



CS6501-003: Datacenter Infrastructure

- Course Overview and Logistics

Qizhe Cai

About me

- **Qizhe Cai**
 - **Assistant Professor**, UVA (started 1 months ago!)
 - **Previously:** Ph.D from Cornell; M.S. from Princeton; Undergrad from Umich
 - **Office:** Rice 102
 - **Office hour:** 2:30pm Monday
- **Research interests**
 - At the intersection of **networking, systems and hardware**
 - Publish in conferences like SIGCOMM, NSDI, OSDI and SOSP
- **Non-research interests**
 - Sports: soccer, football, and gaming (mainly watching these days)

This Course

Motivation: LLMs are part of everyday life

- A LLM is a neural network–based AI trained on vast text corpora to understand and generate human-like language

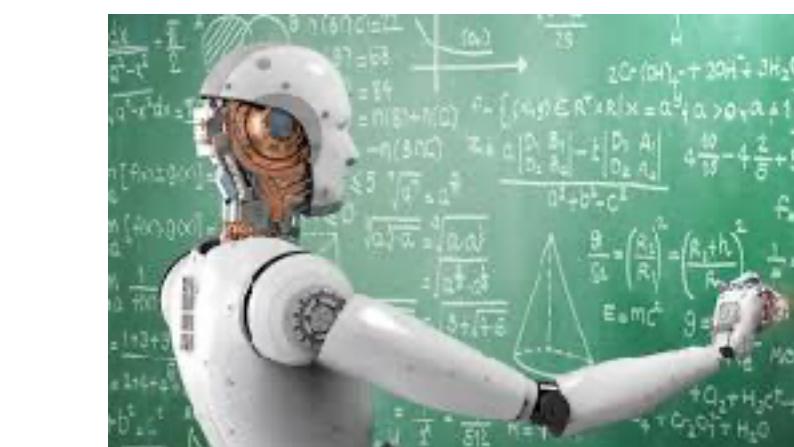
- LLM is on everyday's life



- Chatbots
- Coding
- Education & Learning
- Healthcare & Well-being



GitHub
Copilot

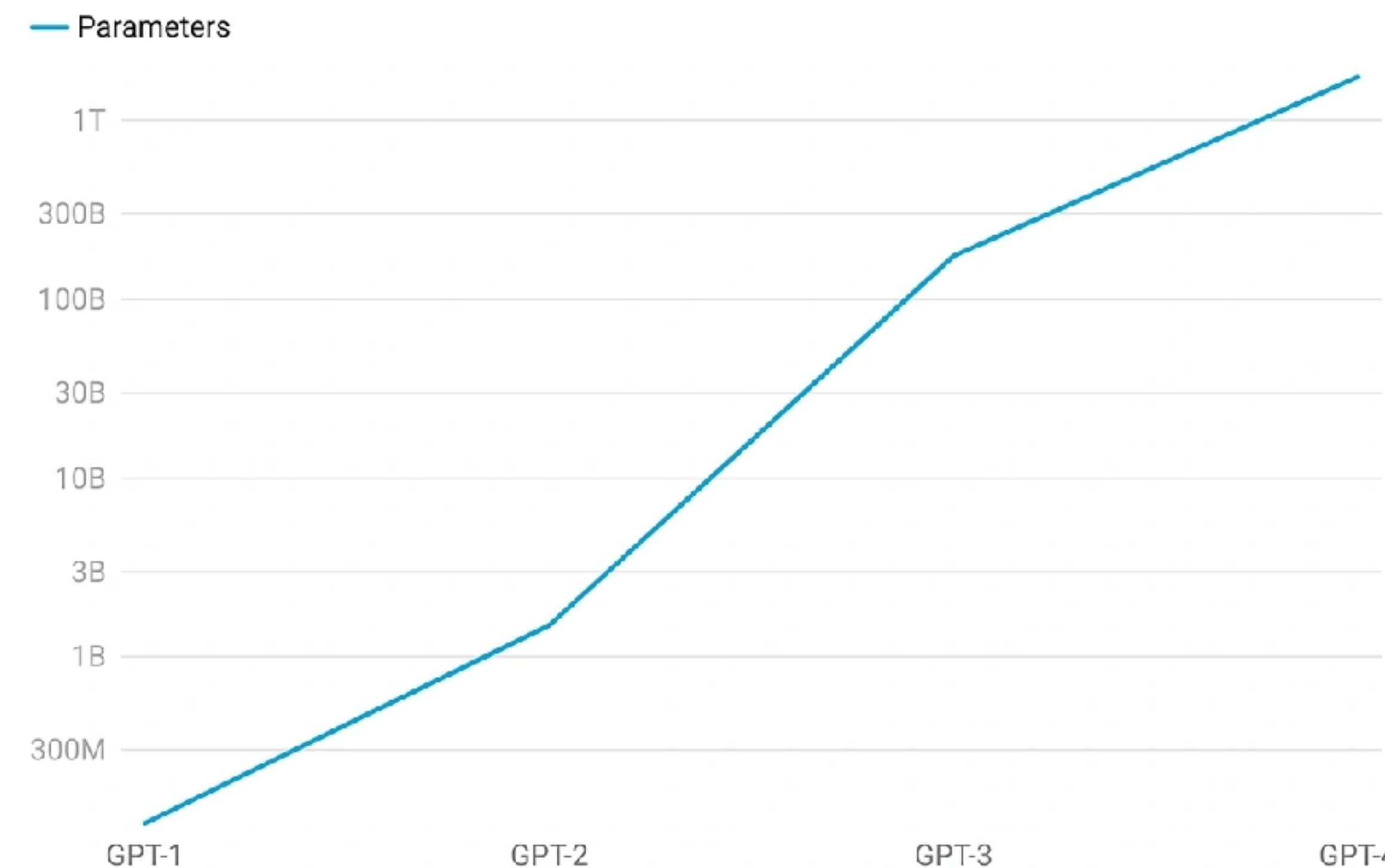


Motivation: Continual Growth in Model Size

- Scaling Law: More parameters, training samples, or compute time => More powerful models
- The model size has increased ~3000x in last 7 years and still increases

