# Data-Aware Caching for Cloud Analytics

DB

Durable Latency

S3 Cost Cache

BI ML ... DS

Object storage

Data-aware Cache

Cache

Data-aware Evict

LiquidCache
10x { cost | latency } reduction
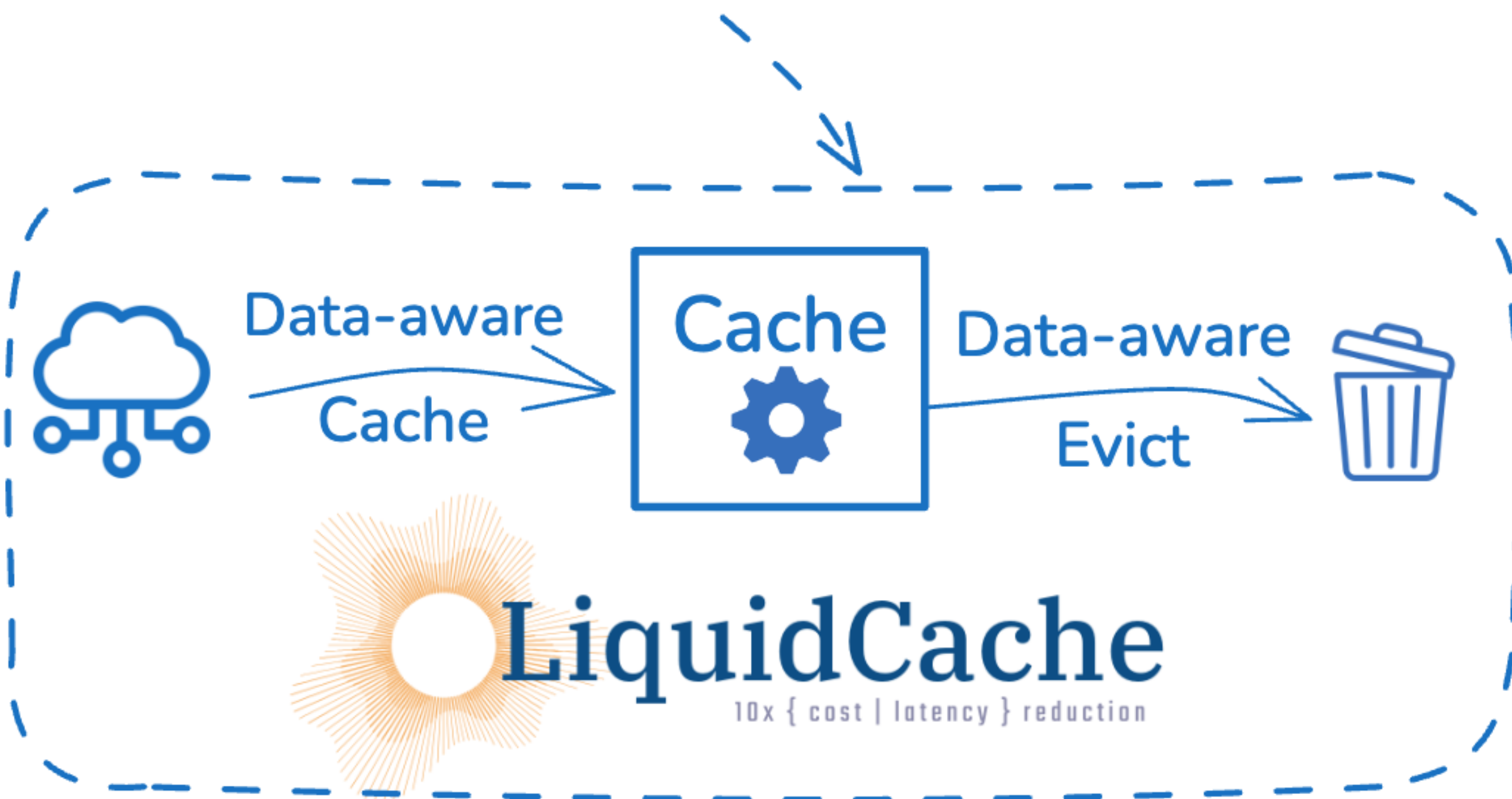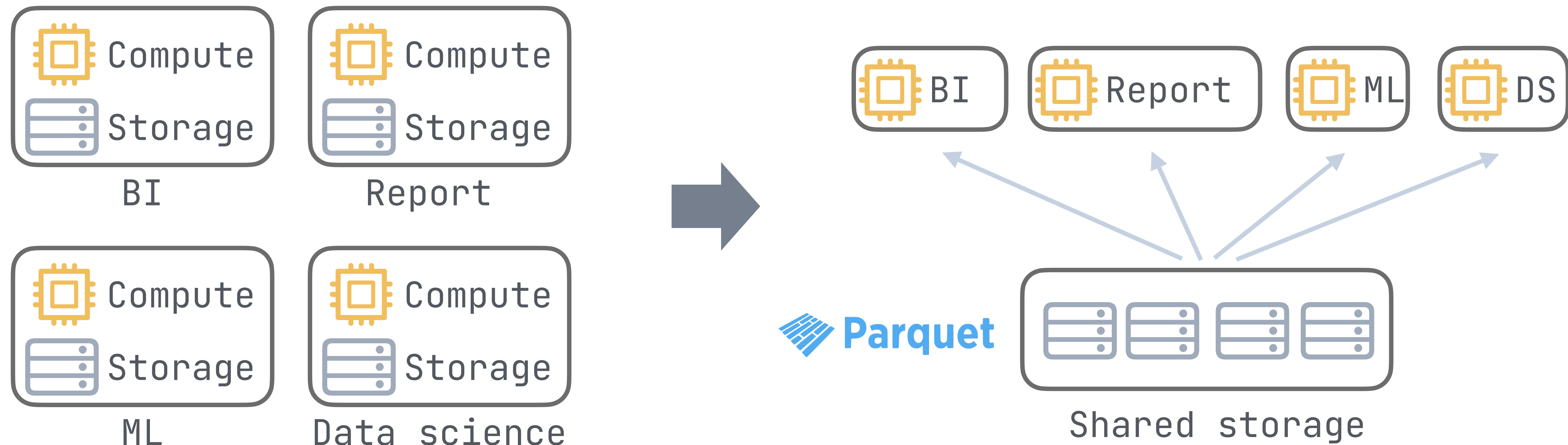
Xiangpeng Hao

