

How Netflix Delivers Key-Value and Time-Series Storage at Any Scale

Joseph Lynch & Vidhya Arvind

December 12, 2023



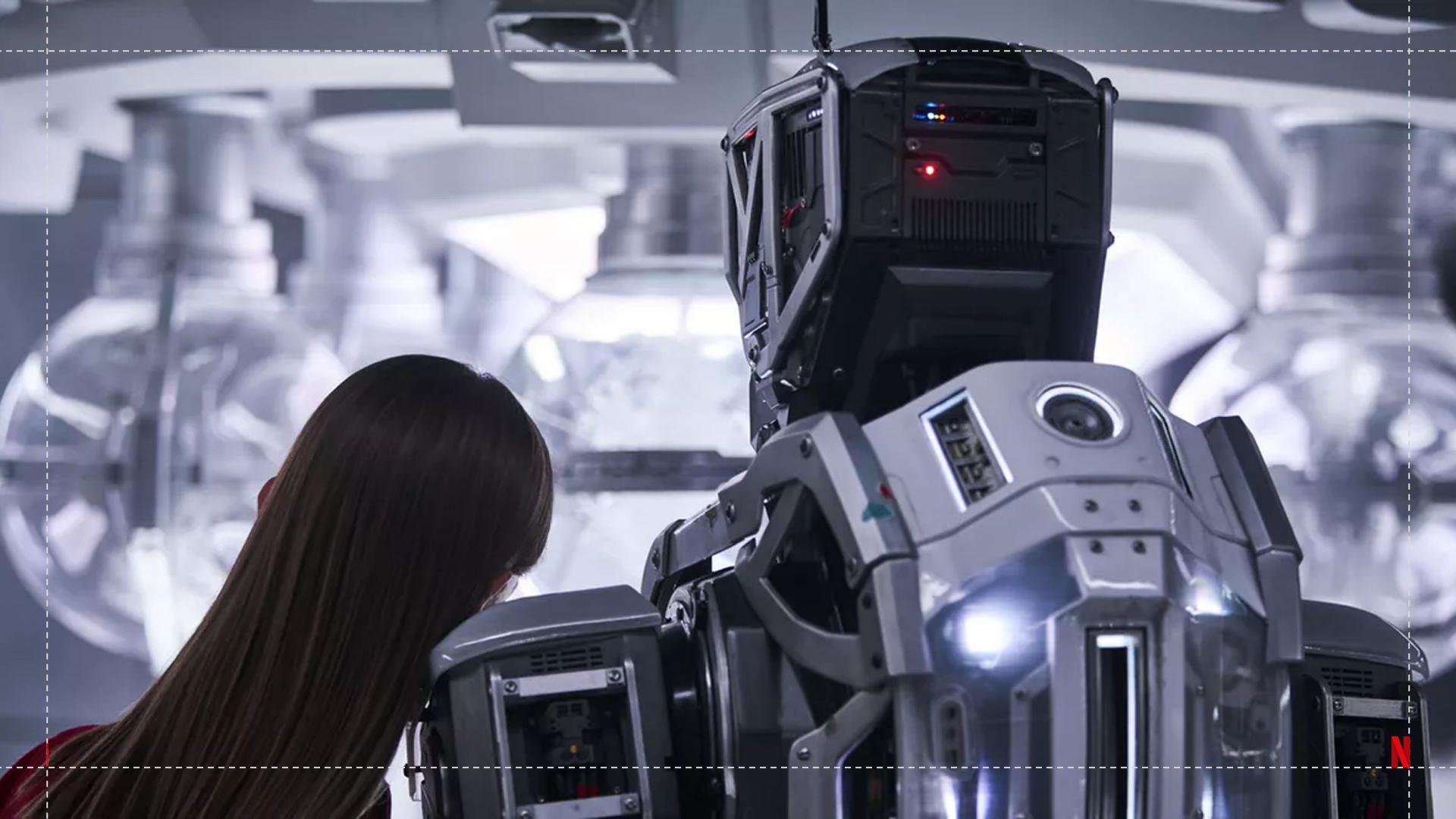
Vidhya Arvind

She/Her
Staff Software Engineer
Data Platform @ Netflix



Joseph (Joey) Lynch

He/Him
Principal Software Engineer
Platform @ Netflix
<https://jolynch.github.io/>



Netflix Storage Problems

All the data

All the scale

All the places

Netflix Online Stateful Scale

Technology	Traffic	Consistency	Latency SLO	Data Size
Near Cache (e.g. Hollow)	0(1 B/s)	Eventual	0(10us)	0(1 TiB)
Remote Cache (e.g. EVCache)	0(50 M/s)	Best-Effort	0(100us)	0(1 PiB)
Databases (e.g. Apache Cassandra)	0(10 M/s)	In-region RyW, Eventual across	0(1000ms)	0(10 PiB)

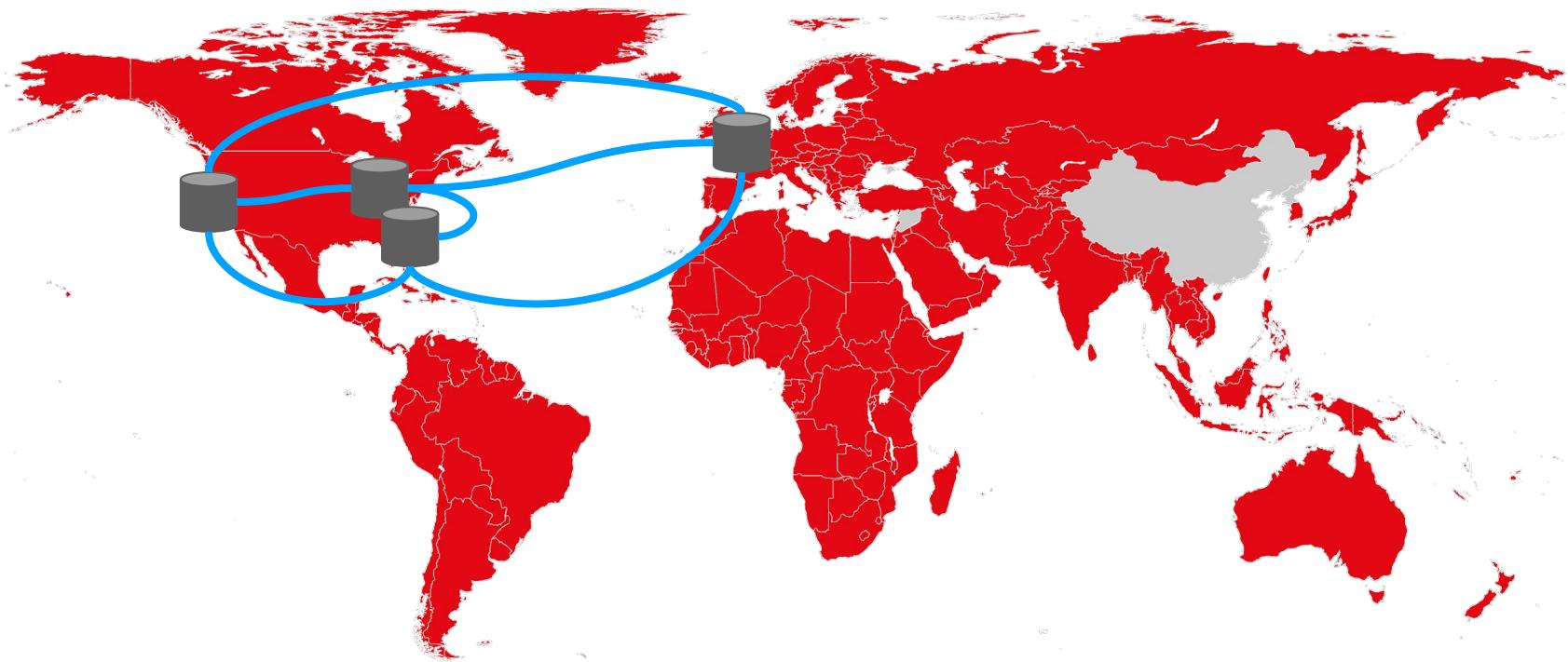
<https://github.com/Netflix/Hollow>

<https://github.com/apache/cassandra>

<https://github.com/Netflix/EVCache>

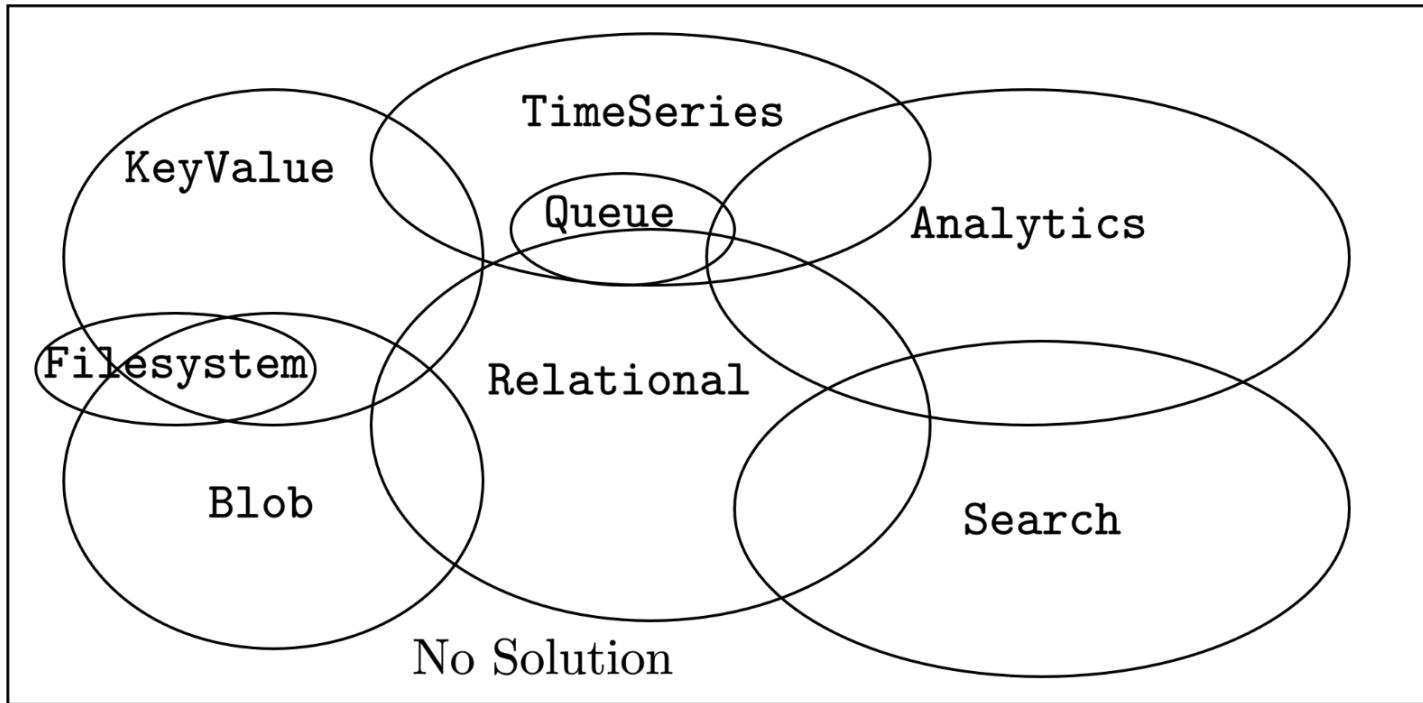


Netflix Online Stateful Scale



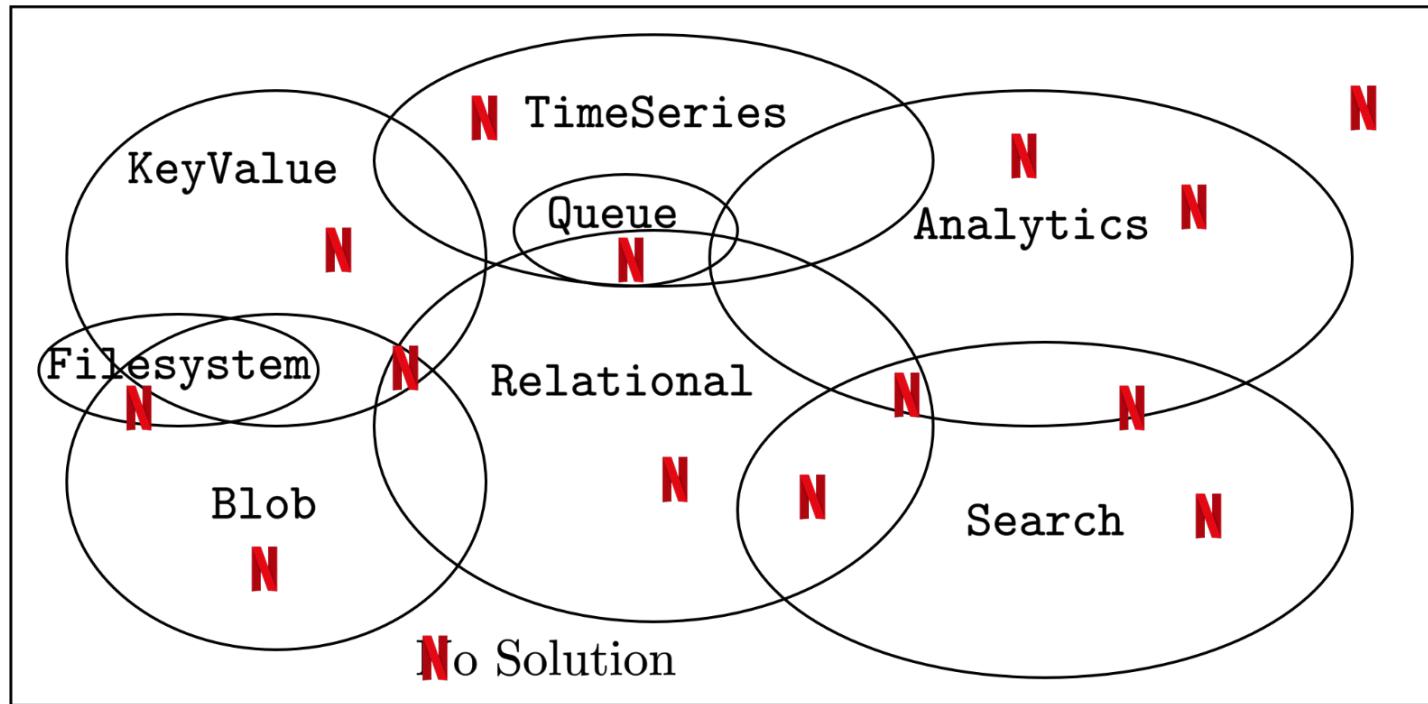
Problem: Varied Requirements

U of database requirements

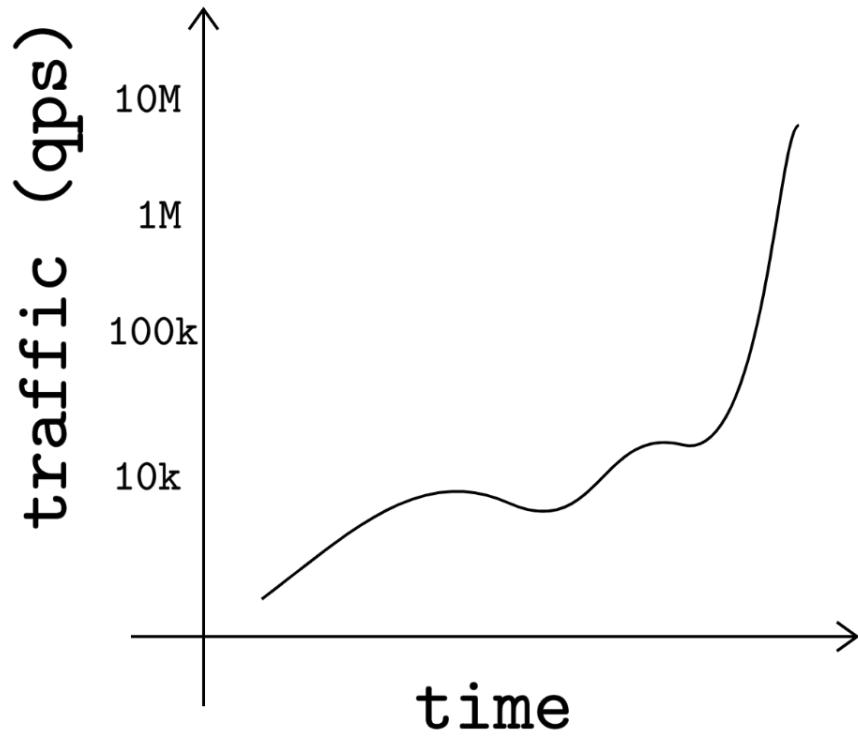
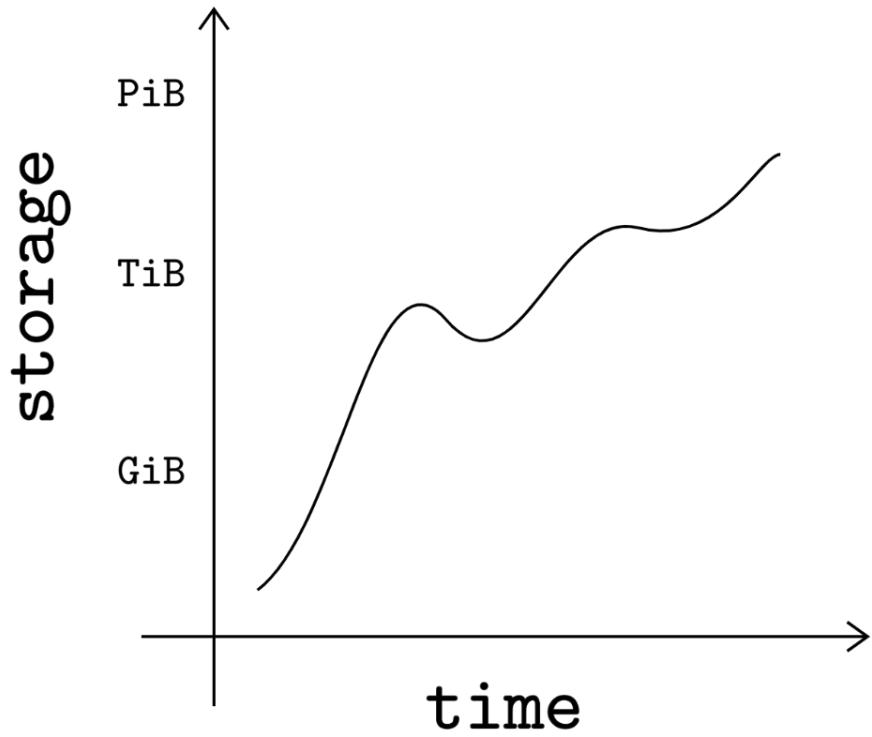