ChatGPT Parameters

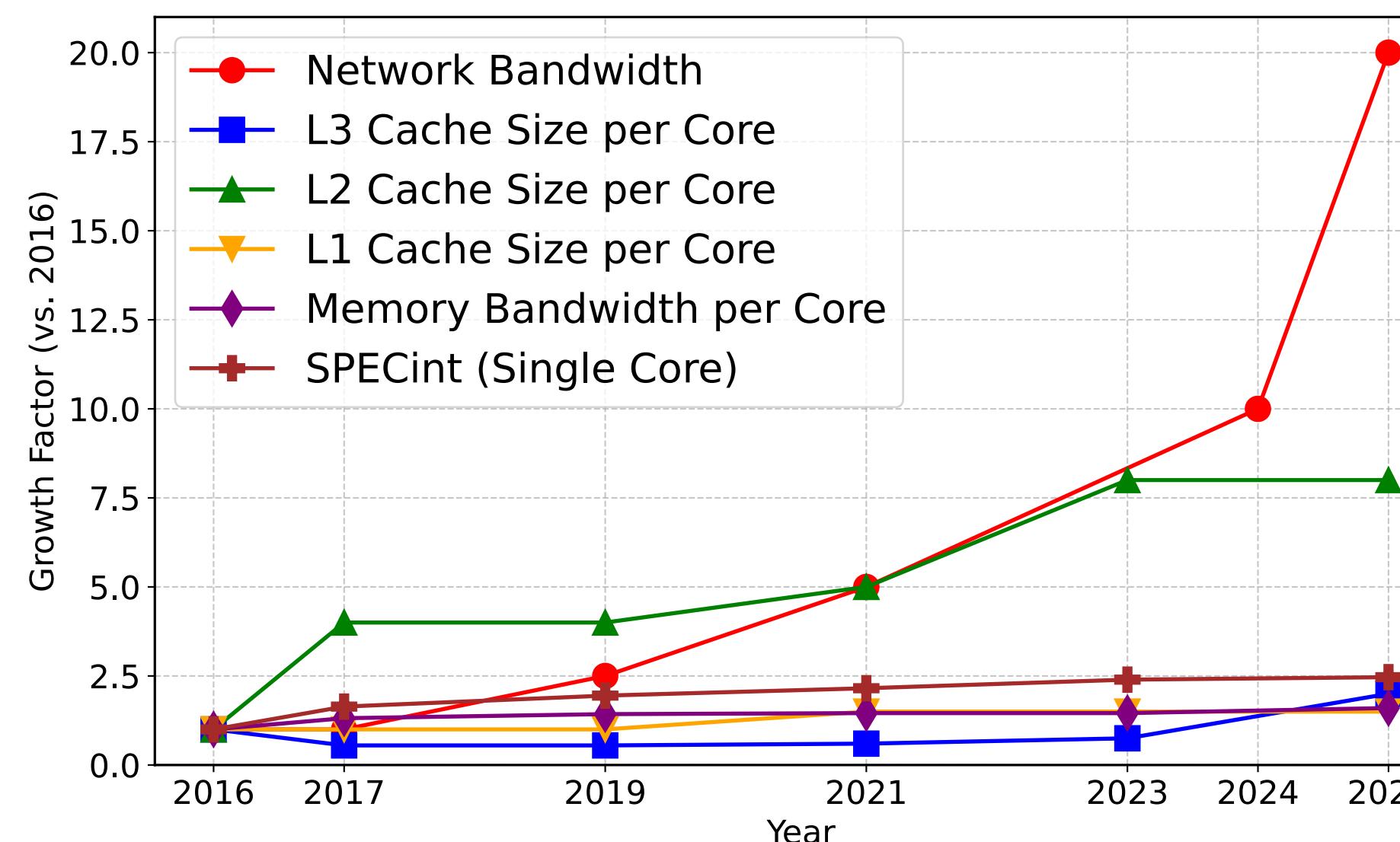
The number of parameters in successive models of ChatGPT has increased massively



