# Exploring Next Token Prediction For Optimizing Databases

**Yeasir Rayhan** and  Walid G. Aref

Purdue University

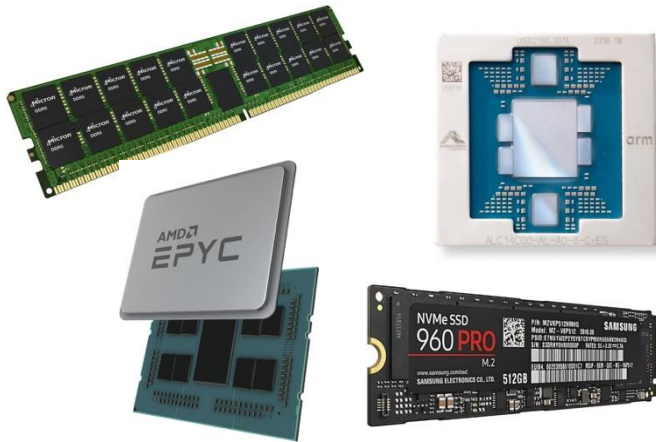West Lafayette, Indiana, USA

PURDUE
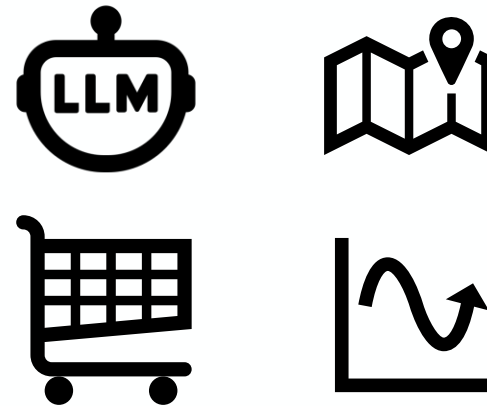UNIVERSITY.

# Roadmap

- The current landscape of Database Systems design space
  - Hardware and Workload
- The LLM recipe
  - Next Token Prediction (NTP): First step towards adopting the LLM recipe
- The building blocks to adopt Next Token Prediction (NTP) in databases
  - Decision Transformers
  - DB-Tokens
- The Probe and Learn (PoLe) framework
  - Preliminary case study
    - Index Scheduling

# The Game Changers in Databases [1]

Hardware

Workload

[1] Anastasia Ailamaki. 2021. Accelerated Data Management Systems Through Real-Time Specialization.
Keynote at MICRO.

# The current state

Hardware                                    Workload

What's the current state?
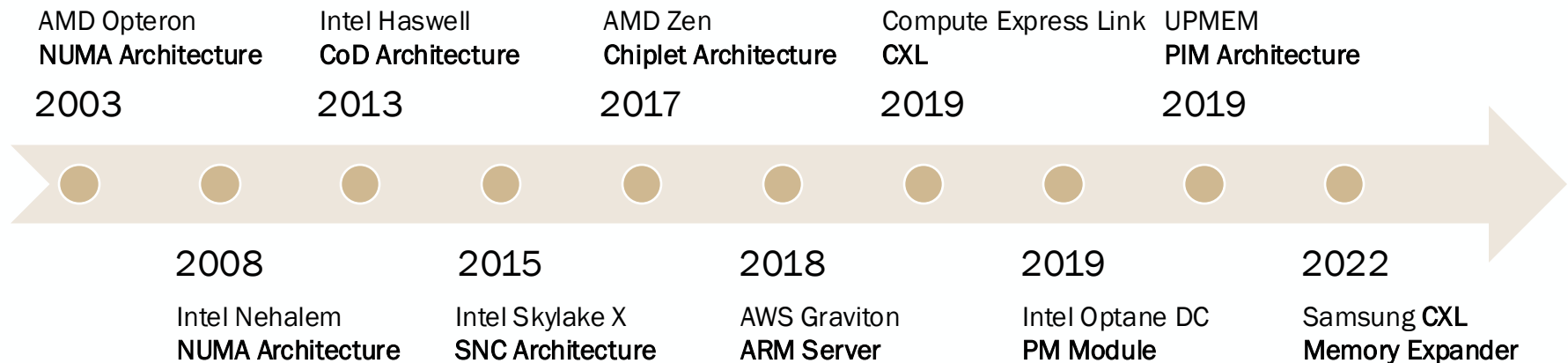
# The Current State: Observation 1

Hardware

Workload

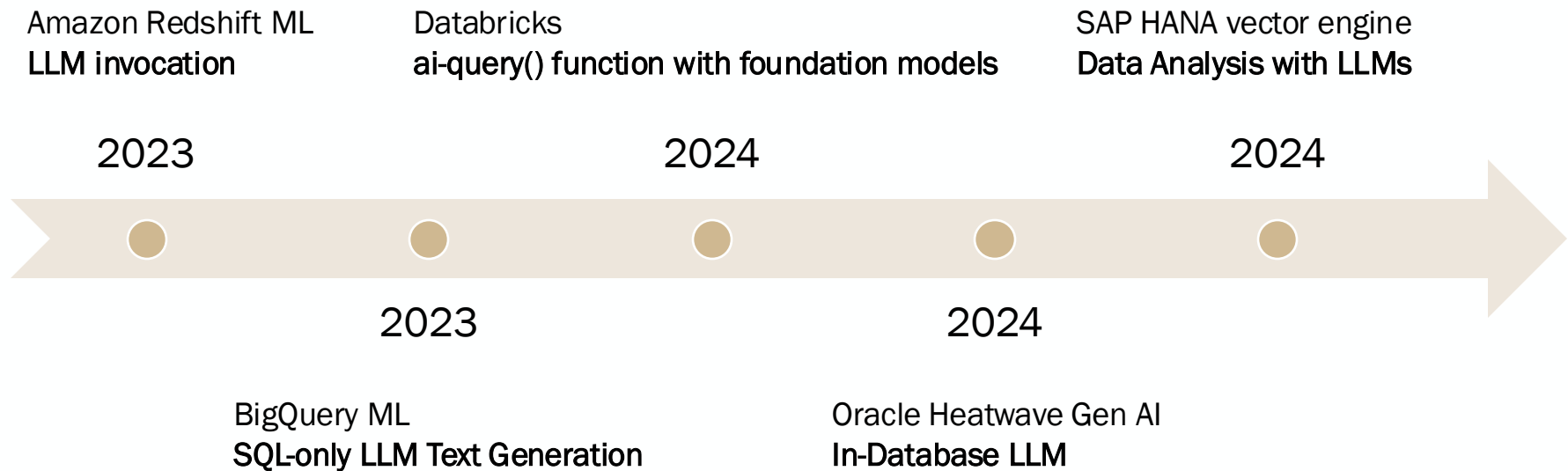1. Both the hardware and workload are **rapidly evolving.**

