



Lesson 01

Introduction to Machine Learning

План лекції

- ▶ Привітання та знайомство
- ▶ Огляд курсу
- ▶ Що таке AI/DS/ML/DL?
- ▶ Задачі, що вирішує ML
- ▶ Етапи реалізації проектів з ML
- ▶ Збір та обробка даних в ML, джерела даних та їх особливості
- ▶ Огляд основних інструментів в ML з використанням Python

Олег Коменчук

Machine Learning Engineer (4 years of commercial experience)

PhD Student in Information Systems and Technologies (2021 -
Present)

Kaggle Notebooks Expert

[LinkedIn](#) | [Kaggle](#) | [GitHub](#)

Знайомство зі студентами

- ▶ Інформація про себе
- ▶ Освіта та сфера інтересів
- ▶ Робочий досвід (якщо є)
- ▶ Попередні курси чи проєкти, досвід програмування на Python
- ▶ Очікування від курсу

Комунікація



Telegram-група: <https://t.me/+RVnrUzqmZZcxMGUy>

Огляд курсу

23 заняття

2 заняття на тиждень в Zoom

Вівторок, П'ятниця (19:15 - 21:15)

Домашні завдання:

- ▶ Практичні завдання
- ▶ Тест
- ▶ Дедлайни
- ▶ Перездачі
- ▶ Оцінювання

- ▶ Introduction to Machine Learning
- ▶ NumPy, Pandas, Matplotlib
- ▶ Linear & Logistic Regression
- ▶ Logistic Regression for Multiclass Classification
- ▶ Regularization
- ▶ Tree Based Models
- ▶ Scikit-Learn Workflow
- ▶ Intro to Deep Learning: PyTorch
- ▶ Intro to Deep Learning: Layers
- ▶ Intro to Deep Learning: Optimization
- ▶ Computer Vision: Intro
- ▶ Computer Vision: Classification Models
- ▶ Computer Vision: Segmentation Models
- ▶ Computer Vision: Object Detection
- ▶ Natural Language Processing: Intro
- ▶ Natural Language Processing: Embeddings
- ▶ Natural Language Processing: RNNs
- ▶ Transformers
- ▶ LLMs, LLM Prompting
- ▶ Recommender Systems
- ▶ Autoencoders
- ▶ Generative Adversarial Networks (GANs)
- ▶ Machine Learning Model Deployment

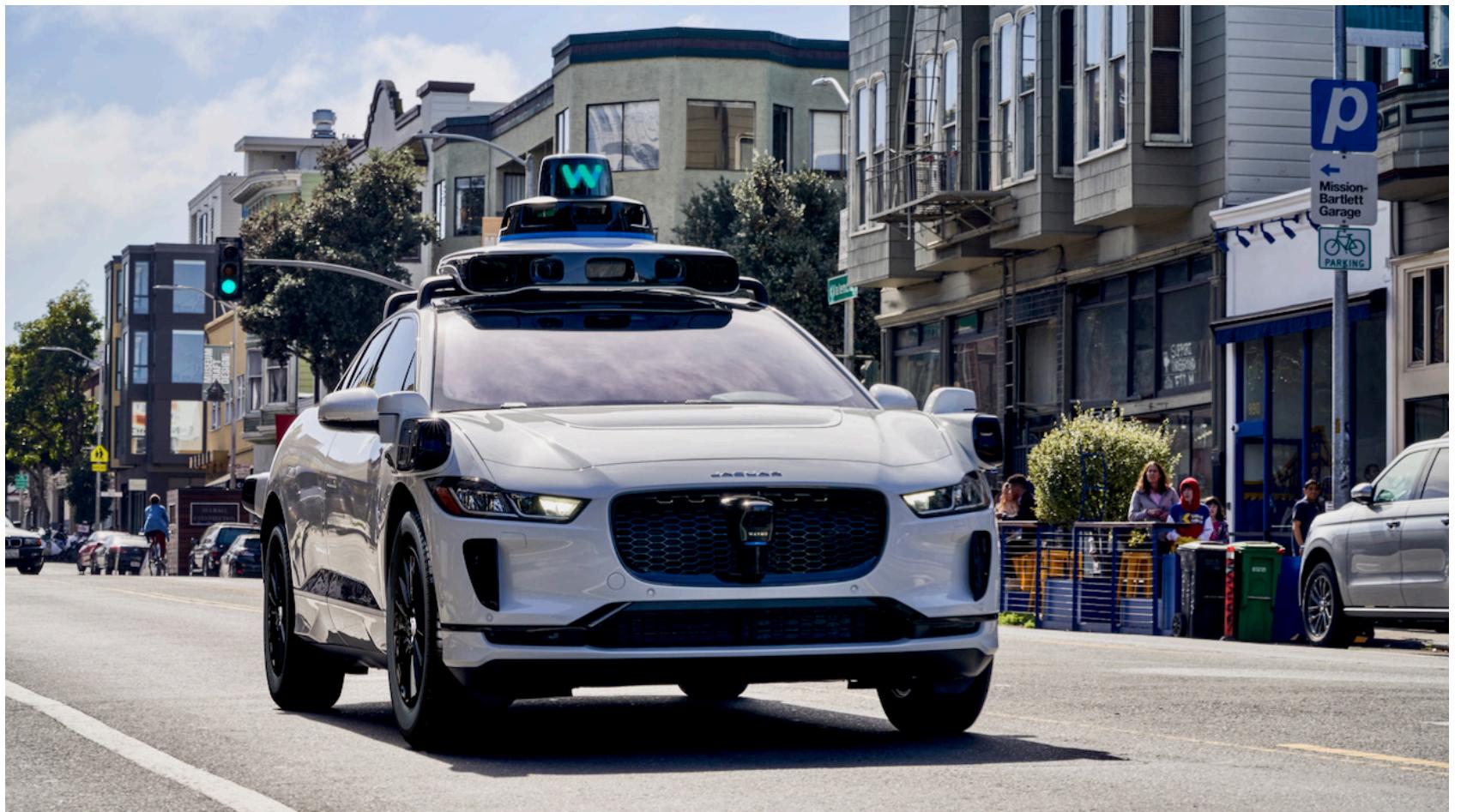
What is Artificial Intelligence, Machine Learning, Deep Learning, and Data Science?

