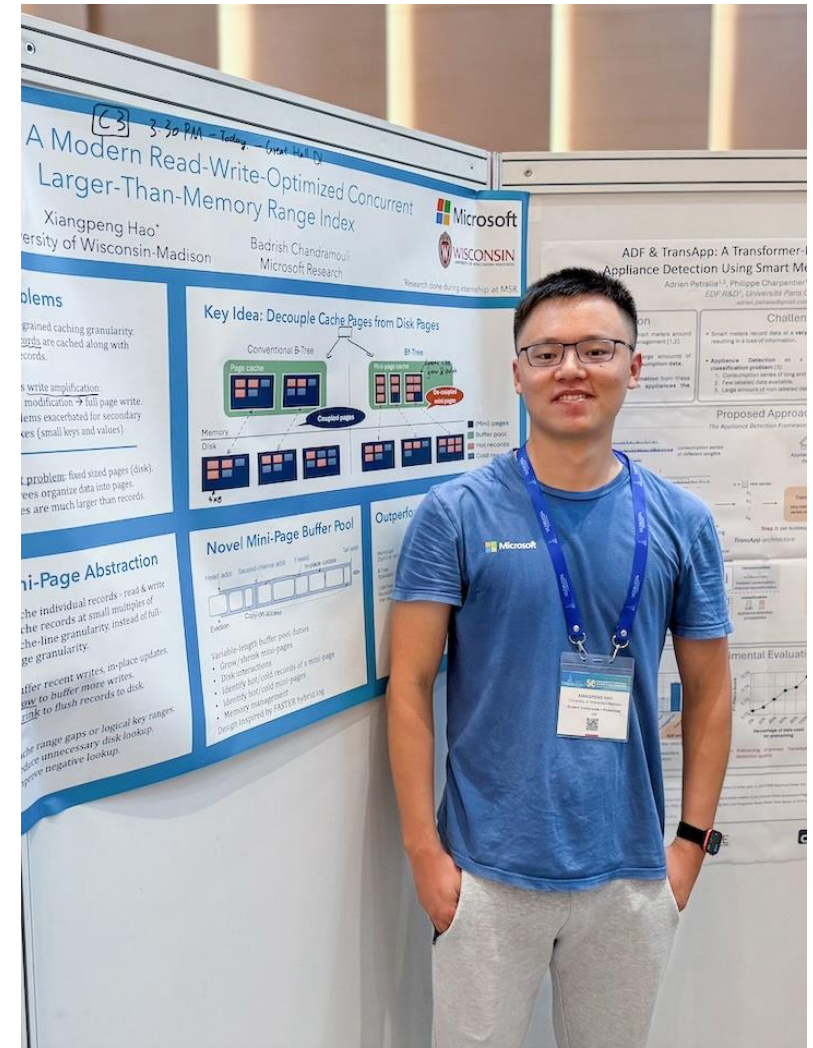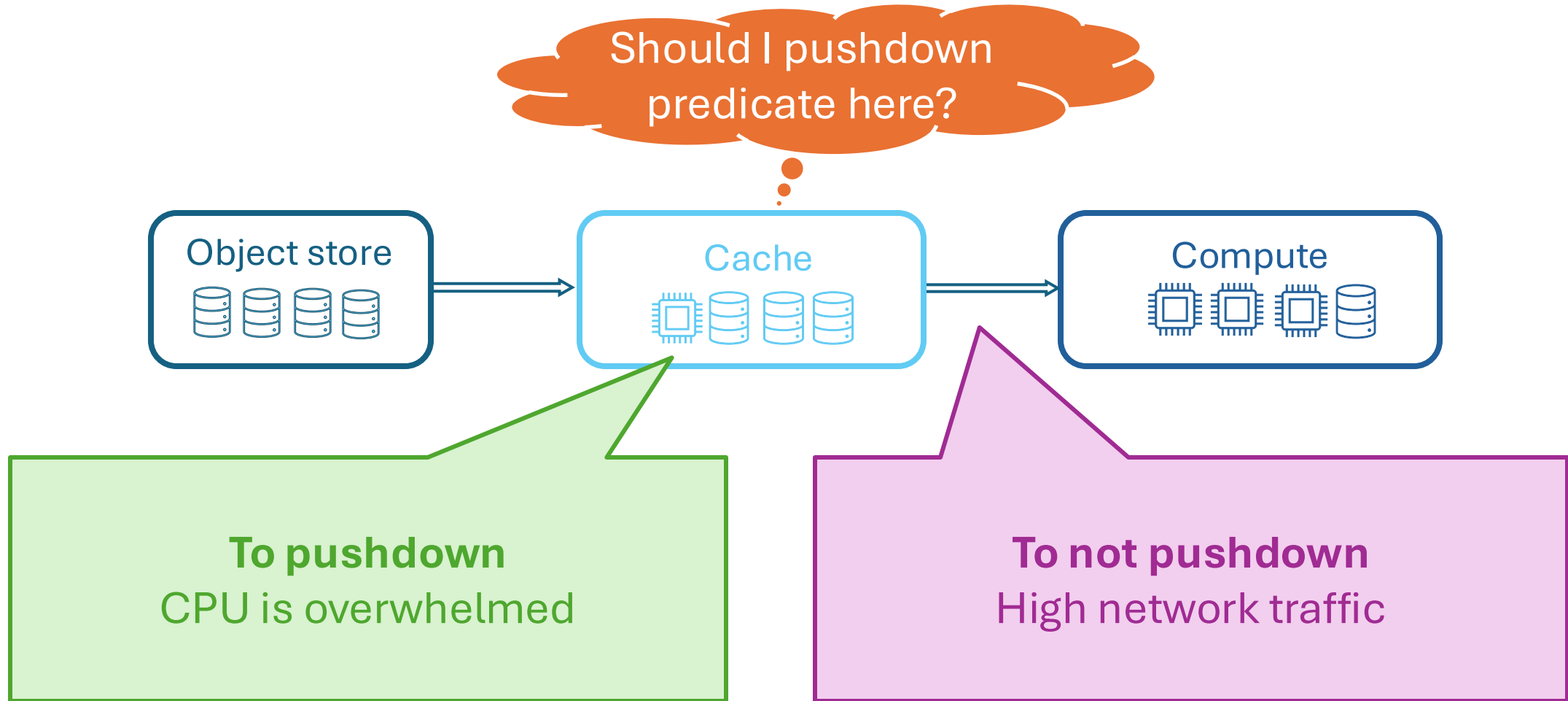# Practical Disaggregated Cache for Apache DataFusion
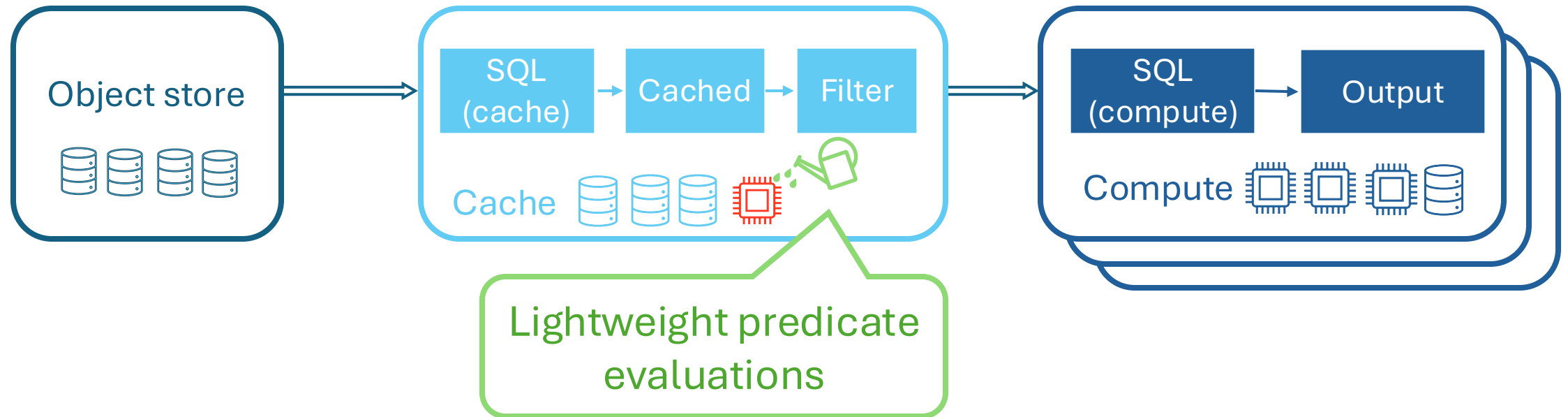
Xiangpeng Hao