# Observation 1: Rapid Evolution of Hardware

- Every few years, there is a new technology:
  - Compute, Memory and Storage
- Within the past 6 years, we have seen the commercial introduction of
  - Processing In Memory (PIM) chips by UPMEM
  - Compute Express Links (CXL)
  - ARM-based server processors

AMD Opteron
NUMA Architecture
2003

Intel Haswell
CoD Architecture
2013

AMD Zen
Chiplet Architecture
2017

Compute Express Link
CXL
2019

UPMEM
PIM Architecture
2019

2008
Intel Nehalem
NUMA Architecture

2015
Intel Skylake X
SNC Architecture

2018
AWS Graviton
ARM Server

2019
Intel Optane DC
PM Module

2022
Samsung CXL
Memory Expander

# Observation 1: Rapid Evolution of Workload

- The applications that databases need to support are continuously growing.
- Within the past 3-4 years, we have seen a drastic shift in the ML workloads.
  - Large Language Models (LLM)

Amazon Redshift ML
**LLM invocation**

Databricks
**ai-query() function with foundation models**

SAP HANA vector engine
**Data Analysis with LLMs**

2023                           2024                           2024

2023                           2024

BigQuery ML
**SQL-only LLM Text Generation**

Oracle Heatwave Gen AI
**In-Database LLM**

PURDUE
UNIVERSITY®

# The current state

Hardware

Workload

What's the current state?

# The Current State: Observation 2

Hardware                                    Workload

2. Both the hardware and workload are **heterogeneous.**

# Heterogeneity in Hardware

intel
NVIDIA
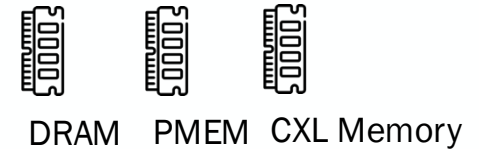AMD
aws

x86
RISC    IBM Power    RISC-v

Different Processor Architecture

P-core    E-core

Compute Cores

Different Vendors

SAMSUNG
ASUS
micron

SSD

Cores    IMC    Chip    NUMA Node

Chiplet

NUMA

NUMA VS Chiplet Architecture

DRAM    PMEM    CXL Memory

Memory

# Heterogeneity in Workload

ML Workload
LLM Queries
Training Data Preparation
Real-time Prediction

Spatial Workload
Location based Services
Geographic Information Systems
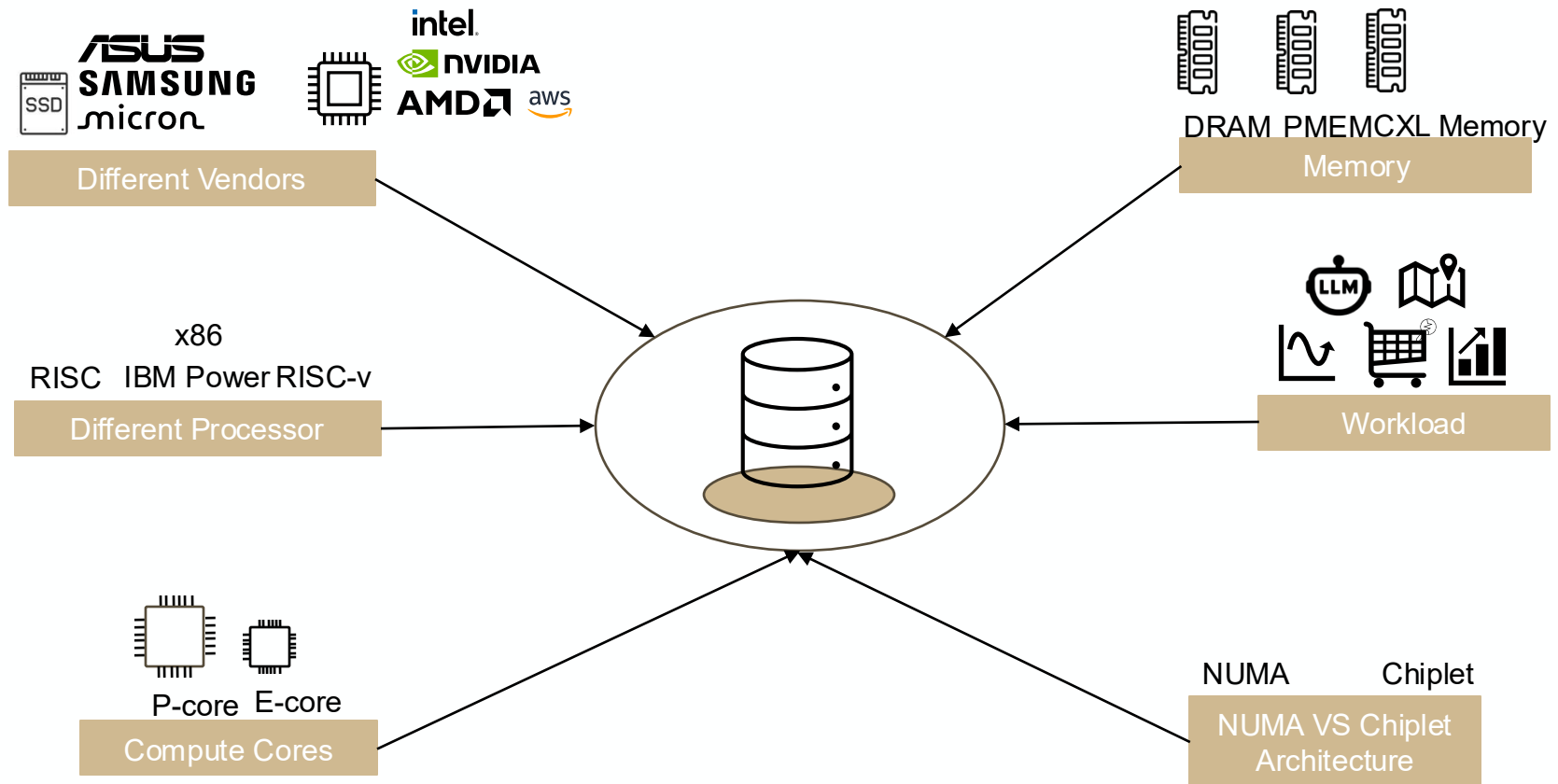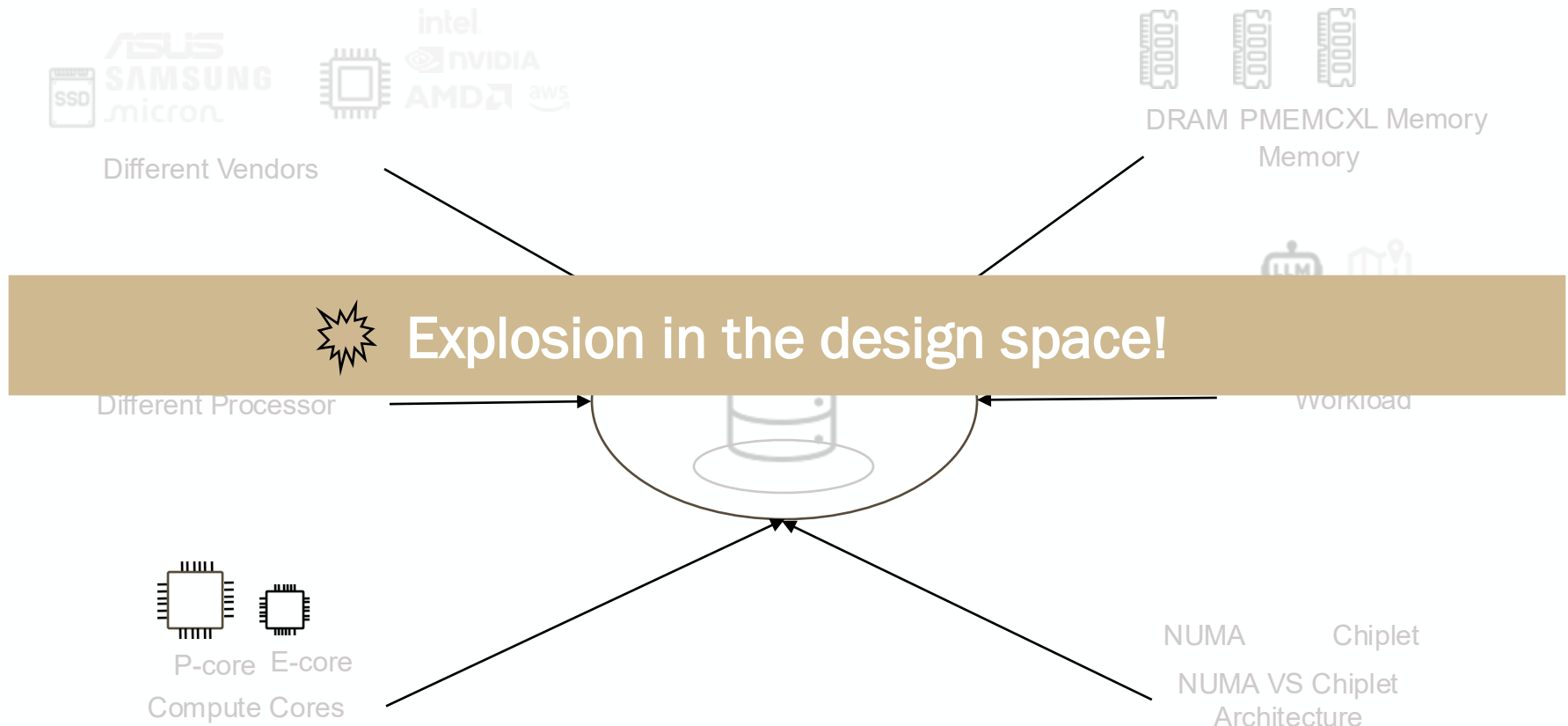Moving Objects