Department of Computer Science
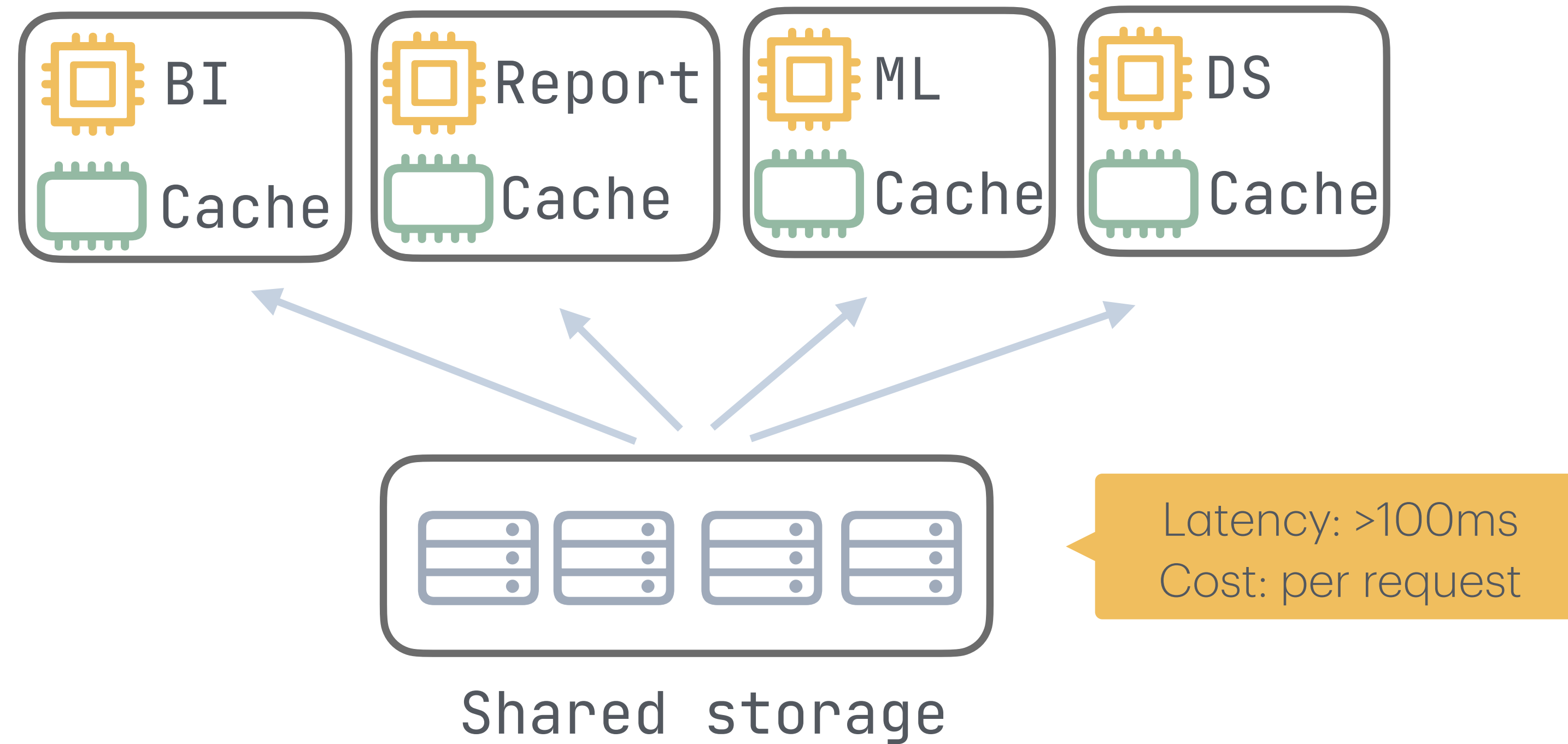University of Wisconsin-Madison

