



# schedule

## presentation schedule

---

although this is a seminar, the first few weeks will include cross-cutting intro/thematic talks to help set the stage for the semester. because we expect students to have a specialization in a sub-area already, the intro lectures will mostly focus on broad ideas, establish common ground, and point to resources. students can use these independently to fill in the gaps based on their own interests, perhaps in consultation with instructors/their advisors.

the (student) seminar presentation schedule will comprise a subset of papers from the list of papers below, determined after students are matched with their topics + assigned seminar dates. students will present a seminar-style lecture on their assigned paper(s), and lead a class discussion following it.

please register your preferences for topics using the form on slack by 11:59pm on 2/23. we'll use this to match students with papers/topics. if you do not submit a preference, you will be assigned a topic based on availability. we can't guarantee that we will be able to accommodate all preferences, but we will do our best to match students with topics that are of interest to them. we reserve the right to make changes to the schedule and seminar assignments based on student preferences and availability.

the presentation schedule will be maintained on under the "resources" tab on slack.

## a tentative list of topics + readings

---

this is an evolving list of papers that we plan to cover, roughly covering a week of material per "block" - mostly for students to read ahead and get a general idea of the topics, or dive deeper into specific topics. please keep an eye out for updates!

---

THEME	TOPICS	SELECTED PAPERS/READINGS/ SYSTEMS
<b>Analytical systems fundamentals</b>	DB architecture Access paths, storage, encodings	<b>Architecture of a Database System</b> (Hellerstein+, '07) <b>Red Book Ch. 4</b> (Stonebraker)
	Query planning & execution	<b>Access path selection in a relational database management system</b> (Selinger+, SIGMOD '79) <b>A Deep Dive into Common Open Formats for Analytical DBMSs</b> (Liu+, VLDB '23) <b>An Empirical Evaluation of Columnar Storage Formats</b> (Zeng+, VLDB '23)
	Indexing & search 1 (Exact search) Optimizers	<b>The Cascades Framework for Query Optimization</b> (Graefe+, IEEE Data Eng. Bull. '95) <b>Apache Calcite: A Foundational Framework for Optimized Query Processing Over Heterogeneous Data Sources</b> (Begoli+, SIGMOD '18)
<b>Modern systems</b>	Data warehouses Streaming engines Lakehouses	<b>The Snowflake Elastic Data Warehouse</b> (Dageville+, SIGMOD '16) <b>The Stratosphere platform for big data analytics</b> (Alexandrov+, SIGMOD '14) <b>Apache Flink</b> (various

THEME	TOPICS	SELECTED PAPERS/READINGS/ SYSTEMS
		<p>publications, documentation)</p> <p><b>Apache Spark: a unified engine for big data processing</b> (Zaharia+, SIGMOD '16 - article)</p> <p><b>Spark SQL: Relational Data Processing in Spark</b> (Ambrust+, SIGMOD '15)</p> <p><b>Delta Lake: High-Performance ACID Table Storage over Cloud Object Stores</b> (Ambrust+, VLDB '20)</p> <p><b>Lakehouse: A New Generation of Open Platforms that Unify Data Warehousing and Advanced Analytics</b> (Armburst+, CIDR '21)</p> <p><b>Photon: A Fast Query Engine for Lakehouse Systems</b> (Behm+, SIGMOD '22)</p>
	Big data systems	<p><b>Dremel: Interactive Analysis of Web-Scale Datasets</b> (Melnik+, VLDB '10)</p> <p><b>Dremel: A Decade of Interactive SQL Analysis at Web Scale</b> (Melnik+, VLDB '20)</p> <p><b>BigLake: BigQuery's Evolution toward a Multi-Cloud Lakehouse</b> (Levandoski+, SIGMOD '24)</p> <p><b>Spanner: Google's Globally-Distributed Database</b> (Corbett+, OSDI '12)</p>

