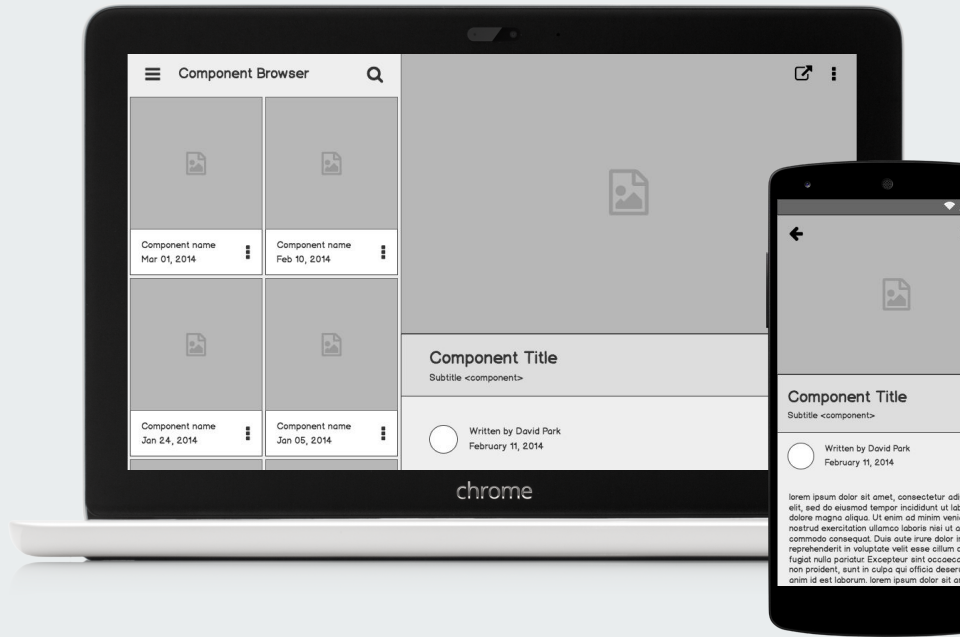
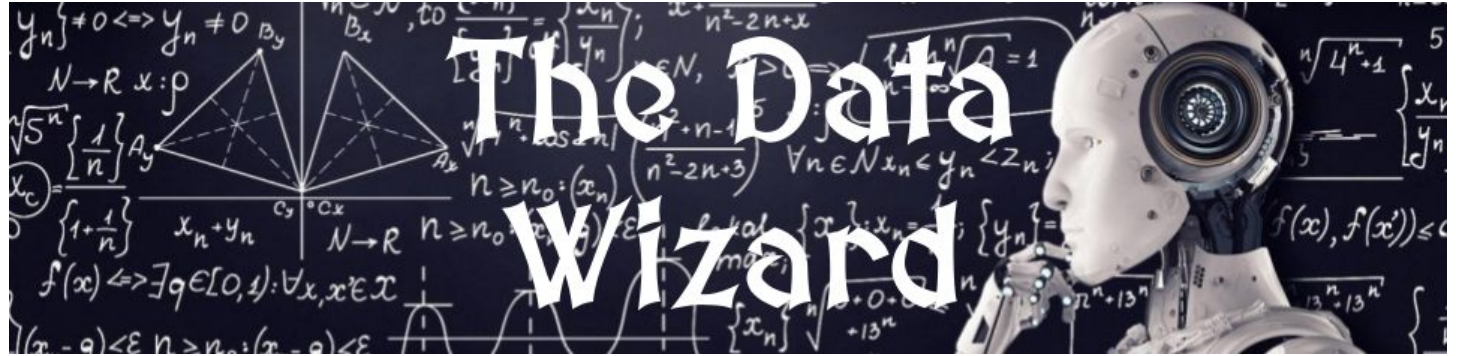


# Enterprise Software Engineering Practices: AI/ML Real-World Applications

Lecture by: Brooks Christensen



# ABOUT ME



- Master's Degree in Physics from CU Boulder (2020)
- Applied Data Science Certificate from MIT PE Program (2022)
- AI / ML for Business Applications Certificate from UT Austin (2023)
- Nielsen Media Ratings (2021 - 2024)
- Freelance Data Scientist (2024 - Present)

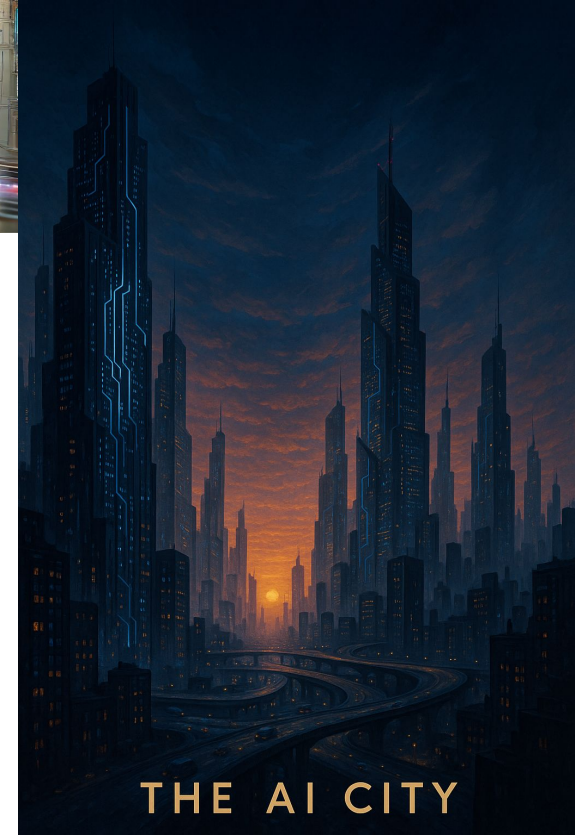


# Outline

- Intro to AI/ML Applications
- AI/ML Types
- AI Models
- Ensembles
- How AI/ML Systems Fail
- Ethical Considerations



“The AI City” - ChatGPT-4o (Oct 1, 2024)



“The AI City” - ChatGPT-5.1 (Nov 30, 2025)

THE AI CITY

---

# Introduction to AI/ML Applications

## U-Net: Convolutional Networks for Biomedical Image Segmentation

Olaf Ronneberger, Philipp Fischer, Thomas Brox

There is large consent that successful training of deep networks requires many thousand annotated training samples. In this paper, we propose an expanding path that enables precise localization. We show that such a network can be trained end-to-end from very few images and output (phase contrast and DIC) we won the ISBI cell tracking challenge 2015 in these categories by a large margin. Moreover, the network is fa

## Extreme-Long-short Term Memory for Time-series Prediction

Sida Xing, Feihu Han, Suiyang Khoo

The emergence of Long Short-Term Memory (LSTM) solves the problems of vanishing gradient and exploding gradient. In traditional Recurrent LSTM, we proposed an advanced LSTM algorithm, the Extreme Long Short-Term Memory (E-LSTM), which adds the inverse matrix part to training rounds, thus reducing the overall training time. In this research, the E-LSTM model is used for the text prediction task. Experimental results showed that the E-LSTM sometimes takes longer to traditional LSTM, whilst also improving the training speed and the overall efficiency of the LSTM.

## Attention Is All You Need

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion

The dominant sequence transduction models are based on complex recurrent and convolutions entirely. Experiments on two machine translation tasks, the WMT 2014 English-to-French translation task, our model establishes a new single-tr

## Generative Adversarial Networks

Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, S

We propose a new framework for estimating generative models via an adversarial process, in which making a mistake. This framework corresponds to a minimax two-player game. In the space of all Markov chains or unrolled approximate inference networks during either training or generation of

## Deep Residual Learning for Image Recognition

Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun

Deeper neural networks are more difficult to train. We present a residual learning framework to ease the training of networks that are substantially deeper showing that these residual networks are easier to optimize, and can gain accuracy from considerably increased depth. On the ImageNet dataset we place on the ILSVRC 2015 classification task. We also present analysis on CIFAR-10 with 100 and 1000 layers. The depth of representations is of central importance for many visual recognition tasks. Solely due to our extremely deep representations, we obtain a 2

## Auto-Encoding Variational Bayes

Diederik P Kingma, Max Welling

How can we perform efficient inference and learning in directed probabilistic models? Our contributions are two-fold. First, we show that a reparameterization of the an approximate inference model (also called a recognition model) to the intractable

## Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms

Han Xiao, Kashif Rasul, Roland Vollgraf

We present Fashion-MNIST, a new dataset comprising of 28x28 grayscale images of 70,000 fashion products from 10 classes. The dataset is designed to be a drop-in replacement for the MNIST dataset, as it shares the same image size, data format and the structure of training and testing splits. The dataset is

## Recurrent Neural Networks (RNNs): A Tutorial

Robin M. Schmidt

State-of-the-art solutions in the areas of "Language Modelling & Generation", "Text-to-Speech", "Machine Translation", "Image Captioning", "Image-to-Text", "Text-to-Image", "Image-to-Image", "Mechanism" or "Pointer Networks". We also give recommendations for further reading in these areas. In this work we give a short overview of "Recurrent Neural Networks" or "RNNs". We also give recommendations for further reading in these areas.