May 19th, 1 PM CDT
CS2310 or Zoom

# On premise → cloud (2010-2020)
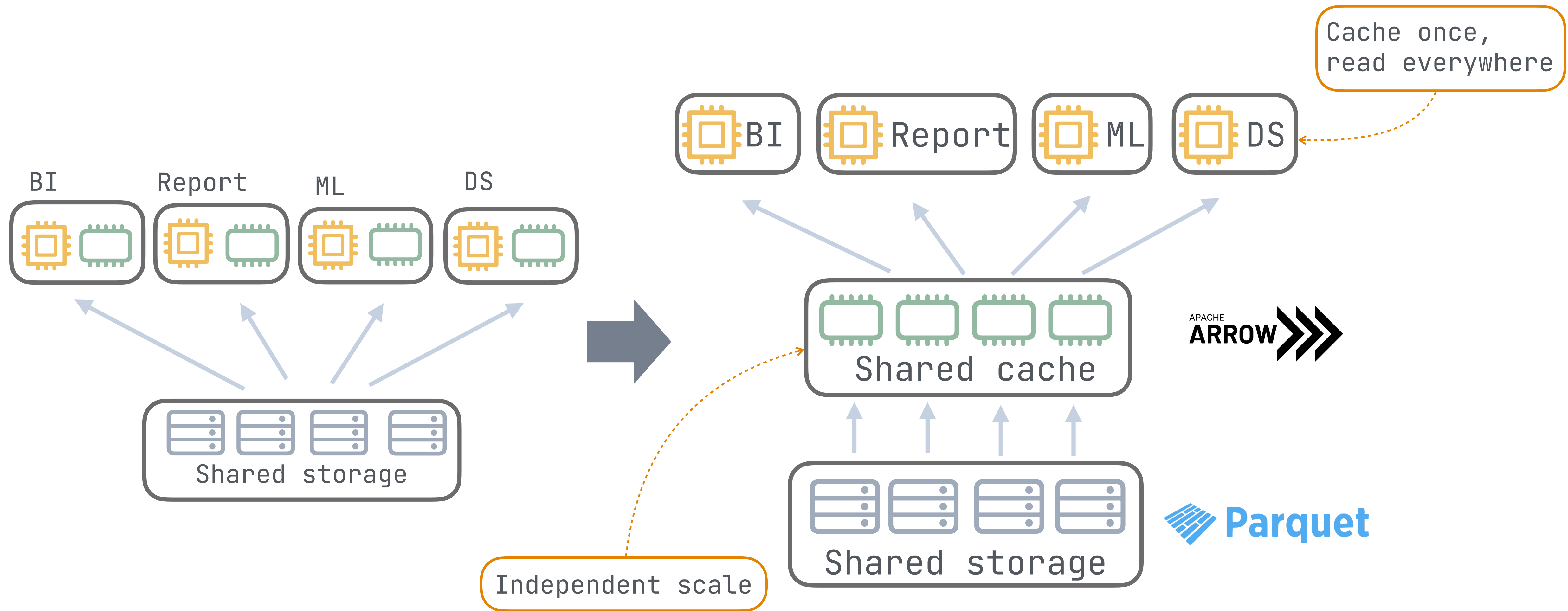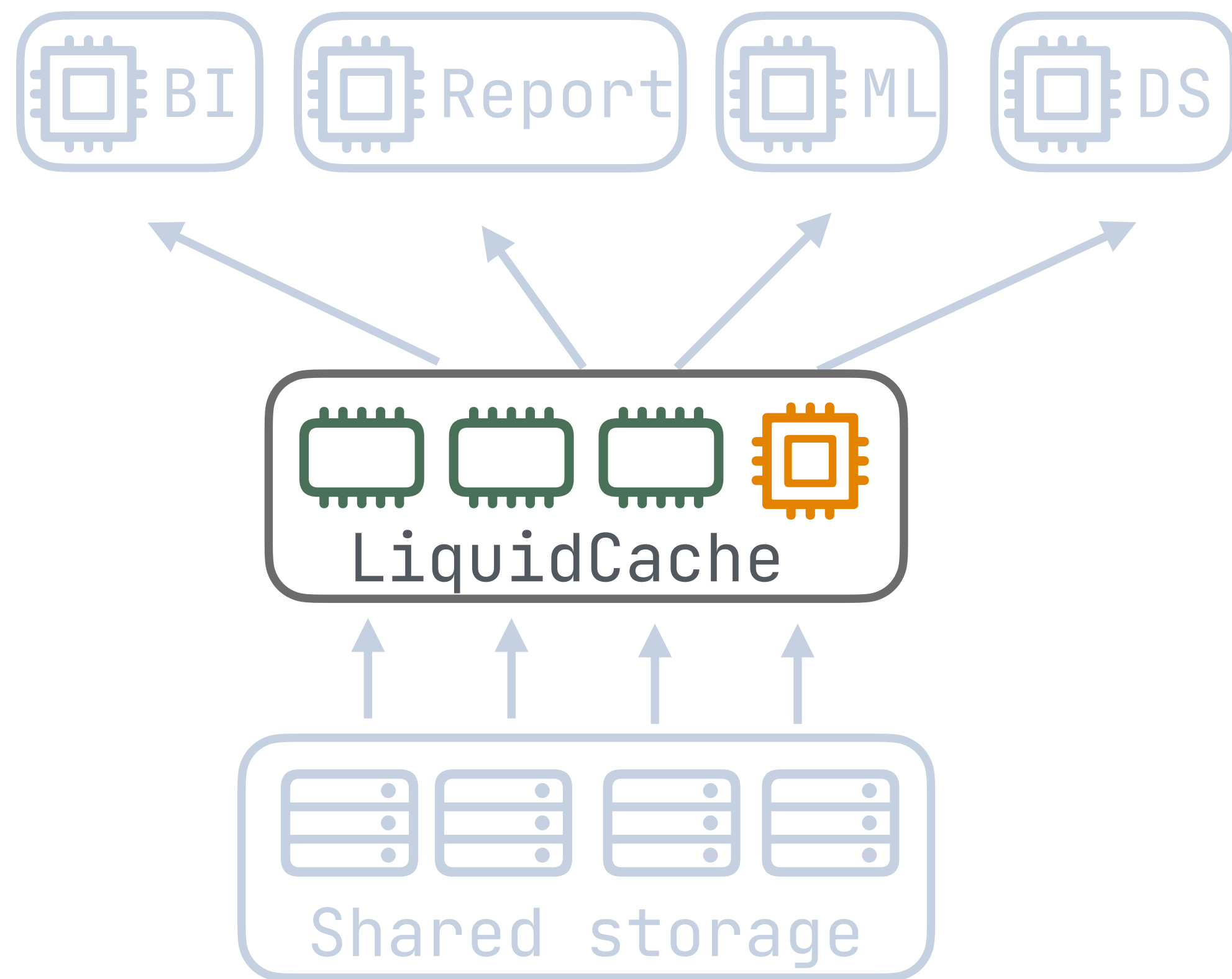
# Every system has a cache (2020-2025)

