# Some large scale use cases

# @ NETFLIX

- Various use cases
  - Example: Stream ingestion, route events to Kafka, ES, Hive
  - Example: Model user interaction sessions

- Mix of stateless / moderate state / large state

- Stream Processing as a Service
  - Launching, monitoring, scaling, updating

# @ NETFLIX

**Keystone** Prod Scale – only Kafka and ES Flink routers are deployed in prod, (Hive output Flink routers are in test, unaccounted below)

- 4000+ Kafka brokers, 50+ clusters

- 100's of Data Streams (Flink Jobs)

- 3700+ Docker containers running Flink

- 1400+ nodes with 22K+ cpu cores

@ **Alibaba** Group

- Blink based on Flink
- A core system in Alibaba Search
  - Machine learning, search, recommendations
  - A/B testing of search algorithms
  - Online feature updates to boost conversion rate

- Alibaba is a major contributor to Flink
- Contributing many changes back to open source

# Blink in Alibaba Production

✓ In production for almost one year

✓ Run on thousands of nodes

- hundreds of jobs
- The biggest cluster is more than 1000 nodes
- the biggest job has 10s TB states and thousands of subtasks

✓ Supported key production services on last Nov 11[th], China Single's Day

- China Single's Day is by far the biggest shopping holiday in China, similar to Black Friday in US
- Last year it recorded $17.8 billion worth of gross merchandise volumes in one day
- Blink is used to do real time machine learning and increased conversion by around 30%
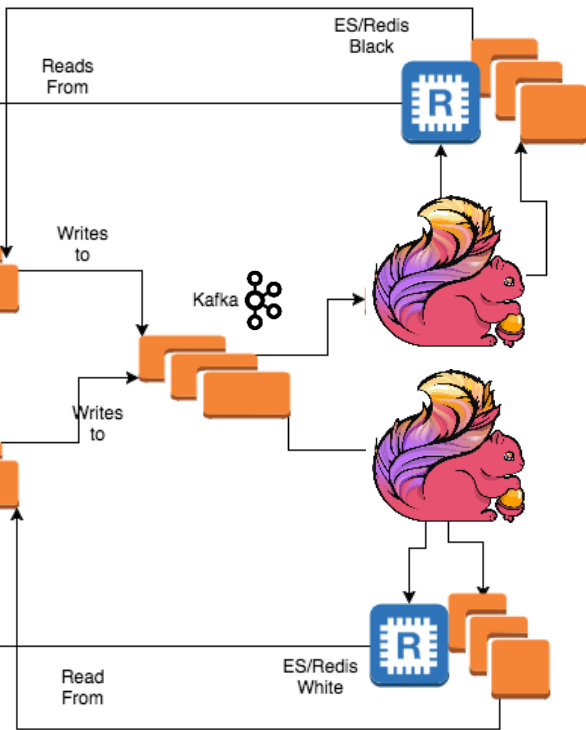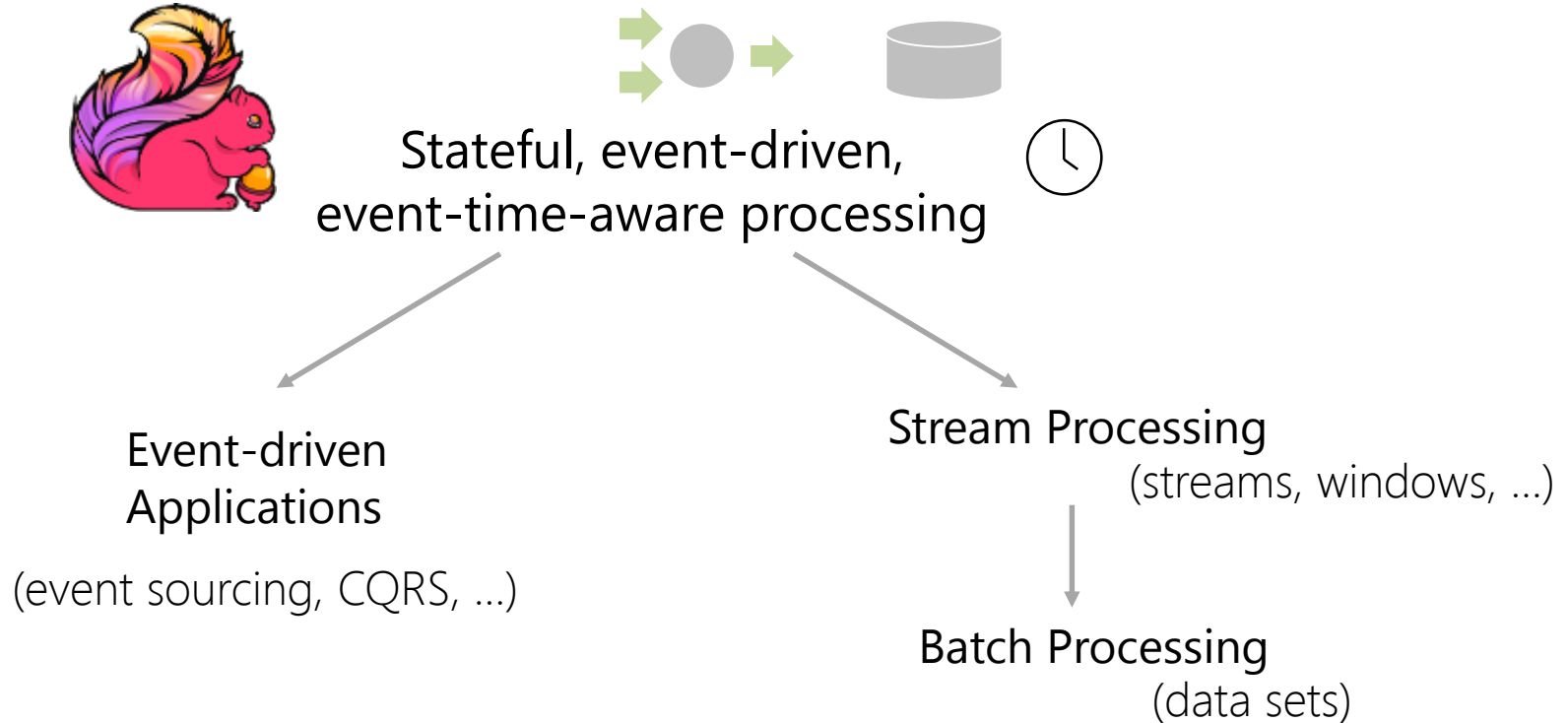
**Social network** implemented using event sourcing and CQRS (Command Query Responsibility Segregation) on Kafka/Flink/Elasticsearch/Redis
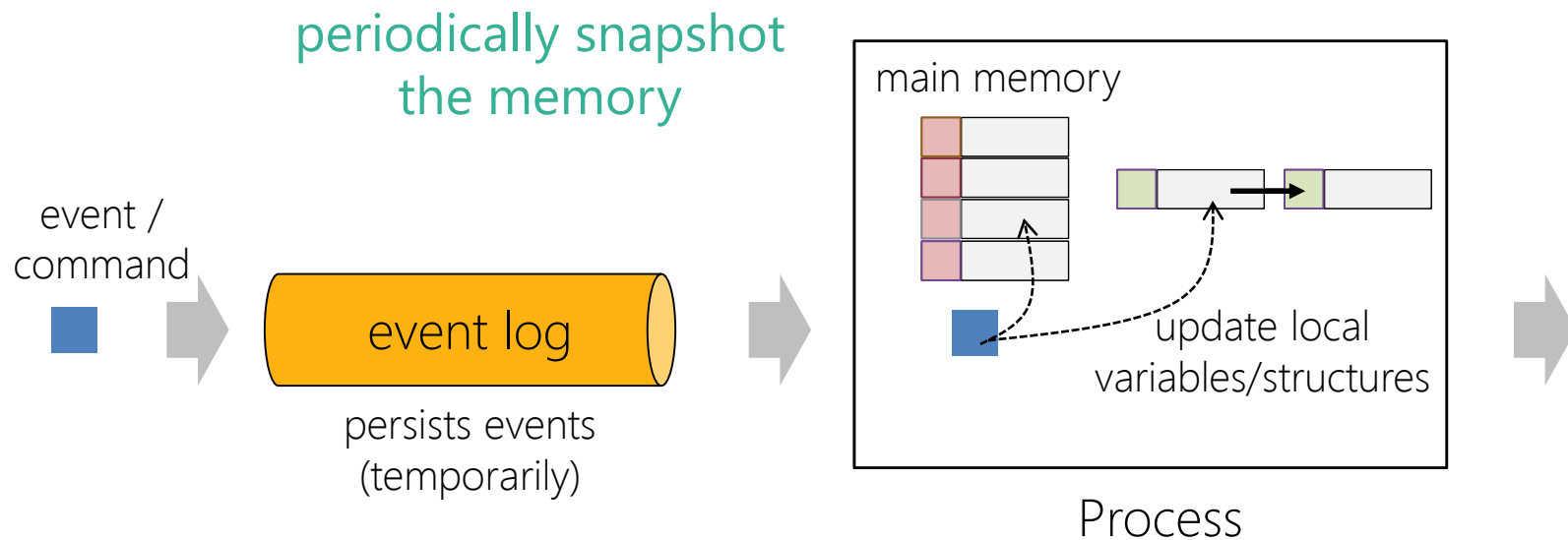
*More: https://data-artisans.com/blog/drivetribe-cqrs-apache-flink*

# How we learned to view Flink through its users

# System for Event–driven Applications



Stateful, event-driven,
event-time-aware processing

Event-driven
Applications

(event sourcing, CQRS, …)

Stream Processing

(streams, windows, …)

Batch Processing

(data sets)

