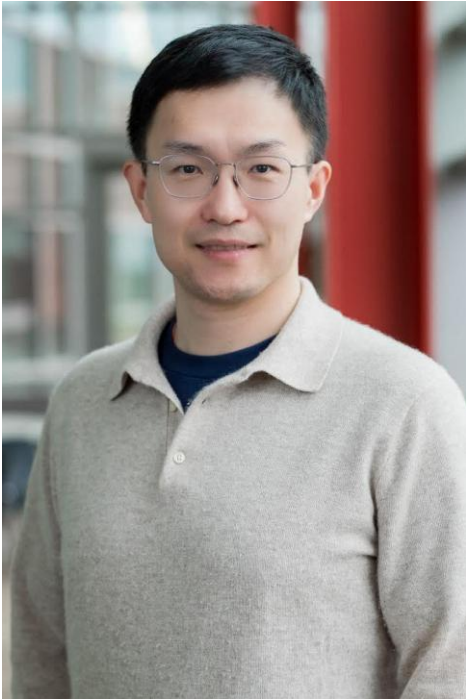




CS 498: Machine Learning System Spring 2026

Minjia Zhang

The Grainger College of Engineering



Minjia Zhang

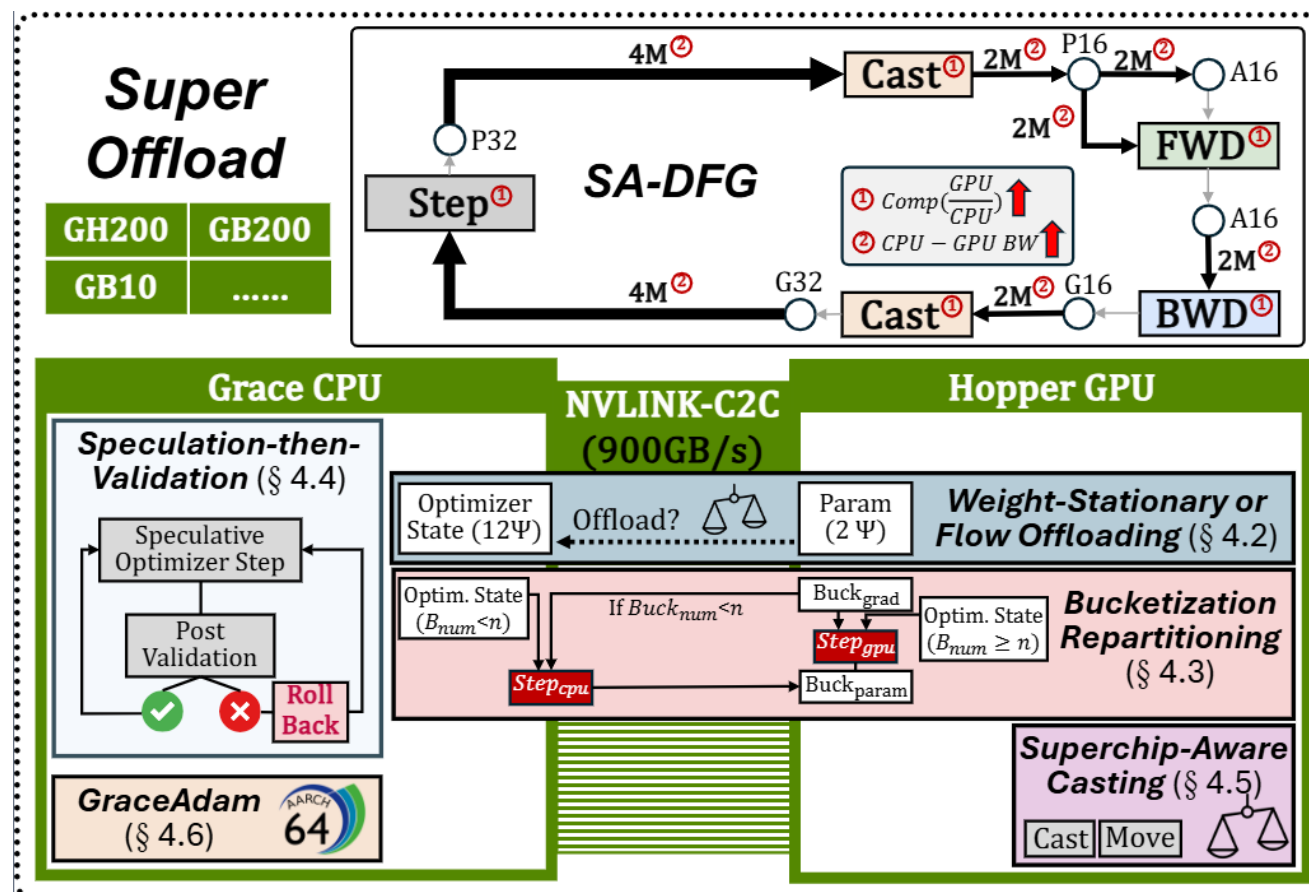
- Now: Assistant Professor at UIUC (2024-Present)
 - Affiliated with ECE and NCSA, UIUC
 - <https://siebelschool.illinois.edu/about/people/faculty/minjiaz>
- Principal researcher at MSR and Microsoft AI (2016-2023)
 - Training and serving LLM at scale, DeepSpeed, Megatron-DeepSpeed

Research area: Systems + Machine Learning

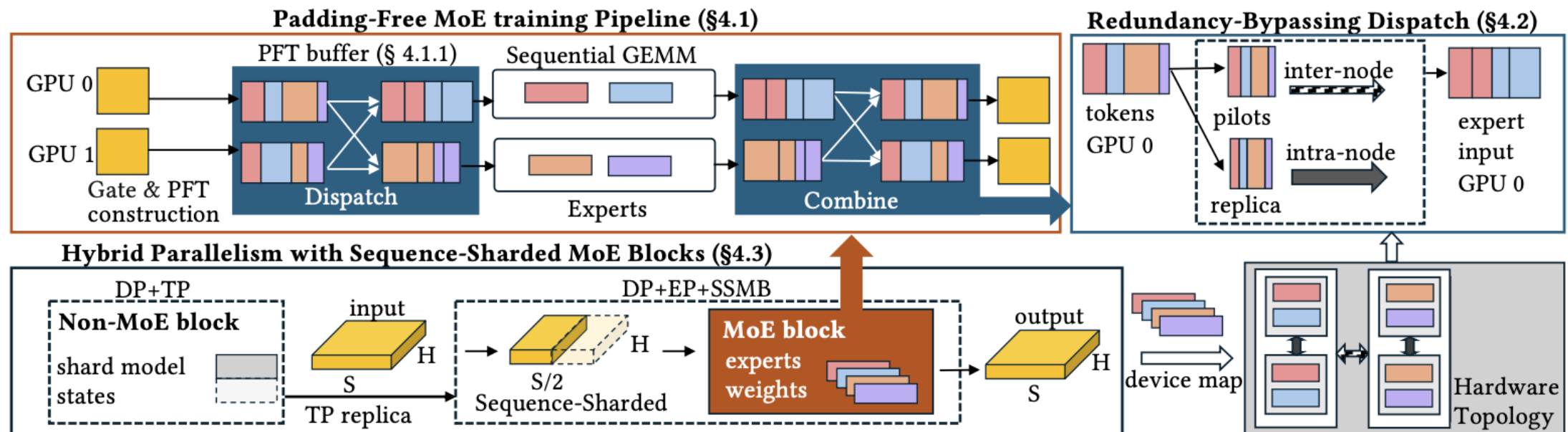
Topics:

- Efficient machine learning systems (training/inference on parallel/distributed/heterogeneous hardware)
- Effective efficiency algorithms (post-training, reasoning model, model compression, etc.)
- Large-scale DL/AI applications (Agentic AI, VLM, Image/Video Generation, DLRM, Vector DB, etc)