Problem: Varied Use-Cases

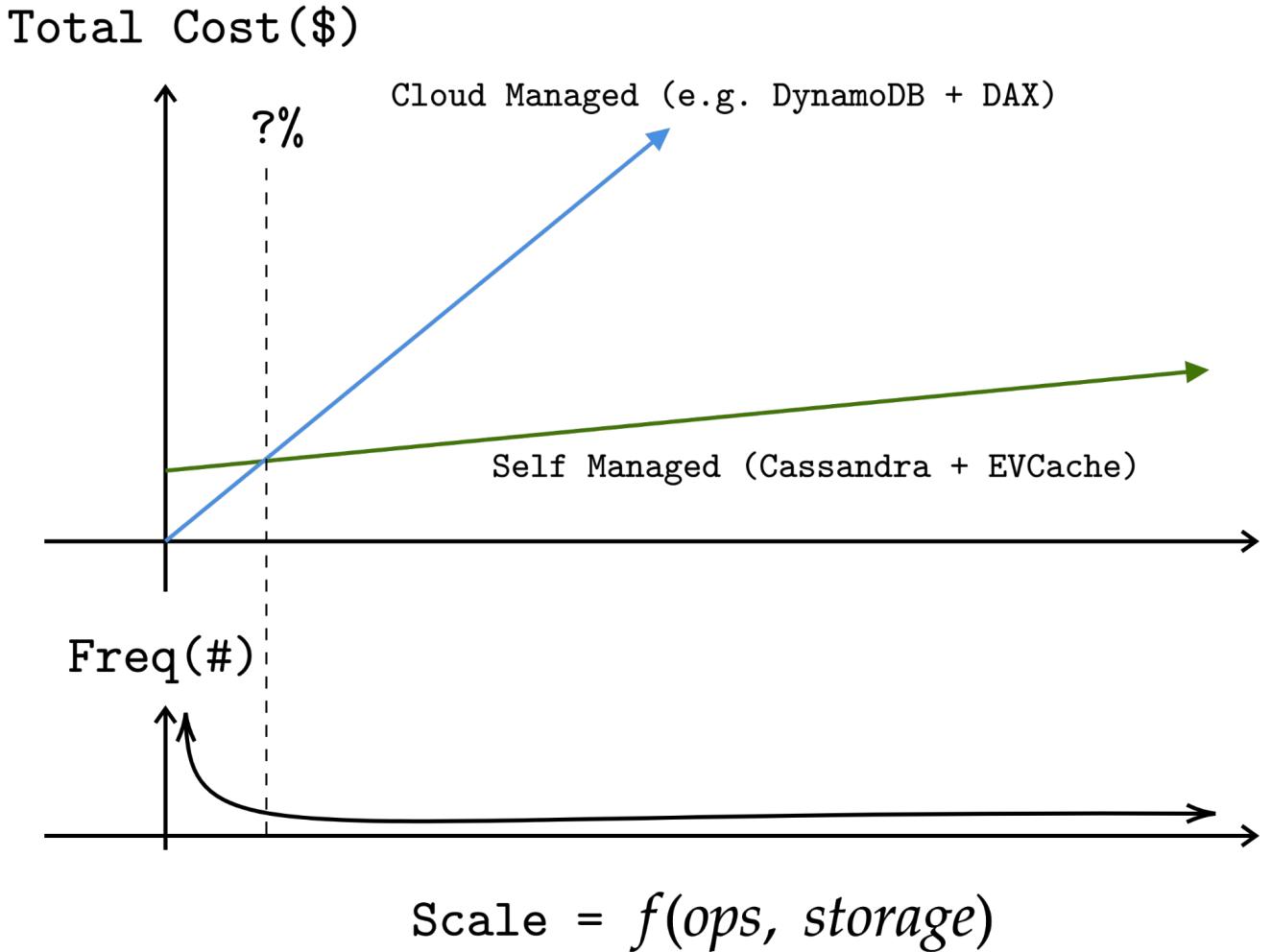
U of database requirements



Problem: Scaling



Problem: Scale-Cost Dynamics

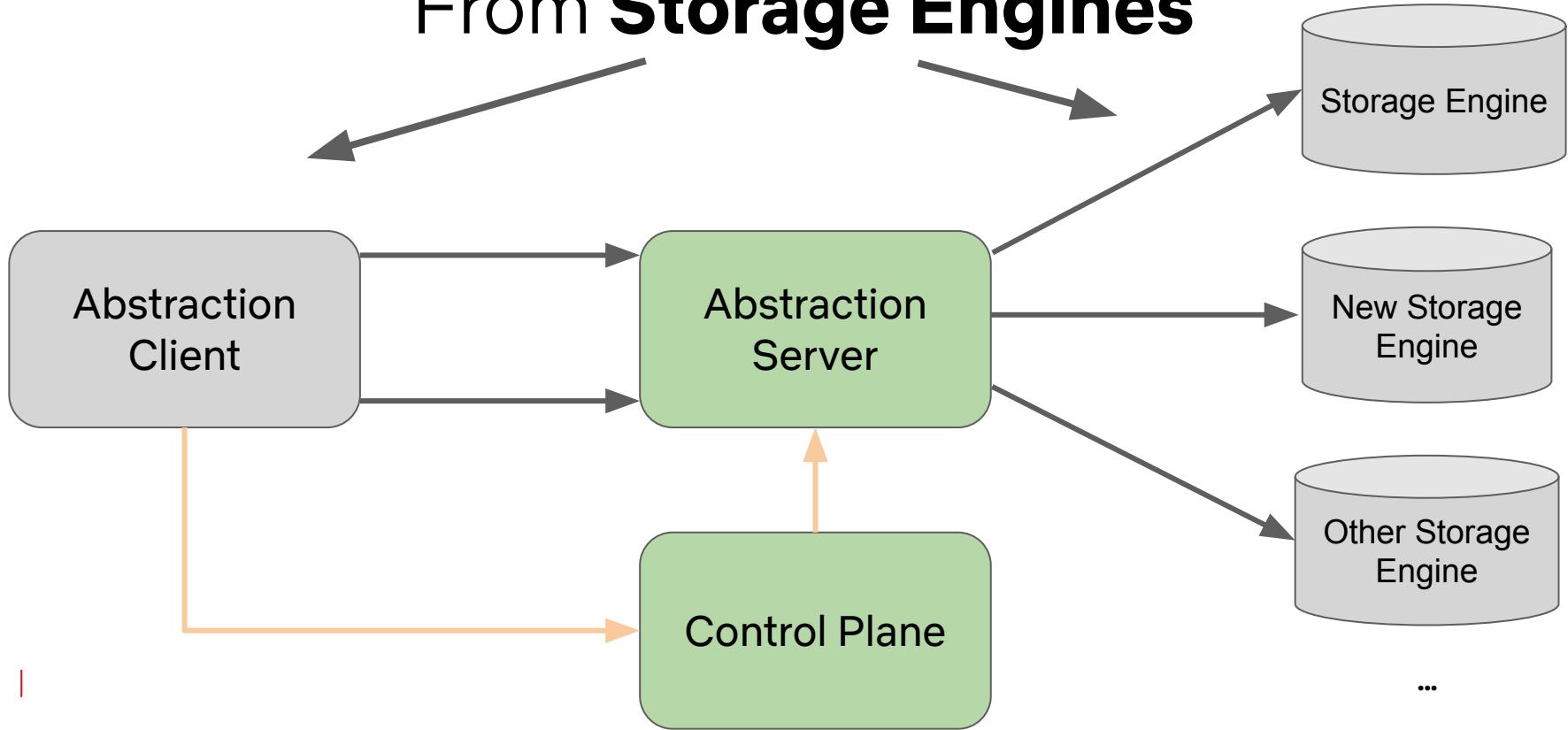


Solution?

Solution?
Data **Abstraction** Layers!

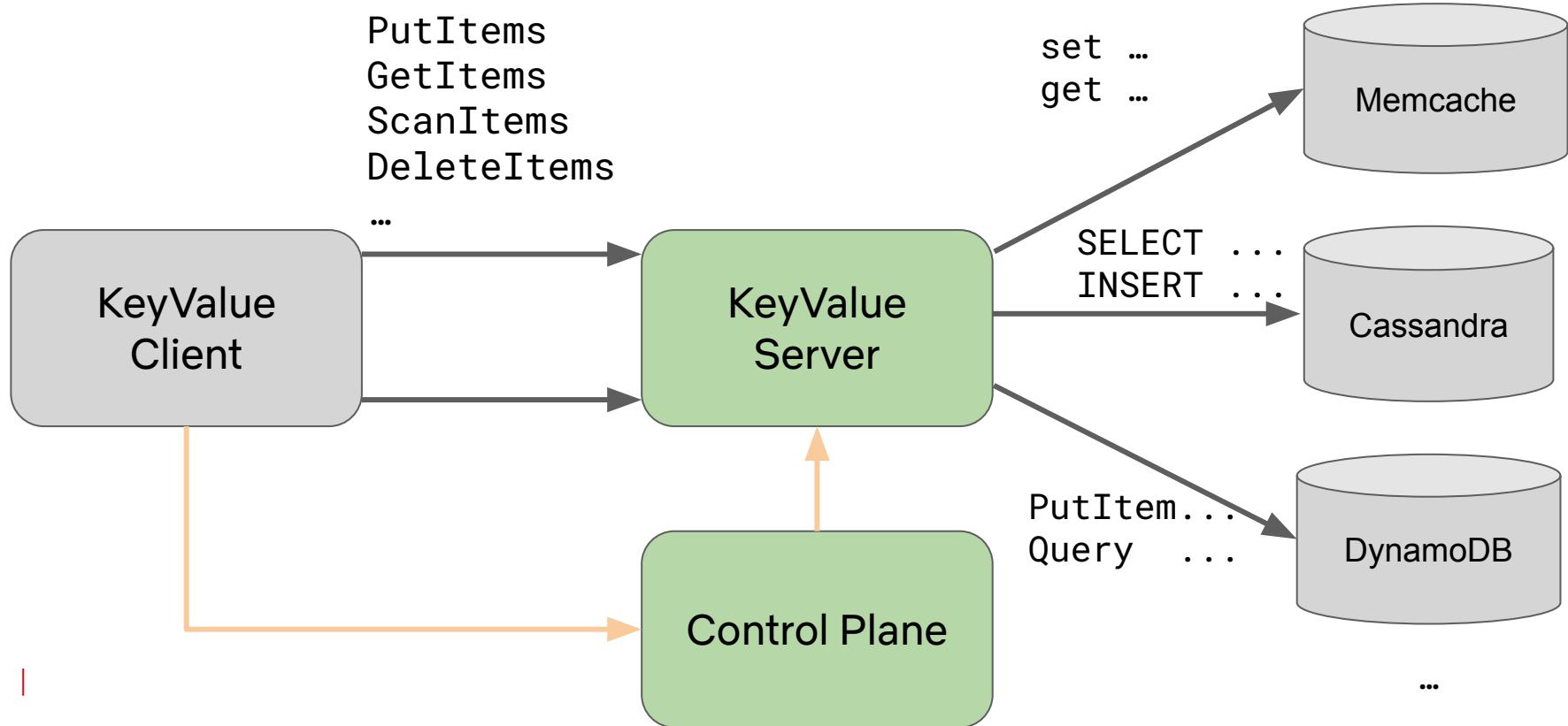
Abstraction

Abstract **Clients** From **Storage Engines**

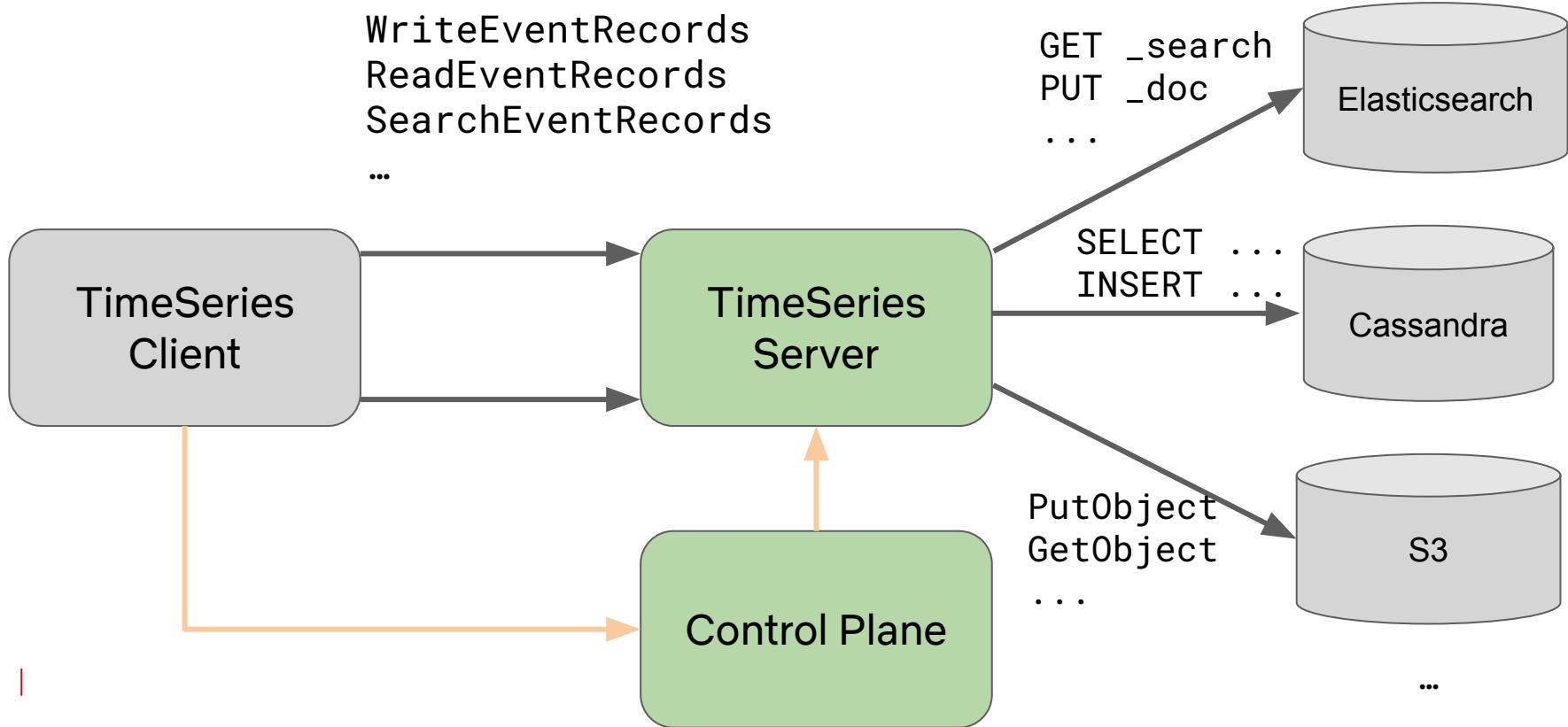


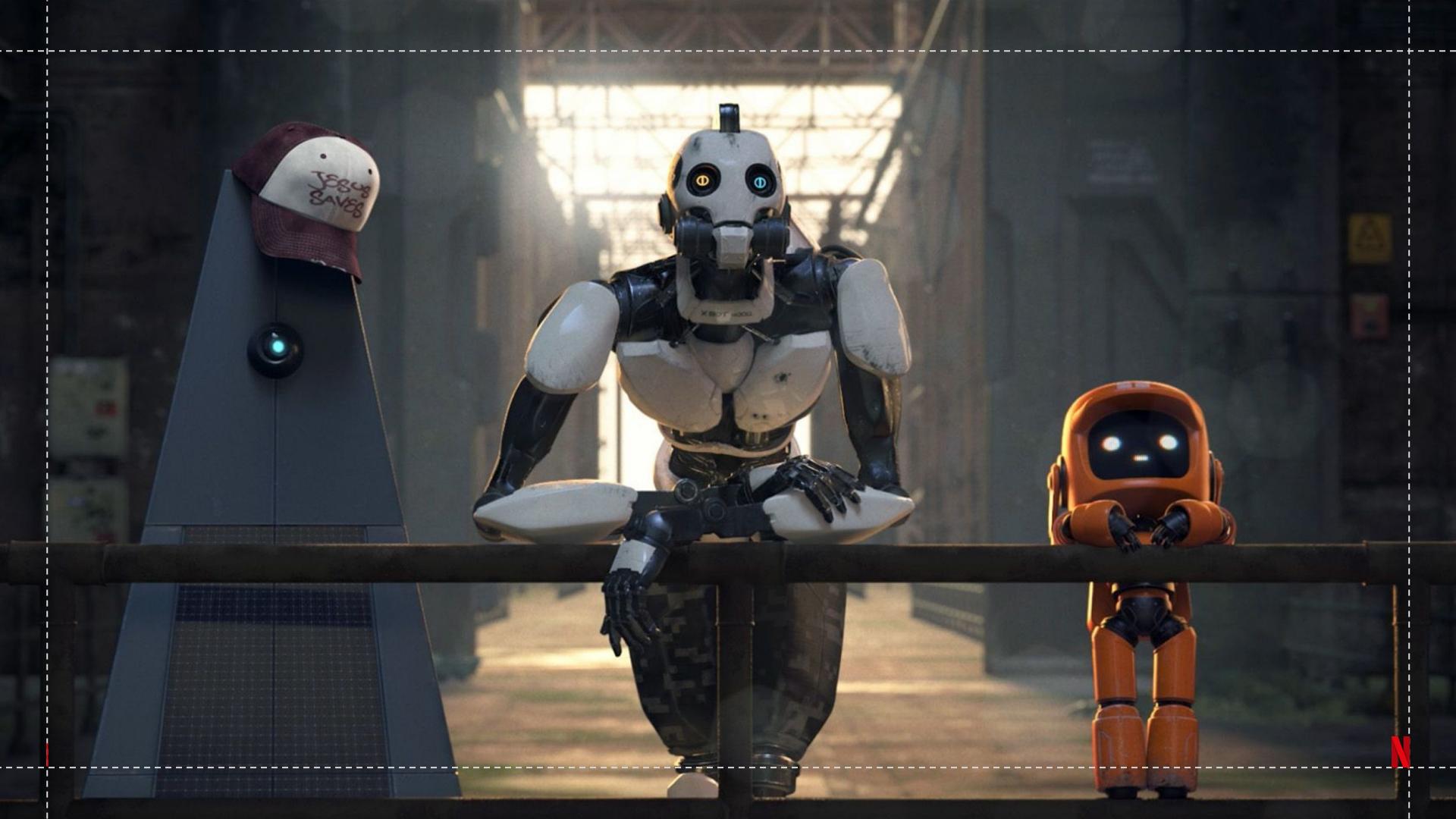
Key-Value

"Two-Level Map as a Service"