AI/ML: Self-Driving Cars

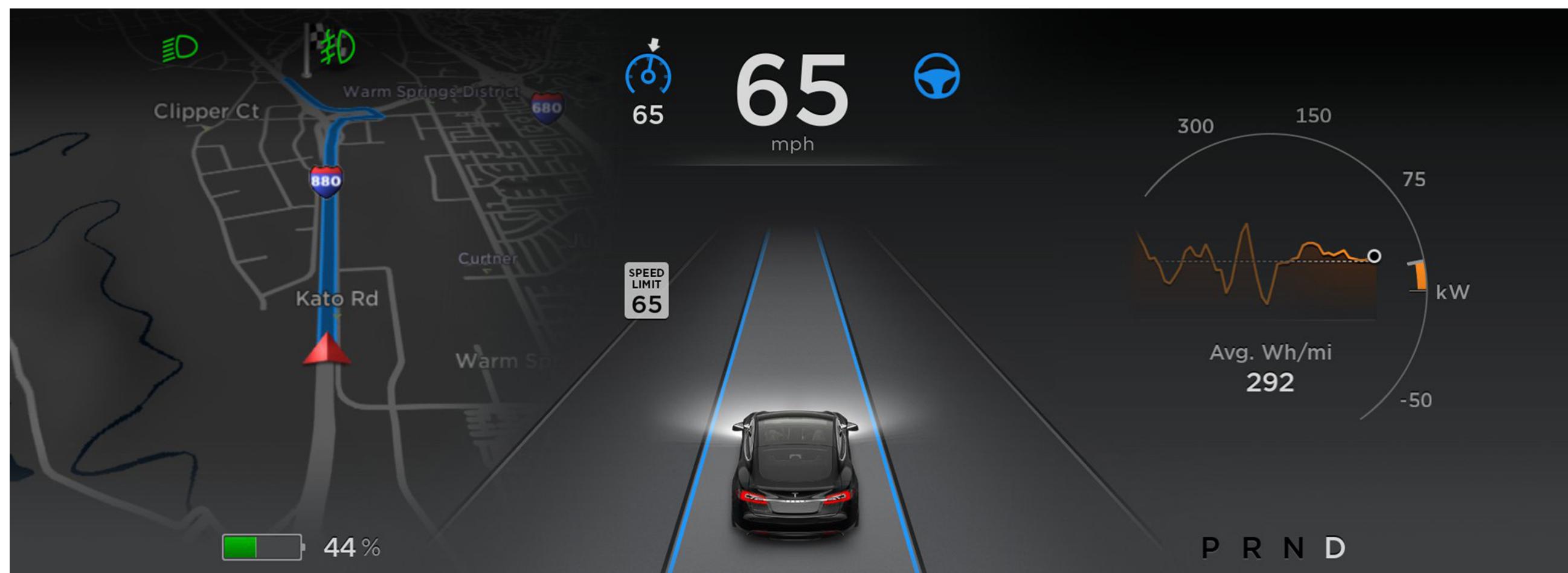
Tesla Autopilot



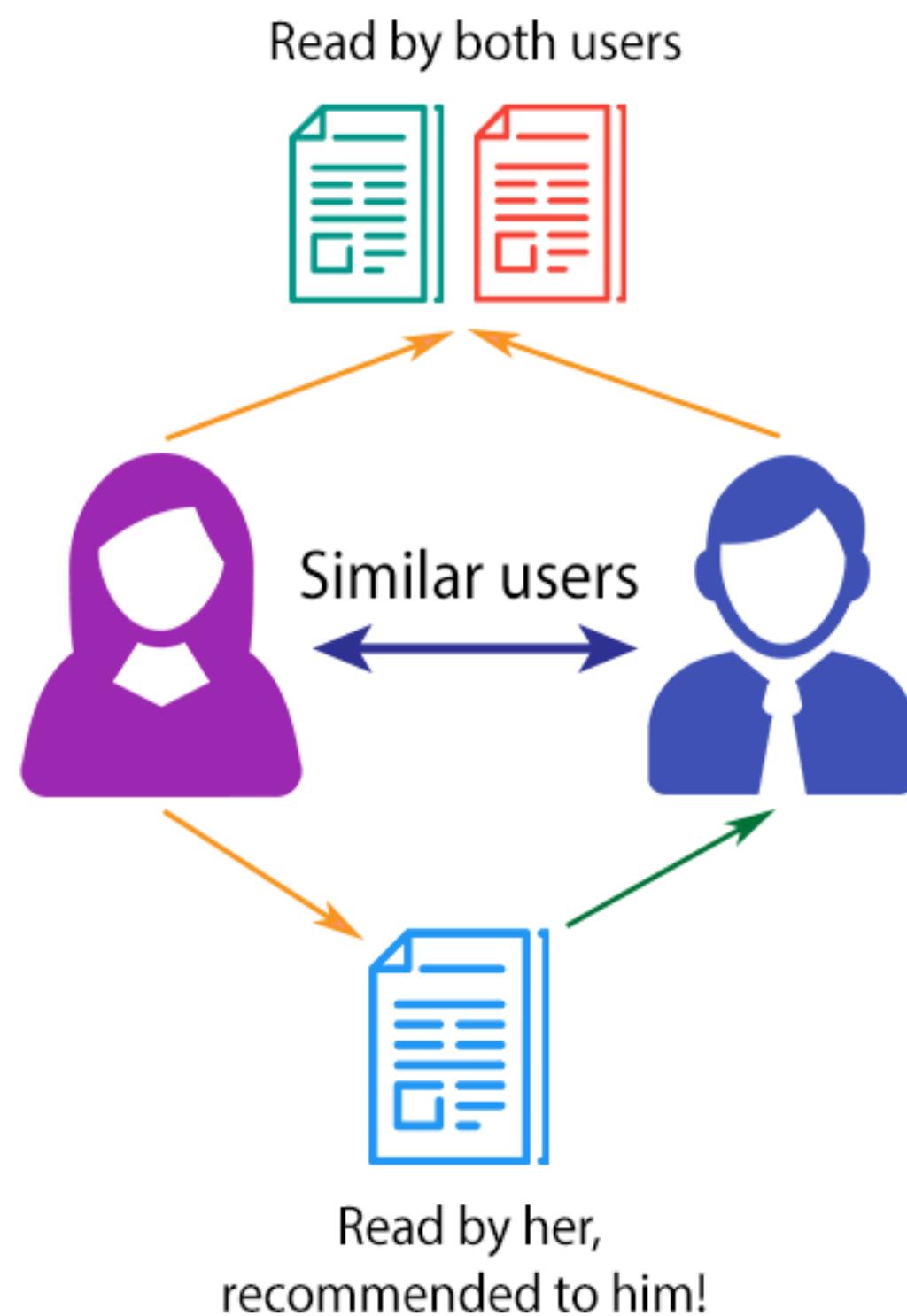
Waymo



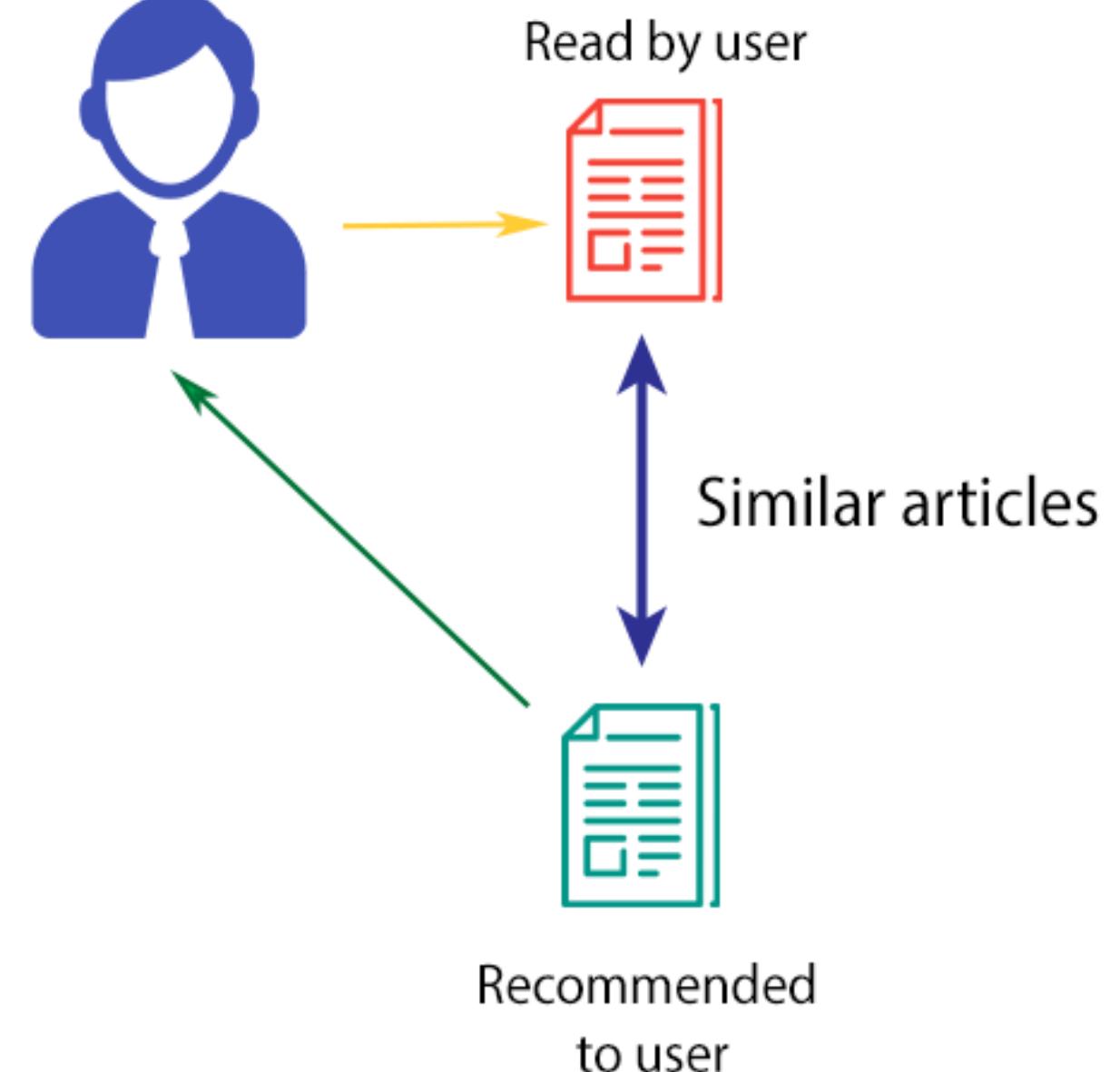
Mobileye



COLLABORATIVE FILTERING



CONTENT-BASED FILTERING



NETFLIX



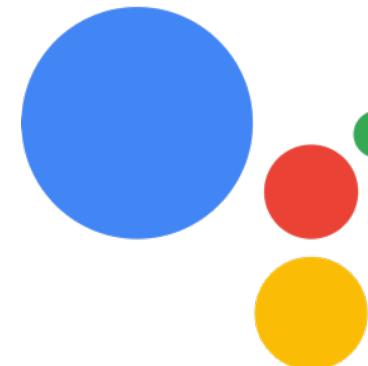
Google





Apple Siri

Siri is a virtual assistant that is part of Apple's iOS, watchOS, macOS, and tvOS operating systems. The assistant uses voice queries and a natural-language user interface to answer questions, make recommendations, and perform actions by delegating requests to a set of Internet services.



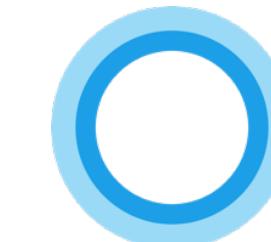
Google Assistant

Google Assistant is an AI-powered virtual assistant developed by Google that is primarily available on mobile and smart home devices, and can engage in two-way conversations.



Amazon Alexa

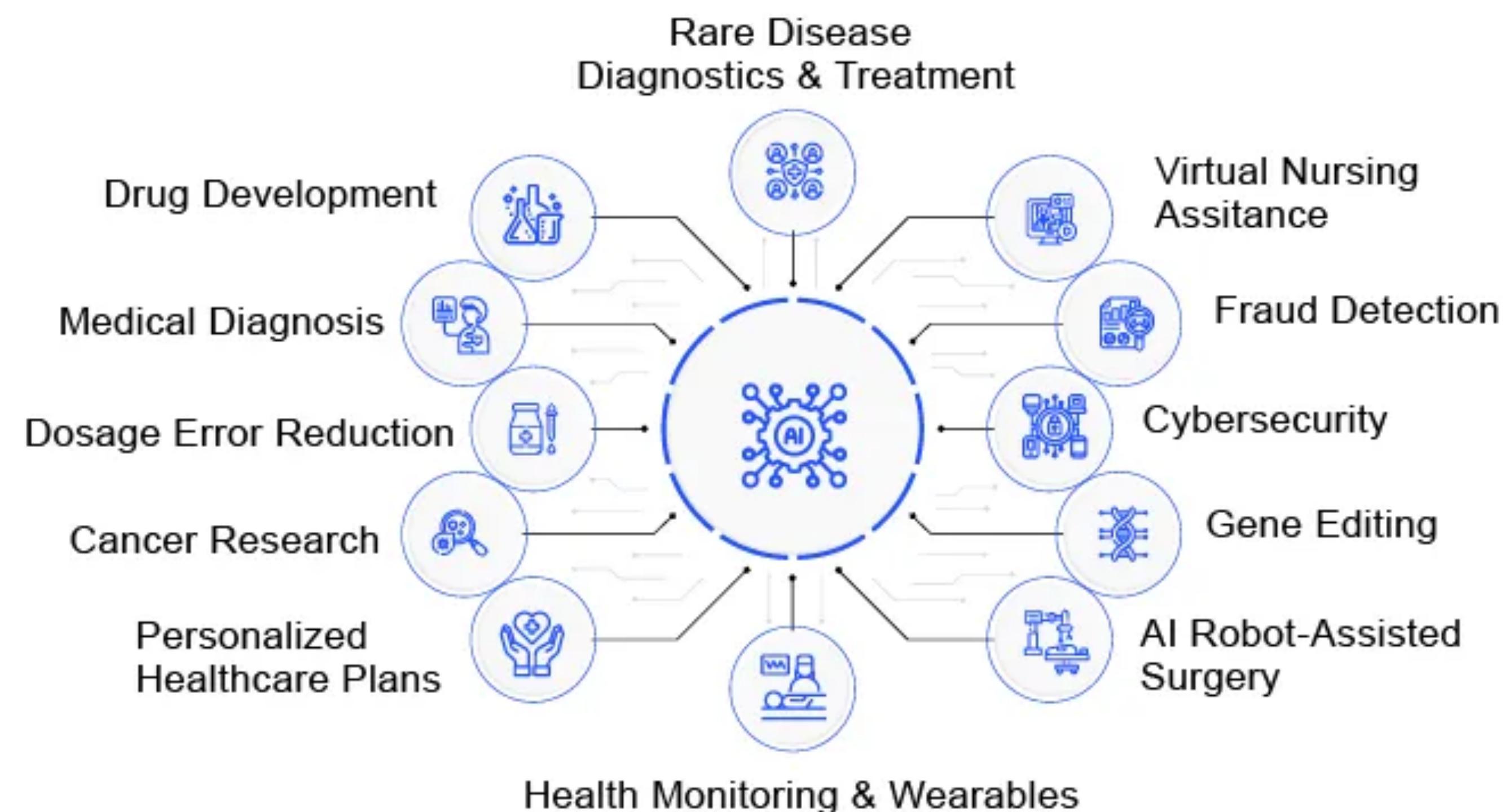
Alexa is a virtual assistant developed by Amazon. It is capable of voice interaction, music playback, making to-do lists, etc.



Microsoft Cortana

Cortana is a virtual assistant created by Microsoft for Windows 10, Windows 10 Mobile, Windows Phone 8.1. It can set reminders, recognize natural voice without the requirement for keyboard input, and answer questions using information from the Bing search engine.

Applications of AI in Healthcare



AI/ML: Generative AI



ChatGPT is a chatbot developed by OpenAI and launched on November 30, 2022. Based on a large language model, it enables users to refine and steer a conversation towards a desired length, format, style, level of detail, and language.



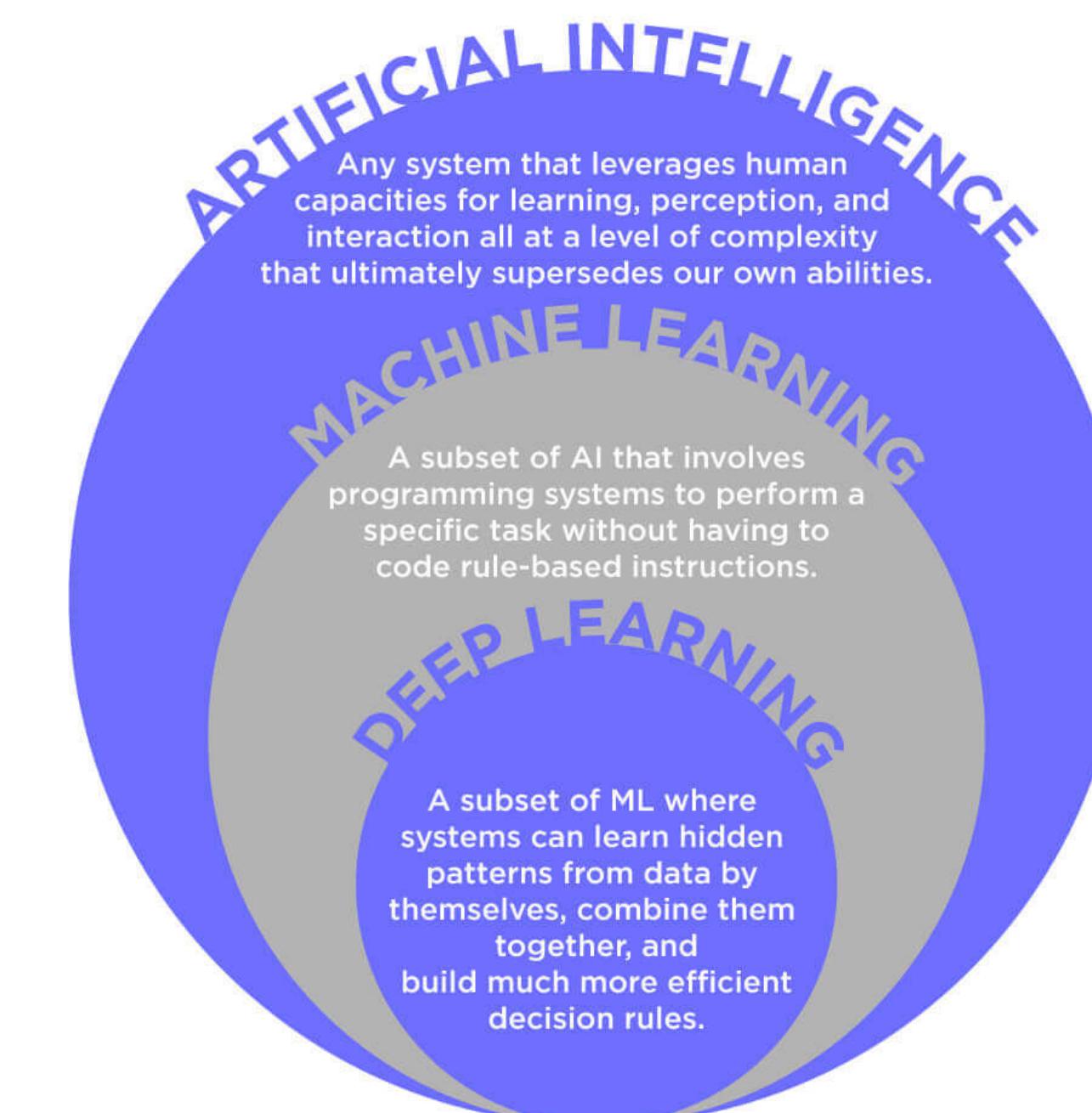
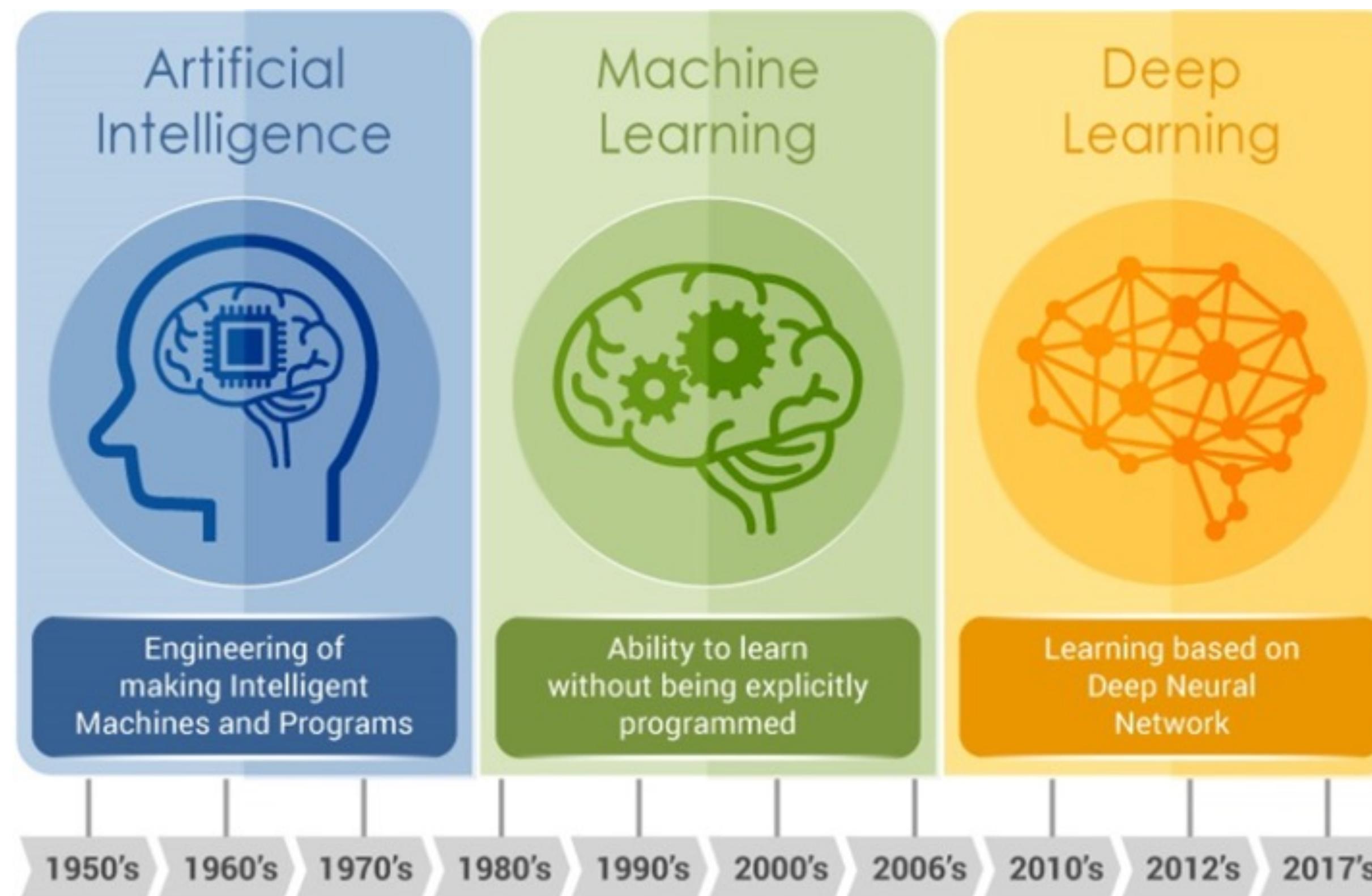
DALL-E, DALL-E 2, and DALL-E 3 are text-to-image models developed by OpenAI using deep learning methodologies to generate digital images from natural language descriptions, called "prompts."



