# *Pravega*: Rethinking storage for streams

Stephan Ewen, *Data Artisans*

Flavio Junqueira, *Pravega*

*Flink Forward – Berlin 2017*

# Outline

- Intro to Pravega

- Flink + Pravega
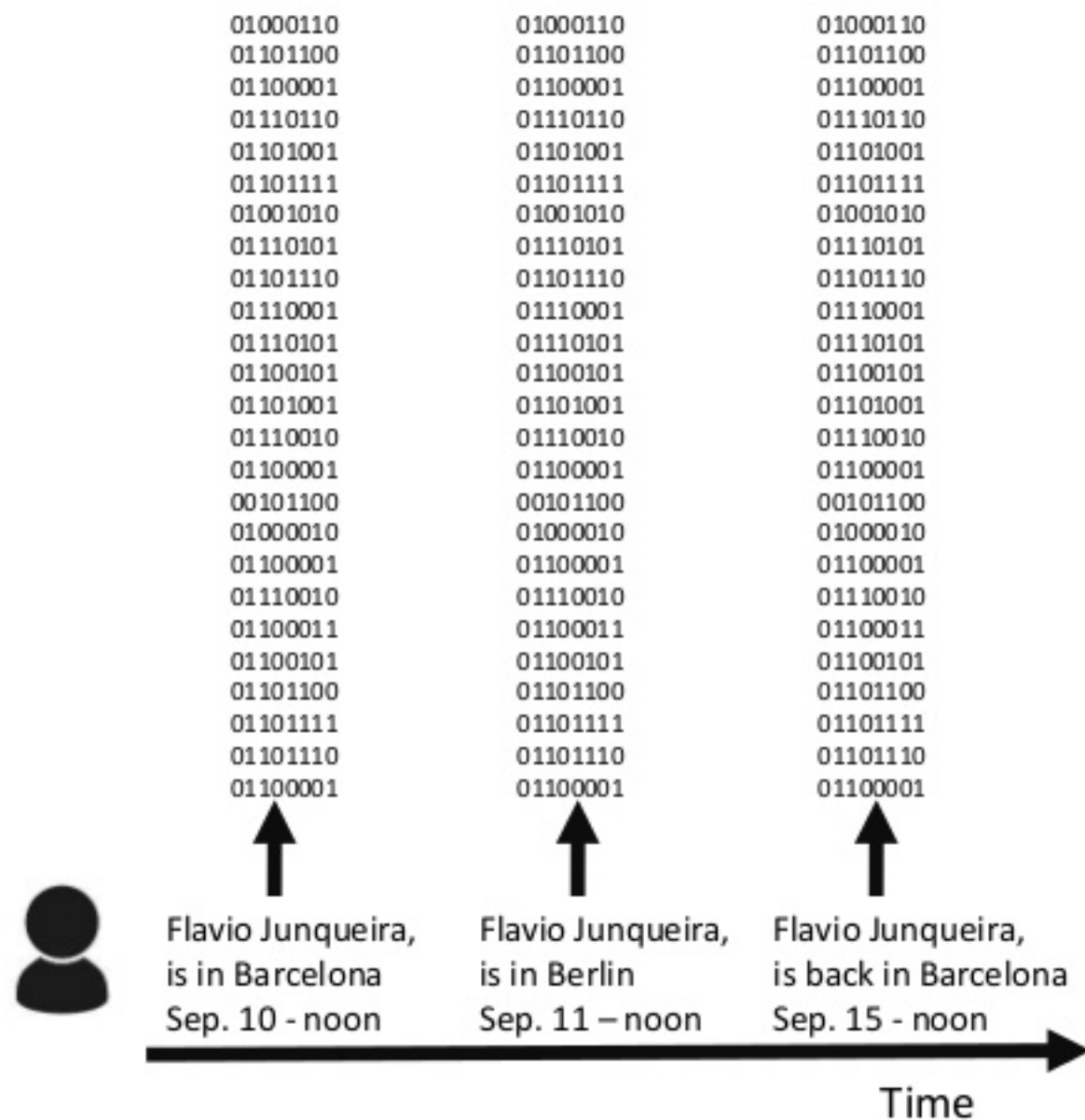
# Pravega

http://pravega.io
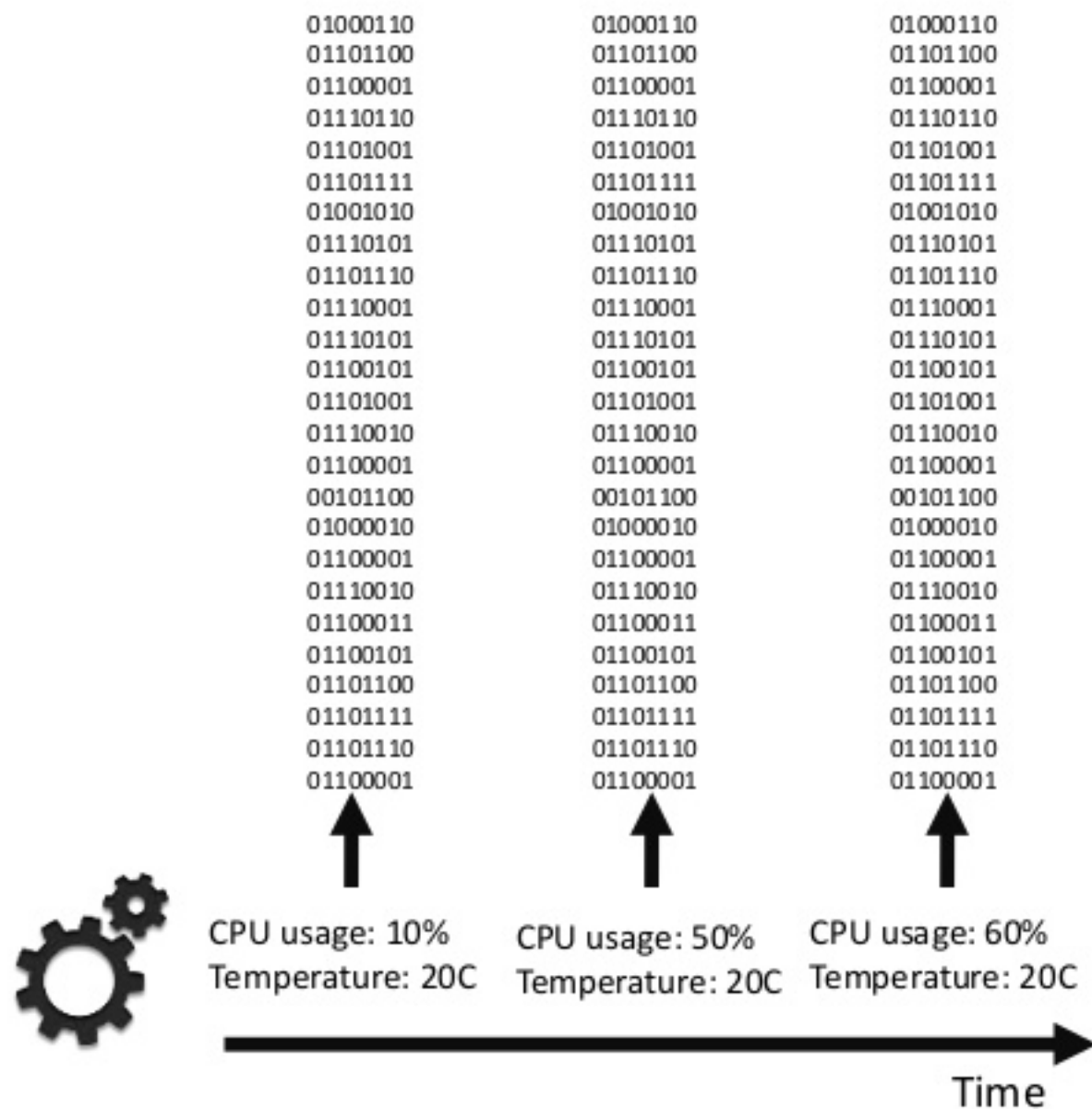
# Streams

Flavio Junqueira,
Lives in Barcelona ———▶

01000110
01101100
01100001
01110110
01101001
01101111
01001010
01110101
01101110
01110001
01110101
01100101
01101001
01110010
01100001
00101100
01000010
01100001
01110010
01100011
01100101
01101100
01101111
01101110
01100001

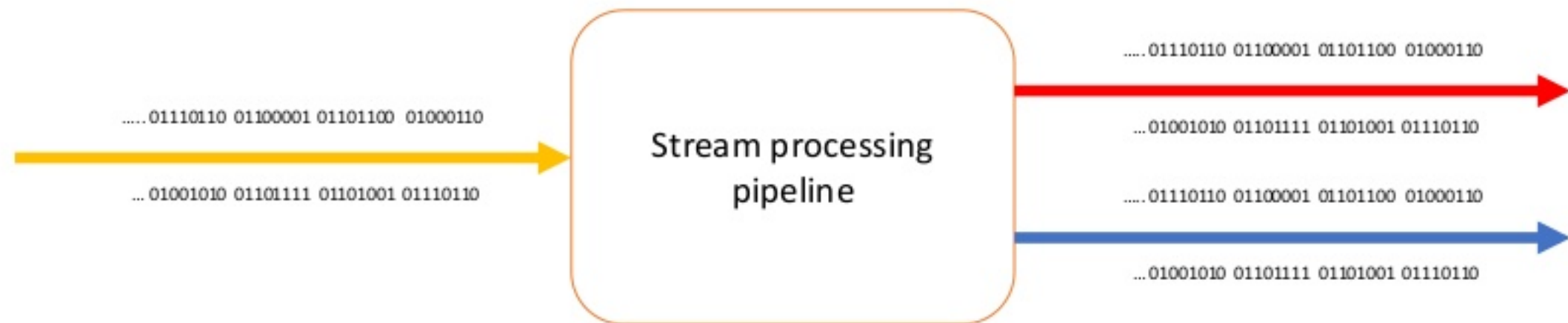- Process bits
- Store bits
- Transmit bits

```
01000110        01000110        01000110
01101100        01101100        01101100
01100001        01100001        01100001
01110110        01110110        01110110
01101001        01101001        01101001
01101111        01101111        01101111
01001010        01001010        01001010
01110101        01110101        01110101
01101110        01101110        01101110
01110001        01110001        01110001
01110101        01110101        01110101
01100101        01100101        01100101
01101001        01101001        01101001
01110010        01110010        01110010
01100001        01100001        01100001
00101100        00101100        00101100
01000010        01000010        01000010
01100001        01100001        01100001
01110010        01110010        01110010
01100011        01100011        01100011
01100101        01100101        01100101
01101100        01101100        01101100
01101111        01101111        01101111
01101110        01101110        01101110
01100001        01100001        01100001
```

Flavio Junqueira,   Flavio Junqueira,   Flavio Junqueira,
is in Barcelona     is in Berlin        is back in Barcelona
Sep. 10 - noon      Sep. 11 – noon      Sep. 15 - noon

Time

- Order matters
- Correlation between events
- Causality maybe?

