# TUNING FLINK FOR ROBUSTNESS AND PERFORMANCE

STEFAN RICHTER

@STEFANRRICHTER

SEPT 4, 2018

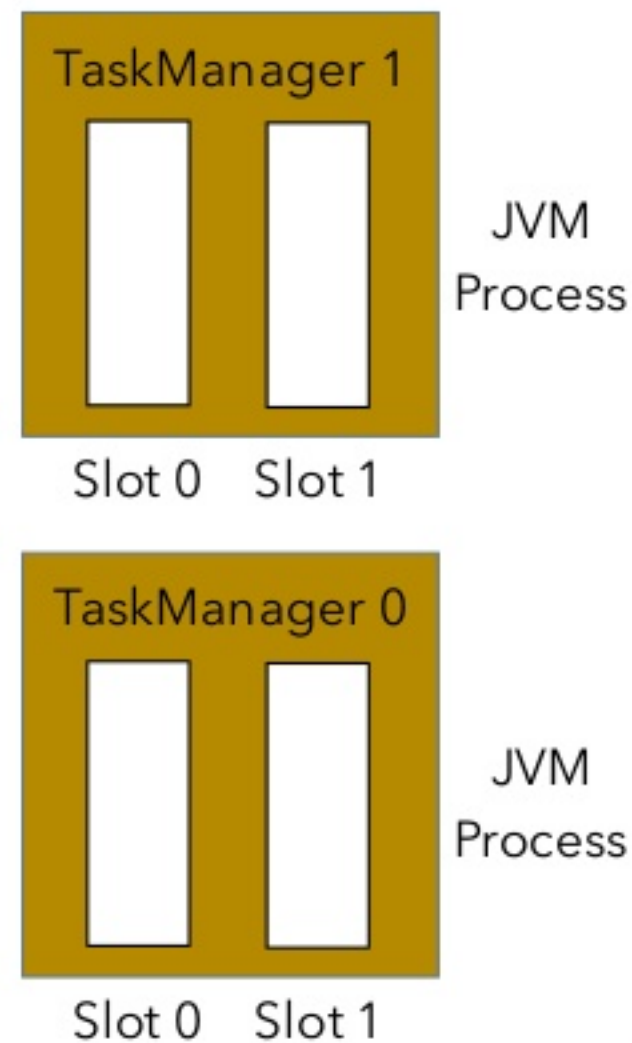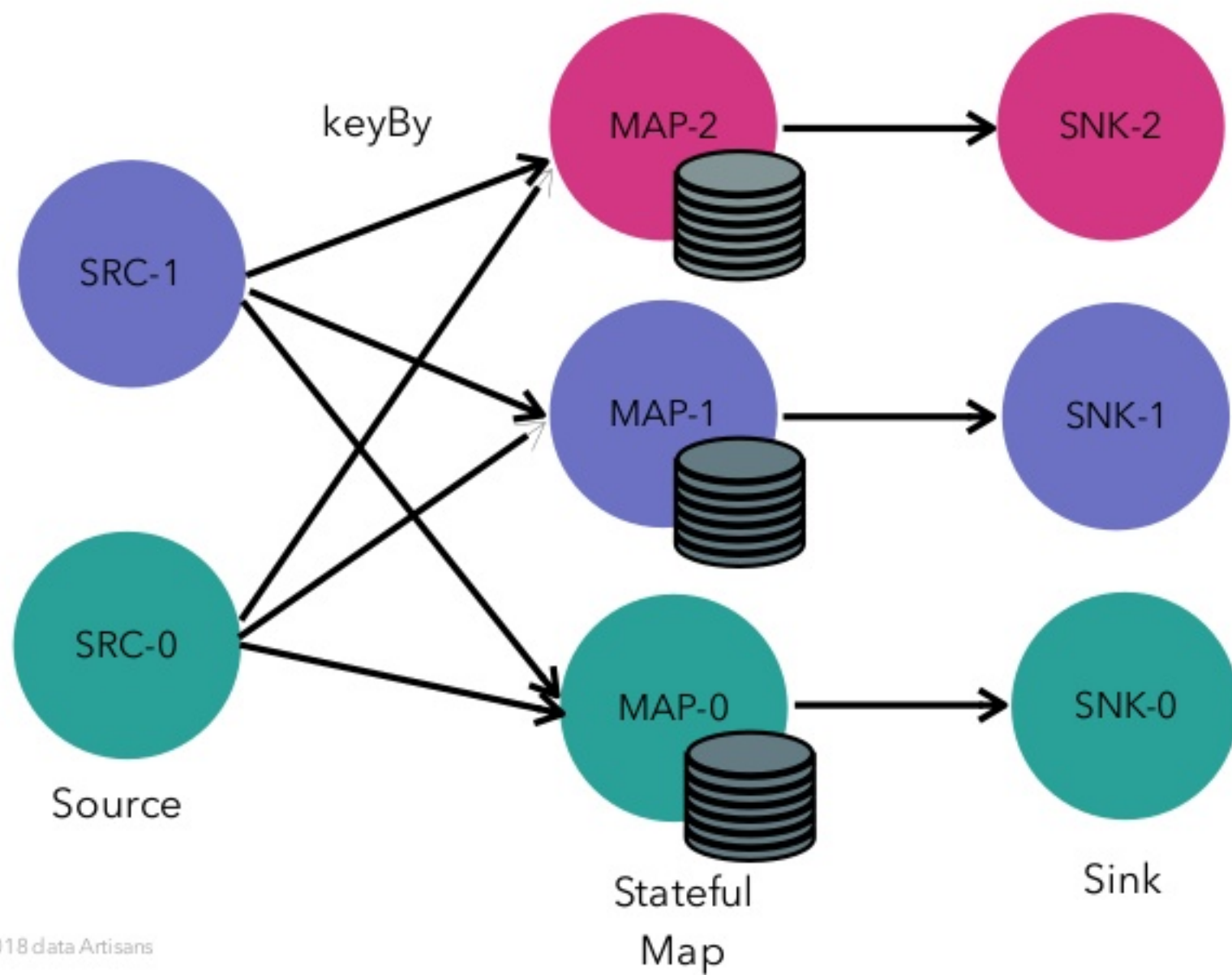data Artisans

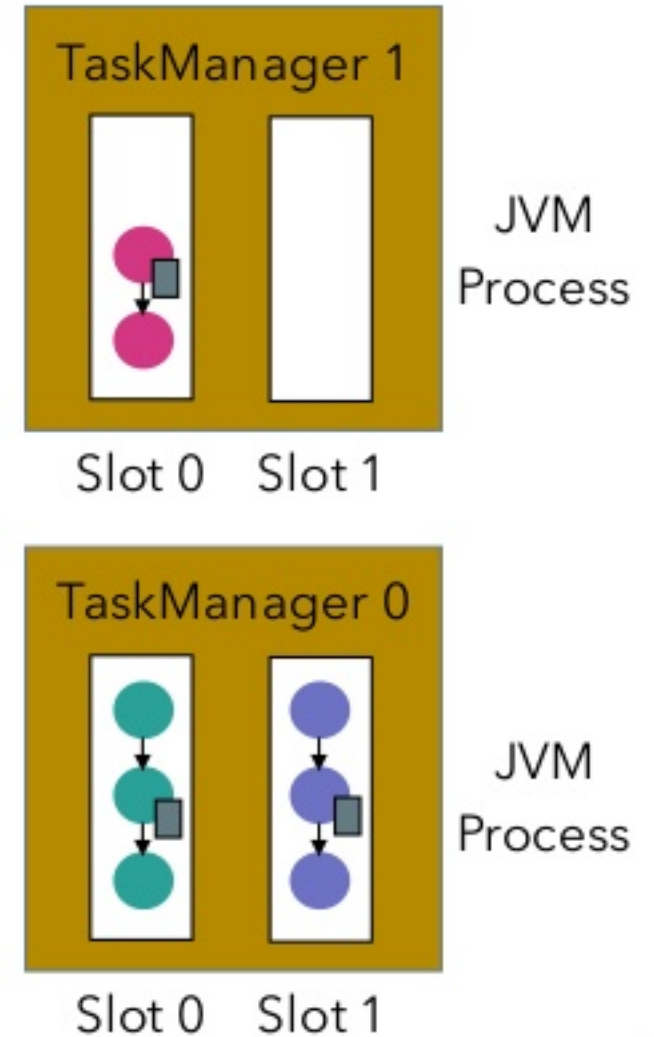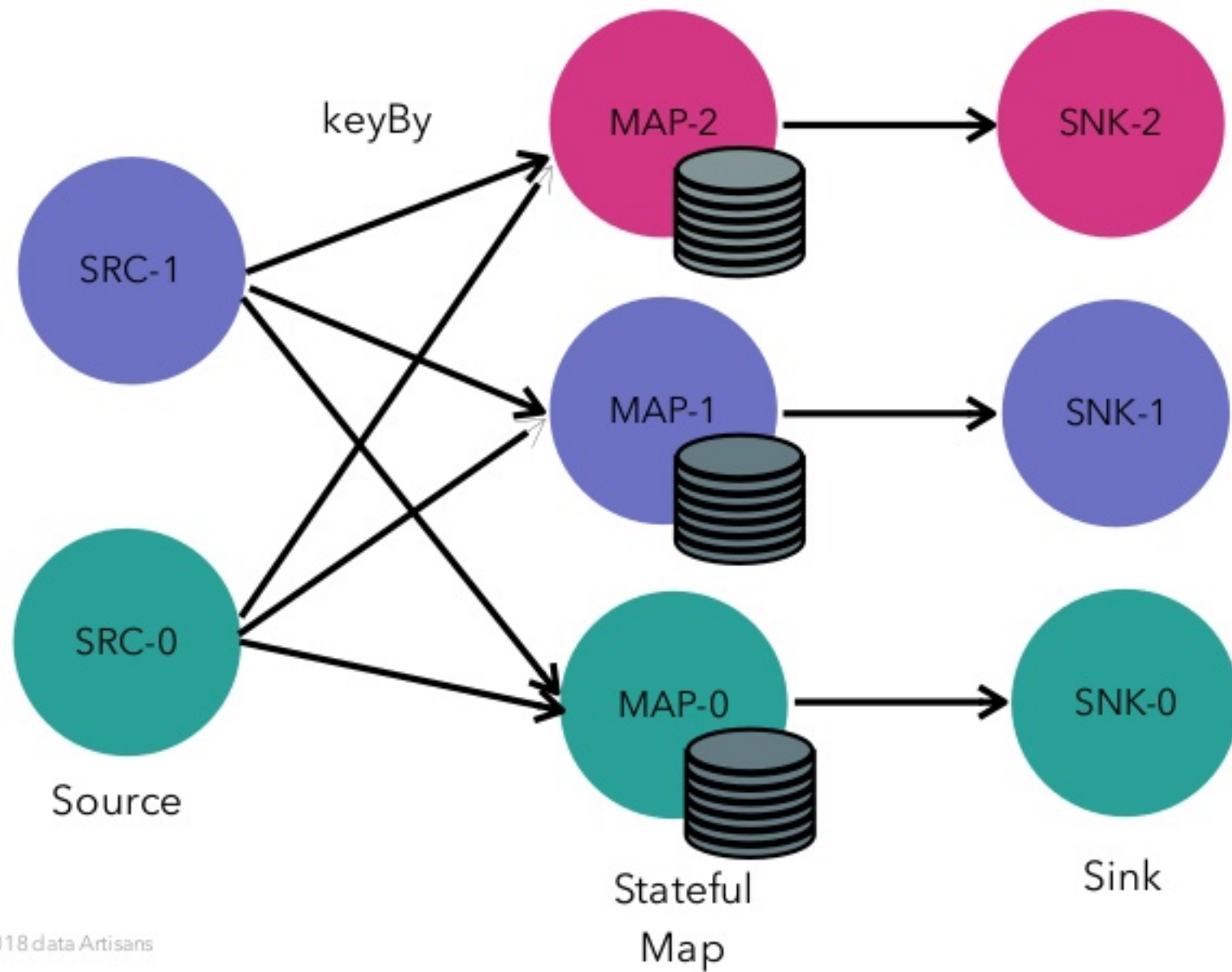# GENERAL MEMORY CONSIDERATIONS

# BASIC TASK SCHEDULING



© 2018 data Artisans

# BASIC TASK SCHEDULING



© 2018 data Artisans

# TASK MANAGER PROCESS MEMORY LAYOUT

Typical Size

Flink Framework etc.