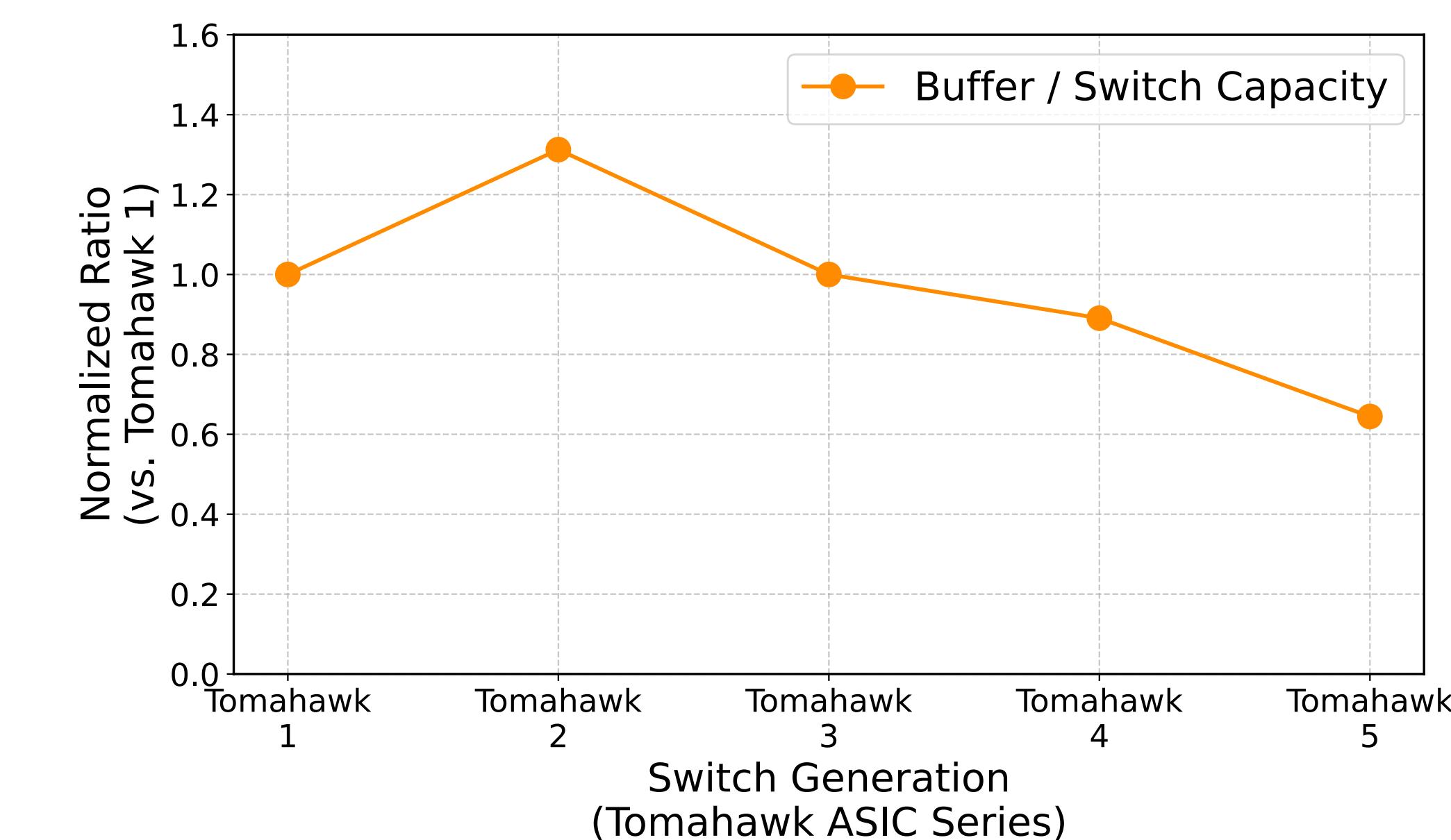
Motivation: Stagnant HW trends within Datacenters

- Datacenters play a central role in hosting LLMs.
- However, resource trends struggles to keep pace with the growing model sizes.