OLTP

HTAP

OLAP

# Designing Databases for the Modern Era



ASUS SAMSUNG micron SSD

intel NVIDIA AMD aws

**Different Vendors**

DRAM PMEM CXL Memory
**Memory**

x86
RISC IBM Power RISC-v
**Different Processor**

LLM
**Workload**

P-core E-core
**Compute Cores**

NUMA Chiplet
**NUMA VS Chiplet Architecture**

# Designing Databases for the Modern Era

Different Vendors

intel
NVIDIA
AMD
aws

DRAM PMEM CXL Memory
Memory

**Explosion in the design space!**

Different Processor

Workload

P-core  E-core
Compute Cores

NUMA          Chiplet

NUMA VS Chiplet
Architecture

PURDUE
UNIVERSITY®

# What's Required?

Databases that are **generalizable** across heterogeneous
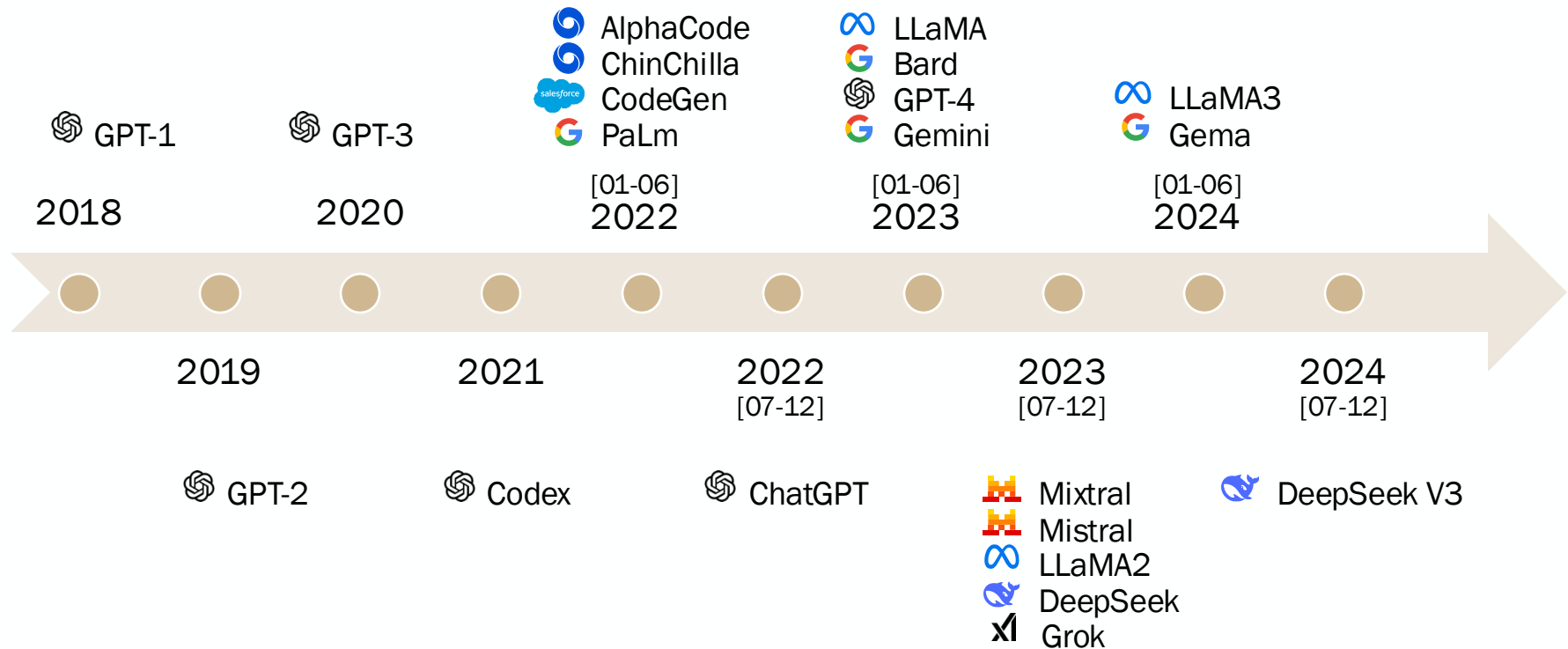hardware and workload
without sacrificing performance.

