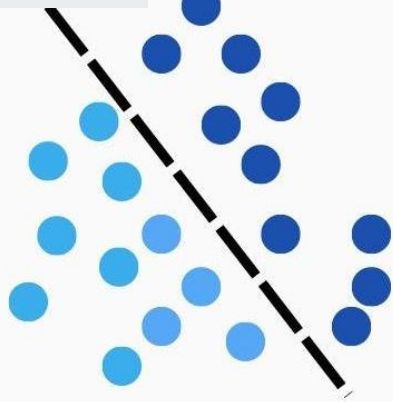
Upon examination of the chest X-ray image:

- The lungs are predominantly clear without any obvious consolidations, masses, or pneumothoraces.
- The cardiac silhouette appears within normal size limits.
- The bony thorax, including the ribs and clavicles, appears intact without any visible fractures.
- The diaphragm and costophrenic angles are well visualized and appear normal.
- There is no visible mediastinal widening or significant lymphadenopathy.

<https://xrayinterpreter.com/resource/how-to-use-chatgpt-to-interpret-xrays>

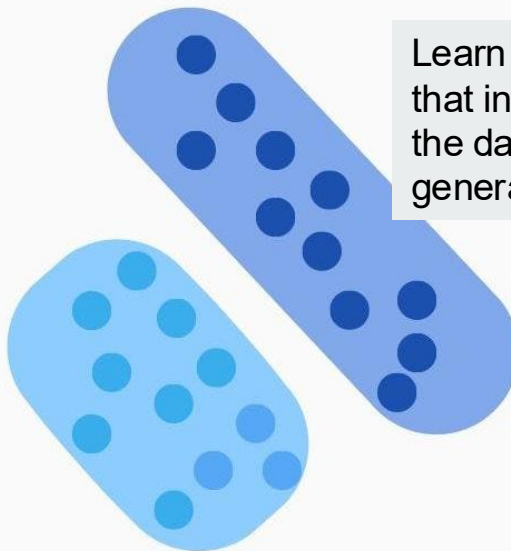
Importance of generative approach

Simple, multiple
equal solutions



Discriminative

Learn factors
that influence
the data during
generation



Generative



VS

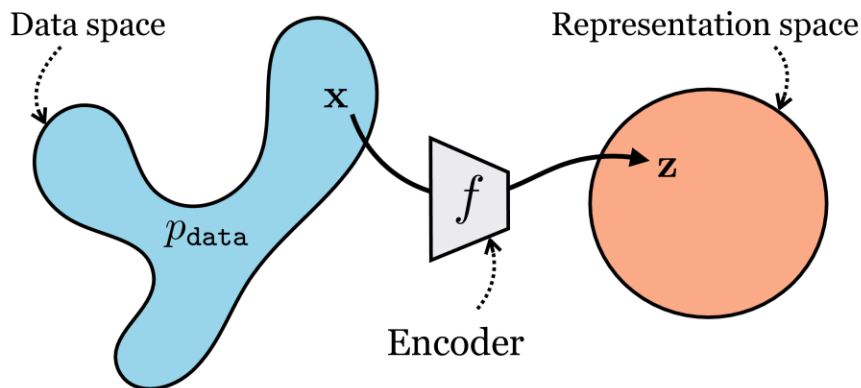


<https://www.turing.com/kb/generative-models-vs-discriminative-models-for-deep-learning>

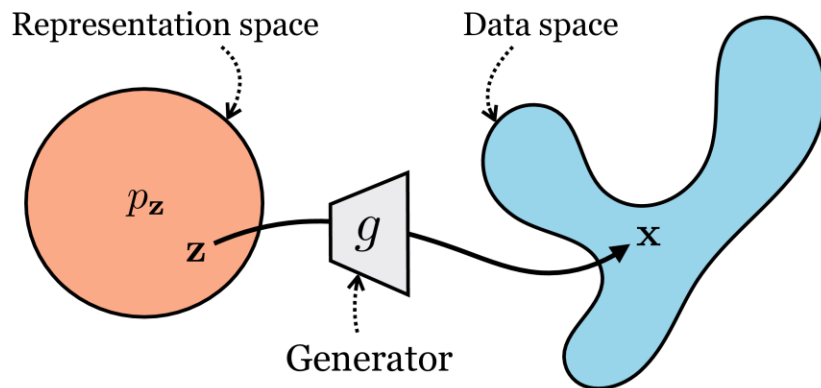
- It takes much more understanding to generate **realistic** data