# Event Sourcing + Memory Image

periodically snapshot
the memory

main memory

event /
command

event log

persists events
(temporarily)

update local
variables/structures

Process

# Event Sourcing + Memory Image

Recovery: Restore snapshot and replay events
since snapshot
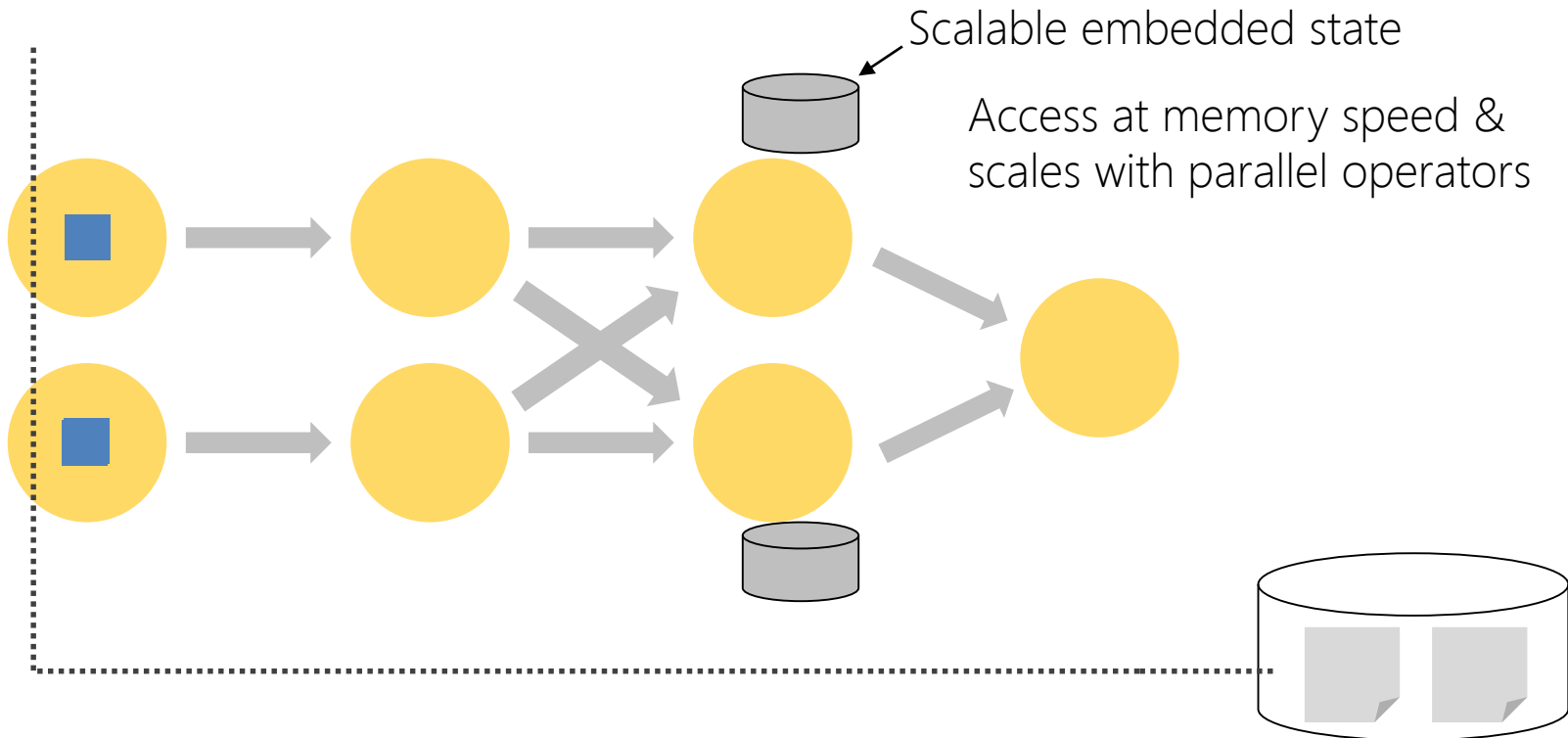


event log

persists events
(temporarily)

Process

# Distributed Memory Image

Distributed application, many memory images.
Snapshots are all consistent together.
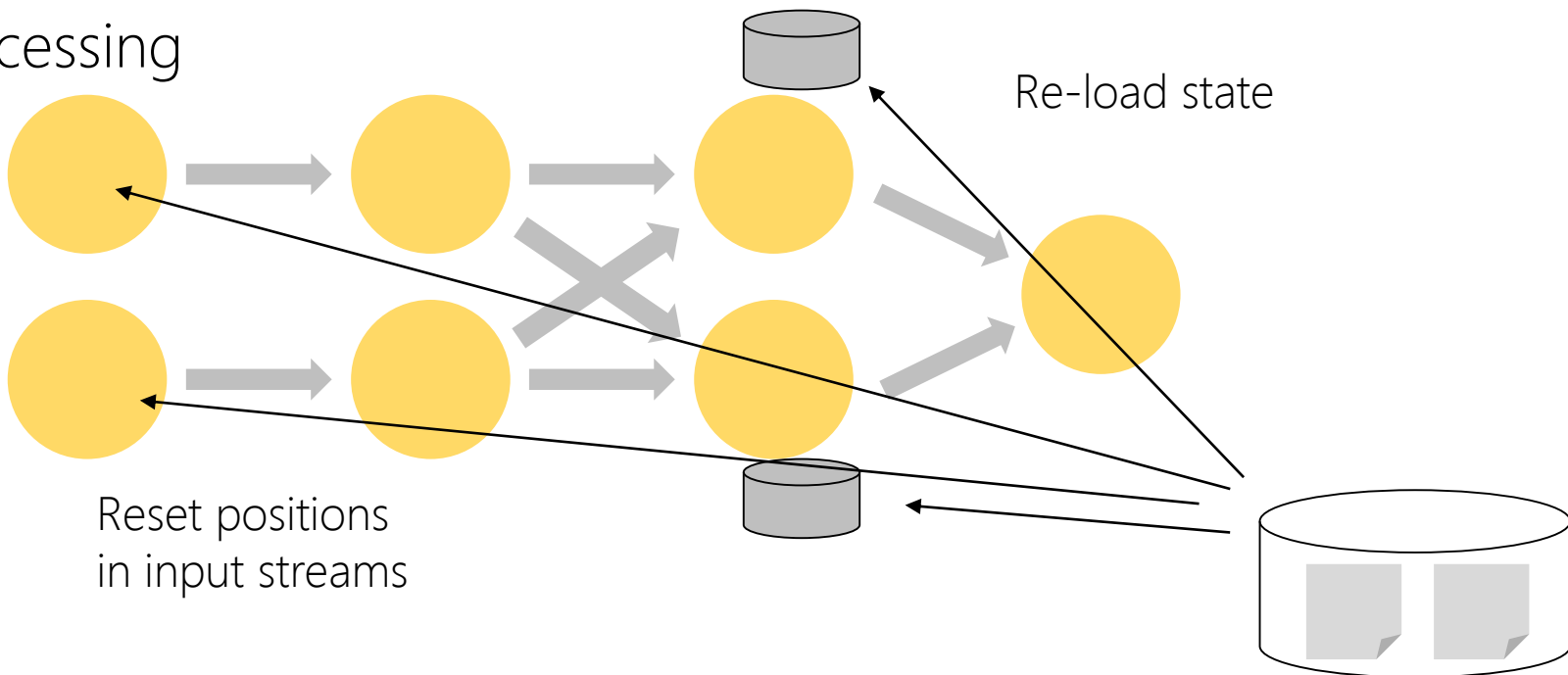
# Stateful Event & Stream Processing

Scalable embedded state

Access at memory speed & scales with parallel operators

# Stateful Event & Stream Processing

Rolling back computation

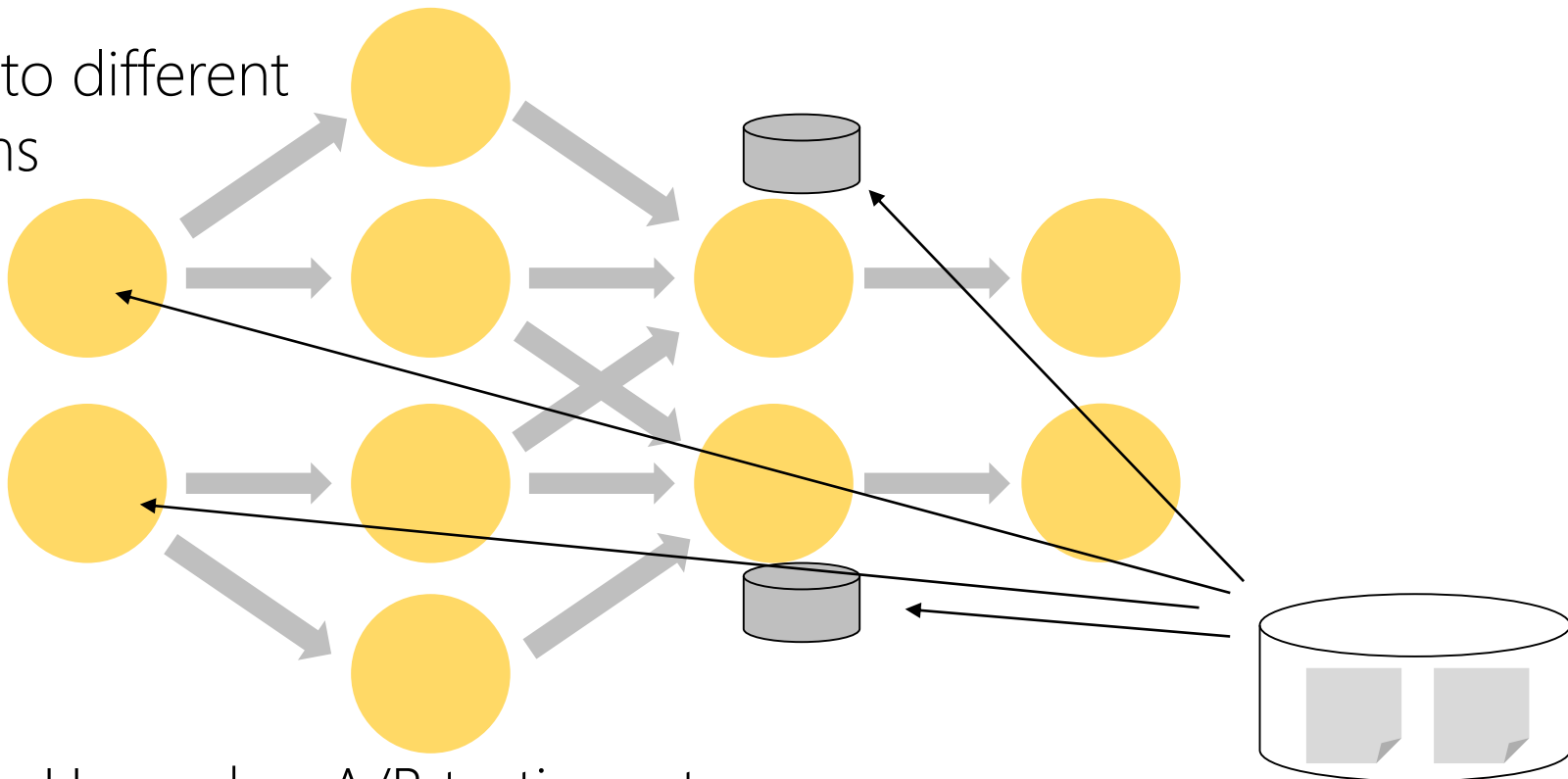Re-processing

Re-load state

Reset positions
in input streams

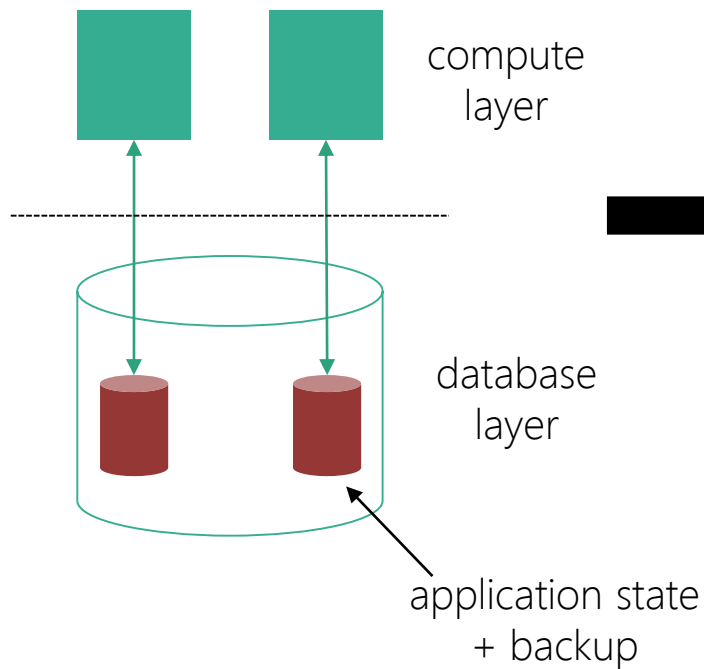# Stateful Event & Stream Processing

Restore to different programs



Bugfixes, Upgrades, A/B testing, etc

# Compute, State, and Storage

Classic tiered architecture

Streaming architecture

compute
layer

database
layer

application state
+ backup

compute
+
application state

stream storage
and
snapshot storage
(backup)

# System for Event–driven Applications

Stateful, event-driven,
event-time-aware processing

Event-driven
Applications

(event sourcing, CQRS, ...)