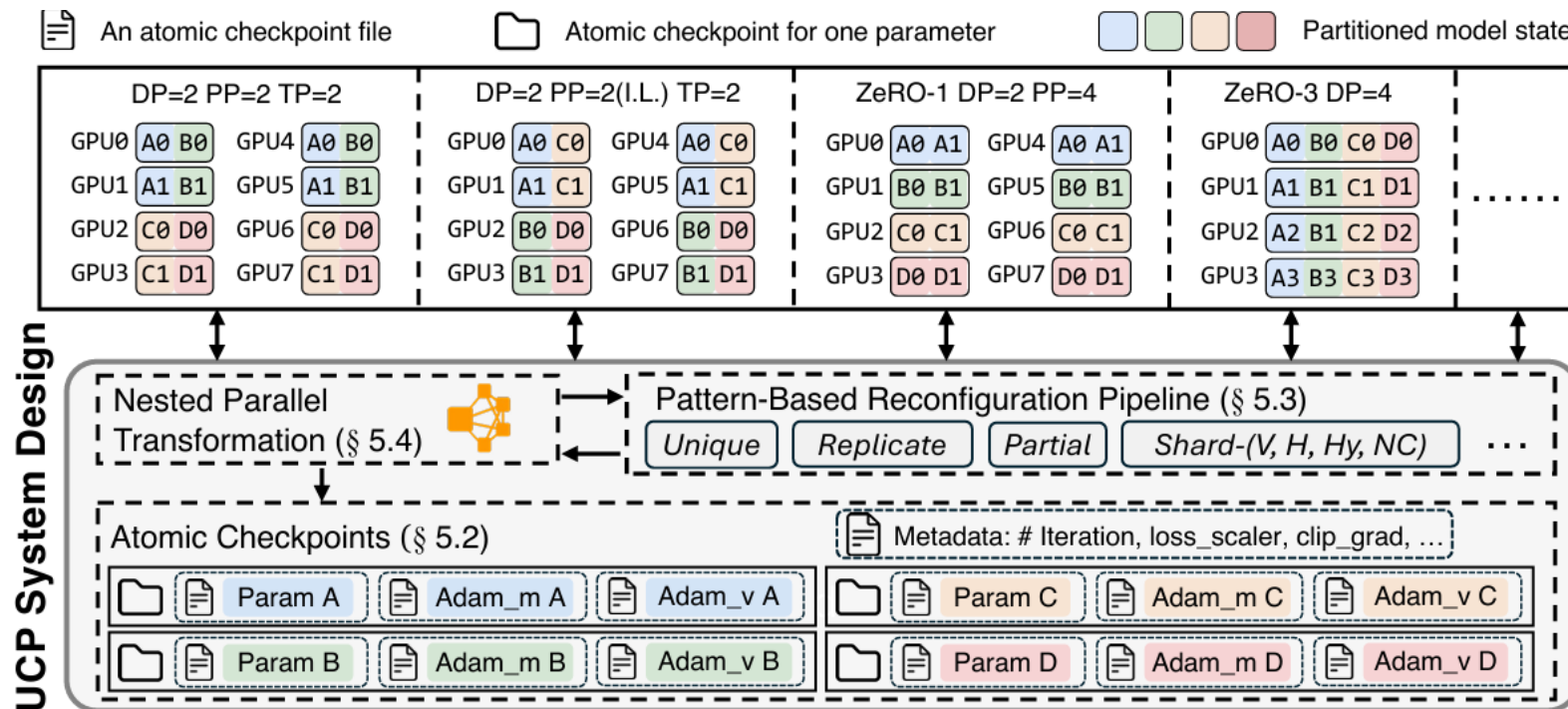
- **Emerging hardware:** Tightly coupled CPU--GPU architectures (e.g., NVIDIA GH200) fundamentally change the cost model of offloading.
- Rethink offloading under high-bandwidth chip-to-chip interconnects for **maximal LLM trainability and speed** on Superchips
- Up to 25B model training on a single GH200, and 1 million sequence length training on 8 GH200 while achieving high MFUs

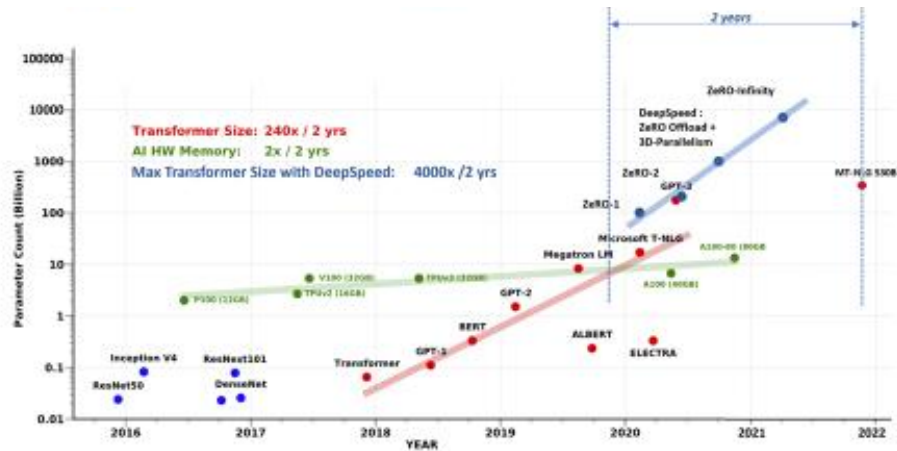


- **Fine-grained MoEs** amplify communication and memory overhead, breaking standard parallelism strategies
- Designed special features to enable DeepSeek-style MoEs training with up to 545 billion parameters across thousand-GPU scale



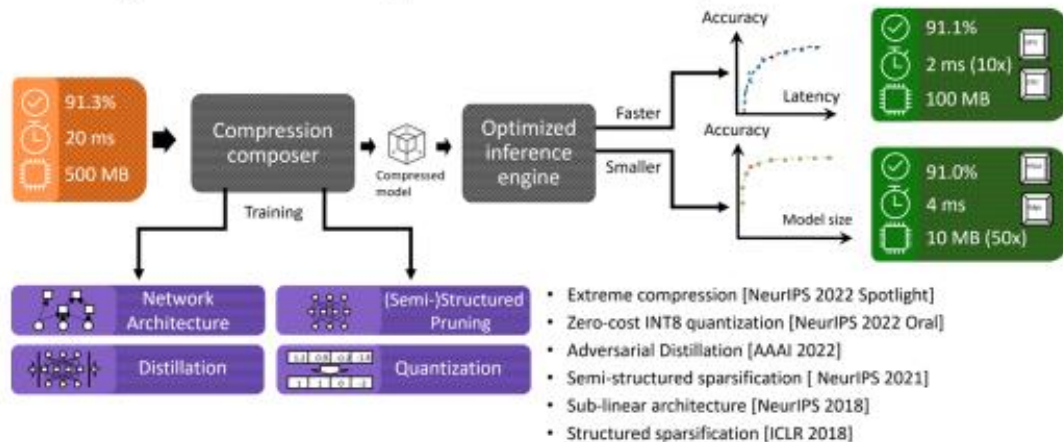
- At scale, failures and reconfiguration are normal, but existing distributed checkpointing assumes static parallelism.
- **Reconfigurable parallelism**: Correctness-preserving transformation pipeline across different parallel training strategies.
- **Failure-tolerant and resilient training** at massive scale under flexible 3D parallel, ZeRO/FSDP, and expert/sequence parallelism



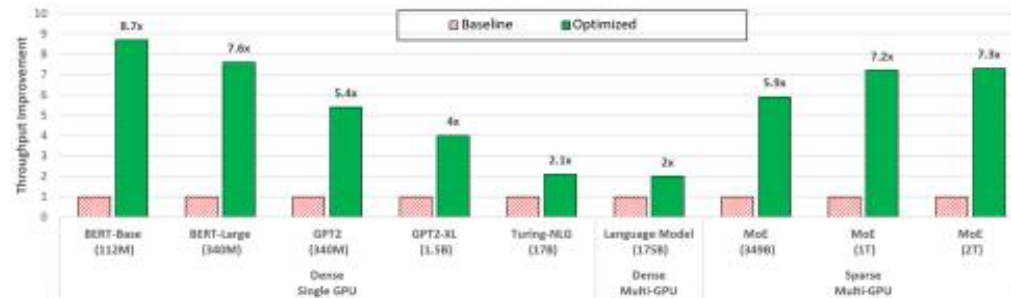


DNN Training at Scale and Speed: Breaking the Memory Wall and Beyond

[NSDI'24, 23, ASPLOS' 23, PPoPP'23, ICLR'23, ICML'22, NeurIPS'22, USENIX ATC'21, HPCA'21,...]