Network Buffers

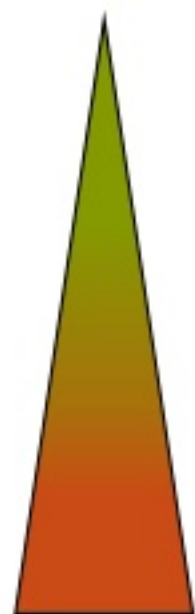Timer State

Keyed State

Task Manager JVM Process

Java Heap

Off Heap / Native

# TASK MANAGER PROCESS MEMORY LAYOUT

Typical Size

Flink Framework etc.

Network Buffers

Timer State

Keyed State

Task Manager JVM Process

Java Heap
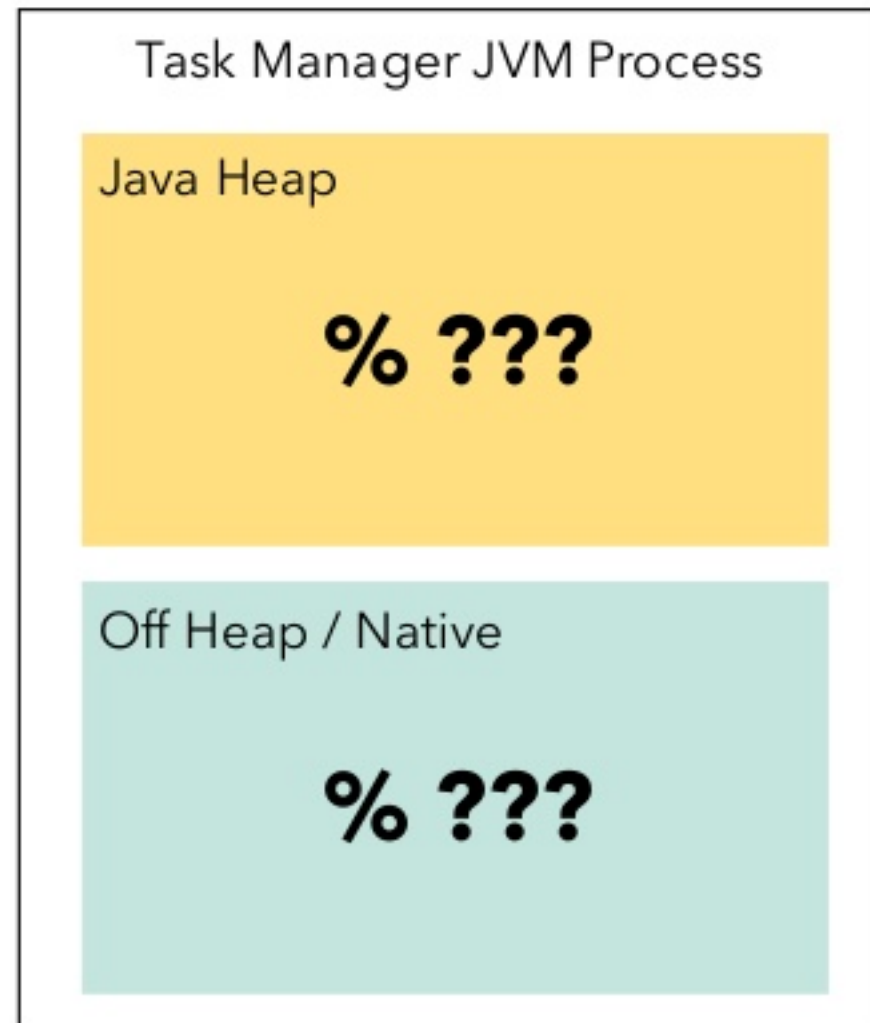
**% ???**

Off Heap / Native

**% ???**

**Σ ???**

# TASK MANAGER PROCESS MEMORY LAYOUT

Typical Size

Timer State

Keyed State

Task Manager JVM Process

Java Heap

Flink Framework etc.

Off Heap / Native

Network Buffers

# TASK MANAGER PROCESS MEMORY LAYOUT



Typical Size

Timer State

Keyed State

Task Manager JVM Process

Java Heap

Flink Framework etc.

Off Heap / Native
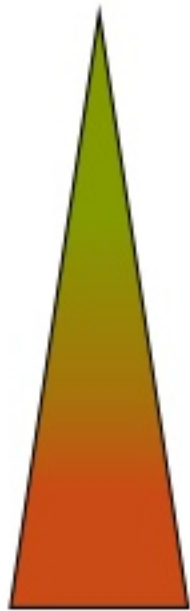
Network Buffers

# STATE BACKENDS

# FLINK KEYED STATE BACKENDS CHOICES

**Based on Java Heap Objects**

**Based on RocksDB**

# HEAP KEYED STATE BACKEND CHARACTERISTICS

- State lives as Java objects on the heap.

- Organized as chained hash table, key $\mapsto$ state.

- One hash table per registered state.

- Supports asynchronous state snapshots through copy-on-write MVCC.

- Data is de/serialized only during state snapshot and restore.


- Highest performance.

- Affected by garbage collection overhead / pauses.

- Currently no incremental checkpoints.

- Memory overhead of representation.

- State size limited by available heap memory.

# HEAP STATE TABLE ARCHITECTURE



- Hash buckets (Object[]), 4B-8B per slot
-  Load factor <= 75%
-  Incremental rehash

# HEAP STATE TABLE ARCHITECTURE

Entry ◯

Entry ◯    Entry ◯

- 4 References:
  - Key
  - Namespace
  - State
  - Next
- 3 int:
  - Entry Version
  - State Version
  - Hash Code

K

N

S

Object sizes and overhead.
Some objects might be shared.

- Hash buckets (Object[]), 4B-8B per slot
- Load factor <= 75%
- Incremental rehash

4 x (4B-8B)
+ 3 x 4B
+ ~8B-16B (Object overhead)

# HEAP STATE TABLE SNAPSHOT MVCC



Copy of hash bucket array is snapshot overhead

# HEAP STATE TABLE SNAPSHOT MVCC

No conflicting modification = no overhead

# HEAP STATE TABLE SNAPSHOT MVCC



Original

Snapshot

Modifications trigger deep copy of entry - only as much as required. This depends on what was modified and what is immutable (as determined by type serializer).
Worst case overhead = size of original state table at time of snapshot.

# HEAP BACKEND TUNING CONSIDERATIONS

- Chose type serializer with efficient copy-method (for copy-on-write).

- Flag immutability of objects where possible to avoid copy completely.

- Flatten POJOs / avoid deep objects. Reduces object overheads and following references = potential cache misses.

- GC choice/tuning can help. Follow future GC developments.

- Scale-out using multiple task manager per node to support larger state over multiple heap backends rather than having fewer and large heaps.

# ROCKSDB KEYED STATE BACKEND CHARACTERISTICS

- State lives as serialized byte-strings in off-heap memory and on local disk.
- Key-Value store, organized as log-structured merge tree (LSM-tree).
    - Key: serialized bytes of <Keygroup, Key, Namespace>.
    - Value: serialized bytes of the state.
- One column family per registered state (~table).
- LSM naturally supports MVCC.
- Data is de/serialized on every read and update.
- Not affected by garbage collection.
- Relative low overhead of representation.
- LSM naturally supports incremental snapshots.
- State size limited by available local disk space.
- Lower performance (~order of magnitude compared to Heap state backend).

# ROCKSDB ARCHITECTURE

Memory

Persistent Store

In Flink:
- disable WAL and sync
- persistence via checkpoints

WriteOp

Active
MemTable

WAL

WAL

# ROCKSDB ARCHITECTURE



Memory

WriteOp

Active MemTable

Full/Switch

ReadOnly MemTable

Persistent Store

WAL

WAL

Flush

Local Disk

SST SST

SST SST

Merge

Compaction

In Flink:
- disable WAL and sync
- persistence via checkpoints

# ROCKSDB ARCHITECTURE



**Memory**

WriteOp

Active MemTable

Full/Switch

ReadOnly MemTable

Flush

Set per column family (~table)

**Persistent Store**

WAL

WAL

**Local Disk**

SST    SST

SST    SST

Merge

Compaction
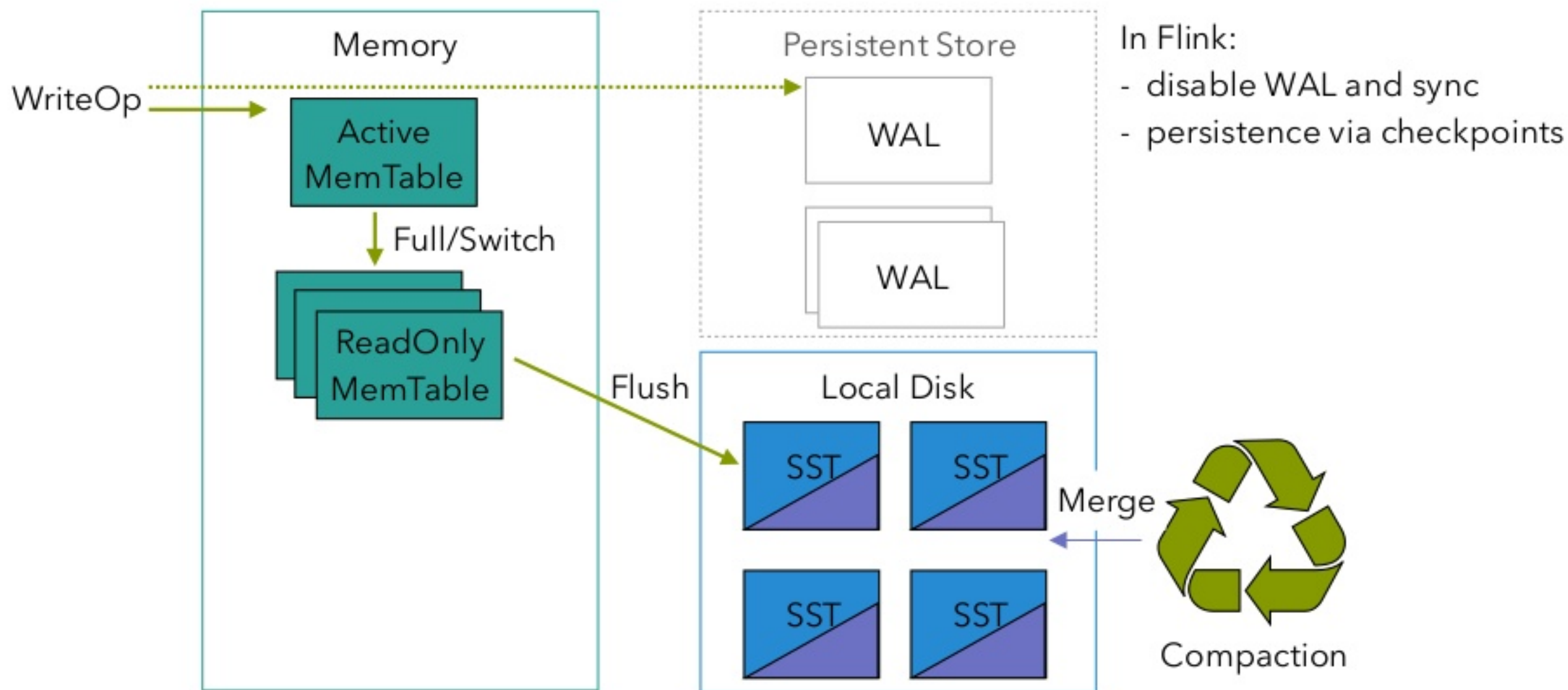
In Flink:
- disable WAL and sync
- persistence via checkpoints