# What's Required?

Databases that are **generalizable** across heterogeneous hardware and workload without sacrificing performance.

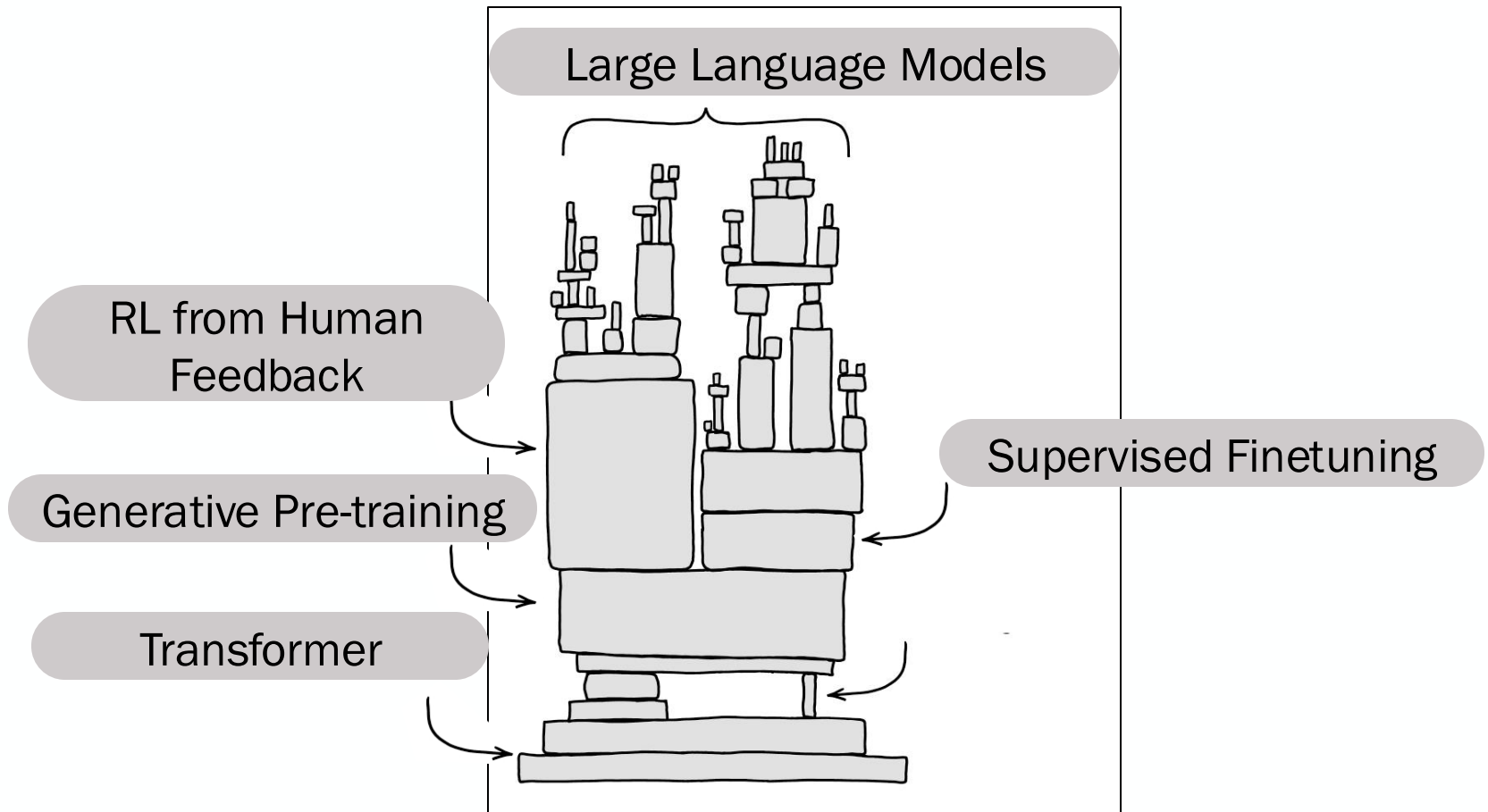## Large Language Models
## (LLM)