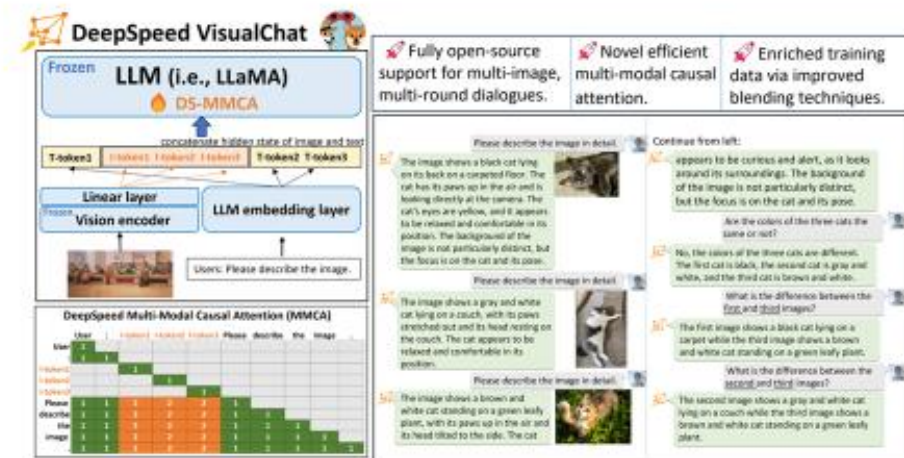


Smaller, Faster, and Cheaper DNN via
Model Compression



Ultra-Fast LLM Inference

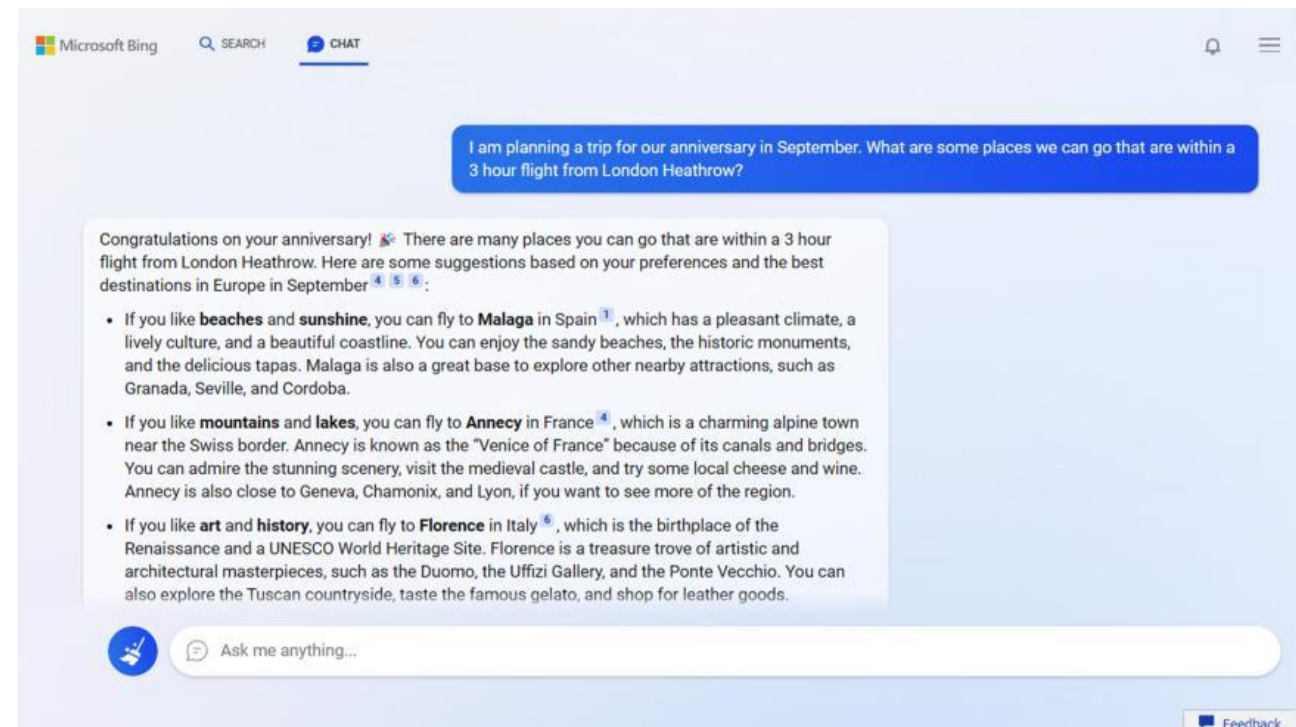
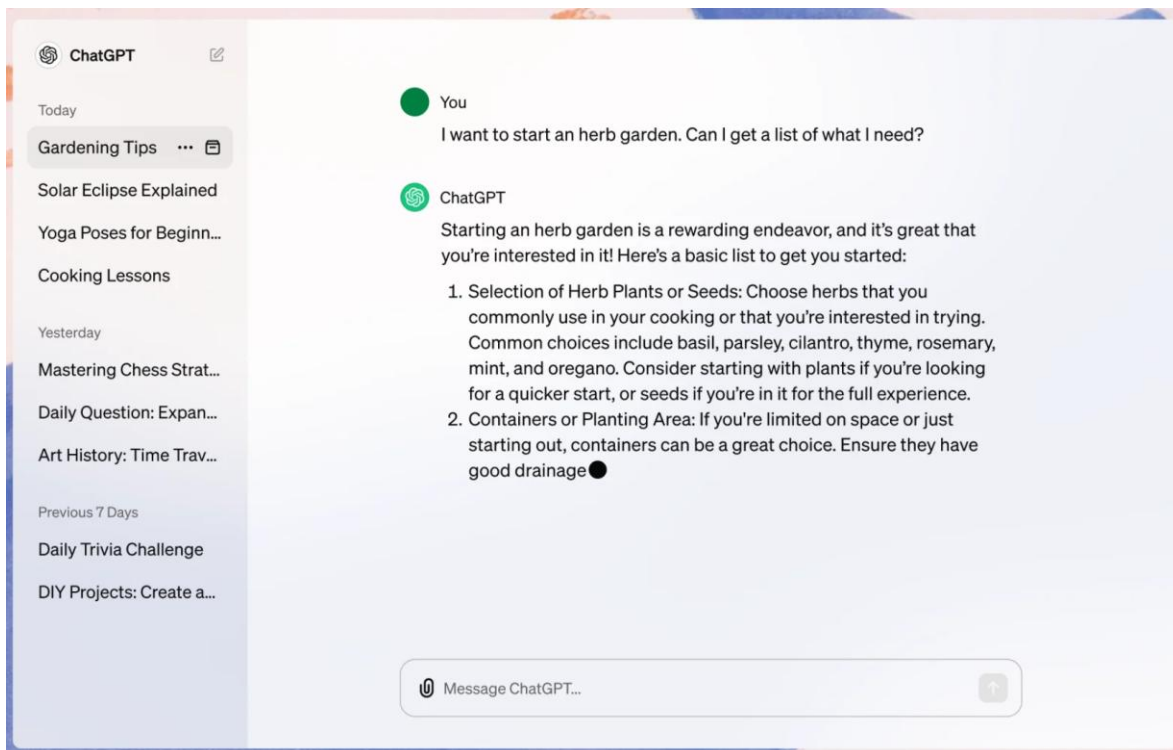
[ICLR'24, SC'22, ICLR'21, NeurIPS'20, USENIX ATC'18,...]



LLM Applications: Chat, Multi-Modal, Diffusion, DiT, Agents, and many more

- **Why study ML Systems?**
- Course overview
- Logistics

The Large Language Model Revolution



ChatGPT: Optimizing Language Models for Dialogue

```
"""
Python 3
Get the current value of a Bitcoin in US dollars using the bitcoincharts api
"""

import requests
import json

def get_bitcoin_price():
    url = 'http://api.bitcoincharts.com/v1/weighted_prices.json'
    response = requests.get(url)
    data = json.loads(response.text)
    return data['USD']['7d']

if __name__ == '__main__':
    print(get_bitcoin_price())
```

[Suggest code and entire function in your editor – Github/OpenAI Codex](#)

Image/Video Generation from Text



TEXT DESCRIPTION

An astronaut Teddy bears A bowl of
soup

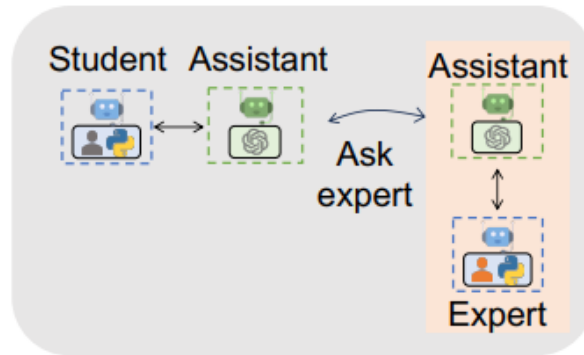
that is a portal to another dimension that
looks like a monster as a planet in the
universe

