

EECS 598: Systems for GenAI

Mosharaf Chowdhury



Today's Agenda

- Administrivia
- Topics
- Projects

About Mosharaf

- Associate Professor of CSE
 - <http://www.mosharaf.com/>
 - <https://symbioticlab.org/>
- Office hours:
 - Appointment-only

ViNEYard

virtual network embedding (2008-2012)



in-memory computing (2009-2014)

Coflow

data-parallel communication (2010-2016)

Infiniswap

software memory disaggregation (2016-2022)

Salus

GPU resource management (2017-2022)



systems for federated learning (2019-)



AI energy optimization (2021-)



About Jiachen Liu (GSI)

- 4th-year PhD student at SymbioticLab
- 2023 ML & Systems Rising Star
- Office hours from next week
 - 4828 BBB, 1230PM-130PM Fridays
 - No office hours this week
 - amberljc@umich.edu



Status

- As of today: ~50 registered or w/ override
- If you are not planning to take the class, drop ASAP
 - Existing overrides that have not converted will be revoked
 - Find me after class if you want override

Course Schedule

- Webpage: <https://github.com/mosharaf/eecs598>
- Meetings
 - 10:30 AM – 12PM (**T/Th** for lectures and seminars)
 - 1:30 PM – 2:30 PM (**Fri** for makeups and projects)
- **Pay attention to the online announcements and schedule**
 - On average, two meetings per week
 - Friday makeups will be added on a need-to-add basis

Prerequisites

- **EECS 482 / 484 / 489 / 491**
 - Equivalent courses are acceptable as well
- **Good programming skills**
 - Build substantial systems for course project

Course Requirements

Paper Summary	15%
Paper Presentation	15%
Participation	10%
Project Report	40%
Project Presentations	20%

Topics (#Lectures)

- GenAI Basics (4)
- Pre-Training (2)
- Fine-Tuning/Alignment (2)
- Inference (2)
- Grounding (2)
- Systems Optimizations (3)
- Special Topics (3)

Group-Based Work

- **ALL activities will be done in groups except for participation**
 - Paper presentation
 - Paper summary
 - Research projects

Form Groups ASAP

- Submit at <https://forms.gle/t8n6V9ewJoDWTaSL9>
 - By January 23 the latest, but **right now** is better
 - We need a group to pickup duties for Jan 18!!!
 - Use piazza to find group members
 - Group size should be 3
 - May allow a few smaller groups if/when students drop off

Readings

- **38 papers/articles across**
 - Primarily from systems venues like SOSP, OSDI, NSDI, EuroSys, and MLSys
 - Some from traditional AI/ML venues but still with systems-y flavor

Paper Presentation

- **This is a seminar-style course**
 - Each group must present at least one lecture (required papers and the rest)
 - Paper presentation account for **15%** of the total grade
- **The entire class will be dedicated to the assigned paper(s)**
 - Aim for 40-minute presentation without interruption
 - But there will be intermittent discussions
- **Lead the discussion**
 - Go through the paper in details, along with its strengths and weaknesses
 - Include companion papers and other related papers

Paper Presentation

- Share your slides to eeecs598-genai-staff@umich.edu 24 hours before the class
 - Use Google Sheets so we can provide in-place comments/feedback
- Prepare early
- Practice a lot
- Also, read
 - [How to Give a Bad Talk](#), by David A. Patterson

Paper Summaries

- **This is a paper-reading course**
 - Paper summaries account for **15%** of the total grade
- **Roughly 1-2 summary per-group (assigned)**
- **Each summary must follow the template and address the following**
 - What is the problem and why is it important?
 - What is the hypothesis of the work?
 - What is the proposed solution, and what key insight guides their solution?
 - What is one (or more) drawback or limitation of the proposal, and how will you improve it?
- **Summary must include the gist of class discussion**

Paper Summaries

- Reviews must be shared to eeecs598-genai-staff@umich.edu within 24 hours of class presentation
 - Use Google docs so we can provide in-place comments/feedback
- **Delayed submission will receive NO CREDIT**
 - There will be NO extensions

Panel Discussion

- **The Authors**
 - Groups that present and write summary
- **The Reviewers**
 - Each group will be assigned to at least one of these slots
- **Rest of the Class**
 - [Submit](#) one insightful question for each presented papers before each class
 - Ask questions directly too
 - [Vote](#) for your favorite group based on their performance