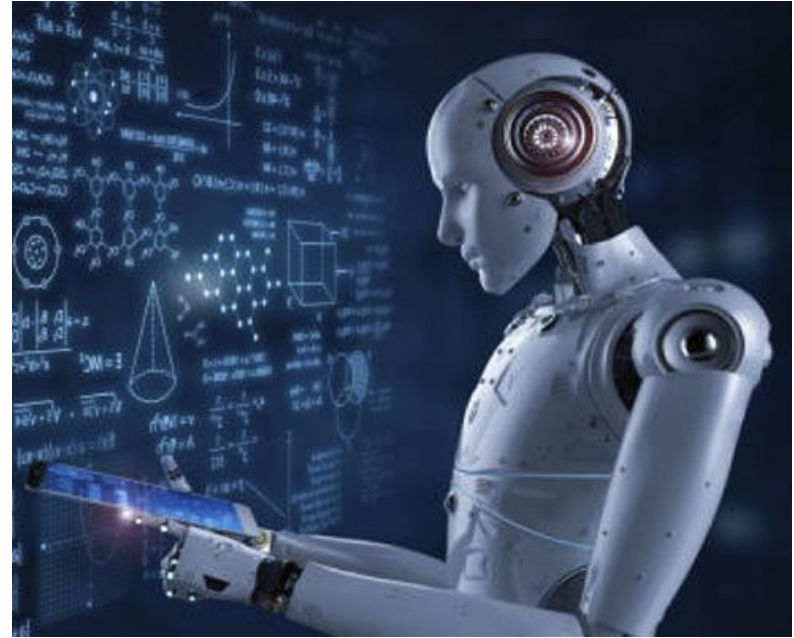
## WaveNet: A Generative Model for Raw Audio

Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex

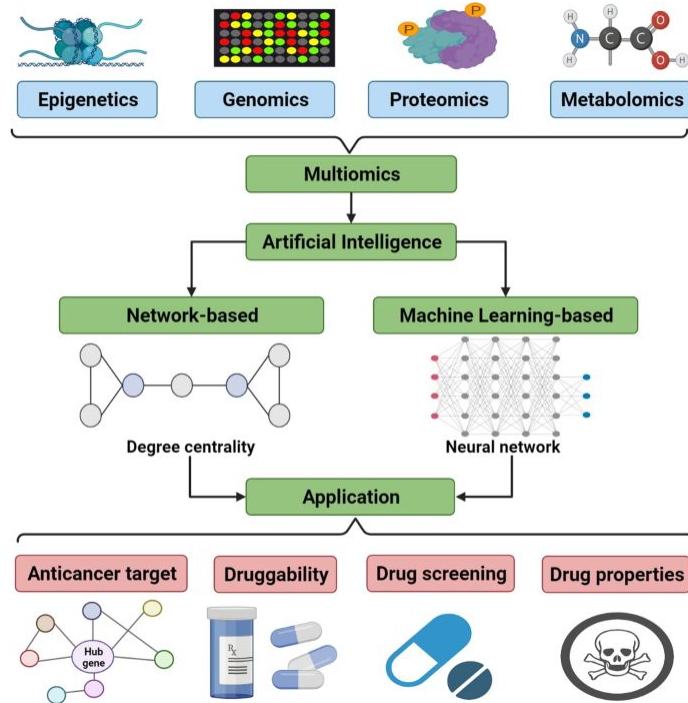
This paper introduces WaveNet, a deep neural network for generating raw audio waveforms. The network is trained on speech, it yields state-of-the-art performance, with human listeners rating it as significantly more natural than model music, we find that it generates novel and often highly realistic musical fragments.

# Examples of AI/ML Applications

- Large Language Models (ChatGPT, etc)
  - Machine Translation
  - Chatbots
  - Text Summarization
  - Sentiment Analysis
  - Content Creation
  - Knowledge Extraction and Question Answering
  - Personalization
  - Research & Development
  - **Retrieval Augmented Generation (RAG)**
- Traffic Management
- Food Production
- **Environmental Modeling and Sustainability**
- Creativity and Art
- Quantum Machine Learning