# About Me

- 4th PhD@Wisconsin-Madison

- Study Database/Storage systems

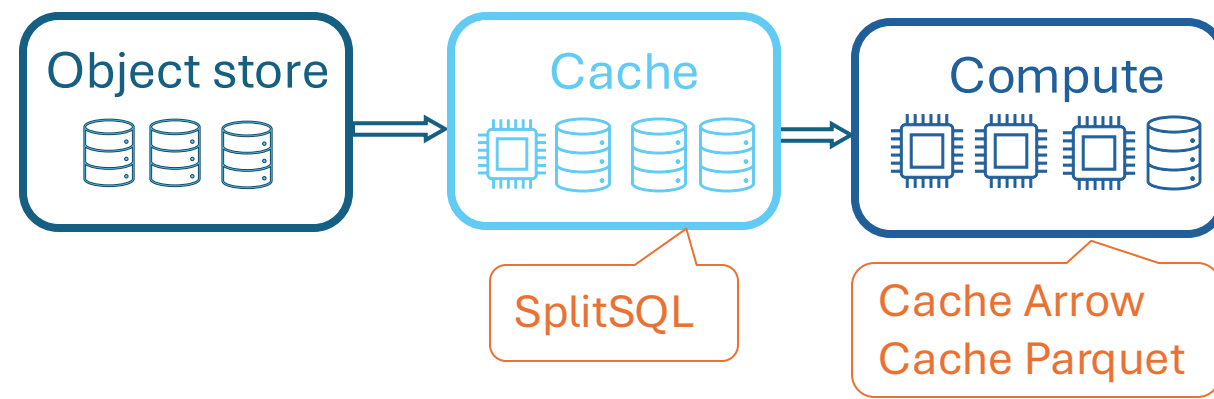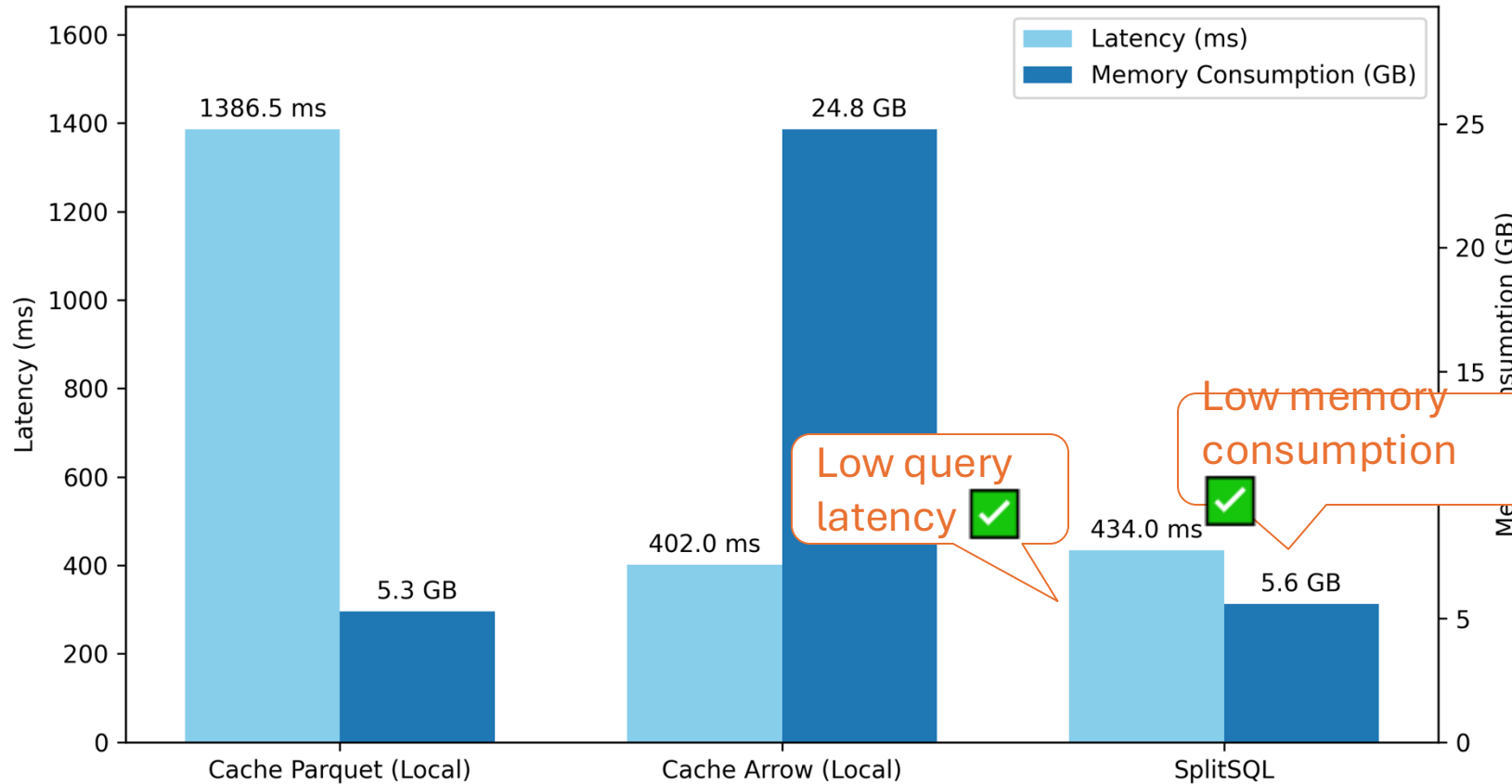- Build high performance/low level systems

# Disaggregated cache

# SplitSQL: Pushdown Done Right



Object store

SQL (cache) → Cached → Filter

Cache

Lightweight predicate evaluations

SQL (compute) → Output

Compute

# Result

Object store → Cache → Compute

SplitSQL (Cache)

Cache Arrow
Cache Parquet (Compute)



ClickBench Q22 (lower is better)

Latency (ms)
Memory Consumption (GB)

- Cache Parquet (Local): 1386.5 ms, 5.3 GB
- Cache Arrow (Local): 402.0 ms, 24.8 GB
- SplitSQL: 434.0 ms, 5.6 GB

Low query latency ✅

Low memory consumption ✅

```sql
SELECT
"SearchPhrase",
MIN("URL"),
MIN("Title"),
COUNT(*) AS c, COUNT(DISTINCT
"UserID")
FROM hits
WHERE
"Title" LIKE '%Google%'
AND
"URL" NOT LIKE '%.google.%' AND
"SearchPhrase" <> ''
GROUP BY "SearchPhrase" ORDER BY c
DESC LIMIT 10;
```
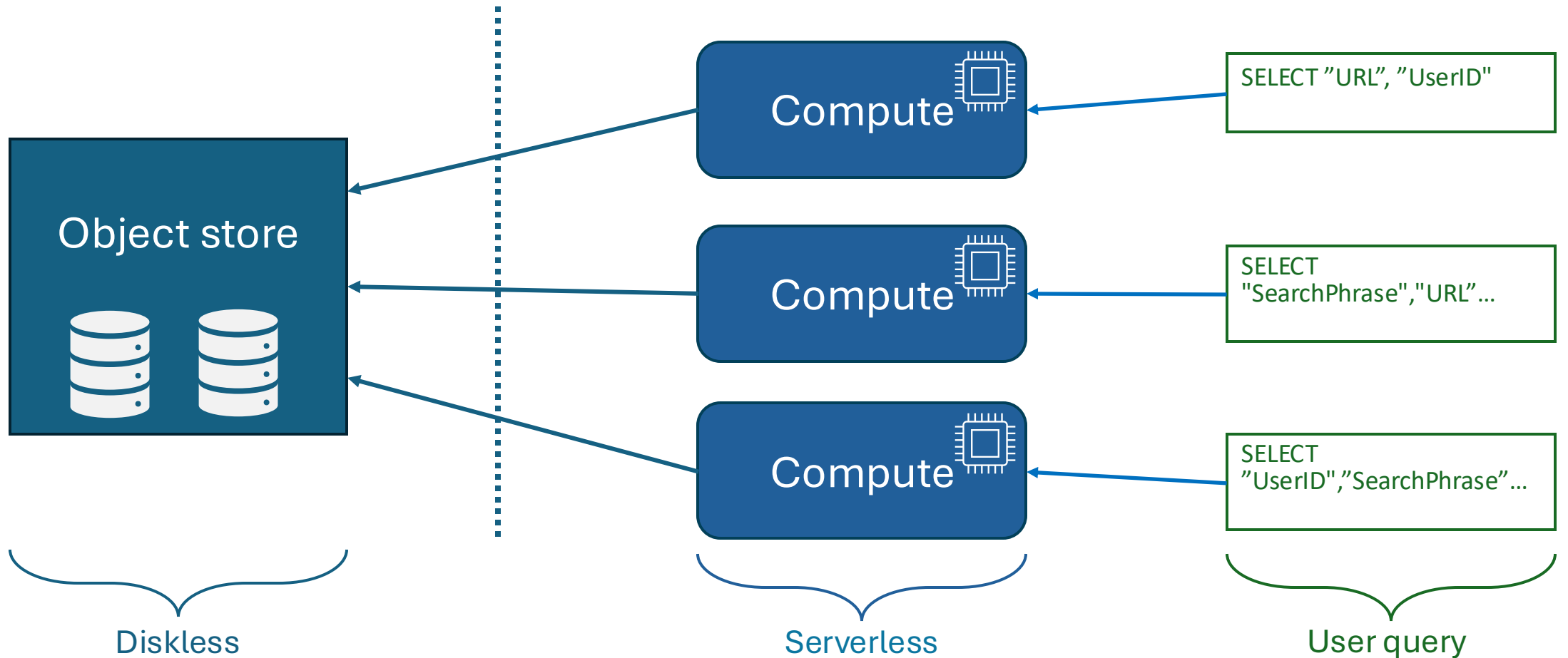
# Outline

Part 1: Disaggregated cache is the future

Part 2: To Pushdown or not to pushdown?

Part 3: SplitSQL: pushdown down right

Part 4: Evaluations

# Data lake architecture



Object store

Compute

Compute

Compute

SELECT "URL", "UserID"

SELECT "SearchPhrase","URL"...

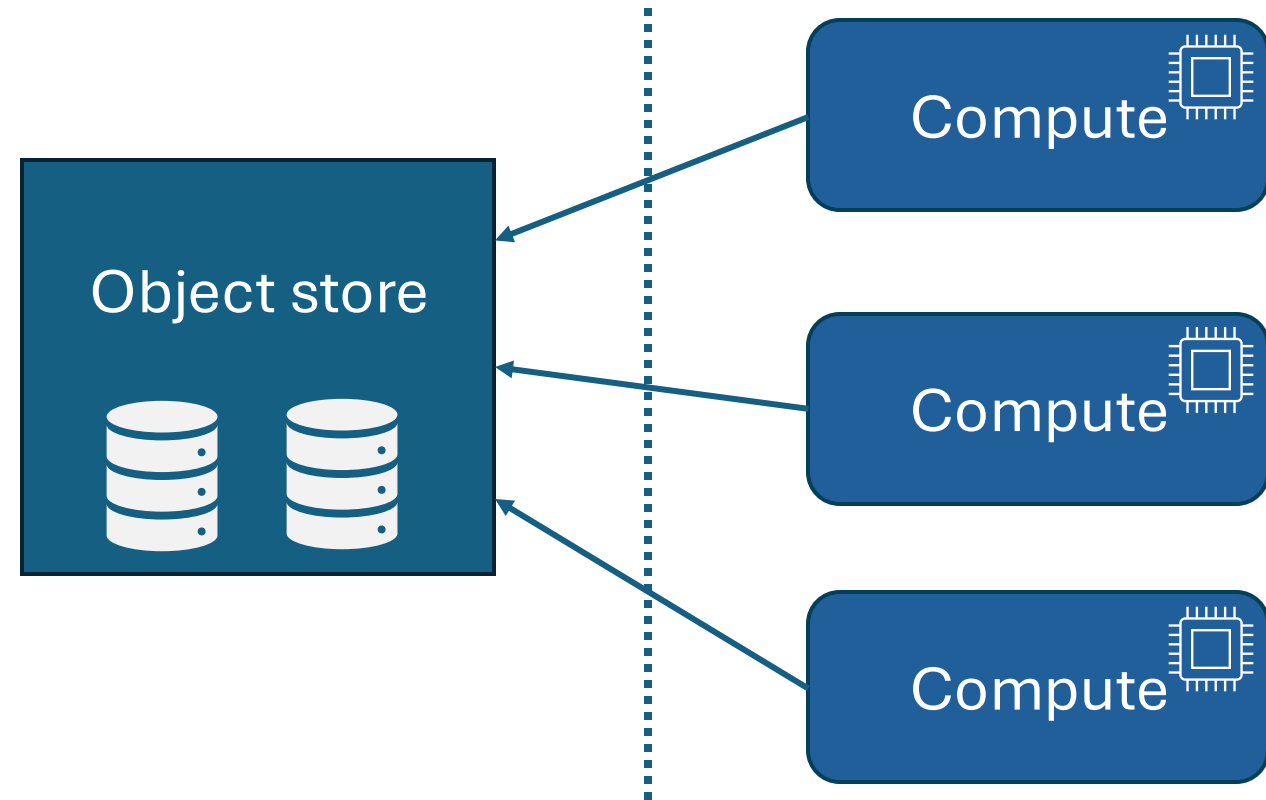SELECT "UserID","SearchPhrase"...