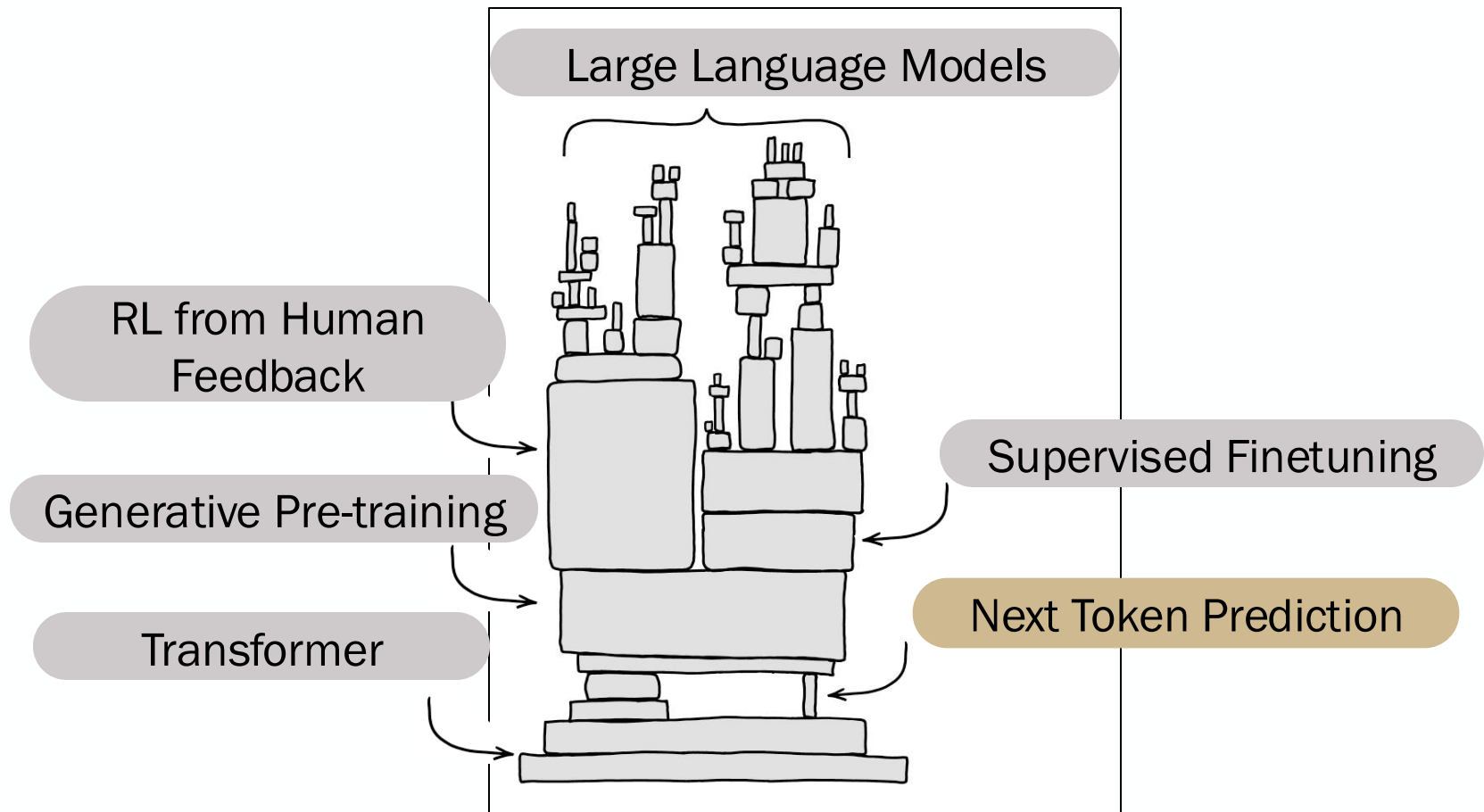# Large Language Models (LLM)

and

# Next Token Prediction (NTP)

PURDUE
UNIVERSITY®

# LLMs are everywhere!



GPT-1

GPT-3

AlphaCode
ChinChilla
CodeGen
PaLm

[01-06]
2022

LLaMA
Bard
GPT-4
Gemini

[01-06]
2023

LLaMA3
Gema

[01-06]
2024

2018

2020

2019

2021

2022
[07-12]

2023
[07-12]

2024
[07-12]

GPT-2

Codex

ChatGPT

Mixtral
Mistral
LLaMA2
DeepSeek
Grok

DeepSeek V3

# The LLM Recipe



Large Language Models

RL from Human Feedback

Generative Pre-training

Transformer

Supervised Finetuning

# First Step Towards Adopting the LLM Recipe

Large Language Models

RL from Human Feedback

Generative Pre-training

Transformer

Supervised Finetuning

Next Token Prediction

# What is Next Token Prediction?

- Next Token Prediction (NTP):
  - Model the probability of the next token in a given sequence based on the past tokens.
  - Token: a discrete unit, i.e., a word, a sub-word, or a character
- For any sequence $\boldsymbol{\tau}$
  - $\tau_i$ denotes the $i$-th token, and
  - $\boldsymbol{\tau_{<i}}$ denotes the $i - 1$ tokens preceding $\tau_i$
- Goal of NTP
  - Estimate the probability distribution of $\tau_i$:
  $$\mathbb{P}(\tau_i \,|\boldsymbol{\tau_{<i}})$$

# Generating a Sequence via NTP

| SIGMOD | is | being | 25 | held | Berlin | in | Germany |

Vocabulary: Possible set of tokens

| SIGMOD |

First Token

# Generating a Sequence via NTP

| SIGMOD | is | being | 25 | held | Berlin | in | Germany |

Vocabulary: Possible set of tokens

SIGMOD

First Token     Next Token = ?

# Generating a Sequence via NTP

SIGMOD is being 25 held Berlin in Germany

Vocabulary: Possible set of tokens

SIGMOD

First Token

Next Token = ?

# Generating a Sequence via NTP

SIGMOD | is | being | 25 | held | Berlin | in | Germany

Vocabulary: Possible set of tokens

SIGMOD

First Token

Next Token = ?

# Generating a Sequence via NTP

- Estimate the probability distribution of each possible next tokens.



SIGMOD

First Token

Next Token = ?

# Generating a Sequence via NTP

- Select the token with the highest probability distribution as the next token.



SIGMOD

First Token

being   is   SIGMOD   25   held

held   Berlin   in   Germany

Next Token = ?

Select

25

# Generating a Sequence via NTP

| SIGMOD | is | being | 25 | held | Berlin | in | Germany |

Vocabulary: Possible set of tokens

| SIGMOD | | 25 |

First Token     Next Token

# Generating a Sequence via NTP

| SIGMOD | is | being | 25 | held | Berlin | in | Germany |
|--------|-----|-------|-----|------|--------|-----|---------|

Vocabulary: Possible set of tokens

| SIGMOD | 25 | is | being | held | in | Berlin | Germany |
|--------|-----|-----|-------|------|-----|--------|---------|

Full Sentence

# Next Token Prediction (NTP)

and

# Database Systems

# Translating NTP into Database Systems

- **Join Order Selection** (A ⋈ B ⋈ C ⋈ D)
  - Tokens: The set of tables: A, B, C, D

Predicting the next table to join

A    C    D    B

- **Transaction scheduling** (T1, T2, T3, T4)
  - Tokens: The set of transactions: T1, T2, T3, T4

Predicting the next transaction to schedule

T4    T3    T2    T1

# Hindrances of Translating NTP into DBMS

1. The mismatch in objectives between the domain of NLP and Databases

| NLP tasks | Database optimization tasks |
|---|---|
| Generative<br>LLMs generate coherent sequences. | Goal-oriented |
| In Generative Pre-training phase, the goal is to compress a significant amount of world knowledge into the LLM by training on a diverse internet-scale corpus. | Improve query performance, scalability, and resource utilization. |

# Hindrances of Translating NTP into DBMS

2. The mismatch in the notion of Tokens in NLP and Databases

| NLP tasks | Database optimization tasks |
|---|---|
| The notion of token is fixed for a particular tokenizer. | The notion of token is diverse and irregular.<br>• In JOS, tables are tokens.<br>• In scheduling, transactions are tokens. |

# NLP Tokens VS Database Tokens

| SIGMOD | 25 | is | being | held | in | Berlin | Germany |

- Syntactic regularity:
  - The verb follows the subject.