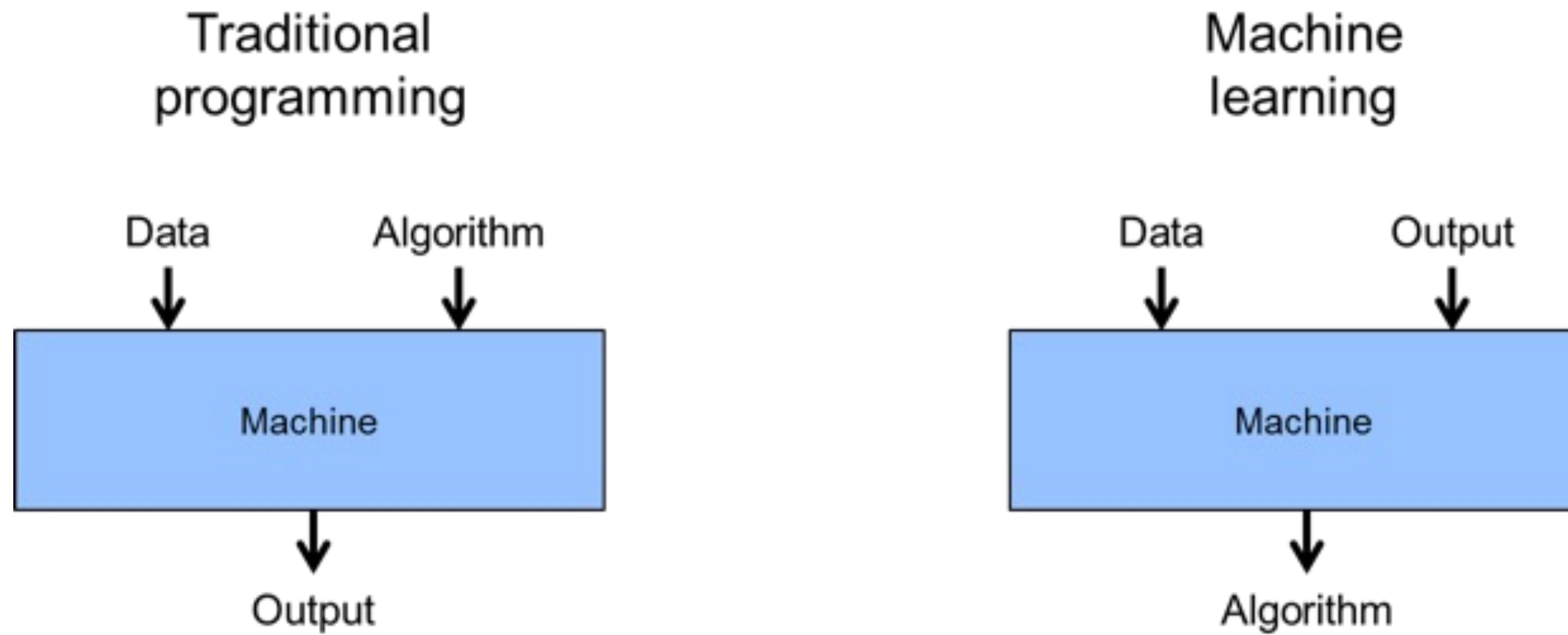
Gemini, formerly known as Bard, is a generative artificial intelligence chatbot developed by Google.



Midjourney is a generative artificial intelligence program and service created and hosted by the San Francisco-based independent research lab Midjourney, Inc. Midjourney generates images from natural language descriptions, called prompts, similar to OpenAI's DALL-E and Stability AI's Stable Diffusion.



The difference between traditional programming and machine learning



Math in ML

$$\begin{bmatrix} a_1 & a_2 & a_3 \\ a_4 & a_5 & a_6 \\ a_7 & a_8 & a_9 \end{bmatrix} \begin{bmatrix} b_1 & b_2 & b_3 \\ b_4 & b_5 & b_6 \\ b_7 & b_8 & b_9 \end{bmatrix} = \begin{bmatrix} c_1 & c_2 & c_3 \\ c_4 & c_5 & c_6 \\ c_7 & c_8 & c_9 \end{bmatrix}$$

- Linear Algebra
- Mathematical Analysis
- Calculus
- Statistics

$$\frac{d(\sin x)}{x} = \cos x$$

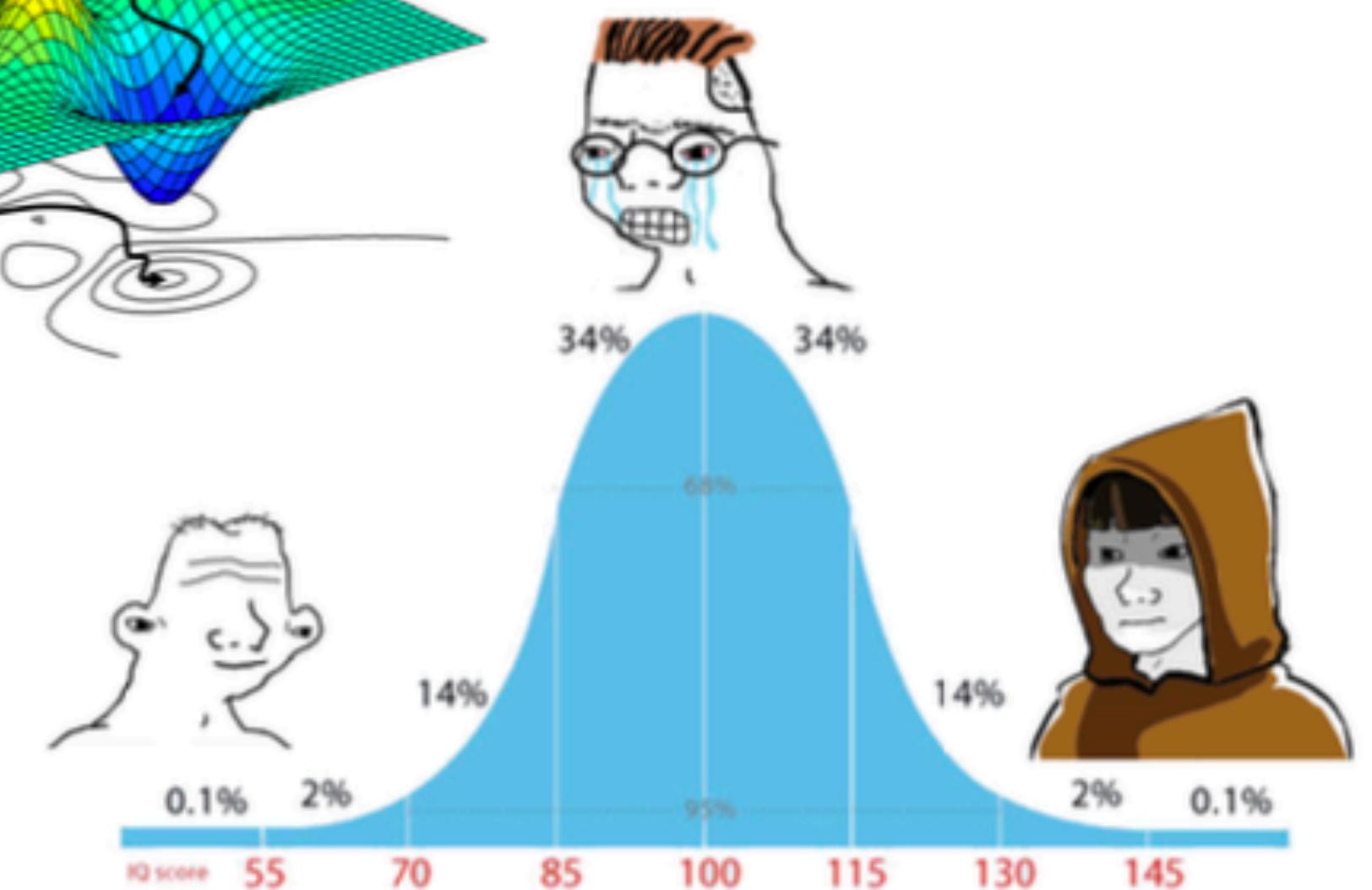
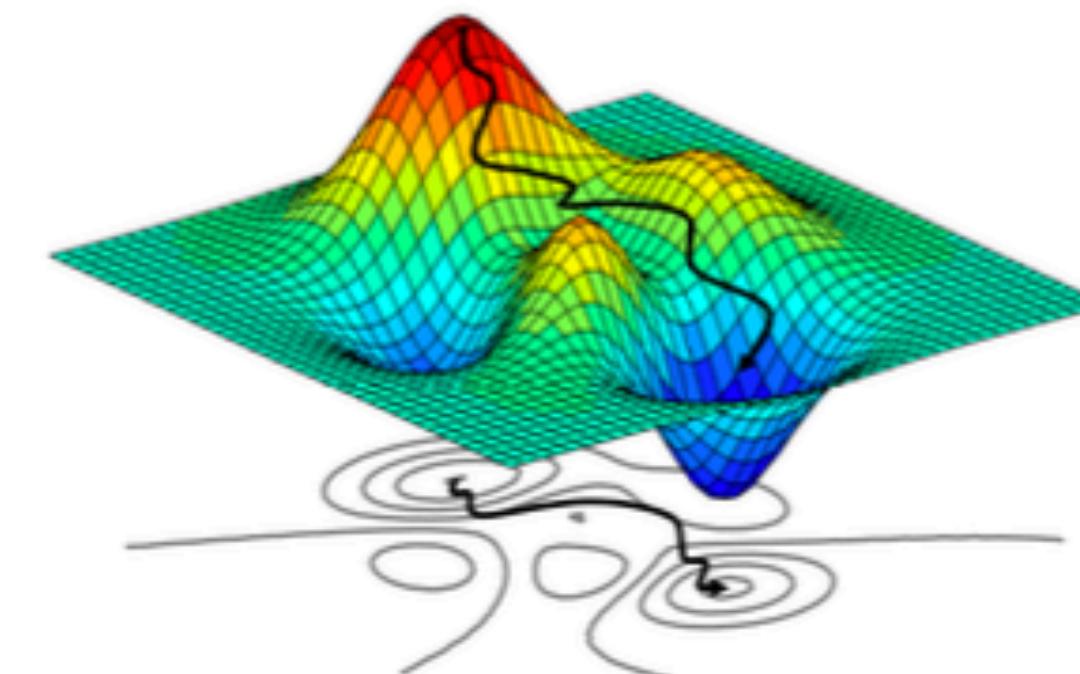
$$\frac{d(\cos x)}{x} = -\sin x$$

