

Deploying AI

Introduction to AI Systems

```
$ echo "Data Science Institute"
```

Introduction

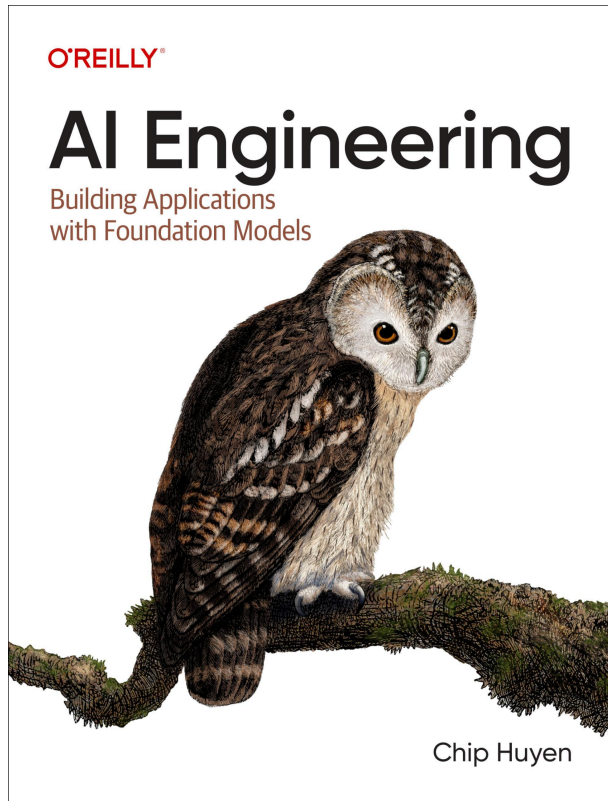
Agenda

Agenda

- What is an AI System?
- Use cases and planning an AI application
- The AI engineering Stack

AI Engineering

We will be covering Chapter 1 of AI Engineering, by Chip Huyen.



Main Points

What is an AI System?

What is an AI System?

- Foundation models
 - Language models
 - Self-supervision
 - From language models to foundation models
- From foundation models to AI engineering

What is an AI System?

- It is a system based on a large-scale machine learning model
- Many principles of productionizing AI applications are similar to those applied in machine learning engineering
- However, the availability of large-scale, readily available models affords new possibilities, and also carries risks and challenges

What Makes AI Different?

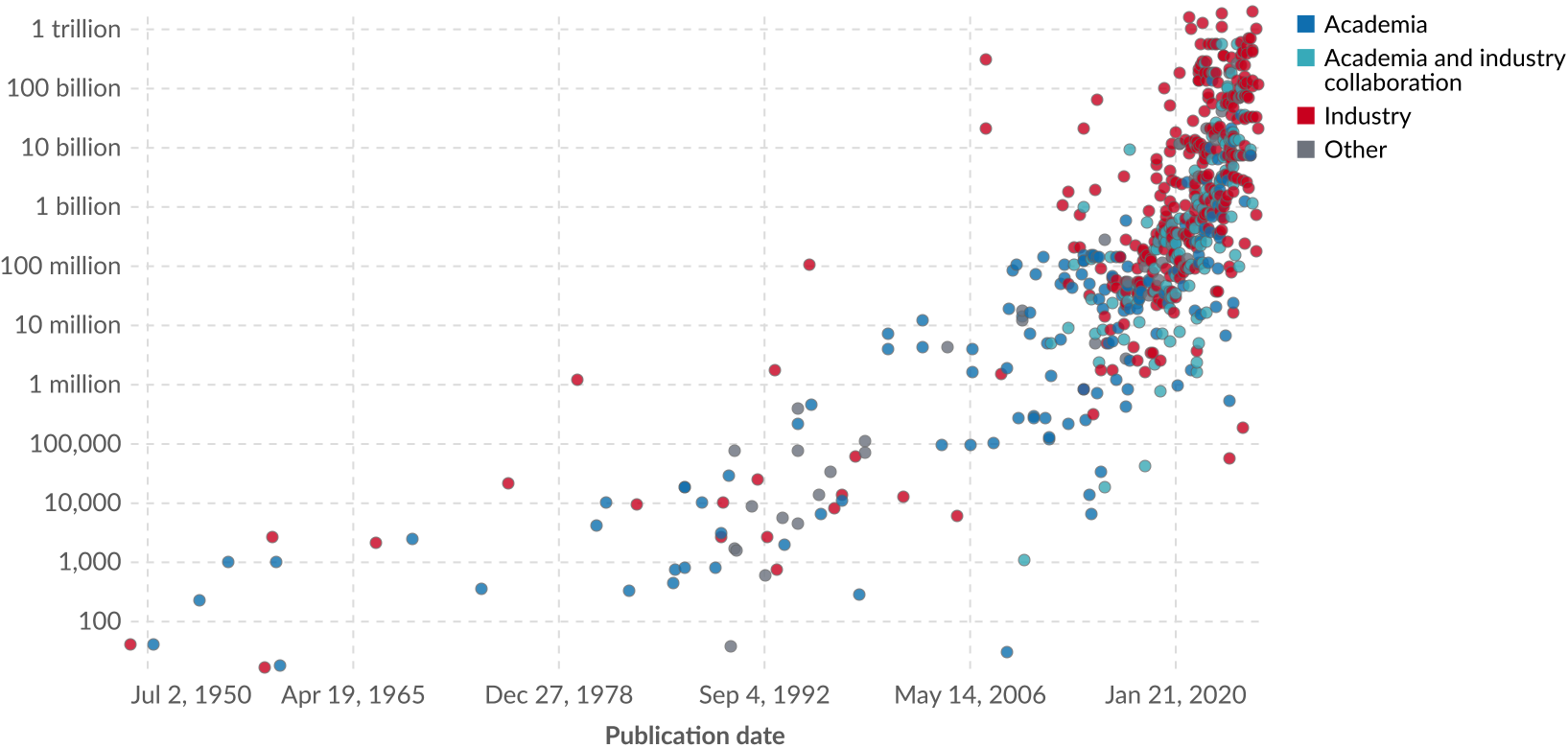
- AI is different because of scale
- Large Language Models (LLMs) and other Foundation Models (FMs) follow a maximalist approach to creating models: more complex models are trained on more data as more compute and storage become available
- FMs are becoming capable of more tasks and therefore they are deployed in more applications and more teams leverage their capabilities
- FMs require more data, compute resources, and specialized talent