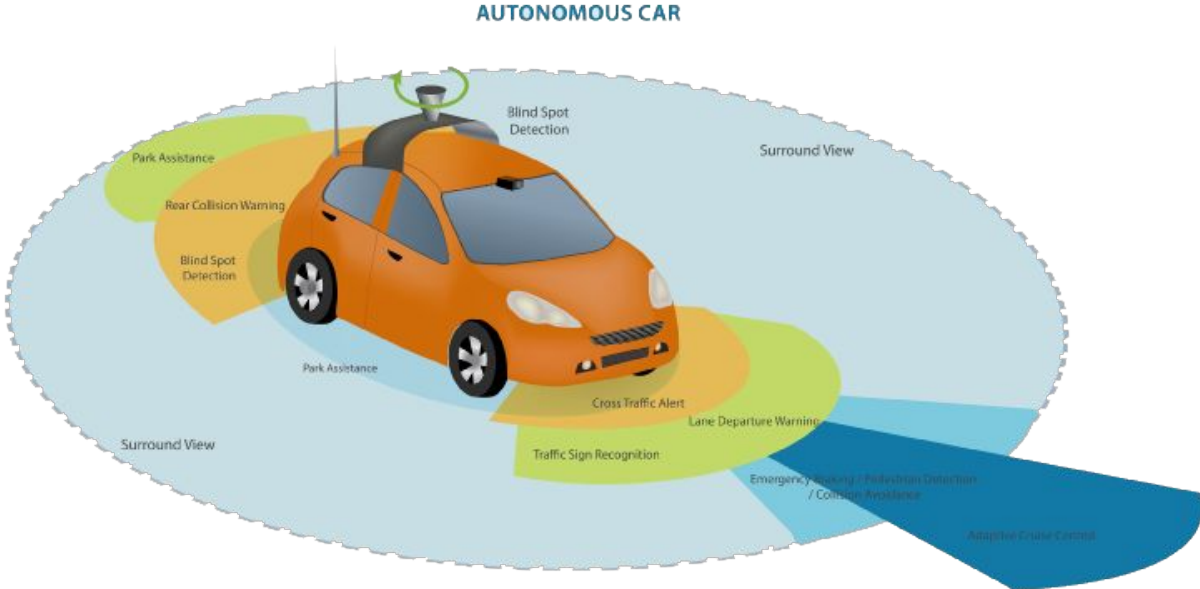
# AI/ML Applications in Health Sciences

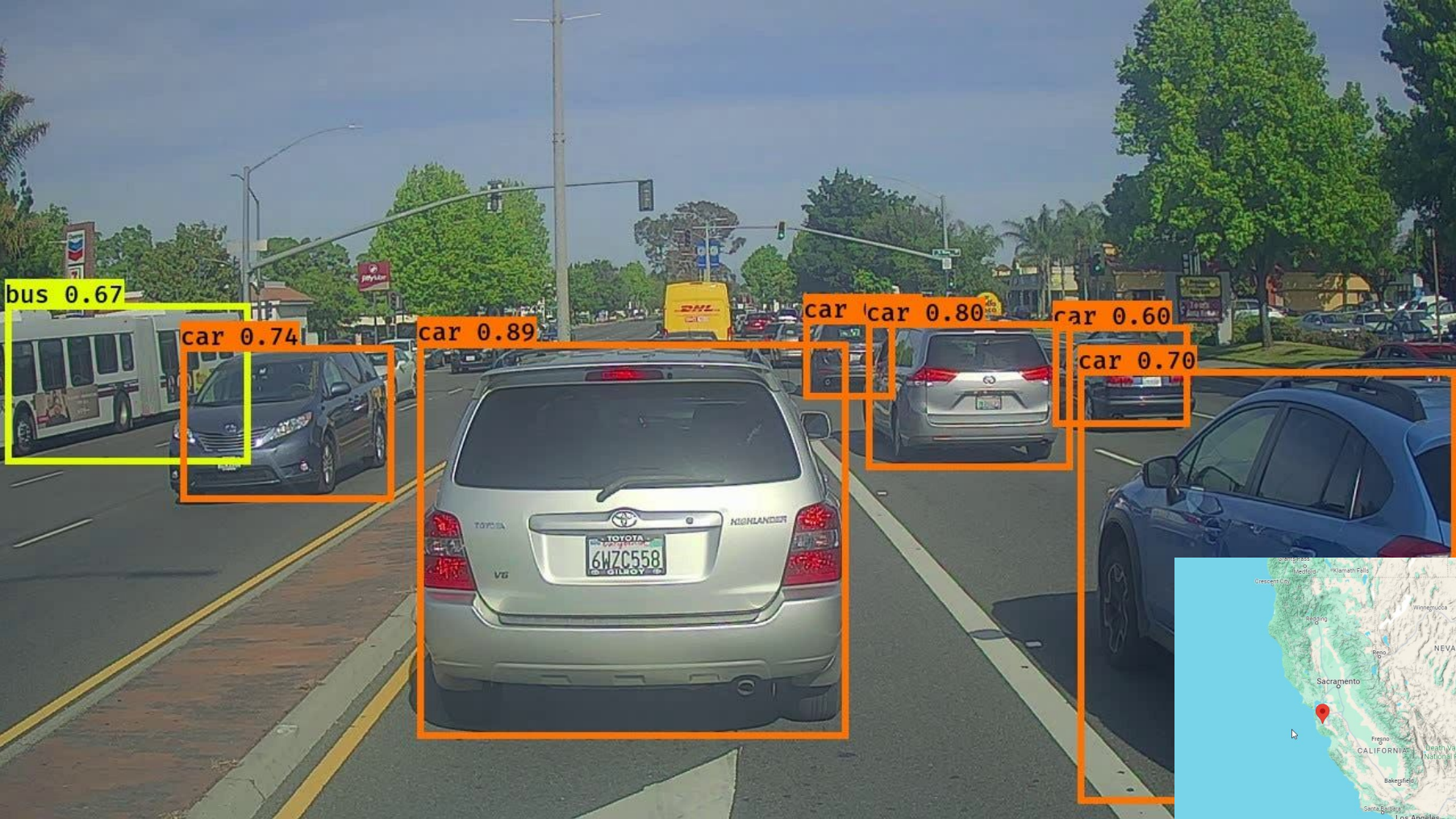


# AI/ML Applications in Business & Finance



# AI/ML Applications in Automation





bus 0.67

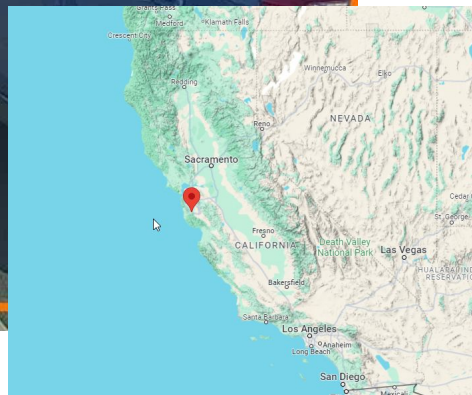
car 0.74

car 0.89

car car 0.80

car 0.60

car 0.70



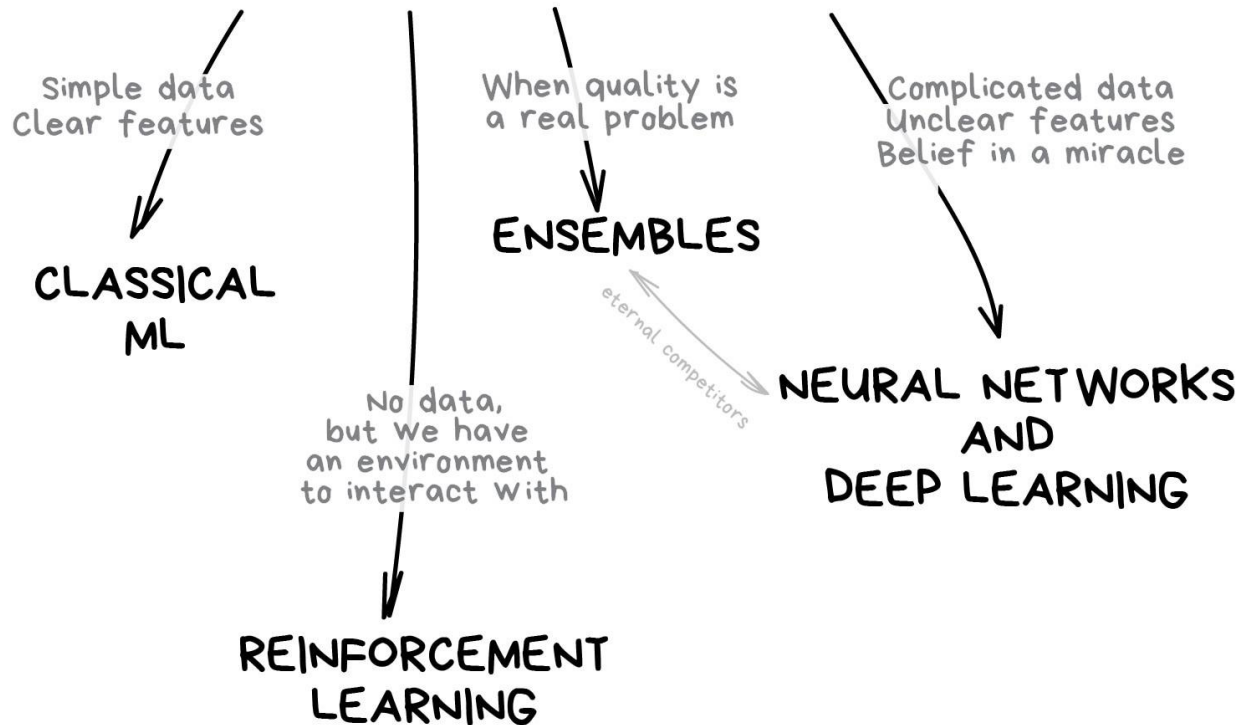




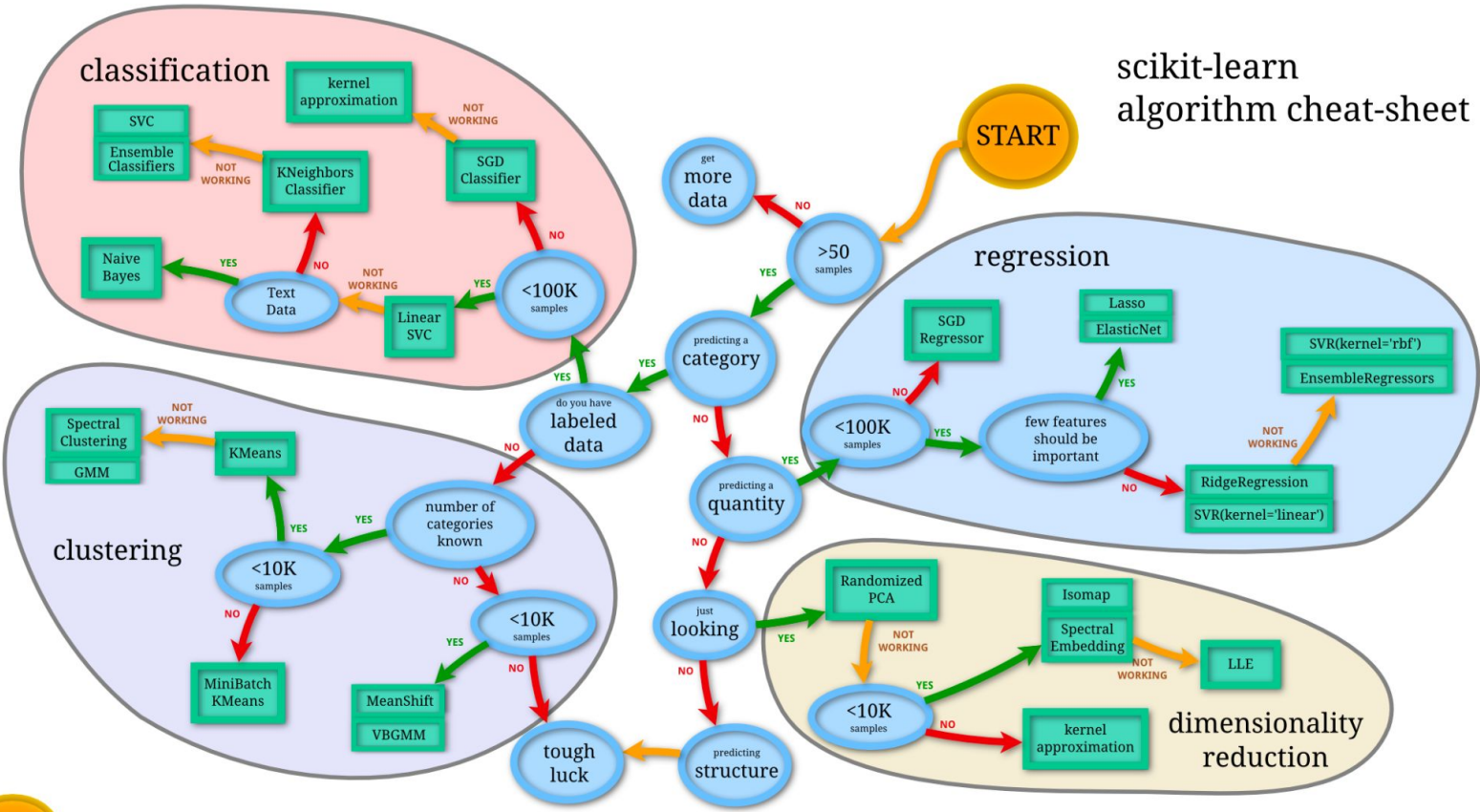