$$\frac{d(\tan x)}{x} = -\sec^2$$

$$\frac{d(\cot x)}{x} = \operatorname{cosec}^2$$

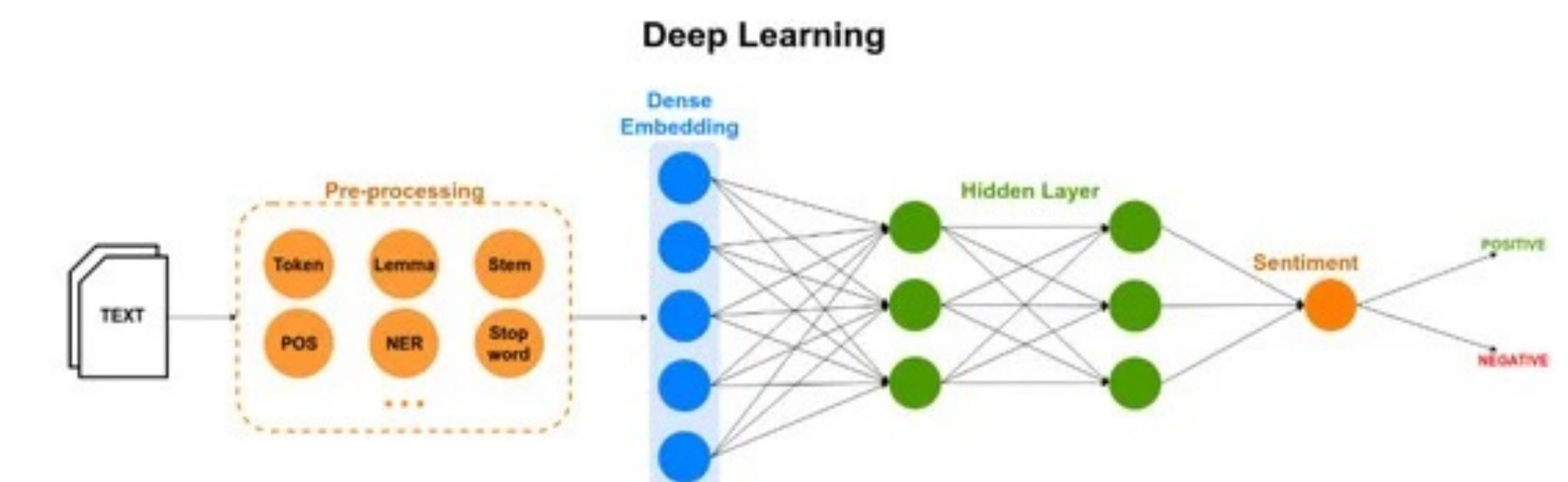
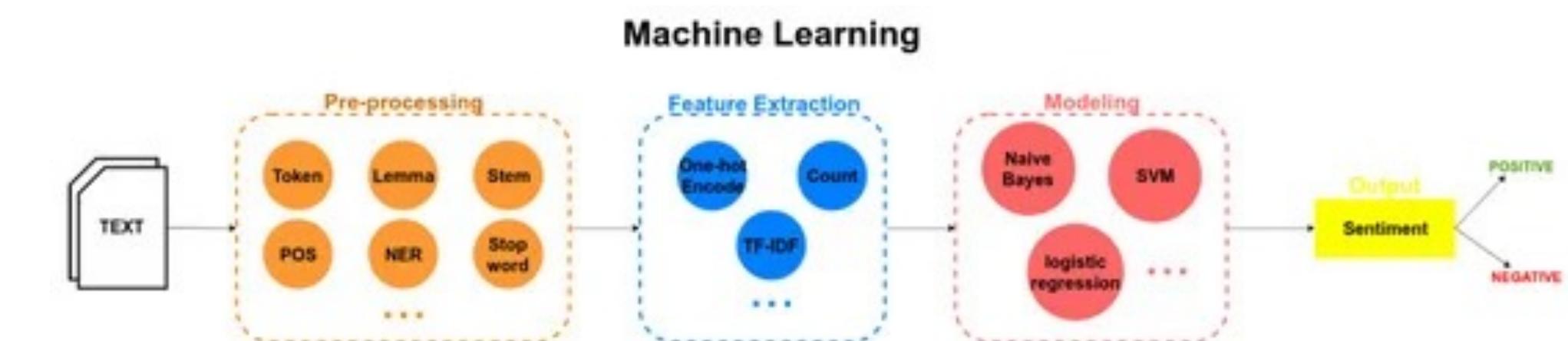
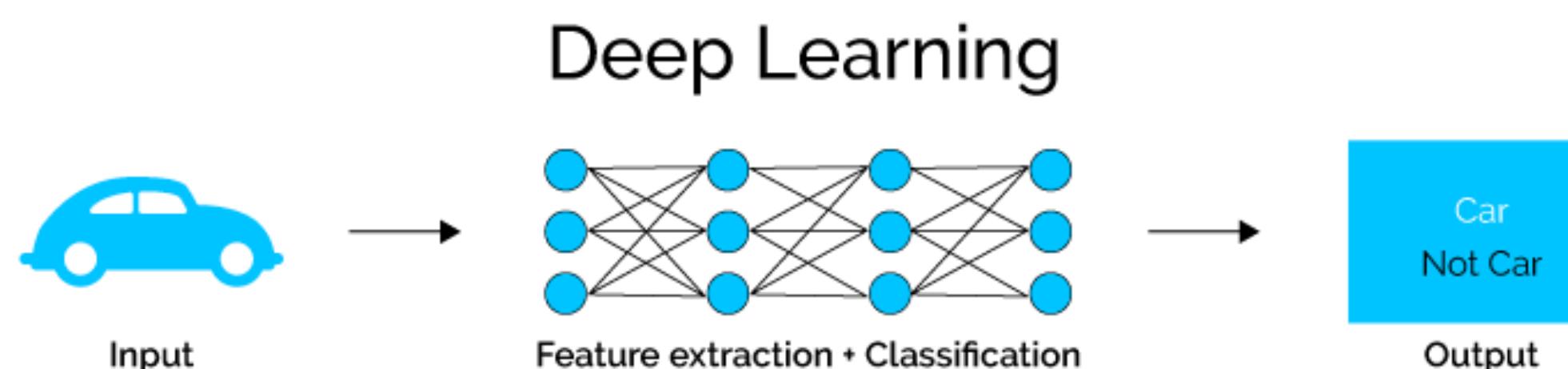
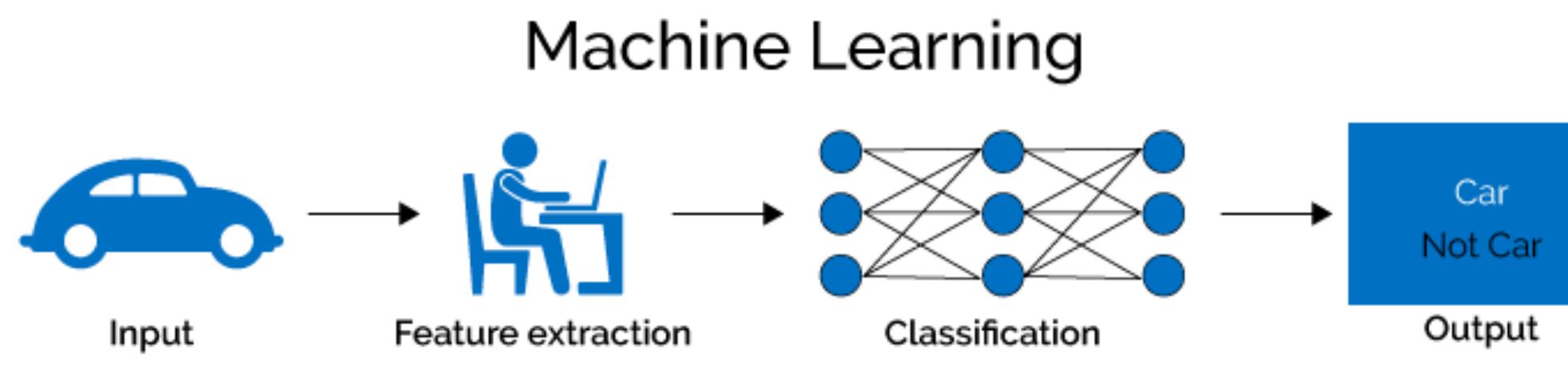
$$\frac{d(\sec x)}{x} = \sec x \cdot \tan x$$

$$\frac{d(\operatorname{cosec} x)}{x} = -\operatorname{cosec} x \cdot \cot x$$

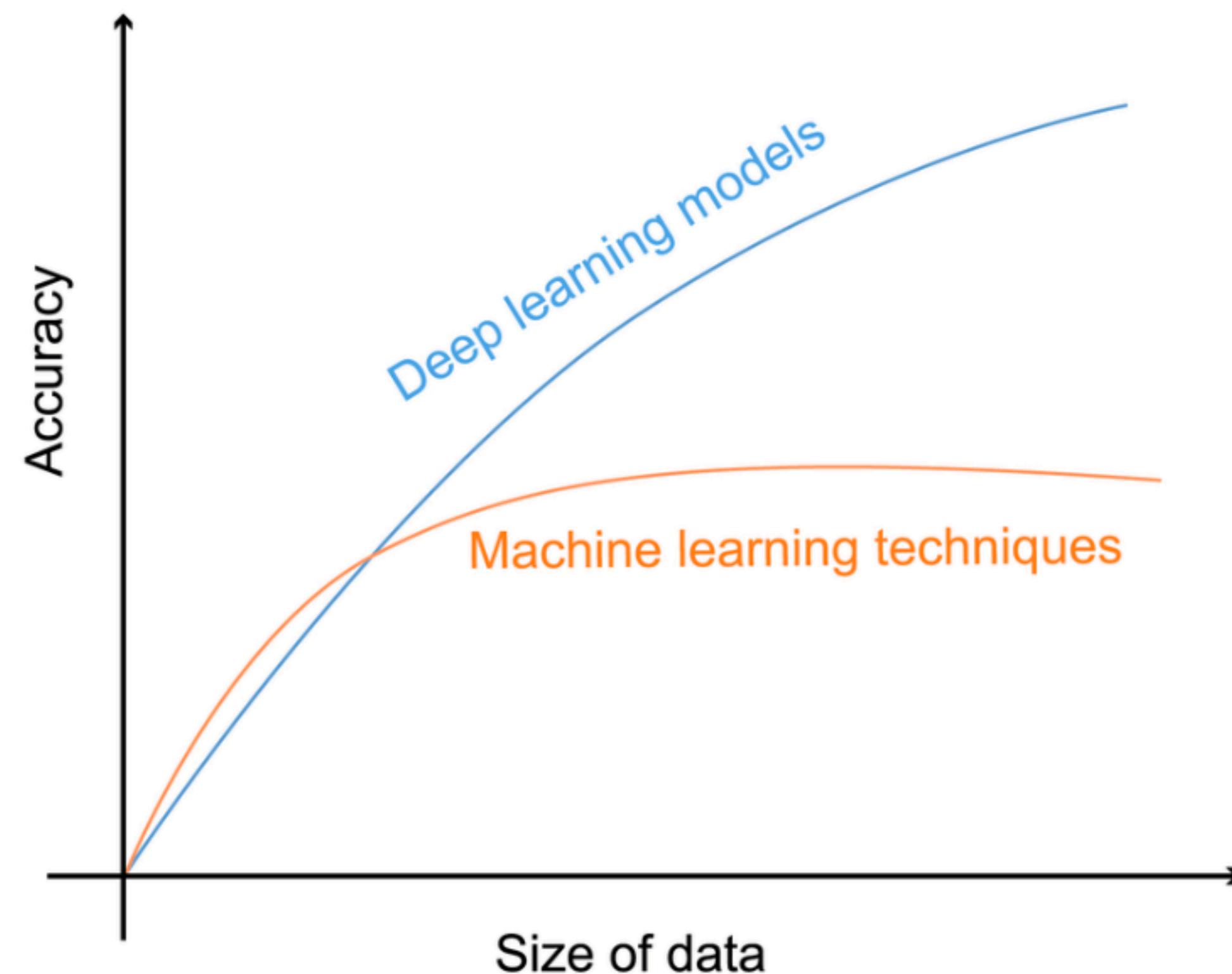


Deep Learning

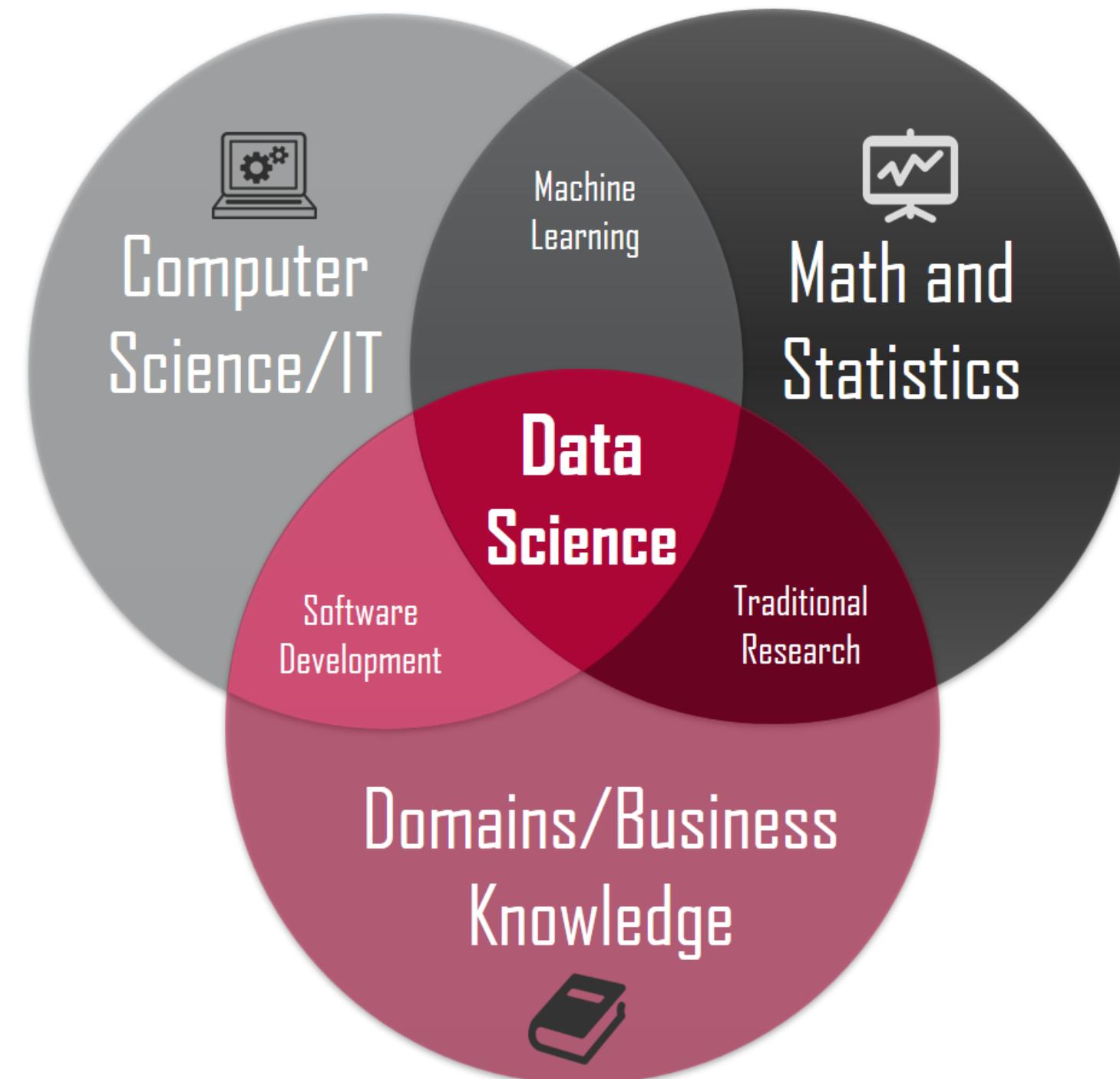
Deep Learning refers to training Artificial Neural Networks, sometimes very large neural networks. Today, the terms Neural Networks and Deep Learning are used almost interchangeably.