Parameters in notable artificial intelligence systems



Parameters are variables in an AI system whose values are adjusted during training to establish how input data gets transformed into the desired output; for example, the connection weights in an artificial neural network.

Number of parameters



Data source: Epoch (2025)

OurWorldinData.org/artificial-intelligence | CC BY

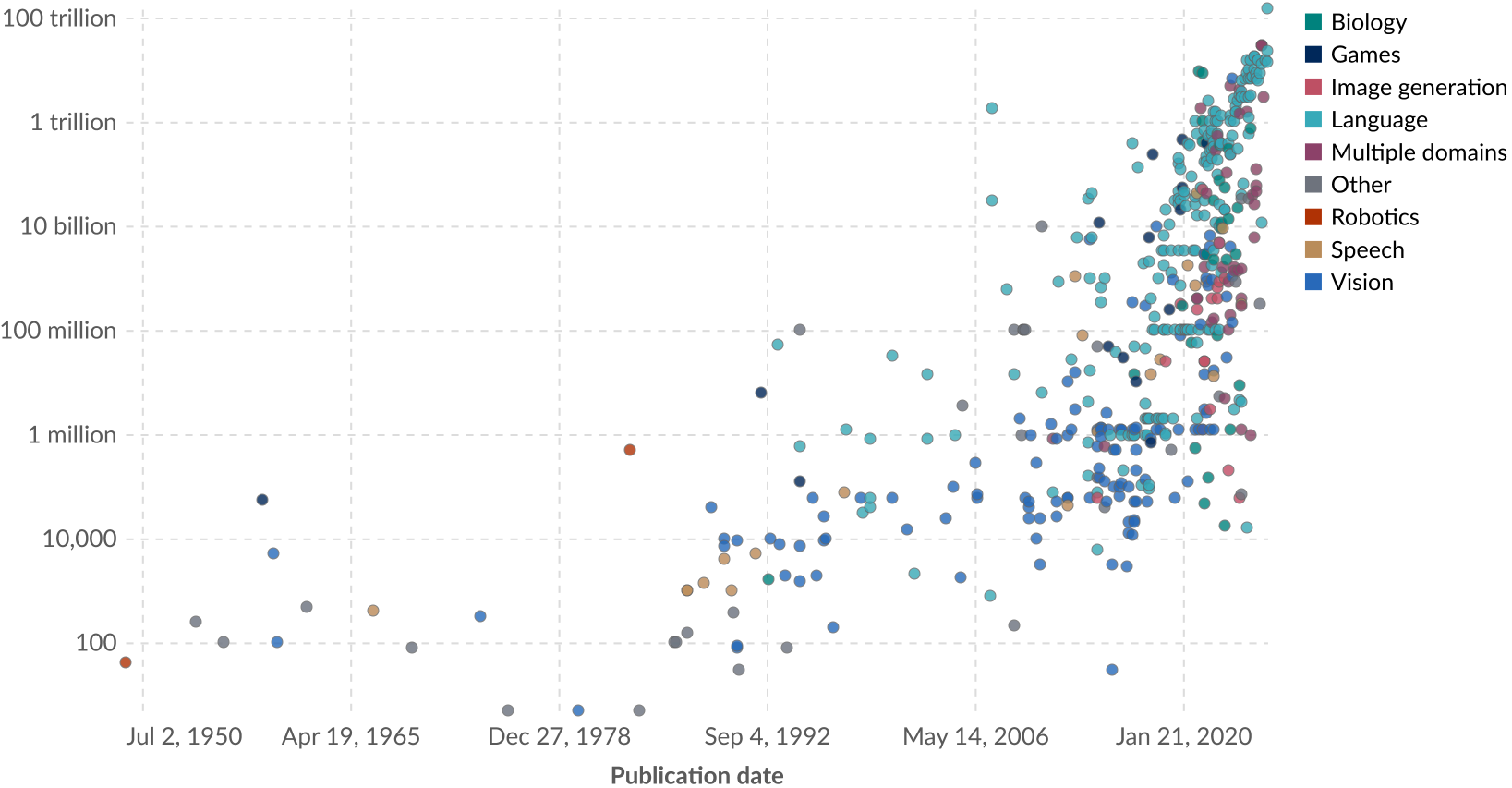
Note: Parameters are estimated based on published results in the AI literature and come with some uncertainty. The authors expect the estimates to be correct within a factor of 10.