# Vision: shared cache (2025+)

BI  Report  ML  DS

BI  Report  ML  DS

Cache once, read everywhere

Shared cache

APACHE ARROW >>>
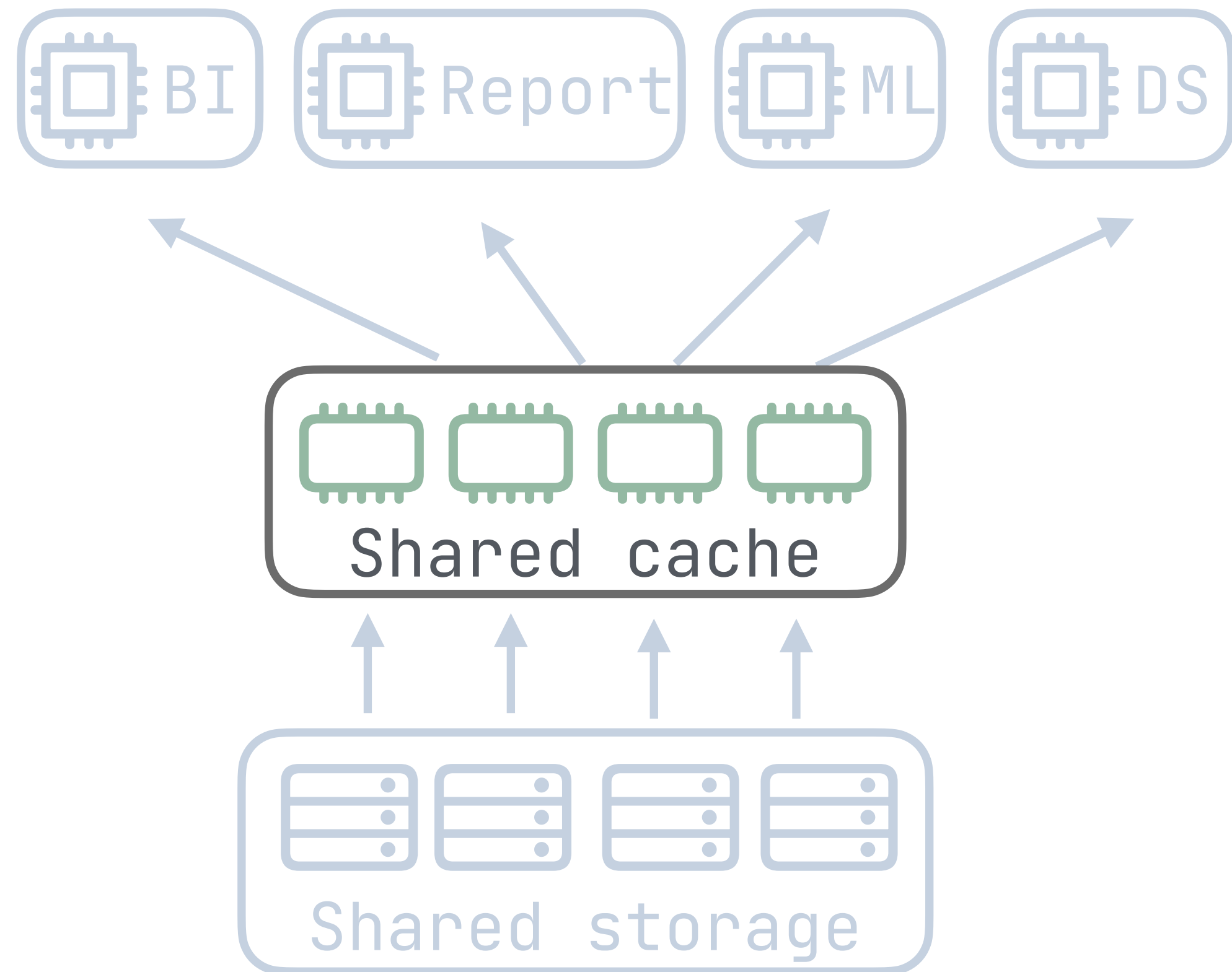
Shared storage

Independent scale

Shared storage

Parquet

# Thesis goal: LiquidCache



Thesis goal:

Design a cost-effective shared cache system by combining compute and data, while preserving ecosystem compatibility.

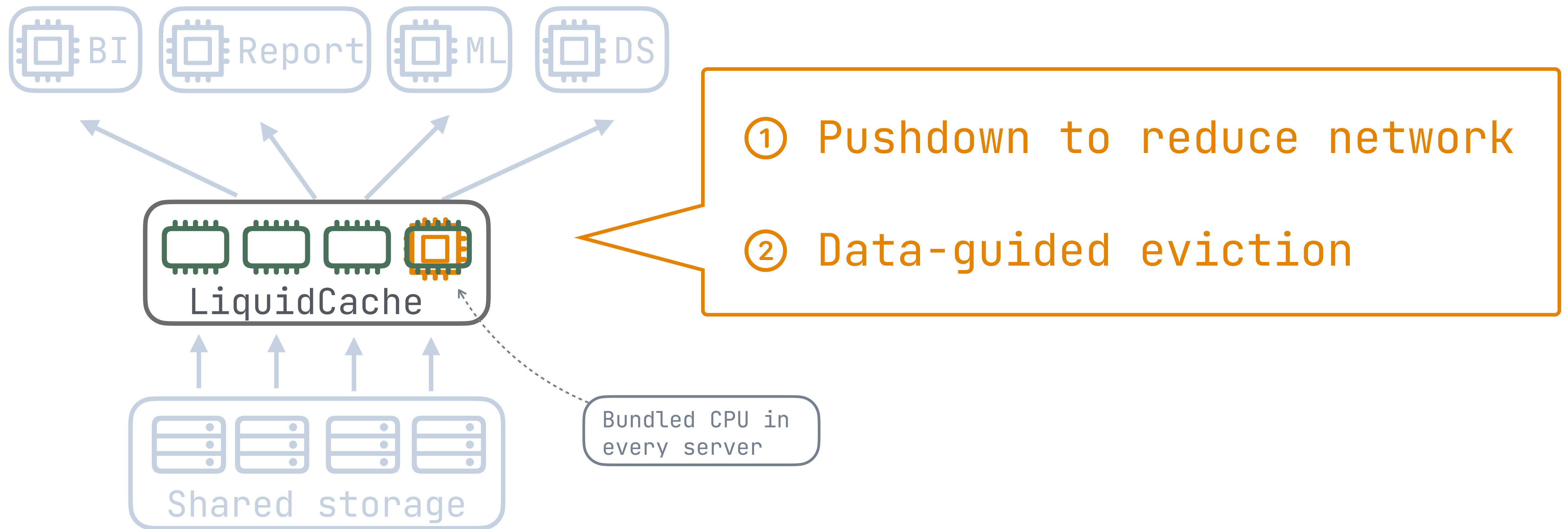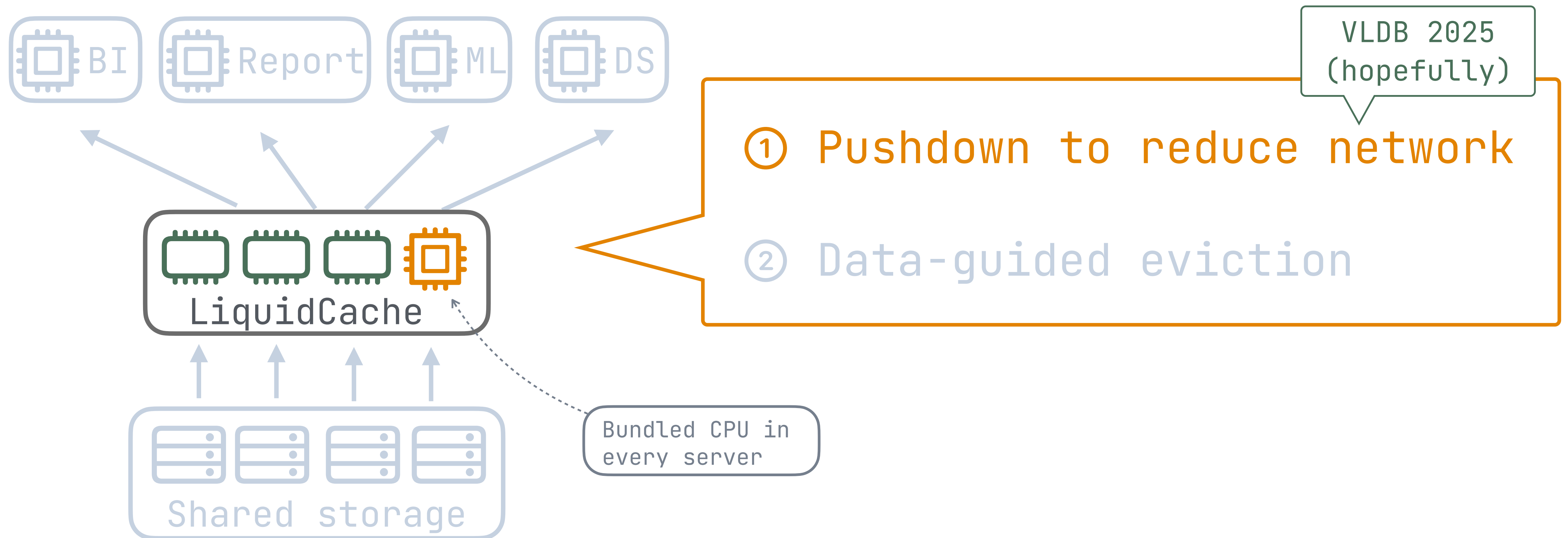# First attempt: byte cache