# AI/ML Types

# Main Types of AI/ML

## THE MAIN TYPES OF MACHINE LEARNING

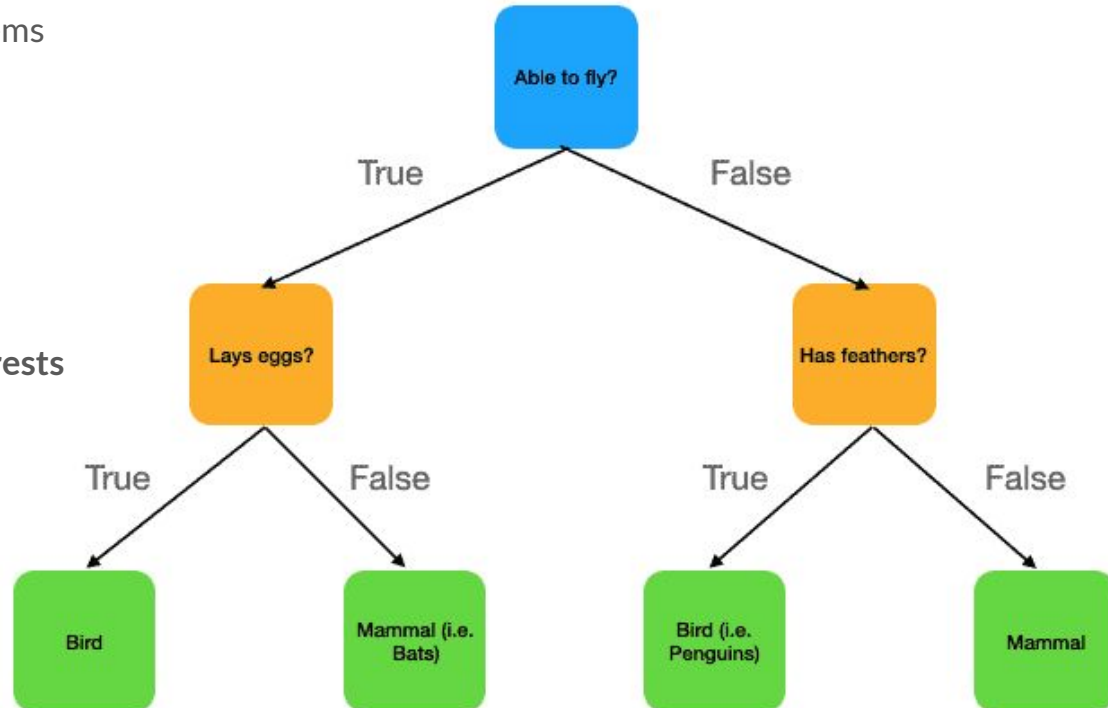


# scikit-learn algorithm cheat-sheet



# Decision Trees

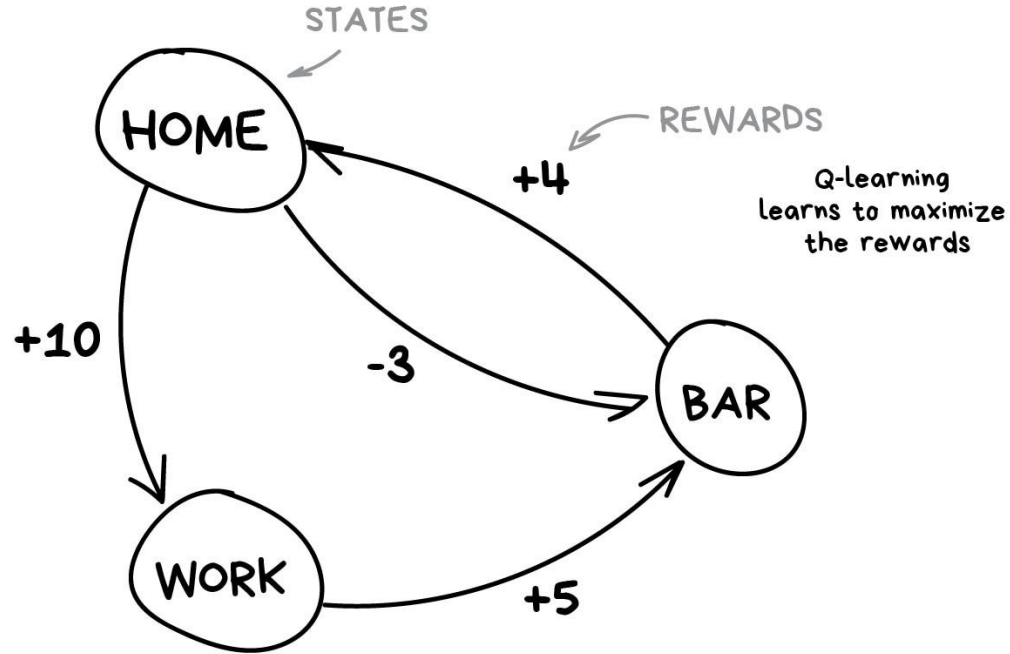
- Classical ML Technique
- Practical Solutions for **Many** Problems
- Useful for Tabular Data
- Not a Neural Network
- Classification or Regression
- Typically Better Performance in **Forests**