Reversed representation learning

Representation learning



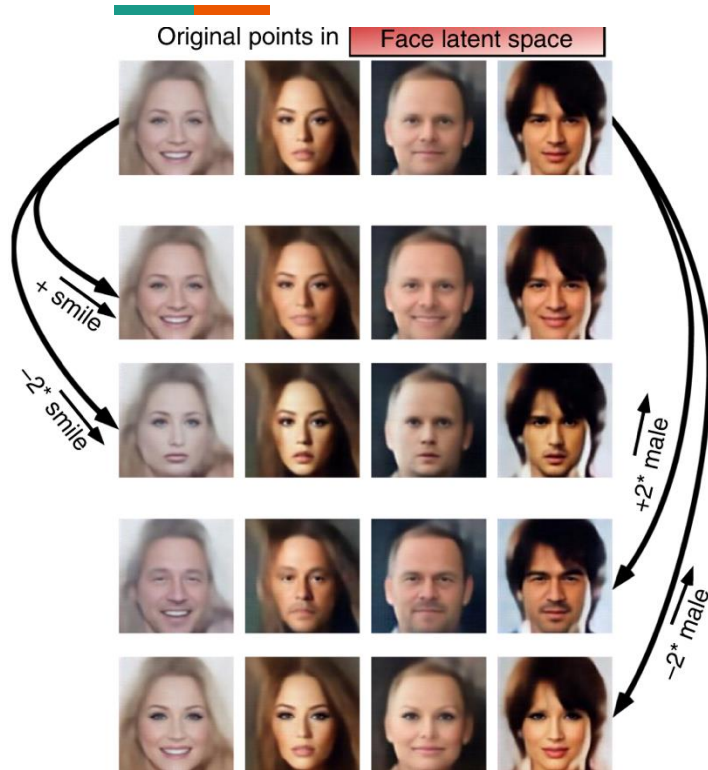
Generative modeling



https://visionbook.mit.edu/generative_modeling_and_rep_learning.html

- Generative model can be thought of as reversed representation learning
 - Compress vs decompress
- How do we train, or “guide”, the generative process?

Interpretable generative process



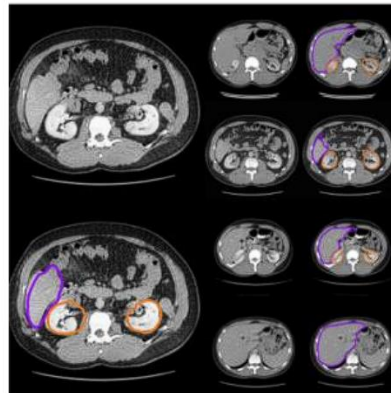
- **Assumption:** Generative model helps us understand the underlying mechanisms and the factors that generated the data
- Disentangle factors from observation
- Counterfactual “what-if” analysis

RadImageGAN

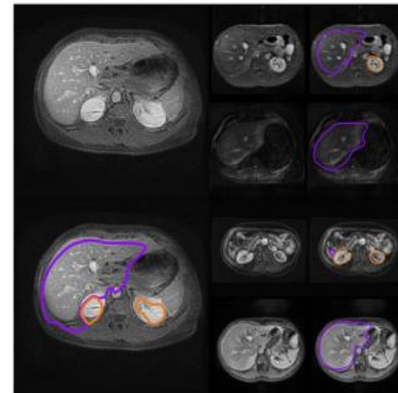


- StyleGAN on radiological images from 100,000 patients
- 12 anatomical regions and 130 pathological classes
- Can generate both images and labeled masks to train future AI models