In general,

- No extensions
- Everyone must come to class after reading the **required** papers of the day

What Do We Talk About When We Talk About “**Systems** for GenAI”

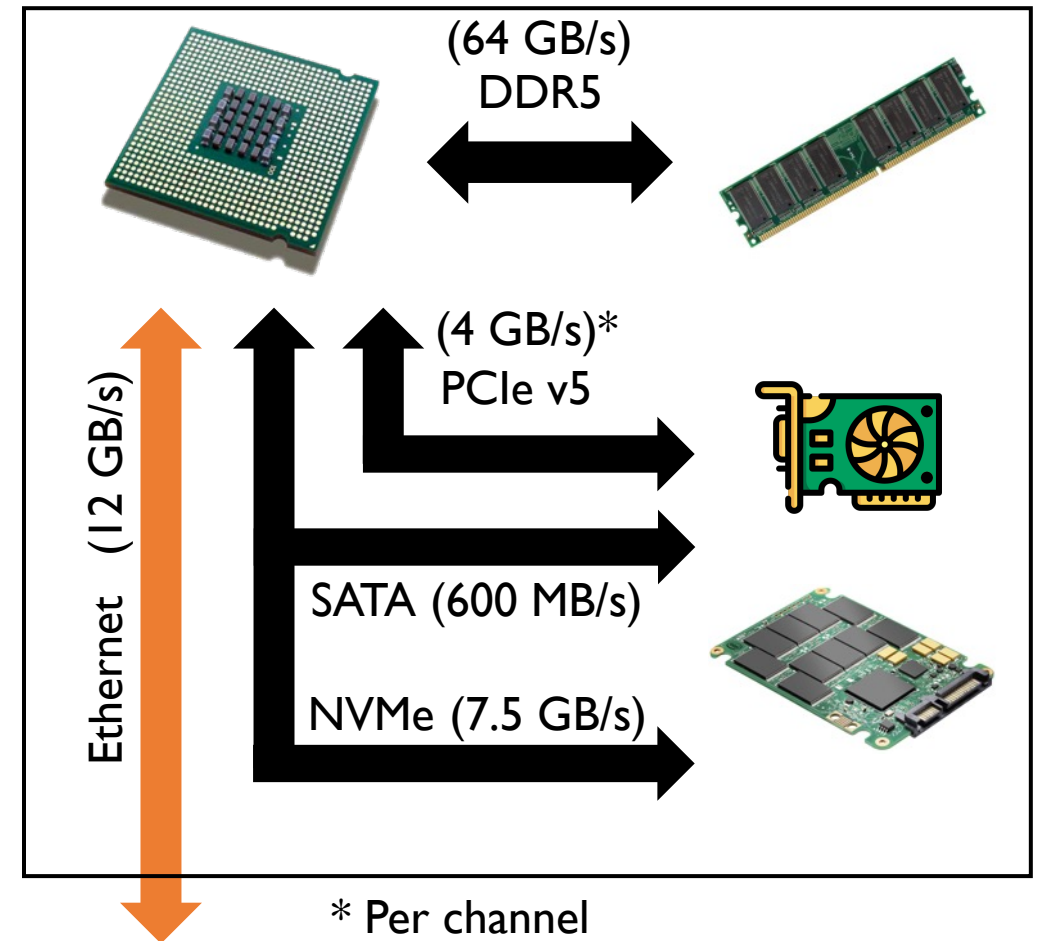
What's in a (Simplified) Server?