# ROCKSDB ARCHITECTURE

Memory

WriteOp

Active
MemTable

Full/Switch

ReadOp

ReadOnly
MemTable

Flush

Persistent Store

WAL

WAL

In Flink:
- disable WAL and sync
- persistence via checkpoints

Local Disk

SST

SST

SST

SST

Merge

Compaction

# ROCKSDB ARCHITECTURE



Memory

WriteOp

Active
MemTable

Full/Switch

ReadOp

ReadOnly
MemTable

Read Only
Block Cache

Persistent Store

WAL

WAL

In Flink:
- disable WAL and sync
- persistence via checkpoints

Flush

Local Disk

SST    SST

SST    SST

Merge

Compaction

# ROCKSDB RESOURCE CONSUMPTION

- One RocksDB instance per keyed operator subtask.
- **block_cache_size**:
  - Size of the block cache.
- **write_buffer_size**:
  - Max. size of a MemTable.
- **max_write_buffer_number**:
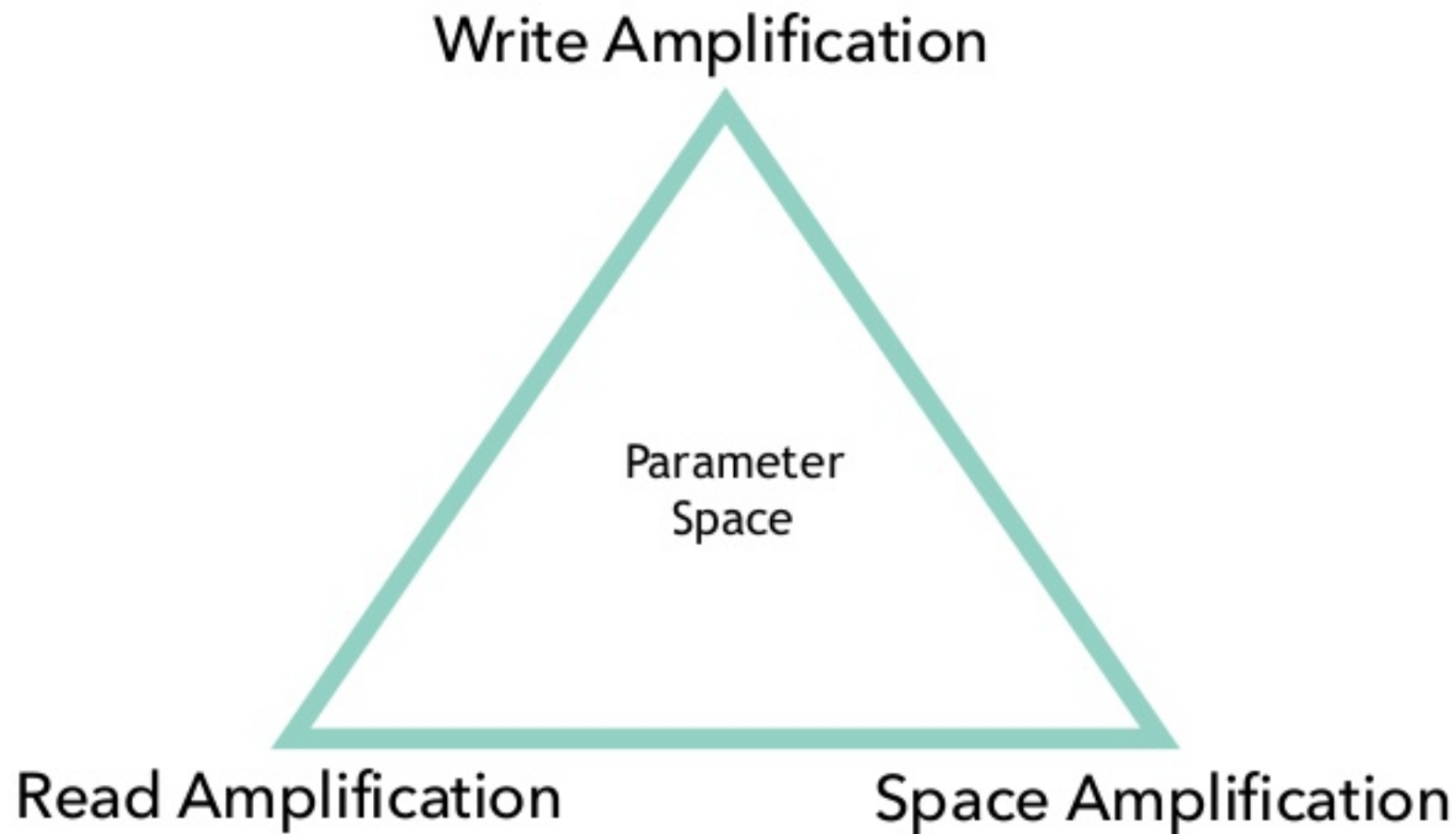  - The maximum number of MemTables in memory before flush to SST files.
- **Indexes and bloom filters** (optional).
- **Table Cache**:
  - Caches open file descriptors to SST files. Default: unlimited!

# PERFORMANCE TUNING - AMPLIFICATION FACTORS

Write Amplification

Parameter
Space

Read Amplification

Space Amplification

More details: https://github.com/facebook/rocksdb/wiki/RocksDB-Tuning-Guide

# PERFORMANCE TUNING - AMPLIFICATION FACTORS

Write Amplification

**Example**: More compaction effort =
increased write amplification
and reduced read amplification

Parameter
Space

Read Amplification

Space Amplification

More details: https://github.com/facebook/rocksdb/wiki/RocksDB-Tuning-Guide

# GENERAL PERFORMANCE CONSIDERATIONS

- Efficient type serializer and serialization formats.

- Decompose user-code objects: business logic / efficient state representation.

- Extreme: „Flightweight Pattern", e.g. wrapper object that interprets/manipulates stored byte array on the fly and uses only byte-array type serializer.

- File Systems:

  - Working directory on fast storage, ideally local SSD. Could even be memory file system because it is transient for Flink. EBS performance can be problematic.

  - Checkpoint directory: Persistence happens here. Can be slower but should be fault tolerant.

# TIMER SERVICE

# HEAP TIMERS

Timer

- 2 References:
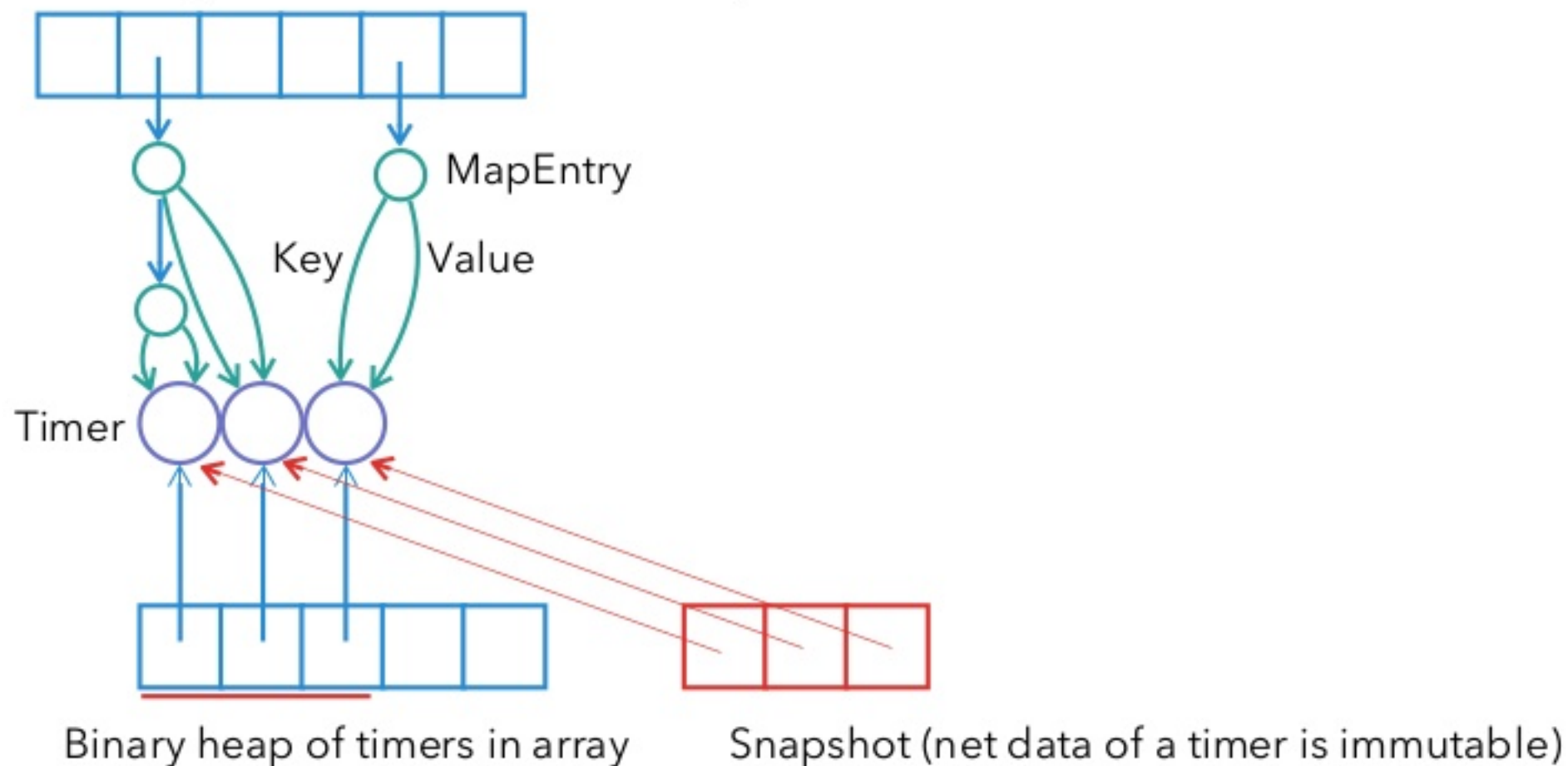  - Key
  - Namespace
- 1 long:
  - Timestamp
- 1 int:
  - Array Index

K

N

Object sizes and overhead.
Some objects might be shared.

Peek: O(1)
Poll: O(log(n))
Insert: O(1)/O(log(n))
Delete: O(n)
Contains O(n)

Binary heap of timers in array

# HEAP TIMERS

HashMap<Timer, Timer> : fast deduplication and deletes

MapEntry

Key    Value

Timer

Binary heap of timers in array

- 2 References:
  - Key
  - Namespace
- 1 long:
  - Timestamp
- 1 int:
  - Array Index

K

N

Object sizes and overhead.
Some objects might be shared.

Peek: O(1)
Poll: O(log(n))
Insert: O(1)/O(log(n))
Delete: O(log(n))
Contains O(1)

# HEAP TIMERS

HashMap<Timer, Timer> : fast deduplication and deletes

MapEntry

Key    Value

Timer

Binary heap of timers in array    Snapshot (net data of a timer is immutable)

# ROCKSDB TIMERS

Column Family - only key, no value

| Key Group | Time stamp | Key | Name space |
|---|---|---|---|
| 0 | 20 | A | X |
| 0 | 40 | D | Z |

...

| Key Group | Time stamp | Key | Name space |
|---|---|---|---|
| 1 | 10 | D | Z |
| 1 | 20 | C | Y |

...

| Key Group | Time stamp | Key | Name space |
|---|---|---|---|
| 2 | 50 | B | Y |
| 2 | 60 | A | X |

...