Host HW resources

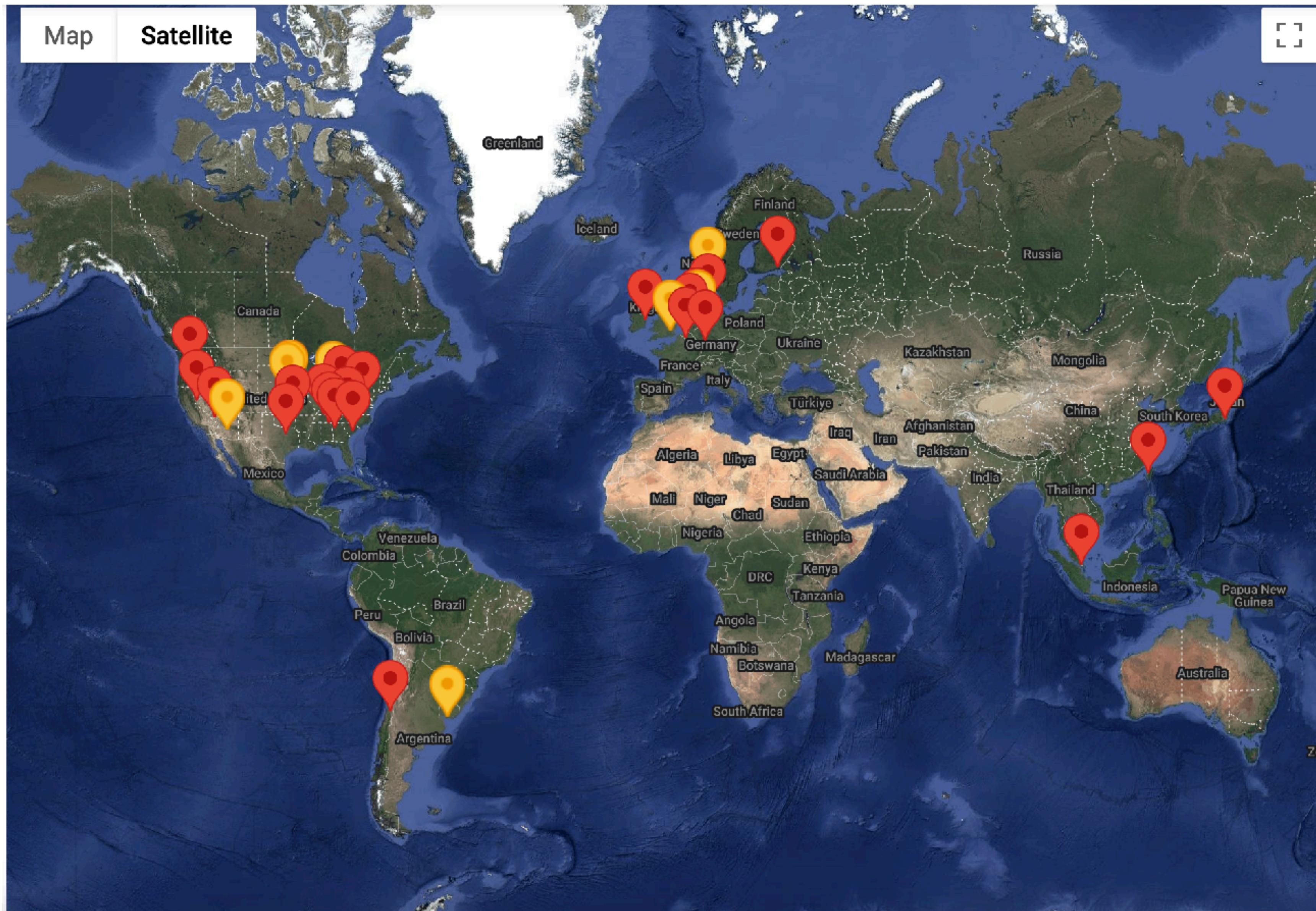


Network switch buffer sizes



Efficient utilization of these hardware resources is critical for achieving high performance in LLMs.