as digital art in the style of
Basquiat drawn on a cave wall

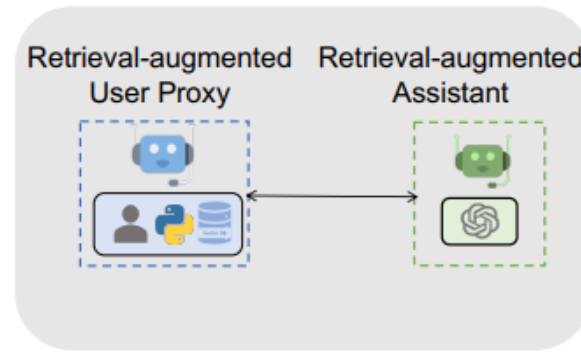


DALL·E 2

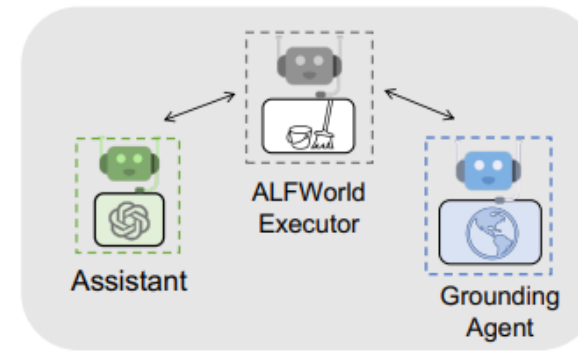




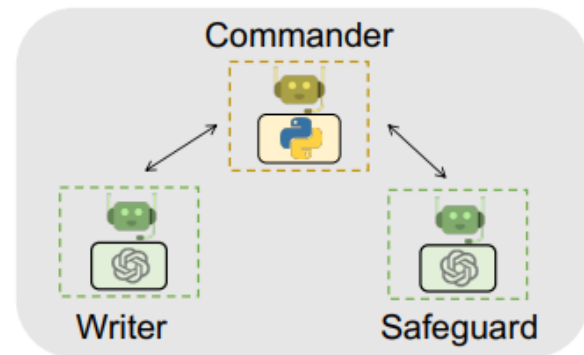
A1. Math Problem Solving



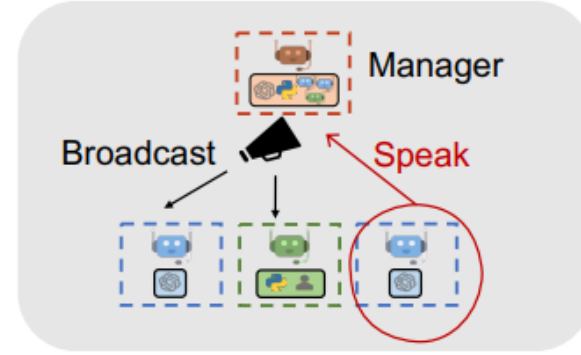
A2. Retrieval-augmented Chat



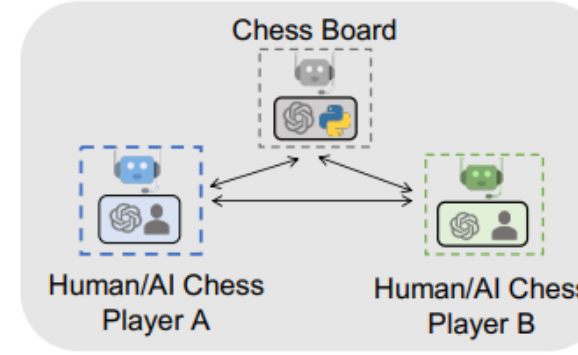
A3. ALF Chat



A4. Multi-agent Coding



A5. Dynamic Group Chat



A6. Conversational Chess

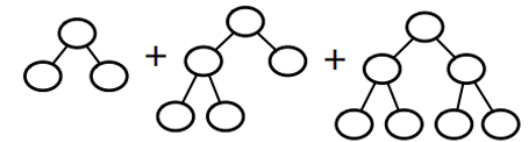
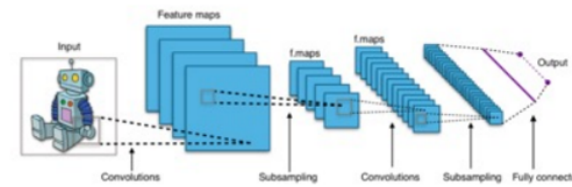
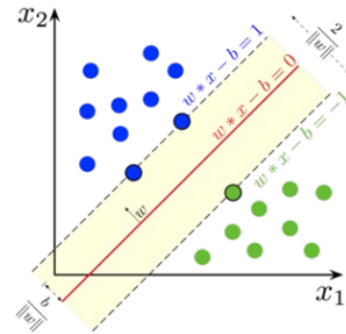
[AutoGen: Enabling Next-Gen LLM Applications via Multi-Agent Conversation](#)



[Using Generative AI to Enable Robots to Reason and Act](#)

How does this Happen?

A key ingredient: ML Systems



Perceptron
Algorithm

Backprop

Support Vector
Machine (SVM)

ConvNet

Gradient Boosting
Machine (GBM)

1958

1986

1992

1998

1999

Many algorithms we use today
are **created before 2000**



2001

flickr

2004

MTurk

2005



2009

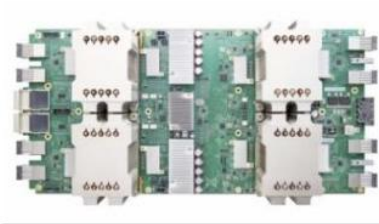
kaggle

IMAGENET

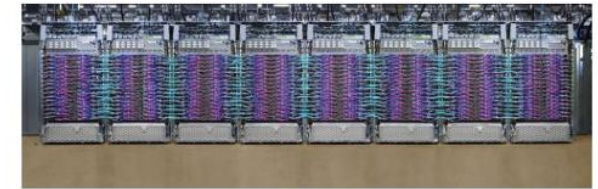
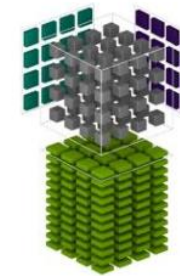
2010