Lexicographically ordered
byte sequences as key, no value

# ROCKSDB TIMERS



Column Family - only key, no value

| Key Group | Time stamp | Key | Name space |
|-----------|-----------|-----|------------|
| 0 | 20 | A | X |
| 0 | 40 | D | Z |

...

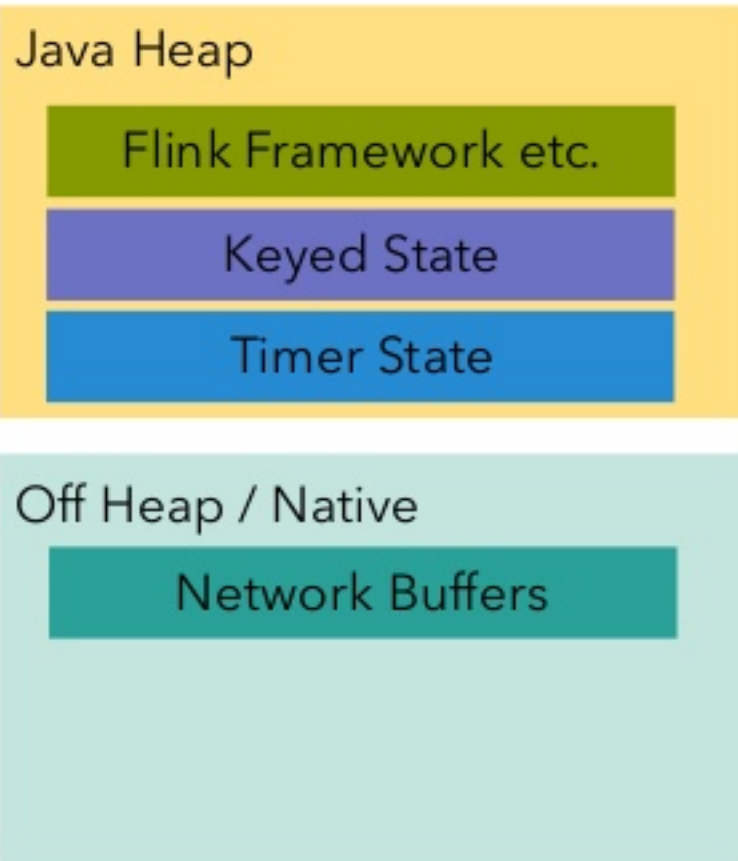| 1 | 10 | D | Z |
| 1 | 20 | C | Y |

...

| 2 | 50 | B | Y |
| 2 | 60 | A | X |

...

Key group queues
(caching first k timers)

Priority queue of
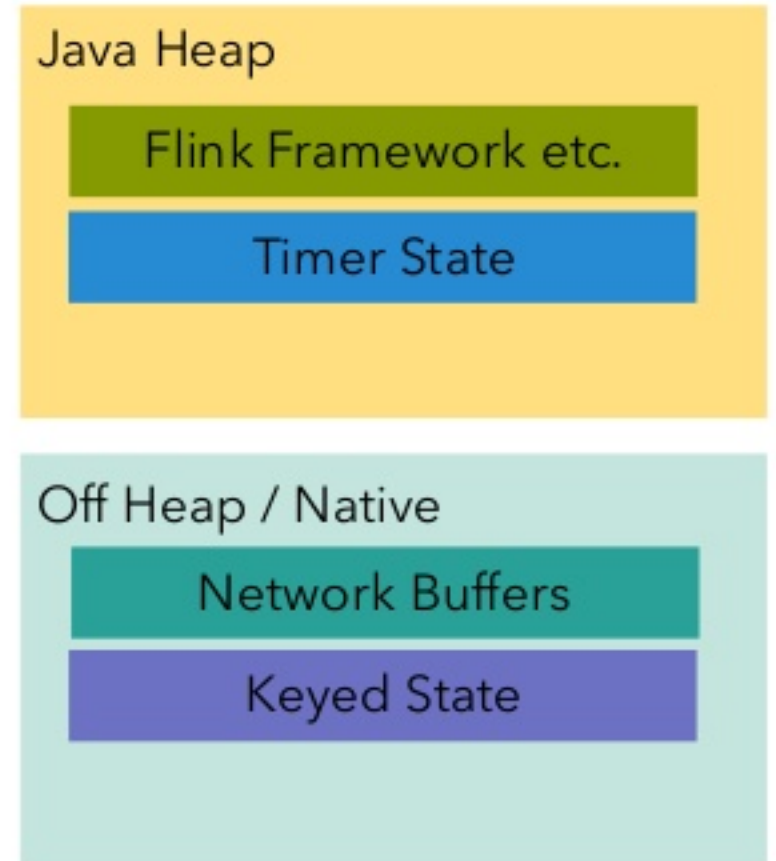key group queues

# 3 TASK MANAGER MEMORY LAYOUT OPTIONS

**Task Manager JVM Process**

Java Heap
- Flink Framework etc.
- Keyed State
- Timer State

Off Heap / Native
- Network Buffers

**Task Manager JVM Process**

Java Heap
- Flink Framework etc.

Off Heap / Native
- Network Buffers
- Keyed State
- Timer State

**Task Manager JVM Process**

Java Heap
- Flink Framework etc.
- Timer State

Off Heap / Native
- Network Buffers
- Keyed State

# FULL / INCREMENTAL CHECKPOINTS

# FULL CHECKPOINT

| | |
|---|---|
| 🟧 | A |
| 🟥 | B |
| ⬜ | C |
| 🟪 | D |

@t₁

| | |
|---|---|
| 🟧 | A |
| 🟥 | B |
| ⬜ | C |
| 🟪 | D |

Checkpoint 1

# FULL CHECKPOINT



Checkpoint 1                    Checkpoint 2

# FULL CHECKPOINT



Checkpoint 1          Checkpoint 2          Checkpoint 3

# FULL CHECKPOINT



Checkpoint 1          Checkpoint 2          Checkpoint 3
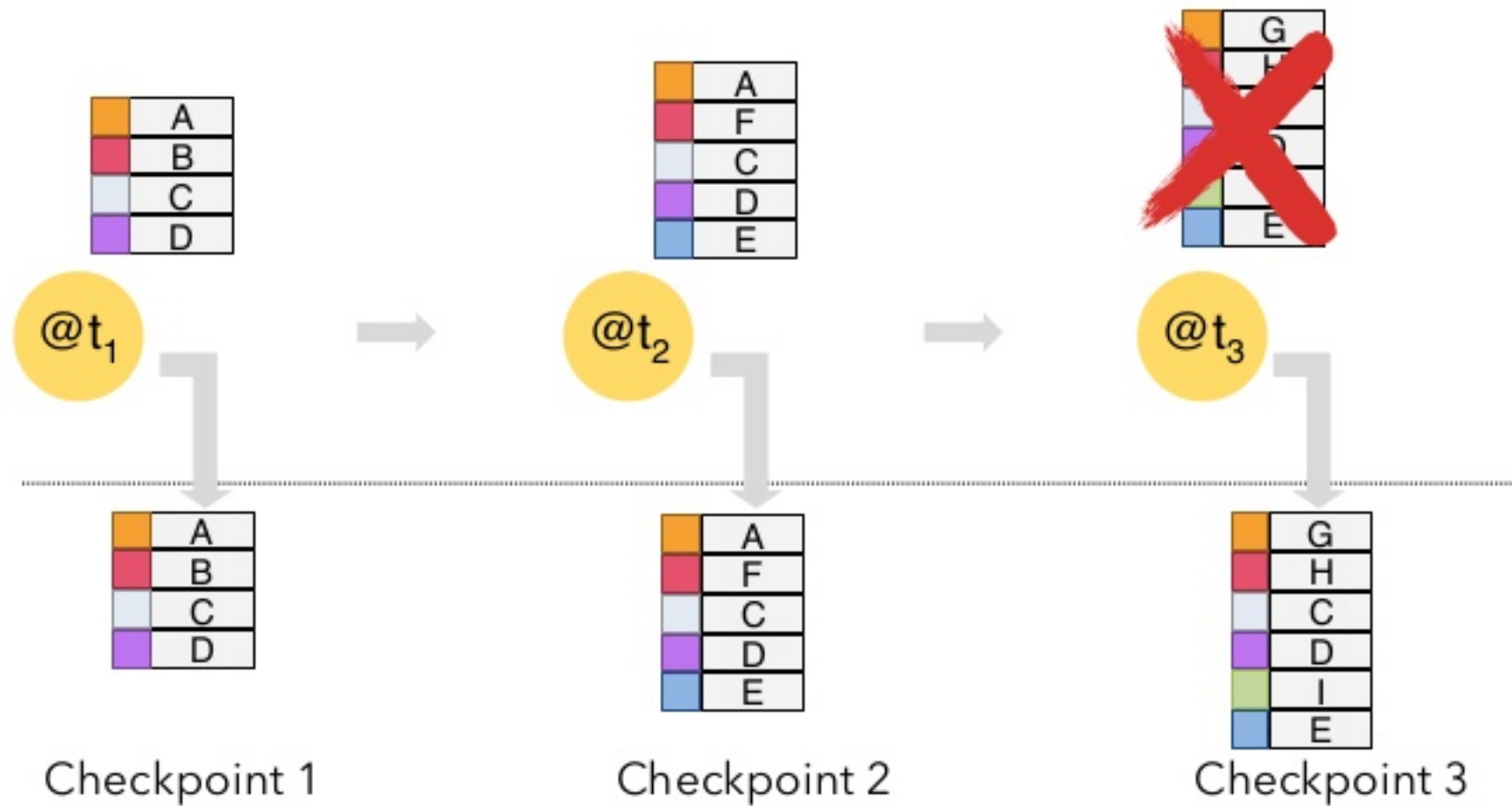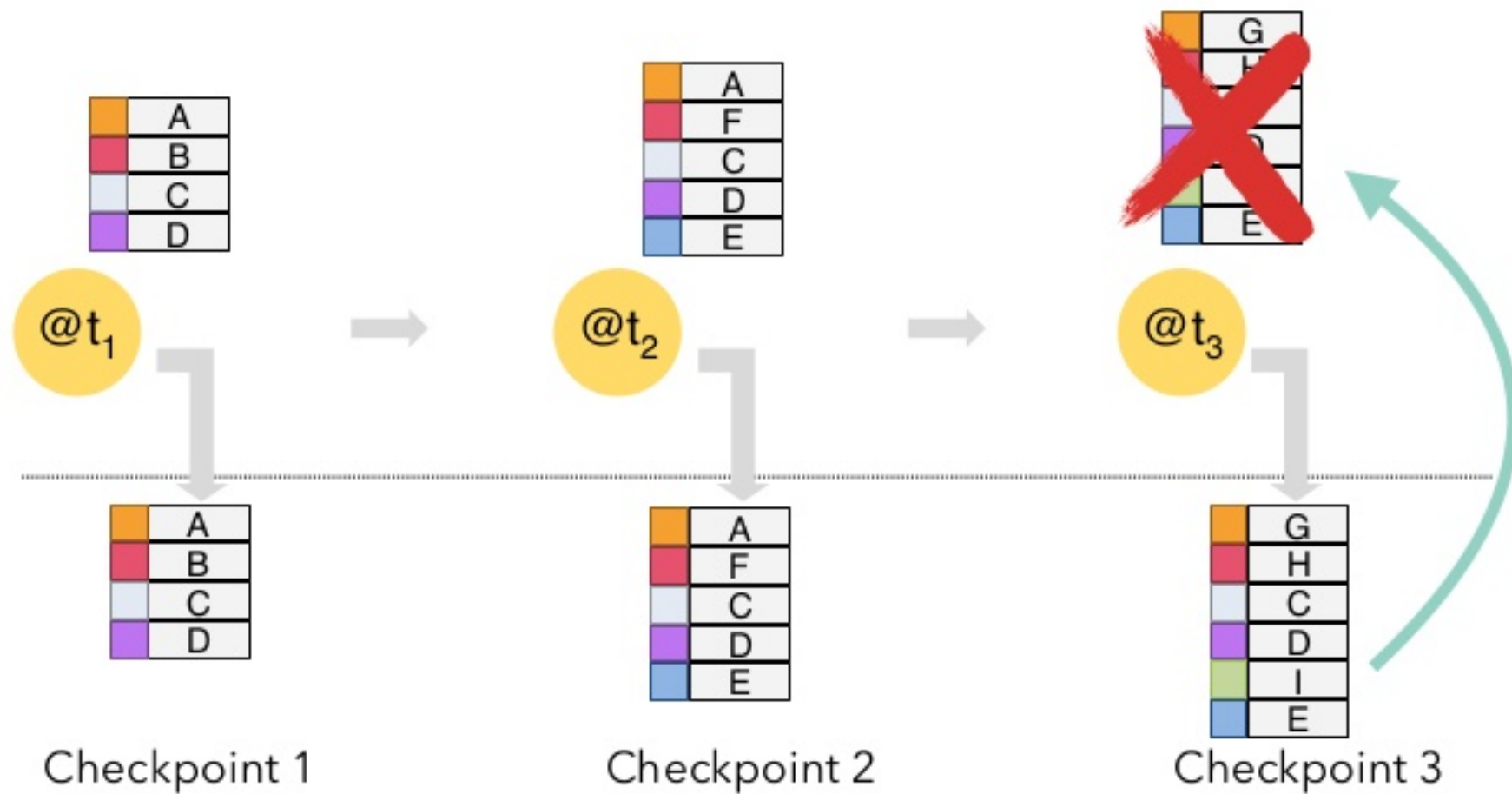
# FULL CHECKPOINT