# Processing data streams

# Processing streams

.....01110110 01100001 01101100 01000110

... 01001010 01101111 01101001 01110110

Stream processing pipeline

.....01110110 01100001 01101100 01000110

...01001010 01101111 01101001 01110110

.....01110110 01100001 01101100 01000110

...01001010 01101111 01101001 01110110

# A typical architecture

Source



One or more stages

..... 01110110 01100001 01101100  01000110

... 01001010 01101111 01101001 01110110

| Messaging substrate | Stream data processor | Messaging substrate | Stream data processor |

..... 01110110 01100001 01101100

... 01001010 01101111 01101001

- Ingests and buffers data
- Decouples source from the engine processing the data

# A typical architecture

Limitations:
- Data stored temporarily
- Not able to store an unbounded amount of stream data

Source

Multiple stages

.....01110110 01100001 01101100 01000110

... 01001010 01101111 01101001 01110110

| Messaging substrate | Stream data processor | Messaging substrate | Stream data processor |

..... 01110110 01100001 01101100

... 01001010 01101111 01101001

- Ingests and buffers data
- Decouples source from the engine processing the data

# A typical architecture

Limitations:
- Data stored temporarily
- Not able to store an unbounded amount of stream data
- Separate bulk store for historical reads

Source

Multiple stages

.....01110110 01100001 01101100  01000110

... 01001010 01101111 01101001 01110110

| Messaging substrate | Stream data processor | Messaging substrate | Stream data processor |

..... 01110110 01100001 01101100

- Ingests and buffers data
- Decouples source from the engine processing the data

Some bulk data store

Target of Pravega is a stream store able to:

- Store stream data permanently
- Preserve order
- Accommodate unbounded streams

# Streams in Pravega

# Pravega and Streams

*What about parallelism?*

# Pravega and Streams

# Segments in Pravega



- Segments are sequences of bytes

# Segments in Pravega



Pravega

Append

..... 01110110 01100001 01101100

01000110

Read

01000110

..... 0110100101110110 01001010

01110110

01110110

01101001

01101001

01101111

01101111

Segments

- Segments are sequences of bytes
- Use routing keys to determine segment

*Segments can be sealed*

# Segments in Pravega

Pravega

Append

..... 01110110 01100001 01101100

01000110

Read

01000110

..... 0110100101110110 01001010

01110110

01110110

01101001

01101001

01101111

01101111

Segments

Once sealed, a segment can't be appended to any longer.

# Segments in Pravega

Pravega



Segments

Pravega is primarily:
- A segment store
- Segments sealed or open

*How is sealing segments useful?*

# Segments in Pravega

Pravega

- Each segment can live in a different server
- Not limited to the capacity of a single server

01101111  01000110

01000110

01000110  01101111

01000110

01000110  01101111

Compose to form a stream

01101111

01101111

01000110

01101111

Segments

01000110  01101111  01101111

Stream

# Segments in Pravega

Pravega

01101111    01000110

01000110

01000110    01101111

01000110

01000110    01101111

01101111

01101111

01000110

01101111    Segments

Compose to form a stream

01000110    01101111

01101111

Stream

*Some useful ways to compose segments*

# Scaling a stream

..... 01110110 01100001 01101100 $\longrightarrow$ | 01000110 |

Stream has one
segment

**1**

- Say input load has increased
- Need more parallelism

..... 01110110 01100001 01101100 $\longrightarrow$

| 01000110 |
| 01000110 | | 01000110 |
| 01000110 |

- Seal current segment
- Create new ones

**2**

Key ranges are not statically assigned to segments

Segment Heat Map

Routing Key Space

10:50 AM  11:00 AM  11:10 AM  11:20 AM  11:30 AM  11:40 AM  11:50 AM  12:00 PM  12:10 PM  12:20 PM  12:30 PM  12:40 PM

*Source*: Virtual cluster - Nautilus Platform

Segment Heat Map

Scale down

Scale up

*Source*: Virtual cluster - Nautilus Platform

# Transactions



Stream has two segments

**1**

Begin txn

**2**

Write to txn

**3**

Write to txn

**4**

Upon commit

**5**

Upon commit

**6**

# Transactions



Stream has two segments

**1**

Begin txn

**2**

Write to txn

**3**

Write to txn

**4**

Upon abort

**5**

*Wait, how are segments manipulated?*

# Controller

- Control plane
- A few of the controller tasks
  - Stream lifecycle
    - Create
    - Delete
    - Scale
  - Txn management
    - Create
    - Commit/Abort

Pravega

# API – *Writer* and *Reader*

# Events

- Internally
  - Pravega is all about bytes
- Current API focused on events
  - Some encapsulation of application bytes
  - `Serializer` interface

```java
public interface Serializer<T> {
  /**
   * Serializes the given event.
   *
   * @param value The event to be serialized.
   * @return The serialized form of the event.
   * NOTE: buffers returned should not exceed {@link #MAX_EVENT_SIZE}.
   */
  ByteBuffer serialize(T value);

  /**
   * Deserializes the given ByteBuffer into an event.
   *
   * @param serializedValue A event that has been previously serialized.
   * @return The event object.
   */
  T deserialize(ByteBuffer serializedValue);
}
```

# EventStreamWriter API

```java
String scope = "myScope";
WriterConfig config = new WriterConfig();
String streamName = "myStream";

ClientFactory factory = ClientFactory.withScope(scope, new URI("//demo.pravega.io:3333"));
EventStreamWriter<String> writer = factory.createEventWriter(streamName, serializer, config);

while(!worldEnd) {
    /* E.g., getNewRecord() reads the next line in a file */
    String record = getNewRecord();
    String key = extractKey(record);
    String event = extractEvent(record);
    writer.writeEvent(key, event);
}
```

# Using the `EventStreamReader` API

```java
String scope = "myScope";
String myReaderId = "myId";
String myReadGroup = "myGroup";
ReaderConfig config = new ReaderConfig(props);
String streamName = "myStream";

ClientFactory factory = ClientFactory.withScope(scope, new URI("//127.0.0.1:3333");
EventStreamReader<String> reader = factory.createEventReader(myReaderId,
                                                             myReadGroup,
                                                             serializer,
                                                             config);
while(!worldEnd) {
  EventRead<String> event = reader.readNextEvent(long timeout);

  // Application consumes event and it is supposed to persist the Position
  // object to guarantee exactly once.
  consumeEvent(event.getEvent(), event.getPosition());
}

reader.close();
```