☒ | SIGMOD | 25 | in | Germany | Berlin | held | is | being |

- Contextual meaning:
  - The mention of Germany implicitly excludes other locations, e.g., USA.

☒ | SIGMOD | 25 | is | being | held | in | Chicago | USA |

**PURDUE** UNIVERSITY.

# NLP Tokens VS Database Tokens

A  C  D  B

Join Order Selection (A ⋈ B ⋈ C ⋈ D)

- Lack of syntactic regularity:
  - Two different database instances can have tables with same name but with different attributes.

- Lack of contextual meaning:
  - A C D B    does not eliminate the possibility of

    D A C B

# Contribution of this Paper



Database Systems 🤝 Next Token Prediction

Database Systems Meet Next Token Prediction

# Building Blocks

Database Systems 🤝 Next Token Prediction

🧱 DB-Tokens    🧱 Decision Transformer

**1. Decision Transformer:** Sequence Modeling
**2. DB-Tokens:** Generalization

# Building Block 1: Decision Transformer

- Goal-directed RL:
    - Treats RL as a supervised sequence modeling problem
    - Predict the next action (token) by conditioning on a desired reward, e.g., query throughput, latency, scalability etc.

$$\mathbb{P}(a_i \,|s_i, \hat{R}_i, a_{<i}, s_{<i}, \hat{R}_{<i})$$

- Each trajectory (policy) is represented as a sequence of (return-to-go, state, action) tuples.

$$\tau \;=\; \boxed{\hat{R}_1} \; \boxed{s_1} \; \boxed{a_1} \quad \boxed{\hat{R}_2} \; \boxed{s_2} \; \boxed{a_2} \qquad \boxed{\hat{R}_T} \; \boxed{s_T} \; \boxed{a_T}$$

Reward-to-go token
Future cumulative reward expected from a given timestep onward.

# Building Block 1: Decision Transformer

- Scalable across large datasets:
  - The underlying model is a Transformer architecture.
  - It follows the Offline Reinforcement Learning Paradigm
    - The model is trained on an offline dataset.
    - The model does not interact with the environment.



Reward, State

Action

Online RL

Reward, State, Action

Logged Interaction

Offline RL

# Building Block 2: DB-Tokens

- Hardware profiles generated from hardware Performance Monitoring Units (PMU, for short)
  - Computationally inexpensive to retrieve from the hardware registers
  - **Generalizable** across different hardware and workload applications
  - Can provide accurate hardware context that the DBMS is running on
  - Can mimic the data distribution and query workload

# The Framework

Database Systems

Next Token Prediction

DB Tokens

Decision Transformer

Probe and Learn

(Po Le)

# Case Study: Index Scheduling

Index: B⁺ Tree

Local: 70 – 80 ns

Remote: 120 – 220 ns

Index node

IMC Tile    Core Tile

NUMA Machine

# Case Study: Index Scheduling



Index: B+ Tree

0  1  2  3

Schedule the
index nodes
to

○ Index node

② a particular NUMA node

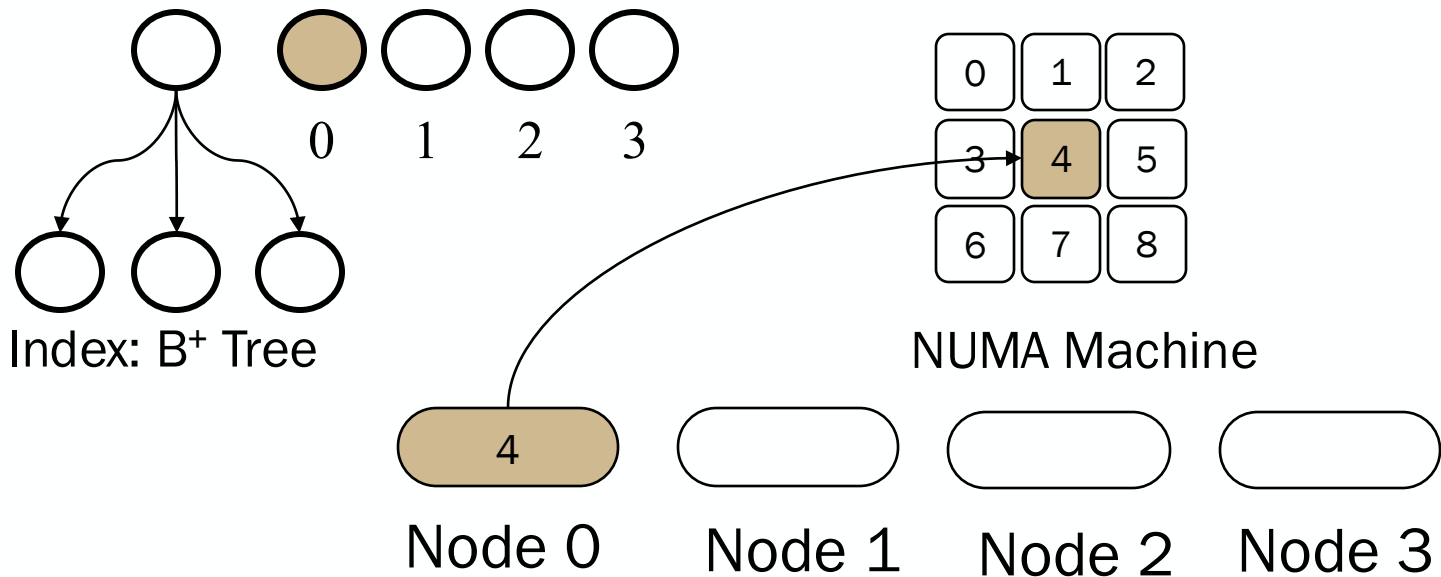① a particular core

△ IMC Tile    ☐ Core Tile

NUMA Machine

# Index Scheduling: NTP Formulation

- Predict the next core to place the $i$-th index node
  - A policy in index scheduling = Sequence of core IDs