Diskless

Serverless

User query

# Modern architecture

Object store is new **disk**
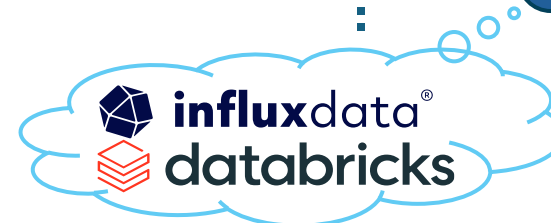
Lambda/EC2 is new **CPU**

Where to **cache** data?
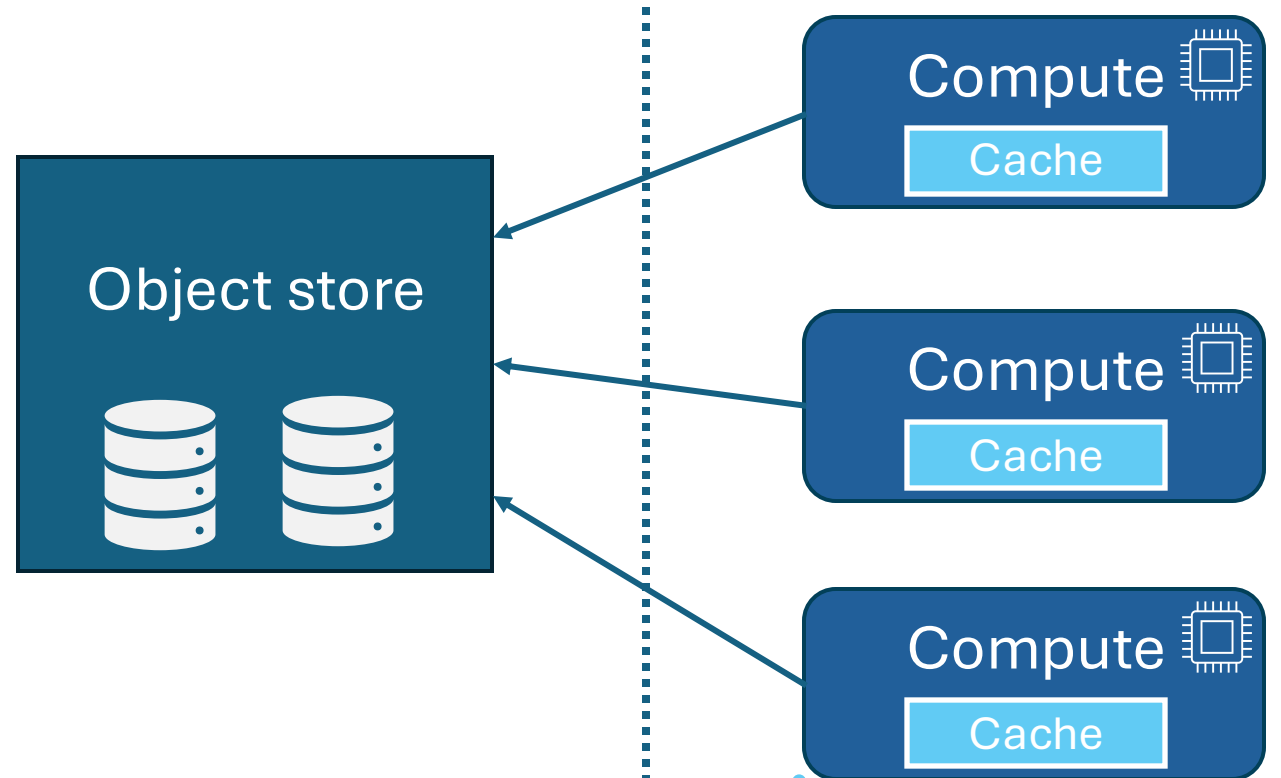
# Option 1: private cache

Simple ✅

Duplicate copies ❌

Can't scale independently ❌

Low resource utilization ❌

Object store

Compute
Cache

Compute
Cache

Compute
Cache

influxdata®
databricks

# Option 2: distributed cache

Complex design ❌

Single copy ✅

Can't scale independently ❌

High resource utilization ✅

Object store

Compute

Compute

Compute

Cache

Snowflake
ALLUXIO

# Option 3: disaggregated cache

Clean architecture 🎯

Single copy ✅

Scale independently ✅

High resource utilization ✅

Object store

Cache

Compute

Compute

Compute

Why not popular?

# Outline

Part 1: Disaggregated cache is the future

**Part 2: To Pushdown or not to pushdown?**

Part 3: SplitSQL: pushdown down right

Part 4: Evaluations

# Disaggregate cache cause high network traffic



Large files send over network!

# Queries filtered out most data

# Evaluate predicates on cache

# This is predicate push down **to cache**

# Outline

Part 1: Disaggregated cache is the future

Part 2: To Pushdown or not to pushdown?

Part 3: SplitSQL: pushdown down right

Part 4: Evaluations

# SplitSQL: Pushdown Done Right

Object store

SQL (cache) → Cached → Filter

Cache

SQL (compute) → Output

Compute

Scale cache:
Provision larger memory,
attach larger elastic storage

Scale compute:
More compute nodes,
same cache node

# Example



How not to overwhelm my CPU?

SELECT DISTINCT "URL"
FROM hits WHERE
"URL" LIKE '%google%'

SELECT "URL"
FROM hits WHERE
"URL" LIKE '%google%'

SplitSQL

SELECT DISTINCT "URL"
FROM cache

Object store

SQL (cache) → Cached → Filter

Cache

SQL (compute) → Output

Compute

# Simple filters can be expensive to evaluate

**Decoding overhead (78%)**

**Useful work (22%)**

Flame Graph

Reset Zoom                                                                                                Search

Decompress Parquet

Parquet -> Arrow

Evaluate filters

snap::decompress::Decoder::decompress

..ssion::snappy_codec::SnappyCodec as parquet::compression::Codec>::decompress

..r::converts::from_utf8    ..y_push    ..    parquet::file::serialized_reader::decode_page

..array::ByteArrayDecoderPlain::read    ..    ..::SerializedPageReader<R> as parquet::column::page::PageReader>::get_next_page    ..air::Finder::find_impl    ..

..r::byte_array::ByteArrayDecoder::read    parquet::column::reader::GenericColumnReader<R,D,V>::read_new_page    ..n::BooleanBuffer::collect_bool

parquet::column::reader::GenericColumnReader<R,D,V>::read_records    ..rray::BooleanArray::from_unary

parquet::arrow::record_reader::GenericRecordReader<V,CV>::read_records    ..ate::Predicate::evaluate_array

<parquet::arrow::array_reader::byte_array::ByteArrayReader<I> as parquet::arrow::array_reader::ArrayReader>::read_records    arrow_string::like::op_scalar

<parquet::arrow::array_reader::struct_array::StructArrayReader as parquet::arrow::array_reader::ArrayReader>::read_records    arrow_string::like::like_op

<parquet::arrow::arrow_reader::ParquetRecordBatchReader as core::iter::traits::iterator::Iterator>::next    arrow_string::like::like

<S as futures_core::stream::TryStream>::try_poll_next    ..l_expr_common::datum::apply_cmp

<futures_util::stream::stream::map::Map<St,F> as futures_core::stream::Stream>::poll_next    ..l_expr::PhysicalExpr>::evaluate

<futures_util::stream::stream::map::Map<St,F> as futures_core::stream::Stream>::poll_next    ..lan::filter::filter_and_project

<datafusion::datasource::physical_plan::file_stream::FileStream<F> as futures_core::stream::Stream>::poll_next

<datafusion_physical_plan::filter::FilterExec_Stream as futures_core::stream::Stream>::poll_next

<datafusion_physical_plan::coalesce_batches::CoalesceBatchesStream as futures_core::stream::Stream>::poll_next

# Predicate push down doesn't like Parquet

**Parquet** is the industry standard
- Rich features, great ecosystem
- Battle tested
- High compression ratio
- De facto file format for big data

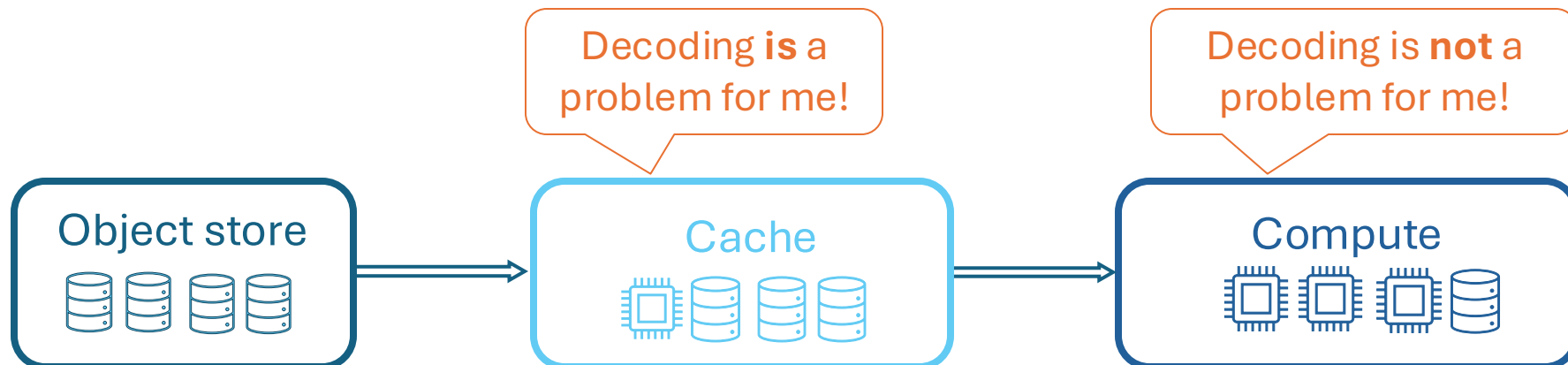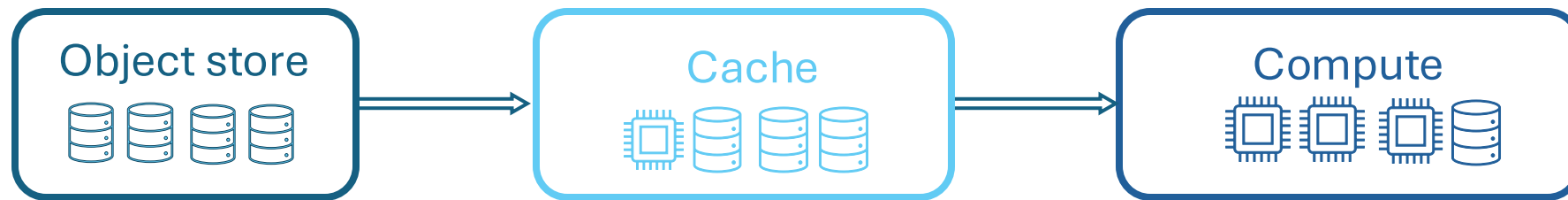Decoding Parquet is CPU-intensive
- Decompression
- Decoding metadata
- Decoding data

# Faster decoding is All-You-Need

Option 1: switch to a different file format
- Small win, big lose – lose all other nice Parquet features
- Significant changes to the ecosystem
- Slow adoption (e.g., >10 years)



Object store → Cache → Compute

Start with a different format!