# Transactions

```
Txn txn = writer.beginTxn();

txn.writeEvent(getKey("Pravega"), "Pravega");
txn.writeEvent(getKey("is"), "is");
txn.writeEvent(getKey("invading"), "invading");
txn.commit();
```
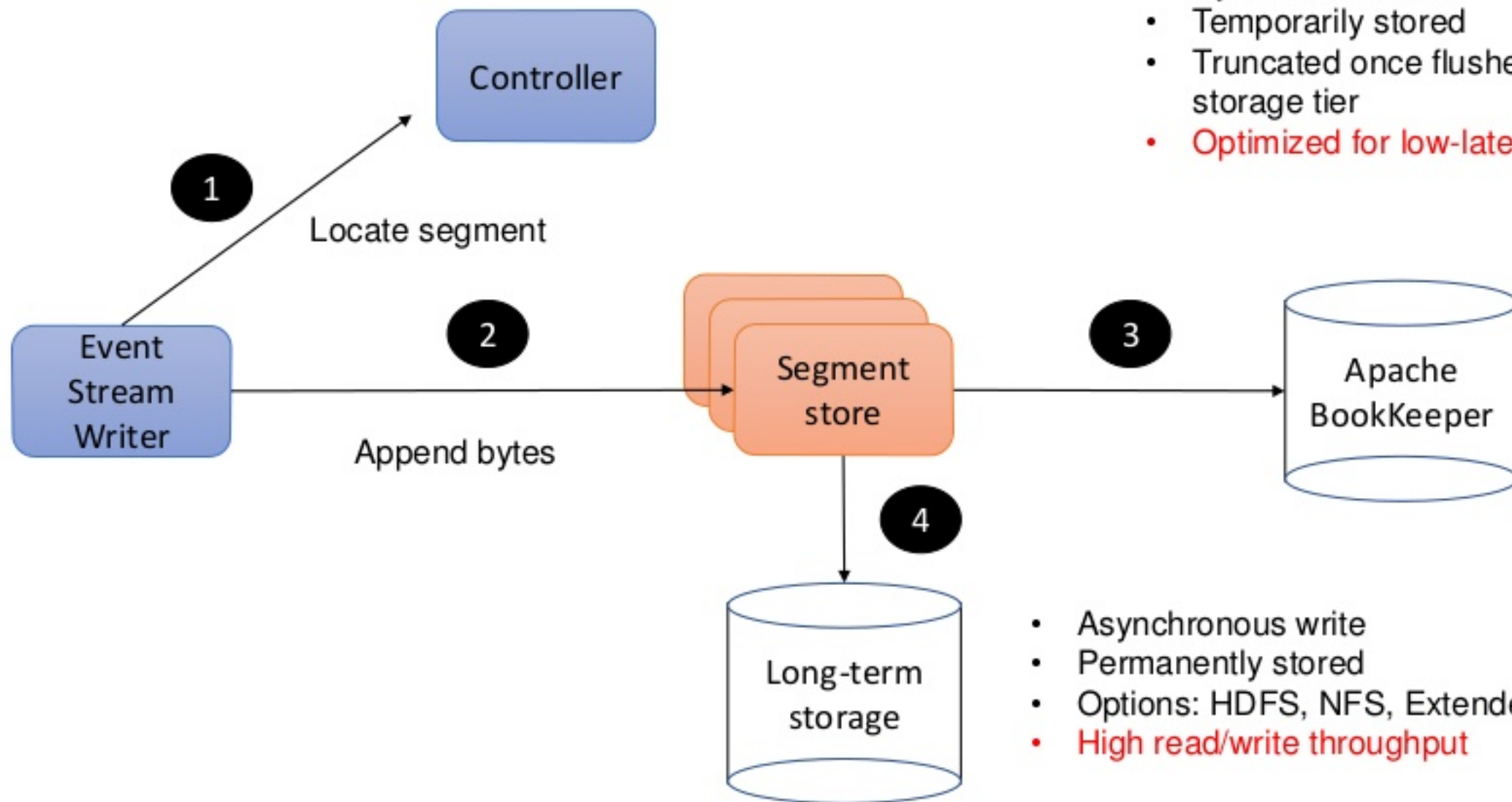
# Pravega semantics

# The write path



Controller

- Synchronous write
- Temporarily stored
- Truncated once flushed to next storage tier
- Optimized for low-latency writes

1 Locate segment

Event Stream Writer

2 Append bytes

Segment store

3

Apache BookKeeper

4

Long-term storage

- Asynchronous write
- Permanently stored
- Options: HDFS, NFS, Extended S3
- High read/write throughput