Stream Processing
(streams, windows, ...)

Batch Processing
(data sets)

# Apache Flink's Layered APIs

Analytics ⟶ **Stream SQL**

Stream- &
Batch Processing ⟶ **Table API** *(dynamic tables)*

⟶ **DataStream API** *(streams, windows)*

Stateful
Event-Driven
Applications ⟶ **Process Function** *(events, state, time)*

# Lessons Learned from Running Flink

The event/stream pipeline
generally just works

☺

# Interacting with the environment

- Dependency conflicts are amongst the biggest problems
  - Next versions trying to radically reduce dependencies
  - Make Hadoop an optional dependency
  - Rework shading techniques

- The deployment ecosystem is crazy complex
  - Yarn, Mesos & DC/OS, Docker & K8s, standalone, …
  - Containers and overlay networks are tricky
  - Authorization and authentication ecosystem complex it itself
  - Continuous work to improve integration

# External systems

- Dependency on any external system eventually causes downtime
  - Mainly: HDFS / S3 / NFS / … for checkpoints

- We plan to reduce dependency on those more and more in the next versions

# Type Serialization

- Type serialization is a harder problem in streaming than in batch
  - The data structure updates require more serialization
  - Types are often more complex than in batch

- State lives long and across jobs
  - Requires to "version" state and serializers
  - Requires a "schema evolution" path
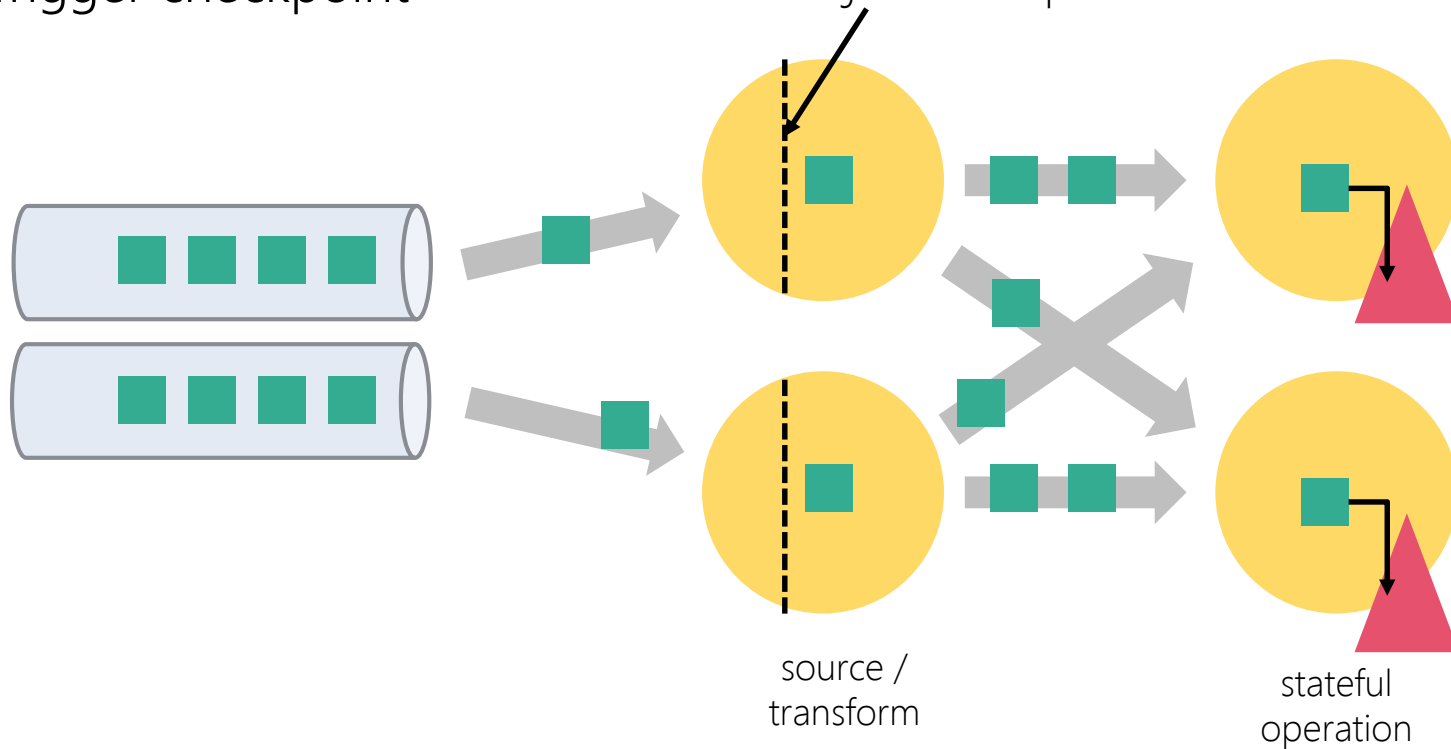  - Much enhanced support in Flink 1.3, more still to come

Robustly checkpointing…

…is the most important part of
running a large scale Flink application

# Review: Checkpoints



Trigger checkpoint
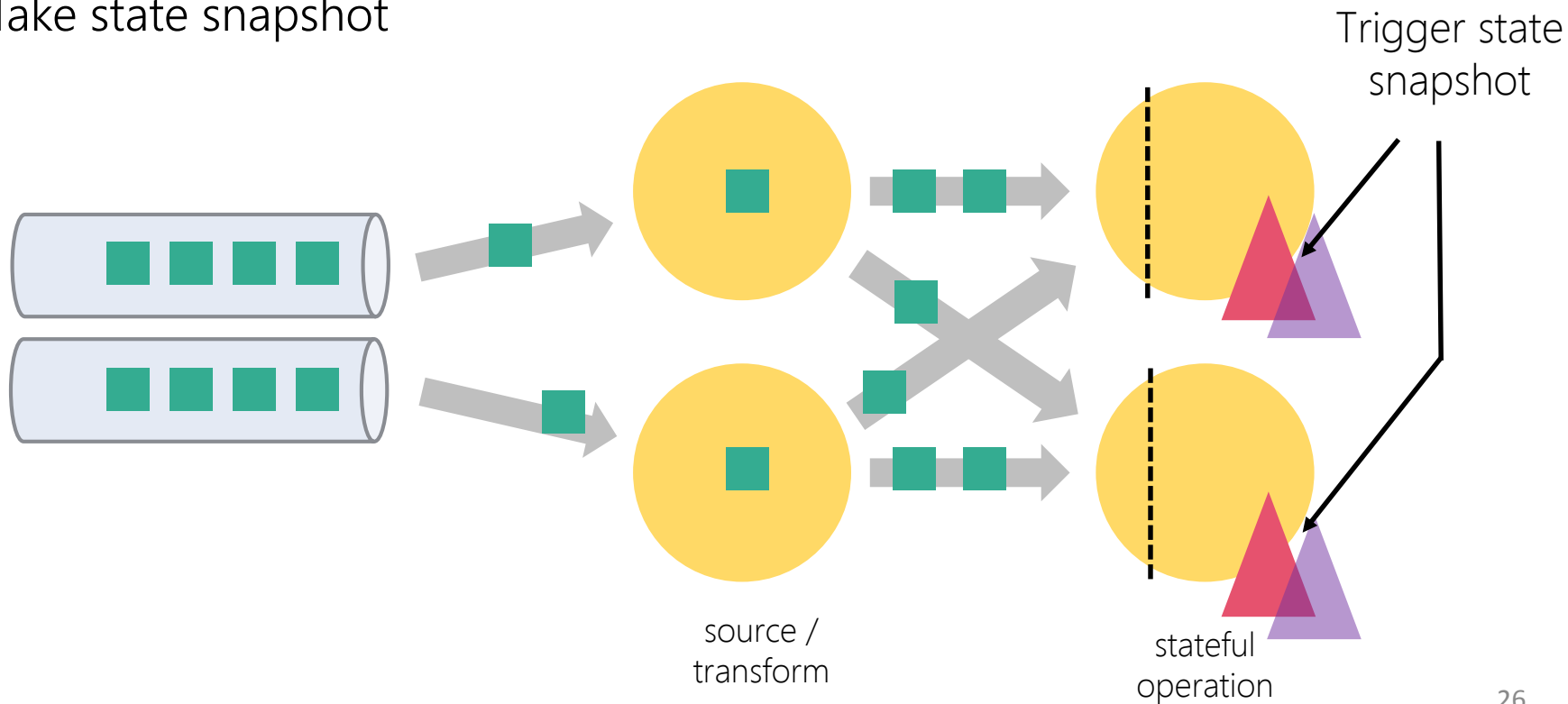
Inject checkpoint barrier

source / transform

stateful operation

# Review: Checkpoints

Take state snapshot

Trigger state snapshot



source / transform

stateful operation

# Review: Checkpoint Alignment



begin aligning

aligning

# Review: Checkpoint Alignment

emit barrier *n*

9 8 7 6 5

4 3 2 1

c b a

operator

checkpoint

input buffer

9 8 7 6 5

4

3 2 1 c

operator

i h g f e d

continue

# Understanding Checkpoints

| Subtasks | TaskManagers | Metrics | Accumulators | **Checkpoints** | Back Pressure |
| --- | --- | --- | --- | --- | --- |

| Overview | History | Summary | Configuration | Details for Checkpoint 4 |
| --- | --- | --- | --- | --- |

| ID | Status | Acknowledged | Trigger Time | Latest Acknowledgement | End to End Duration | State Size | Buffered During Alignment | Discarded |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| 4 | ✔ Completed | 8/8 (100%) | 15:42:26 | 15:42:26 | 15ms | 96.2 KB | 26.7 KB | No |

## Operators

| Name | Acknowleged | Latest Acknowledgment | End to End Duration | State Size | Buffered During Alignment | |
| --- | --- | --- | --- | --- | --- | --- |
| Source: Custom Source | 4/4 (100%) | 15:42:26 | 14ms | 48.5 KB | 0 B | Show Subtasks ⌄ |
| Flat Map -> Sink: Unnamed | 4/4 (100%) | 15:42:26 | 15ms | 47.7 KB | 26.7 KB | Show Subtasks ⌄ |

# Understanding Checkpoints

delay =
end_to_end – sync – async

How long do
snapshots take?