"Massively Scalable Event Store"





N

Storage At Any Scale

Key Concepts

Idempotency

Chunking

Compression

Pagination

Adaptivity

Signaling

Key-Value

API

Storage

Time-Series

API

Storage

Storage At Any Scale

Key Concepts

Idempotency

Chunking

Compression

Pagination

Adaptivity

Signaling

Key-Value

API

Storage

Time-Series

API

Storage

Retriable Writes

Assume you have to retry, how do you make your mutable API's safe?

Retriable Writes

Assume you have to retry, how do you make your mutable API's safe?

Idempotency Tokens!

```
def do_write():
    token = get_token()
    result = write(
        data=...,
        idempotency_token=token
    )
    # When hedging or retrying
    # Send the same token
    if result.status != SUCCESS:
        retry_write(
            data=...,
            idempotency_token=token)
```

Idempotency Tokens

Uniquely identify the mutations with an idempotency token.

```
// Mutation operations should be idempotent,  
// and idempotency tokens are how you can  
// achieve that. Writes retry with the same  
// token - then storage engines de-duplicate.  
  
// Typically a LWW clock or a version clock.  
message IdempotencyToken {  
    google.protobuf.Timestamp generation_time = 1;  
    string token = 2;  
}
```

```
last_token = (0, uuid4())

def get_token() -> IdempotencyToken:
    timestamp_micros = max(
        now_in_micros(),
        last_token.generation_time + 10
    )
    nonce = uuid4()
    return (mix_last_digit(timestamp, nonce), nonce)
```

Isolated LWW Tokens

Servers make monotonic timestamps
within same region

```
last_token = (0, uuid4())

def do_isolated_write(mutations):
    # Share random nonce
    nonce = uuid4()
    # Lock keys, send writes, release
    with acquire_clock(mutations.keys) as timestamp:
        for write in mutations:
            perform_write((timestamp.bump(), nonce), write)
    # close releases at last timestamp + 1
```

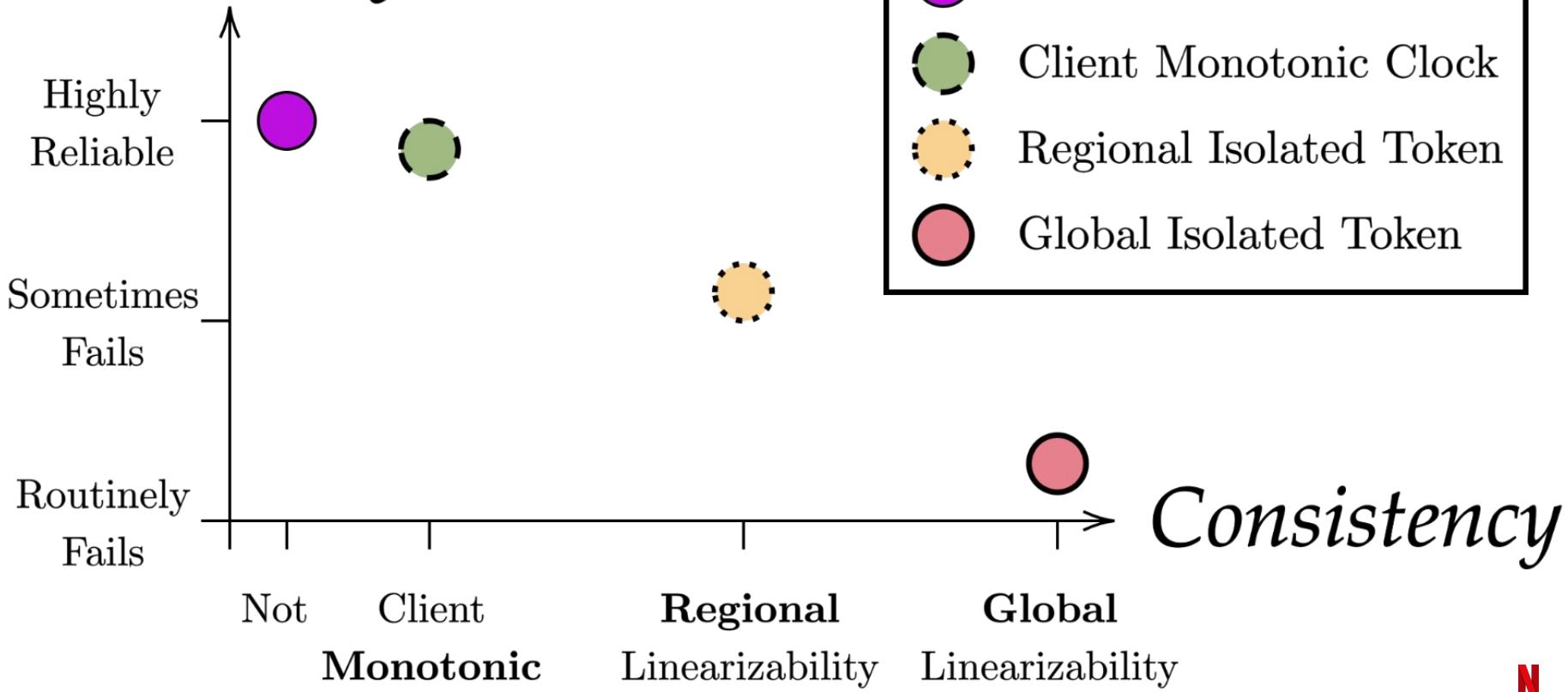
Isolated Global Tokens

Servers make global monotonic version

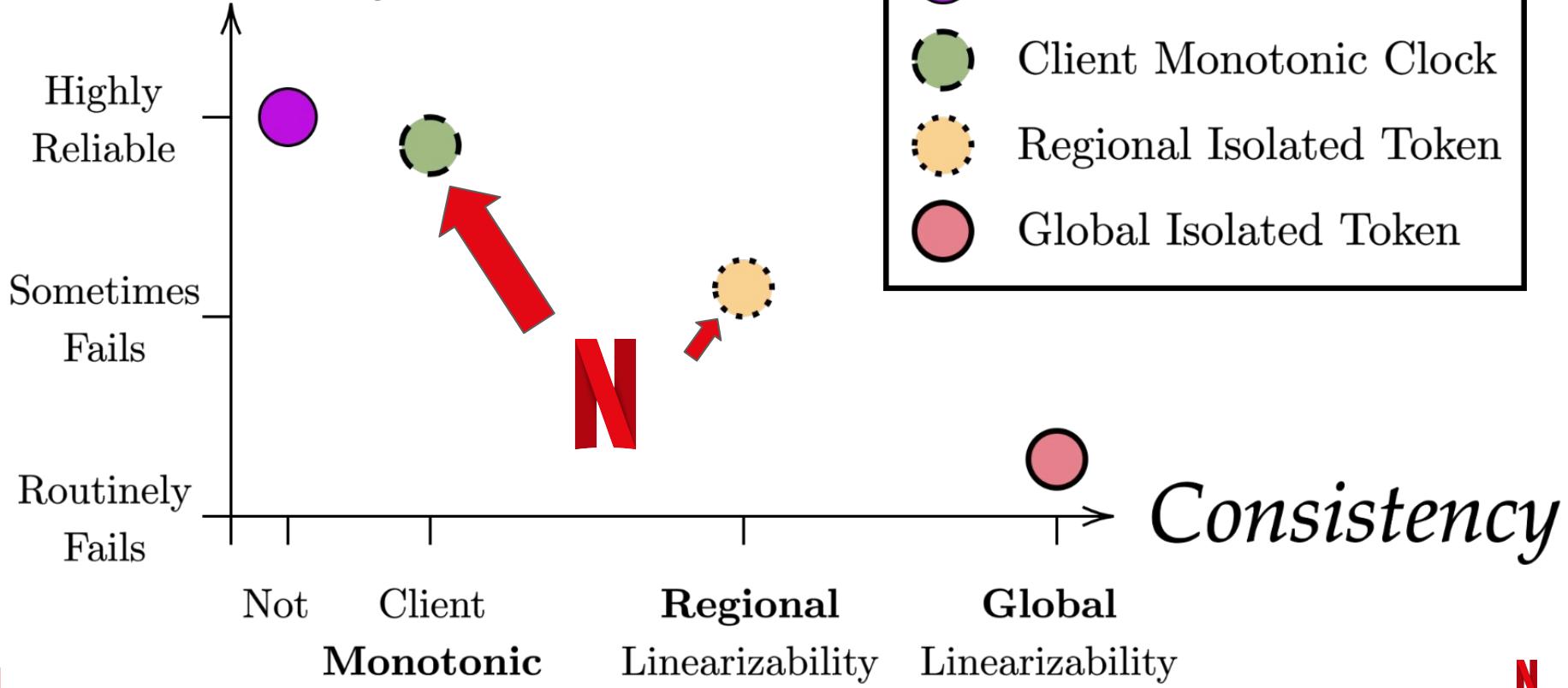
```
last_token = (0, uuid4())

def do_isolated_write(mutations):
    # Share generation timestamp
    timestamp = now()
    # Acquire global transaction id and perform transaction
    with acquire_version(mutations) as token:
        for write in mutations:
            perform_write((timestamp, token), write)
```

Reliability

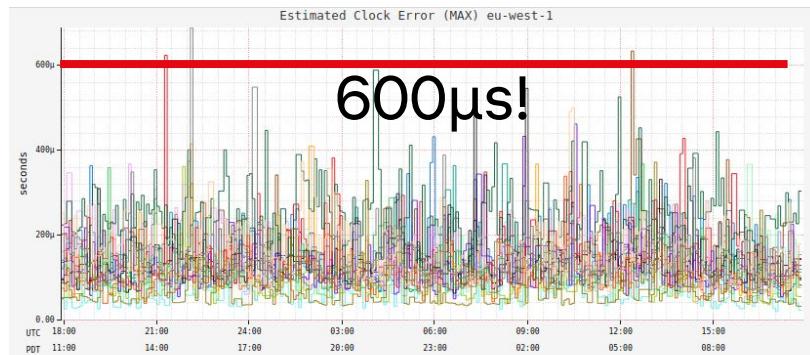


Reliability

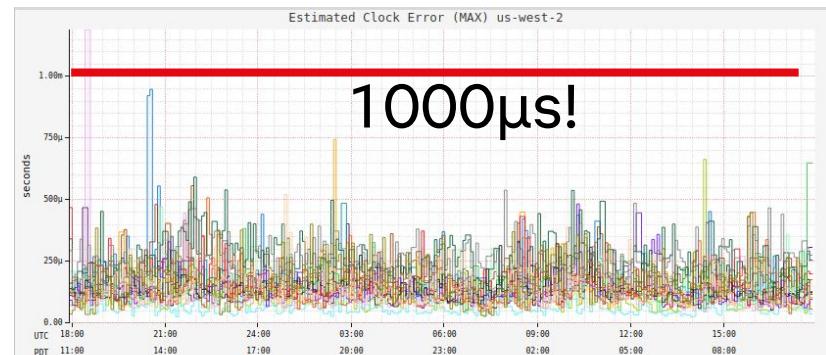


N

Clocks Work



Min	Max ▾	Average	Total	Name	Visibility
49.8 μ	688 μ	137 μ	39.4m	r5d.2xlarge	
81.8 μ	633 μ	167 μ	48.0m	i3.4xlarge	
65.8 μ	623 μ	151 μ	43.4m	m5d.12xlarge	
105 μ	588 μ	245 μ	70.3m	r5.2xlarge	
62.8 μ	545 μ	136 μ	39.1m	r5d.8xlarge	
38.8 μ	499 μ	149 μ	42.9m	r5.4xlarge	
49.8 μ	461 μ	125 μ	36.0m	m5d.4xlarge	
66.8 μ	431 μ	154 μ	44.2m	i3.2xlarge	
57.8 μ	414 μ	122 μ	35.1m	i3en.6xlarge	
48.8 μ	373 μ	146 μ	42.0m	m5.2xlarge	
32.8 μ	371 μ	69.2 μ	19.9m	r5d.large	
55.8 μ	367 μ	126 μ	36.2m	m5d.2xlarge	
53.8 μ	367 μ	125 μ	36.0m	i3en.3xlarge	
54.8 μ	366 μ	121 μ	34.6m	m5.4xlarge	
35.8 μ	361 μ	83.2 μ	23.9m	i3.xlarge	
57.8 μ	341 μ	131 μ	37.7m	i3en.2xlarge	
38.8 μ	338 μ	108 μ	31.0m	m5d.xlarge	
19.8 μ	222 μ	70.5 μ	20.2m	i3en.xlarge	



Min	Max ▾	Average	Total	Name	Visibility
42.8 μ	1.19m	123 μ	35.3m	m5.2xlarge	
88.8 μ	947 μ	183 μ	52.4m	i3.2xlarge	
57.8 μ	742 μ	134 μ	38.5m	i3en.6xlarge	
36.8 μ	662 μ	89.6 μ	25.7m	i3.xlarge	
55.8 μ	648 μ	160 μ	45.9m	m5d.2xlarge	
114 μ	590 μ	266 μ	76.2m	r5.2xlarge	
89.8 μ	555 μ	194 μ	55.6m	i3.4xlarge	
111 μ	537 μ	260 μ	74.6m	r5d.2xlarge	
50.8 μ	519 μ	133 μ	38.2m	r5.4xlarge	
36.8 μ	506 μ	143 μ	41.0m	m5d.xlarge	
74.8 μ	484 μ	138 μ	39.7m	i3en.3xlarge	
59.8 μ	481 μ	144 μ	41.3m	m5d.4xlarge	
53.8 μ	468 μ	141 μ	40.3m	m5.4xlarge	
67.8 μ	466 μ	164 μ	47.0m	i3d.12xlarge	
50.8 μ	466 μ	123 μ	35.4m	i3en.2xlarge	
25.8 μ	446 μ	173 μ	49.8m	r5d.large	
56.8 μ	312 μ	136 μ	39.2m	m5d.8xlarge	
22.8 μ	207 μ	58.8 μ	16.9m	i3en.xlarge	

|Clock Drift Memo: <https://docs.google.com/document/d/1XLBjQ9scZCy-xlo51Rs--CSdFV781fnp5hXdXTBAk1k/edit>

Storage At Any Scale

Key Concepts

Idempotency

Chunking

Compression

Pagination

Adaptivity

Signaling

Key-Value

API

Storage

Time-Series

API

Storage

Small data

Payload size < O(1MiB)

```
# Chunk after payload size 1 MiB
chunk_after_size = 1024 * 1024

# Send data to storage engine directly in one
# stage + commit step
if sizeof(payload) < chunk_after_size:
    perform_write(payload)
```

Small data

Payload size < O(1MiB)