# Guarantees on the write path
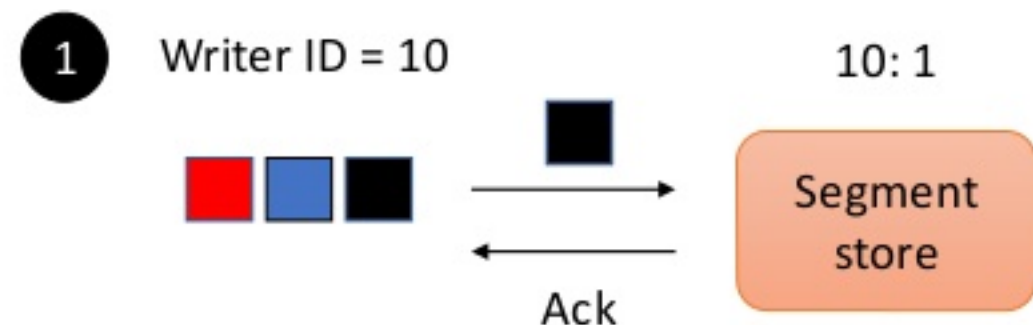
- **Order**
  - Writer appends in application order
- **Duplicates**
  - Writer IDs
  - Maps to last appended data on the segment store
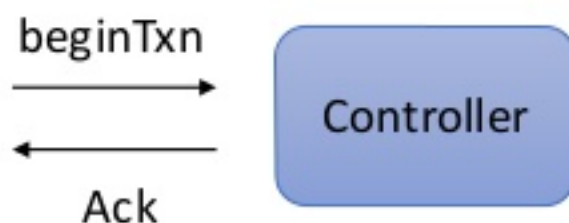  - Writer does not persist ID to tolerate a crash
- Txns
  - Atomicity at the stream level
  - If anything goes wrong with the writes, either abort or let it time out
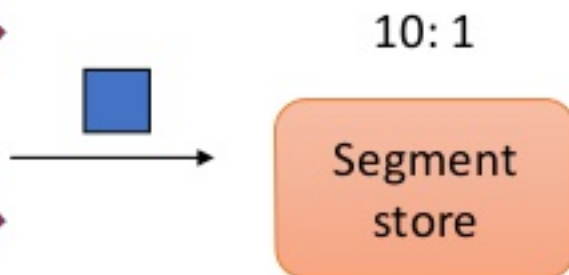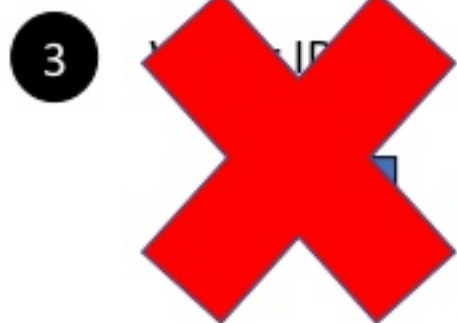
# Avoiding duplicates – Reconnect

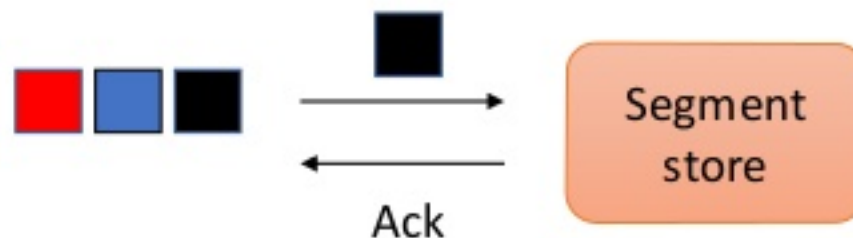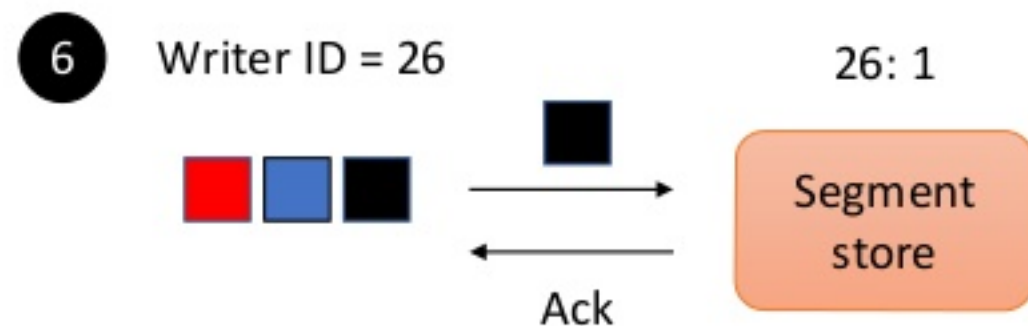# Avoiding duplicates – Transactional writes