Nimble
Vortex    DuckDB
LanceDB

# Faster decoding is All-You-Need

**Option 1: switch to a customized file format**
- Small win, big lose – lose all other nice Parquet features
- Significant changes to the ecosystem
- Slow adoption (e.g., >10 years)

❌

**Option 2: cache decoded values (e.g., cache Arrow)**

❓

Cache Arrow here!

Object store → Cache → Compute

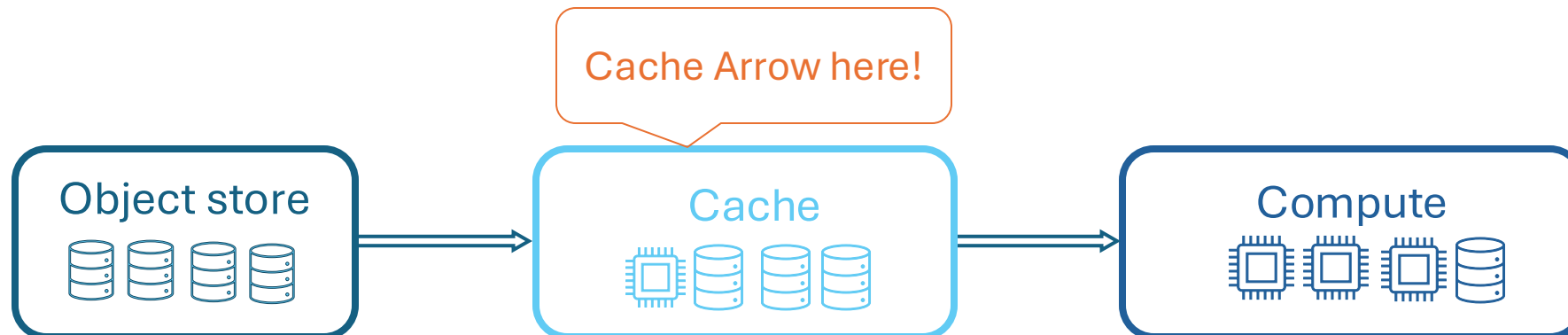# Cache Arrow speeds up some queries



ClickBench

# Faster decoding is All-You-Need

**Option 1: switch to a customized file format**
- Small win, big lose – lose all other nice Parquet features
- Significant changes to the ecosystem
- Slow adoption (e.g., >10 years)
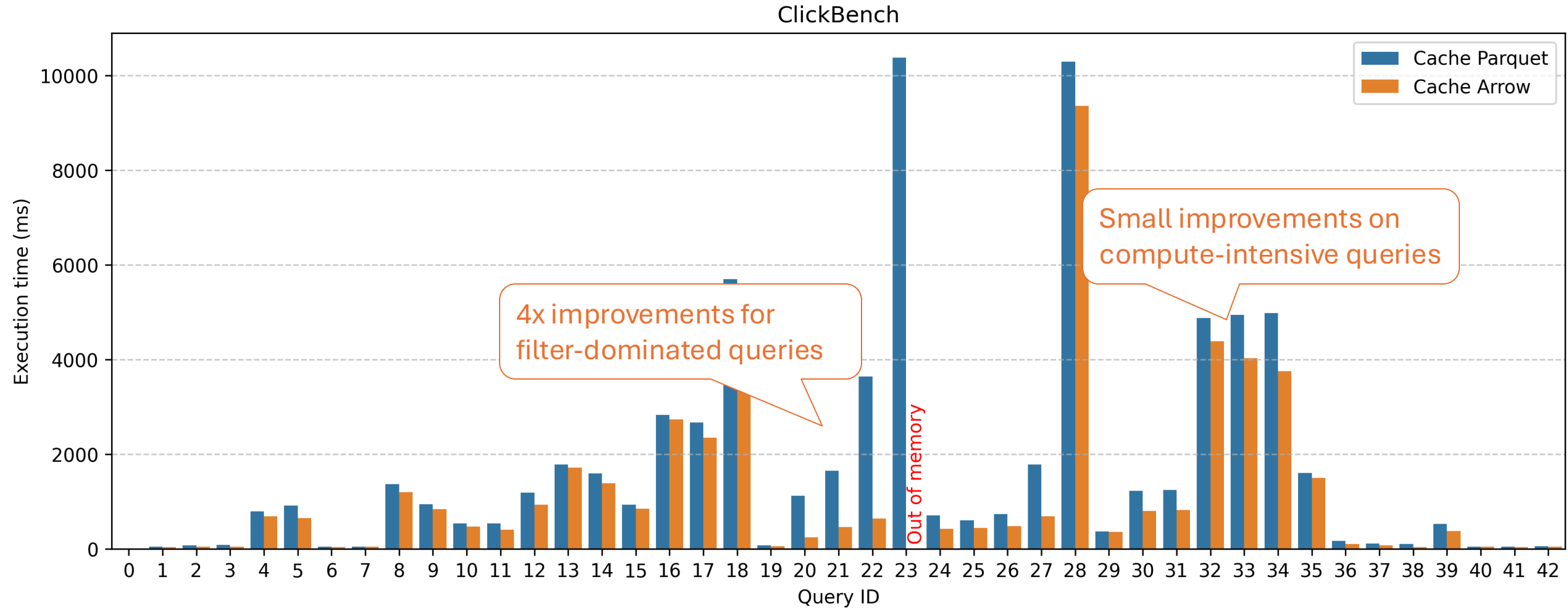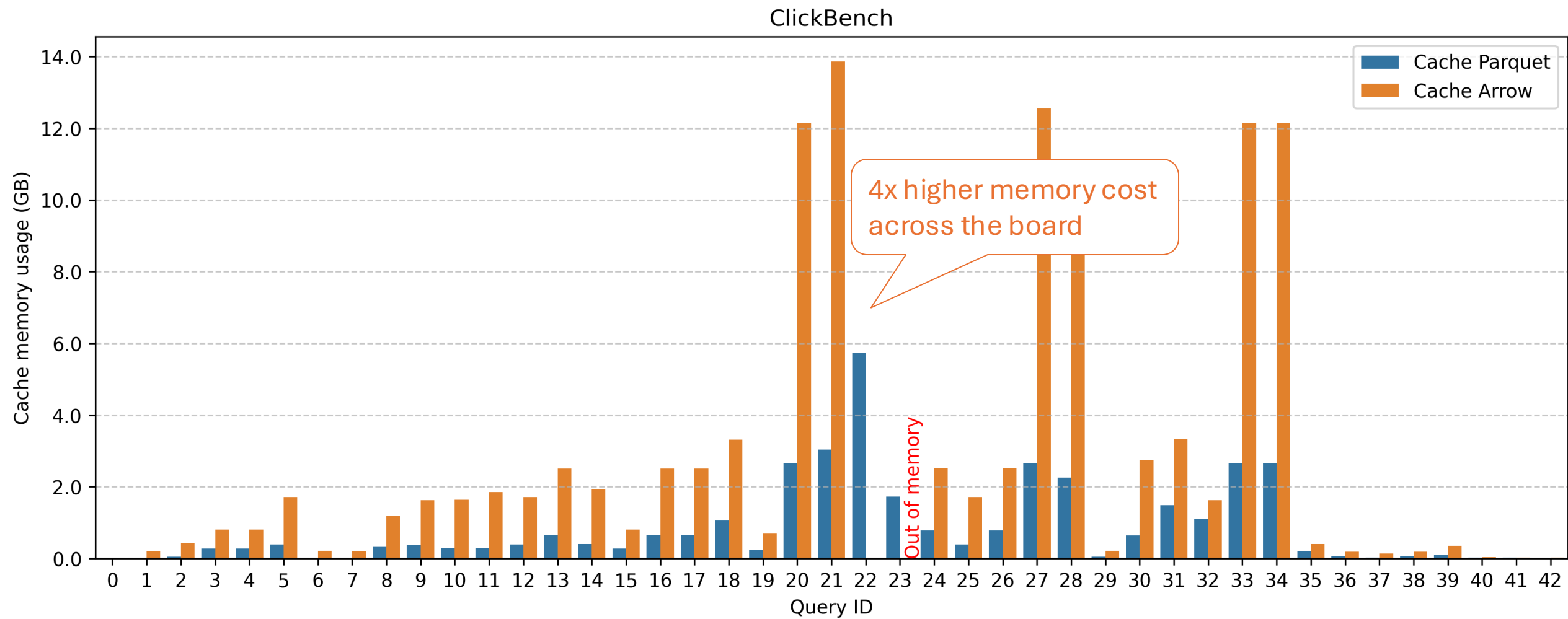
☒

**Option 2: cache decoded values (e.g., cache Arrow)**
- Performance improvement varies
- 4x more memory usage

☒

**We need: cache-specific file format**
- Transparent – what happens in cache, stays in cache
- Unlocks new opportunities

☑

# Cache-specific format?

- Leverage modern encoding algorithms
  - SIMD friendly
  - Fine-grained decoding – decode only relevant data
  - Evaluate predicates on encoded data
- Make modern trade-offs
  - IO time vs decode time
  - Compression ratio vs access cost

Cache-specific format here!