Interconnected compute and storage resources

- Different bandwidth and latency constraints

Simplified diagram

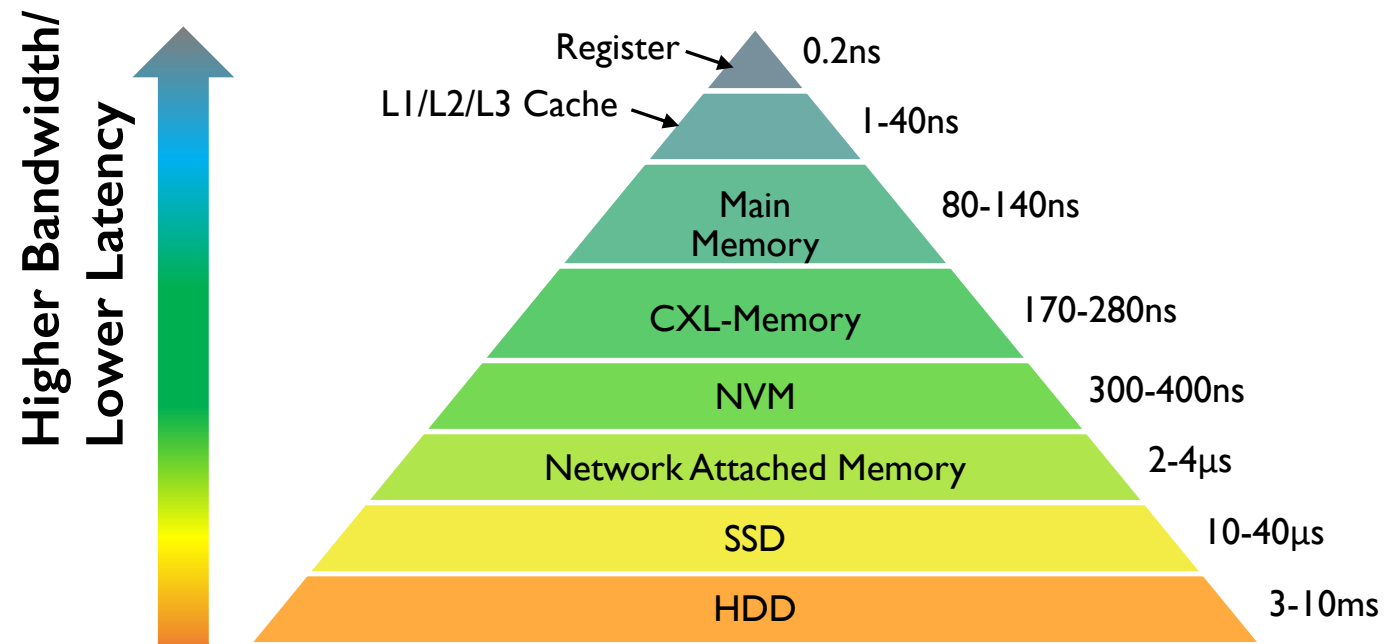
- Doesn't include faster networks such as RDMA, dedicated GPU interconnects such as NVlink, etc...