**1** Writer ID = 10

beginTxn →

← Ack

Controller

**2** Writer ID = 10

10: 1

→

← Ack

Segment store

**3** 10: 1

→

Segment store

**4** Controller

Eventually times out and aborts txn

# Avoiding duplicates – Transactional writes

**5**    Writer ID = 26

beginTxn →

← Ack

Controller

**6**    Writer ID = 26

26: 1

→

← Ack

Segment store

# The read path

Controller

**1** Locate segment

Event Stream Reader

**2** Read bytes

Segment store

**4** Bytes read

**3**

Long-term storage

- Used for recovery alone
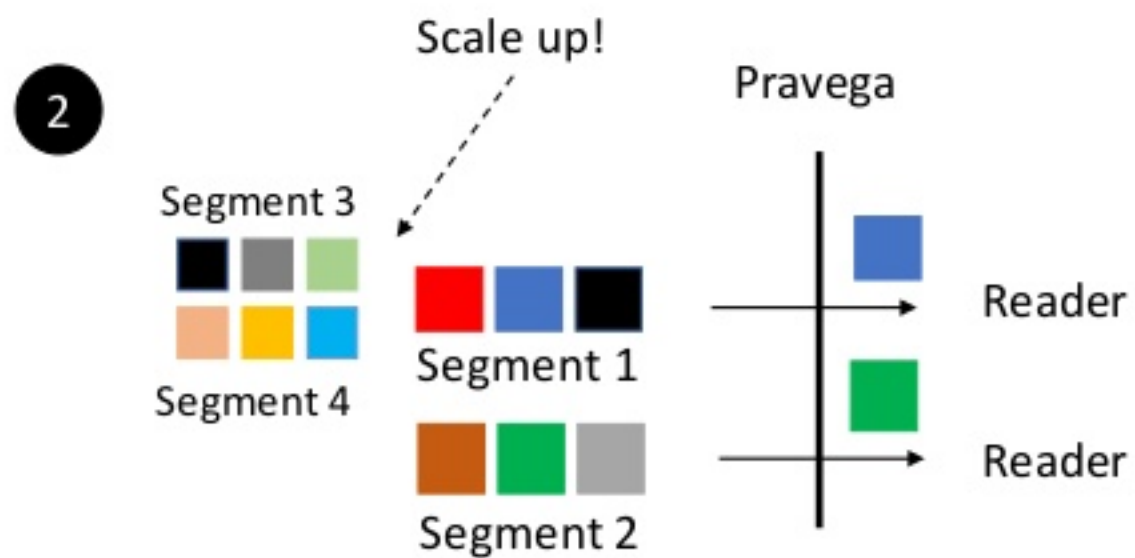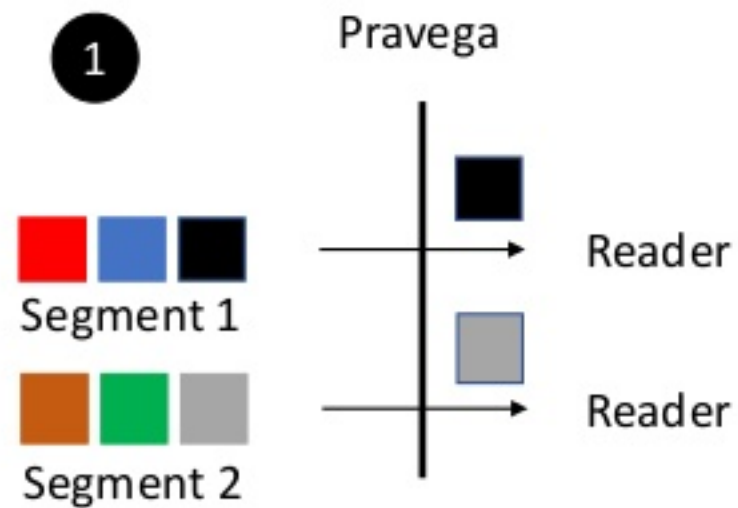- Not used to serve reads

Apache BookKeeper

- Bytes read from memory
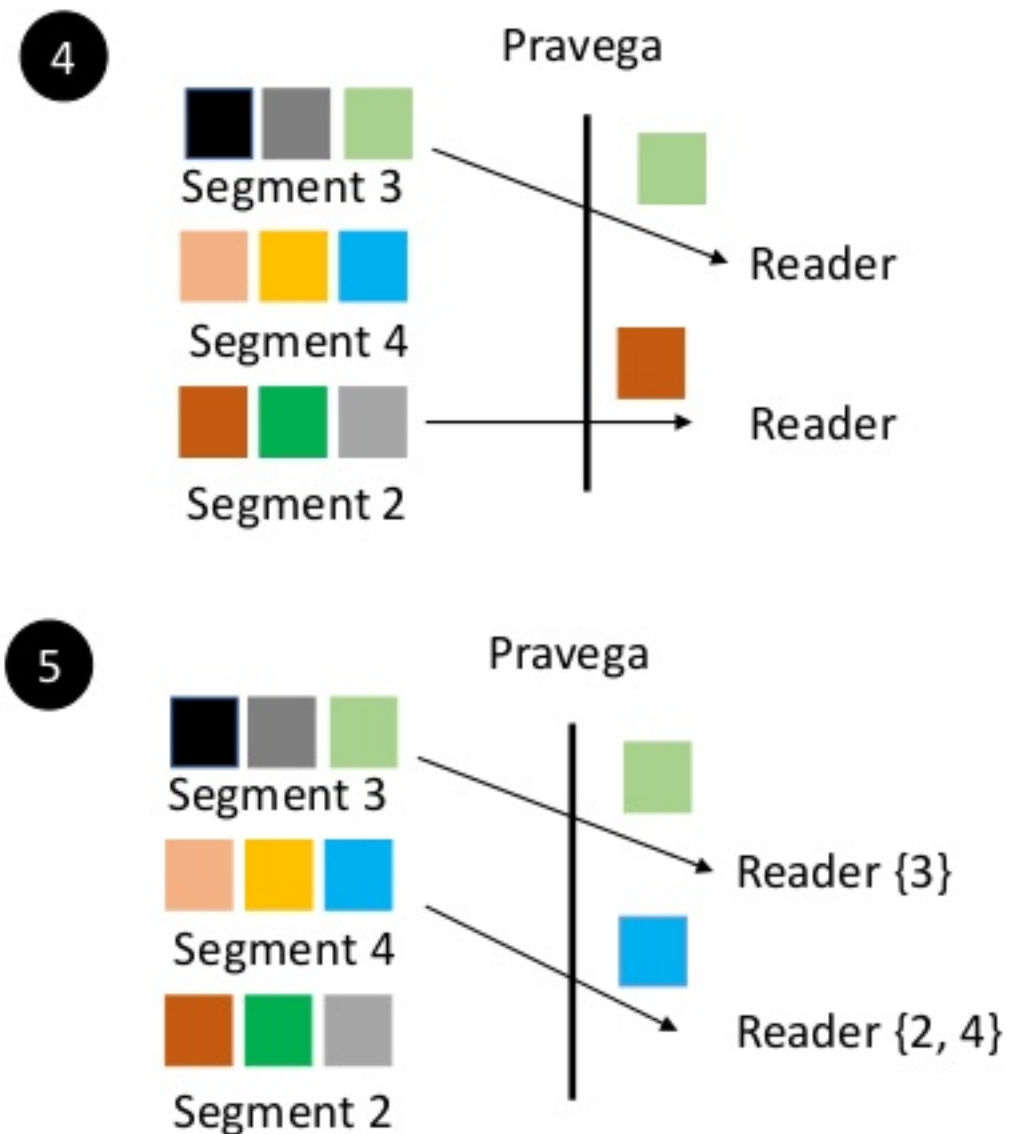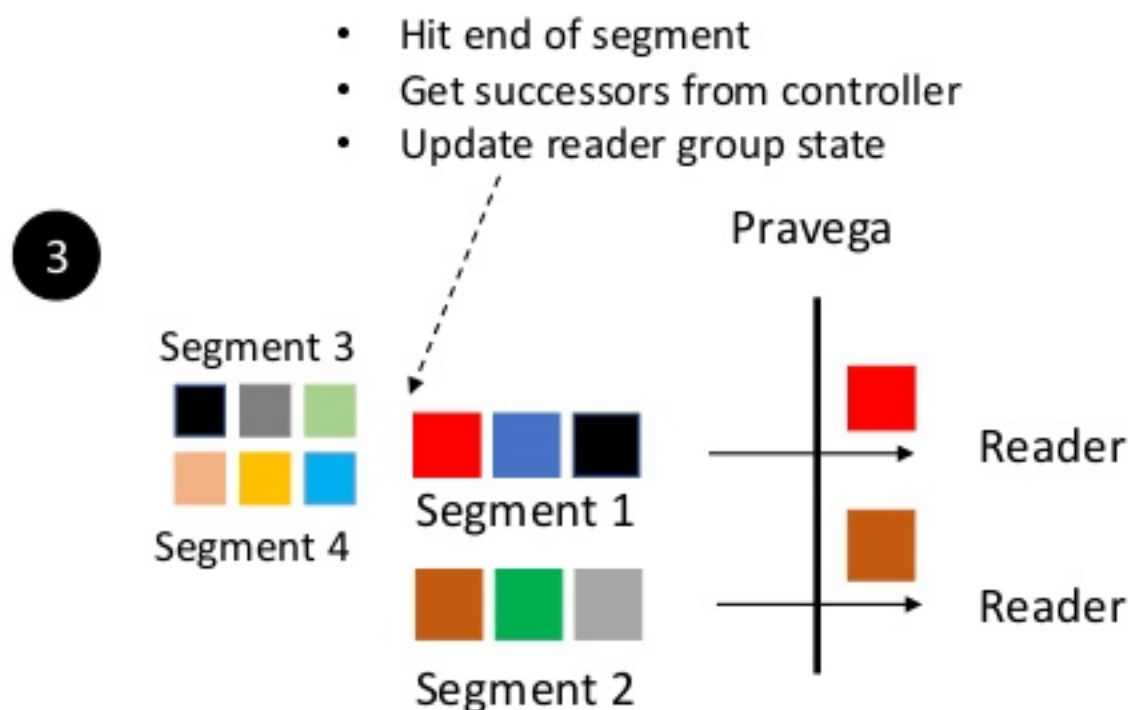- If not present, pull data from Tier 2