Checkpoint 1          Checkpoint 2          Checkpoint 3
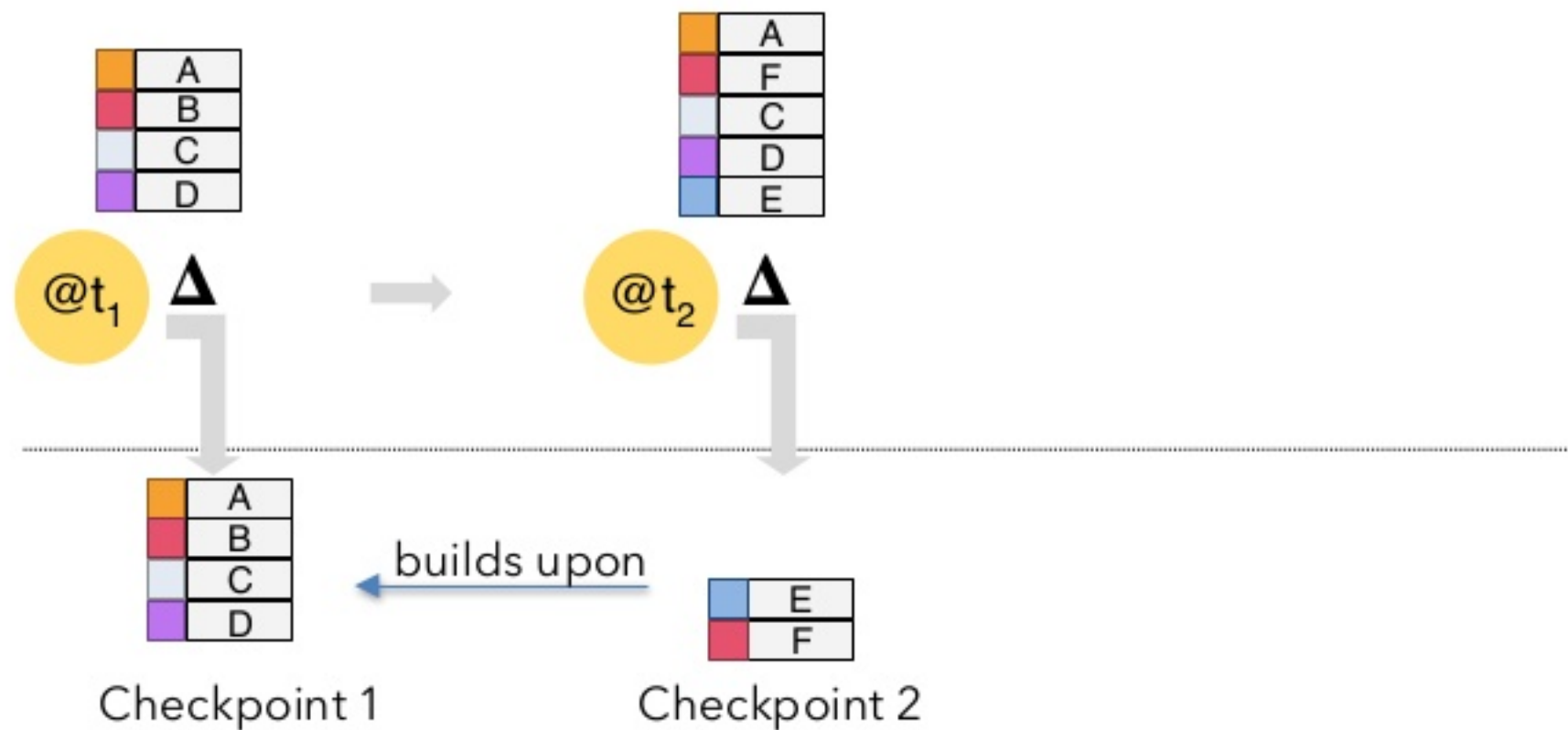
# FULL CHECKPOINT OVERVIEW

- Creation iterates and writes full database snapshot as stream to stable storage.

- Restore reads data as stream from stable storage and re-inserts into backend.

- Each checkpoint is self contained, size is proportional to size of full state.
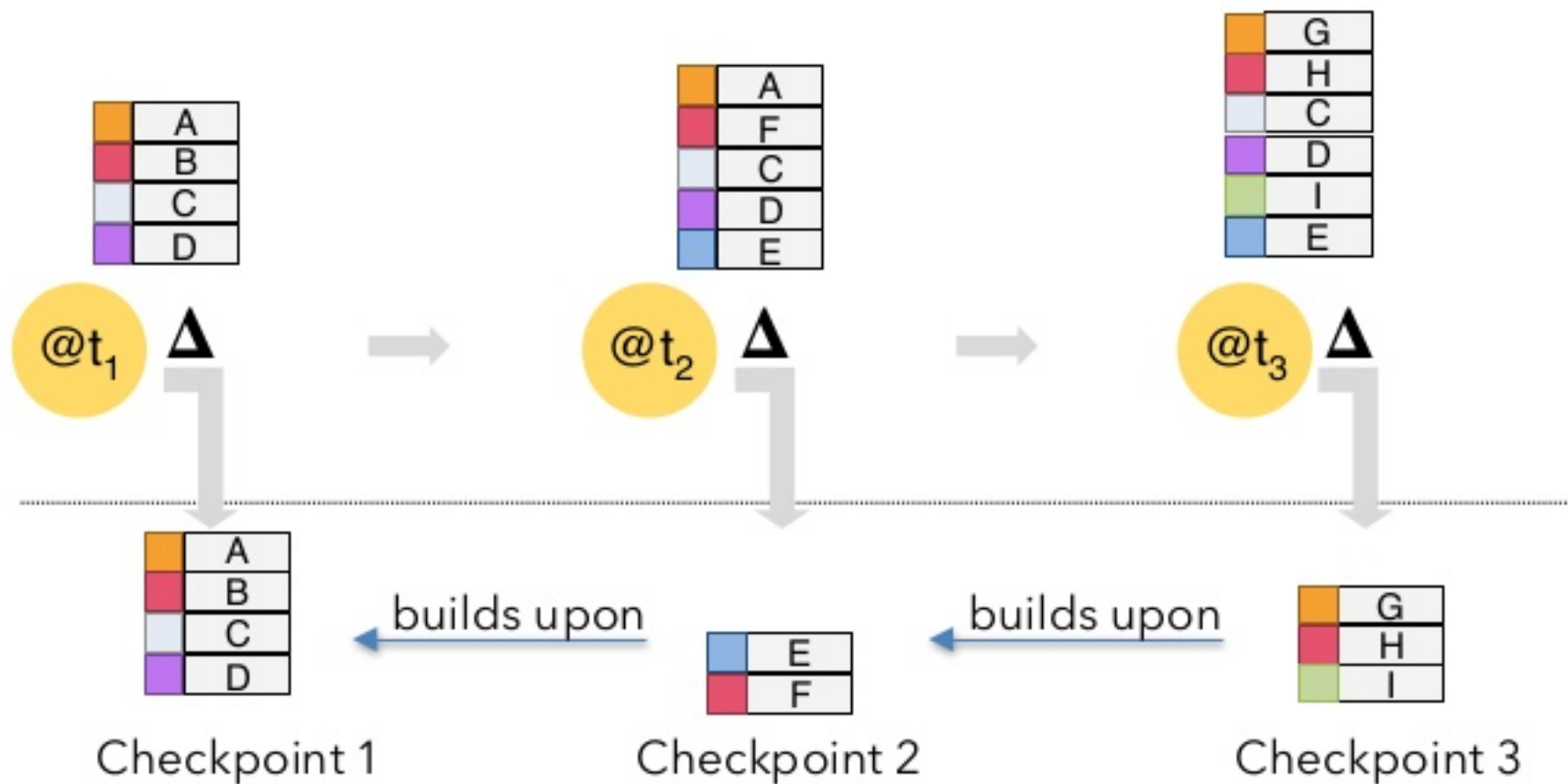
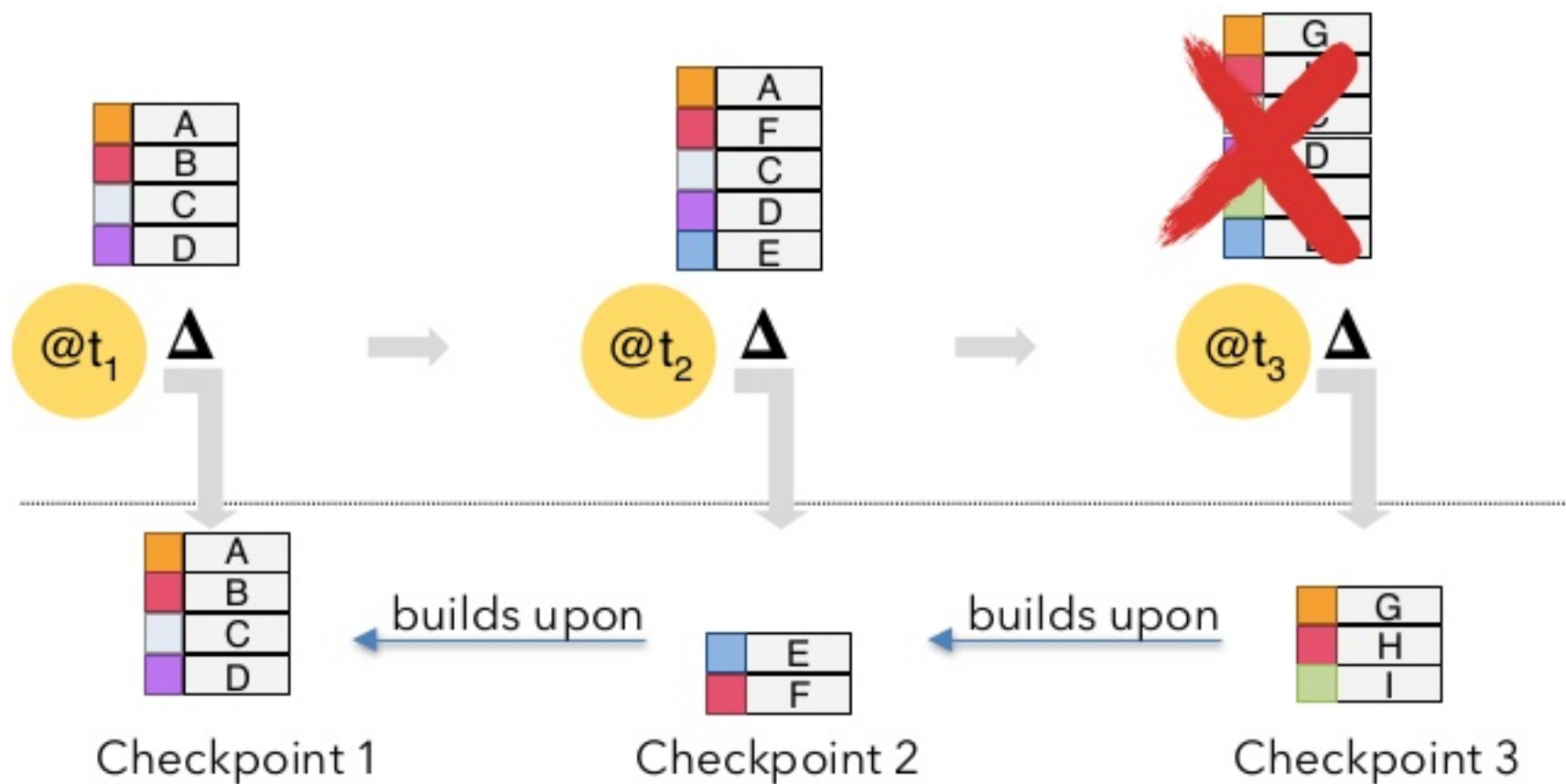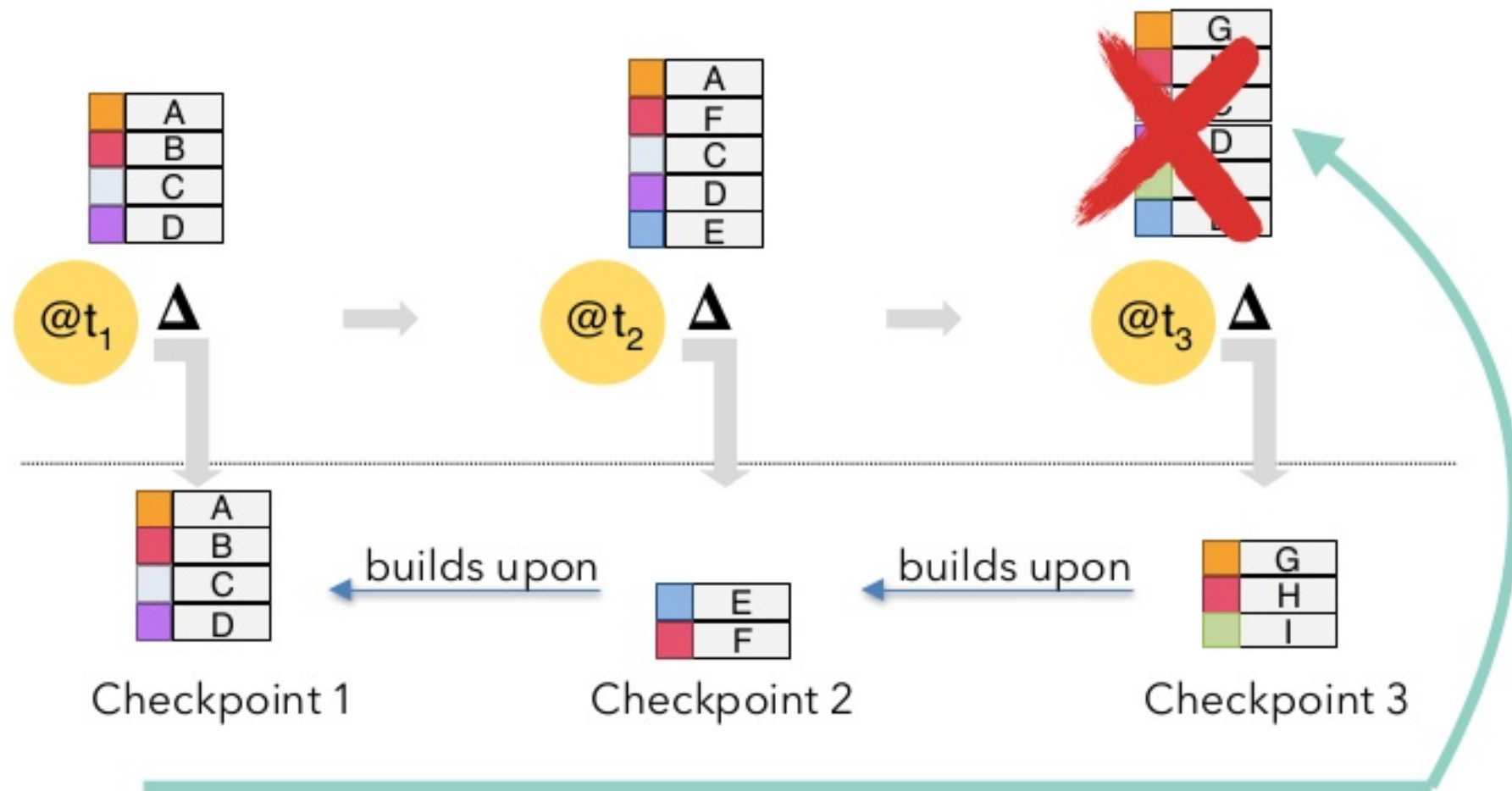- Optional: compression with Snappy.