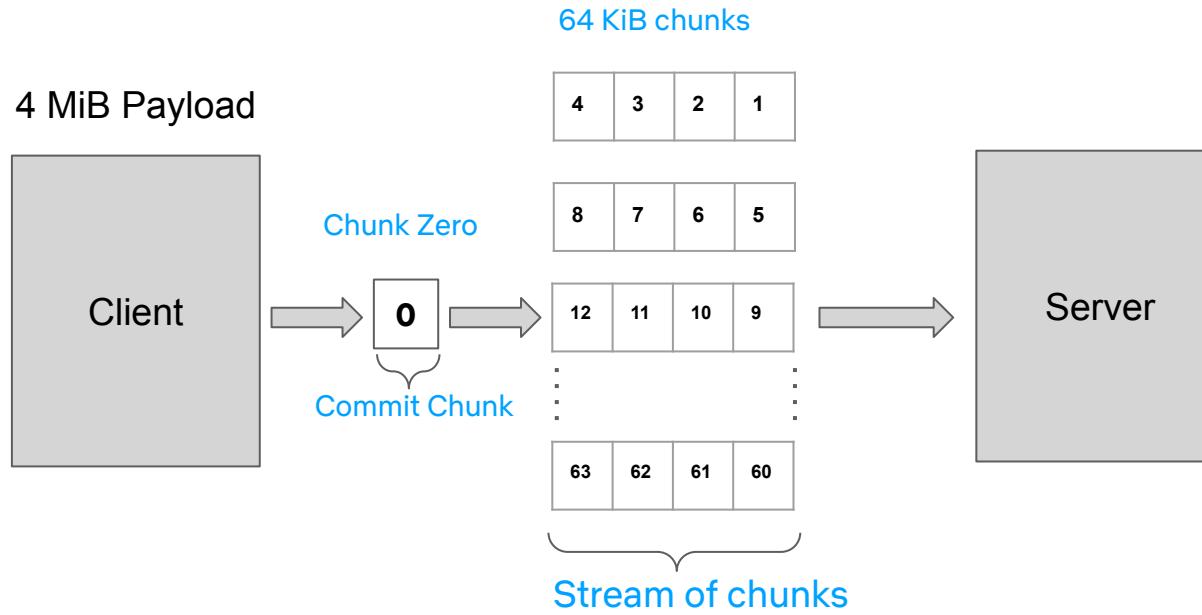
```
# Chunk after payload size 1 MiB
chunk_after_size = 1024 * 1024

# Send data to storage engine directly in one
# stage + commit step
if sizeof(payload) < chunk_after_size:
    perform_write(payload)

else ?
```

Chunk Data

Write Path: Payload size: O(> 1MiB)



chunk-after-bytes = 1 MiB

chunk-size-bytes = 64 KiB

Chunk Data

Write Path: Payload size: O(> 1MiB)

```
# Step 1: Stage data chunks and use the idempotency
# token to tie chunks together
page, writes = page(), []
for idx, chunk in enumerate(chunk(payload)):
    item = Item(
        key=key, value=chunk,
        chunk=(idx + 1),
        token=idempotency_token
    )
    page.add(item)
    if page.size > 2 * 1024 * 1024:
        writes.add(perform_write(page))
        page.clear()
writes.add(perform_write(page))
```

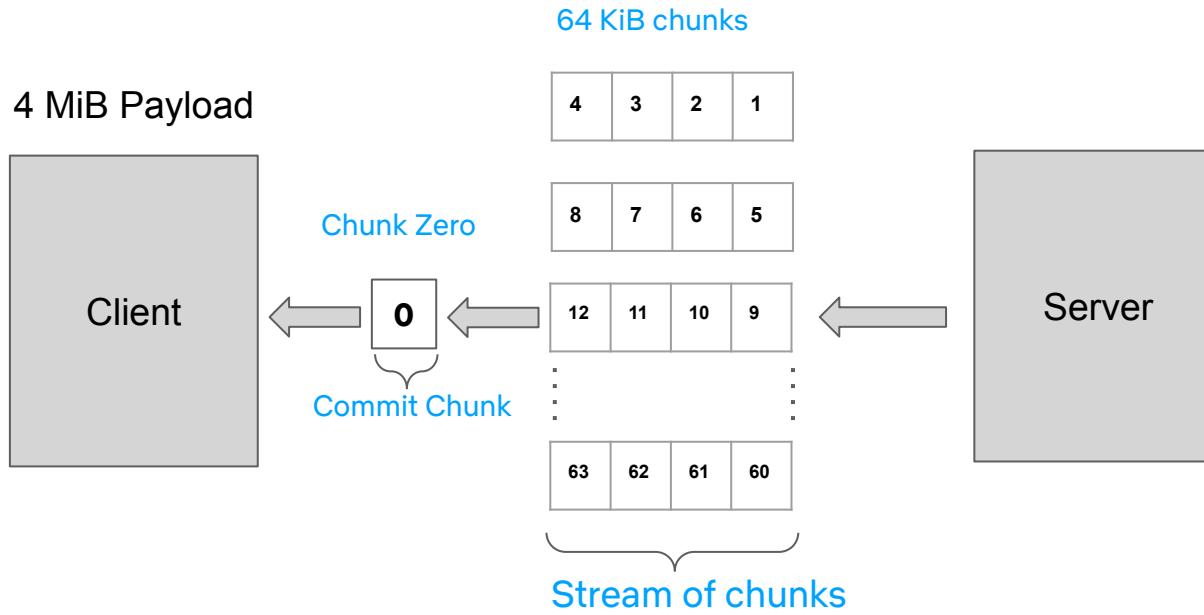
Chunk Data

Write Path: Payload size: O(> 1Mib)

```
# Step 2: Commit the write by writing chunk zero
perform_write(
    Item(
        key      = key,
        chunk    = 0,
        token    = idempotency_token,
        metadata = {
            chunk_count      = 43,
            chunk_size_bytes = 64KiB,
            # Filled by server
            chunk_token_xxh64 = xxh64(token.nonce),
            chunk_offset       = consistent_hash(key)
        }
    )
))
```

Chunk Data: Read Path

Read Path: Payload size: $O(> 1\text{Mib})$



chunk-after-bytes = 1 MiB

chunk-size-bytes = 64 Kib

Chunk Data

Read Path: Payload size: O(> 1Mib)

```
# Server-side: Read chunk 0s value_metadata for
# number of chunks to retrieve
chunks, commit = [], read_chunk_zero(id, key)
version         = commit.metadata.version
num_chunks     = commit.metadata.chunks
chunk          = 1
while sizeof(chunks.value) < page_limit:
    last_chunk = perform_chunk_read(chunk)
    chunks.add(last_chunk)
page_token = last_chunk.chunk < num_chunks
            ? construct_page_token(last_chunk)
            : null
return (chunks, page_token)
```

```
# Client-side: retrieve pages and stitch chunks
for page in pages:
    for item in page.items:
        current_item = select(current_item, item)
        # SMALL ITEM: Read the entire item from the single RPC
        if !current_item.metadata.has_chunk_metadata():
            yield current_item
        else
            # LARGE ITEM: append the chunk to the item
            current_item.value.append(item.value)
            if item.chunk == current_item.metadata.chunk_count:
                yield current_item
                current_item = null
```

Storage At Any Scale

Key Concepts

Idempotency

Chunking

Compression

Pagination

Adaptivity

Signaling

Key-Value

API

Storage

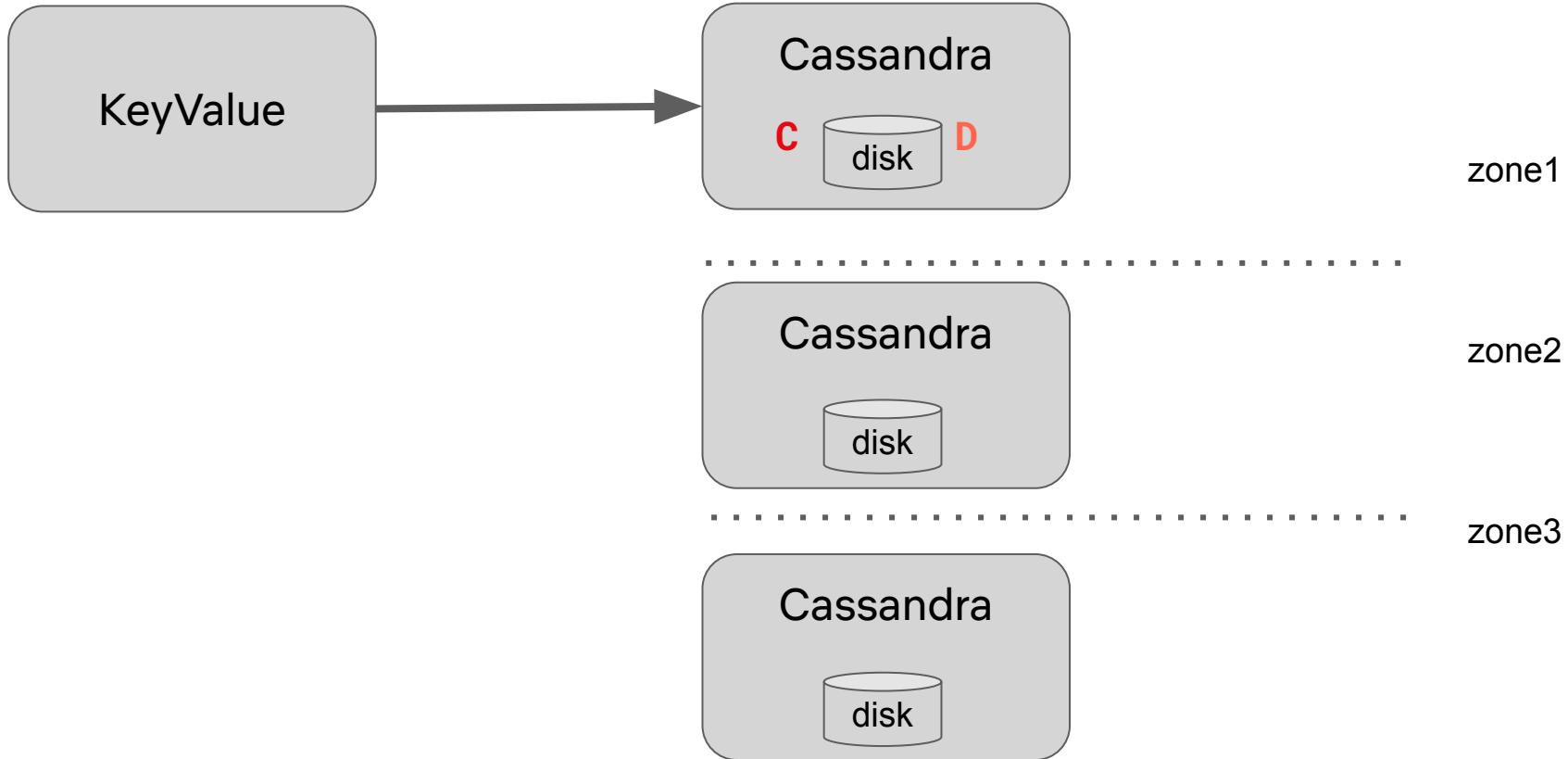
Time-Series

API

Storage

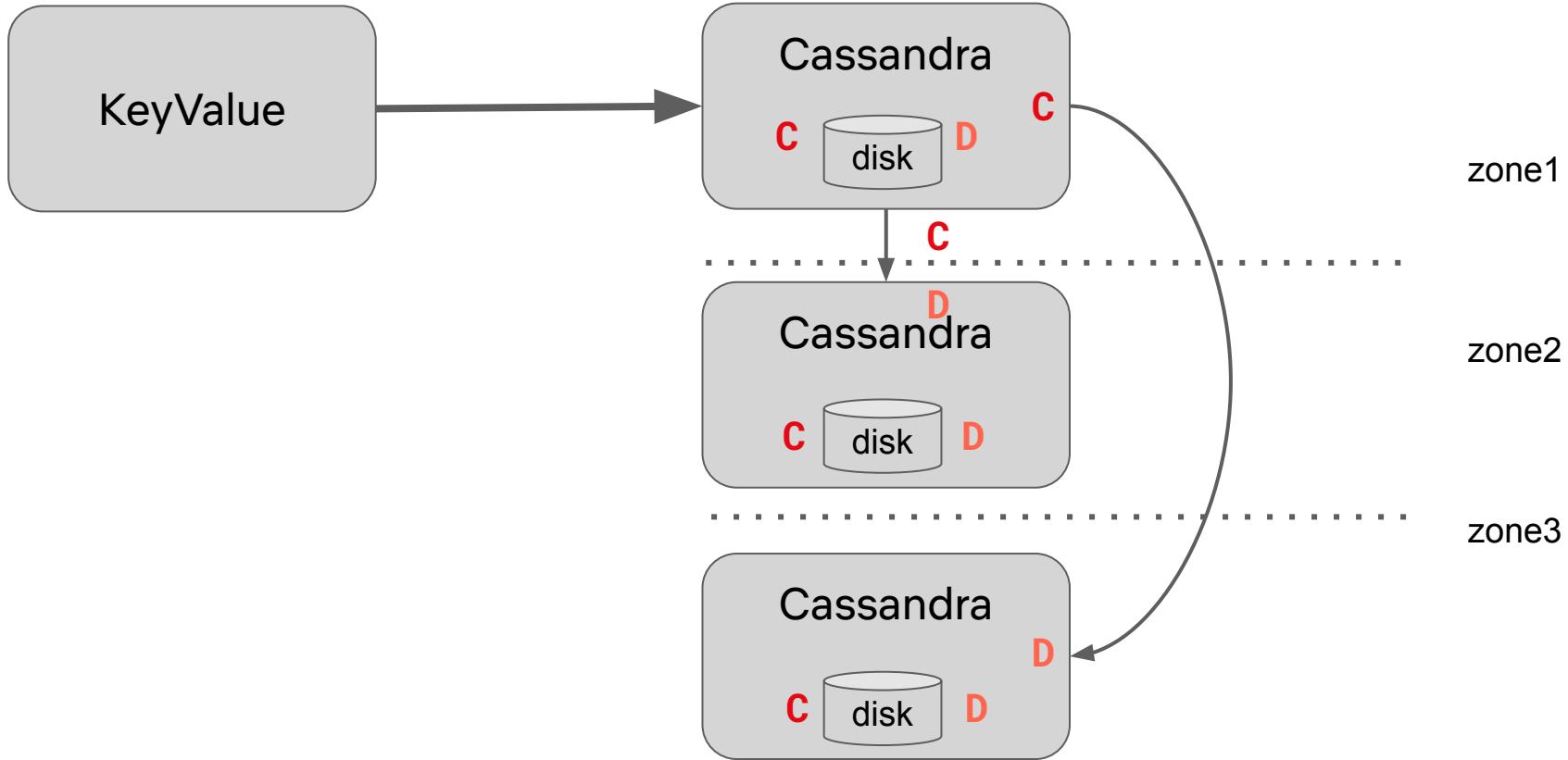
Compression

Server Traditionally Compresses



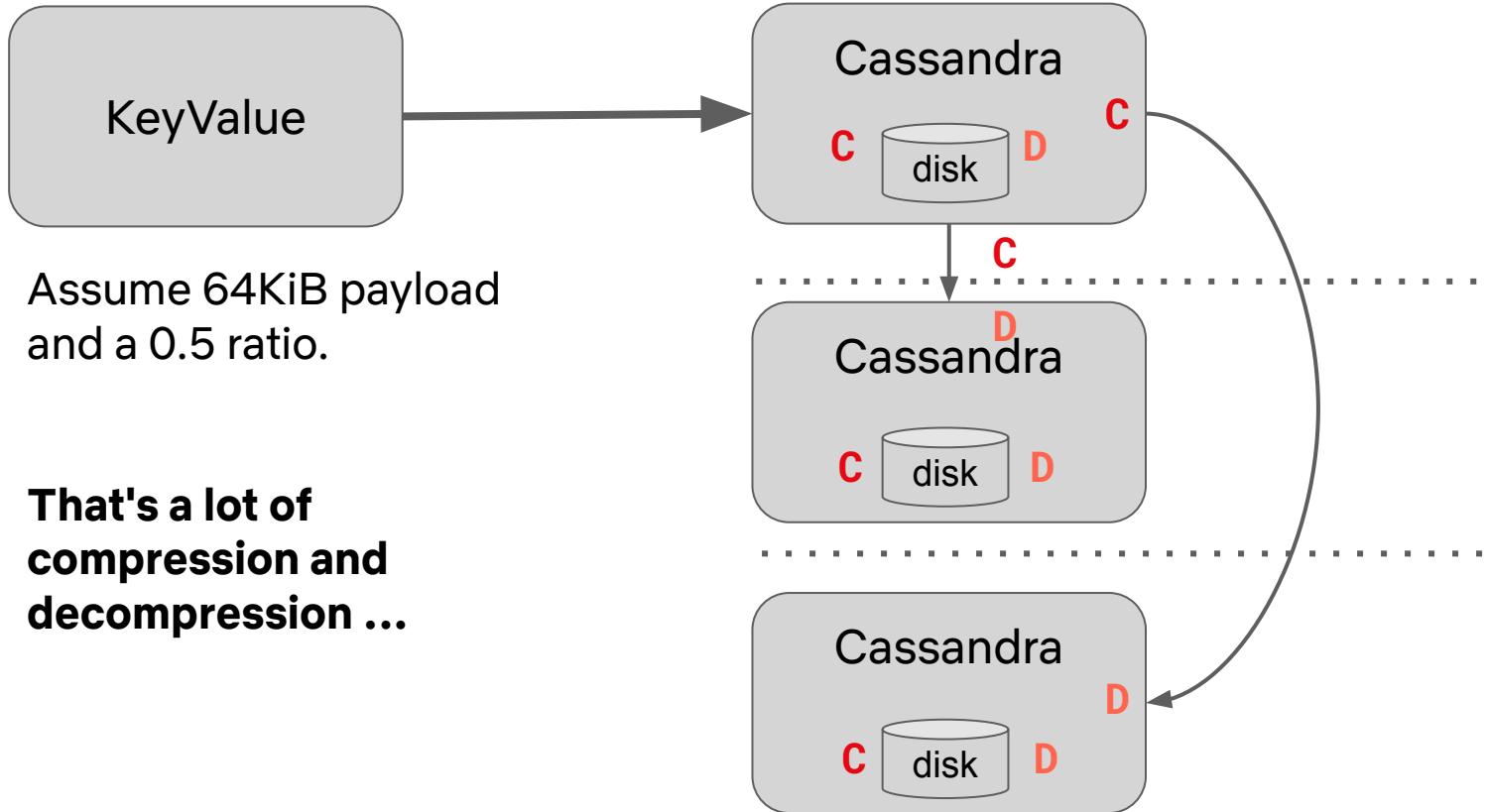
Compression

Server Traditionally Compresses



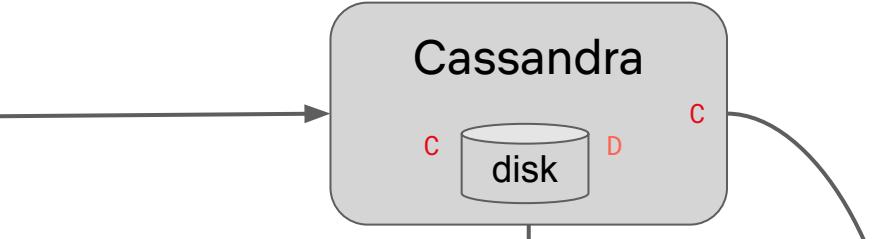
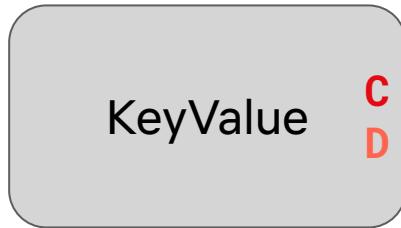
Compression