BI   Report   ML   DS

Shared cache

Shared storage

Challenges:

1. Network bottleneck

2. Inefficient cache eviction

# LiquidCache = compute + data



BI  Report  ML  DS

**LiquidCache**

Bundled CPU in every server

Shared storage

① Pushdown to reduce network

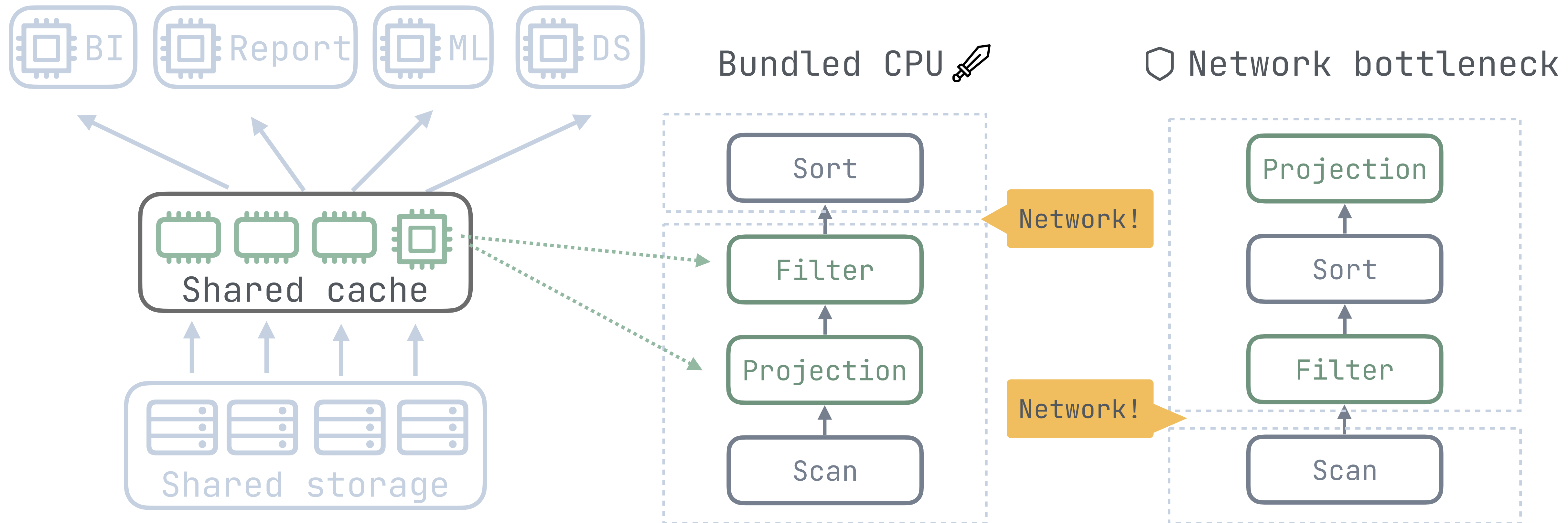② Data-guided eviction

# LiquidCache = compute + data

# Pushdown to reduce network

# Pushdown overwhelms cache CPU



Previously believed bottleneck:
**Filter evaluation**

Our findings:
**Data decoding**

# LiquidCache



**Part 1:**

co-designed format to skip decoding

**Part 2:**

progressive, selective, asynchronous transcoding

# Part 1:
## co-designed format to skip decoding

**Read one row**

**Decode all**

**Parquet page**

**Extract**

**Parquet page**

**Read one row**

**Decode one**

**Desired page**

**Co-design principle:**
Each row must be independently decodable

# Each row must be independently decodable (string example)
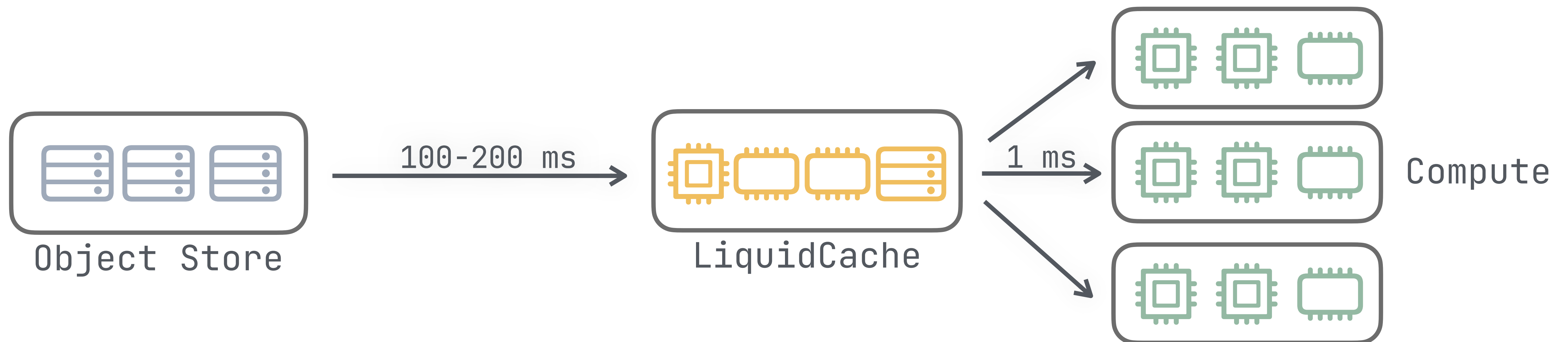


Liquid (Arrow) → Liquid (Dictionary) → Liquid (BitPacked) → Liquid

No general purpose compression