# INCREMENTAL CHECKPOINT



Checkpoint 1

# INCREMENTAL CHECKPOINT



@$t_1$ Δ → @$t_2$ Δ

Checkpoint 1     builds upon     Checkpoint 2

# INCREMENTAL CHECKPOINT

# INCREMENTAL CHECKPOINT

# INCREMENTAL CHECKPOINT



Checkpoint 1     builds upon     Checkpoint 2     builds upon     Checkpoint 3

# INCREMENTAL CHECKPOINTS WITH ROCKSDB

Memory

Active
MemTable

Full/Switch

ReadOnly
MemTable

Flush

**Incremental checkpoints:**
**Observe created/deleted**
**SST files since last checkpoint**

Local Disk

SST
SST
SST
SST

Merge

Compaction

# INCREMENTAL CHECKPOINT OVERVIEW

- Expected trade-off: faster* checkpoints, slower* recovery

- Creation only copies deltas (new local SST files) to stable storage.

- Write amplification because we also upload compacted SST files so that we can prune checkpoint history.

- Sum of all increments that we read from stable storage can be larger than the full state size. Deletes are also explicit as tombstones.

- But no rebuild required because we simply re-open the RocksDB backend from the SST files.
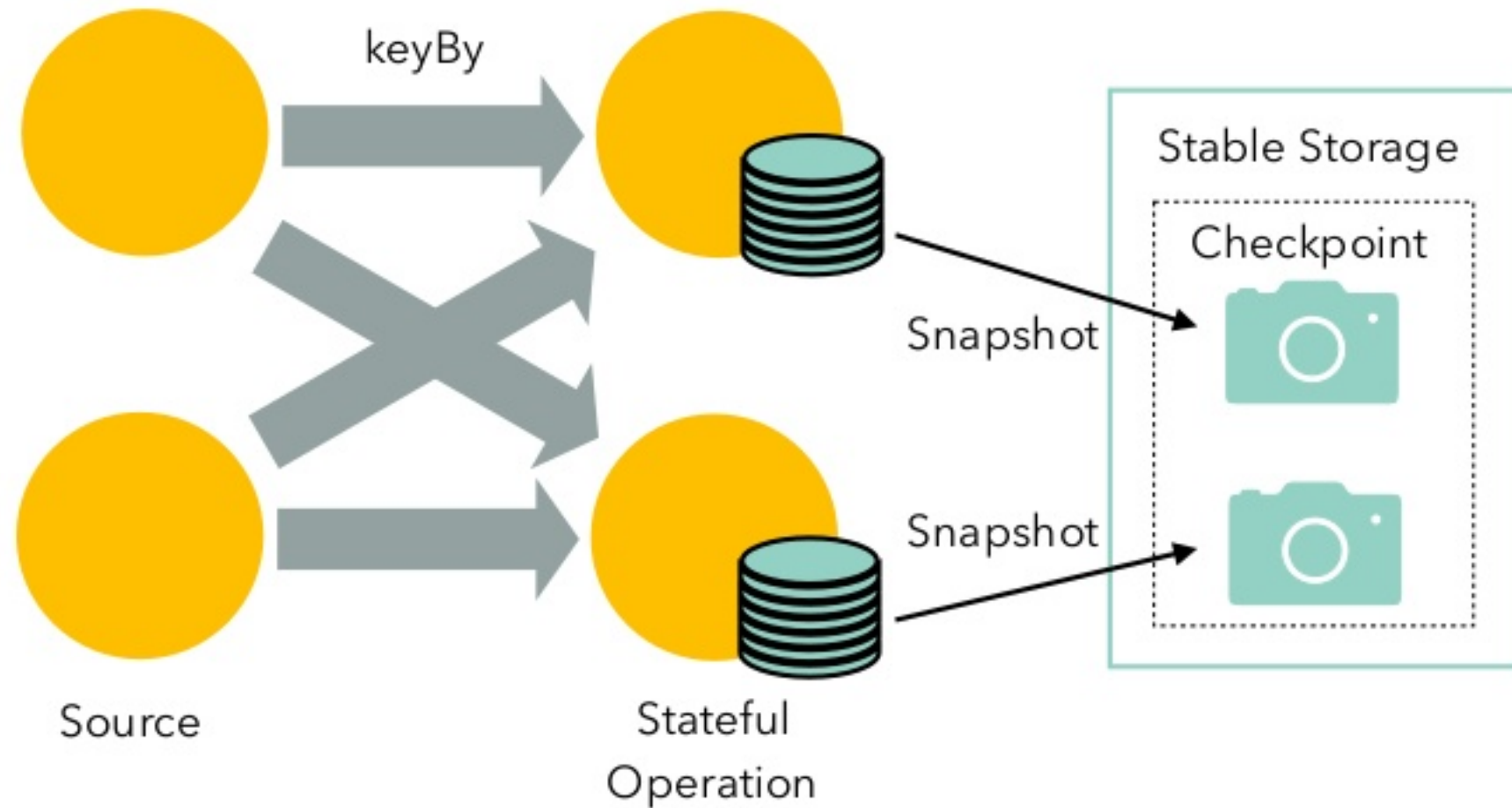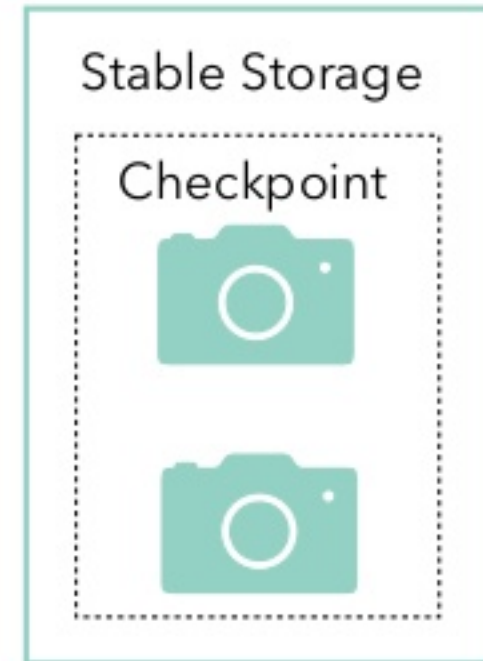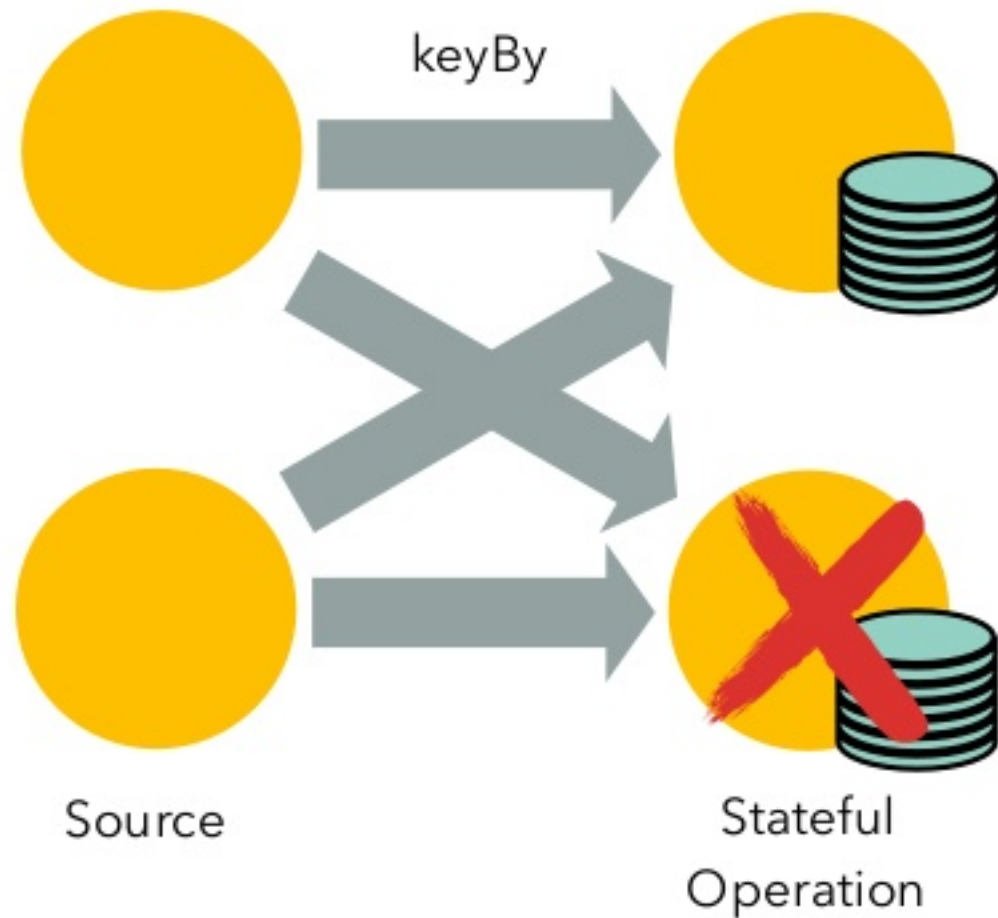
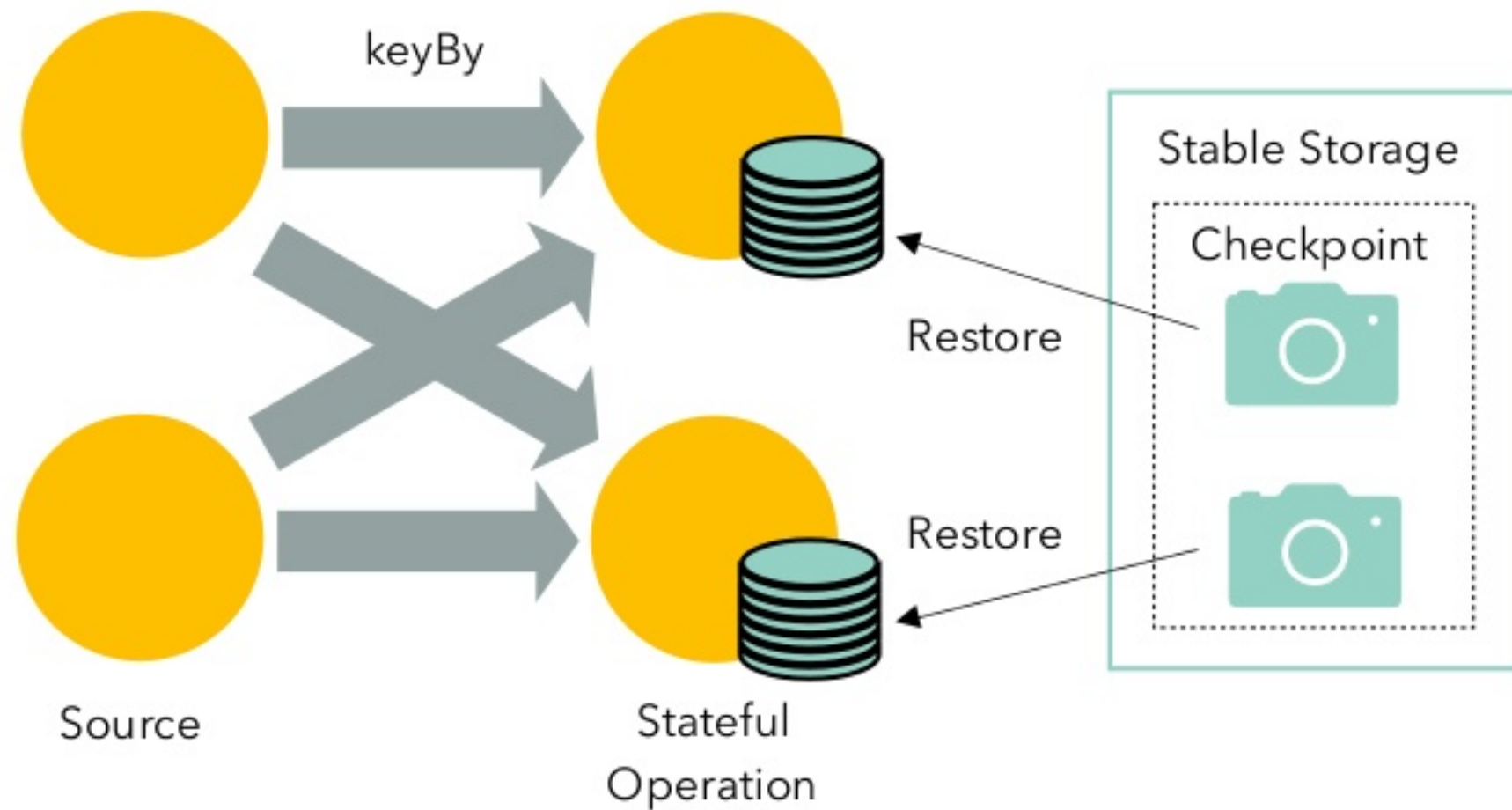- SST files are snappy-compressed by default.

# LOCAL RECOVERY

# CHECKPOINTING WITHOUT LOCAL RECOVERY



keyBy

Source

Stateful
Operation

Snapshot

Snapshot

Stable Storage

Checkpoint

# RESTORE WITHOUT LOCAL RECOVERY



keyBy

Source

Stateful
Operation

Stable Storage

Checkpoint

# RESTORE WITHOUT LOCAL RECOVERY



keyBy

Restore

Restore

Stable Storage

Checkpoint

Source

Stateful
Operation

# CHECKPOINTING WITH LOCAL RECOVERY



keyBy

Snapshot

Snapshot

Stable Storage

Checkpoint

Source

Stateful
Operation

# RESTORE WITH LOCAL RECOVERY

Scenario 1: No task manager failures, e.g. user code exception



keyBy

Restore