Object store

Cache

Compute

# Example: Encode String arrays



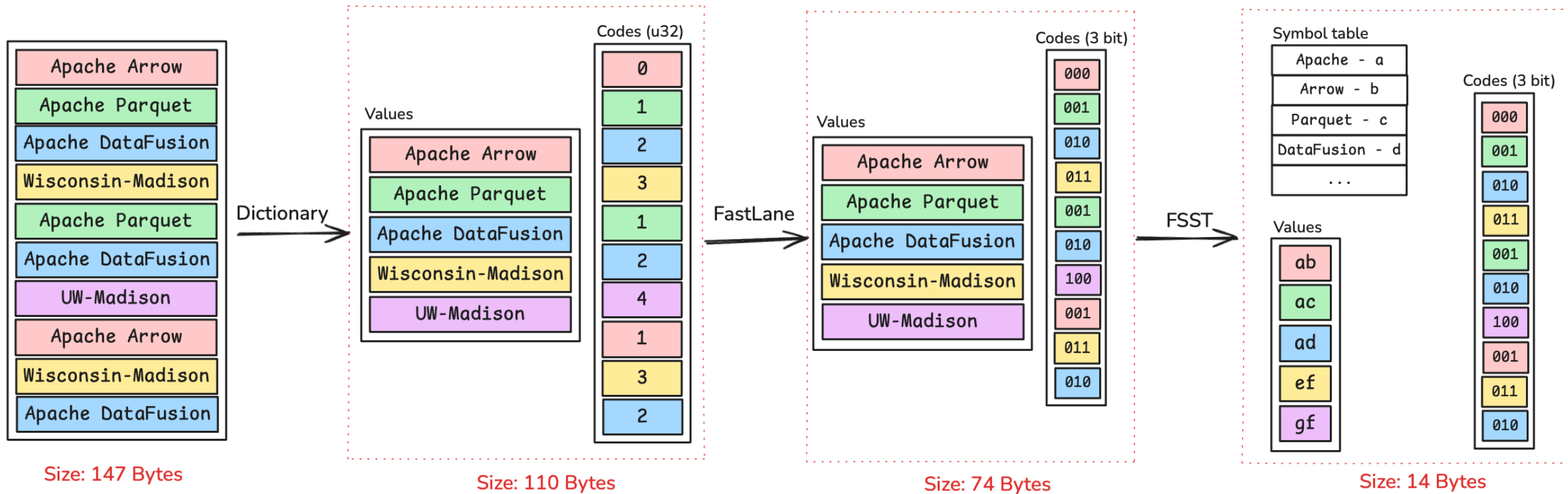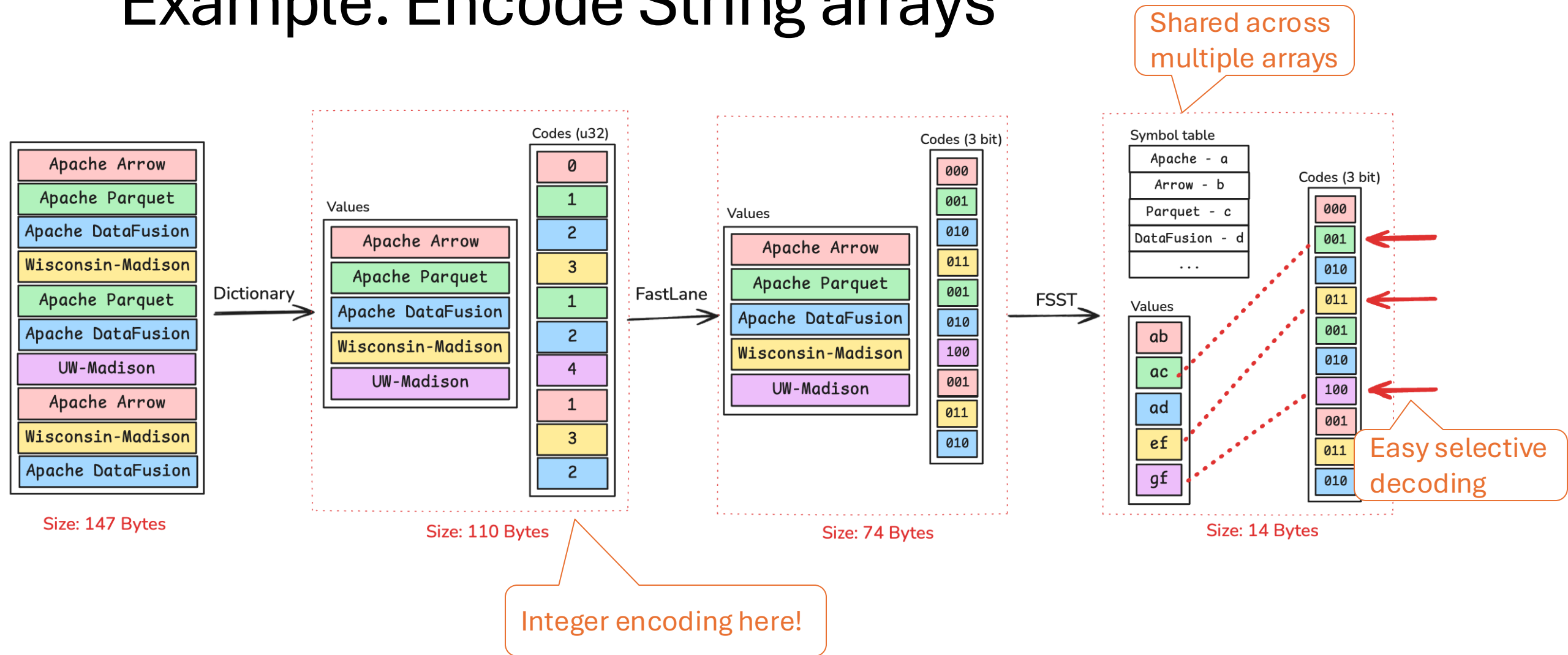| | | | |
|---|---|---|---|
| Apache Arrow | | | |
| Apache Parquet | | | |
| Apache DataFusion | | | |
| Wisconsin-Madison | | | |
| Apache Parquet | | | |
| Apache DataFusion | | | |
| UW-Madison | | | |
| Apache Arrow | | | |
| Wisconsin-Madison | | | |
| Apache DataFusion | | | |

Size: 147 Bytes

**Dictionary**

**Values**

| | Codes (u32) |
|---|---|
| Apache Arrow | 0 |
| Apache Parquet | 1 |
| Apache DataFusion | 2 |
| Wisconsin-Madison | 3 |
| | 1 |
| | 2 |
| UW-Madison | 4 |
| | 1 |
| | 3 |
| | 2 |

Size: 110 Bytes

**FastLane**

**Values**

| | Codes (3 bit) |
|---|---|
| Apache Arrow | 000 |
| Apache Parquet | 001 |
| Apache DataFusion | 010 |
| Wisconsin-Madison | 011 |
| | 001 |
| | 010 |
| UW-Madison | 100 |
| | 001 |
| | 011 |
| | 010 |

Size: 74 Bytes

**FSST**

**Symbol table**

| |
|---|
| Apache - a |
| Arrow - b |
| Parquet - c |
| DataFusion - d |
| ... |

**Values**

| | Codes (3 bit) |
|---|---|
| ab | 000 |
| ac | 001 |
| ad | 010 |
| ef | 011 |
| gf | 001 |
| | 010 |
| | 100 |
| | 001 |
| | 011 |
| | 010 |

Size: 14 Bytes

# Example: Encode String arrays



Size: 147 Bytes

Dictionary

Values

Codes (u32)

Size: 110 Bytes

Integer encoding here!

FastLane

Values

Codes (3 bit)

Size: 74 Bytes

FSST

Shared across multiple arrays

Symbol table

Apache - a
Arrow - b
Parquet - c
DataFusion - d
...

Values

Codes (3 bit)

Easy selective decoding

Size: 14 Bytes

# Evaluate predicates on encoded data



Size: 147 Bytes

Size: 110 Bytes

Size: 74 Bytes

Size: 14 Bytes

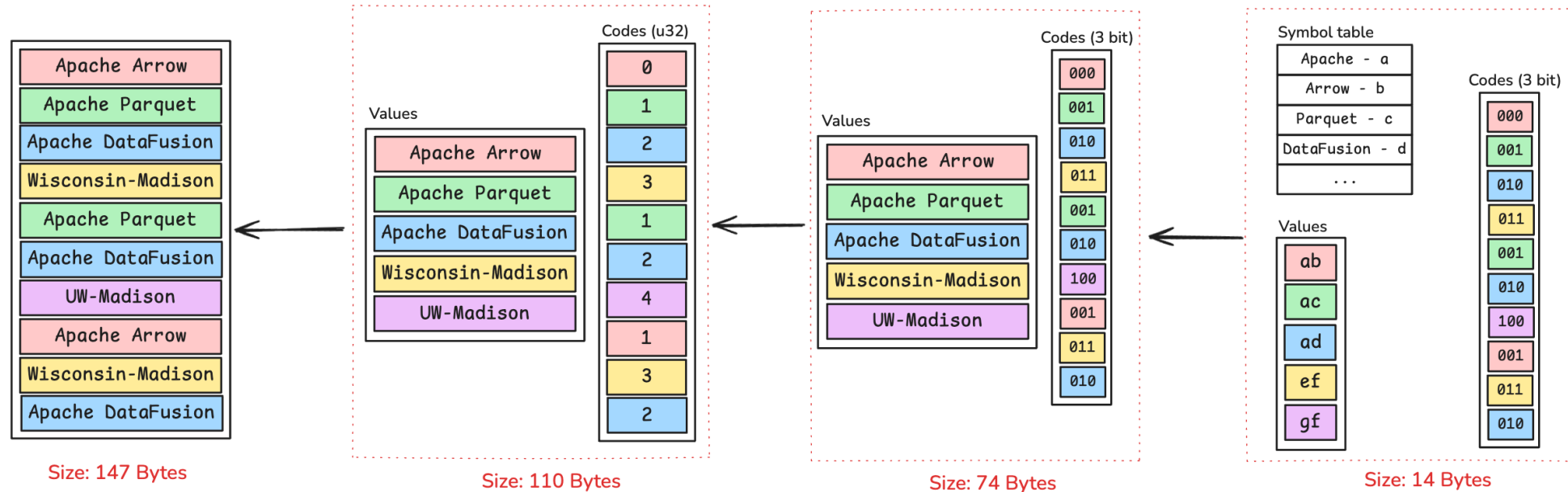How find strings that contains "DataFusion"?