How well behaves
the alignment?
(lower is better)

| | End to End Duration | State Size | Checkpoint Duration (Sync) | Checkpoint Duration (Async) | Alignment Buffered | Alignment Duration |
|---|---|---|---|---|---|---|
| Source: Custom Source 4/4 (100%) 15:42:26 | | | 14ms | 48.5 KB 0 B | | Hide Subtasks ⌃ |
| Minimum | 8ms | 11.9 KB | 0ms | 0ms | 0 B | 0ms |
| Average | 10ms | 12.1 KB | 0ms | 0ms | 0 B | 0ms |
| Maximum | 14ms | 12.3 KB | 0ms | 1ms | 0 B | 0ms |

# Understanding Checkpoints

delay =
end_to_end − sync − async

How long do
snapshots take?

How well behaves
the alignment?
(lower is better)

long delay = under backpressure

too long means
→ too much state
per node
→ snapshot store cannot
keep up with load
(low bandwidth)

most important
robustness metric

under constant backpressure
means the application is
under provisioned

vastly improved with
incremental checkpoints in Flink 1.3

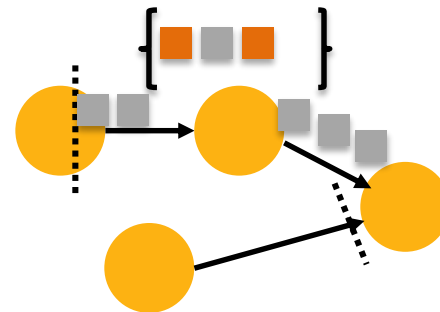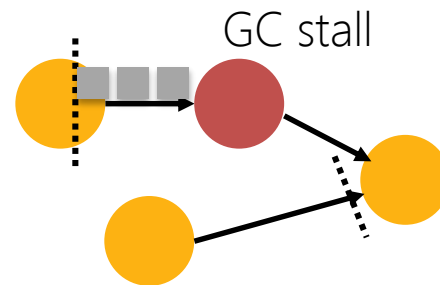| | End to End Duration | State Size | Checkpoint Duration (Sync) | Checkpoint Duration (Async) | Alignment Buffered | Alignment Duration |
|---|---|---|---|---|---|---|
| Source: Custom Source 4/4 (100%) 15:42:26 | 14ms | 48.5 KB | 0 B | | | |
| | | 1.9 KB | 0ms | 0ms | 0 B | 0ms |
| Average | 10ms | 12.1 KB | 0ms | | 0 B | 0ms |
| Maximum | 14ms | 12.3 KB | 0ms | 1ms | 0 B | 0ms |

# Heavy alignments

- A heavy alignment typically happens at some point
  → Different load on different paths

- Skewed window emission
  (lots of data on one node)

- Stall of one operator on the path

# Heavy alignments

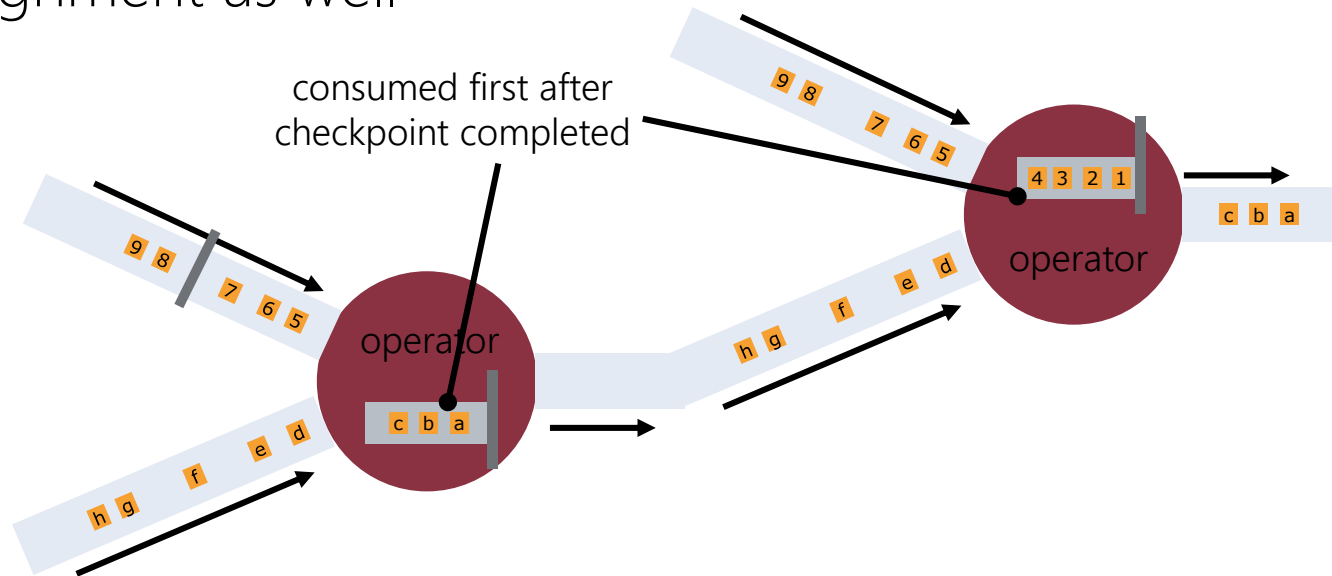- A heavy alignment typically happens at some point
  → Different load on different paths

- <span style="color:#c0504d">Skewed window emission
  (lots of data on one node)</span>

- Stall of one operator on the path

# Heavy alignments

- A heavy alignment typically happens at some point
  → Different load on different paths

- Skewed window emission
(lots of data on one node)

GC stall

- Stall of one operator on the path

# Catching up from heavy alignments

- Operators that did heavy alignment need to catch up again
- Otherwise, next checkpoint will have a heavy alignment as well



consumed first after checkpoint completed
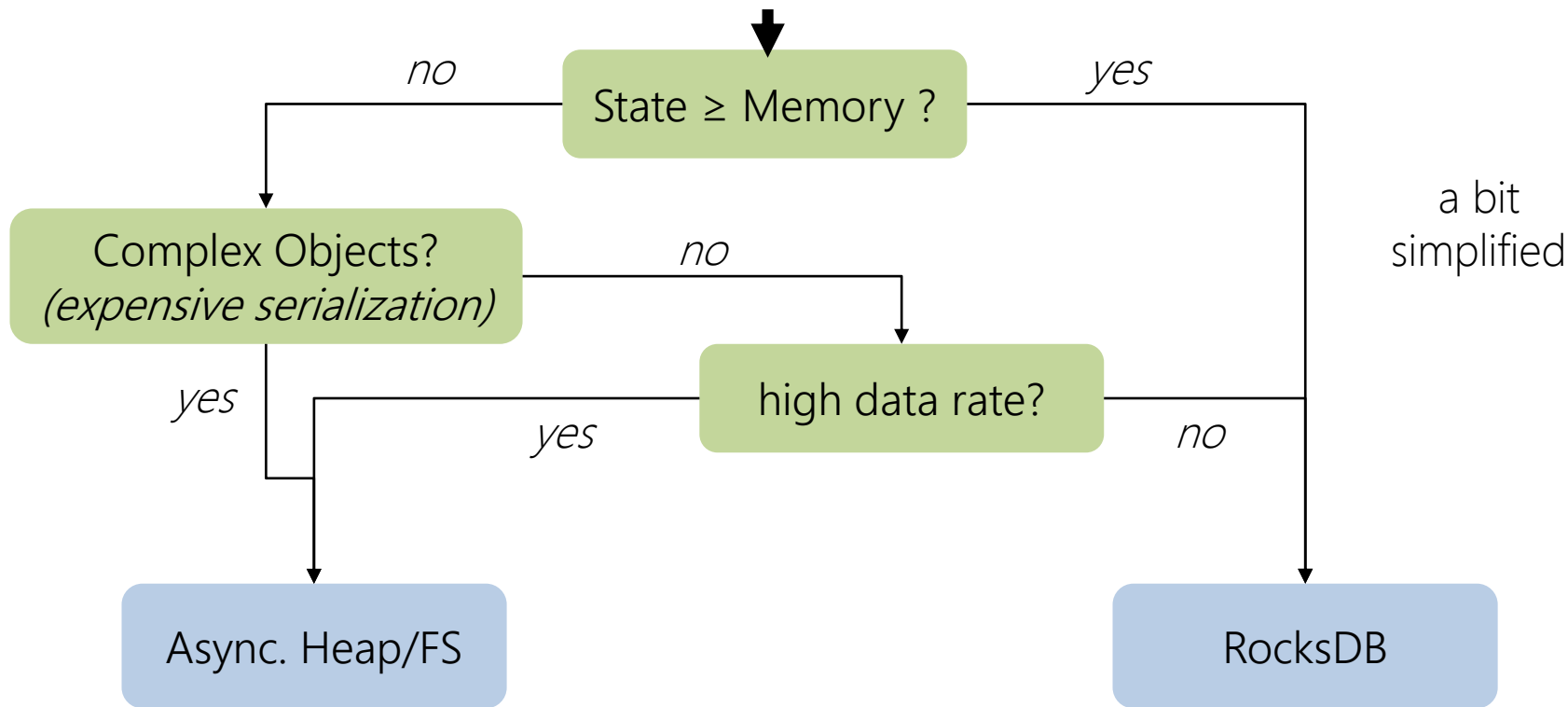
# Catching up from heavy alignments

- Giving the computation time to catch up before starting the next checkpoint
  - Set the min-time-between-checkpoints
  - Ideas to change checkpoints to policy based (spend x% of capacity on checkpoints)

- Asynchronous checkpoints mitigate most of problem
  - Very short stalls in the pipelines means shorter alignment phase
  - Catch up already happens concurrently to state materialization

# Asynchrony of different state types

| State | Flink 1.2 | Flink 1.3 | Flink 1.4 |
|---|---|---|---|
| Keyed state RocksDB | ✓ | ✓ | ✓ |
| Keyed State on heap | ✗ (✓) (hidden in 1.2.1) | ✓ | ✓ |
| Timers | ✗ | ✗ | ✓ (PR) |
| Operator State | ✗ | ✓ | ✓ |

# When to use which state backend?



State ≥ Memory ?

no

yes

Complex Objects?
*(expensive serialization)*

no

high data rate?

yes

yes

no

a bit
simplified

Async. Heap/FS

RocksDB

Berlin
11-13 Sep 2017
Flink Forward, the premier conference on Apache Flink®, is coming back to Berlin

Call for Submissions is open

dataArtisans
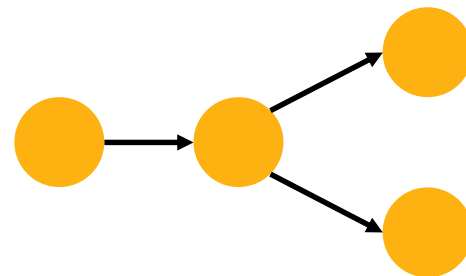
We are hiring!