THEME	TOPICS	SELECTED PAPERS/READINGS/ SYSTEMS
	<p>ML for systems</p> <p>Optimization techniques: bandits, mathematical programming (linear, integer, dynamic)</p> <p>RL (part 1): MDPs, Regret, Q error...</p> <p>Misc.: Heuristic/metaheuristic algorithms (MCTS + DRL, ...), weak supervision</p>	<p><b>The Case for Learned Index Structures</b> (Kraska+, SIGMOD '18)</p> <p><b>Neo: A Learned Query Optimizer</b> (Marcus+, VLDB '19)</p> <p><b>SageDB: A Learned Database System</b> (Kraska+, CIDR '19)*</p>
<b>Vectors, Embeddings</b>	<p>Vector search</p> <p>Embedding models</p> <p>Indexing &amp; search 2</p> <p>Index optimization: Quantization, parameter tuning, rebuilding indices, dim. reduction (PCA, t-SNE)</p> <p>Hybrid, Multivector, Cross-modal retrieval</p> <p>Retrieval Augmented Generation</p>	<p><b>Efficient Estimation of Word Representations in Vector Space</b> (Mikolov+, ICLR '13)</p> <p><b>Dense Passage Retrieval for Open-Domain Question Answering</b> (Karpukhin+, EMNLP '20)</p> <p><b>ColBERT: Efficient and Effective Passage Search via Contextualized Late Interaction over BERT</b> (Khattab+, SIGIR '20)</p>
<b>AI/ML systems</b>	<p>Computation and Numerics</p> <p>ML compilation overview</p> <p>Computational graphs</p> <p>Automatic differentiation</p> <p>GLO</p> <p>Kernel fusion</p>	<p><b>TASO: Optimizing Deep Learning Computation with Automatic Generation of Graph Substitutions</b> (Jia+, SOSP '19)</p>

THEME	TOPICS	SELECTED PAPERS/READINGS/ SYSTEMS
	Distributed ML Data parallelism (Parameter servers, All-reduce) Model, Pipeline parallelism	<b>GPipe: Efficient Training of Giant Neural Networks Using Pipeline Parallelism</b> (Huang+, NIPS '19) <b>Scaling Distributed Machine Learning with the Parameter Server</b> (Li+, OSDI '14) <b>Large Scale Distributed Deep Networks</b> (Dean+, NIPS '12)
	Hardware for AI Architecture-aware design GPU and systems architecture intro	<b>QLORA: Efficient Finetuning of Quantized LLMs</b> (Dettmers+, NIPS '23) <b>FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness</b> (Dao+, NIPS '22)
	Inference + serving at scale (1) Part 1:	<b>DeepSpeed Inference: Enabling Efficient Inference of Transformer Models at Unprecedented Scale</b> (MSR) <b>AlpaServe: Statistical Multiplexing with Model Parallelism for Deep Learning Serving</b> (Li+, OSDI '23) <b>Orca: A Distributed Serving System for Transformer-Based Generative Models</b>

THEME	TOPICS	SELECTED PAPERS/READINGS/ SYSTEMS
<b>Systems for Reasoning</b>	Reasoning frameworks and techniques: CoT, ReAct, Tree of Thoughts, Self-reflection, RL for reasoning (Q*, ...) Additional topics: Self and path consistency, Neuro-symbolic systems, Diffusion models	<p>(Yu+, OSDI '22) <b>Ray: A Distributed Framework for Emerging AI Applications</b></p> <p><b>Outrageously Large Neural Networks: The Sparsely-Gated Mixture-of-Experts Layer</b> (Shazeer+, ICLR '17) <b>Chain-of-Thought Prompting Elicits Reasoning in Large Language Models</b> (Wei+, NIPS '22) <b>ReAct: Synergizing Reasoning and Acting in Language Models</b> (Yao+, ICLR '23) <b>Tree of Thoughts: Deliberate Problem Solving with Large Language Models</b> (Yao+, NIPS '23) <b>Hierarchical Reasoning Model</b> (Sapient Intelligence, '25) <b>On the role of planning in model-based deep reinforcement learning</b> (Hamrick+, '21)</p>
<b>Compound AI systems</b>	Prompt optimization frameworks: GRPO, MIPRO, GEPA Orchestrated systems: Deep research, ...	<p><b>The Shift from Models to Compound AI Systems</b> (BAIR) <b>AlphaEvolve: A coding agent for scientific and algorithmic discovery</b> (Google DeepMind)</p>