Locations of Google's datacenters



Inside a datacenter



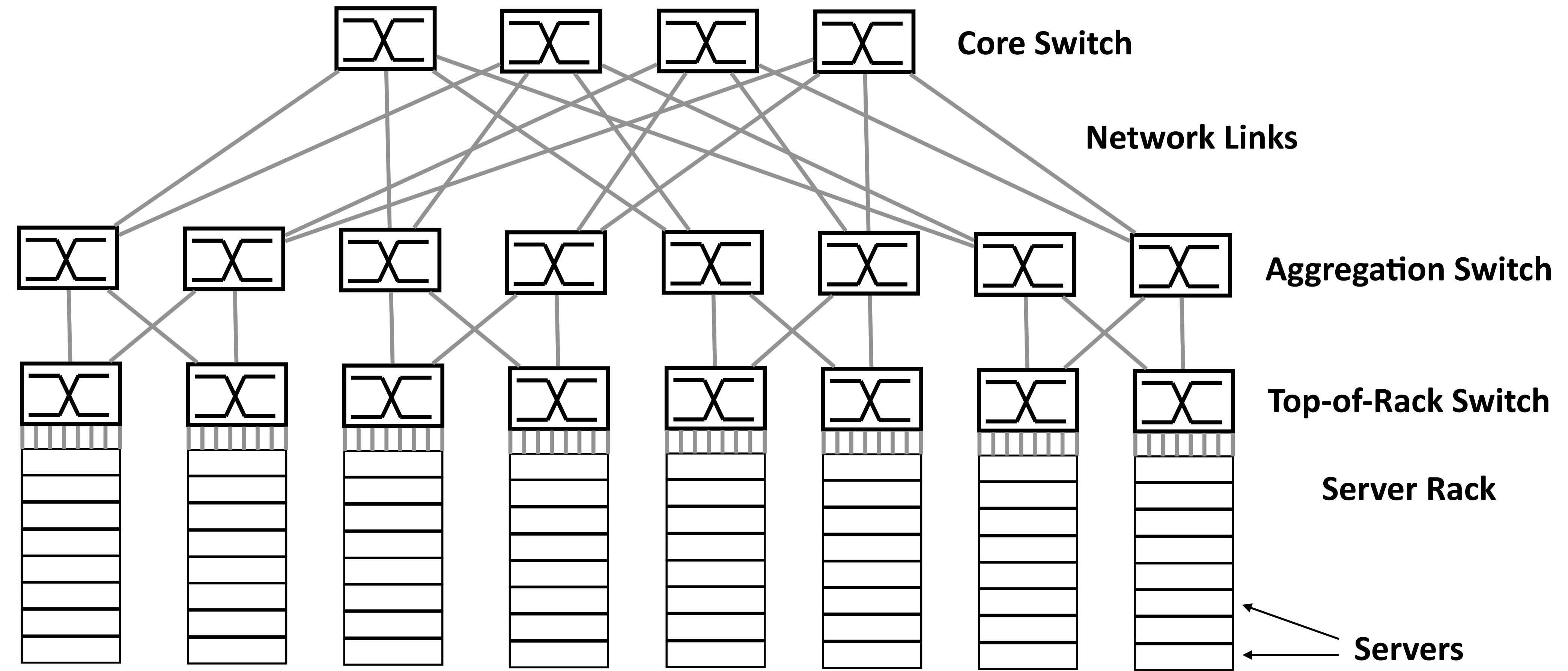
Inside a datacenter



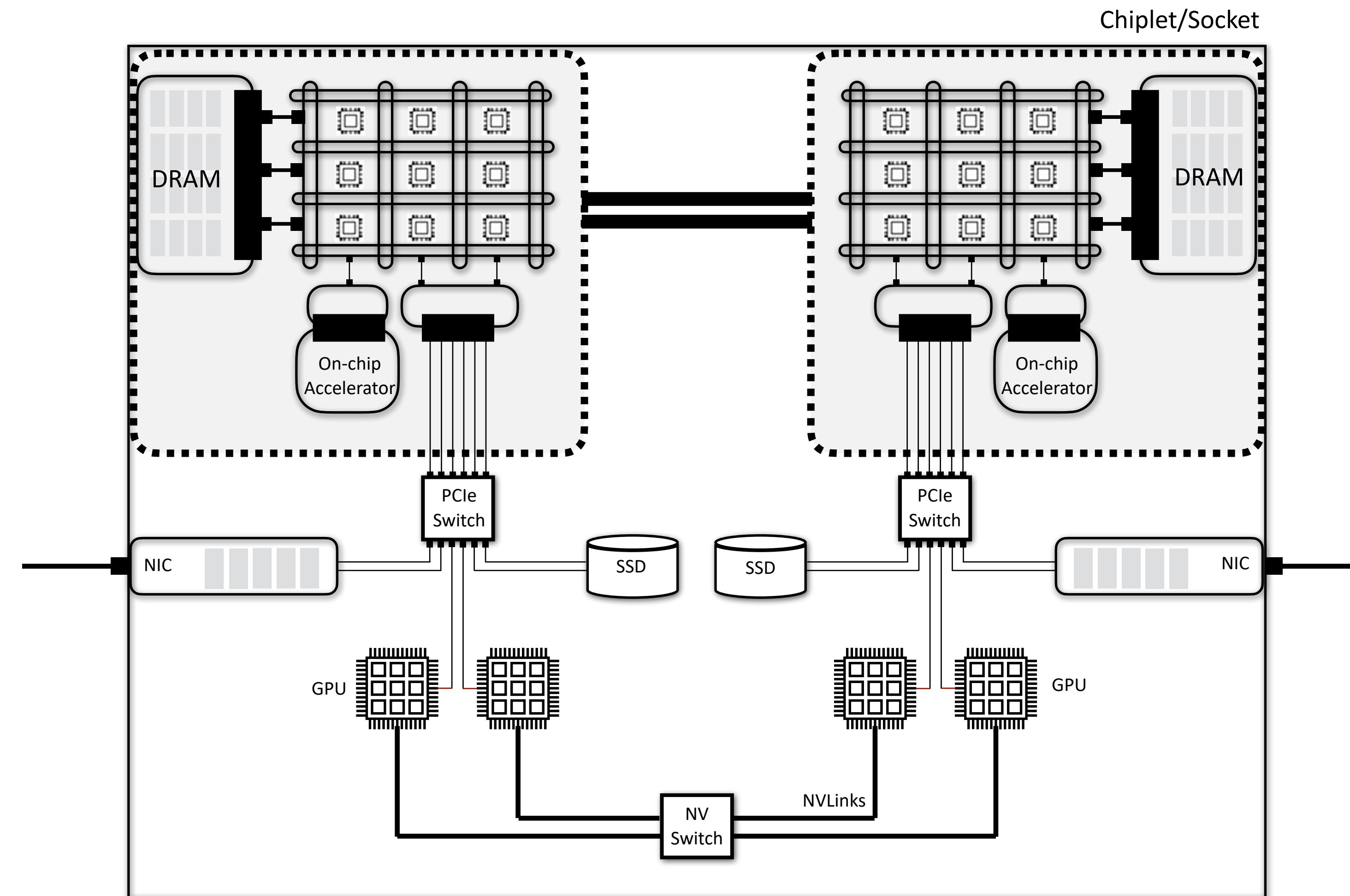
Inside a datacenter



Datacenter network architecture



Datacenter server architecture



How are datacenter resources shared/managed?

How are datacenter resources shared/managed?

- Network resources — switch buffers and link bandwidth
 - (Anything between the sender and receiver NICs)
- The mechanisms used to manage network resources
 - Traffic engineering
 - Load balancing
 - Traffic Shaping
 - Traffic Prioritization
 - Transport protocols

How are datacenter resources shared/managed?

- Server/host resources
 - Compute
 - Memory (DRAM and caching hierarchy)
 - Peripherals (disks or SSDS, GPUs, FPGAs)
- **Various components within the operating system (OS) used to allocate host resources**
 - CPU schedulers
 - Memory management
 - Network stack
 - Storage stack
 - ...

How are datacenter resources shared/managed?

- Additional (often overlooked) server/host resources
 - Memory interconnect
 - Peripheral interconnect
 - Processor interconnect
- **Hardware components/protocols** allocate these host resources
 - Memory controller and DDR
 - DMA engine, Root Complex, and PCIe/NVLink
 - Caching agents, home agents, ...