Server Traditionally Compresses



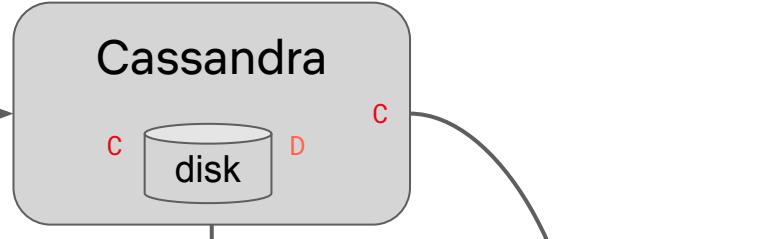
Compression

Better to Compress Chunks at Client

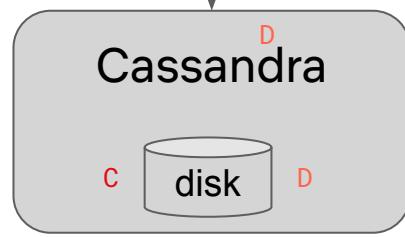


Storage still
compresses, but saves:

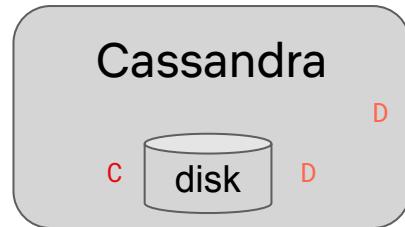
- Commit Log
- Allocations
- Disk IO
- Network IO
- Overall ratio



zone1



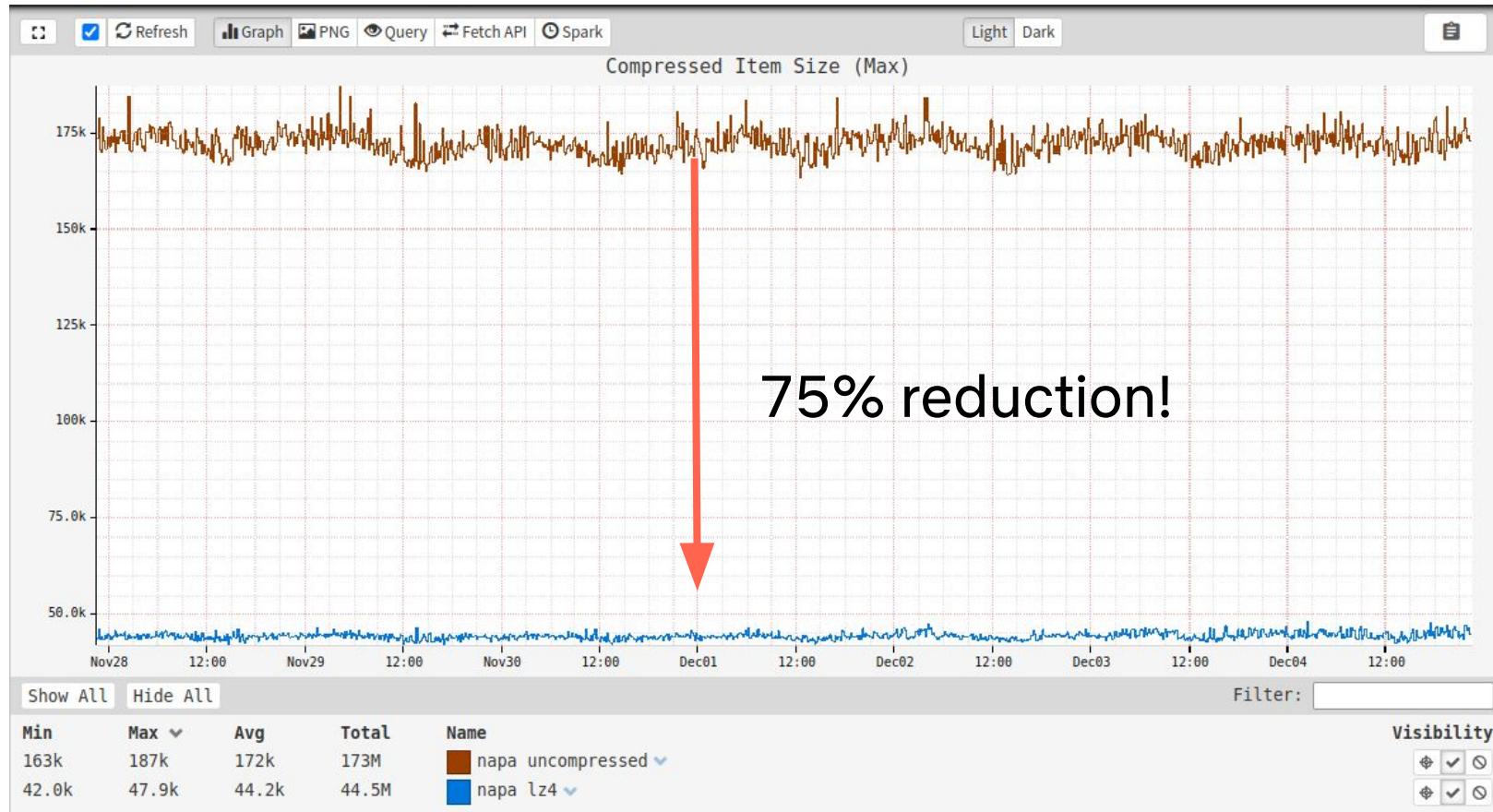
zone2



zone3

Compression

Better to Compress Chunks at Client



Storage At Any Scale

Key Concepts

Idempotency

Chunking

Compression

Pagination

Adaptivity

Signaling

Key-Value

API

Storage

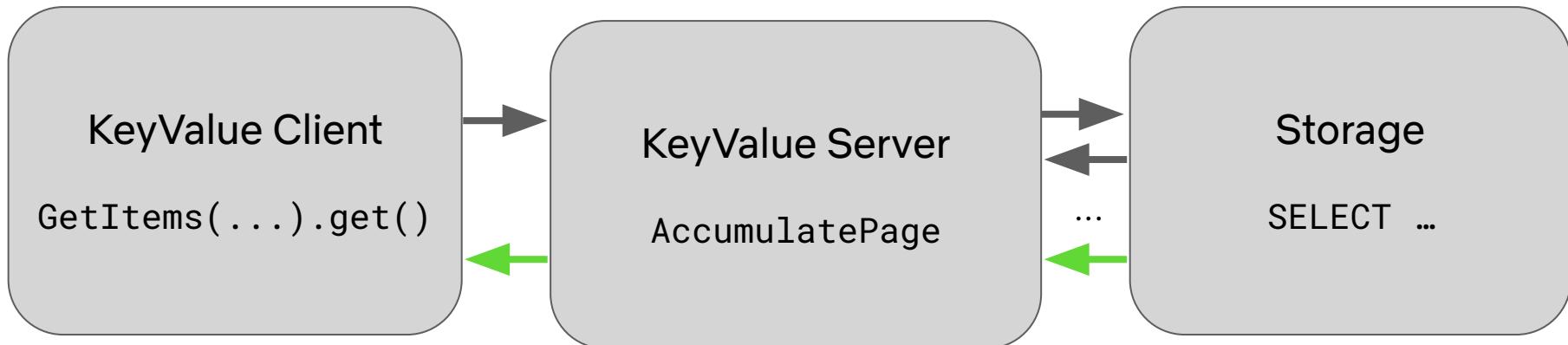
Time-Series

API

Storage

Pagination

Page Limit of 4MiB



```
{  
  "items": [item1, item2, ...],  
  "next_page_token": <page_token>  
}
```

Pagination

Accumulate pages of fixed size work (MiB), clients must ask for more

Storage engines almost always paginate by row count

Robust APIs paginate by *size not count*. Translation required.

Storage At Any Scale

Key Concepts

Idempotency

Chunking

Compression

Pagination

Adaptivity

Signaling

Key-Value

API

Storage

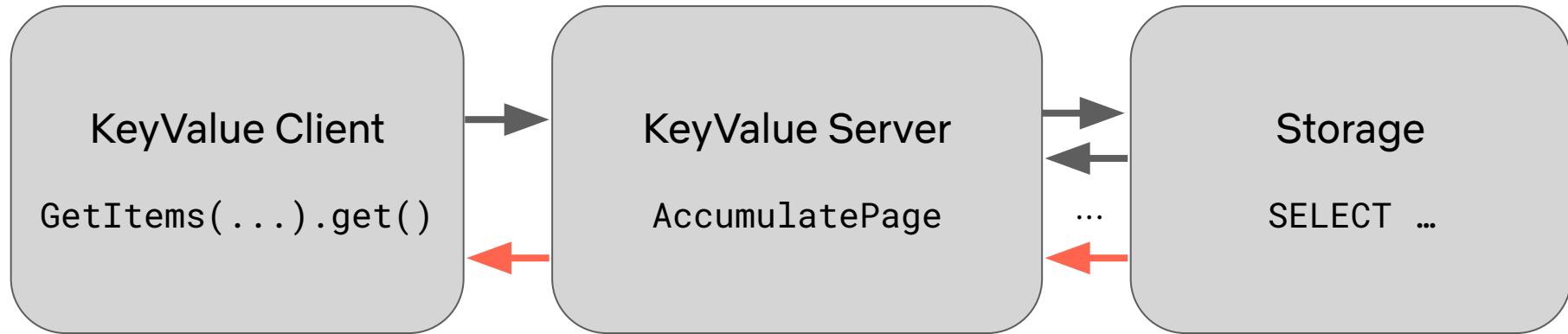
Time-Series

API

Storage

Adaptivity: Target Page size

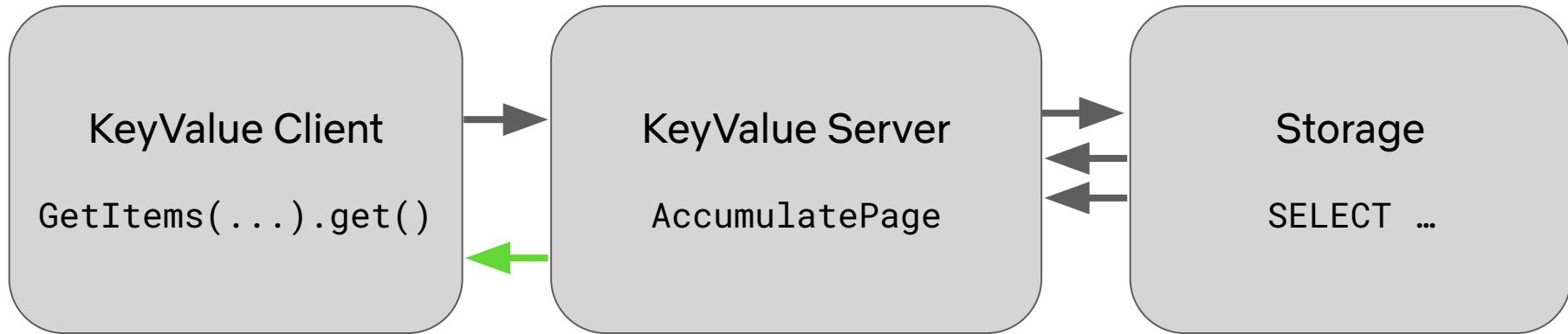
KeyValue Read Target page size
HashMap[String, SortedMap[Bytes,
Bytes]]



Server might do many round trips
to DB to accumulate a full page

**Adaptivity:
Target Page size**

KeyValue Read Target page size
HashMap[String, SortedMap[Bytes,
Bytes]]



4MiB Page size reached → change the page size via
`fetch_size` to limit DB calls

Adaptivity: Early SLO Return

KeyValue Read SLO

HashMap[String, SortedMap[Bytes, Bytes]]

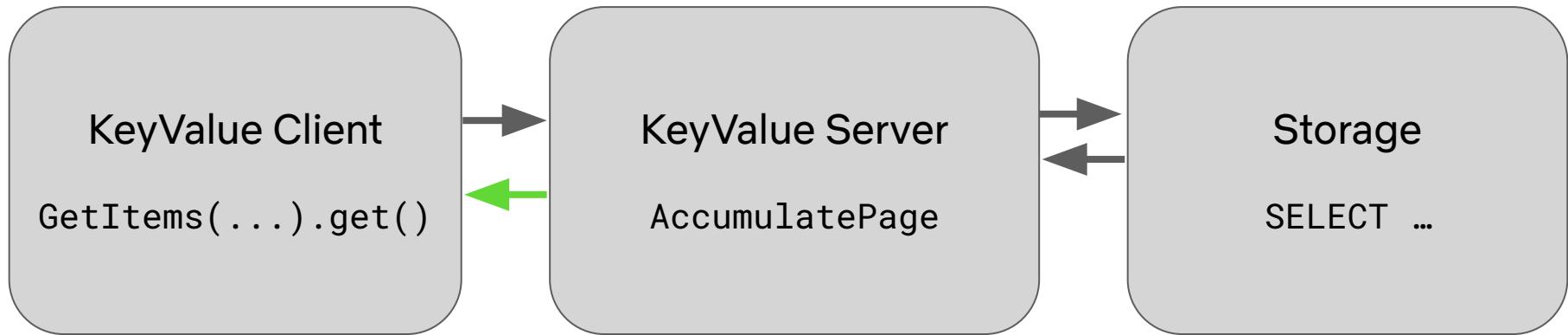


DB might have more data to fetch but SLO is reached before a page limit is reached

Adaptivity: Early SLO Return

KeyValue Read SLO

HashMap[String, SortedMap[Bytes, Bytes]]



80% SLO_{max} reached → return early

Configuration?