Find "Apache DataFusion"

1. Encode the needle – "ad"

2. Find the index in values – 2

2. Return the index in codes

# Evaluate predicates on **partially** encoded data



Size: 147 Bytes

Size: 110 Bytes

Size: 74 Bytes

Size: 14 Bytes

Apply filters along the decoding path
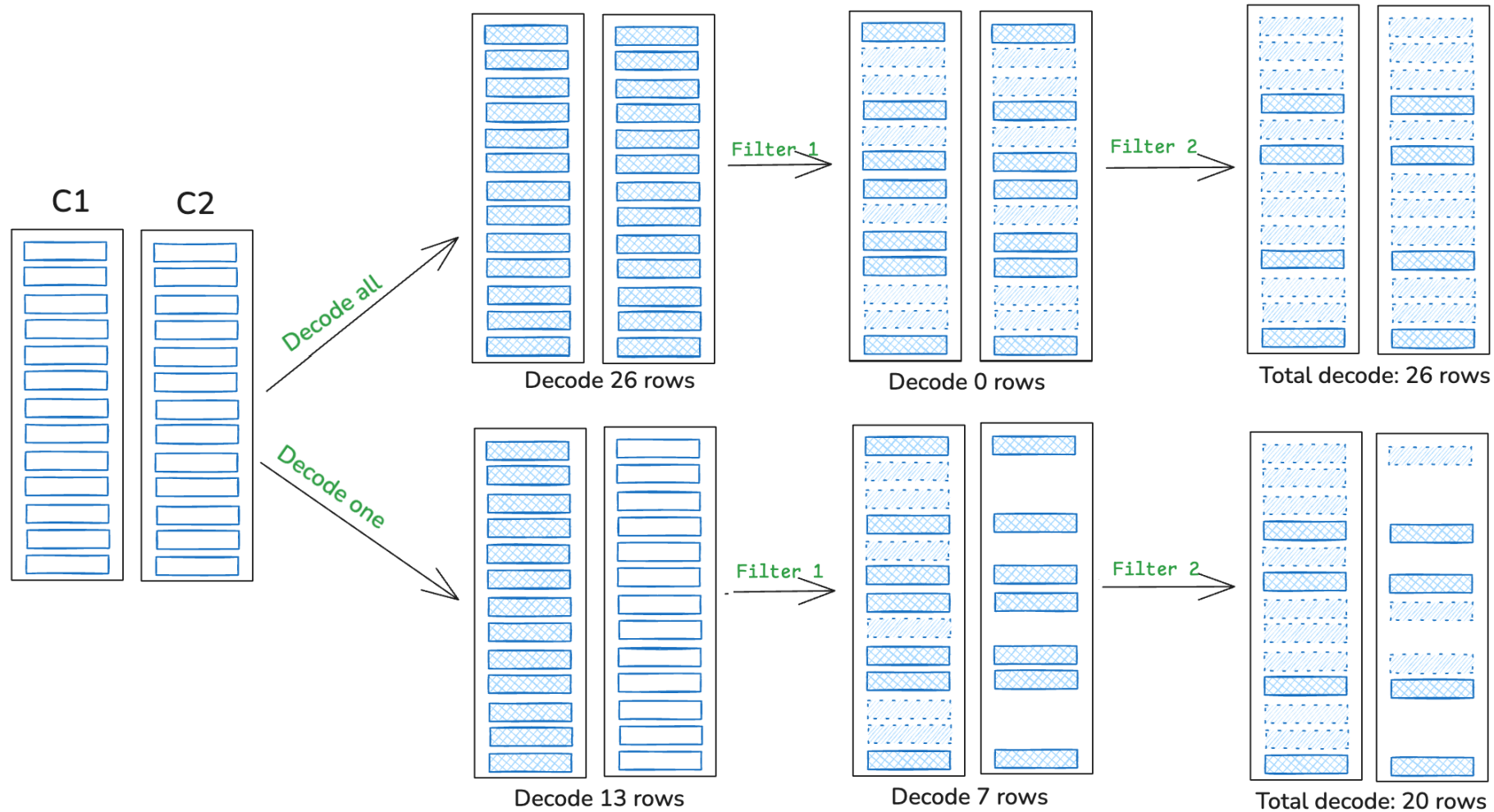
Find contains "DataFusion"

1. Decode the values

2. Find the index in values – 2

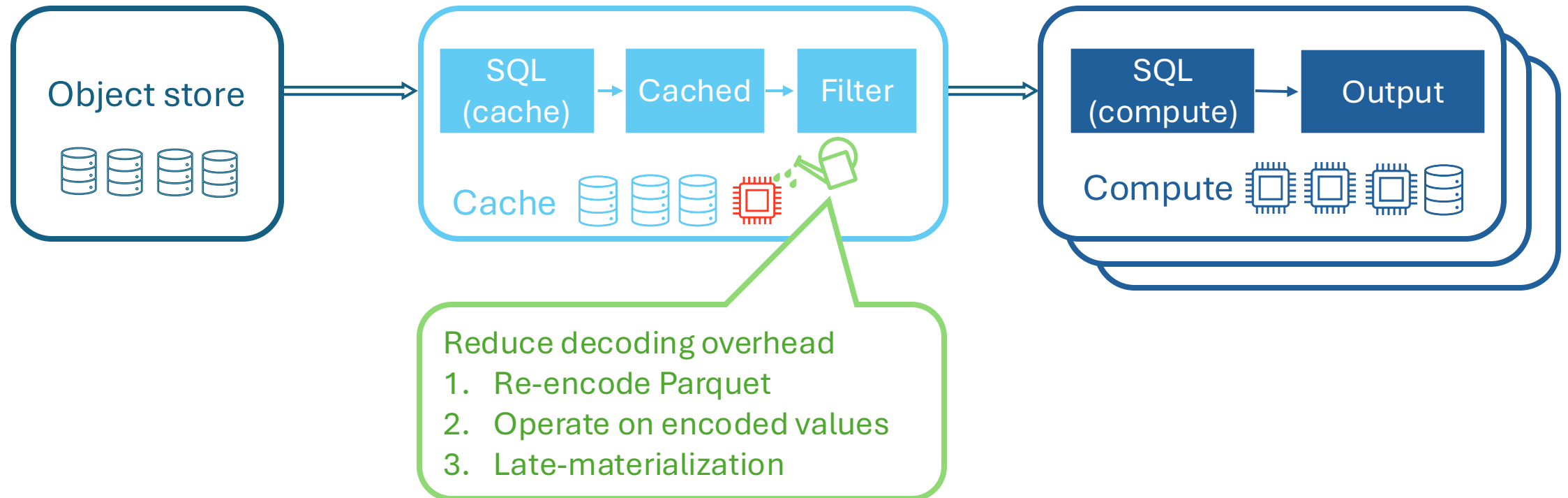2. Return the index in codes

Partial decoding

# Random access for late-materialization



C1    C2

Decode all

Decode 26 rows → Filter 1 → Decode 0 rows → Filter 2 → Total decode: 26 rows

Decode one

Decode 13 rows → Filter 1 → Decode 7 rows → Filter 2 → Total decode: 20 rows

Decode only the necessary
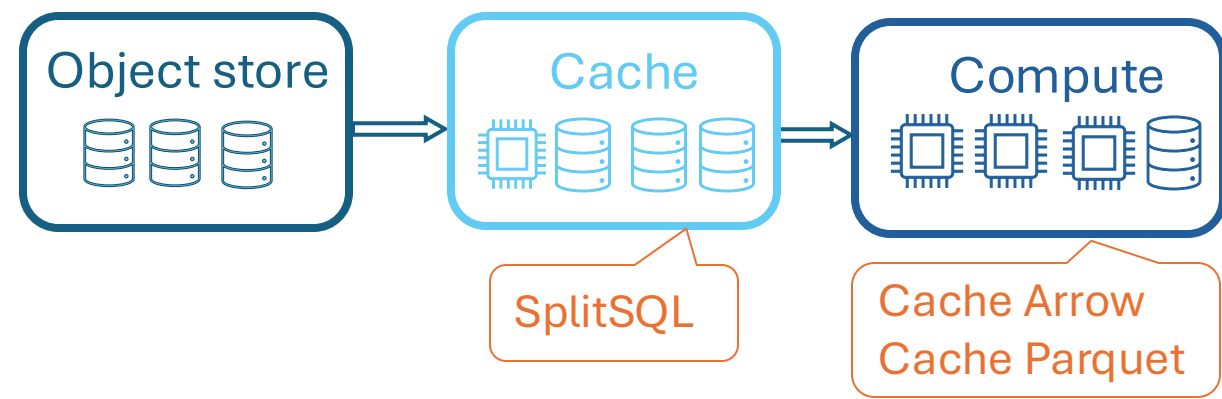
# SplitSQL: Practical Disaggregated Cache

# Outline

Part 1: Disaggregated cache is the future

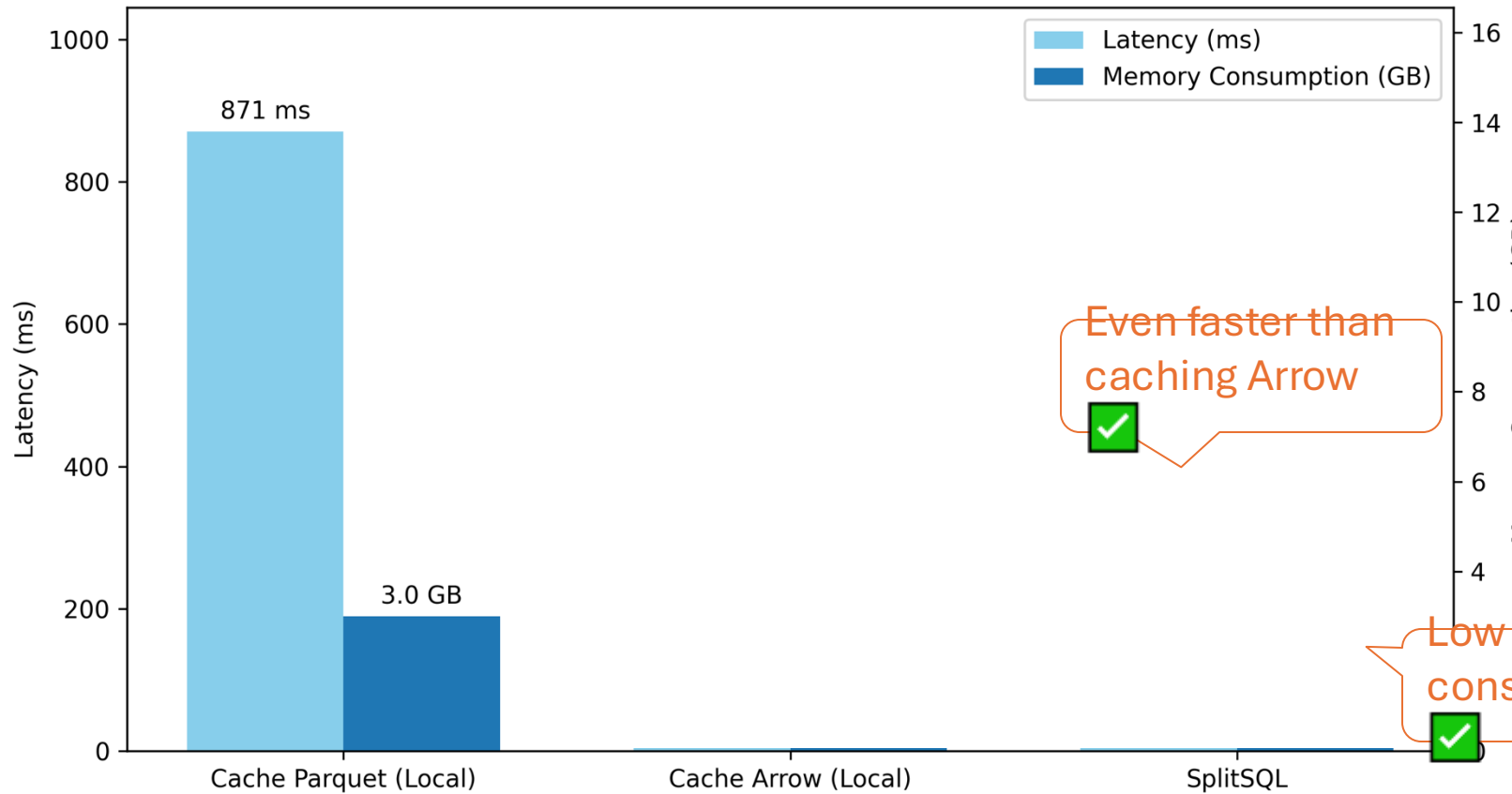Part 2: To Pushdown or not to pushdown?