data-artisans.com/careers

# Backup Slides

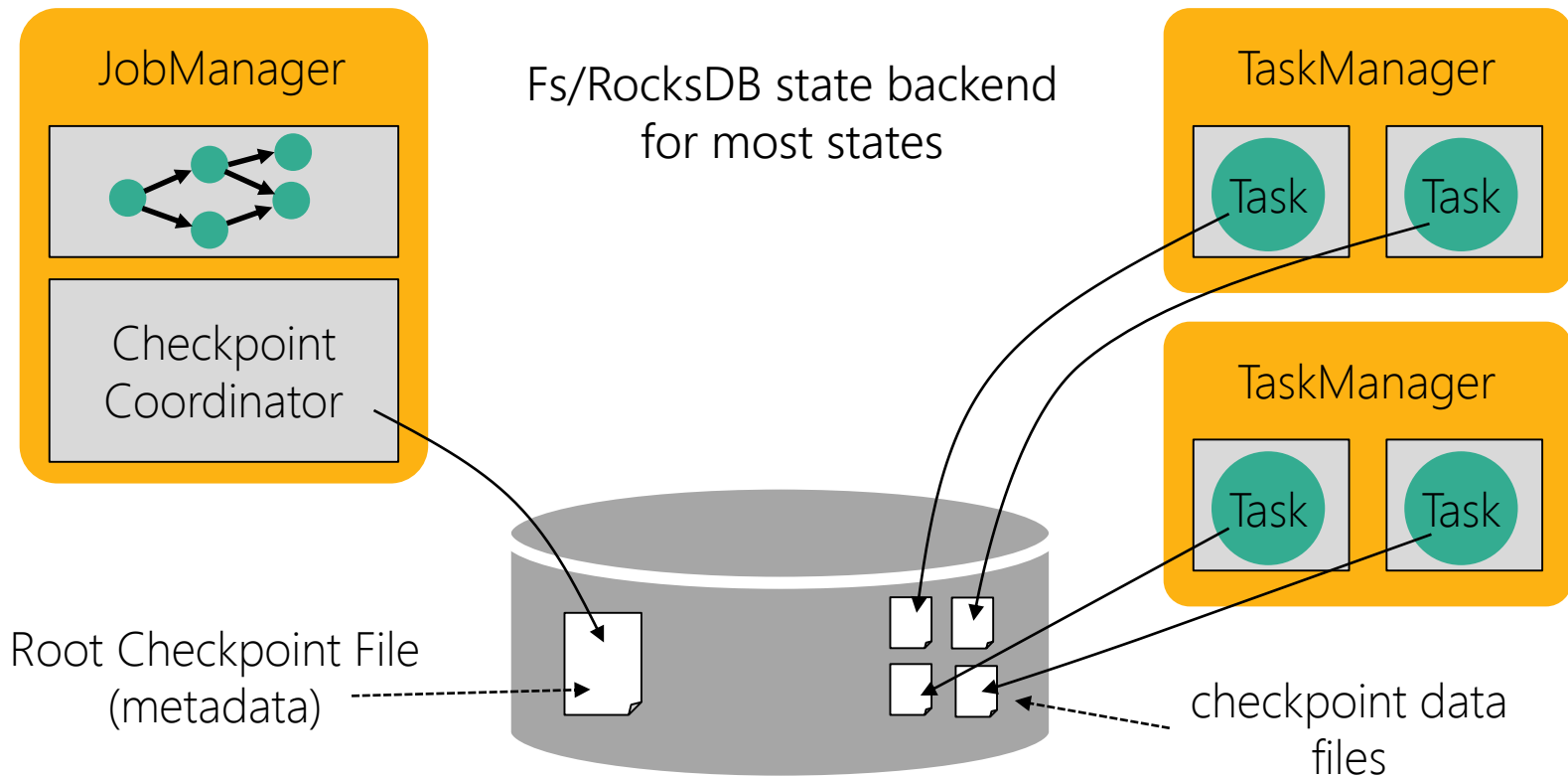# Avoiding DDOSing other systems

# Exceeding FS request capacity

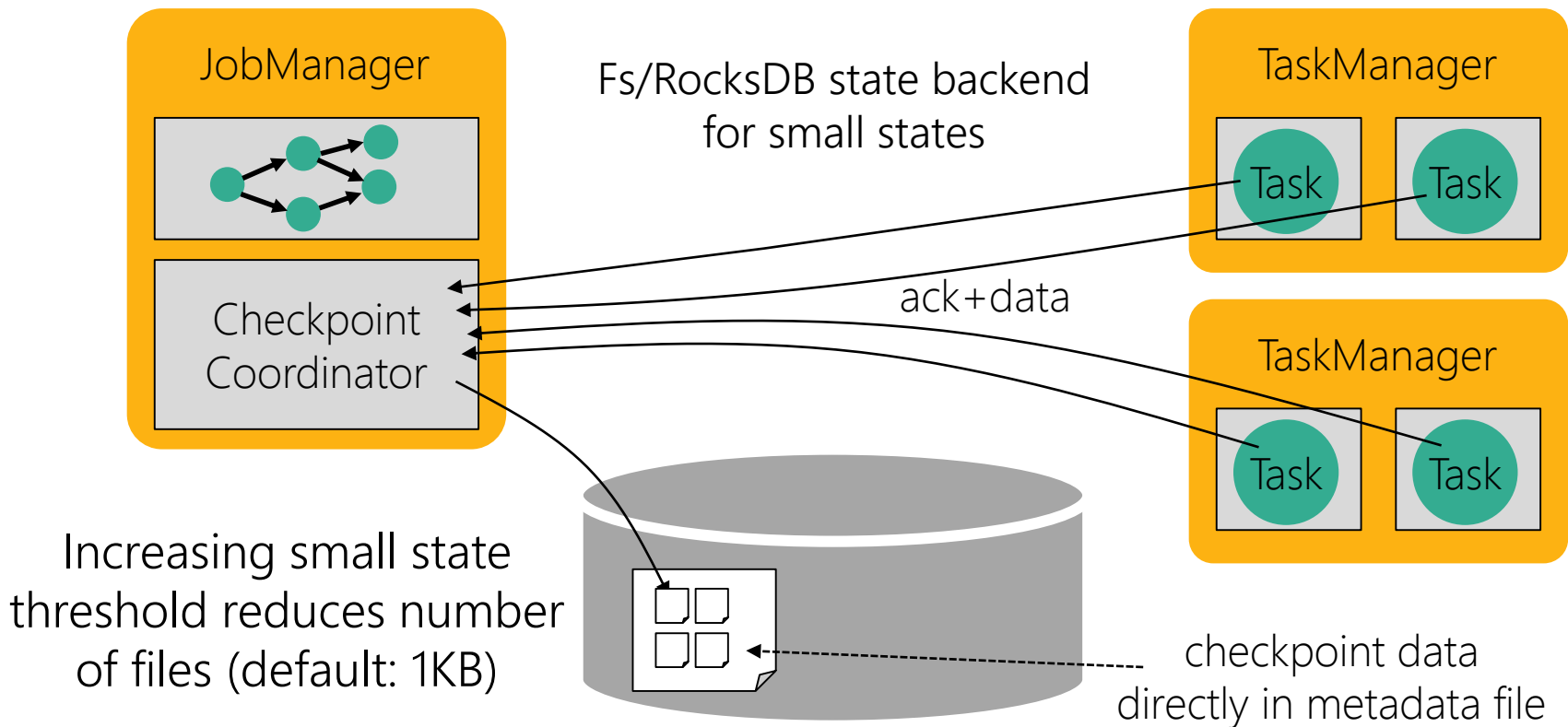- Job size: multiple 1000 operators
- Checkpoint interval: few secs

- State size: KBs per operator, 1000 of state chunks

- Via the S3 FS (from Hadoop), writes ensure "directory" exists, 2 HEAD requests
- Symptom: S3 blocked off connections after exceeding 1000s HEAD requests / sec

# Reducing FS stress for small state



JobManager

Fs/RocksDB state backend for most states

TaskManager

Task    Task

Checkpoint
Coordinator

TaskManager

Task    Task

Root Checkpoint File
(metadata)

checkpoint data
files

47

# Reducing FS stress for small state



JobManager

Fs/RocksDB state backend for small states

TaskManager

Task  Task

Checkpoint Coordinator

ack+data

TaskManager

Task  Task

Increasing small state threshold reduces number of files (default: 1KB)

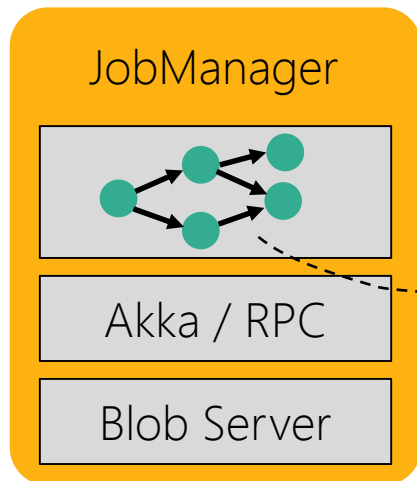checkpoint data directly in metadata file

48

# Distributed Coordination

# Deploying Tasks

Happens during initial deployment and recovery



Contains
- Job Configuration
- Task Code and Objects
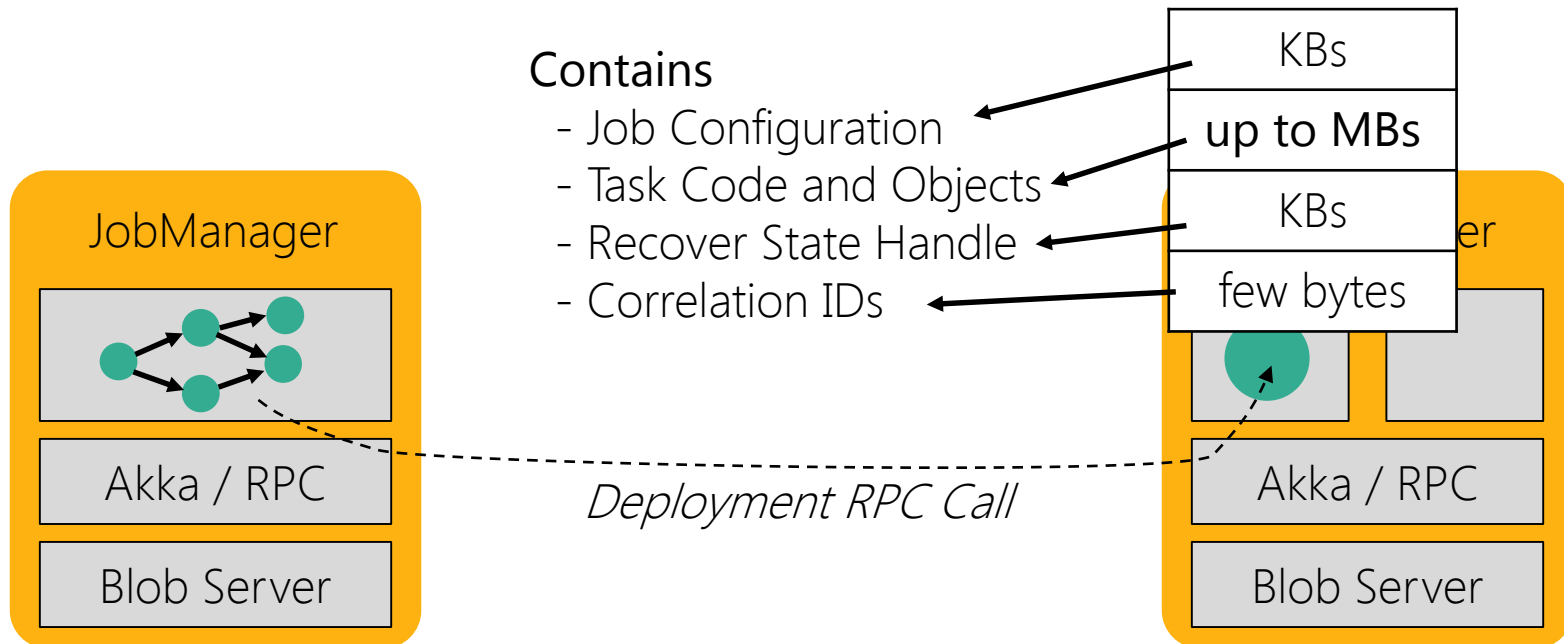- Recover State Handle
- Correlation IDs