Storage At Any Scale

Key Concepts

Idempotency

Chunking

Compression

Pagination

Adaptivity

Signaling

Key-Value

API

Storage

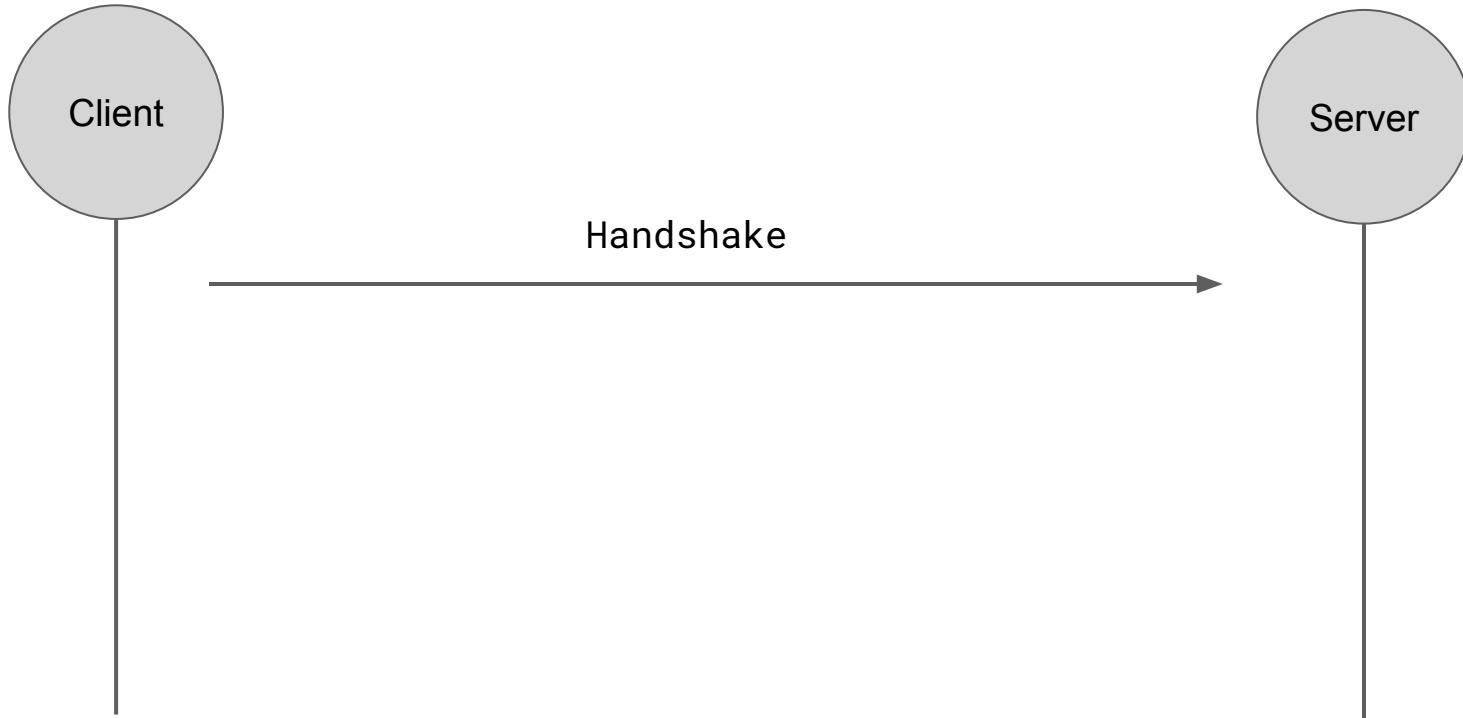
Time-Series

API

Storage

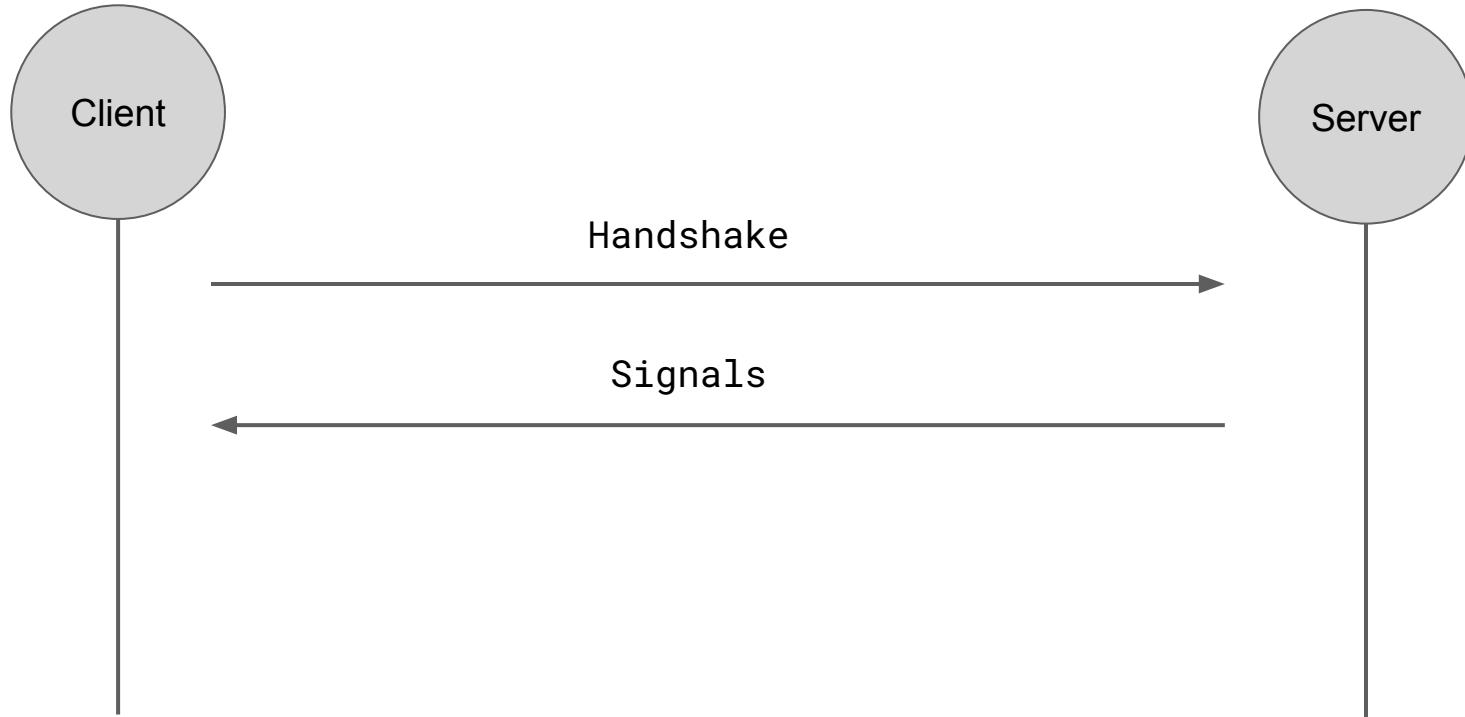
Signals

Clients Introspect Servers



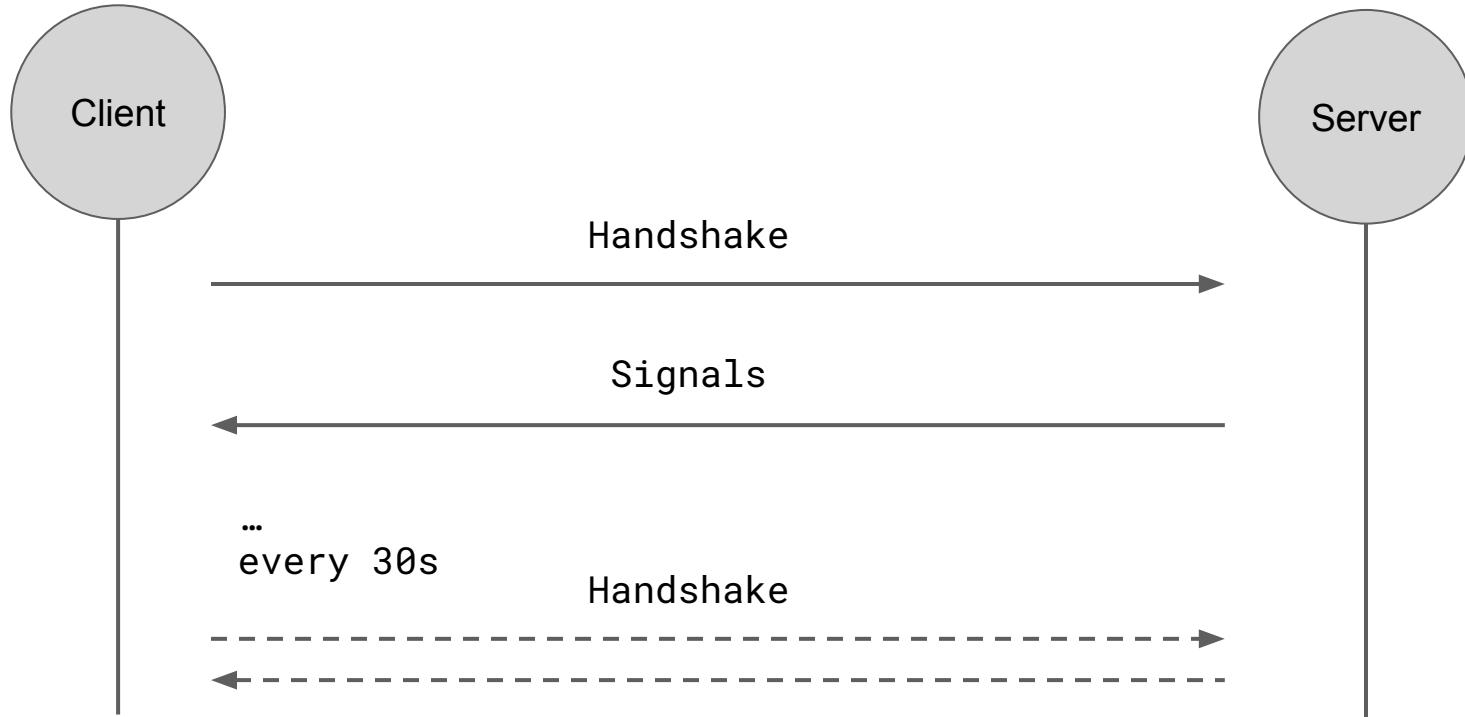
Signals

Clients Introspect Servers



Signals

Clients Introspect Servers



Servers Signal to Clients

Client tuning is error prone
Prefer to send signals from Servers

```
$ grpc -a dgwkv.napa
    -e prod
    -r us-east-1
    com.netflix.dgw.kv.v2.KeyValueServiceV2/Handshake
| jq .signals
```

```
$ grpc -a dgwts.sessionlogs
    -e prod
    -r us-west-2
    com.netflix.dgw.ts.TimeSeriesService/Handshake
| jq .signals
```

Servers Signal to Clients

Client tuning is error prone
Prefer to send signals from Servers

```
"payload": {  
    "chunked": {  
        "chunk-after        "chunk-size": 65536  
    },  
    "stats": {  
        "record-count": "39371484"  
    }  
}
```

Clients can transparently
chunk payloads to make
them fixed amount of work

Large payloads must be
broken down into chunks

Servers Signal to Clients

Client tuning is error prone
Prefer to send signals from Servers

```
"compression": [  
    {  
        "algorithm": "LZ4",  
        "parameters": {  
            "min-length": 128  
        }  
    }  
]
```

Clients can transparently
compress payloads

Fewer bytes = more reliable
(also checksumming)

Servers Signal to Clients

Client tuning is error prone
Prefer to send signals from Servers

```
"slo_by_endpoint": {  
    "access": {  
        "endpoint": "access",  
        "latency": {  
            "target": "0.010s",  
            "max": "0.500s"  
        },  
        "concurrency_limits": {  
            "normal": 1000,  
            "hedge": 50  
        }  
    }  
}
```

```
"slo_by_endpoint": {  
    "access": {  
        "endpoint": "access",  
        "latency": {  
            "target": "0.050s",  
            "max": "0.500s"  
        },  
        "concurrency_limits": {  
            "normal": 1000,  
            "hedge": 50  
        }  
    }  
}
```

Each Namespace is Different

```
"namespaces": {  
    "eventual_ns": {  
        "consistency_scope":  
            "LOCAL", [REDACTED]  
        "consistency_target":  
            "EVENTUAL", [REDACTED]  
    }  
    ...  
    "physical_storage": {  
        "type": "CASSANDRA",  
        "cluster": "clstr1",  
        "dataset": "ks1",  
        "table": "tbl1",  
    }  
}
```

Lower consistency provides better SLO

```
"namespaces": {  
    "strong_ns": {  
        "consistency_scope":  
            "GLOBAL", [REDACTED]  
        "consistency_target":  
            "READ_YOUR_WRITES", [REDACTED]  
    }  
    ...  
    "physical_storage": {  
        "type": "CASSANDRA",  
        "cluster": "clstr1",  
        "dataset": "ks1",  
        "table": "tbl1",  
    }  
}
```



N

Storage At Any Scale

Key Concepts

Idempotency

Chunking

Compression

Pagination

Adaptivity

Signaling

Key-Value

API

Storage

Time-Series

API

Storage

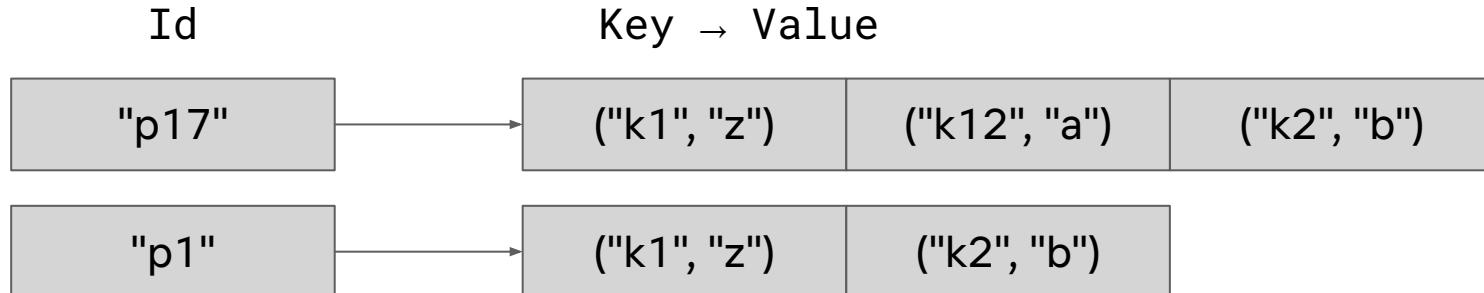
Key-Value Abstraction APIs

Server gRPC Abstraction API

Map[String, List[Item]]

Client High-Level API

HashMap[String, SortedMap[Bytes, Bytes]]



KeyValue

HashMap[String, SortedMap[Bytes, Bytes]]

```
rpc PutItems(  
    # Allows retry  
    IdempotencyToken token,  
    # Primary key  
    String namespace,  
    String id,  
    # Key-Value pairs  
    List<Item> items  
) -> PutItemsResponse
```

```
PutItemsResponse {  
    Trilean durable;  
    Trilean visible; ...  
}  
  
Item {  
    bytes key;  
    bytes value;  
    # Zero implies commit  
    int32 chunk;  
}
```

KeyValue

HashMap[String, SortedMap[Bytes, Bytes]]

```
rpc MutateItems(  
    IdempotencyToken token,  
    String namespace,  
    String id,  
    List[MutationRequest]  
    mutations, ...  
) -> MutationResponse
```

```
MutationRequest {  
    oneof {  
        PutItems put;  
        DeleteItems delete;  
    }  
}
```

KeyValue

HashMap[String, SortedMap[Bytes, Bytes]]

```
rpc GetItems(  
    # Primary key  
    String namespace,  
    String id,  
  
    # What matches  
    Predicate predicate,  
  
    # Pagination/Selection  
    Selection selection  
) -> GetItemsResponse  
  
GetItemsResponse {  
    # Page of results  
    # 1 MiB pages  
    List[Item] items  
  
    # If set there is  
    # more data to read  
    String next_page  
}
```

KeyValue

HashMap[String, SortedMap[Bytes, Bytes]]

```
rpc ScanItems(  
    # Table name  
    String namespace,  
  
    # What matches  
    Predicate predicate,  
  
    # Pagination & selection  
    Selection selection  
) -> ScanItemsResponse  
  
ScanItemsResponse {  
    # Page of results  
    List[ScanResult] results;  
  
    # Multiple concurrent  
    # pages consumed  
    List[String] next_page;
```

Writes

- All writes sent with **IdempotencyToken** so we can deduplicate
- Large values **chunked**, staged, and finally committed

Reads

- All **reads return pages** within the SL0
- No SL0 across all pages
- **Fixed size work**, not count!

Storage At Any Scale

Key Concepts

Idempotency

Chunking

Compression

Pagination

Adaptivity

Signaling

Key-Value

API

Storage

Time-Series

API

Storage

Key-Value Storage Layout

Base Table

id	key	value	value_metadata
1234	key1	value1	{"metadata": ...}
1234	key2	<empty>	{"metadata": ... }

```
CREATE TABLE IF NOT EXISTS <ns>.<ns> (
    id          text,
    key         blob,
    value       blob,
    value_metadata blob,
PRIMARY KEY (id, key))
WITH CLUSTERING ORDER BY (key <ASC|DESC>)
AND compaction = LCS(256MiB)
AND compression = LZ4(16KiB)
```

Key-Value Storage Layout

Data Table

id	bucket	key	chunk	version	value	value_metadata
1234	64200	key2	1	timeuuid1	chunk1	{ "metadata": ... }
1234	64201	key2	2	timeuuid1	chunk2	{ "metadata": ... }
1234	64201	key2	3	timeuuid1	chunk2	{ "metadata": ... }

```
CREATE TABLE IF NOT EXISTS <ns>.data_<ns> (
    id          text,
    bucket      int,
    key         blob,
    chunk       int,
    version     timeuuid,
    value        blob,
    value_metadata blob,
    PRIMARY KEY ((id, bucket), key, version, chunk))
    WITH CLUSTERING ORDER BY (key <ASC|DESC>, version DESC, chunk ASC)
```

Data Model

Key-Value bucketing per key

id	key	value	value_metadata
1234	key1	<empty>	chunk_count: 2 chunk_offset: 64201
6789	key2	<empty>	chunk_count: 3 chunk_offset: 34001
12	key1	"small"	{}

id	bucket	key	chunk	version	value	value_meta data
1234	64201	key1	1	timeuuid1	chunk value 1	{"metadata": ... }
1234	64201	key1	2	timeuuid1	chunk value 2	{"metadata": ... }
3443	64202	key1	1	timeuuid3	chunk value 1	{"metadata": ... }
6789	34001	key2	1	timeuuid2	chunk value 1	{"metadata": ... }
6789	34001	key2	2	timeuuid2	chunk value 2	{"metadata": ... }
6789	34001	key2	3	timeuuid2	chunk value 3	{"metadata": ... }

```
"chunks": {  
    "chunk_size_bytes": 65536,  
    "chunk_count": 2,  
    "chunk_offset": 124431,  
    "chunk_token_xxh64": "12595968139585192830",  
    "version": "3d52c5c0-6229-11ee-aecd-de49621d037e"  
}
```

Chunking: Spread the Load

```
offset = consistentHash(xxHash(key), NUM_BUCKETS)  
bucket = offset + (chunk / 8)
```