Typical Memory/Storage Hierarchy

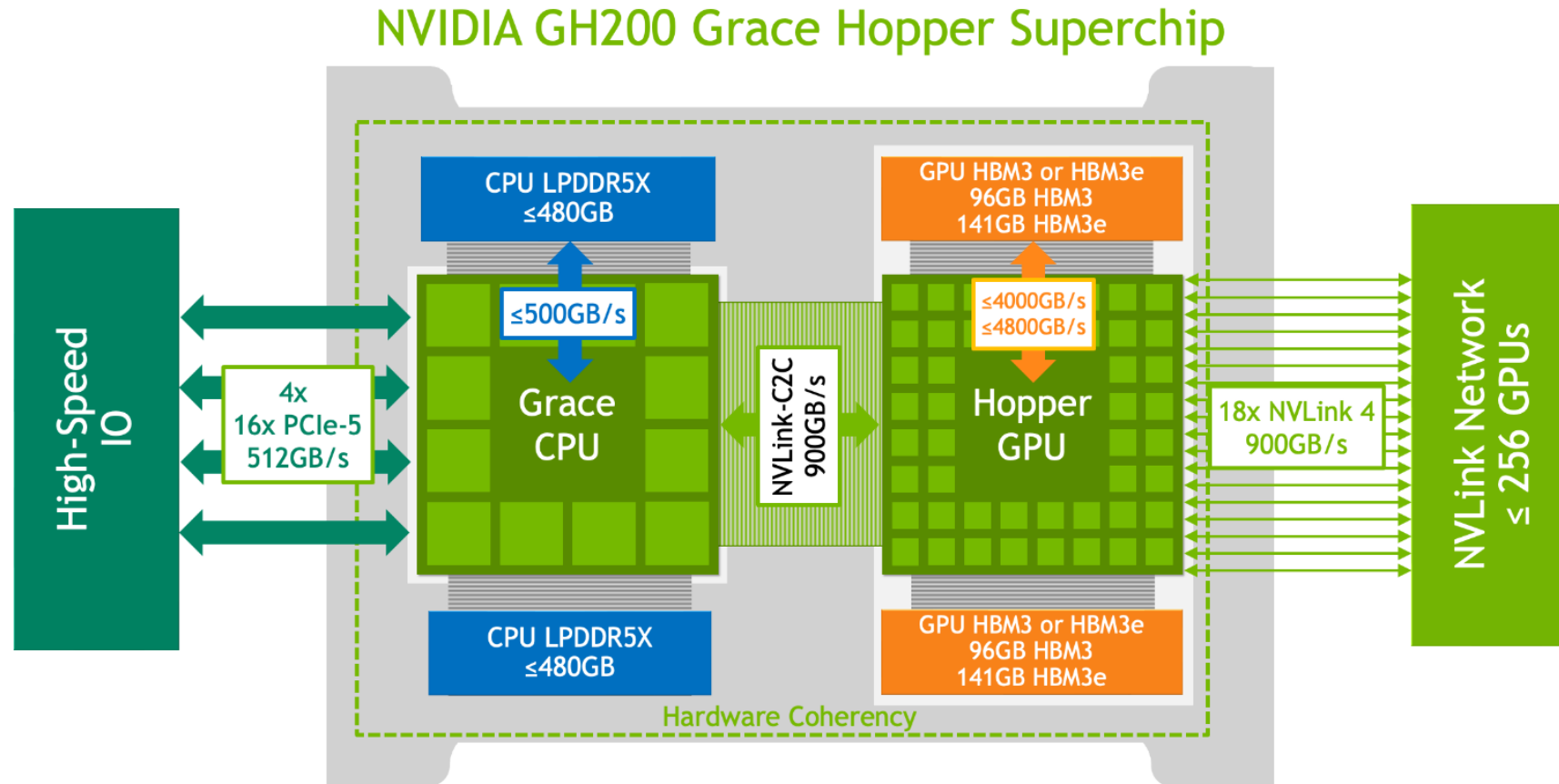
Fundamental Goals of (SW/HW) System Design

- Minimize time to access data
- Maximize compute utilization
- **Balanced System**



Maruf et al, SIGMETRICS 2023

What's in a Modern AI Server?



<https://resources.nvidia.com/en-us-grace-cpu/nvidia-grace-hopper>

Scale Out: Warehouse-Scale Computer (WSC)

Single organization

Homogeneity (to some extent)