Datapoints used to train notable artificial intelligence systems



Each domain has a specific data point unit; for example, for vision it is images, for language it is words, and for games it is timesteps. This means systems can only be compared directly within the same domain.

Training datapoints



Data source: Epoch (2025)

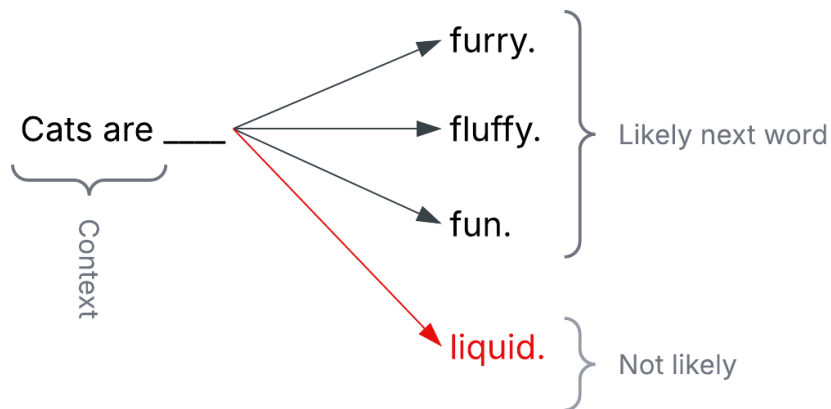
OurWorldinData.org/artificial-intelligence | CC BY

What Makes AI Engineering Different?

- FMs are costly to create, develop, deploy, and maintain. Only a few organizations have the capabilities to do so and typical applications are built upon Models-as-a-Service
- AI Engineering is the process of building applications on top of readily available models

Language Models

- FMs emerged from LLMs which developed from language models
- Language models are not new, but have recently developed greatly through *self-supervision*
- A language model encodes statistical information about one or more languages. Intuitively, we can use this information to know how likely a word is to appear in a given context



Foundation model use cases

- Coding
- Image and Video Production
- Writing
- Education
- Conversational Bots
- Information Aggregation
- Data Organization
- Workflow Automation

Planning an AI application

- Use Case Evaluation
- Setting Expectations
- Milestone Planning
- Maintenance

The AI engineering Stack

- Three layers of the AI Stack
- AI Engineering vs ML Engineering
- AI Engineering vs Full-Stack Engineering

References

References

- Huyen, Chip. Designing machine learning systems. O'Reilly Media, Inc., 2022