Why DL ?



What is Data Science?



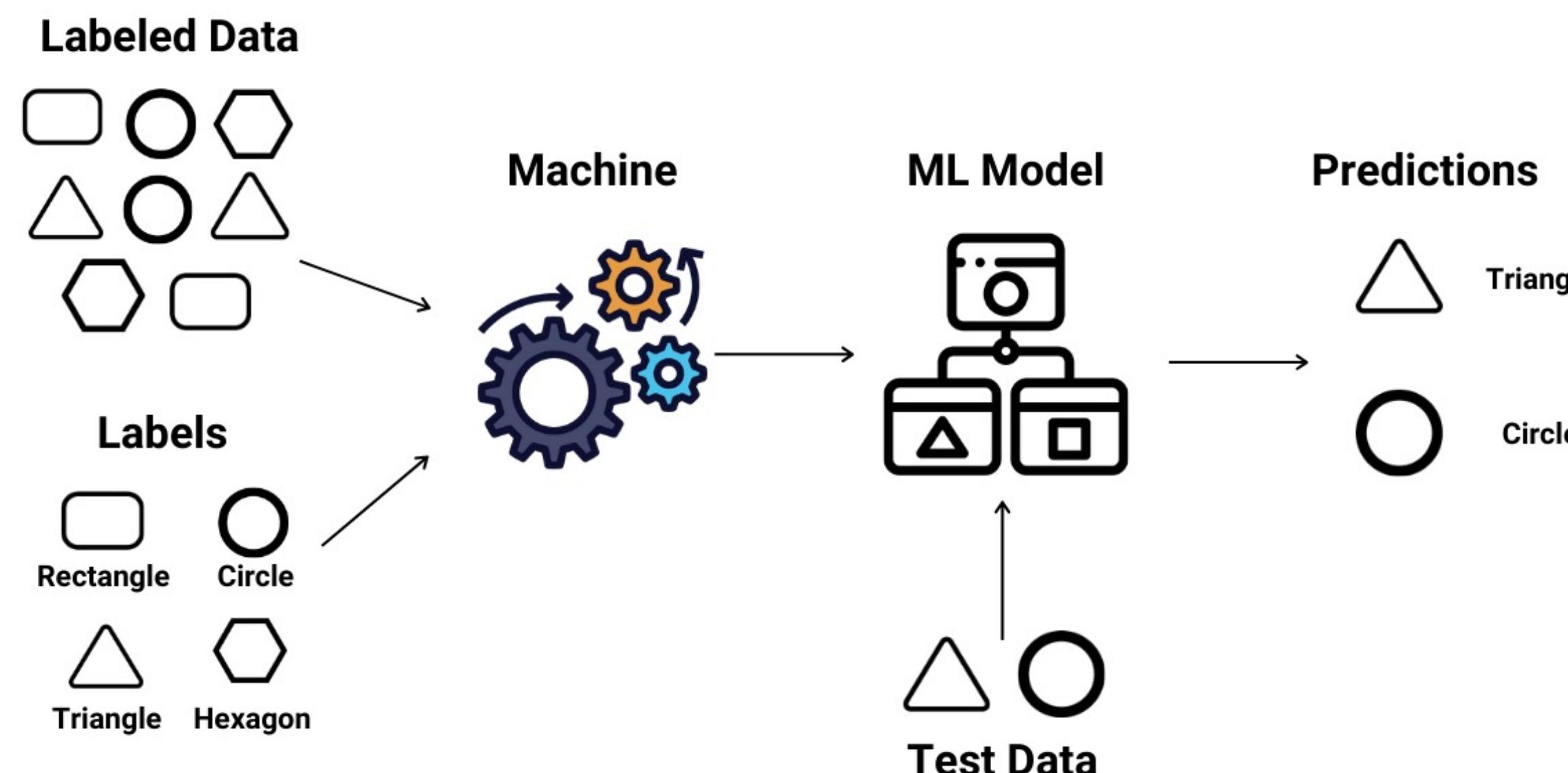
Data science is the domain of study that deals with vast volumes of data using modern tools and techniques to find unseen patterns, derive meaningful information, and make business decisions. Data science uses complex machine learning algorithms to build predictive models. The data used for analysis can come from many different sources and be presented in various formats.

In general, any machine learning problem can be assigned to one of three broad classifications:

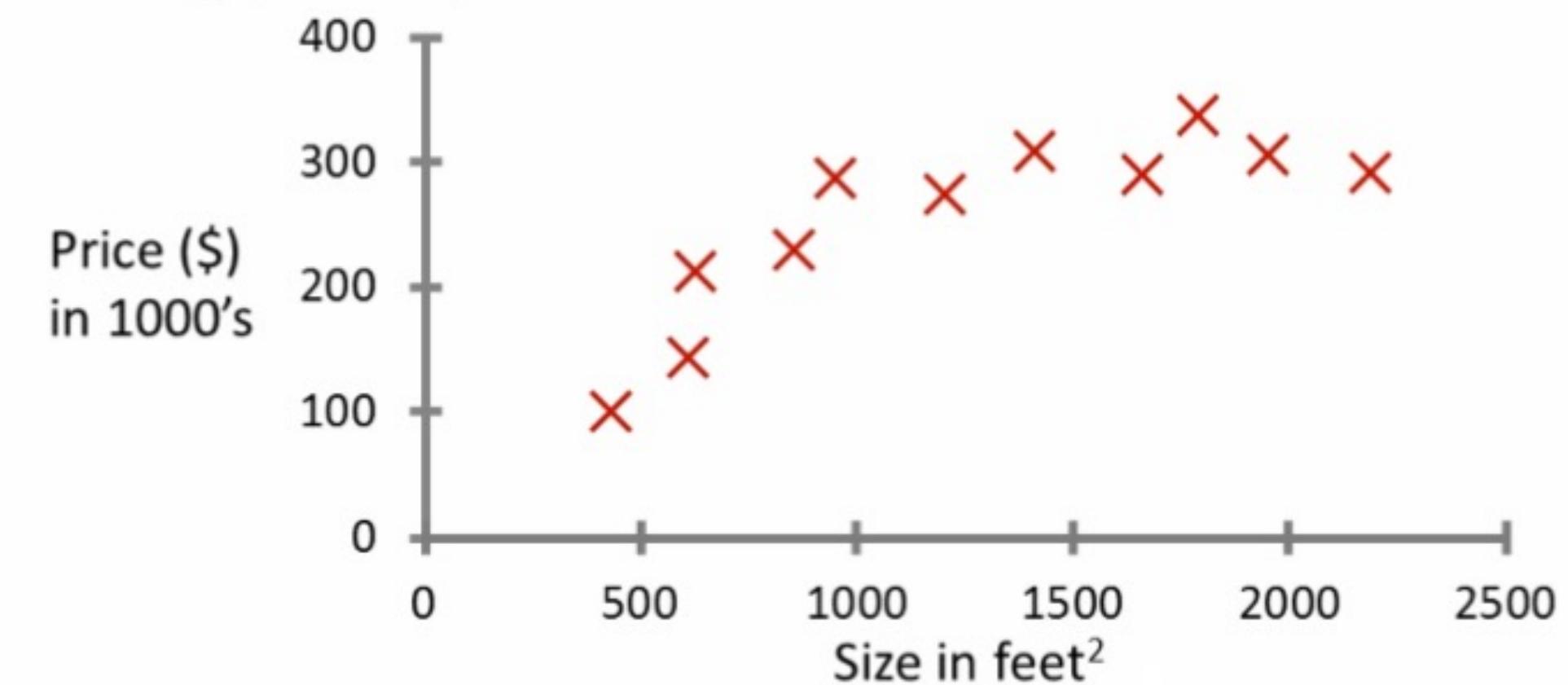
- Supervised Learning
- Unsupervised Learning
- Reinforcement Learning

Supervised Learning

In **supervised learning**, we are given a data set and **already know what our correct output should look like**, having the idea that there is a relationship between the input and the output.
Supervised learning problems are categorized into **regression** and **classification** problems:

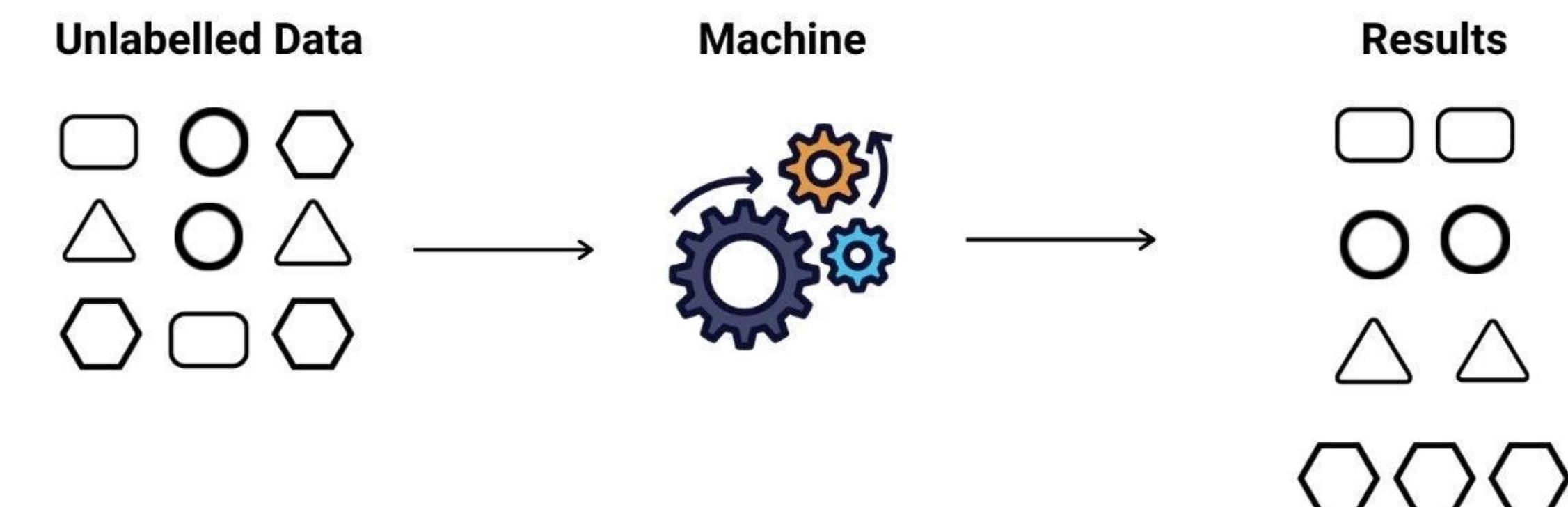
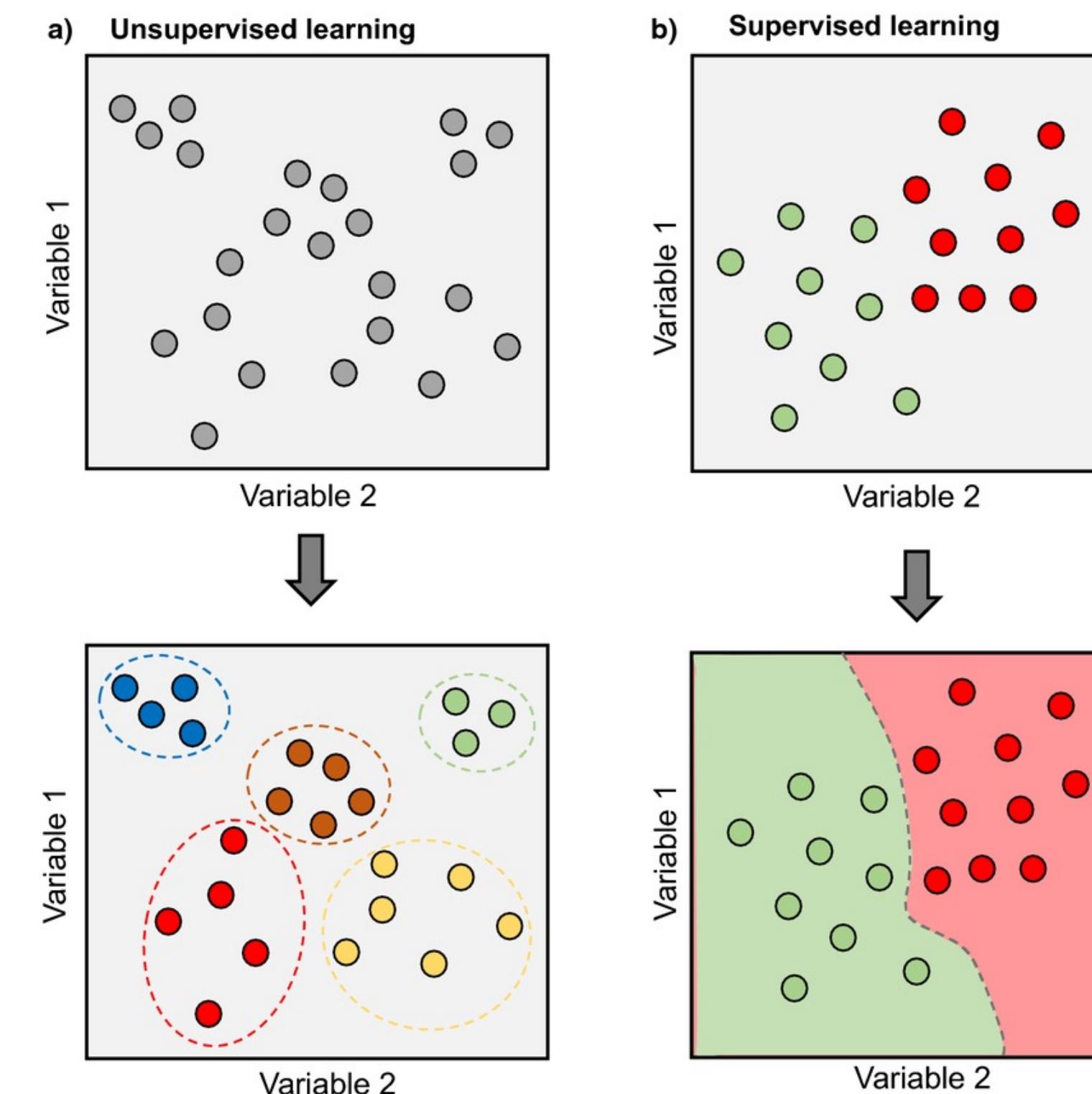


Housing price prediction.