(id, key) -> value

("1", "a") -> 1MiB

("1", "b") -> 2MiB

("2", "c") -> 10MiB



((id, bucket), key, chunk)

Chunking: Spread the Load

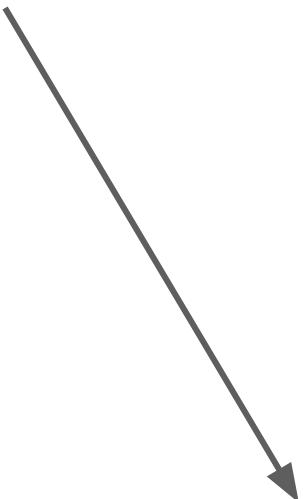
```
offset = consistentHash(xxHash(key), NUM_BUCKETS)  
bucket = offset + (chunk / 8)
```

(id, key) -> value

("1", "a") -> 1MiB

("1", "b") -> 2MiB

("2", "c") -> 10MiB



((id, bucket), key, chunk)

(("1" , 1072), "a" , 1)
...
(("1" , 1074), "a" , 16)

Chunking: Spread the Load

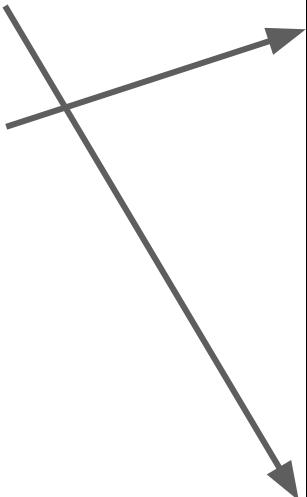
```
offset = consistentHash(xxHash(key), NUM_BUCKETS)
bucket = offset + (chunk / 8)
```

(id, key) -> value

("1", "a") -> 1MiB

("1", "b") -> 2MiB

("2", "c") -> 10MiB



((id, bucket), key, chunk)

((("1", 2), "b", 1)
...
(("1", 6), "b", 32))

((("1", 1072), "a", 1)
...
(("1", 1074), "a", 16))

Chunking: Spread the Load

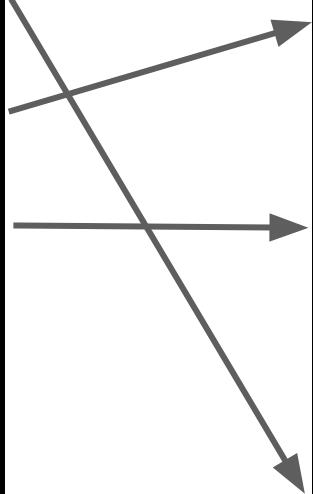
```
offset = consistentHash(xxHash(key), NUM_BUCKETS)
bucket = offset + (chunk / 8)
```

(id, key) -> value

("1", "a") -> 1MiB

("1", "b") -> 2MiB

("2", "c") -> 10MiB



((id, bucket), key, chunk)

(("1" , 2), "b" , 1)

...
(("1" , 6), "b" , 32)

(("2" , 12), "c" , 1)

...
(("2" , 32), "c" , 160)

(("1" , 1072), "a" , 1)

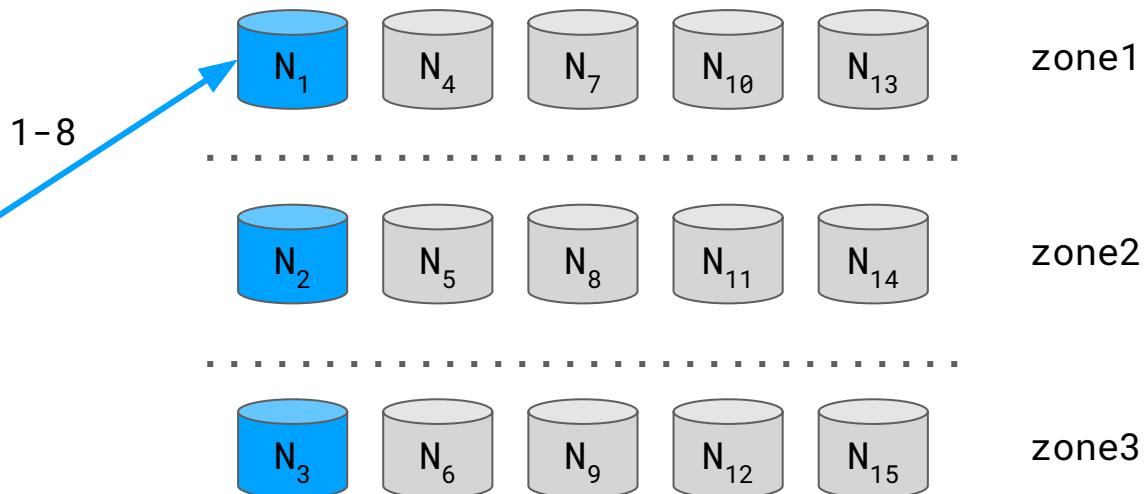
...
(("1" , 1074), "a" , 16)

Chunking: Spread the Load

Chunks Distributed Across the Cluster

(id, key) -> value

("2", "c") -> 10MiB

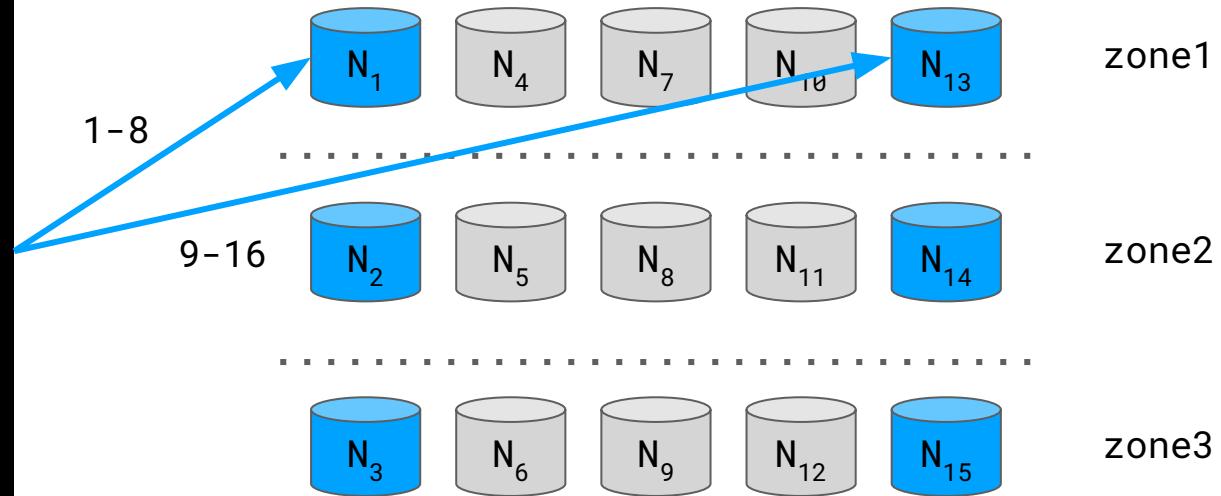


Chunking: Spread the Load

Chunks Distributed Across the Cluster

(id, key) -> value

("2", "c") -> 10MiB

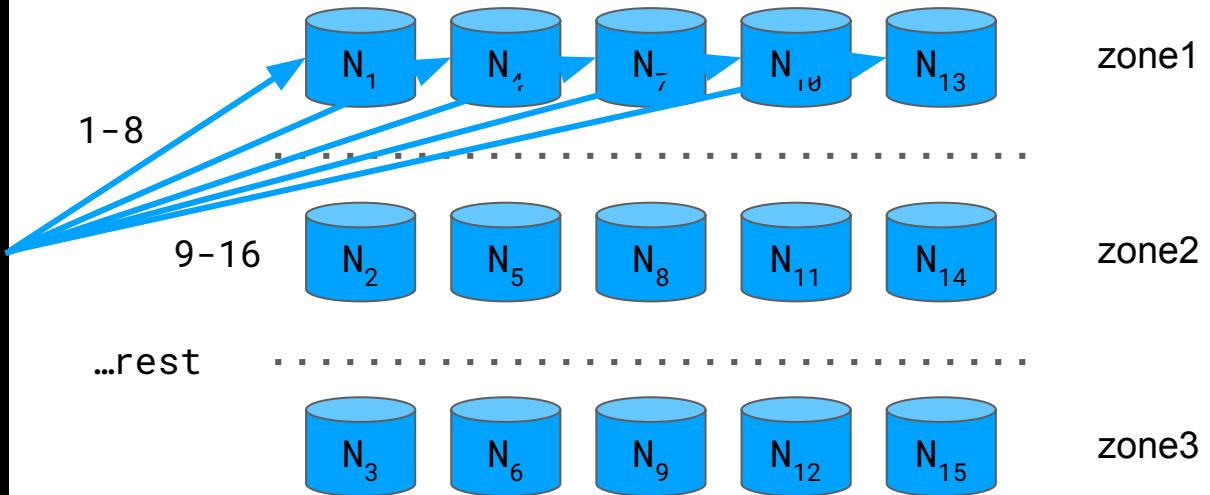


Chunking: Spread the Load

Chunks Distributed Across the Cluster

(id, key) -> value

("2", "c") -> 10MiB



Versioning

Key-Value concurrency through versions

id	key	value	value_metadata
1234	key1	<empty>	{version: v2}
6789	key2	<empty>	{version: timeuuid()}

id	bucket	key	chunk	versio n	value	value_meta data
1234	64201	key1	1	V1	chunk value 1	{"metadata": ...}
1234	64201	key1	2	V1	chunk value 2	{"metadata": ...}
⋮						
1234	34001	key1	1	V2	chunk value 1	{"metadata": ...}
1234	34001	key1	2	V2	chunk value 2	{"metadata": ...}
1234	34002	key1	3	V2	chunk value 3	{"metadata": ...}

```
"chunks": {  
    "chunk_size_bytes": 65536,  
    "chunk_count": 2,  
    "chunk_offset": 124431,  
    "chunk_token_xxh64": "12595968139585192830",  
    "version": "3d52c5c0-6229-11ee-aecd-de49621d037e"  
}
```

Storage At Any Scale

Key Concepts

Idempotency

Chunking

Compression

Pagination

Adaptivity

Signaling

Key-Value

API

Storage

Time-Series

API

Storage

TimeSeries Abstraction APIs

Server gRPC Abstraction API

Map[String, List[EventRecord]]

Client High-Level API

**HashMap[String,
SortedMap[[Timestamp, UUID], Event]]**

Time Series ID

(Event Time, Event ID) -> Event

"p17"	(t_4, U_1, v_1)	(t_2, U_2, v_2)	(t_1, U_3, v_3)
"p1"	(t_{272}, U_4, v_4)	(t_4, U_1, v_1)	

TimeSeries `HashMap[String,
 SortedMap[[Timestamp, String], Event]]`

`rpc WriteEventRecords[Sync](rpc ReadEventRecords(
 List[EventRecord] events
) -> (Durable, Visible)`

`message EventRecord {
 String namespace,
 String time_series_id,
 Timestamp time,
 String event_id,
 List[EventItem] items,
 }`

`) -> Page`

`rpc SearchEventRecords(
 String namespace,
 SearchQuery query
) -> Page`

TimeSeries

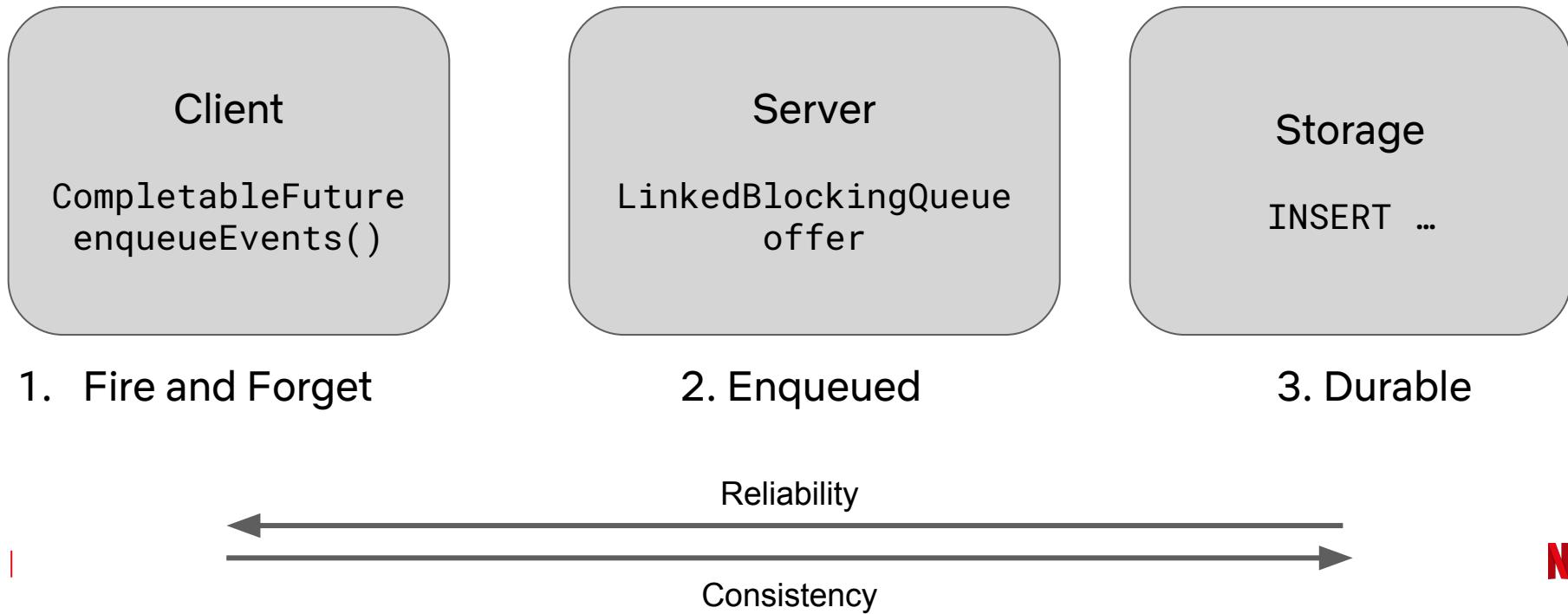
```
HashMap[String,  
        SortedMap[[Timestamp, String], Event]]
```

```
message EventRecord {  
    String namespace,  
    String time_series_id,  
    # Idempotency Token  
    Timestamp time,  
    String event_id,  
    List[EventItem] items,  
}
```

```
message  
ReadEventRecordResponse {  
    StringValue  
    next_page_token;  
  
    # The page of results  
    List[EventRecord] records;  
}
```

TimeSeries Three-Ack Options

HashMap[String,
SortedMap[[Timestamp, String], Event]]



TimeSeries

HashMap[String,
SortedMap[[Timestamp, String], Event]]

Writes

- All writes sent with components to construct an **IdempotencyToken**
- Three modes of acknowledging
- Large events prohibited (>4MiB)

Reads

- All **reads return pages** of events within the SL0
- No SL0 across all pages
- **Fixed size work**, not count!

Storage At Any Scale

Key Concepts

Idempotency

Chunking

Compression

Pagination

Adaptivity

Signaling

Key-Value

API

Storage

Time-Series

API

Storage

Why take on the Complexity?

Save millions of dollars on computers

Flexible operations

Simplify massive scale event datasets
for developers

Why take on the Complexity?