# Reinforcement Learning



Reinforcement Learning



ROUTINE MARKOV PROCESS

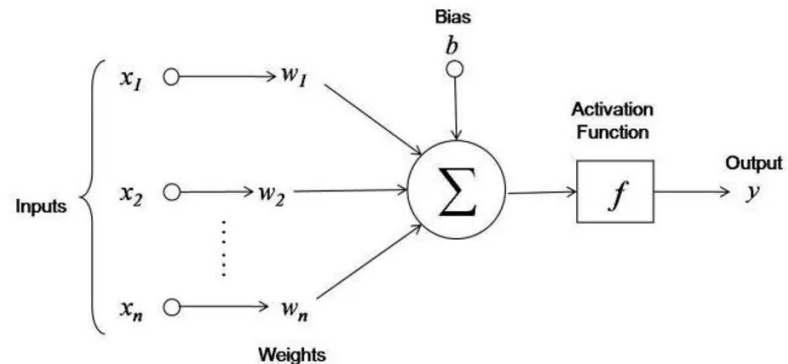
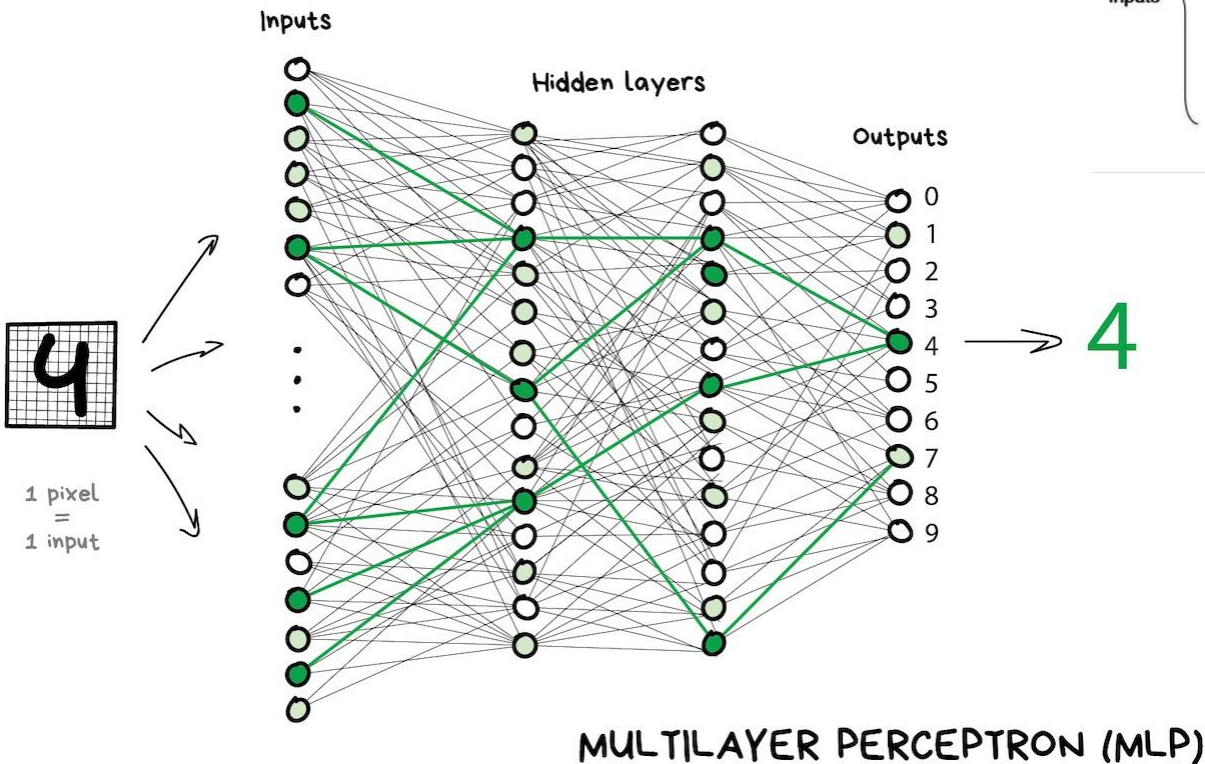
MICROSOFT / WEB / TL;DR

# Twitter taught Microsoft's AI chatbot to be a racist asshole in less than a day



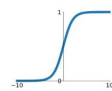
# AI Models

# AI Models - MLP

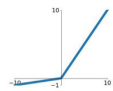


## Activation Functions

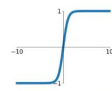
**Sigmoid**  
 $\sigma(x) = \frac{1}{1+e^{-x}}$



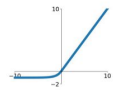
**Leaky ReLU**  
 $\max(0.1x, x)$



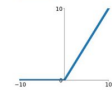
**tanh**  
 $\tanh(x)$



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



**ReLU**  
 $\max(0, x)$



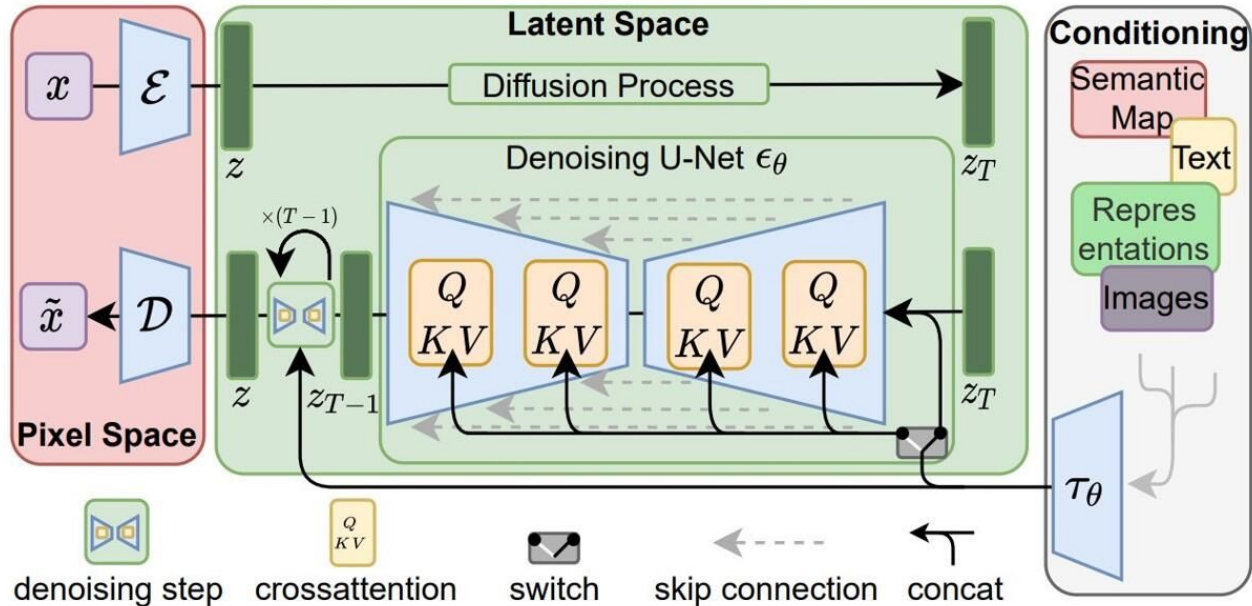
**ELU**  
 $\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$

## Cross-Entropy Loss Function

$$H(P^* | P) = - \sum_i \underbrace{P^*(i)}_{\text{TRUE CLASS DISTRIBUTION}} \log \underbrace{P(i)}_{\text{PREDICTED CLASS DISTRIBUTION}}$$