Leverages state-of-the-art encoding schemes

Carefully designed encoding/layout for each data types

# Co-design with filter pushdown (selective decoding)

```sql
SELECT val, location
FROM sensor_data
WHERE location = 'office';
```

# Co-design with filter pushdown
# (filter late materialization)

```sql
SELECT val, location
FROM sensor_data
WHERE location = 'office', date > '2025-03-21';
```

# Co-design with filter pushdown (evaluate on encoded data)

```sql
SELECT val, location
FROM sensor_data
WHERE name = 'Apache Arrow'
```



In paper:
Evaluate on **partially** encoded data

Baseline                                                    Evaluate on encoded

# LiquidCache



**Part 1:**

co-designed format to skip decoding

**Part 2:**

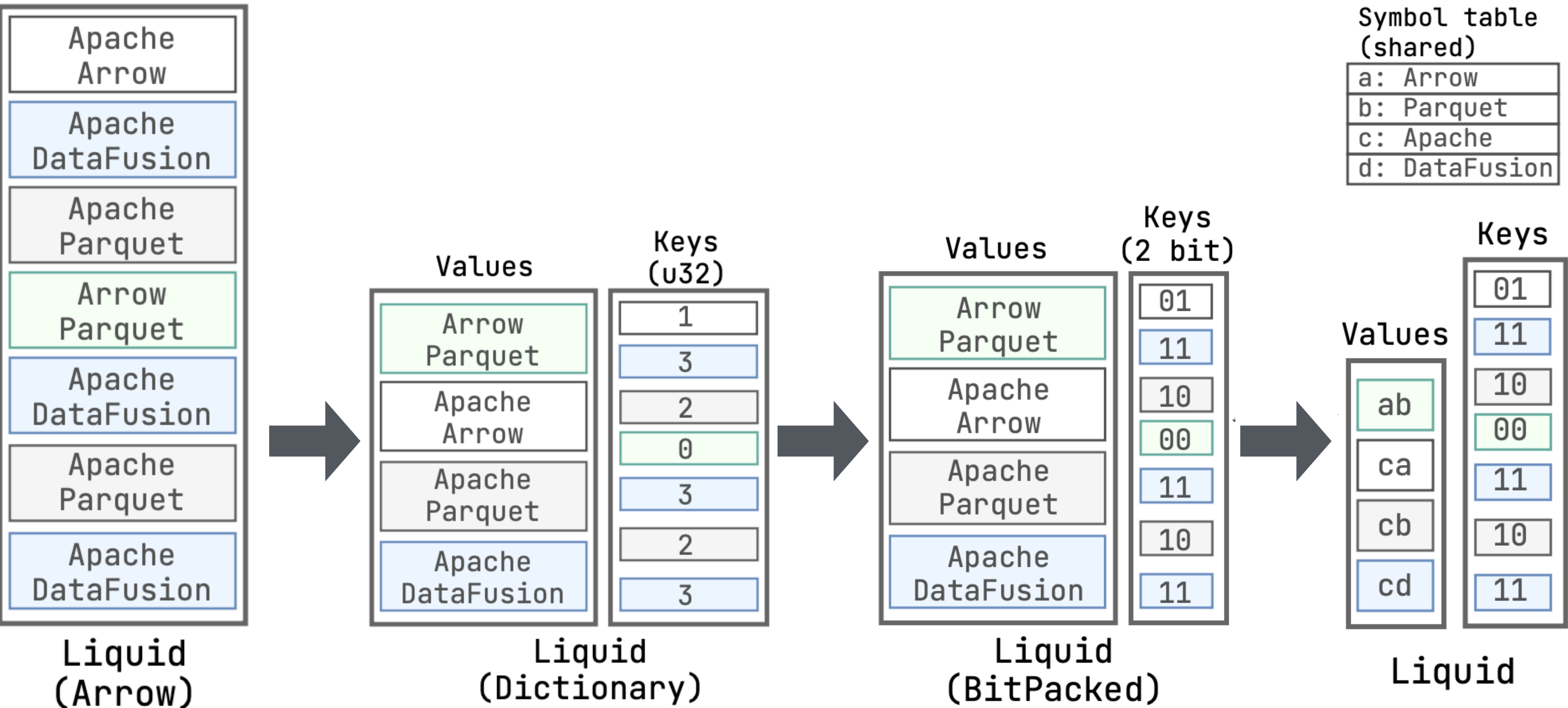progressive, selective, asynchronous transcoding

# Part 2:
# progressive, selective, asynchronous transcoding

## Yet another file format?

∞ **Nimble**  **LanceDB**

ⓢ **Vortex**  **ClickHouse**

**DuckDB**

## No!

Technical: they are not much different
Organizational: license, governance, ecosystem
LiquidCache: progressively bend the world