Restore

Source

Stateful
Operation

Stable Storage

Checkpoint

# RESTORE WITH LOCAL RECOVERY

Scenario 2: With task manager failure, e.g. disk failure

keyBy

Restore

Stable Storage

Checkpoint

Restore

Source

Stateful
Operation

# LOCAL RECOVERY TAKEAWAY POINTS

- Works with both state backends, for full and incremental checkpoints.

  - Keeps a local copy of the snapshot. Typically, this comes at the cost of mirroring the snapshot writes to remote storage also to local storage.

  - Restore with LR avoids the transfer of state from stable to local storage.

- LR works particularly well with RocksDB incremental checkpoints.

  - No new local files created, existing files might only live a bit longer.

  - Opening database from local, native table files - no ingestion / rebuild.

- Under TM failure recovery still bounded by slowest restore, but still saves a lot of resources!

# REINTERPRET STREAM AS KEYED STREAM

# REINTERPRETING A DATASTREAM AS KEYED

```
env.addSource(new InfiniteTupleSource(1000))
      .keyBy(0)
      .map((in) -> in)
      .timeWindow(Time.seconds(3));
```

Problem: Will not compile because we can no longer ensure a keyed stream!

# REINTERPRETING A DATASTREAM AS KEYED

```
env.addSource(new InfiniteTupleSource(1000))
    .keyBy(0)
    .map((in) -> in)
    .timeWindow(Time.seconds(3));
```

Problem: Will not compile because we can no longer ensure a keyed stream!

```
KeyedStream<T, K> reinterpretAsKeyedStream(
    DataStream<T> stream,
    KeySelector<T, K> keySelector)
```

Solution: Method to explicitly give (back) „keyed" property to any data stream
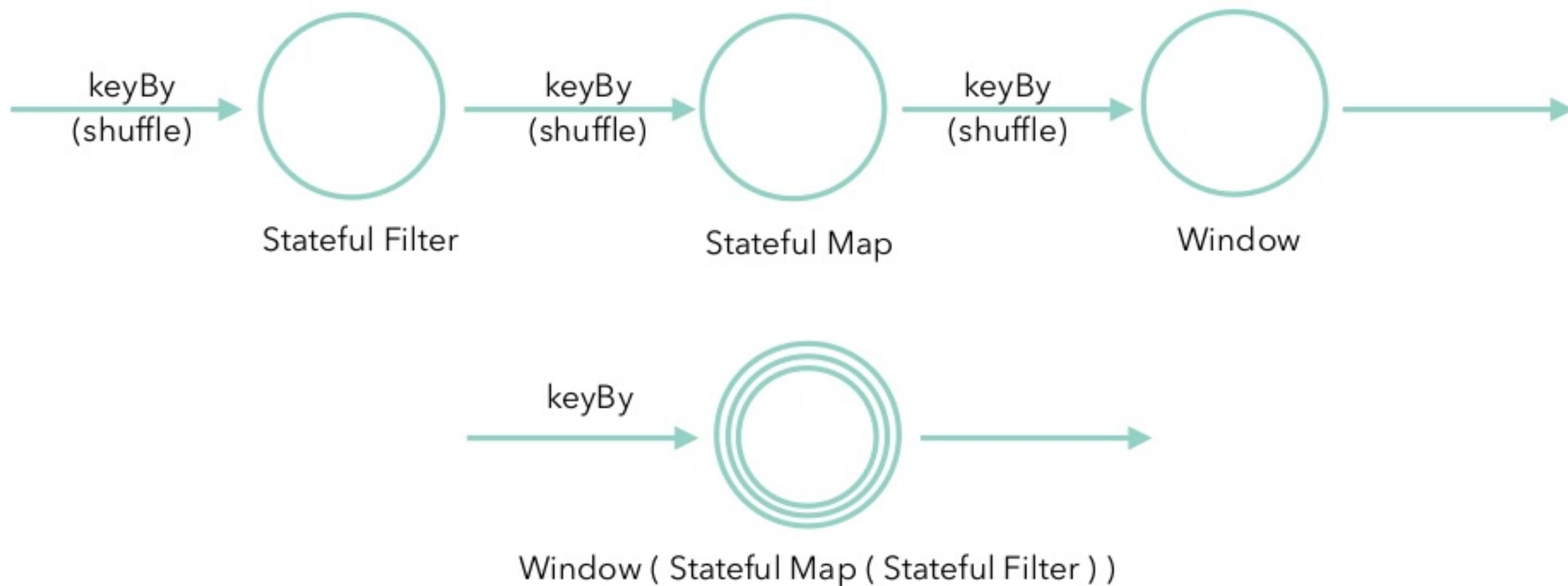
```
DataStreamUtils.reinterpretAsKeyedStream(
    env.addSource(new InfiniteTupleSource(1000))
        .keyBy(0)
        .filter((in) -> true), (in) -> in.f0)
    .timeWindow(Time.seconds(3));
```