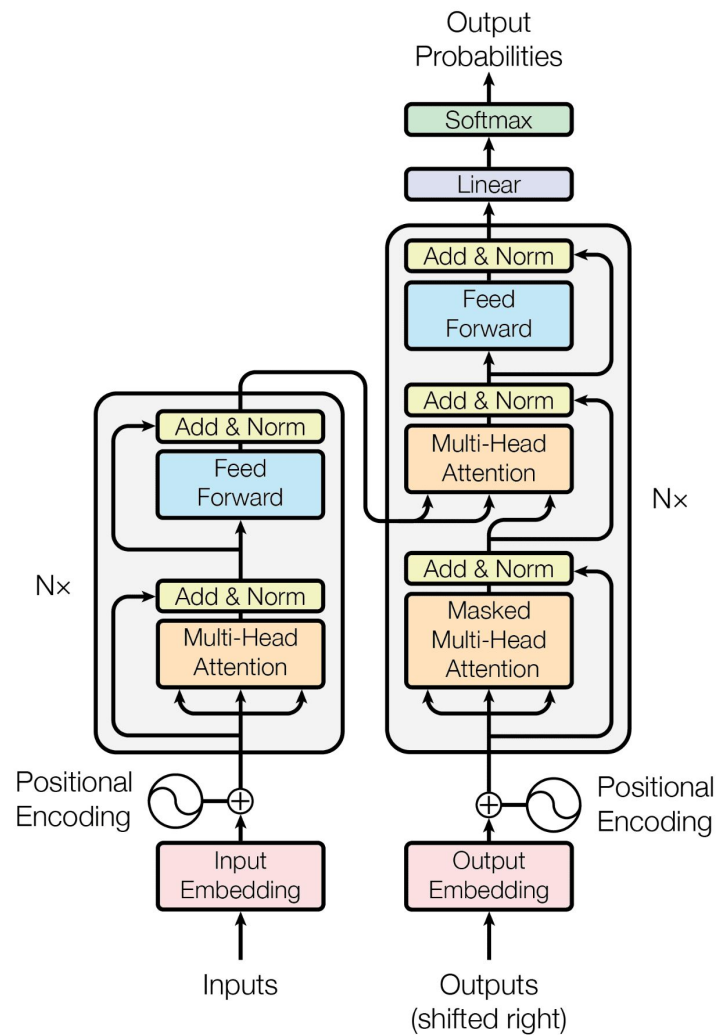
# AI Models - Stable Diffusion

- Previous Image Generation Systems used SD-type Algorithms



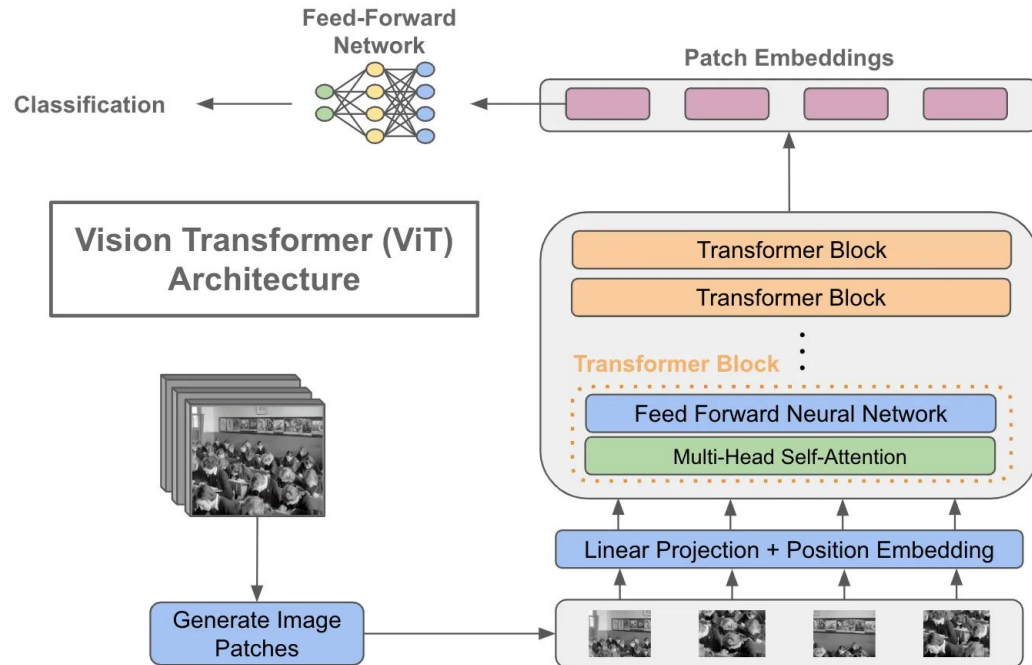
# AI Models - Transformer

- “Attention is All You Need”, Vaswani, Et. al (2017)
- Original GPT Architecture
- Key, Query, Value
- Encoder/Decoder
- Positional Encoding
- Masking



# AI Models - Vision Transformer (ViT)

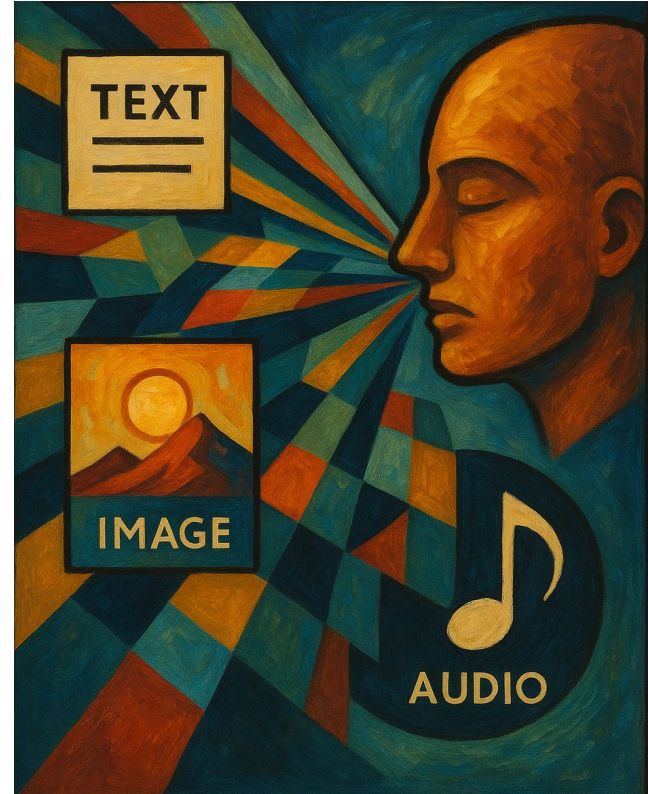
- Creates patches out of image
- Linear projection & position embedding
- Multiple transformer blocks:
  - Multi-head self-attention
  - FFNN (aka MLP)
- Patch embeddings act like tokens
- Final classification with FFNN (aka MLP)



# AI Models - Multimodal Transformer-Based Generators

- Transformer-based Architectures that **Additionally Support**:
  - Vision
  - Language
  - Audio
  - Action
- Current GPT Architecture

“Multimodal Possibility” - ChatGPT-5.1 (Nov 30, 2025)



# AI Models - Recommendation Systems

- Sparse matrix calculations
- Predicts which users will be interested in which products or services
- Individual user history is important to predicting future interest
- Examples:
  - Netflix
  - Amazon



# AI Models - Retrieval Augmented Generation (RAG)

- Allows users to make queries about a **specific set** of existing data
- **Vectorization** allows for fast comparison between query and relevant text in corpus
  - FAISS is widely used
- Examples:
  - Search your emails
  - Search medical records
  - Make queries about legal documents