# Progressive transcoding

**Progressive**:
Transcode as needed,
no upfront cost

Parquet
subset

Partial
Liquid

Fully
Liquid

Parquet file

(Disk)

Object Store

100-200 ms

LiquidCache

1 ms

Compute

# Selective transcoding

# Asynchronous transcoding



location    date    val

**Selective**:
Transcode only touched data

100

Transcode

0

6 am    12 pm    6 pm    12 am    6 am

**Asynchronous**:
Transcode when less busy

With same memory: 10x lower latency

With same CPU: 10x lower CPU time

Decoding cost: close to theoretical optimal

Compression ratio: comparable to Parquet

Transcoding cost: negligible, no latency spike

# LiquidCache = compute + data

# Motivating example



FILO eviction

| Apache DataFusion |
| InfluxDB |
| Arrow Parquet |
| UW-Madison |

LLM-based eviction

| Apache DataFusion |
| InfluxDB |
| Arrow Parquet |
| UW-Madison |

Evict

| Apache DataFusion |
| InfluxDB |
| Arrow Parquet |
| UW-Madison |

Can we do better?

LRU eviction

| Apache DataFusion |
| InfluxDB |
| Arrow Parquet |
| UW-Madison |

Evicted
**Retained**

Count( name = "**Apache DataFusion**")

| | Cache miss |
|---|---|
| FILO | 2 |
| LRU | 2 |
| LLM-based | 2 |
| Theoretical opt. | 2 |

# Data-guided eviction

# Insight: evict partial data

BI  Report  ML  DS

**LiquidCache**

Shared storage

Some parts of data are
more important than others

| String prefix | Bit width | Dictionary |
| String length | Min/max/avg | Nullable mask |

r_name = 'EUROPE'

SearchPharse ≠ ''    p_brand = 'Brand#23'

MIN("URL")

MIN("TITLE")    **Real-world queries**

MobilePhoneModel ≠ ''    n_name = 'SAUDI ARABIA'

UserID = 43509093289640449    l_quantity < 24

# Evicts data, keep summaries

| | 0x98105 |
|---|---|
| | 0x53703 |
| | 0x94040 |
| | 0x15214 |

Evict →

**Conventional**

~~0x98105~~

**0x53703**

~~0x94040~~

~~0x15214~~

**LiquidEvict**

Bloom filter    ~~Evicted~~    **Retained**

Find "0x80309"?

| | Cache miss |
|---|---|
| Conventional | 3 |
| LiquidEvict | 0? |

# Real world example: StringView eviction

## Logical content

| |
|---|
| Apache DataFusion |
| InfluxDB |
| Arrow Rust Impl |
| Parquet pushdown |
| Apache DataFusion |

## Physical representation

| | | | |
|---|---|---|---|
| 17 | 0 | 0 | Apac |
| 8 | InfluxDB | | |
| 15 | 0 | 17 | Arro |
| 16 | 1 | 0 | Parq |
| 17 | 0 | 0 | Apac |

Views

String length

Buffer id

Buffer offset

Prefix

| A | p | a | c | h | e | | D |
|---|---|---|---|---|---|---|---|
| a | t | a | F | u | s | i | o |
| n | A | r | r | o | w | | R |
| u | s | t | | I | m | p | l |
| | | | | | | | |

Buffer 0

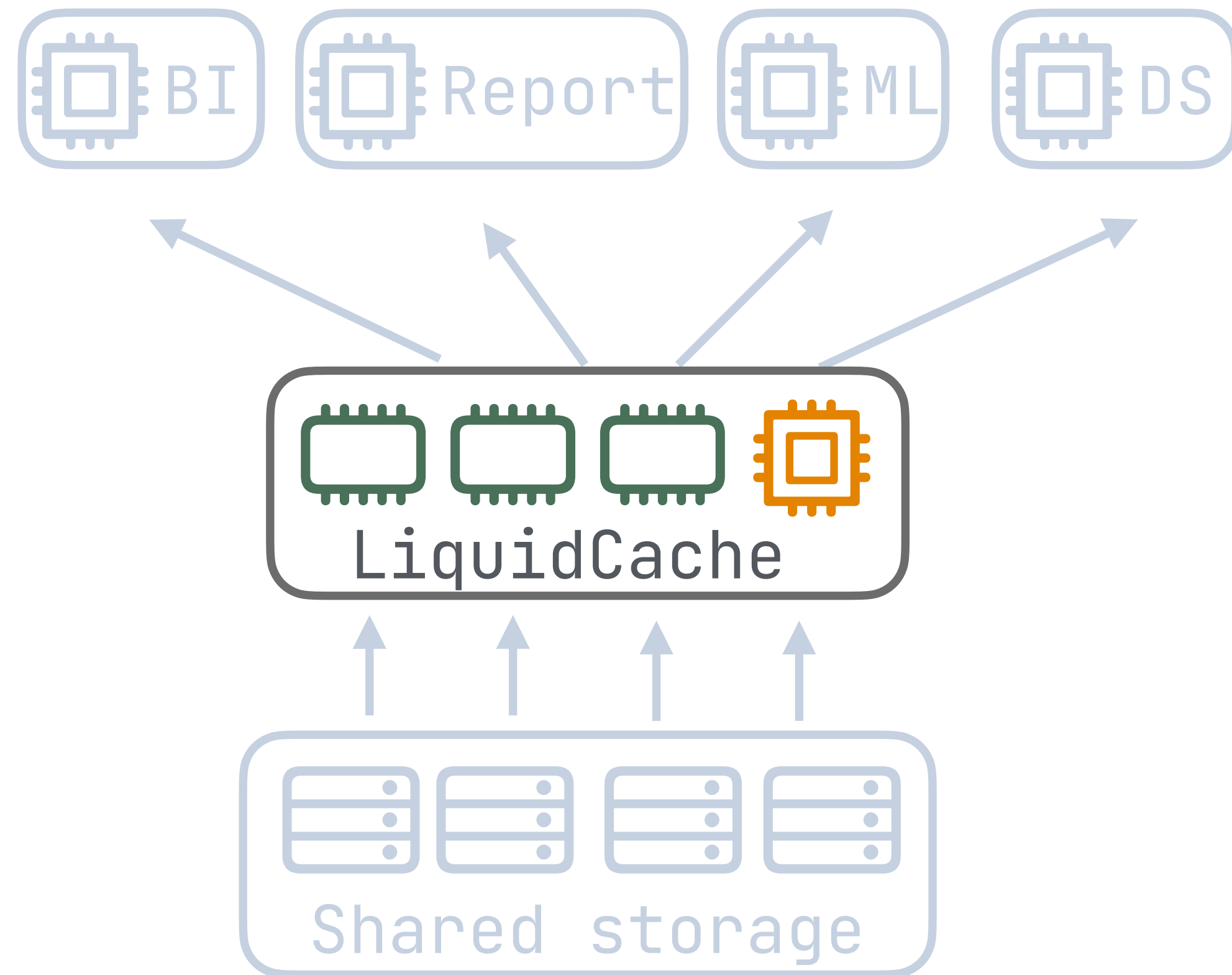| P | a | r | q | u | e | t |
|---|---|---|---|---|---|---|
| p | u | s | h | d | o | w | n |
| | | | | | | |
| | | | | | | |
| | | | | | | |