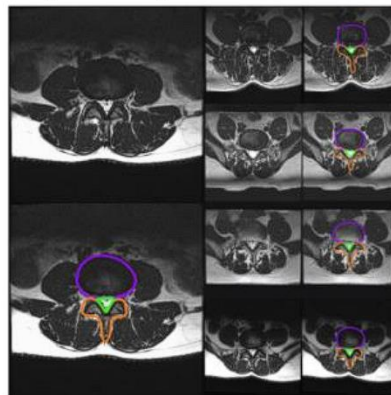
A) CT Abdomen



B) MRI Abdomen



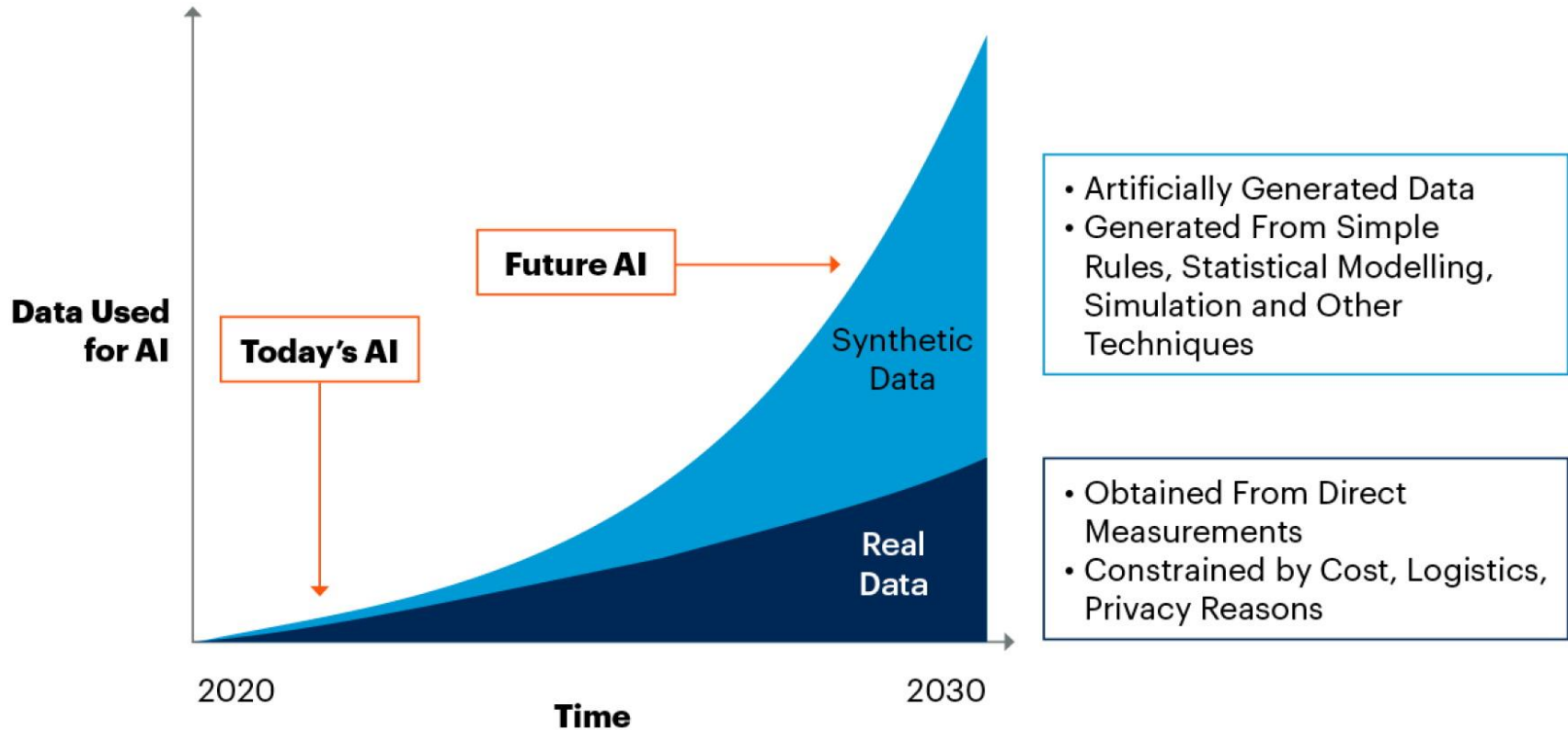
C) MRI Spine



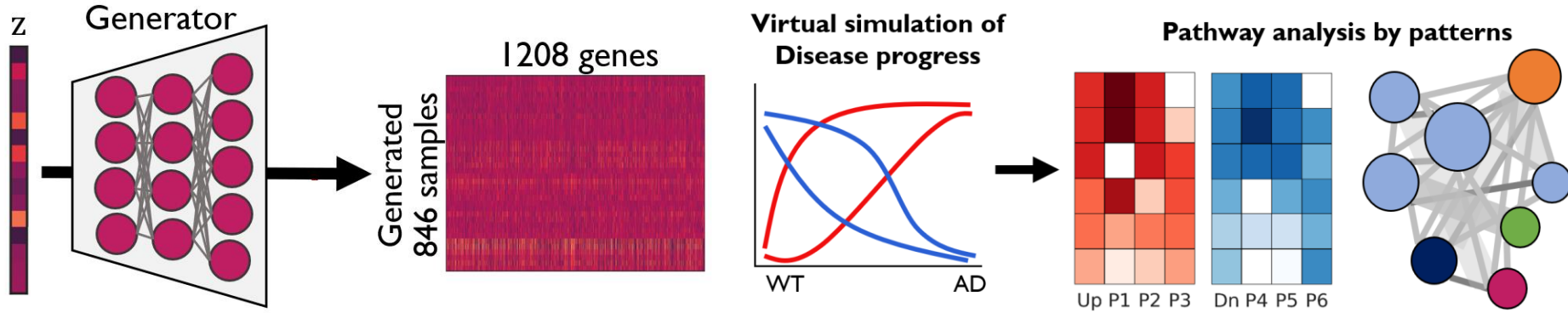
D) Colonoscopy Polyp



By 2030, Synthetic Data Will Completely Overshadow Real Data in AI Models



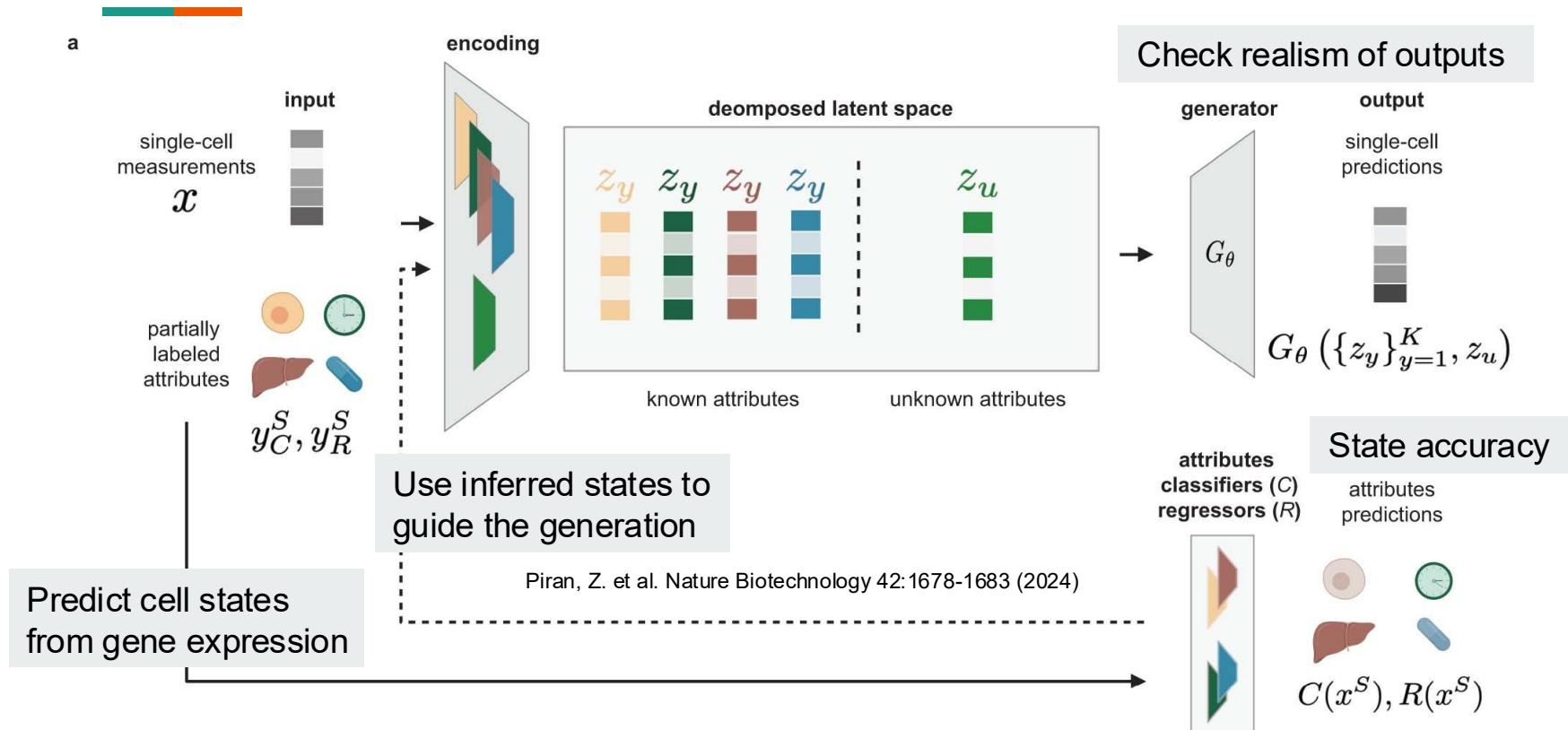
Knowledge from synthetic data



Park, J. et al. PLoS Computational Biology 16:e1008099 (2020)

- Train a generative model with data from small-scale experiment
- Simulate time-course transcriptomics data
- Analyze simulated data to gain insights into disease progression
- Validate biomarkers in external datasets

Disentangle factors influencing gene expression



Summary



- Early AI were driven by human knowledge and fixed rules
- The emergence of digital computer and internet enabled data-driven, machine learning approach to AI
- New hardware and training techniques gave rise to the ability of AI to learn representations by themselves
- Learning through generation allows modern AI to mimic complex mechanisms behind the data and disentangle important factors

Any question?



- See you next time