NEW  
MODELS  
EVERY DAY

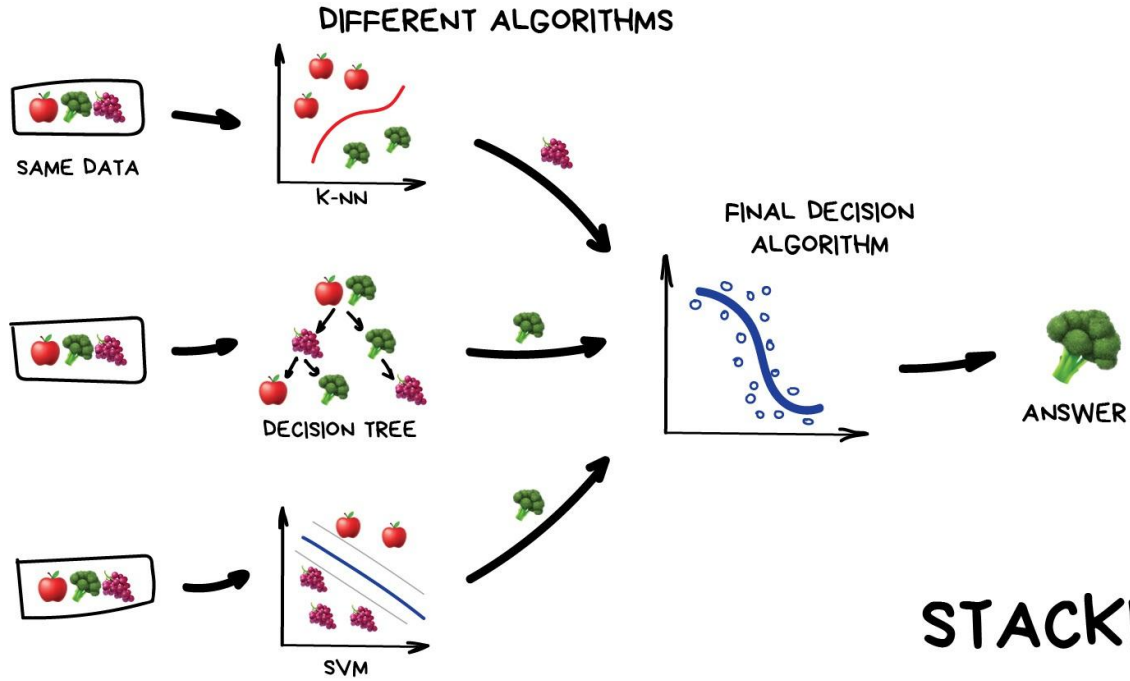


# Ensembles

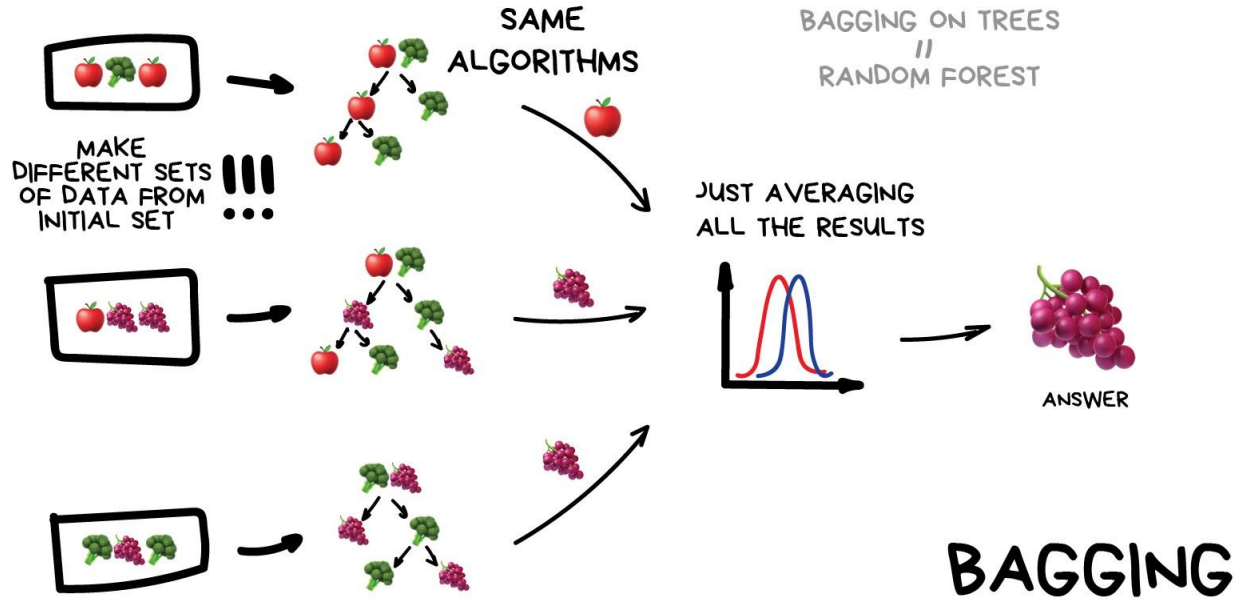
# Ensembles



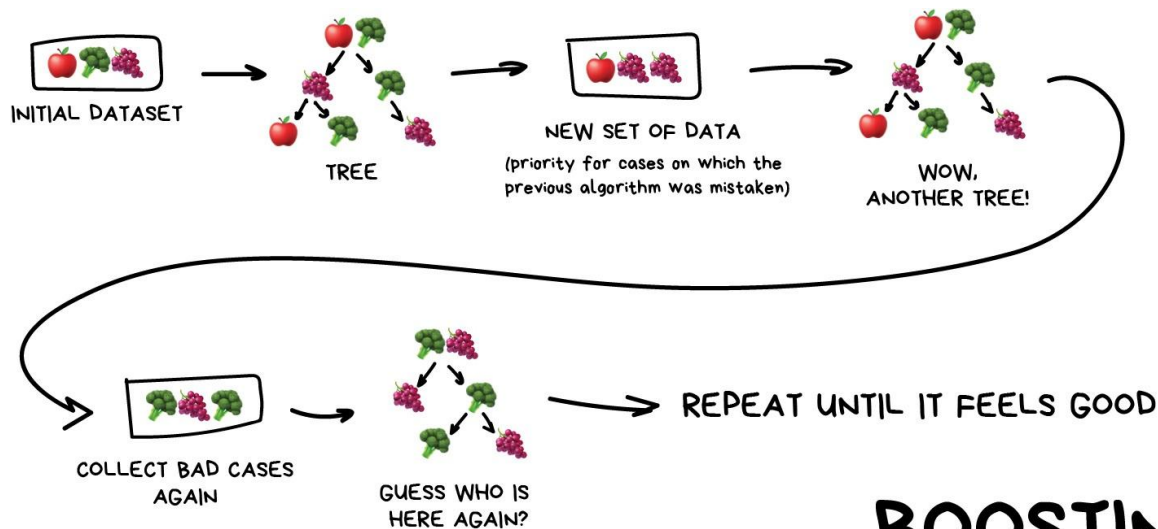
# Ensembles - Stacking



# Ensembles - Bagging



# Ensembles - Boosting



# BOOSTING

---

# How AI/ML Systems Fail

ARTIFICIAL INTELLIGENCE

# Hundreds of AI tools have been built to catch covid. None of them helped.

## What went wrong

Many of the problems that were uncovered are linked to the poor quality of the data that researchers used to develop their tools.

# Traditional Software System Failures

- **Dependency failure** – A third party package no longer maintained
- **Hardware failure** – Hard-drive failure, CPU overheating, losing network access, etc
- **Downtime or Crashing** – Backend infrastructure outage

TECH

## **Dead Roombas, stranded packages and delayed exams: How the AWS outage wreaked havoc across the U.S.**

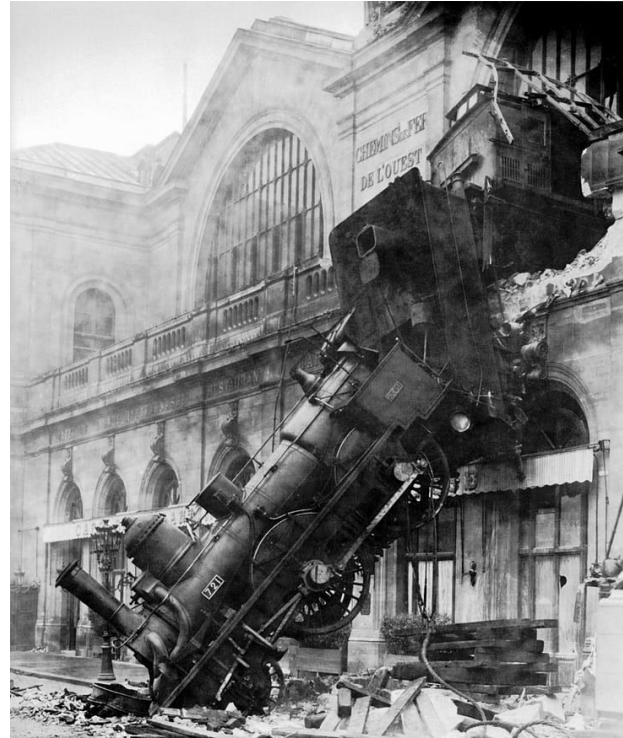
PUBLISHED THU, DEC 9 2021-8:00 AM EST | UPDATED THU, DEC 9 2021-10:51 AM EST



# ML-Specific Failures

They account for a smaller proportion of failures but are **harder to detect and fix**:

- **Edge cases**
- **Data Distribution Shifts**
- **Degenerate feedback loops**



# ML-Specific Failures – Edge Cases



# ML-Specific Failures – Data Distribution Shifts

Consider a model estimating the likelihood of pets being adopted in a shelter. Let's call  $X$  the characteristics of a pet,  $Y$  whether its adopted or not, and  $P(Y|X)$  the likelihood of a pet being adopted

- **Covariate shift** –  $P(X)$  changes,  $P(Y|X)$  remains

A shelter takes in more older pets, but the likelihood of adoption remains the same by age.