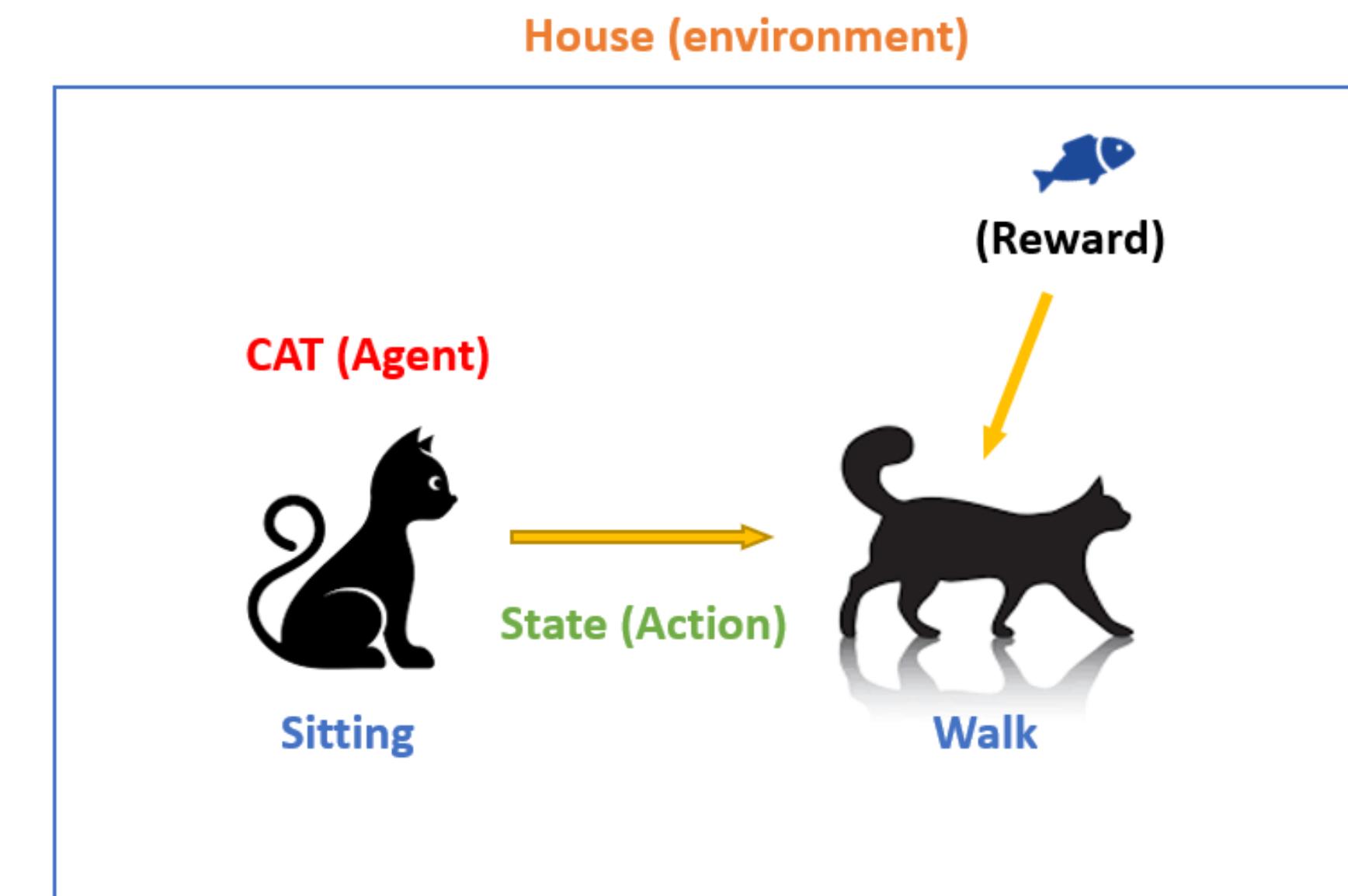
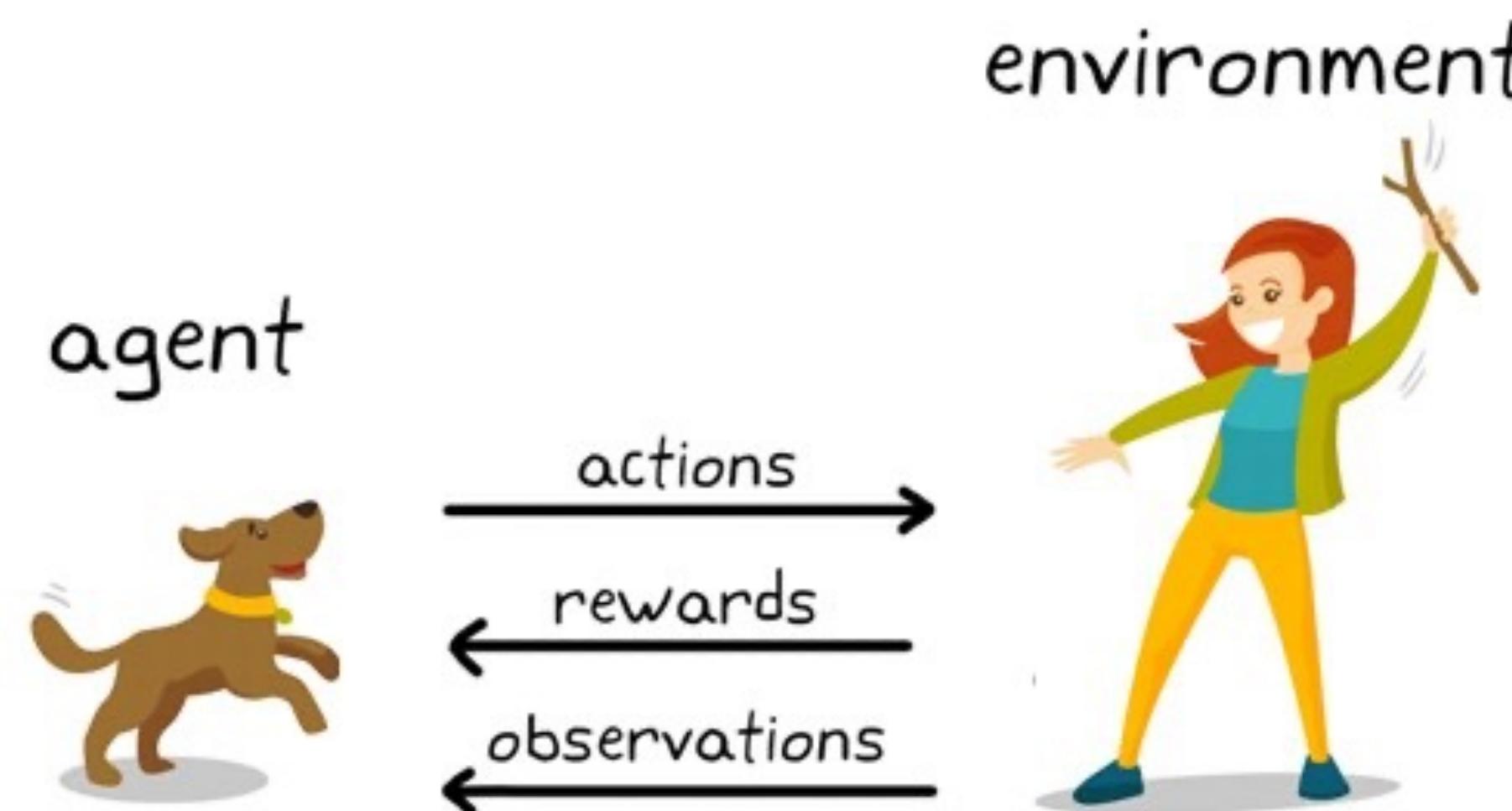


Classification - predict a class label
Regression - predict a number

Unsupervised learning, on the other hand, allows us to approach problems with little or no idea what our results should look like. We can derive structure from data where we don't necessarily know the effect of the variables. We can derive this structure by clustering the data based on relationships among the variables in the data.

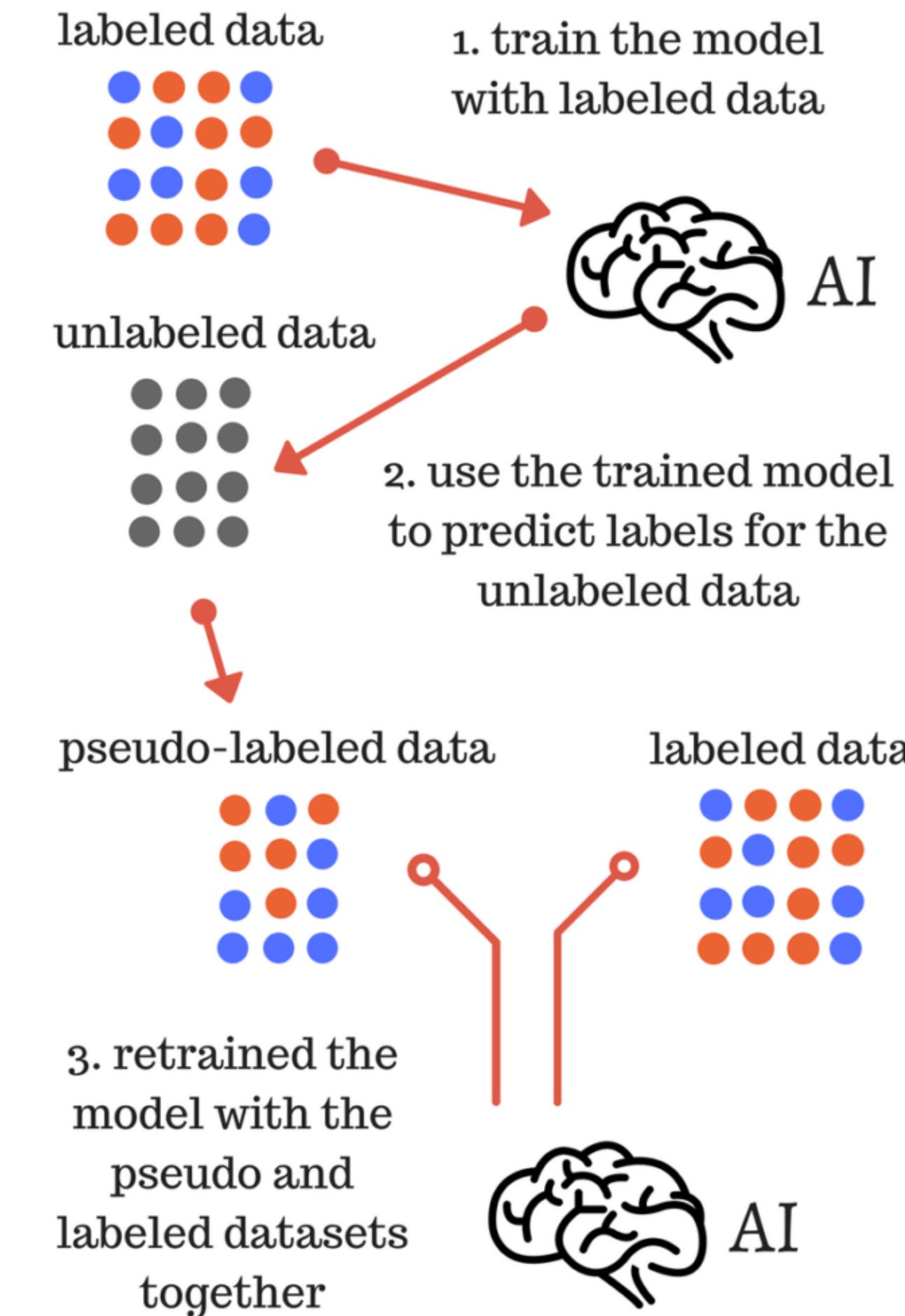


Reinforcement learning is learning what to do—how to map situations to actions—so as to maximize a numerical reward signal. The learner is not told which actions to take, but instead must discover which actions yield the most reward by trying them. In the most interesting and challenging cases, actions may affect not only the immediate reward but also the next situation and, through that, all subsequent rewards. These two characteristics—trial-and-error search and delayed reward—are the two most important distinguishing features of reinforcement learning.