JobManager

Akka / RPC

Blob Server

TaskManager

Akka / RPC

Blob Server

*Deployment RPC Call*

# Deploying Tasks

Happens during initial deployment and recovery

Contains
- Job Configuration
- Task Code and Objects
- Recover State Handle
- Correlation IDs

KBs
up to MBs
KBs
few bytes

JobManager

Akka / RPC

Blob Server

*Deployment RPC Call*

Akka / RPC

Blob Server

# RPC volume during deployment

(back of the napkin calculation)

| number of tasks | | x | parallelism | x | size of task objects | = | RPC volume |
|---|---|---|---|---|---|---|---|
| 10 | | x | 1000 | x | 2 MB | = | 20 GB |

~20 seconds on full 10 GBits/s net

> 1 min with avg. of 3 GBits/s net

> 3 min with avg. of 1GBs net

# Timeouts and Failure detection

~20 seconds on full 10 GBits/s net

> 1 min with avg. of 3 GBits/s net

> 3 min with avg. of 1GBs net

Default RPC timeout: **10 secs**

default settings lead to **failed deployments with RPC timeouts**

Solution:      Increase RPC timeout
Caveat:      Increasing the timeout makes failure detection slower
Future:      Reduce RPC load (next slides)
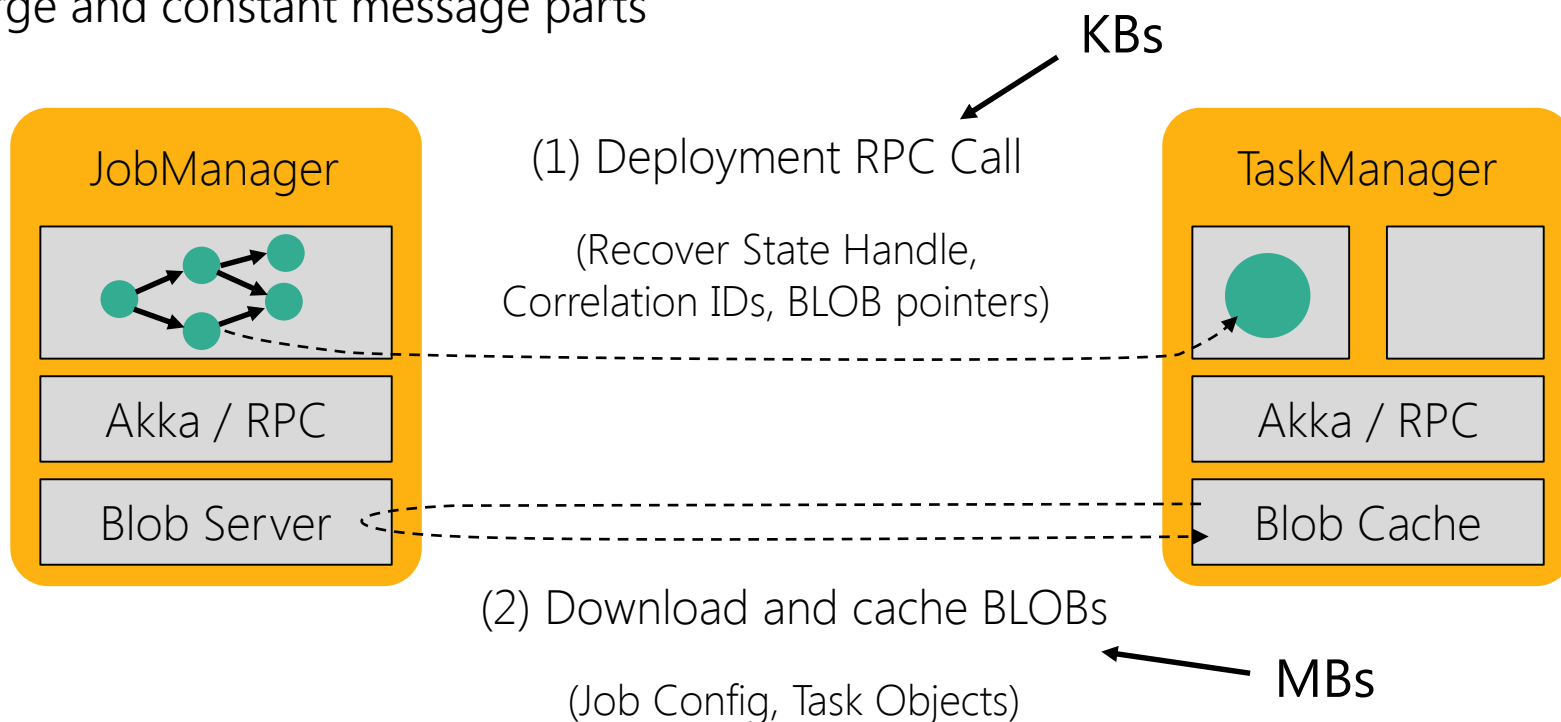
# Dissecting the RPC messages

| Message part | Size | Variance across subtasks and redeploys |
|---|---|---|
| Job Configuration | KBs | constant |
| Task Code and Objects | **up to MBs** | constant |
| Recover State Handle | KBs | variable |
| Correlation IDs | few bytes | variable |

# Upcoming: Deploying Tasks

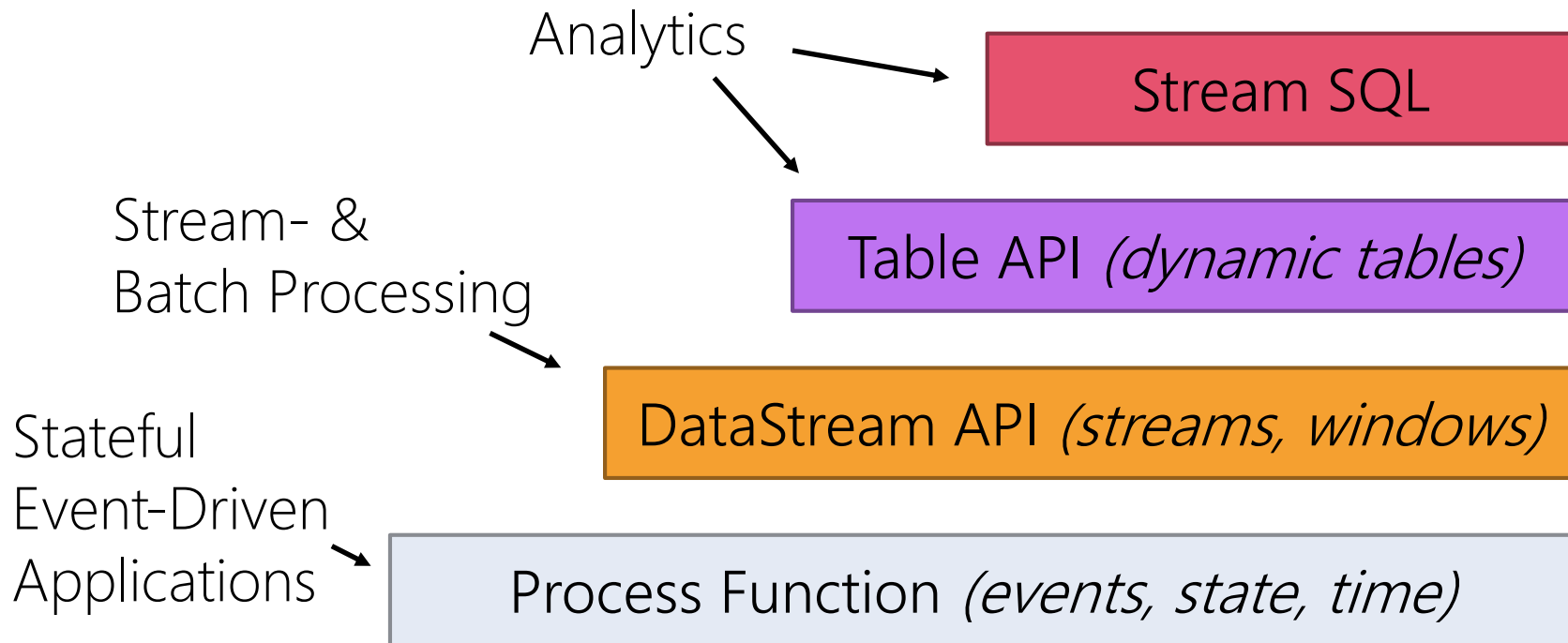Out-of-band transfer and caching of
large and constant message parts

Layers of abstraction

# Apache Flink's Layered APIs

Analytics → **Stream SQL**

Stream- & Batch Processing → **Table API** *(dynamic tables)*

**DataStream API** *(streams, windows)*

Stateful Event-Driven Applications → **Process Function** *(events, state, time)*

# Process Function

```scala
class MyFunction extends ProcessFunction[MyEvent, Result] {

    // declare state to use in the program
    lazy val state: ValueState[CountWithTimestamp] = getRuntimeContext().getState(…)

    def processElement(event: MyEvent, ctx: Context, out: Collector[Result]): Unit = {
        // work with event and state
        (event, state.value) match { … }

        out.collect(…) // emit events
        state.update(…) // modify state

        // schedule a timer callback
        ctx.timerService.registerEventTimeTimer(event.timestamp + 500)
    }

    def onTimer(timestamp: Long, ctx: OnTimerContext, out: Collector[Result]): Unit = {
        // handle callback when event-/processing- time instant is reached
    }
}
```

# Data Stream API

```scala
val lines: DataStream[String] = env.addSource(
                                new FlinkKafkaConsumer09<>(…))

val events: DataStream[Event] = lines.map((line) => parse(line))

val stats: DataStream[Statistic] = stream
    .keyBy("sensor")
    .timeWindow(Time.seconds(5))
    .sum(new MyAggregationFunction())

stats.addSink(new RollingSink(path))
```
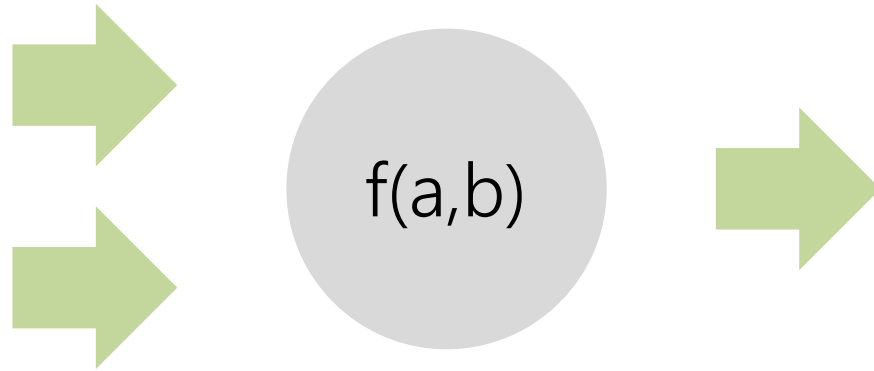
# Table API & Stream SQL

```
// Table API
val tapiResult: Table = tEnv.scan("sensors")          // scan sensors table
  .window(Tumble over 1.hour on 'rowtime as 'w)       // define 1-hour window
  .groupBy('w, 'room)                                 // group by window and room
  .select('room, 'w.end, 'temp.avg as 'avgTemp)       // compute average temperature
```

```
|SELECT room, TUMBLE_END(rowtime, INTERVAL '1' HOUR), AVG(temp) AS avgTemp
|FROM sensors
|GROUP BY TUMBLE(rowtime, INTERVAL '1' HOUR), room
```