# Reader groups

- Group of event readers
  - Read events from a set of streams
  - Load distributed across readers of the group

- Segments
  - A given reader reads from a set of segments
  - Coordination of segment assignment done via a **state synchronizer**

- State synchronizer
  - General facility for synchronizing state across processes
  - Uses a revisioned Pravega segment
  - Advanced topic: not explained further in this talk
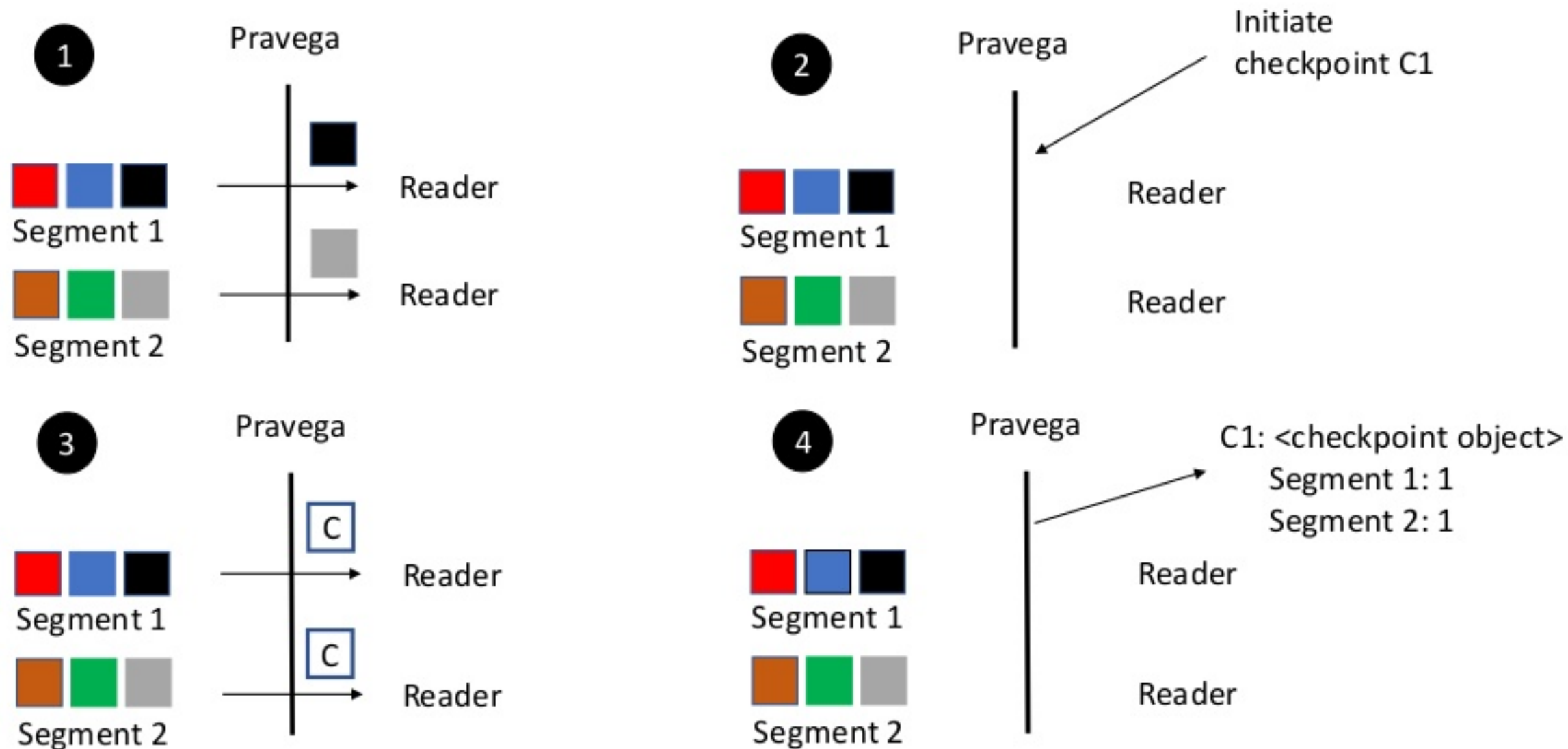
# Reader groups + Scaling