Buffer 1

Conventional: evict the entire array

LiquidEvict (structure-aware):
1. Evict buffers
2. Retain only prefix and str len
3. Retain only prefix

# LiquidCache + LiquidEvict



**Data-aware cache**
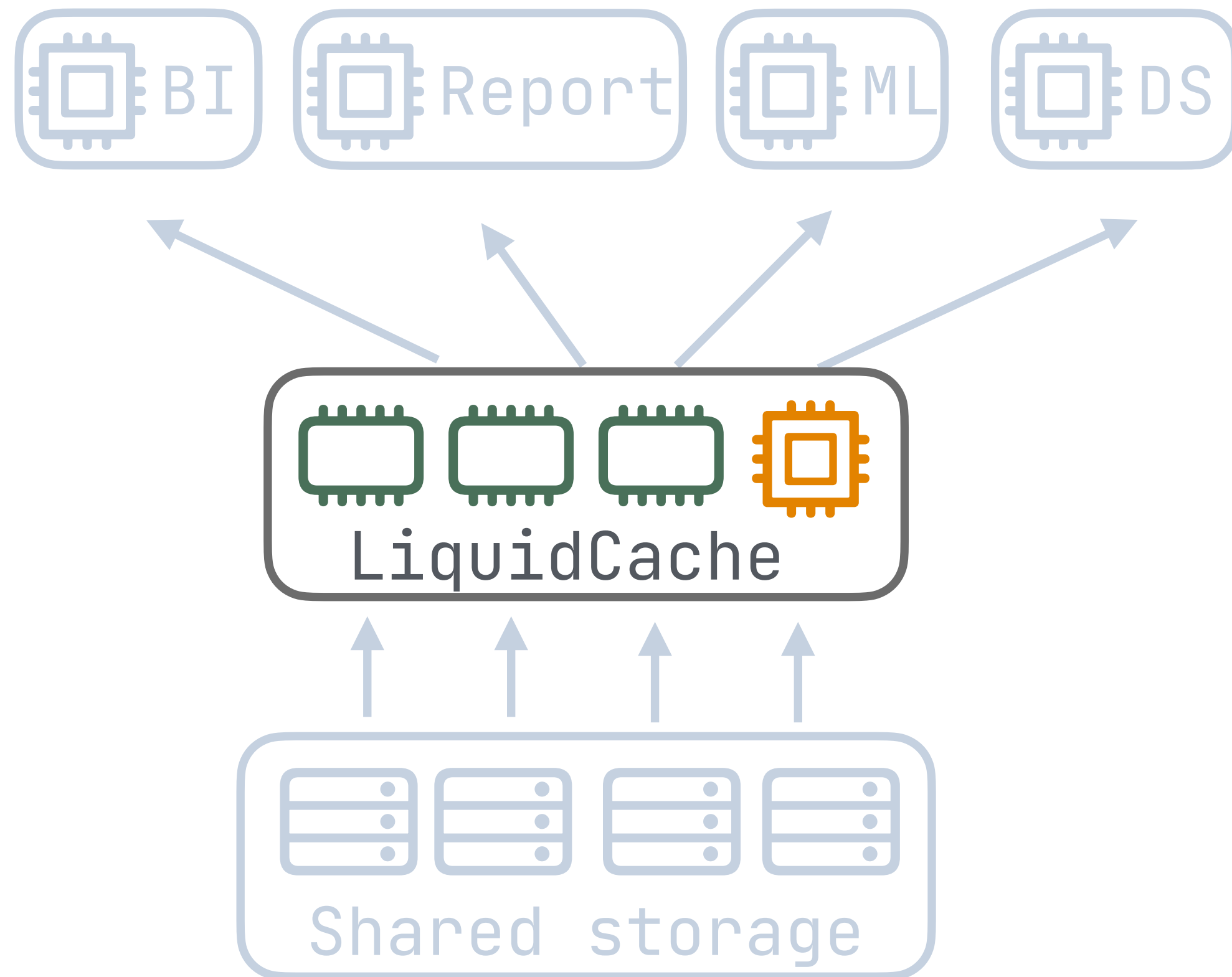
1. Pushdown to reduce traffic
2. Efficient decoding

**Data-aware eviction**

1. Evicts unimportant parts of data
2. Evicts data, keep summaries

**Shared, pushdown cache system**
10x lower latency, 10x lower cost

# Timeline



Jun. 1 - LiquidCache VLDB revision

July 1 - LiquidEvict prototype

Aug 1 - LiquidEvict bench & polish

Sep 1 - LiquidEvict VLDB submission

Dec 1 - LiquidEvict revision

Jan - Mar 2026 - Polish & present

May - Defense + graduate

Other first-author papers:

Bf-Tree: A Modern Read-Write-Optimized Concurrent Larger-Than-Memory Range Index.
Xiangpeng Hao, Badrish Chandramouli. (VLDB 2024)

Towards Buffer Management with Tiered Main Memory.
Xiangpeng Hao, Xinjing Zhou, Xiangyao Yu, Michael Stonebraker. (SIGMOD 2024)