Save millions of dollars on computers
(compute efficiency)

Flexible operations
(operator efficiency)

Simplify massive scale event datasets
for developers
(developer efficiency)

Time-Series Storage Layout

Metadata table

namespace name	version	start time	end time	metadata	status
ns1	2022_01_01	2023-01-01	2023-01-02	{<config>}	ACTIVE
ns1	2022_01_02	2023-01-02	2023-01-03	{<config>}	ACTIVE
ns2	2022_01_02	2022-01-02	2022-01-07	{<config>}	CLOSED

```
CREATE TABLE IF NOT EXISTS <ns>.metadata (
    namespace_name text,
    version        text,
    start_time     timestamp,
    end_time       timestamp,
    // how many time or event buckets to make
    metadata       map<text, text>,
    status         text,
PRIMARY KEY (namespace_name, version)
) WITH CLUSTERING ORDER BY (version DESC)
```

Time-Series Storage Layout

Metadata table

namespace_name	version	start_time	end_time	status
_history_service	data_20240721_00_42	1721522555888000	1729298555888000	ACTIVE
_history_service	data_20240422_00_42	1713746555888000	1721522555888000	ACTIVE
_history_service	data_20240123_00_42	1705970555888000	1713746555888000	ACTIVE
_history_service	data_20231025_00_42	1698194555888000	1705970555888000	ACTIVE
_history_service	data_20230727_00_42	1690418555888000	1698194555888000	ACTIVE
_history_service	data_20230428_00_42	1682642555888000	1690418555888000	ACTIVE
_history_service	data_20230128_00_42	1674866555888000	1682642555888000	ACTIVE
_history_service	data_20221030_00_42	1667090555888000	1674866555888000	ACTIVE
_history_service	data_20220801_00_42	1659314555888000	1667090555888000	ACTIVE
_history_service	data_20220503_00_42	1651538555888000	1659314555888000	DELETED
_history_service	data_20220202_00_42	1643762555888000	1651538555888000	DELETED
_history_service	data_20211104_00_42	1635986555888000	1643762555888000	DELETED
_history_service	data_20210806_00_42	1628210555888000	1635986555888000	DELETED
_history_service	data_20210508_00_42	1620434555888000	1628210555888000	DELETED
_history_service	data_20210207_00_42	1612658555888000	1620434555888000	DELETED
_history_service	data_20201109_00_42	1604882555888000	1612658555888000	DELETED
_history_service	data_20200811_00_42	1597106555888000	1604882555888000	DELETED
_history_service	data_20200513_00_42	1589330555888000	1597106555888000	DELETED
_history_service	data_20200213_00_42	1581554555888000	1589330555888000	DELETED
_history_service	FIRST_SLICE	1651538555888000	1651538555888000	ACTIVE

Time Slice - Shard by Table

```
persistence_prod cass_dgw_ts_ [REDACTED]_history_service--useast1d us-east-1  
(nfsuper) /mnt/data/cassandra/data/impression_history_service $ du -shc *  
0      data_20200213_00_42-0b91ad00ca7a11ec9d3efb8c1dedee89  
0      data_20200513_00_42-08216570ca7a11ecbadeb1d99c6b1cff  
0      data_20200811_00_42-04200170ca7a11ecbb66ed9e09af8275  
0      data_20201109_00_42-00b8ba90ca7a11ecad1f2944ed95fca9  
0      data_20210207_00_42-fd41bc40ca7911ec9b39db076ee72b0d  
0      data_20210508_00_42-f947ab40ca7911ecb296e56a17d142b4  
0      data_20210806_00_42-f621da30ca7911ecaa9759e45316053e  
0      data_20211104_00_42-f239a380ca7911eca7b30d5f8869e00b  
0      data_20220202_00_42-ee1993f0ca7911ec95395fdea5a24fa8  
0      data_20220503_00_42-5a5d8ee0e35c11eca081256038c5d02f  
1.5T    data_20220801_00_42-5dc216a0e35c11ecb8acff3beb89a17d  
1.2T    data_20221030_00_42-62011fe0e35c11ecaca0153c136468d2  
1.2T    data_20230128_00_42-90270f431f4435869fc376edf07e0538  
1.1T    data_20230428_00_42-530d343cd8b73fb79344af761425f827  
1.1T    data_20230727_00_42-04b376573be33ab683a8d17bc6daf455  
572G   data_20231025_00_42-5e9ffebc21413ce8b094802352417830  
0      data_20240123_00_42-afb6482796343a2bb6e81be80ad37c17  
0      data_20240422_00_42-7e17ccc11656396a893cffa94586bcce  
0      data_20240721_00_42-4878c769892433ee8c6674683b7821aa  
0      metadata-8dfc79b0ca7911ec901035819e30d34e  
6.5T    total
```

Slice State Machine

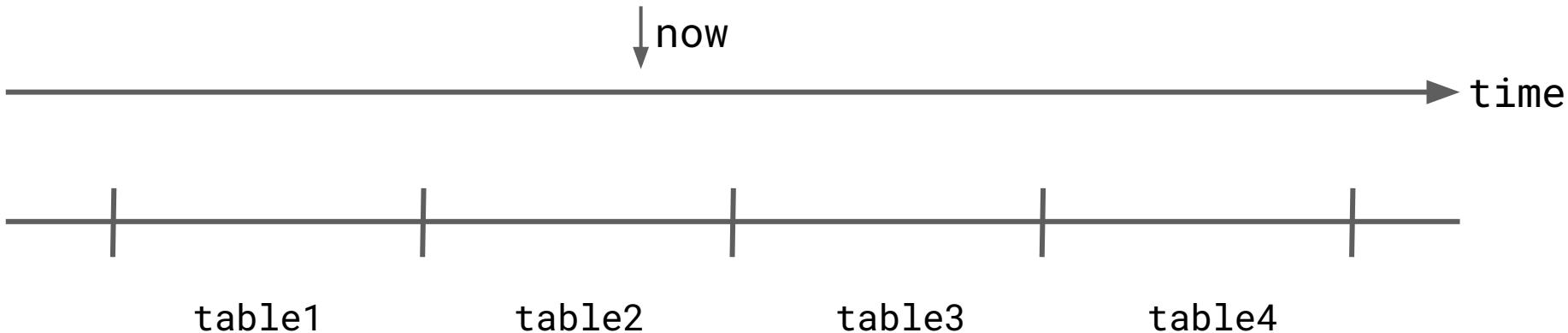
PENDING

ACTIVE

CLOSED

DELETED

Time Slice - Shard by Table



Advantages over TWCS?

- Use up to 80-90% disk
- Decouple retention from TTL
- Decouple background operations: compression, compaction

Pro-Tips

- Create "runway" in future
- Idempotent CREATE WITH id = <hash>

Time-Series Storage Layout

Data table

time series id	time bucket	event bucket	event time	event id	value	value_metadata
1234	19682	1	ts1	uuid1	event1	{ metadata: ... }
1234	19681	2	ts2	uuid2	event2	{ metadata: ... }
5678	89681	1	ts1	uuid3	event3	{ metadata: ... }

```
CREATE TABLE IF NOT EXISTS <ns>.<version> (
    time_series_id text,
    time_bucket     bigint,
    event_bucket   int,
    event_time      timestamp,
    event_id        text,
    value           blob,
    value_metadata blob,
    PRIMARY KEY(
        (time_series_id, time_bucket, event_bucket), event_time, event_id))
    WITH CLUSTERING ORDER BY (event_time DESC, event_id DESC, key DESC)
```

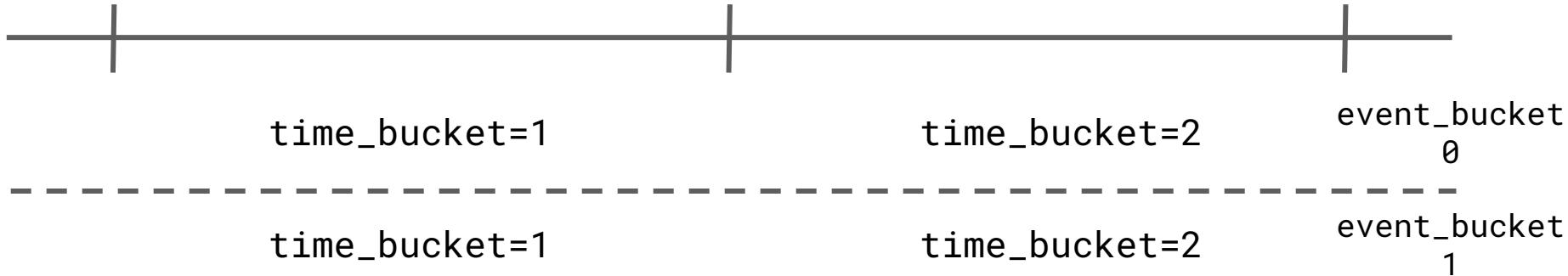
Time-Series Storage Configuration

metadata

```
{'buckets_per_id': '2', 'buffer_capacity': '4194304', 'buffer_coalesce': '1', 'seconds_per_bucket': '5', 'seconds_per_interval': '7776000'}  
{'buckets_per_id': '2', 'buffer_capacity': '4194304', 'buffer_coalesce': '1', 'seconds_per_bucket': '5', 'seconds_per_interval': '7776000'}  
{'buckets_per_id': '2', 'buffer_capacity': '4194304', 'buffer_coalesce': '1', 'seconds_per_bucket': '5', 'seconds_per_interval': '7776000'}  
{'buckets_per_id': '2', 'buffer_capacity': '4194304', 'buffer_coalesce': '1', 'seconds_per_bucket': '5', 'seconds_per_interval': '7776000'}
```

now

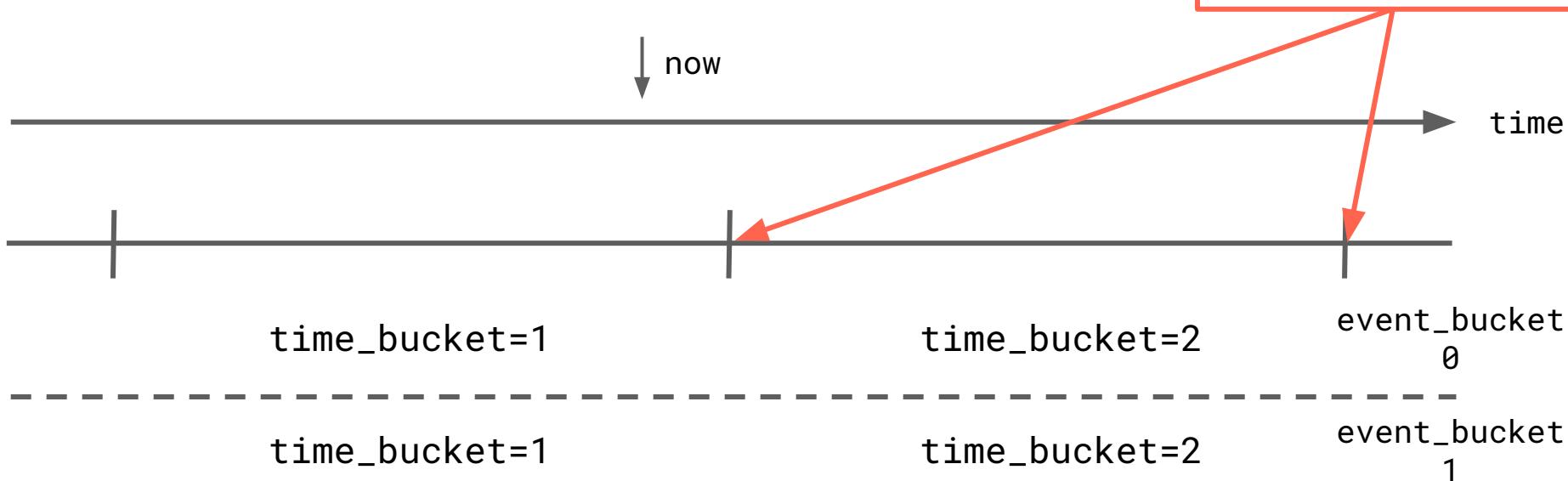

time

Time-Series Storage Configuration

metadata

```
{'buckets_per_id': '2', 'buffer_capacity': '4194304', 'buffer_coalesce': '1', 'seconds_per_bucket': '5', 'seconds_per_interval': '7776000'}  
{'buckets_per_id': '2', 'buffer_capacity': '4194304', 'buffer_coalesce': '1', 'seconds_per_bucket': '5', 'seconds_per_interval': '7776000'}  
{'buckets_per_id': '2', 'buffer_capacity': '4194304', 'buffer_coalesce': '1', 'seconds_per_bucket': '5', 'seconds_per_interval': '7776000'}  
{'buckets_per_id': '2', 'buffer_capacity': '4194304', 'buffer_coalesce': '1', 'seconds_per_bucket': '5', 'seconds_per_interval': '7776000'}
```



Time-Series Storage Configuration

metadata

```
{'buckets_per_id': '2', 'buffer_capacity': '4194304', 'buffer_coalesce': '1', 'seconds_per_bucket': '5', 'seconds_per_interval': '7776000'}  
{'buckets_per_id': '2', 'buffer_capacity': '4194304', 'buffer_coalesce': '1', 'seconds_per_bucket': '5', 'seconds_per_interval': '7776000'}  
{'buckets_per_id': '2', 'buffer_capacity': '4194304', 'buffer_coalesce': '1', 'seconds_per_bucket': '5', 'seconds_per_interval': '7776000'}  
{'buckets_per_id': '2', 'buffer_capacity': '4194304', 'buffer_coalesce': '1', 'seconds_per_bucket': '5', 'seconds_per_interval': '7776000'}
```

now

time

time_bucket=1

time_bucket=2

event_bucket
=0

time_bucket=1

time_bucket=2

event_bucket
=1

Time-Series Storage Configuration

metadata

```
{'buckets_per_id': '2', 'buffer_capacity': '4194304', 'buffer_coalesce': '1', 'seconds_per_bucket': '5', 'seconds_per_interval': '7776000'}  
{'buckets_per_id': '2', 'buffer_capacity': '4194304', 'buffer_coalesce': '1', 'seconds_per_bucket': '5', 'seconds_per_interval': '7776000'}  
{'buckets_per_id': '2', 'buffer_capacity': '4194304', 'buffer_coalesce': '1', 'seconds_per_bucket': '5', 'seconds_per_interval': '7776000'}  
{'buckets_per_id': '2', 'buffer_capacity': '4194304', 'buffer_coalesce': '1', 'seconds_per_bucket': '5', 'seconds_per_interval': '7776000'}
```

now

time

time_bucket=1

time_bucket=2

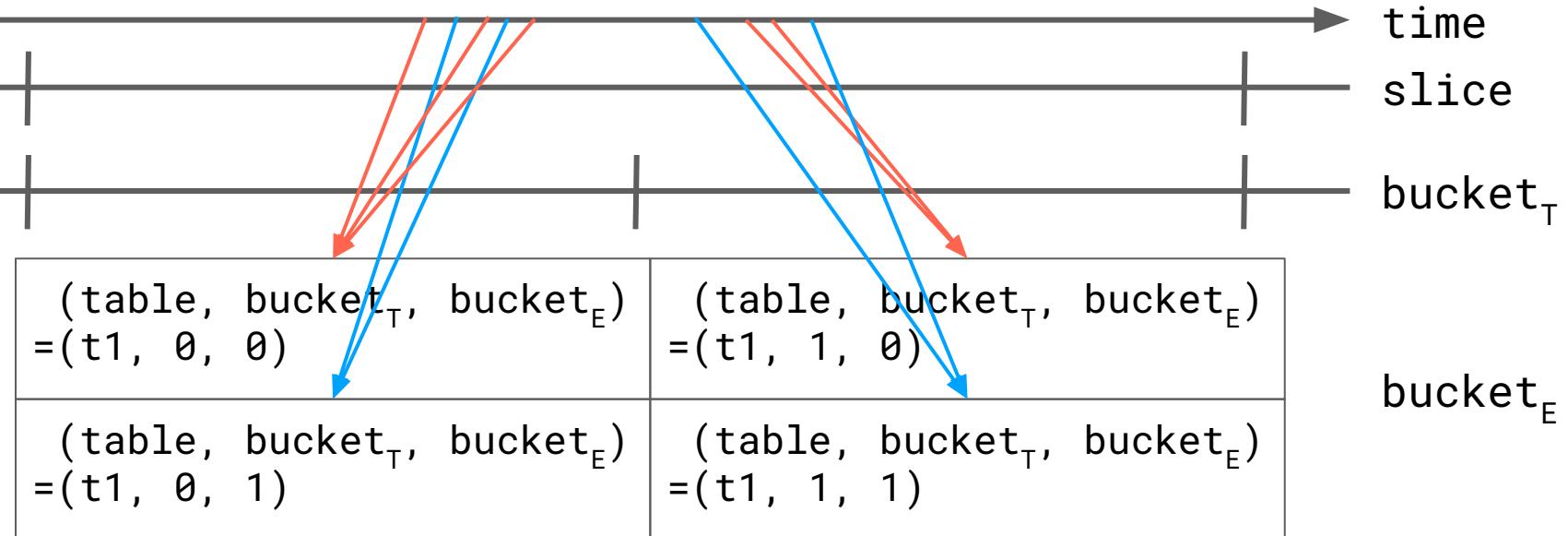
event_bucket
=0

time_bucket=1

time_bucket=2

event_bucket
=1

Time Partition - Tunable



Why So Many Buckets?

Flexible Buckets

Tune for your use case!

Tunable	Primary Purpose	Tradeoff
Tables (slice_interval)	Retention	DDL in Critical Path
Time Buckets (time_bucket)	Read amplification	Lookback latency
Random Event Buckets (event_bucket)	Sharding wide partitions	Parallel fan-out

Key-Value Storage Layout

Data Table - Metadata buckets

id	bucket	key	chunk	version	value	value_metadata
1234	-2	key2	6	timeuuid1	<empty>	{ "metadata": ... }
1234	-2	key2	12	timeuuid2	<empty>	{ "metadata": ... }

Data Table

id	bucket	key	chunk	version	value	value_metadata
1234	64200	key2	1	timeuuid1	chunk1	{ "metadata": ... }
1234	64201	key2	2	timeuuid1	chunk2	{ "metadata": ... }
1234	78902	key2	1	timeuuid2	chunk1	{ "metadata": ... }

Key-Value Storage Layout

Writes

- **Stage** chunks to data table and **commit** to base table after chunks are staged.
- **IdempotencyToken** is used for concurrent writes
- Small values are committed to **base table only**

Reads

- Read base table for version and number of chunks
- Retrieve chunks adaptively
- Fallback to older version if chunks data is not replicated

Success Stories

Cloud Saves



Tracing



Profile Level User Interaction (PLUID)



Cloud Saves

Netflix **Game Data**
streaming scale Key-Value
Data

30 - 300MB of game
progress metadata per call

Enabled by **chunking and
compression** features in
Key-Value

| **Key-Value deploys - 400+ shards, 3000+ use-cases**
Time-Series has 10+ PiB of events stored across 10s of use-cases

Tracing

Netflix **Distributed Tracing**
data

1000 of services sending
IPC traces

Stores 1-2 PiB of
uncompressed data

Immutable events

Profile Level User Interaction (PLUID)

Video + Profile
play/pause/stop events

Key-Value for active
sessions and Time-Series
for session history

800k + 10k Writes
150k + 80k Reads
8 hour TTL + 1 year
retention

Tuned for heavy write
operations

Further Work

Summarization



Re-sharding



Compare and Swap



Thank You.

Vidhya Arvind
varvind@netflix.com

Joey Lynch
josephl@netflix.com

N