Index: B⁺ Tree

NUMA Machine

0  1  2
3  4  5
6  7  8

Predict the next core to place

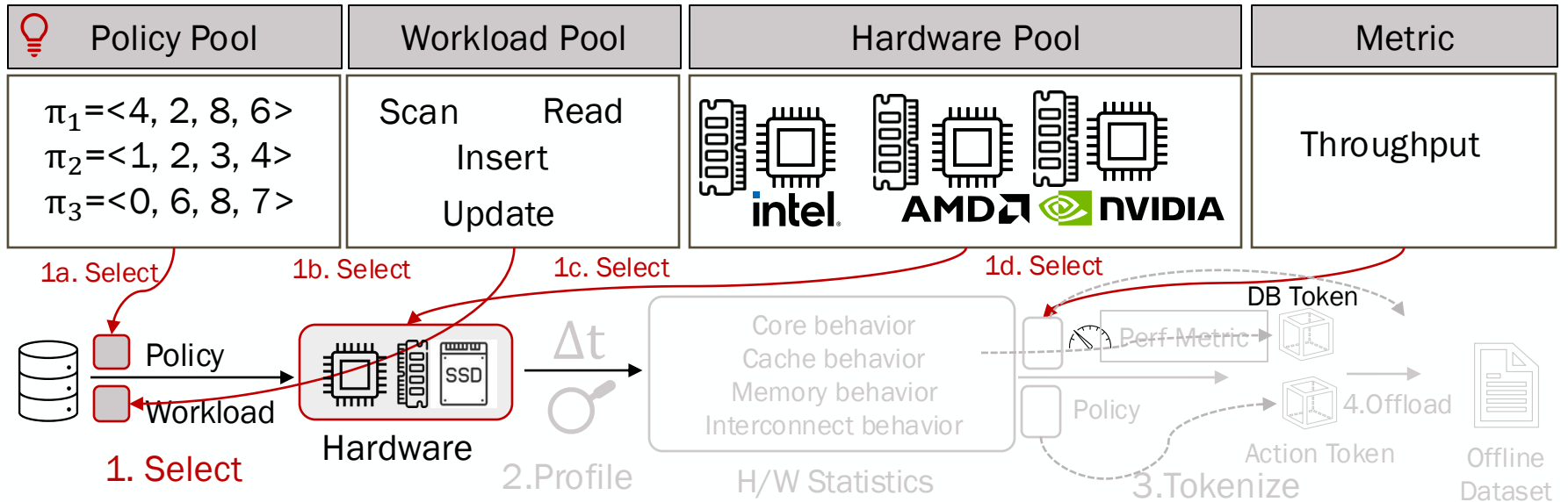Node 0    Node 1    Node 2    Node 3

# Index Scheduling: NTP Formulation

- Predict the next core to place the $i$-th index node
  - A policy in index scheduling = Sequence of core IDs



Index: B⁺ Tree

NUMA Machine

Predict the next
core to place

Node 0    Node 1    Node 2    Node 3

# Index Scheduling: NTP Formulation

- Predict the next core to place the $i$-th index node
  - A policy in index scheduling = Sequence of core IDs



Index: B⁺ Tree

NUMA Machine

**Predict the next core to place**

Node 0    Node 1    Node 2    Node 3

# Index Scheduling: NTP Formulation

- Predict the next core to place the $i$-th index node
  - A policy in index scheduling = Sequence of core IDs



Index: B+ Tree

NUMA Machine

**Predict the next core to place**

| 4 | 2 | 8 | 6 |
|---|---|---|---|
| Node 0 | Node 1 | Node 2 | Node 3 |

# Index Scheduling: NTP Formulation

■ Predict the next core to place the $i$-th index node

   ■ A policy in index scheduling = Sequence of core IDs



Index: B⁺ Tree

NUMA Machine

**Predict the next core to place**

Node 0   Node 1   Node 2   Node 3

# Index Scheduling: DB-Tokens

- Predict the next core to place the $i$-th index node
  - A policy in index scheduling = Sequence of core IDs



Index: B+ Tree

NUMA Machine

Node 0    Node 1    Node 2    Node 3

**DB-Tokens:** The L1, L2, LLC cache misses, branch misses, TLB misses, local and remote memory accesses of <Core 4>, query throughput of <Core 4>.

# Probe Phase of PoLe Framework Index Scheduling

# Probe Phase: Select (Step 1)



| Policy Pool | Workload Pool | Hardware Pool | Metric |
|---|---|---|---|
| $\pi_1 = <4, 2, 8, 6>$ <br> $\pi_2 = <1, 2, 3, 4>$ <br> $\pi_3 = <0, 6, 8, 7>$ | Scan   Read <br> Insert <br> Update | | Throughput |

1a. Select   1b. Select   1c. Select   1d. Select

Policy
Workload
1. Select

Hardware
2.Profile

Δt

Core behavior
Cache behavior
Memory behavior
Interconnect behavior
H/W Statistics

Perf-Metric
Policy

DB Token

4.Offload
Action Token
3.Tokenize

Offline
Dataset

**1. Select:** Execute different policies across various hardware configurations and workloads, drawn from a diverse set of policy, hardware, and workload pools.

# Probe Phase: Profile (Step 2)

| | | | |
|---|---|---|---|
| $\pi_1$=<4, 2, 8, 6><br>$\pi_2$=<1, 2, 3, 4><br>$\pi_3$=<0, 6, 8, 7> | Scan    Read<br>Insert<br>Update | | Throughput |

1a.Choose    1b.Choose    1c.Choose    1d.Choose

DB Token

Policy

Workload

**Hardware**

1. Select

Δt

**2.Profile**

Core behavior
Cache behavior
Memory behavior
Interconnect behavior

H/W Statistics

Perf-Metric

Policy

4.Offload

Action Token

3.Tokenize

Offline
Dataset

**2. Profile:** During query execution, periodically profile the hardware to capture the behavior of crucial hardware components.