**Warning: Only use this when you are absolutely sure that the elements in the reinterpreted stream follow exactly Flink's keyBy partitioning scheme for the given key selector!**
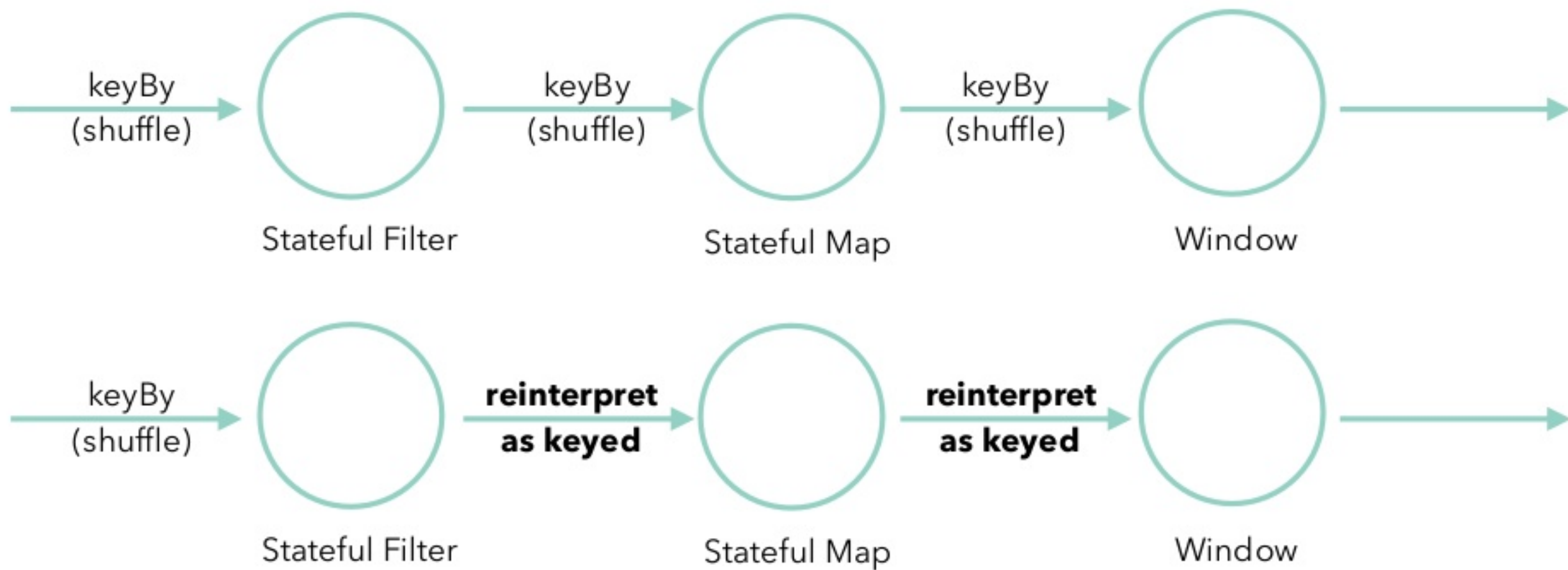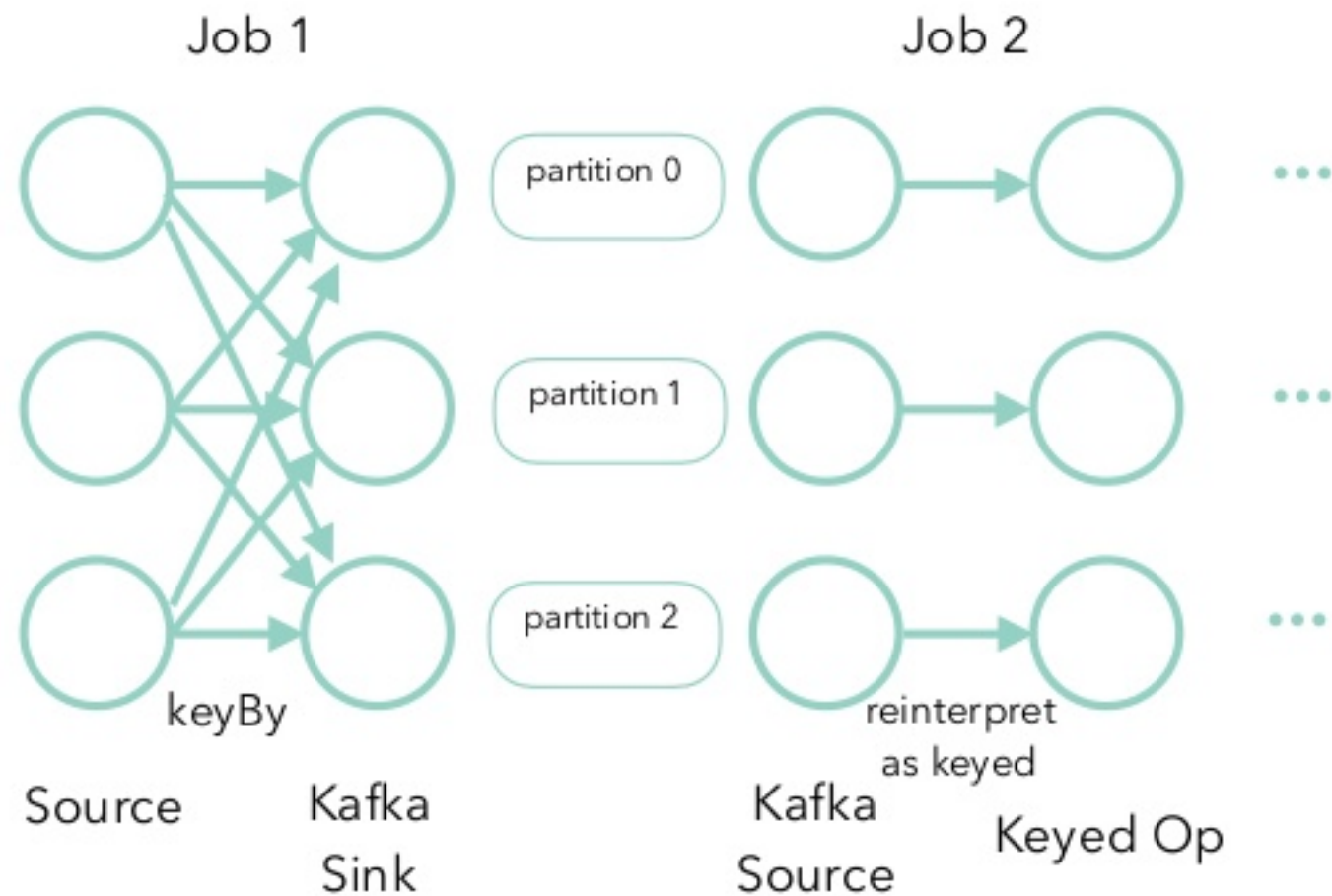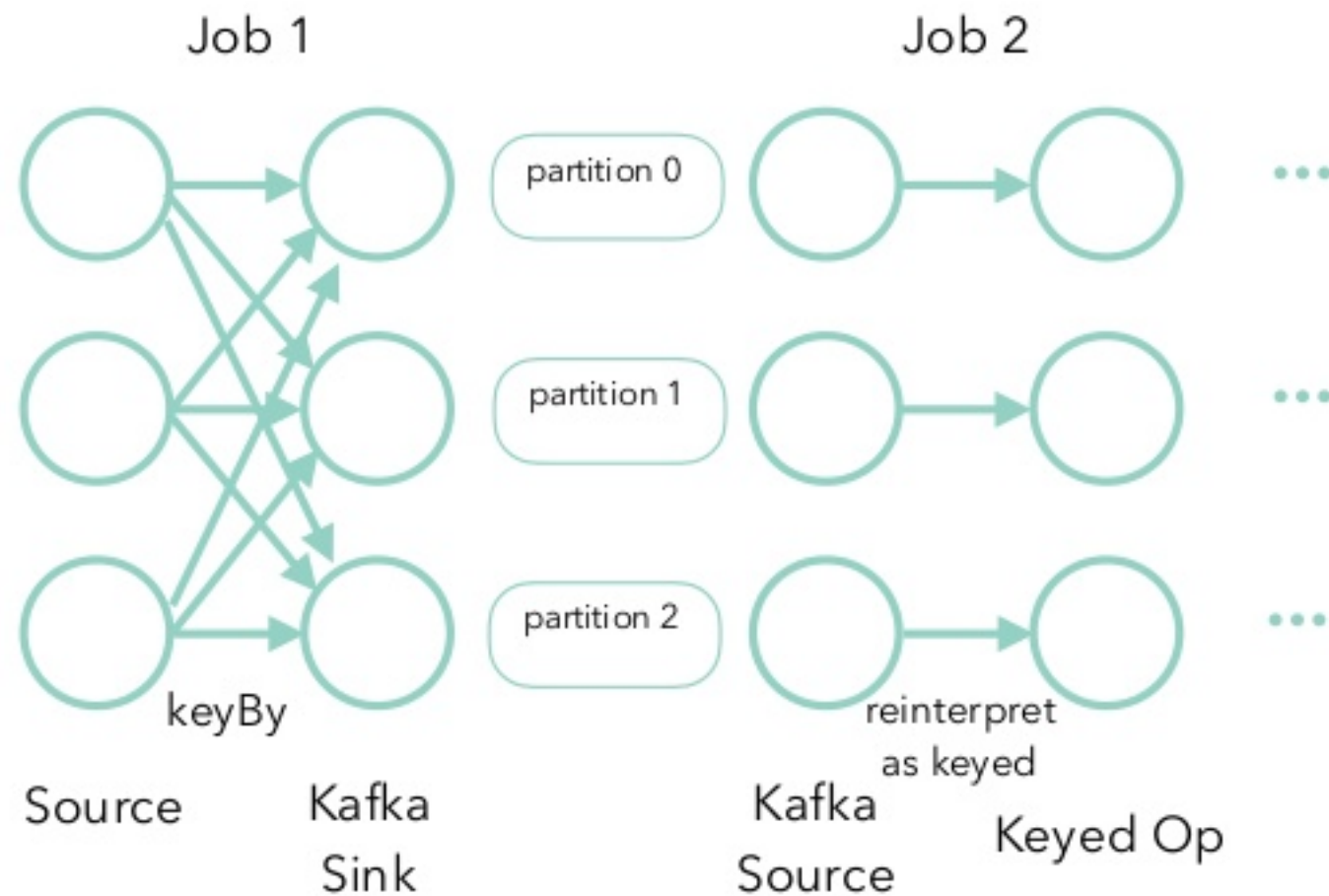
# IDEA 1 - REDUCING SHUFFLES



Window ( Stateful Map ( Stateful Filter ) )

# IDEA 1 - REDUCING SHUFFLES



© 2018 data Artisans

# IDEA 2 - PERSISTENT SHUFFLE



Job 1                                    Job 2

partition 0

partition 1

partition 2

keyBy

Source          Kafka          Kafka          Keyed Op
                Sink           Source

reinterpret
as keyed

# IDEA 2 - PERSISTENT SHUFFLE



Job 1

Job 2

partition 0

partition 1

partition 2

keyBy

reinterpret
as keyed

Source

Kafka
Sink

Kafka
Source

Keyed Op

Job 2 becomes
embarrassingly
parallel and can use
fine grained recovery!

# THANK YOU!

@StefanRRichter
@dataArtisans
@ApacheFlink

WE ARE HIRING

data-artisans.com/careers

dataArtisans