Data serves as fuel for machine learning models

Public
cloud



TensorCore



2006

2007

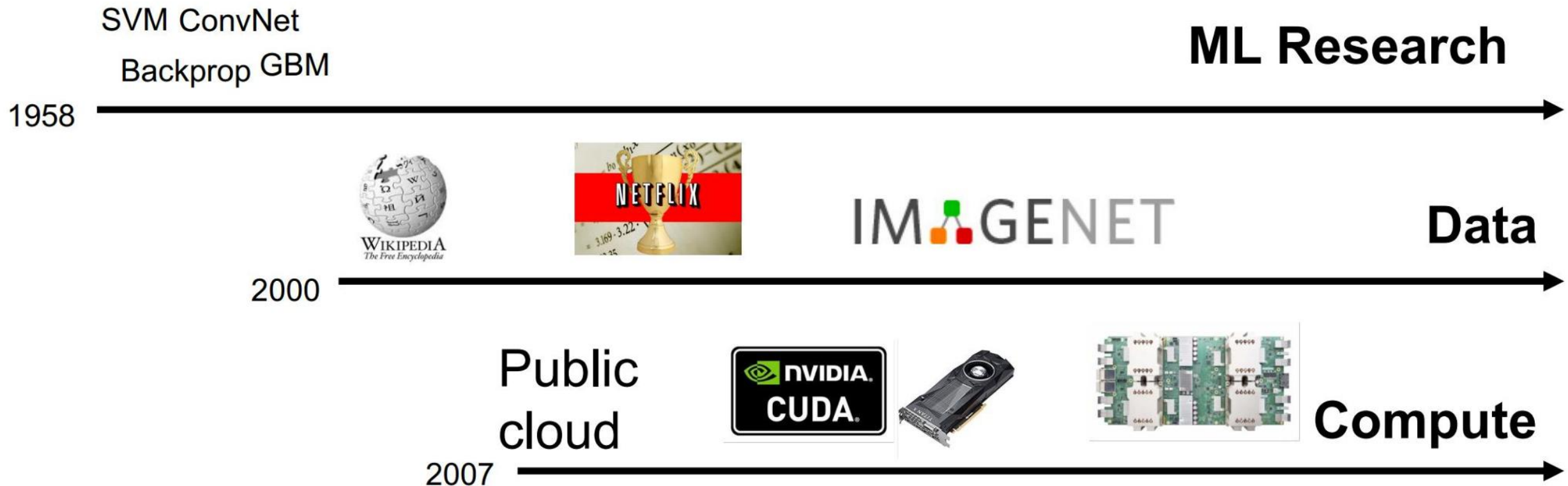
2016

2017

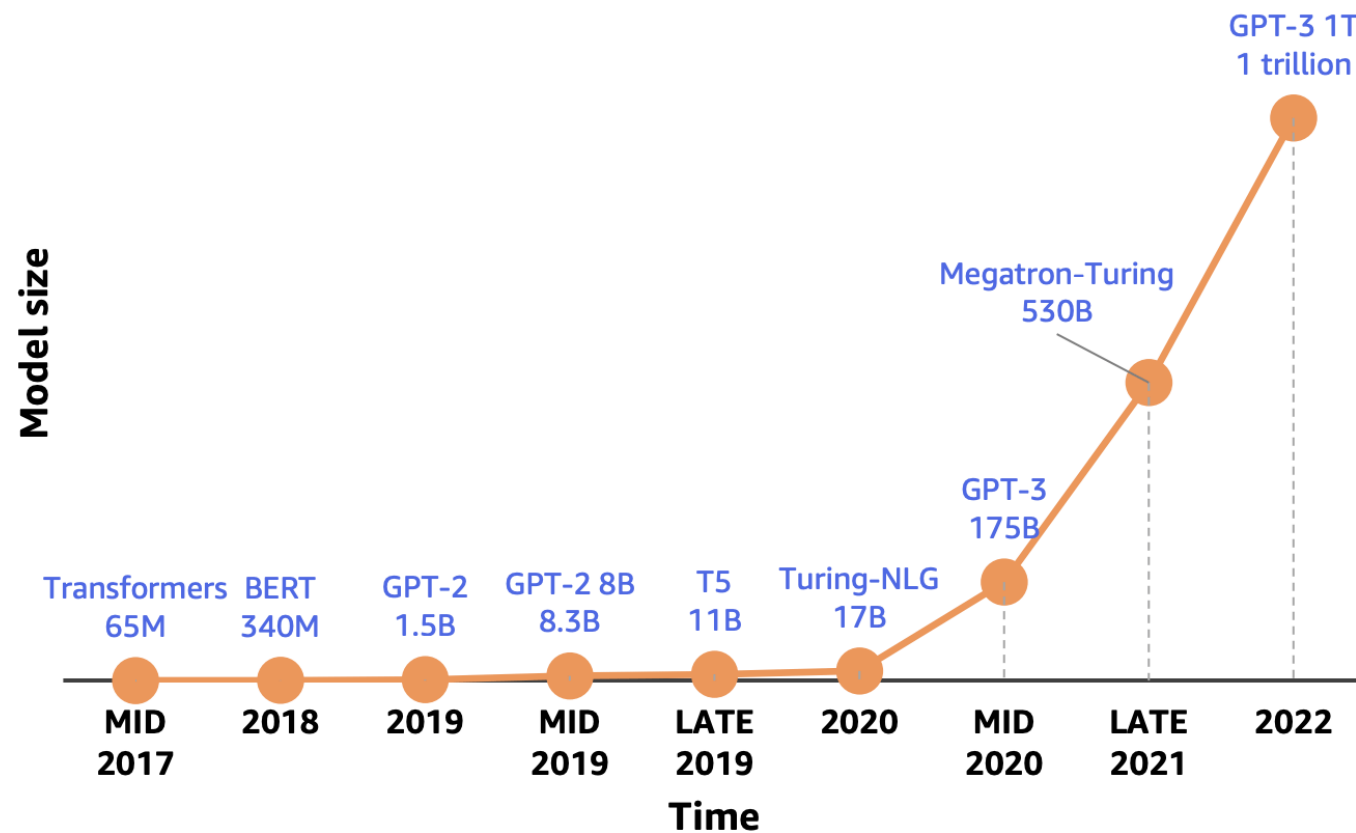
2019

Compute scaling

When three things come together and ready



15,000x increase in 5 years

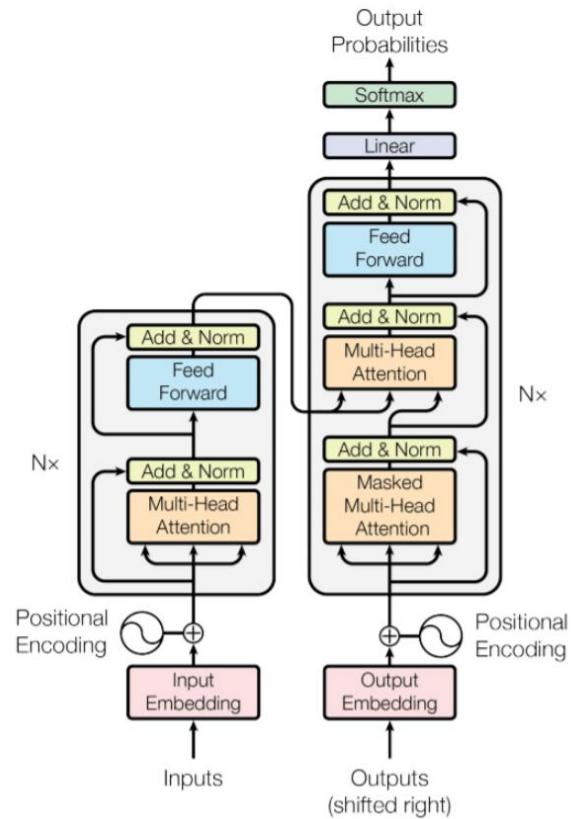


Larger models → better accuracy

Model size is still growing

Not reached the accuracy limit yet

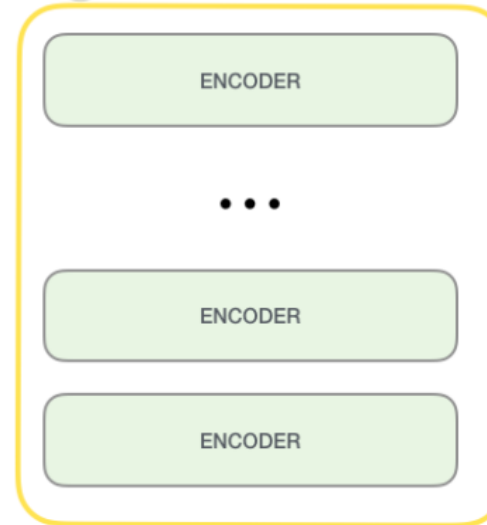
More compute-efficient to train larger models than smaller ones to same accuracy



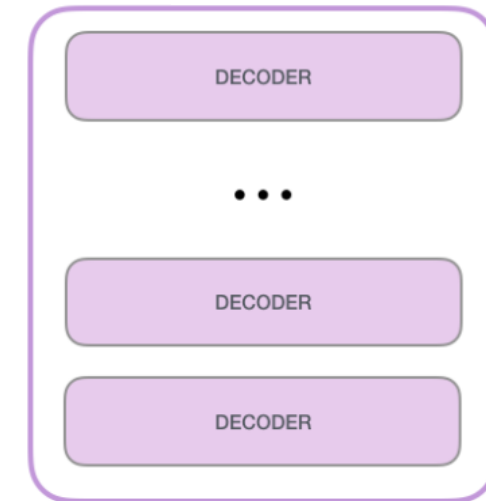
Attention Is All You Need, NeurIPS 2017



BERT



GPT



BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding, ACL 2019

Language Models are Few-Shot Learners, NeurIPS 2020

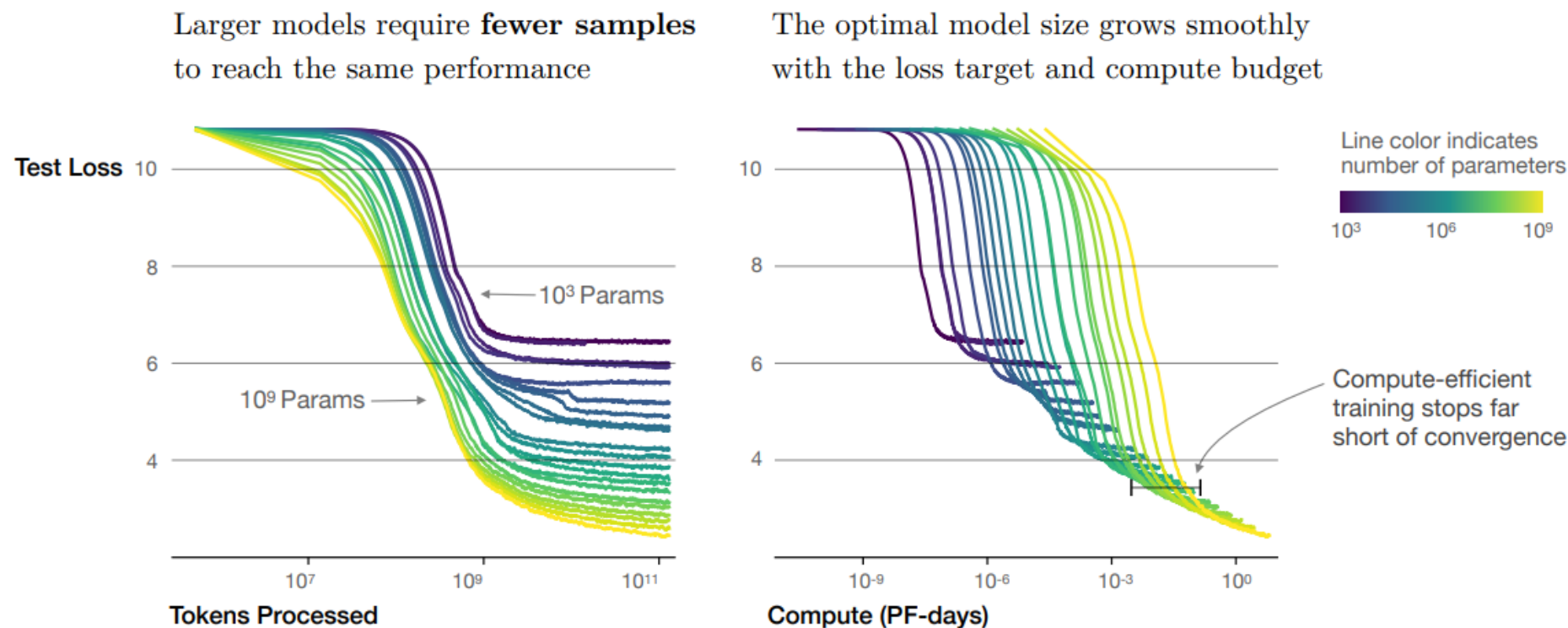


Figure 2 We show a series of language model training runs, with models ranging in size from 10^3 to 10^9 parameters (excluding embeddings).