In this course, we will use the term “**DC infrastructure**” to refer to the
network/host resources + protocols/stacks used to manage these resources

Designing an Efficient DC Infrastructure is a Challenging Problem

- **Plethora of protocols, stacks and hardware**
 - Often designed/developed independent from each other
 - But intricate interactions can significantly impact end-to-end performance
- **Increasingly heterogeneous and complex hardware**
 - Different hardware resources often have different technology trends
 - In terms of capacity, performance (bandwidth/latency), energy efficiency, cost, etc
 - Bottlenecks can keep shifting with time;
 - Protocols/stacks need to either be resilient, or adapt to changing hardware/trends
- **Application workloads and user demands keep evolving over time**
 - Protocol/stacks often “optimized” to better serve the specific “average case” workloads

Topics covered in this course

- Four modules:
- Datacenter networking
- Host hardware & interconnects
- OS layers
- ML systems

Topics covered in this course: Datacenter Networking

- Datacenter topology
- Transport design
- Load balancing
- Networking infrastructure for ML
- Network communication for ML
- Optical networks

Topics covered in this course: Host Hardware & Interconnects

- PCIe
- CXL
- Silicon Photonics
- Host network
- FPGA-based computing/networking

Topics covered in this course: OS

- Network stacks
- Storage stacks
- CPU schedulers
- Memory management
- Memory protection
- Virtualization

Topics covered in this course: ML Systems

- Systems for inference
- Systems for training

Topics not covered in this course

- Course barely scratches the surface in discussion on cloud infrastructure
- **Goal:** preliminary insights into reasoning about end-to-end performance of datacenter applications
- Course does not talk about many more and important topics
 - Security
 - Power/energy efficiency
 - Telemetry
 - Monitoring
 - Debuging
 - Verification
 - ...

Questions?

Course Logistics

Course website: <https://www.qizhecai.com/cs6501-fall25/>

Classwork

- Paper reviews (15%)
- Paper presentation (15%)
- Class participation (20%)
- Research project or Survey (50%)

Paper Reviews

- Read and submit reviews for 2 papers for each lecture (starting from 09/03)
 - Required readings in course schedule
 - Reviews must be short (each with less than 2 * 200 words) and constructive
 - Suggested review outline
 - **Problem:** What is the **problem** being solved?
 - **Motivation:** Why is it interesting or important?
 - **Key ideas:** What are the key technical **insights** of the solution?
 - **Limitations:** What are few potential **limitations** of the current solution?
 - **Next steps:** What are few potential **next problems** to solve in this space?
- **Deadline to submit reviews:** 10am EST on the day of the lecture
- Check course logistics page for link to submit reviews
- You may skip up to **five classes** of reviews without penalty.

Paper Presentation

- Students will take turn giving 30 mins presentation on one paper each
- Presentation should cover the relevant related work for the presented paper/topic
 - Suggest taking a look at recommended readings, in addition to required readings
- **If you present a paper in class, you do not submit the review for that class, and it does not count as a missed review.**
- Suggested outline for presentation (no longer than 20 slides)
 - What is the **problem** being solved? (1-2 slides)
 - Why is it an **interesting** problem? (1-2 slides)
 - What is the existing **solution space** (related work)? (3-4 slides)
 - What are the key technical **insights** of the solution? (3-4 slides)
 - What are the **techniques** used to solve the problem? (3-4 slides)
 - What are a few potential **limitations** of the current solution? (1-2 slides)
 - What are a few potential **next problems** to be solved in this direction? (1-2 slides)
- **Deadline to submit slides: 10am Friday of the week before your presentation**
 - Exception: Presentations for 09/03 may be submitted at the last minute.

Class Participation

- We will all come to class prepared, having read the papers that will be presented
- Everyone is expected to actively participate in discussions
 - Discussions should ideally also incorporate the end-to-end picture
 - Based on the concepts accumulated over the semester

Research Project or Survey

- Two tracks — research project or survey
 - Any topic relevant to the course
- Research track: Motivate and solve a new research problem
- Survey track: Explore existing solution space and open questions for any existing research problem
- **Checkpoints** (see course website for more details about what to include in each report)
 - Proposal: due 10/15
- Final report: due 12/8
- Project presentations (optional)
 - For those who wish to present and get feedback from the class on their project

Questions?