Part 3: SplitSQL: pushdown down right

Part 4: Evaluations
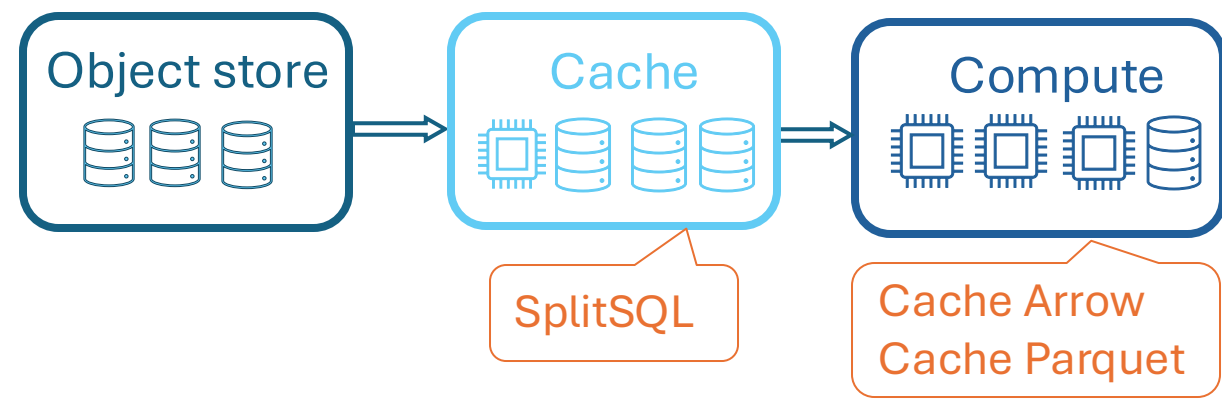
# Evaluations – Q21



ClickBench Q21 (lower is better)

Chart legend:
- Latency (ms)
- Memory Consumption (GB)

Cache Parquet (Local): 871 ms, 3.0 GB

Even faster than caching Arrow ✅

Low memory consumption ✅

X-axis: Cache Parquet (Local), Cache Arrow (Local), SplitSQL

Diagram boxes: Object store → Cache → Compute

SplitSQL (pointing to Cache)
Cache Arrow / Cache Parquet (pointing to Compute)
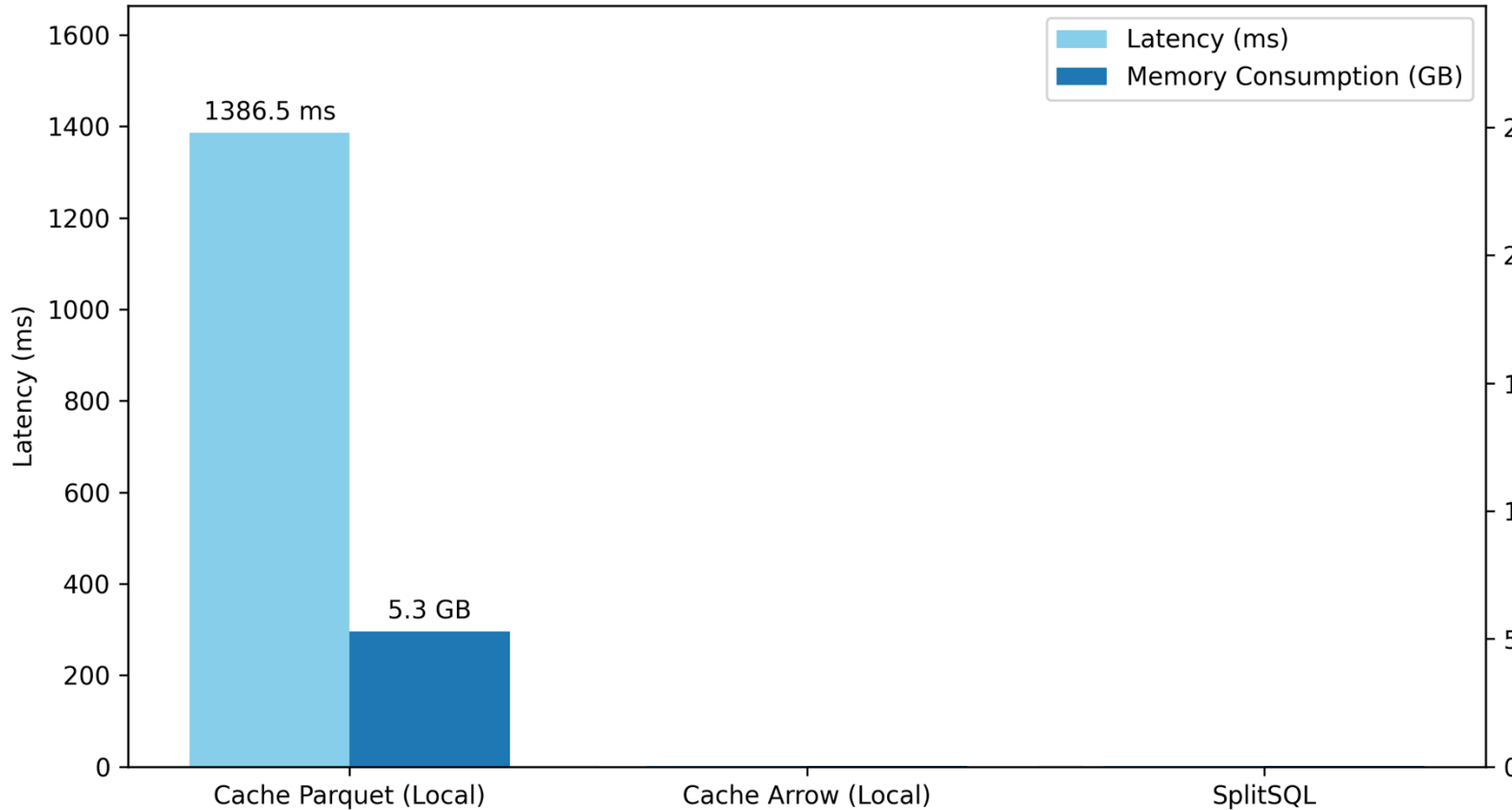
```sql
SELECT
"SearchPhrase", MIN("URL"),
COUNT(*) AS c
FROM hits
WHERE "URL" LIKE '%google%'
AND "SearchPhrase" <> ''
GROUP BY "SearchPhrase"
ORDER BY c DESC LIMIT 10;
```
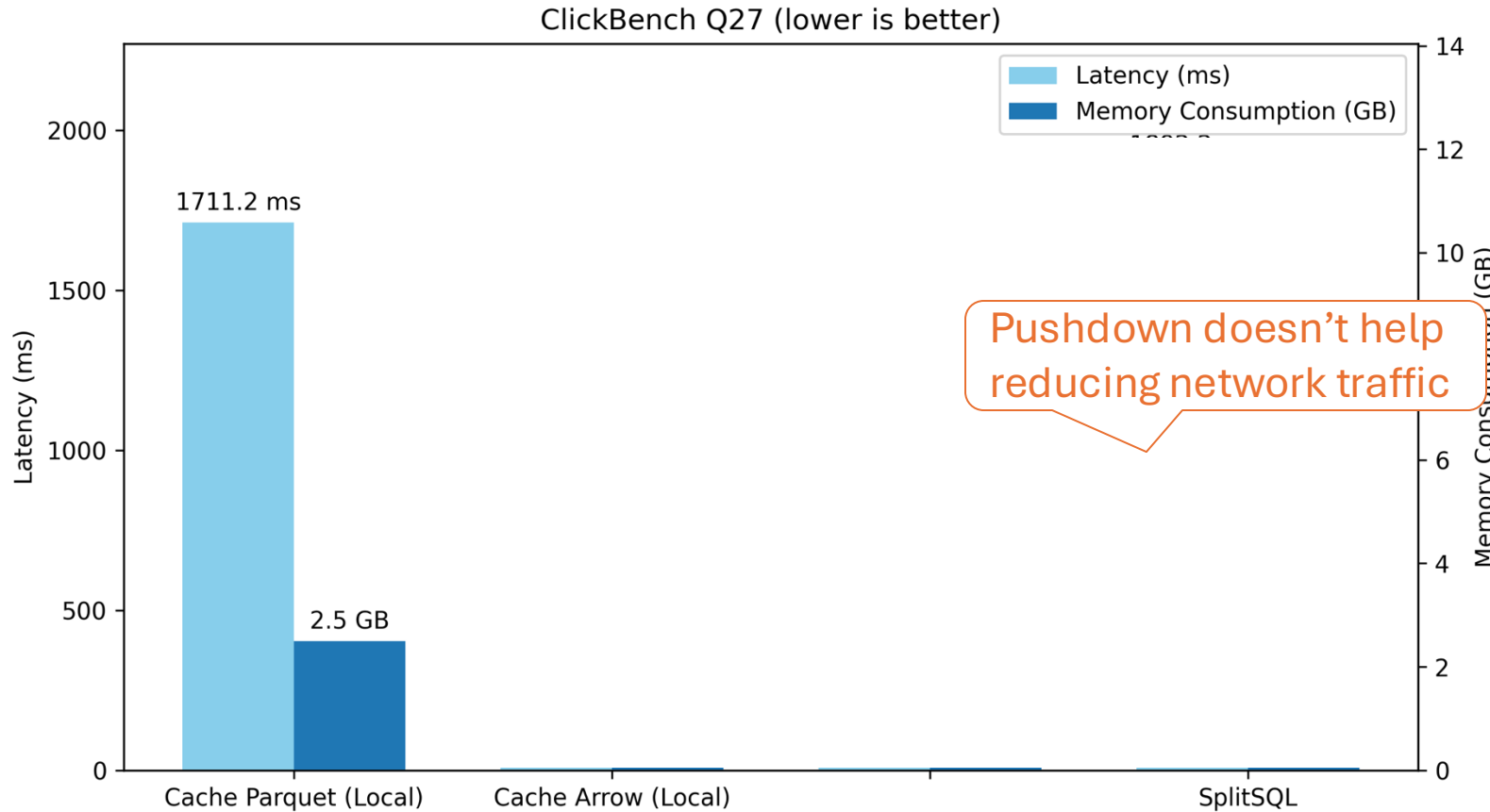
# Evaluations – Q22

Object store → Cache → Compute

SplitSQL

Cache Arrow
Cache Parquet

ClickBench Q22 (lower is better)