THEME	TOPICS	SELECTED PAPERS/READINGS/ SYSTEMS
		<p><b>DSPy: Compiling Declarative Language Model Calls into Self-Improving Pipelines</b> (Khattab+, ICLR '24)</p> <p><b>Automated Design of Agentic Systems</b> (Hu+, ICLR '25)</p>
	<p>Inference + serving at scale (2)</p>	<p><b>vLLM: Efficient Memory Management for Large Language Model Serving with PagedAttention</b></p> <p><b>CacheGen: KV Cache Compression and Streaming for Fast Large Language Model Serving</b></p> <p><b>CacheBlend: Fast Large Language Model Serving for RAG with Cached Knowledge Fusion</b></p> <p><b>Autellix: An Efficient Serving Engine for LLM Agents as General Programs</b> (Luo+, '25)</p>
<p><b>AI-native analytical systems</b></p>	<p><b>CHASE-SQL: Multi-Path Reasoning and Preference Optimized Candidate Selection in Text-to-SQL</b></p> <p><b>DocETL: Agentic Query Rewriting and Evaluation for Complex Document Processing</b> (Shankar+, '24)</p> <p><b>LOTUS: Enabling Semantic Queries with LLMs Over Tables of Unstructured and Structured Data</b> (Patel+, '24)</p> <p><b>A Declarative System for Optimizing AI Workloads</b> (Liu+, CIDR</p>	

THEME	TOPICS	SELECTED PAPERS/READINGS/ SYSTEMS
	'24)	

## important dates

---

- project proposals: Mar 6
- project milestone: Apr 10
- project submission: Apr 29
- presentation preferences: Feb 23
- short summaries: before the start of class (1:45p) on the day of the seminar



# policies

## grading

---

COMPONENT	WEIGHT
short summaries	10%
presentation + seminar	20%
seminar report	15%
project <ul style="list-style-type: none"><li>• proposal: 5%</li><li>• milestone: 20%</li><li>• presentation: 5%</li><li>• final submission: 70%</li></ul>	35%
attendance and participation	20%

## attendance and participation

---

all classes will be held in person barring extreme circumstances. since seminar courses are centered around discussion, the course will not be recorded so that students do not need to fear their questions or comments being immortalized.

as a result, (1) participation and discussion is highly encouraged, and we suggest you

bring a pen(cil) and paper to take notes and leave your electronics in your bag (still bring them – some lectures will have online interactive components), although this is not mandatory and you are welcome to take electronic notes if you prefer. (2) there is essentially no way to "make ud class, so i will expect students miss no more than 4 class periods throughout the semester.

*policy adopted from ryan's offering of cis6500.*

## short summaries

---

when it is *not* your turn to present, you will submit a ~500 word summary that thematically connects and discusses the papers chosen for the seminar. summaries are graded on completion. consider only selected "main" papers for the summary (marked with an asterisk \* on the presentation schedule on slack).

your summary is due before the start of class (1:45p) on the day of the seminar. late summaries will receive no credit. *you can turn in your summaries for the first week of the seminar until 11:59pm on 2/23.*

the submission form can be found here.

## research, academic integrity and ai policies

---

due to the unique positioning of this class, which places it closely with publishable research undertakings, we take academic integrity very seriously. you are responsible for knowing penn's code of academic integrity. we will aggressively pursue all cases of suspected violations.

unless otherwise noted, all assignments are individual work.

**ai tools:** judicious ai tool use is acceptable and encouraged for learning concepts and literature review in assignments. ideally, ai tools should not be used as a substitute for thinking or completing assignments end-to-end. you can also try using it to automate boilerplate or simple code that is low stakes.

however, for any work you turn in or present towards a grade, we will assume a yardstick where you are the sole creator of your work. that is, you should both understand this

work thoroughly and be fully responsible for what you turn in, since it is "your work".

more concretely, you're expected to be able to explain all important details of the work you turn in, except for only very minor boilerplate/setup parts. this ensures that you can engage in discussions and contribute to projects effectively, justifying and iterating on work. this also prevents you from potentially misleading your peers, advisors, and other stakeholders (even if unconsciously) about your progress, work, and decision-making.

moreover, this also means that should your work demonstrate signs of violations resulting from inadvertent careless ai tool use, you (and not the ai tool) would be held solely responsible. for example, if an experimental result is found to be fabricated, a claim of an ai tool implementing the code will still land you in grave danger.

finally, certain assignments may contain a section that requires you to disclose any ai tool use. you must answer this honestly and transparently.

you should also save your conversations with ai tools for future reference.