Scaling Laws for Neural Language Models, OpenAI, 2020

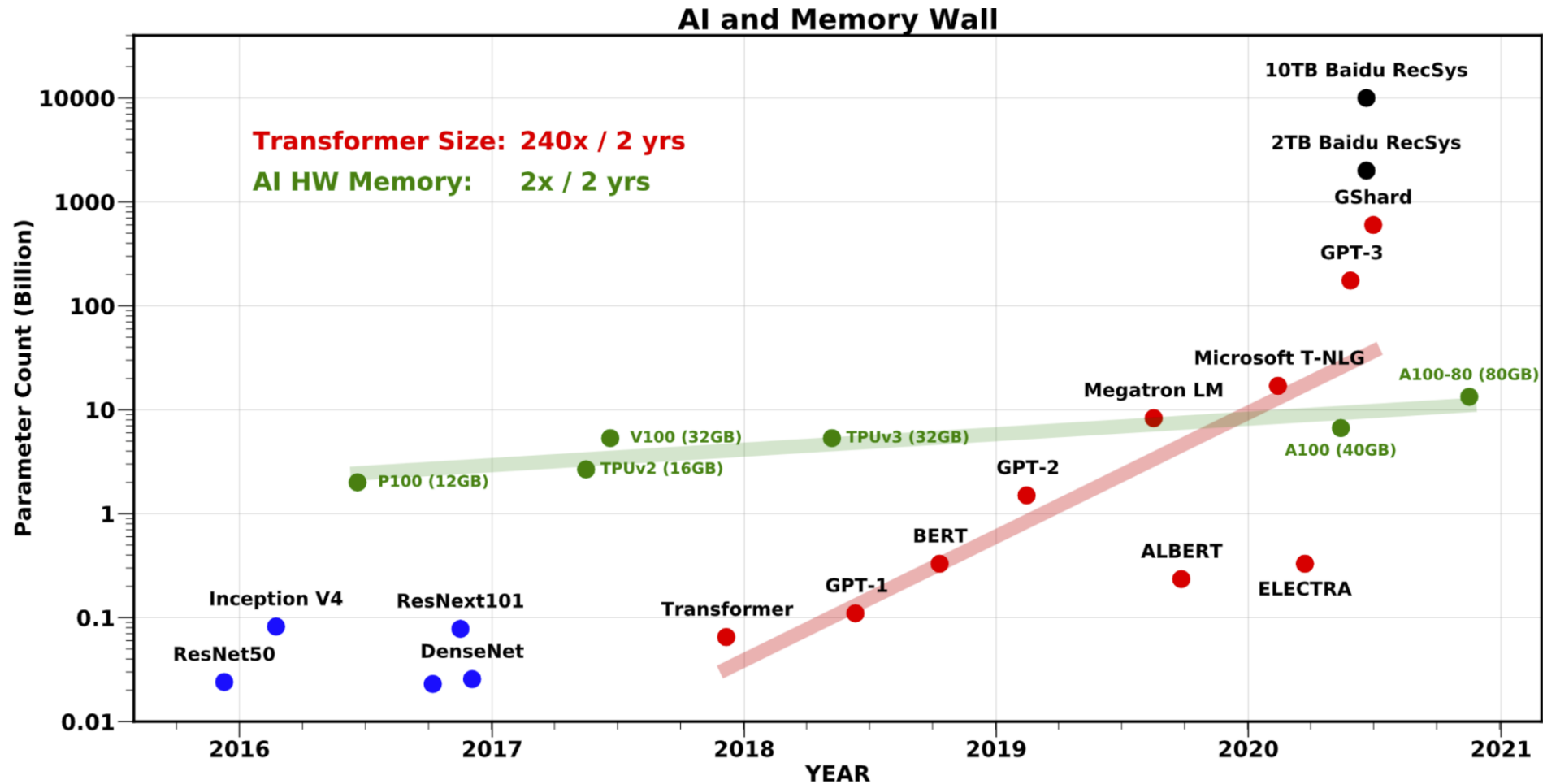
- Too slow to train high-quality models on massive data
 - More hardware \neq higher throughput, bigger model
 - Higher throughput \neq better accuracy, faster convergence, lower cost
 - Better techniques \neq handy to use
- Slow and expensive to deploy the models

Efficiency: Efficient use of hardware for high scalability and throughput

Effectiveness: High accuracy and fast convergence, lowering cost

Easy to use: Improve development productivity of model scientists

DNN Training Hits the Memory Wall

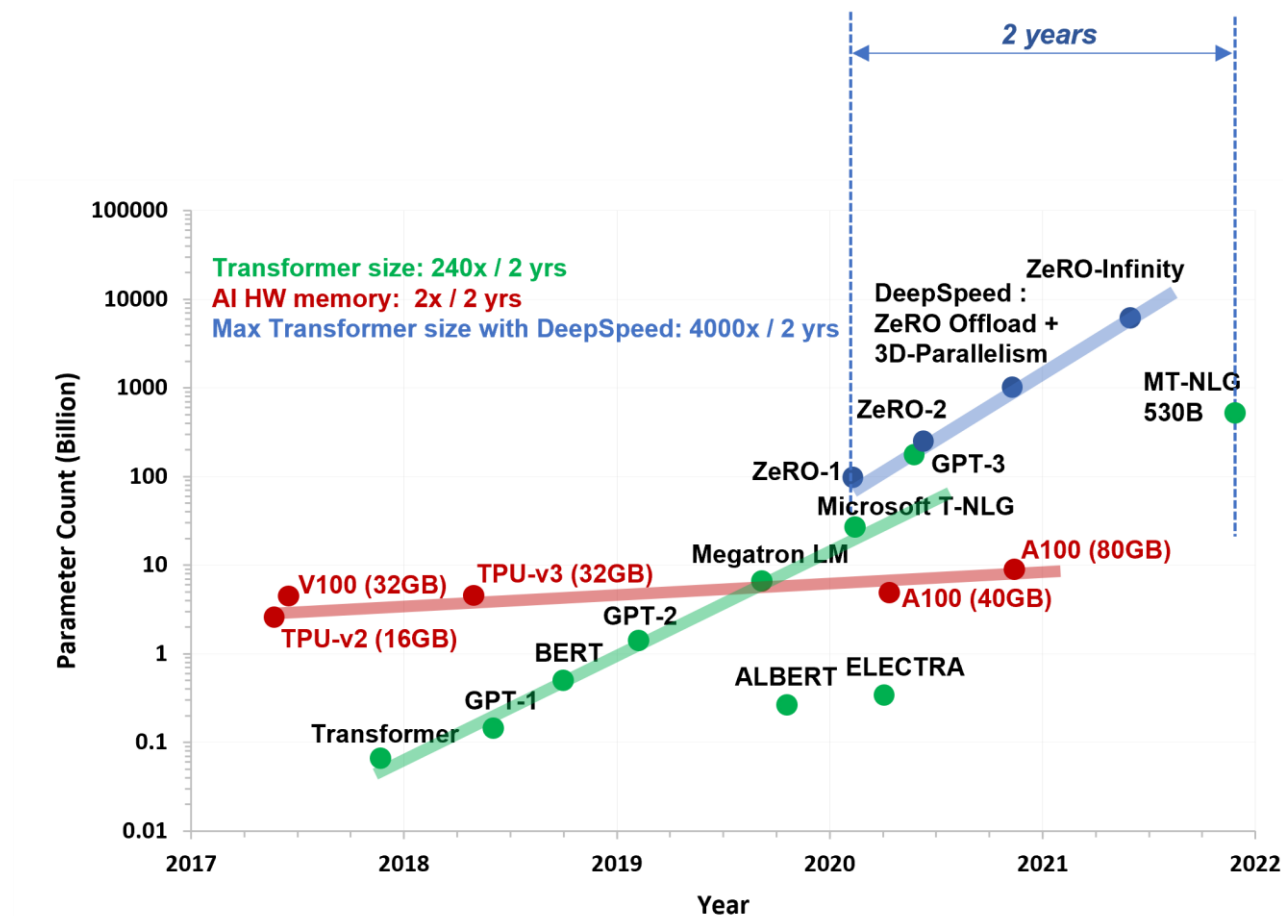


DeepSpeed Powered Massive Models:

- Z-code MoE **10B**
- Microsoft-Turing NLG **17B**
- GPT Neo-X **20B**
- Jurassic-1 **178B**
- Big Science **176B**
- Megatron-Turing NLG **530B**

Key training technologies:

- ☐ ZeRO redundancy optimizer
- ☐ 3D parallelism
- ☐ Optimized CUDA/ROCm/CPU kernels
- ☐ Optimized communication libraries
- ☐ Mixed precision training
- ☐ Communication efficient Adam
- ☐ Sparse Attention
- ☐ Mixture of quantization
- ☐ Curriculum learning
- ☐ ...



Year 2012

Methods

SGD
Dropout
ConvNet
Initialization

Data

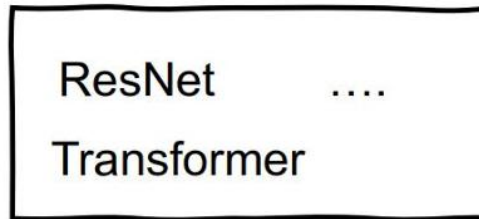
IMGENET

1M labeled
images

Compute

Two GTX 580

Six days



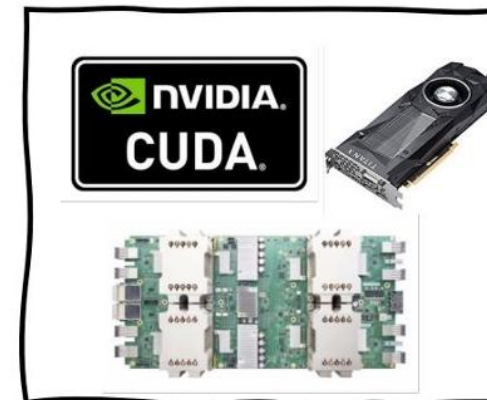
ML Research

44k lines of code

Six months



Data



Compute



ResNet
Transformer

ML Research

100 lines of python

A few hours

System Abstractions

Systems (ML Frameworks)