# Probe Phase: Tokenize (Step 3)



| 1a.Choose | 1b.Choose | 1c.Choose | 1d.Choose |
| --- | --- | --- | --- |

$\pi_1 = <4, 2, 8, 6>$
$\pi_2 = <1, 2, 3, 4>$
$\pi_3 = <0, 6, 8, 7>$

Scan    Read
Insert
Update

Throughput

Policy
Workload
1. Select

Hardware

$\Delta t$

2.Profile

**Core behavior**
**Cache behavior**
**Memory behavior**
**Interconnect behavior**

H/W Statistics

Perf Metric

DB Token

Policy

3.Tokenize

Action Token

4.Offload

Offline Dataset

**3. Tokenize:** Tokenize the hardwire profiles, alongside the desired performance metric and the current policy.

# Probe Phase: Offload (Step 4)



$\pi_1 = <4, 2, 8, 6>$
$\pi_2 = <1, 2, 3, 4>$
$\pi_3 = <0, 6, 8, 7>$

1a.Choose

Scan     Read
Insert
Update

1b.Choose

1c.Choose

intel.  AMD  NVIDIA

1d.Choose

Throughput

Policy
Workload

1. Select

Hardware

$\Delta t$

2.Profile

Core behavior
Cache behavior
Memory behavior
Interconnect behavior

H/W Statistics

Perf-Metric

Policy

3.Tokenize

DB Token

4.Offload

Action Token

Offline
Dataset

**4. Offload:** Offload the action tokens along with the associated DB-tokens to an offline dataset. This is provided to the Decision Transformer during the learning phase.

# Learn Phase of PoLe Framework Index Scheduling

# Learn Phase: Training (Step 1)



1. **Training:** Train a Decision Transformer (DT) on the collected offline dataset in a supervised manner.

# Learn Phase: Inference (Step 2)



2. **Inference:** Infer a new policy via auto-regression using the trained Decision Transformer.

# Experiment Settings

- Index:  A main-memory B$^+$-Tree

- Workload:  YCSB-A workload

- Baselines
  - OS Scheduling Policies
    - Default, Local, Interleaved
    - OS handles core scheduling.
  - Heuristics:
    - Shared Everything NUMA
      - Nearby index nodes are placed on the same NUMA node.
      - OS handles scheduling.
    - Shared Nothing
      - Nearby index nodes are placed on the same NUMA node.
      - Core Scheduling follows the data placement strategy.

# Preliminary Results

- On seen hardware (Intel Sandy Bridge and NVIDIA Grace Hopper)
  - PoLe outperforms the baselines by up to 2.78×

- On unseen hardware (Intel Skylake X)
  - PoLe outperforms the baselines by up to 3×

# PoLe's Learned Scheduling Policies



Intel Sandy Bridge ($\pi_l$)   NVIDIA Grace Hopper ($\pi_l$)   Intel Skylake X ($\pi_l$)

Cells sharing the same color indicate that the associated index chunks are scheduled on the same NUMA server.

# PoLe's Learned Scheduling Policies

Intel Sandy Bridge ($\pi_l$)

| 32 | 28 | 24 | 44 | 16 | 28 | 48 | 48 | 29 | 29 | 29 | 37 | 25 | 49 | 45 | 51 |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| 45 | 41 | 21 | 25 | 29 | 51 | 51 | 13 | 33 | 42 | 34 | 45 | 30 | 34 | 27 | 41 |
| 34 | 45 | 30 | 34 | 29 | 34 | 34 | 30 | 30 | 13 | 36 | 12 | 12 | 12 | 12 | 14 |
| 26 | 13 | 35 | 20 | 34 | 43 | 39 | 39 | 28 | 26 | 26 | 33 | 26 | 33 | 35 | 50 |
| 16 | 16 | 16 | 16 | 14 | 18 | 38 | 38 | 47 | 47 | 47 | 50 | 39 | 16 | 21 | 47 |
| 47 | 28 | 39 | 47 | 47 | 12 | 39 | 39 | 28 | 12 | 28 | 41 | 12 | 16 | 28 | 41 |
| 31 | 31 | 31 | 31 | 48 | 22 | 13 | 13 | 13 | 23 | 23 | 29 | 29 | 31 | 31 | 31 |
| 33 | 21 | 45 | 45 | 33 | 44 | 14 | 46 | 46 | 30 | 36 | 36 | 45 | 36 | 17 | 17 |
| 21 | 43 | 21 | 44 | 51 | 19 | 19 | 15 | 26 | 19 | 29 | 15 | 49 | 49 | 40 | 42 |
| 30 | 30 | 49 | 49 | 13 | 35 | 17 | 17 | 17 | 39 | 39 | 29 | 29 | 26 | 30 | 29 |
| 29 | 14 | 29 | 42 | 41 | 26 | 41 | 37 | 37 | 37 | 37 | 30 | 30 | 39 | 13 | 15 |
| 15 | 39 | 26 | 13 | 17 | 46 | 17 | 34 | 13 | 13 | 13 | 35 | 26 | 31 | 31 | 22 |
| 22 | 22 | 46 | 46 | 42 | 33 | 33 | 46 | 46 | 46 | 37 | 33 | 33 | 33 | 46 | 37 |
| 33 | 33 | 37 | 39 | 39 | 22 | 24 | 39 | 21 | 13 | 26 | 26 | 17 | 24 | 24 | 30 |
| 30 | 16 | 24 | 19 | 19 | 43 | 24 | 24 | 43 | 31 | 43 | 15 | 21 | 21 | 21 | 21 |
| 50 | 16 | 15 | 31 | 31 | 38 | 43 | 38 | 38 | 38 | 31 | 50 | 38 | 50 | 38 | 38 |

**32** → Index Chunk (A collection of index nodes)

Core ID
where the chunk is scheduled

→ Scheduling Policy

Cells sharing the same color indicate that the associated index chunks are scheduled on the same NUMA server.

**PURDUE** UNIVERSITY.

# PoLe's Learned Scheduling Policies



Cells sharing the same color indicate that the associated index chunks are scheduled on the same NUMA server.

# PoLe's Scheduling Policies

## PoLe's Learned Scheduling Policies



Intel Sandy Bridge ($\pi_l$)     NVIDIA Grace Hopper ($\pi_l$)     Intel Skylake X ($\pi_l$)

## Scheduling Policies in the training set

# PoLe's Learned Scheduling Policies

## PoLe's Learned Scheduling Policies



Intel Sandy Bridge ($\pi_l$)  NVIDIA Grace Hopper ($\pi_l$)  Intel Skylake X ($\pi_l$)

1. Learned Policies are different from the observed policies in the training set.
2. Different hardware ➡ Different learned policies.

# Conclusion

- PoLe brings **Next Token Prediction** into the world of database optimization.

- It leverages **offline RL and DB-Tokens** to learn generalizable scheduling strategies.

- What's next?
  - Categorize the optimization tasks that can benefit from the Next Token Prediction (NTP) paradigm
  - Assess to what extent the PoLe framework can provide consistent performance and adaptivity guarantees

# Thank You!
## Questions?



Read our paper!



Get in touch!