# Events, State, Time, and Snapshots

# Events, State, Time, and Snapshots



f(a,b)

Event-driven function
executed distributedly

# Events, State, Time, and Snapshots

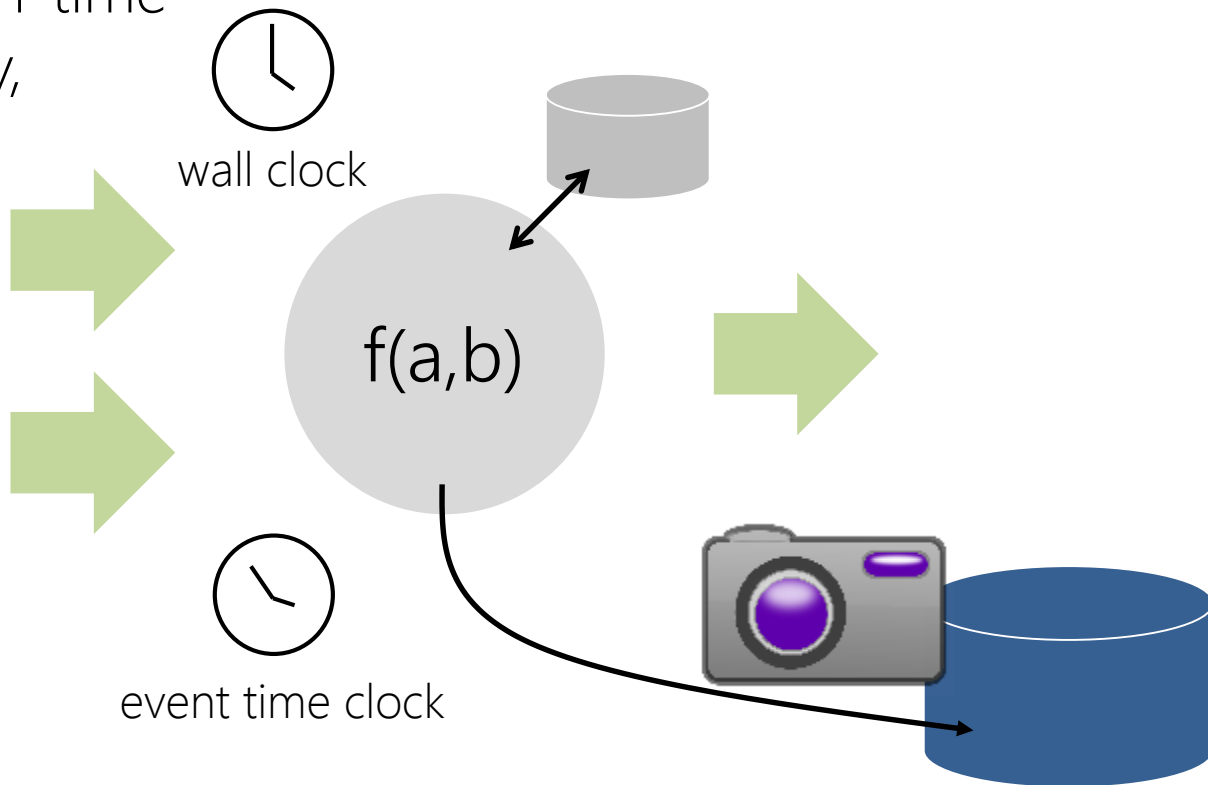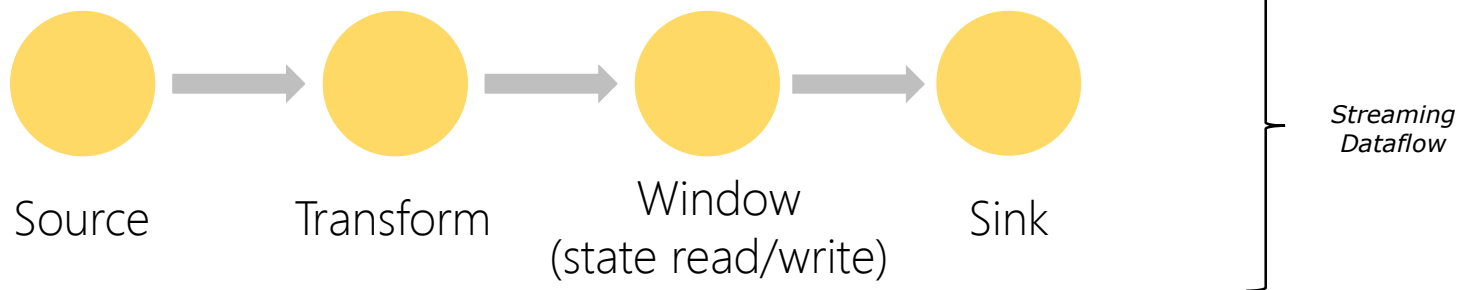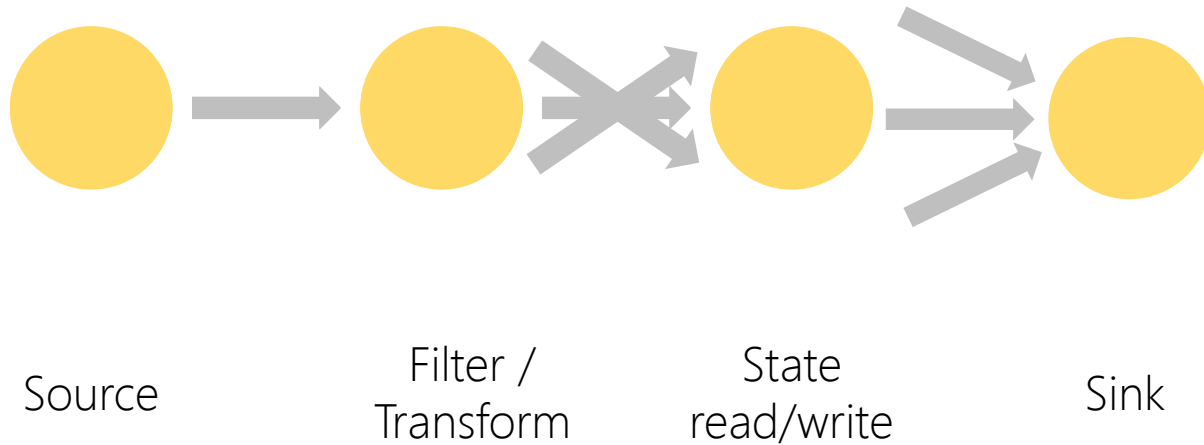Maintain fault tolerant local state similar to any normal application

Main memory +
out of core (for maps)

f(a,b)

# Events, State, Time, and Snapshots



wall clock

f(a,b)

event time clock

Access and react to notions of time and progress, handle out-of-order events

# Events, State, Time, and Snapshots

Snapshot point-in-time view for recovery, rollback, cloning, versioning, etc.

wall clock

f(a,b)

event time clock

# Stateful Event & Stream Processing

```scala
val lines: DataStream[String] = env.addSource(new FlinkKafkaConsumer09(…))
```
*Source*

```scala
val events: DataStream[Event] = lines.map((line) => parse(line))
```
*Transformation*

```scala
val stats: DataStream[Statistic] = stream
    .keyBy("sensor")
    .timeWindow(Time.seconds(5))
    .sum(new MyAggregationFunction())
```
*Transformation*

```scala
stats.addSink(new RollingSink(path))
```
*Sink*