IMAGENET

Data

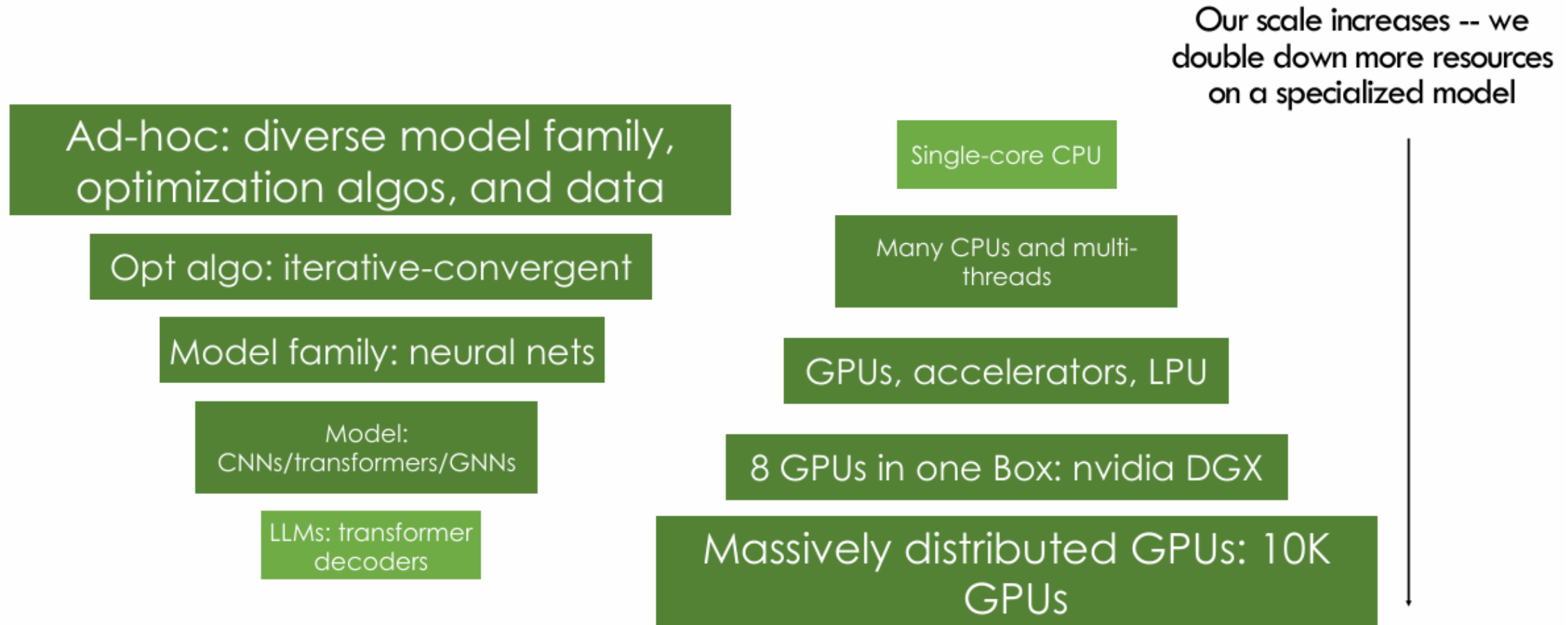


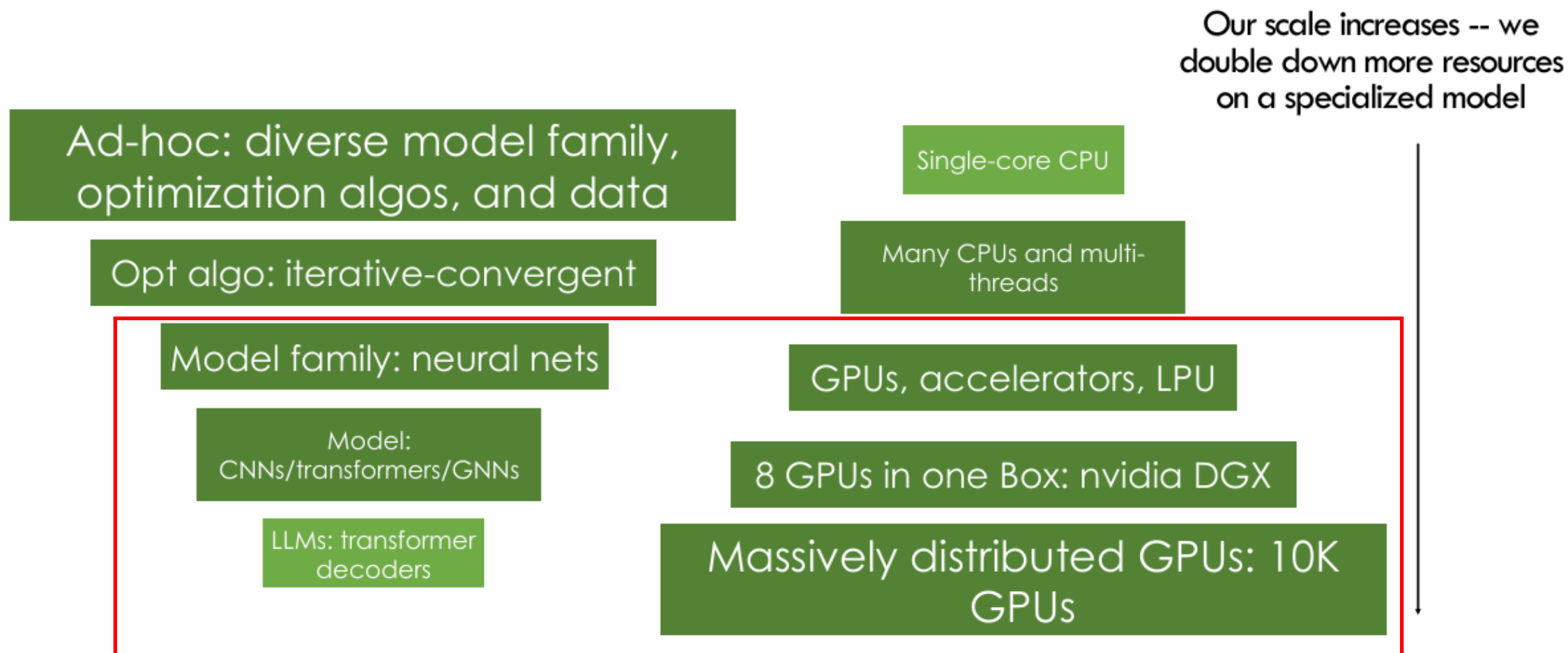
Compute



- Enable new system capabilities to break the memory wall
 - Accelerate ML research
 - Reduce deployment cost
 - Democratize AI to everyone
 - ML \leftrightarrow System codesign
 - ...
-
- In summary: ML System is becoming an essential skill

- Why study ML Systems?
- **Course overview**
- Logistics





- Distributed ML, ML parallelization
- Efficient model adaption
- ML inference optimizations
- Compression algorithms
- AI applications

Machine learning system basic

- Deep Learning workloads
- Computation graph
- ML frameworks

Distributed training strategies to break the memory wall

- Data parallelism
- Tensor parallelism
- Pipeline parallelism
- Sequence parallelism
- Gradient checkpointing
- Auto parallelism

Ultra-fast inference optimizations

- CUDA kernels
- Kernel fusion
- Flash attention
- Paged attention

Compression algorithms to make model smaller, faster, and cheaper

- Quantization
- Sparsification
- Low rank decomposition
- Distillation

- Mixture-of-Experts
- Speculative decoding
- LoRA
- RAG
- Agents
- ...

By the end of this course, you will

- Understand the basic functioning of modern DL libraries, including concepts like compute operators, automatic differentiation, etc.
- Understand the full pipeline of modern ML systems, starting from pre-training all the way to serving...
- Understand scaling-up, why and how? All sorts of machine learning parallelization techniques, and latest research in the area ...
- Understand hardware acceleration/CUDA/GPUs, and can program/debug a little accelerator programs ...
- Ground all you learning in the context of LLMs, understand the L of LLM, how it is optimized, scaled, trained, served...
- Have fun!

Questions?

CS 498 Machine Learning System, Spring 2026

Basic Information

Instructor: [Minjia Zhang](#)

Schedule: Tuesdays and Thursdays 3:30-4:45pm CST

Location: 2018 Campus Instructional Facility

Instructor Email: minjiaz AT illinois.edu

TA: TBD

TA Email: TBD

TA Office hours: TBD LMS: [Canvas](#)

Recommended Prerequisites: [CS 425 - Distributed Systems](#) [CS 484 - Parallel Programming](#), [CS 533 - Parallel Computer Architecture](#), [CS 446 - Machine Learning](#)

Course Description

Welcome to the Spring 2026 offering of CS 498: Machine Learning System!

This is a new undergraduate course offered for the second time at UIUC. Therefore, we may still adjust the schedule and content depending on your learning progress.

The goal of this course is to provide students with an in-depth understanding of various elements of modern machine learning systems, ranging from the performance characteristics of ML models such as transformers and diffusers, performance optimization techniques that reduce the compute, memory, and communication for training and inference of large ML models, and compression algorithms that make ML models smaller and cheaper. The course will also conduct case studies on modern large language model training and serving and cover the design rationale behind state-of-the-art machine learning frameworks.

Course Schedule

Course Policy

Grading

The course assignments consist of (i) attendance and class participation, (ii) lab assignments, (iii) reading summary, (iv) final project presentation, and (v) completing an open-ended research project. The breakdown is as follows.