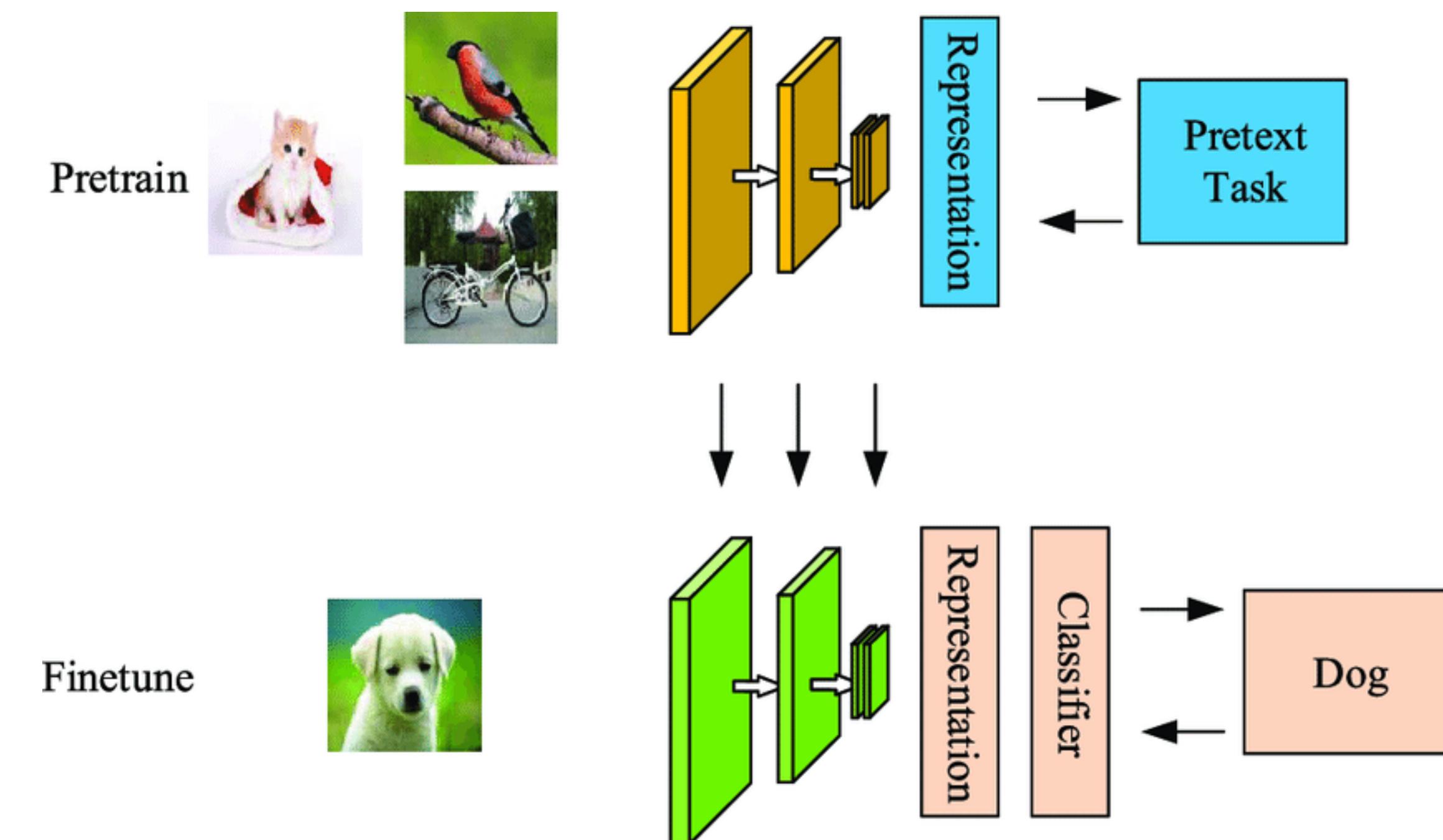
Semi-supervised learning

- ▶ Semi-supervised learning (you have some amount of data with labels and other data without labels)



Self-supervised learning

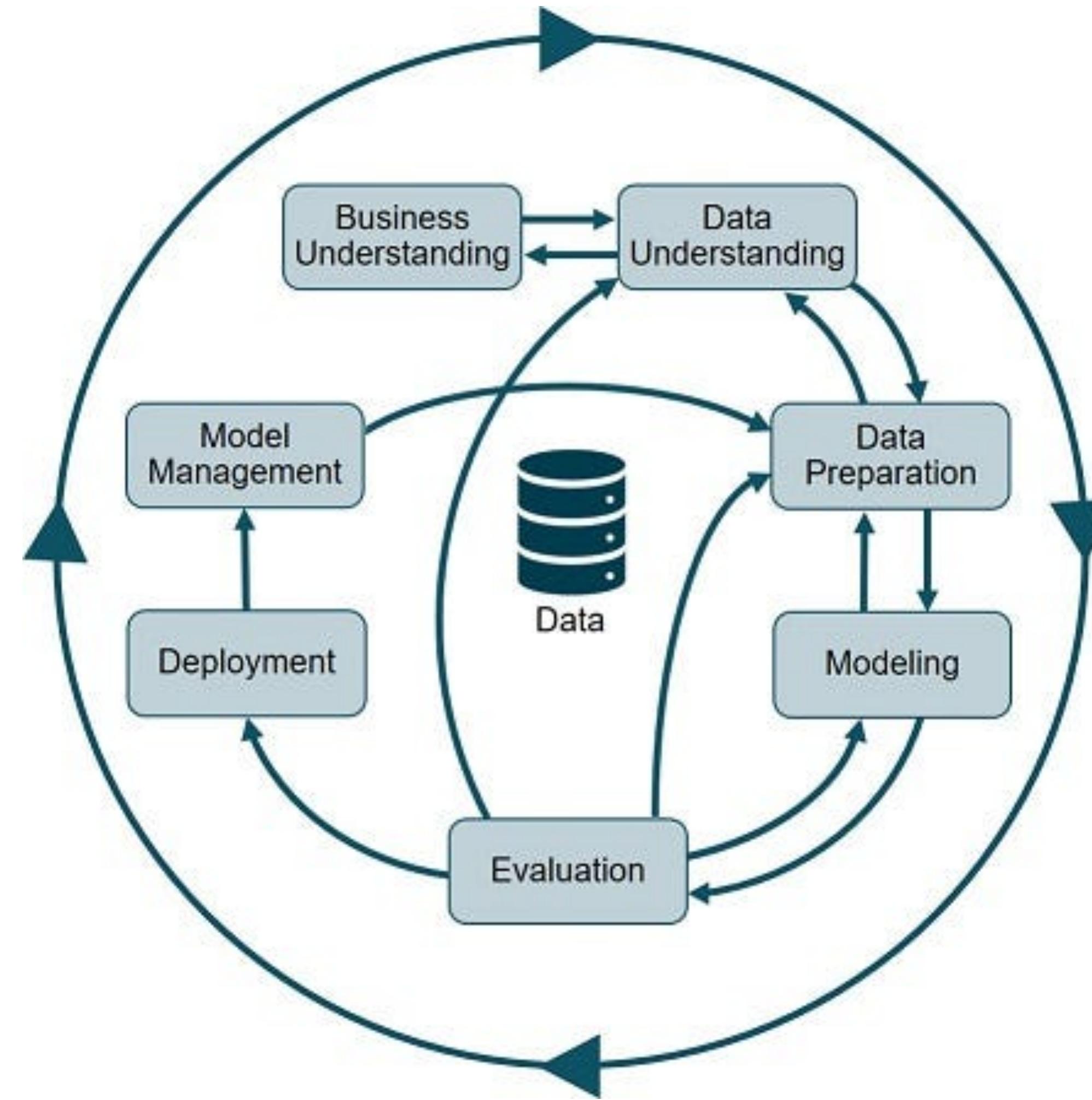
- ▶ Self-supervised learning
(unsupervised learning problem that is framed as a supervised learning problem)



ML Project Stages

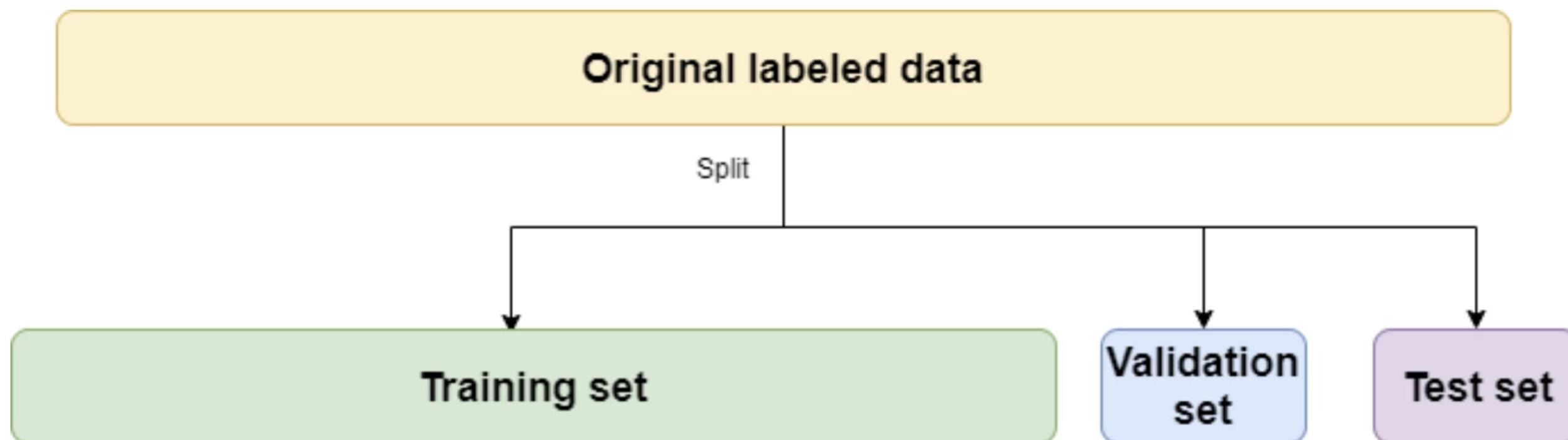
- Initiation – problem understanding, vision, goals
 - Data collection and exploratory data analysis
 - Data cleaning & pre-processing
 - Model development, training and evaluation
 - Model deployment and maintenance

CRISP-DM in ML

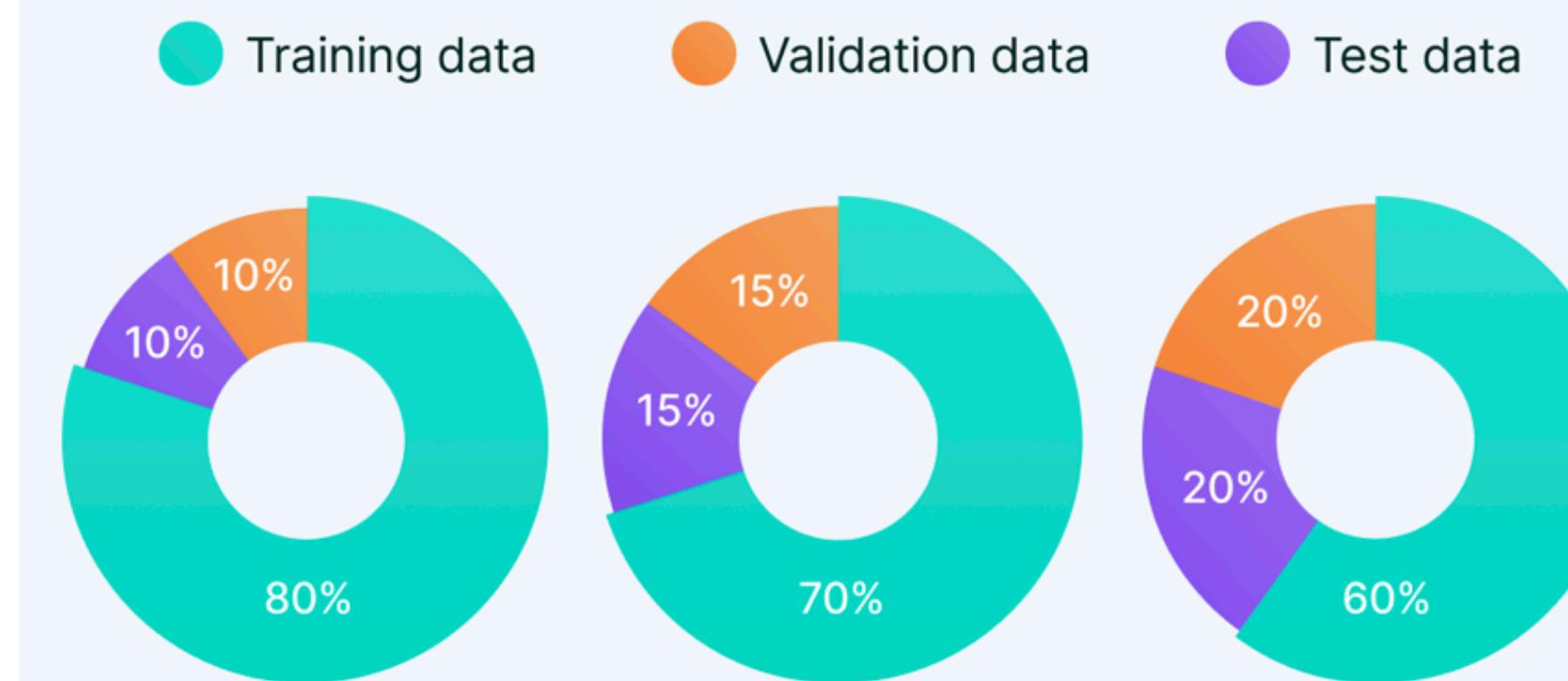


Cross Industry Standard Process for Data Mining (CRISP-DM)