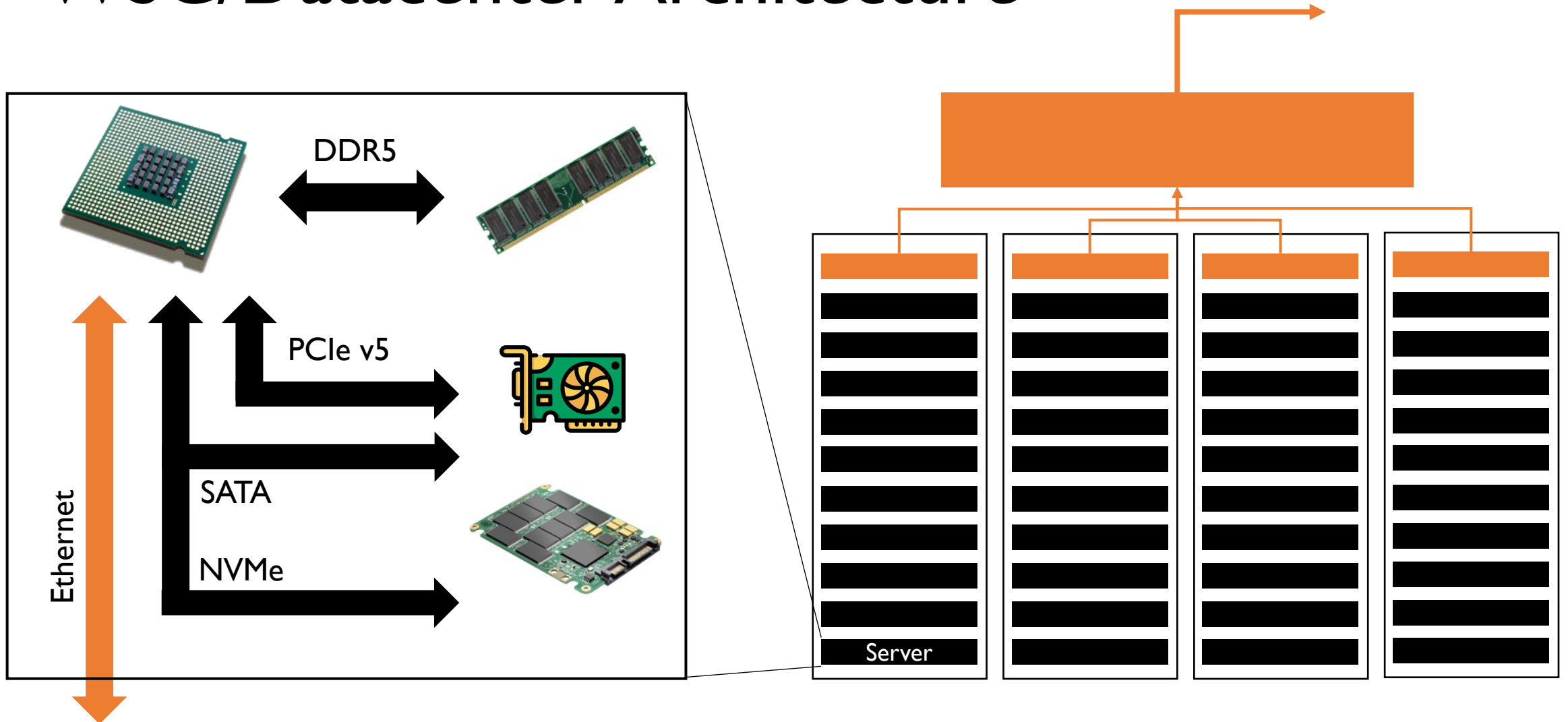
Cost efficiency at scale

- Multiplexing across applications and services
- Rent it out!

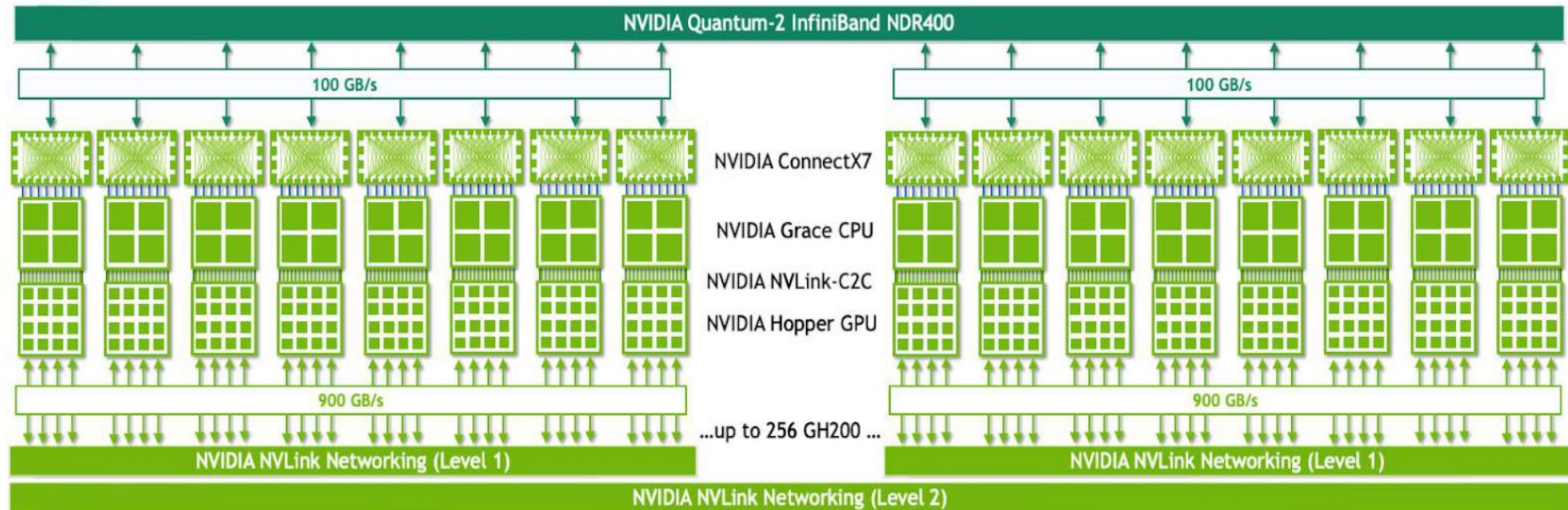
Many concerns

- Infrastructure
- Networking
- Storage
- Software
- Power/Energy
- Failure/Recovery
- ...