Grading Breakdown

Attendance and class participation	20%
Lab assignments	21% (3 lab assignments, 7% each)
Reading summary	20% (10 readings, 2% each)
Final research project presentation	14%
Project report	25% (5% + 5% + 15%)

Paper Reviews

- Attendance and class participation 20%
- Lab assignments 21% (3 lab assignments, 7% each)
- Reading summary 20% (10 readings, 2% each)
- Final research project presentation 15%
- Project report 25% (5% + 5% + 14%)

Grading Scheme (grade is the better of the two)



Grade	Absolte Cutoff (\geq)	Relative Bin
A+	95	Highest 5%
A	90	Next 10%
A-	85	Next 15%
B+	80	Next 15%
B	75	Next 15%
B-	70	Next 15%
C+	65	Next 5%
C+	60	Next 5%
C-	55	Next 5%
D	50	Next 5%
F	<50	Lowest 5%

Grading Scheme (grade is the better of the two)



Grade	Absolte Cutoff (\geq)	Relative Bin
A+	95	Highest 5%
A	90	Next 10%
A-	85	Next 15%
B+	80	Next 15%
B	75	Next 15%
B-	70	Next 15%
C+	65	Next 5%
C	60	Next 5%
C-	55	Next 5%
D	50	Next 5%
F	<50	Lowest 5%

Example, 82
and 33%,
Abs: B+; Rel: B-;
Final: B+

Structure of the Course (Tentative)



Week	Part 1: Basics	
1	Course intro	DL Workloads
2	DL frameworks	AutoDiff
	Part 2: Distributed ML	
3	Overview of training	Communication collectives
4	Data parallelism, tensor parallelism	Pipeline parallelism
5	Zero-style data parallelism	Heterogeneous GPU-CPU
6	3D parallelism	Auto parallelism
7	Mixed precision training	Communication compression
	Part 3: Inference optimizations	
8	CUDA basics	
9	FlashAttention	PagedAttention
10	Continuous batching	Efficient scaling of transformer inference
11	TVM and DL compiler	
	Part 4: Compression	
11	Quantization 1	Quantization 2
12	Sparsification 1	Sparsification 2
13	Distillation	Low-rank decomposition
14	KV cache compression 1	KV cache compression 2
	Part 5: Misc	
15	MoE 1	MoE2
16	Vector db	RAG

- 3 lab assignments
 - Likely will use NCSA clusters for GPUs
 - The instructor needs to figure out some details
- Topics
 - AutoDiff
 - Inference optimizations
 - Compress ML models

- Required reading:
 - The instruction will select 10 highly relevant papers in MLSys.
 - One paper per week (starting from Jan 27), submit your reading by the end of day of each Friday.
 - The reading summary should be done independently and include the following content:
 - The problem the paper is trying to tackle.
 - What's the impact of the work, e.g., why is it an important problem to solve?
 - The main proposed idea(s).
 - A summary of your understanding of different components of the proposed technique, e.g., the purpose of critical design choices.
 - Your perceived strengths and weaknesses of the work, e.g., novelty, significance of improvements, quality of the evaluation, easy-to-use.
 - Is there room for improvement? If so, what idea do you have for improving the techniques?
 - The reading summary length should be around 4-5 paragraphs.

Grading criteria, each summary has 12 points in total. For each review item above, you get:

- 2: The summary item demonstrates a clear understanding of the paper.
- 1: The summary item misses the point of the paper.
- 0: The summary item is missing.

Paper Reviews

The instructor will select 10 highly relevant papers in MLSys (mostly under 12 pages). One paper per week (starting from Jan 27), submit your :

Reading List

The reading summary should be done independently and include the following contents:

- The problem the paper is trying to tackle.
- What's the impact of the work, e.g., why is it an important problem to solve?
- The main proposed idea(s).
- A summary of your understanding of different components of the proposed technique, e.g., the purpose of critical design choices.
- Your perceived strengths and weaknesses of the work, e.g., novelty, significance of improvements, quality of the evaluation, easy-to-use.
- Is there room for improvement? If so, what idea do you have for improving the techniques?

ML System Reading List

[\(3D Parallelism\) Efficient large-scale language model training on GPU clusters using megatron-LM](#)
SC 2021

[\(ZeRO-style Data Parallelism\) ZeRO: Memory Optimizations Toward Training Trillion Parameter Models](#)
SC 2020

[Alpa: Automating Inter- and Intra-Operator Parallelism for Distributed Deep Learning](#)
OSDI 2022

[FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness](#)
NeurIPS 2022

[Orca: A Distributed Serving System for Transformer-Based Generative Models](#)
OSDI 2022

[Efficiently Scaling Transformer Inference](#)
MLSys 2023

[\(vLLM\) Efficient Memory Management for Large Language Model Serving with PagedAttention](#)
SOSP 2023

[SGLang: Efficient Execution of Structured Language Model Programs](#)
NeurIPS 2024

[GPTQ: Accurate Post-Training Quantization for Generative Pre-trained Transformers](#)
ICLR 2023

[Fast Inference from Transformers via Speculative Decoding](#)
ICML 2023

- The course also includes proposing and completing a course project. The project can involve, but is not limited to, any of the following tasks:
 - Benchmark and analyze important DL workloads to understand their performance gap and identify important angles to optimize their performance.
 - Apply and evaluate how existing solutions work in the context of emerging AI/DL workloads.
 - Design and implement new algorithms that are both theoretically and practically efficient.
 - Design and implement system optimizations, e.g., parallelism, cache-locality, IO-efficiency, to improve the compute/memory/communication efficiency of AI/DL workloads.
 - Offer customized optimization for critical DL workloads where latency is extremely tight.
 - Build library/tool/framework to improve the efficiency of a class of problems.
 - Integrate important optimizations into existing frameworks (e.g., DeepSpeed), providing fast and agile inference.
 - Combine system optimization with modeling optimizations.
 - Combine and leverage hardware resources (e.g., GPU/CPU, on-device memory/DRAM/NVMe/SSD) in a principled way.
 - ...
- The project will be done in groups of 2-3 people, which consists of a proposal, mid-term report, final presentation, and final report. The tentative timeline for the project is as follows.

- 15%
- Please spend a significant amount of time on working your project and making this presentation nice and clear.
- Graded by instruction team (50%) and your classmates (50%)
 - Instructor: based on format, correctness, depth, clarity, insights
 - Peers: make sure your classmates feel they indeed learn something after listening to your presentation
- Happening in the end of the semester

- Final report: The final report will be in the style of a research paper describing your project. The recommended length is about **6-8 pages** long (excluding references) and a potential division can be: An abstract, which summarizes the project (0.25 pages).
 - An introduction, which describes and motivates the problem and summarizes the main results of the work (0.5 pages).
 - A brief discussion of related work (0.25 pages).
 - A brief overview of preliminary and background knowledge needed to understand the paper (0.25 pages).
 - Analysis and characterization to show the existence and severity of the problem (1 page).
 - Main design and implementation (1 pages).
 - Evaluation methodology and experiment results (1 page).
 - Concluding remarks, which can include a discussion of open questions or directions for future work (0.25 pages).

Formatting Instructions For NeurIPS 2020

David S. Hippocampus*
Department of Computer Science
Cranberry-Lemon University
Pittsburgh, PA 15213
hippo@cs.cranberry-lemon.edu

Abstract

The abstract paragraph should be indented 1/2 inch (3 picas) on both the left- and right-hand margins. Use 10 point type, with a vertical spacing (leading) of 11 points. The word **Abstract** must be centered, bold, and in point size 12. Two line spaces precede the abstract. The abstract must be limited to one paragraph.

1 Submission of papers to NeurIPS 2020

NeurIPS requires electronic submissions. The electronic submission site is

<https://cm3.research.microsoft.com/NeurIPS2020/>

Please read the instructions below carefully and follow them faithfully.

L1 Style

Papers to be submitted to NeurIPS 2020 must be prepared according to the instructions presented here. Papers may only be up to eight pages long, including figures. Additional pages containing only a section on the broader impact, acknowledgments and/or cited references are allowed. Papers that exceed eight pages of content will not be reviewed, or in any other way considered for presentation at the conference.

The margins in 2020 are the same as those in 2007, which allow for ~15% more words in the paper compared to earlier years.

Authors are required to use the NeurIPS L^AT_EX style files obtainable at the NeurIPS website as indicated below. Please make sure you use the current files and not previous versions. Tweaking the style files may be grounds for rejection.

L2 Retrieval of style files

The style files for NeurIPS and other conference information are available on the World Wide Web at

<http://www.neurips.cc/>

The file `neurips_2020.pdf` contains these instructions and illustrates the various formatting requirements your NeurIPS paper must satisfy.

The only supported style file for NeurIPS 2020 is `neurips_2020.sty`, rewritten for L^AT_EX 2_ε. Previous style files for L^AT_EX 2.09, Microsoft Word, and RTF are no longer supported!

*Use footnote for providing further information about author (webpage, alternative address)—not for acknowledging funding agencies.

34th Conference on Neural Information Processing Systems (NeurIPS 2020), Vancouver, Canada.

- All assignments are due on the respective due date. Only on-time assignments will be accepted.

Questions?