# EPFL
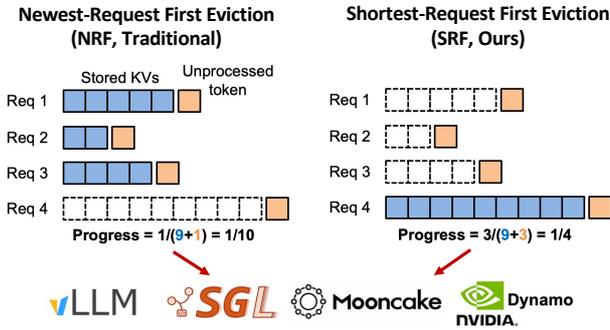
# Harmonizing Data Systems, LLMs, and Vectors :)
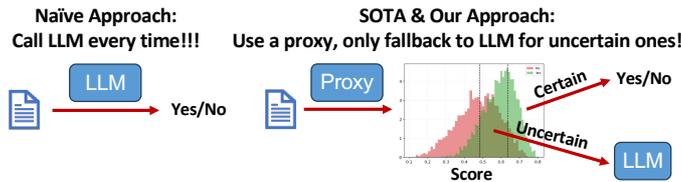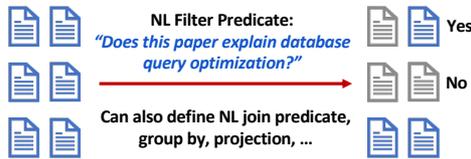
Data-Intensive Applications and Systems Laboratory

**AiAS** — Data-Intensive Applications and Systems

## LLM Inference System = Data System

**Newest-Request First Eviction (NRF, Traditional)**

**Shortest-Request First Eviction (SRF, Ours)**



Stored KVs — Unprocessed token

Req 1, Req 2, Req 3, Req 4

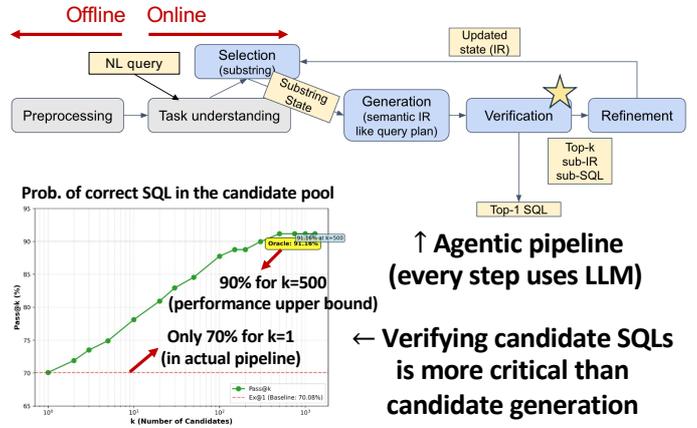**Progress = 1/(9+1) = 1/10**

**Progress = 3/(9+3) = 1/4**

vLLM · SGL · Mooncake · NVIDIA Dynamo

**Data management is crucial (~20% GPU hour saving in LLM inference by doing better KV cache replacement)**

## Agentic LLM Pipeline for Translating NL to SQL

Offline | Online



NL query → Selection (substring) → Substring State → Updated state (IR) → Generation (semantic IR like query plan) → Verification → Refinement

Preprocessing → Task understanding

Top-k sub-IR sub-SQL → Top-1 SQL

**↑ Agentic pipeline (every step uses LLM)**

**Prob. of correct SQL in the candidate pool**



**90% for k=500 (performance upper bound)**

**Only 70% for k=1 (in actual pipeline)**

**← Verifying candidate SQLs is more critical than candidate generation**

## Bulk LLM Calls for Semantic Operators

**NL Filter Predicate:** *"Does this paper explain database query optimization?"* → Yes / No

**Can also define NL join predicate, group by, projection, ...**

**Naïve Approach: Call LLM every time!!!**

LLM → Yes/No

**SOTA & Our Approach: Use a proxy, only fallback to LLM for uncertain ones!**

Proxy → Certain → Yes/No ; Uncertain → LLM

Score

**Saving # LLM calls by ~99%, and 50% on average**

## Why Not Vector *Joins* for Batch Query Proc.?



Vector (point) — $n$ — $\bowtie_{dist(X, Y) < \theta}$ — Data Index

X — Queries / Vector columns (possibly with indexes) — Y — Data

**IndexNestedLoopJoin(X, Y, θ)**
For each query x in X:
  Y.IndexSearch(x, θ)

Search seed

**Cache query results, reuse for similar queries**

**Work Sharing**

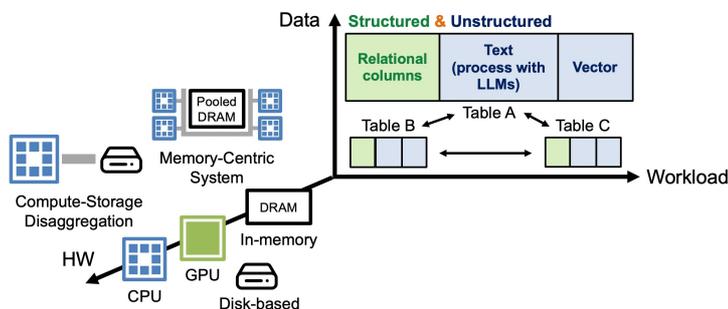**Offload search to offline phase**

**Work Offloading = HashTableBuild**

**Work sharing/offloading: ~33x speedup, ~43% higher recall**

## In-memory Similarity Caching for Disk-based Vector Search



**Fast Local Search** (Scans the hottest mini-indexes first)

**Semantic Threshold Check** (Checks if the distance falls within the dynamically learned spatial threshold)

**Incoming Query Vector** → **Cache Search** (MI-2, MI-1, MI-0, MI-3) Eviction Policy (MRU to LRU) → **Is Hit?**

**Sub-millisecond Hit** (Served entirely from local in-memory mini-indexes) → **Return IDs + Distances**

**Fallback & Execution** (Exact search executed on disk/memory backend) → **Backend Vector Database** → **Return IDs + Distances**

**In-memory mini-index**

**Async Cache Fill** (Fetches new vectors and inserts them into the hottest MRU mini-index, evicting the LRU if full)

**Async Threshold Learning** (Updates region-specific similarity thresholds based on ground-truth backend distances)

Cache Hit Ratio / Latency (P50) / Query Throughput / 10-Recall@10 / Vectors In Cache

**DB-agnostic caching delivers 40-1000x lower query latencies on cache hits**

## Our End Goal: HW-Optimized System Supporting Hybrid Relational/Vector/Semantic Queries



Data — **Structured & Unstructured**

Relational columns | Text (process with LLMs) | Vector

Table B — Table A — Table C — Workload

Pooled DRAM — Memory-Centric System

Compute-Storage Disaggregation

HW — CPU — GPU — DRAM — In-memory — Disk-based

**Lots of open challenges to solve!!!**
- One-size-fits-all solution for filtered vector search
- Declarative control of accuracy-efficiency trade-off
- Query optimization of vector/semantic queries
- LLM for query optimization and direct plan writing

**Under large & complex data, queries (100's ops), and HW configurations**