Source    Transform    Window (state read/write)    Sink

*Streaming Dataflow*

66

# Stateful Event & Stream Processing



Source

Filter / Transform

State read/write

Sink

# Stateful Event & Stream Processing

Scalable embedded state

Access at memory speed & scales with parallel operators

# Stateful Event & Stream Processing

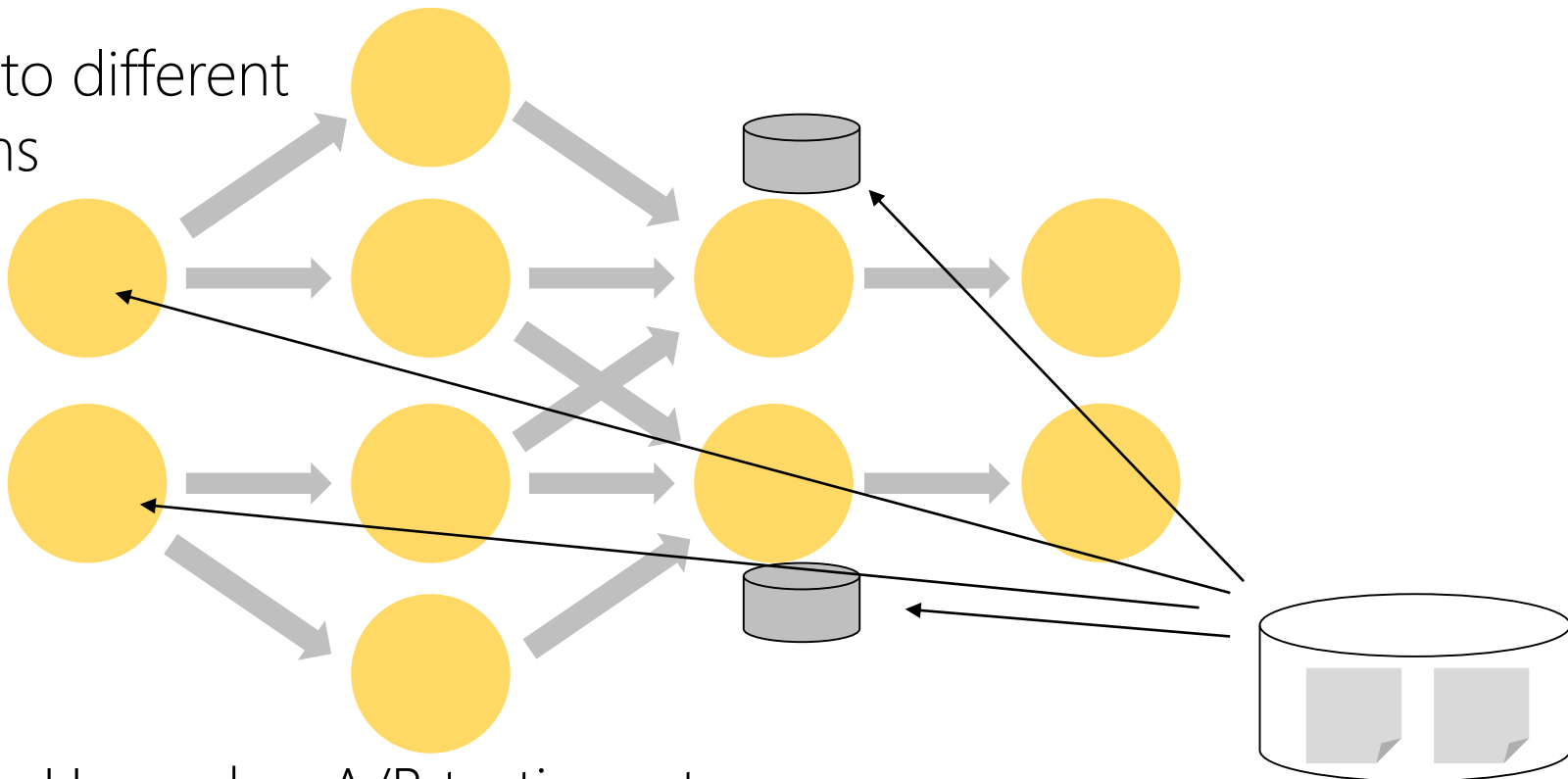Rolling back computation

Re-processing

Re-load state

Reset positions
in input streams

# Stateful Event & Stream Processing

Restore to different programs



Bugfixes, Upgrades, A/B testing, etc
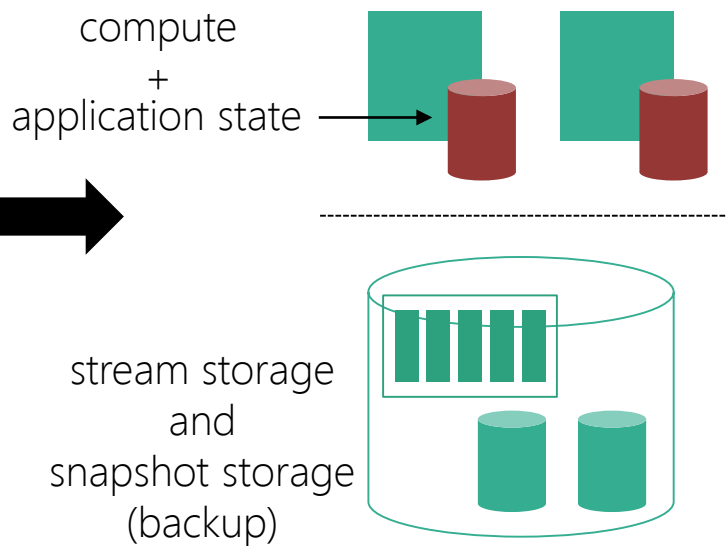
# "Classical" versus Streaming Architecture

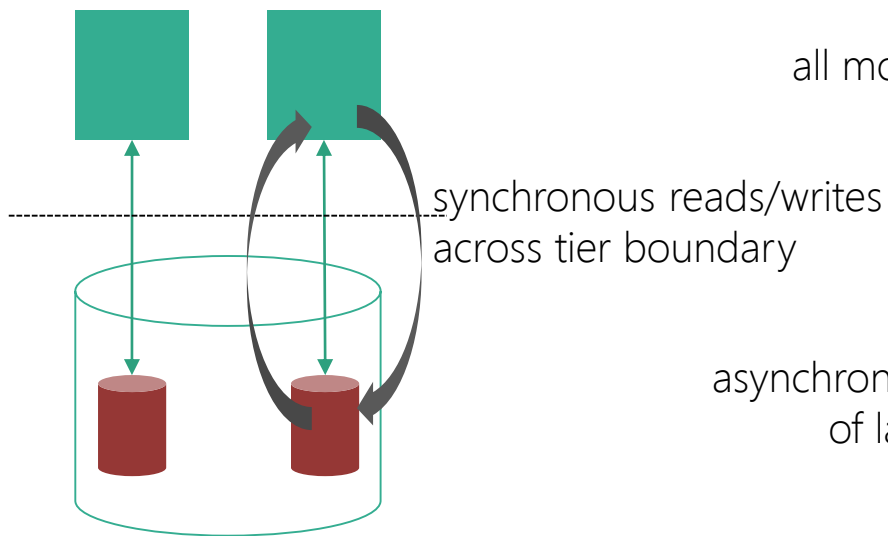# Compute, State, and Storage



Classic tiered architecture

compute layer

database layer

application state + backup

Streaming architecture

compute + application state

stream storage and snapshot storage (backup)

# Performance

Classic tiered architecture

Streaming architecture

all modifications
are local

synchronous reads/writes
across tier boundary

asynchronous writes
of large blobs

# Consistency

Classic tiered architecture

Streaming architecture

exactly once
per state

snapshot consistency
across states

=1

=1

distributed transactions

at scale typically
at-most / at-least once

# Scaling a Service

Classic tiered architecture

Streaming architecture

provision compute

provision compute
and state together

separately provision additional
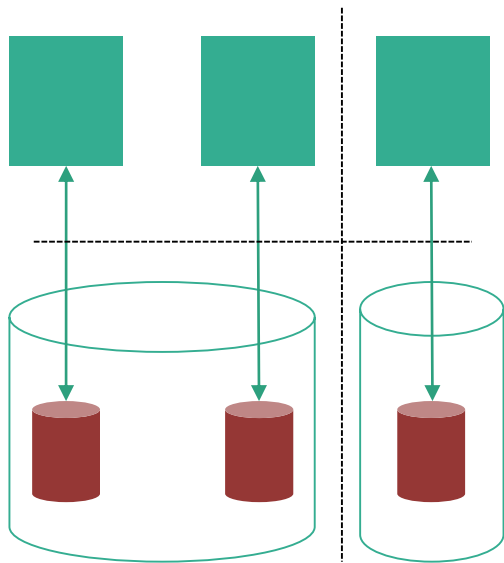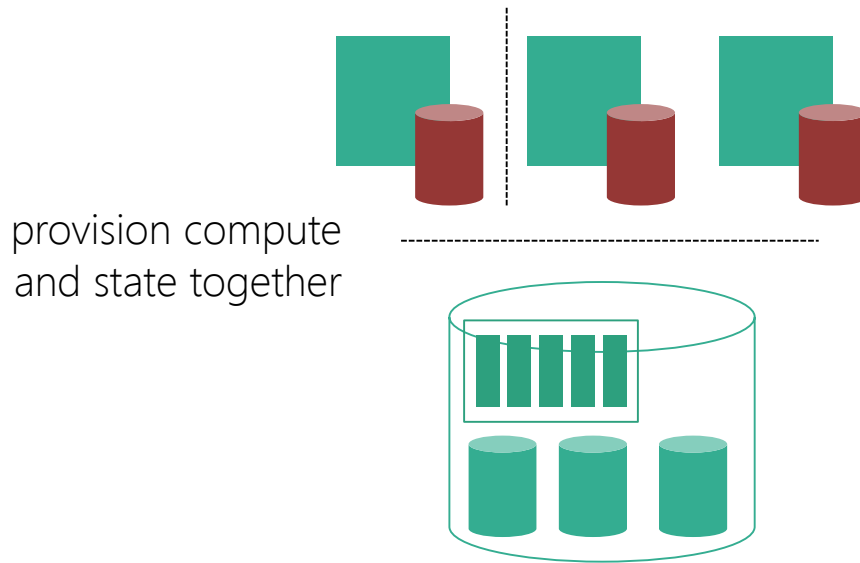database capacity

# Rolling out a new Service



Classic tiered architecture

provision a new database
(or add capacity to an existing one)

Streaming architecture
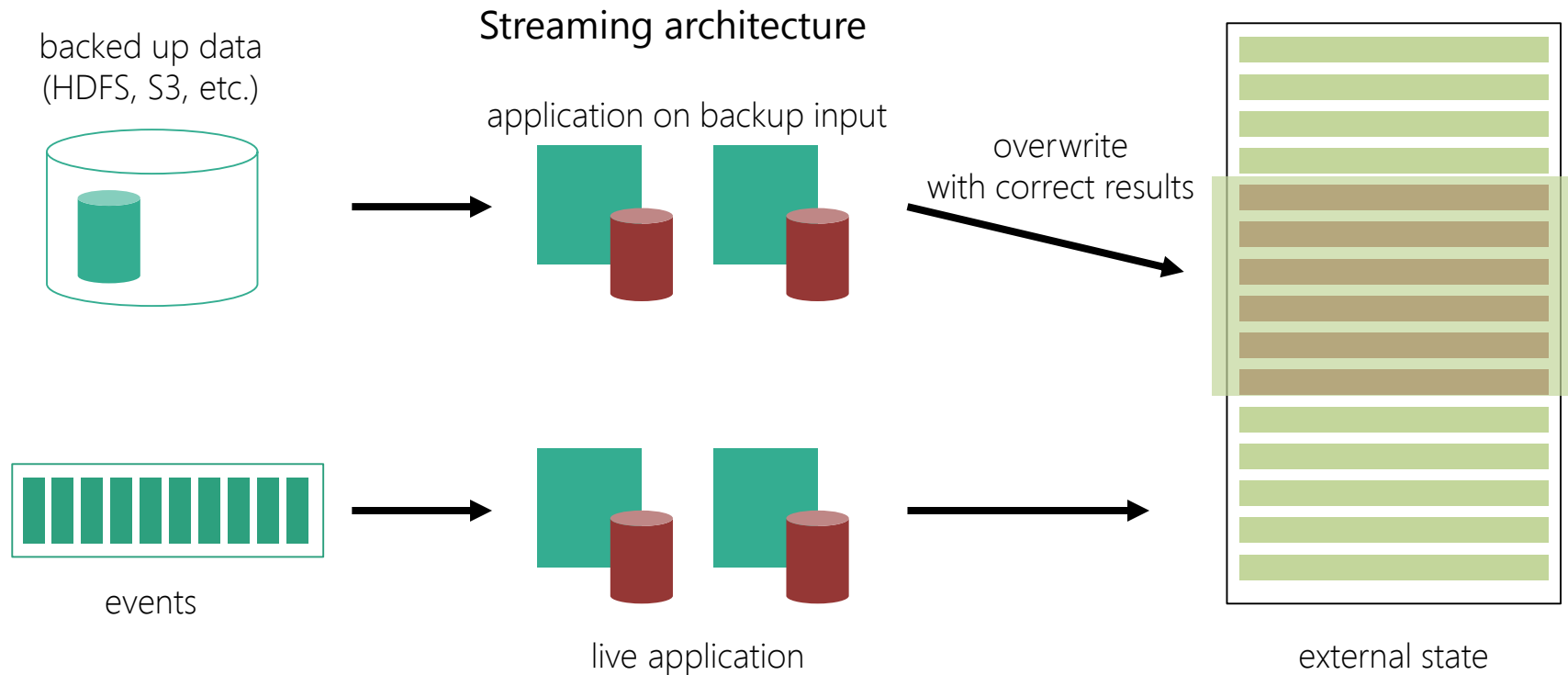
provision compute
and state together

simply occupies some
additional backup space

76

# Repair External State



Streaming architecture

backed up data
(HDFS, S3, etc.)

wrong results

events

live application

external state

# Repair External State



Streaming architecture

backed up data
(HDFS, S3, etc.)

application on backup input

overwrite
with correct results

events

live application

external state

# Repair External State



Streaming architecture

backed up date
(HDFS, S3, etc.)

application on backup input

overwrite
with correct results

Each application doubles as
a batch job!

events

live application

external state