- **Label shift** –  $P(Y)$  changes,  $P(X|Y)$  remains

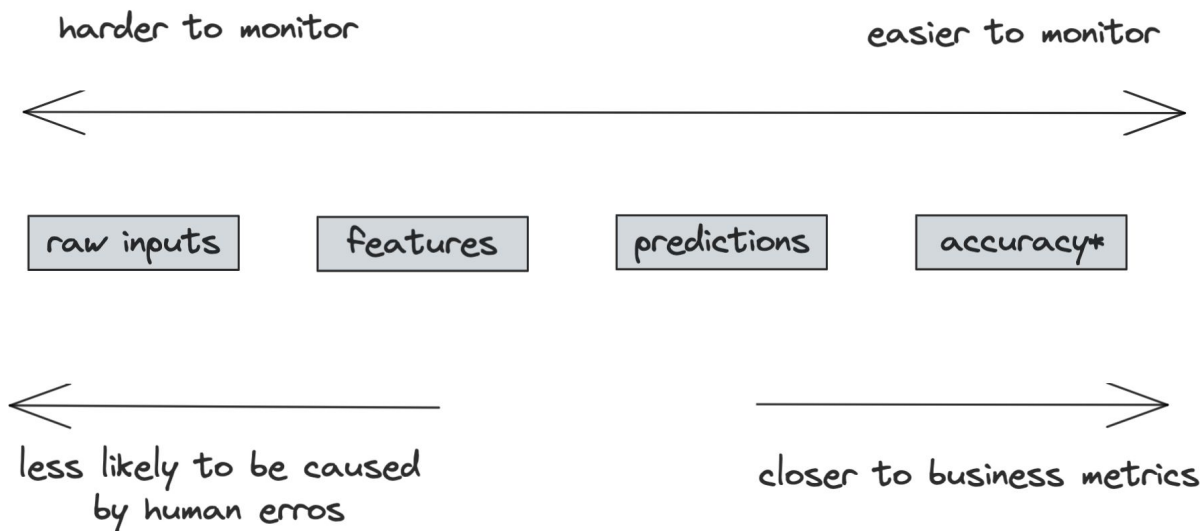
People want to adopt more pets, regardless of age.

- **Concept drift** –  $P(Y|X)$  changes,  $P(X)$  remains

People want to adopt more older pets, but the age of pets in the shelter remains the same.

# ML-Specific Failures – Monitoring Data Distributions

Monitor data all along the inference pipeline over time

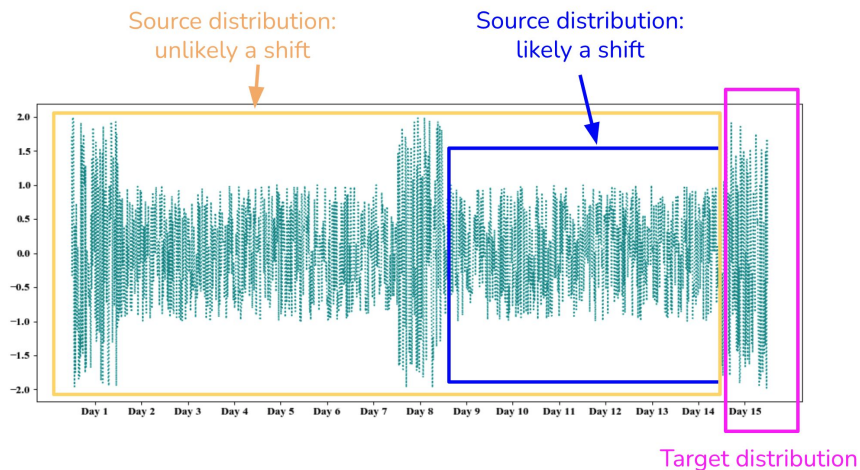


\* if natural labels available

# ML-Specific Failures – Detecting Data Distribution Shifts

A couple of monitoring pitfalls:

## Consider seasonal variations

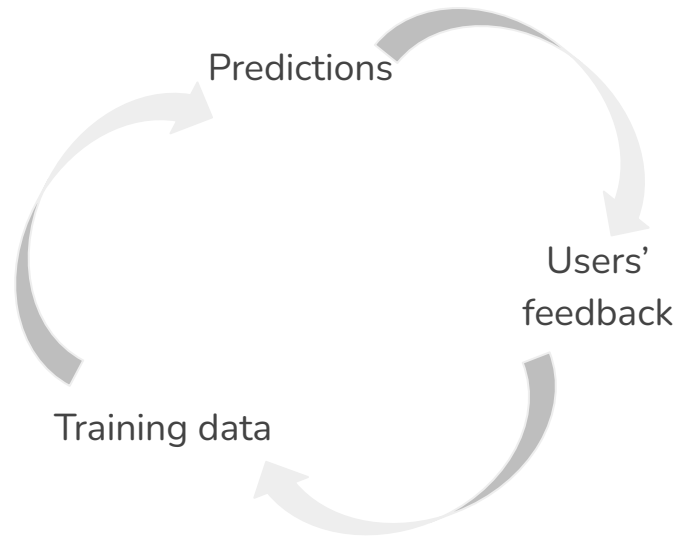


## Cumulative statistics hide sudden shifts



# ML-Specific Failures – Degenerate Feedback Loops (pt. 1)

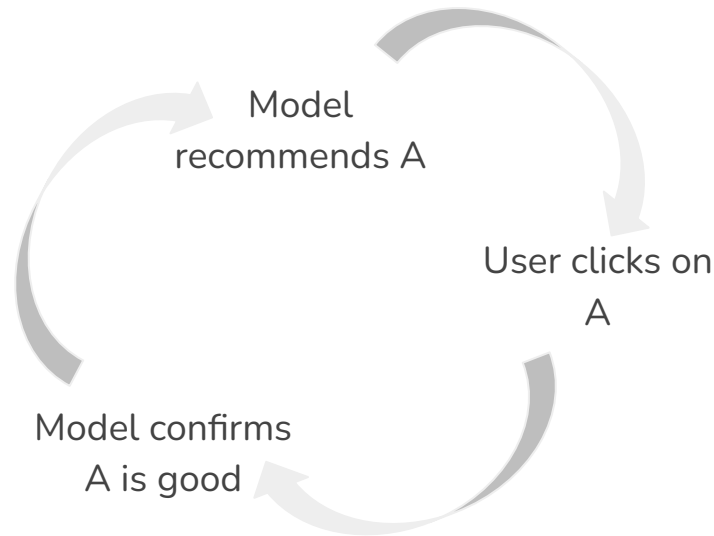
- Typical feedback loop:
  - User feedback from predictions used to train the next iteration of the model to create new predictions



**Only arise once models are in production: hard to detect during development!**

# ML-Specific Failures – Degenerate Feedback Loops (pt. 2)

- Originally, A is ranked marginally higher than B: model recommends A;
- After a while, A is ranked much higher than B



**Only arise once models are in production: hard to detect during development!**

---

# Ethics in AI

# General Ethical Considerations

- Bias and Fairness
- Transparency and Explainability
- Privacy and Data Security
- Accountability
- Autonomy and Control
- Job Displacement and Economic Impact
- Weaponization of AI
- Environmental Impact

# Case Study 1: Automated Grader's Biases

- Summer 2020 - Pandemic
- UK A-level exams
- Failures:
  - Setting the Wrong Objective
  - Insufficient Evaluation
  - Lack of Transparency

## Case Study 2: The Danger of “Anonymized” Data

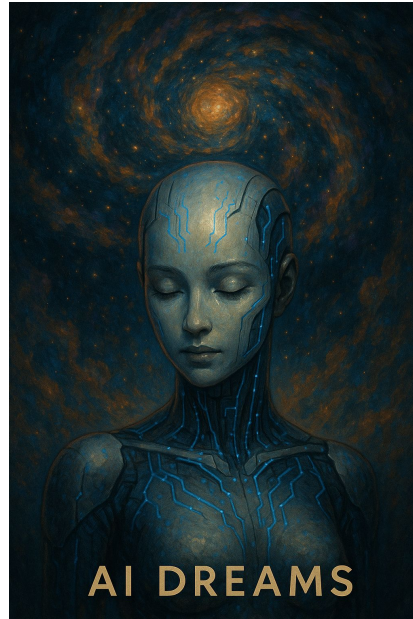
- Strava Fitness App (2018)
- PII was “Anonymized”
- Identified US Military Patrol Routes
- “opt-out” vs “opt-in” Confusion
- Further Potential for PII Misuse

# Acknowledgement

- Vaswani, A. "Attention is all you need." *Advances in Neural Information Processing Systems* (2017).
- Huyen, Chip. *Designing machine learning systems*. "O'Reilly Media, Inc.", 2022.
- ChatGPT



"AI Dreams" - ChatGPT-4o (Oct 1, 2024)



"AI Dreams" - ChatGPT-5.1 (Nov 30, 2025)

O'REILLY®

## Designing Machine Learning Systems

An Iterative Process  
for Production-Ready  
Applications



Chip Huyen

Q&A