Legend:
- Latency (ms)
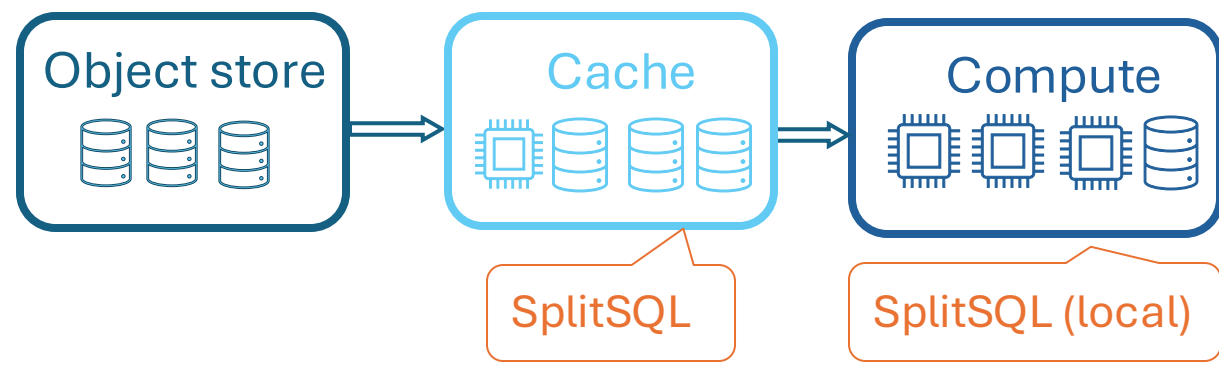- Memory Consumption (GB)

Cache Parquet (Local): 1386.5 ms, 5.3 GB

Cache Arrow (Local)

SplitSQL

```sql
SELECT
"SearchPhrase",
MIN("URL"),
MIN("Title"),
COUNT(*) AS c, COUNT(DISTINCT
"UserID")
FROM hits
WHERE
"Title" LIKE '%Google%'
AND
"URL" NOT LIKE '%.google.%' AND
"SearchPhrase" <> ''
GROUP BY "SearchPhrase" ORDER BY c
DESC LIMIT 10;
```
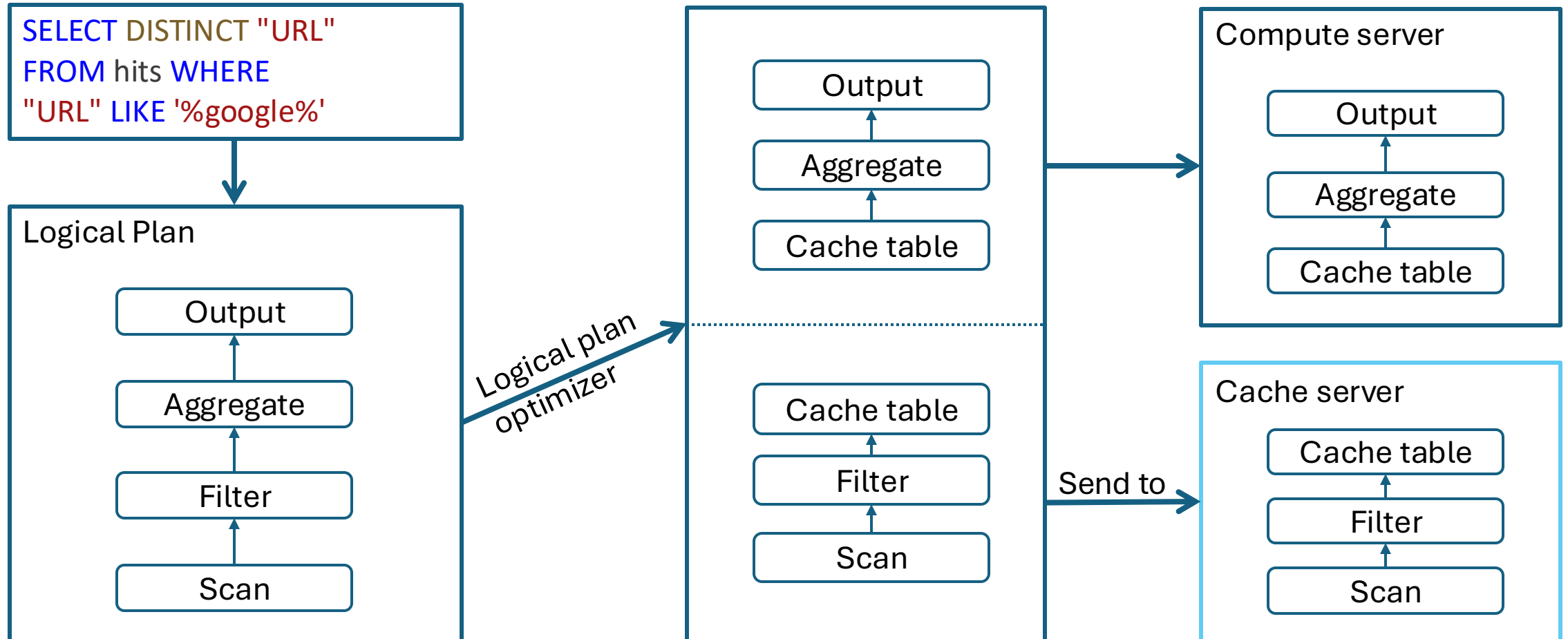
# Evaluations – Q27

Object store

Cache

Compute

SplitSQL

SplitSQL (local)

ClickBench Q27 (lower is better)

Latency (ms)
Memory Consumption (GB)

1711.2 ms

2.5 GB

Latency (ms)

Memory Consumption (GB)

Cache Parquet (Local)    Cache Arrow (Local)    SplitSQL

Pushdown doesn't help reducing network traffic

```
SELECT
"CounterID", AVG(length("URL")) AS
l, COUNT(*) AS c
FROM hits
WHERE "URL" <> ''
GROUP BY "CounterID" HAVING
COUNT(*) > 100000 ORDER BY l
DESC
LIMIT 25;
```

```
SELECT "CounterID", "URL"
FROM "hits"
WHERE "URL" <> ''
```
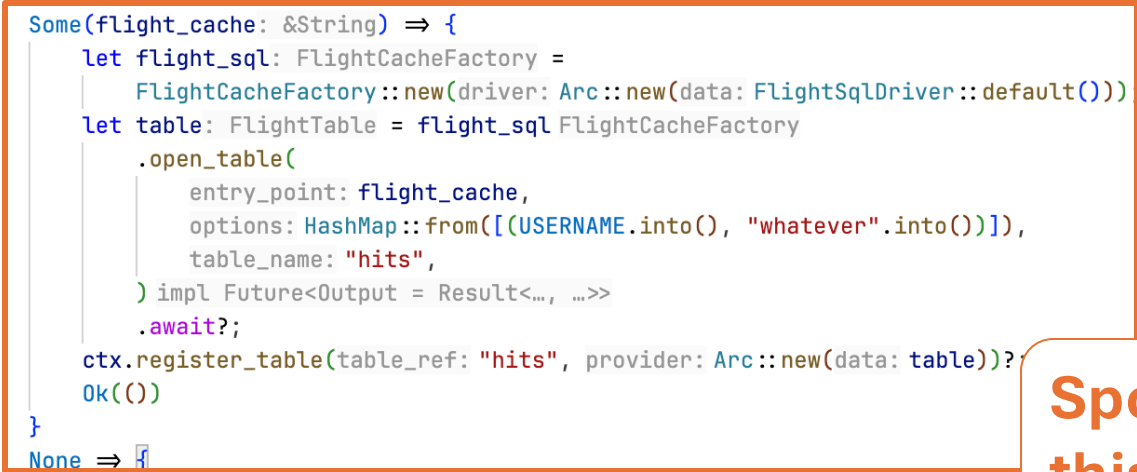
# Implementation in DataFusion

SELECT DISTINCT "URL"
FROM hits WHERE
"URL" LIKE '%google%'



Logical Plan

Output

Aggregate

Filter

Scan

Logical plan optimizer

Output

Aggregate

Cache table

Cache table

Filter

Scan

Send to

Compute server

Output

Aggregate

Cache table

Cache server

Cache table

Filter

Scan

# Easy integration to DataFusion universe

```rust
/// Registers the `hits.parquet` as a table named `hits`
async fn register_hits(
    &self,
    ctx: &SessionContext,
    flight_cache: &Option<String>,
) -> Result<()> {
    let path: &str = self.path.as_os_str().to_str().unwrap();

    match flight_cache {
        Some(flight_cache: &String) => {
            let flight_sql: FlightCacheFactory =
                FlightCacheFactory::new(driver: Arc::new(data: FlightSqlDriver::default()))
            let table: FlightTable = flight_sql FlightCacheFactory
                .open_table(
                    entry_point: flight_cache,
                    options: HashMap::from([(USERNAME.into(), "whatever".into())]),
                    table_name: "hits",
                ) impl Future<Output = Result<…, …>>
                .await?;
            ctx.register_table(table_ref: "hits", provider: Arc::new(data: table))?;
            Ok(())
        }
        None => {
            ctx.register_parquet(table_ref: "hits", table_path: &path, options: Def
                .await Result<()
                .map_err(op
                    DataFus
                        for
                        Box::new(e);
                    )
                })
        }
    }
}

} fn register_hits
```
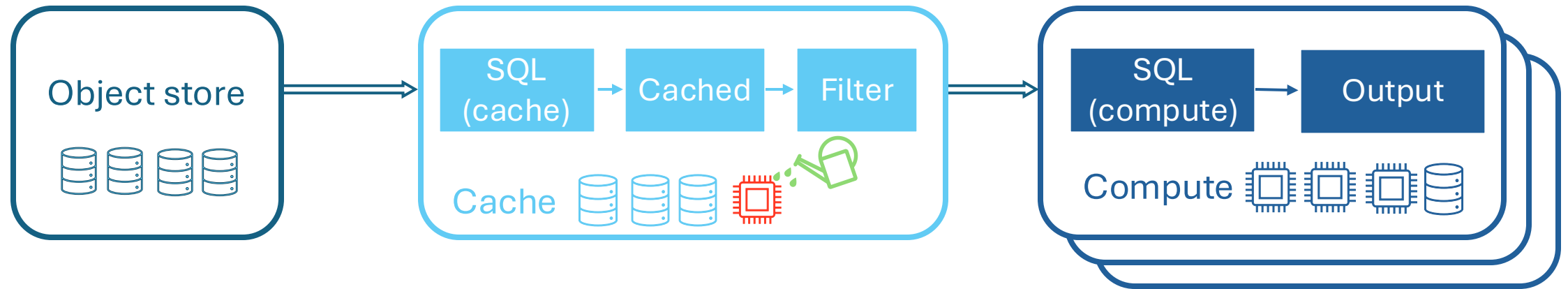
**Sponsor of this work!**

**10 loc change**

# Conclusions & Future work



**Disaggregated**
- Independently scale
- Well-suited for query with filters

**Practical**
- Low CPU overhead
- Compatible with FDAP ecosystem
- Works on commodity hardware

**Even lower network traffic**
- For high-cardinality queries (Q27)
- Aggregate and join push down

**Even faster decoding**
- Storage-aware encodings – different encodings for memory, SSD, HDD