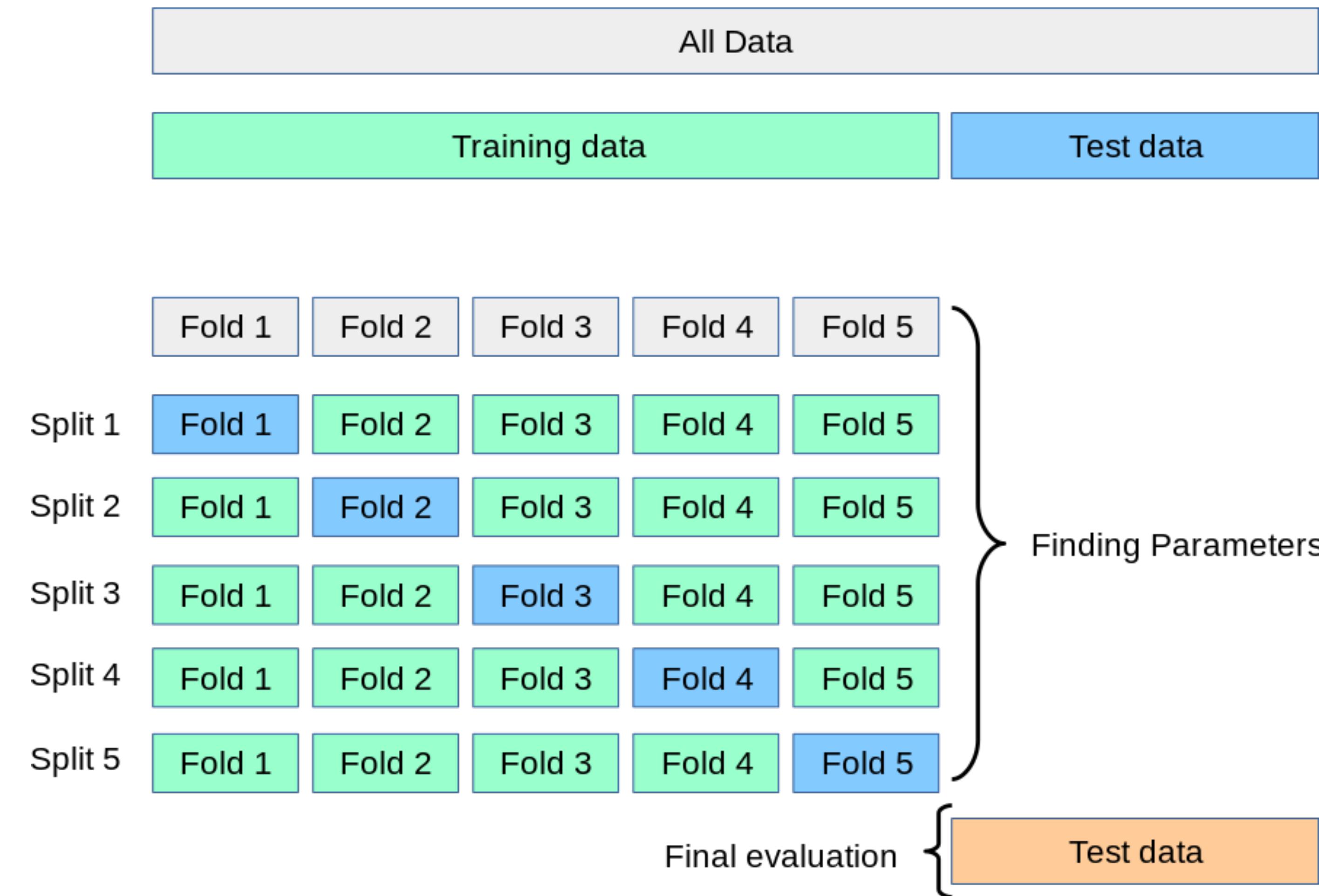
Data Split



Data Training Needs



Evaluation Strategies



Datasets

- Kaggle
- Papers With Code
- 😊 HF Datasets

Python ML Tools Basics Review. NumPy

Released as Numeric in 1995, **NumPy** (since 2005) is a library for the Python programming language, adding support for large, multi-dimensional arrays and matrices, along with a large collection of high-level mathematical functions to operate on these arrays.

Install Documentation Learn Community About Us News Contribute English ▾

NumPy

The fundamental package for scientific computing with Python

LATEST RELEASE: NUMPY 1.26. [VIEW ALL RELEASES](#)

NumPy 1.26.0 released 2023-09-16

POWERFUL N-DIMENSIONAL ARRAYS
Fast and versatile, the NumPy vectorization, indexing, and broadcasting concepts are the de-facto standards of array computing today.

NUMERICAL COMPUTING TOOLS
NumPy offers comprehensive mathematical functions, random number generators, linear algebra routines, Fourier transforms, and more.

OPEN SOURCE
Distributed under a liberal [BSD license](#), NumPy is developed and maintained [publicly on GitHub](#) by a vibrant, responsive, and diverse [community](#).

INTEROPERABLE
NumPy supports a wide range of hardware and computing platforms, and plays well with distributed, GPU, and sparse array libraries.

PERFORMANT
The core of NumPy is well-optimized C code. Enjoy the flexibility of Python with the speed of compiled code.

EASY TO USE
NumPy's high level syntax makes it accessible and productive for programmers from any background or experience level.

Try NumPy

Use the interactive shell to try NumPy in the browser

Python Basics Review. Pandas

pandas is a fast, powerful, flexible and easy to use open source data analysis and manipulation tool, built on top of the Python programming language.

 pandas

About us ▾ Getting started Documentation Community ▾ Contribute

pandas

pandas is a fast, powerful, flexible and easy to use open source data analysis and manipulation tool, built on top of the Python programming language.

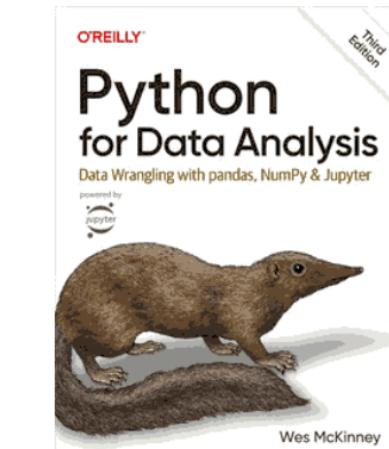
[Install pandas now!](#)

Latest version: 2.1.1

- [What's new in 2.1.1](#)
- Release date:
Sep 20, 2023
- [Documentation \(web\)](#)
- [Download source code](#)

Follow us

Get the book



Previous versions

- 2.1.0 (Aug 30, 2023)
[changelog](#) | [docs](#) | [code](#)
- 2.0.3 (Jun 28, 2023)
[changelog](#) | [docs](#) | [code](#)
- 2.0.2 (May 28, 2023)
[changelog](#) | [docs](#) | [code](#)
- 2.0.1 (Apr 24, 2023)
[changelog](#) | [docs](#) | [code](#)

Show more

Getting started

- [Install pandas](#)
- [Getting started](#)

Documentation

- [User guide](#)
- [API reference](#)
- [Contributing to pandas](#)
- [Release notes](#)

Community

- [About pandas](#)
- [Ask a question](#)
- [Ecosystem](#)

With the support of:



OPEN CODE = BETTER SCIENCE



TWO SIGMA



VOLTRON DATA



Coiled
A Dask Company



Quansight Labs



NVIDIA.



TIDELIFT

Chan
Zuckerberg
Initiative

bodo.ai

Python Basics Review. Matplotlib

Initially released in 2003, **Matplotlib** is a plotting library for the Python programming language and NumPy.

[matplotlib](https://matplotlib.org) Plot types User guide Tutorials Examples Reference Contribute Releases

Search icon | Help icon | GitHub icon | Q&A icon | Twitter icon

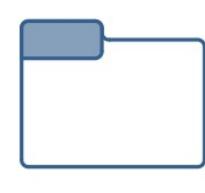
Matplotlib: Visualization with Python

Matplotlib is a comprehensive library for creating static, animated, and interactive visualizations in Python. Matplotlib makes easy things easy and hard things possible.

- Create publication quality plots.
- Make interactive figures that can zoom, pan, update.
- Customize visual style and layout.
- Export to many file formats.
- Embed in JupyterLab and Graphical User Interfaces.
- Use a rich array of third-party packages built on Matplotlib.

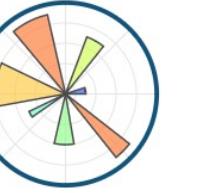
[Try Matplotlib \(on Binder\) →](#)

 Getting Started

 Examples

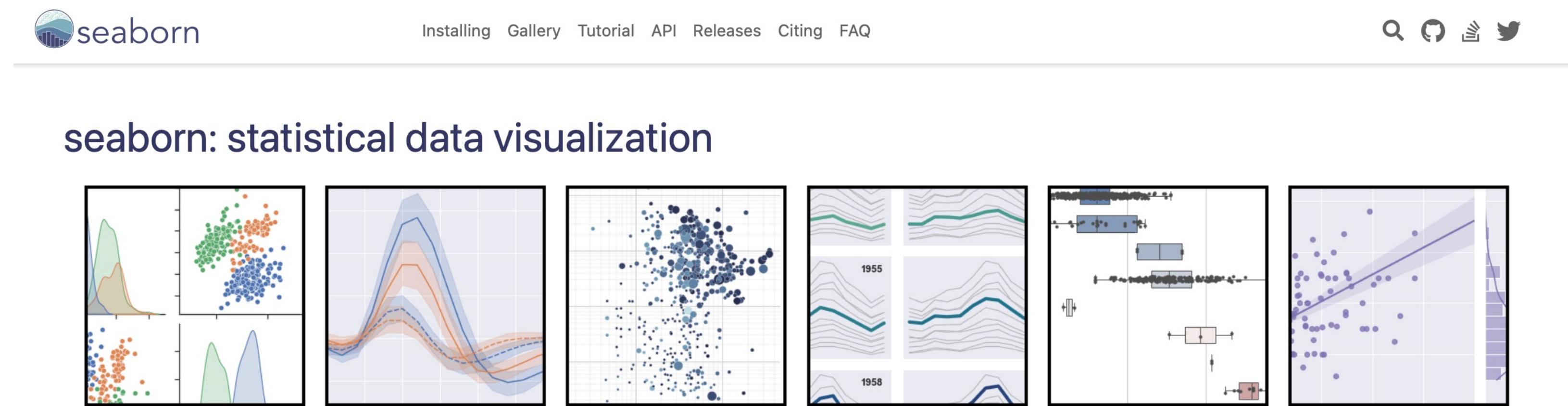
 Reference

 Cheat Sheets

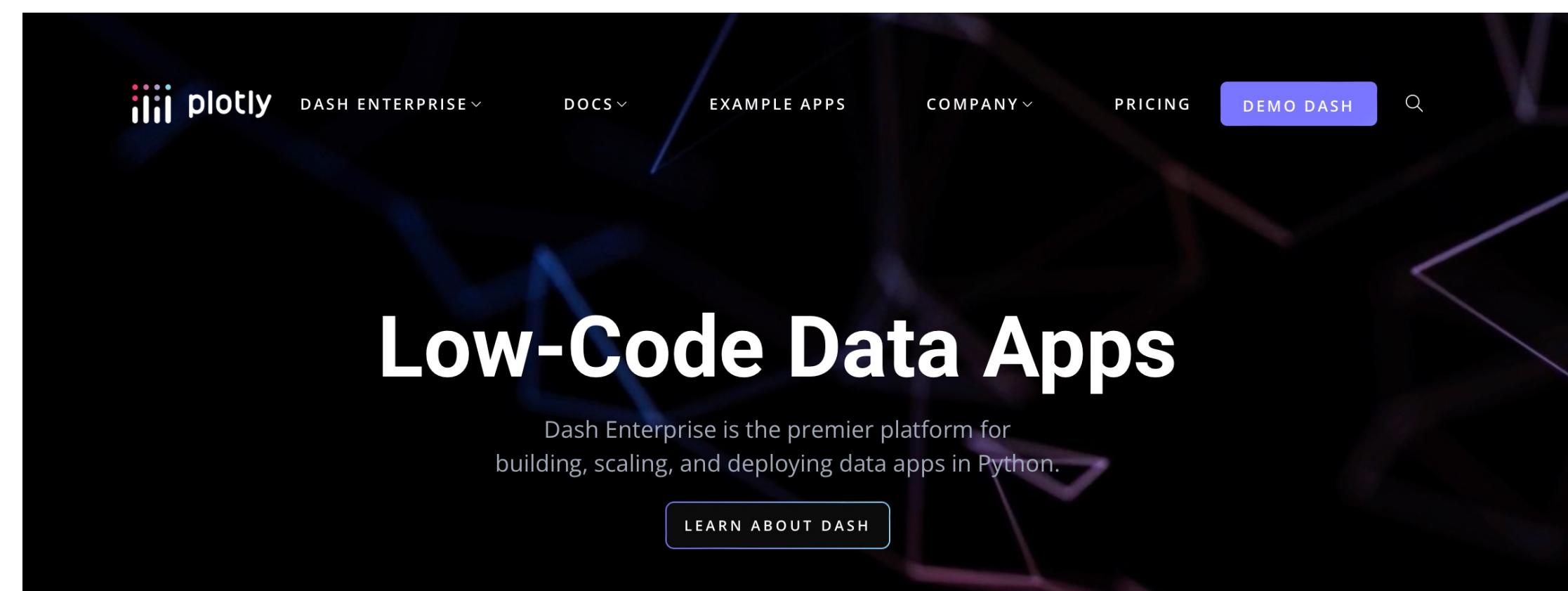
 Documentation

Seaborn & Plotly

Seaborn is a Python data visualization library based on matplotlib. It provides a high level interface for drawing attractive and informative statistical graphics.



Plotly makes interactive, publication-quality graphs: line plots, scatter plots, area charts, bar charts, error bars, box plots, histograms, heatmaps, subplots, multiple-axes, polar charts, bubble charts, etc.





Thanks for your attention