WSC/Datacenter Architecture



Example: Scaling Out Using NVIDIA GH200



<https://resources.nvidia.com/en-us-grace-cpu/nvidia-grace-hopper>

Datacenter Needs an Operating System

Datacenter is a collection of

- Compute
- Memory
- All connected by an interconnect

Not unlike a computer

Some differences

1. VERY high level of parallelism
2. VERY large scale
3. Diversity of workloads
4. Resource heterogeneity
5. Failure is the norm

Three Categories of Software

1. Platform-level

- Software firmware that are present in every machine

2. Cluster-level

- Distributed systems to enable everything

3. Application-level

- User-facing applications built on top

Common “Systems” Techniques

Technique	Performance/Efficiency	Availability/Resilience
Replication & Erasure coding	X	X
Sharding/partitioning	X	X
Scheduling & Load balancing	X	
Health & Integrity checks		X
Compression & Quantization	X	
Centralized controller	X	
Canaries		X
Speculation & Redundant execution	X	

Break!

Systems for GenAI Projects

Research-Oriented Course!

- The final project accounts for **60%** of total grades
- What can and cannot be a project?
 - Just surveys are not allowed
 - Measurements of new environments or of existing solutions on new environments are acceptable
 - Reproducing results from existing solutions is also acceptable
- An ideal project should answer the questions you asked during paper reviews and points you cared about for presentations

How to Approach it?

1. Find a problem and motivate why this is worth solving
2. Quickly survey background and related work
 - Might require you to go back to the first step
3. Form/update your hypothesis
4. Test your hypothesis
 - Go back to 3 until you are happy
5. Present your findings on poster and in writing
 - Discuss known limitations

Milestones

Date	Milestone	Details
01/23/24	Form Group	Find 3 like-minded students
02/09/24	Submit Proposal	Send your proposal by email to receive feedback either via email or in-person or both
03/12/24 03/14/24	Mid-Semester Presentations	Define and motivate a problem, overview related work, and form initial hypothesis and idea
04/16/24 04/18/24	In-Class or Poster Presentations	Present your findings
05/01/24	Research paper	Submit a report like the papers you read

Draft Proposal (Feb 9)

- **Two pages including references that **must** include**
 - What is the problem?
 - Why is it important to solve?
 - Any initial thoughts on what you want to do?
 - How would you evaluate your solution?
- **Include team members**
 - Meaning, form a group ASAP
- **Approved by the instructor and agreed upon by you**
 - Forms the basis of expectation

Mid-Semester Checkpoint (Mar 12,14)

- **In-class short presentation over two days**
 - This is to make sure you are making progress
- **Must include**
 - What is the problem?
 - Why is it important?
 - What are the most related work?
 - What's your hypothesis so far?
 - How are/will you evaluate it?

Presentation & Paper (Apr 16, 28, May 1)

- **Research paper**
 - The key part
 - Should be written like the papers you've read
 - As if you'd submit it to a workshop with ~3 more months of work or to a conference after ~6 more months of work
 - [How to Write a Great Research Paper](#) by Simon Peyton Jones
- **Extended from the mid-semester checkpoint writeup**

Project Ideas

Some project suggestions

- <https://docs.google.com/document/d/1GcYee3HUqjhSkfLkQX7enickHQ00yZFI71Vvn9Kbb-M/edit>

You can propose your own projects too!

Next Class...

Read the required readings

Form groups of 3 and fill out <https://forms.gle/t8n6V9ewJoDWTaSL9> by *Jan 23*

- Decide if you'll drop, **before** you fill it
- If you are to drop, drop immediately

Sign up for Jan 18 slot for extra benefits!