# Reader groups + Scaling

- Hit end of segment
- Get successors from controller
- Update reader group state

**3**

Segment 3

Segment 4

Segment 1

Segment 2

Pravega

Reader

Reader

**4**

Pravega

Segment 3

Segment 4

Segment 2

Reader

Reader

**5**

Pravega
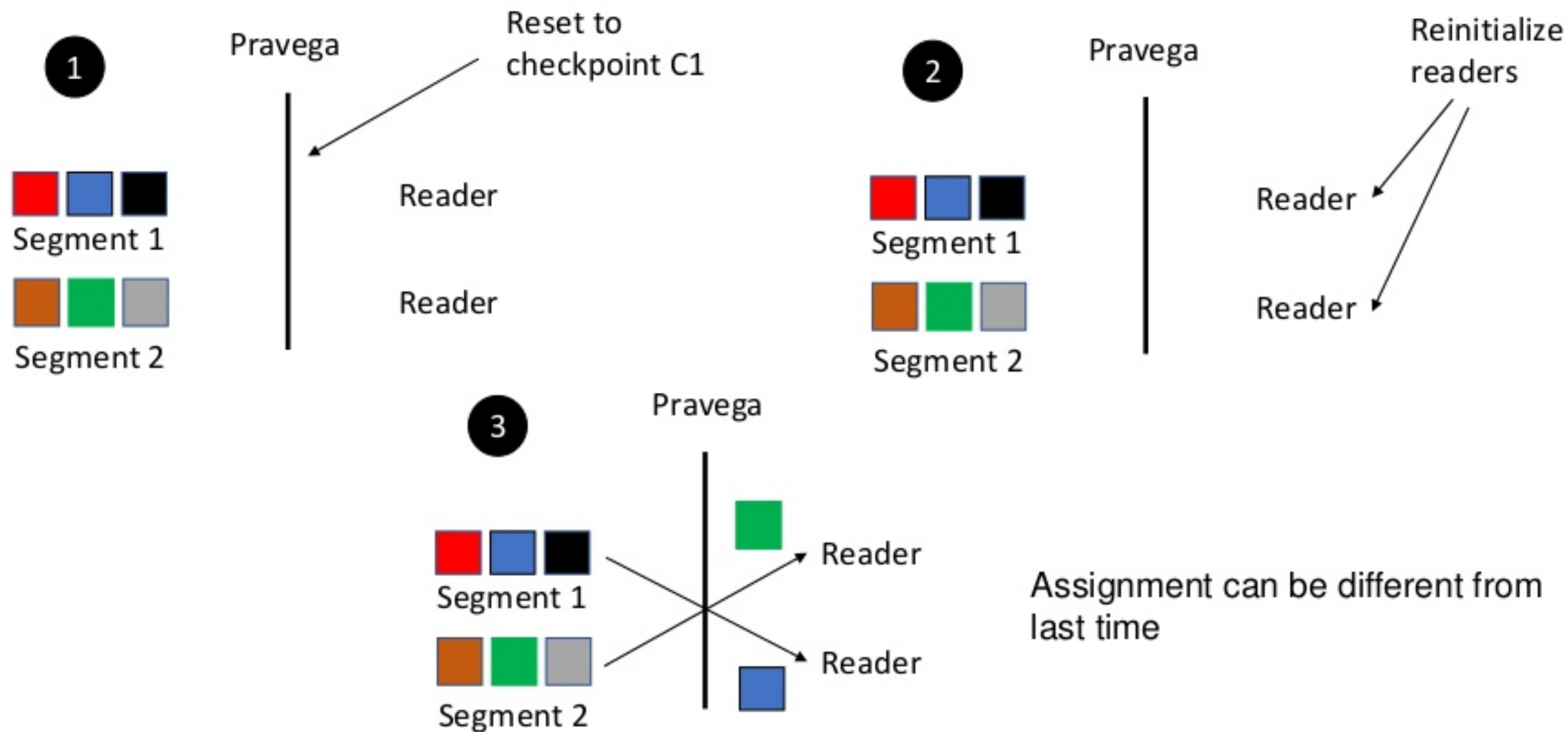
Segment 3

Segment 4

Segment 2

Reader {3}

Reader {2, 4}

# `Checkpoint` object

- Maps segments to corresponding offsets

- Opaque to the application

# Getting a checkpoint

# Resetting to a checkpoint

**①** Pravega

Reset to checkpoint C1

Segment 1

Segment 2

Reader

Reader

**②** Pravega

Reinitialize readers

Segment 1

Segment 2

Reader

Reader

**③** Pravega

Segment 1

Segment 2

Reader

Reader

Assignment can be different from last time

# Wrap up

# Take-away messages

- Pravega is all about
  - Unbounded stream data
  - Permanently stored
  - Elasticity for streams
  - Scaling producers and consumers independently
- Under active development
- Looking at first use cases

# Ongoing work

- Performance tuning

- Scaling support

- Event-time support

- Geo-distribution
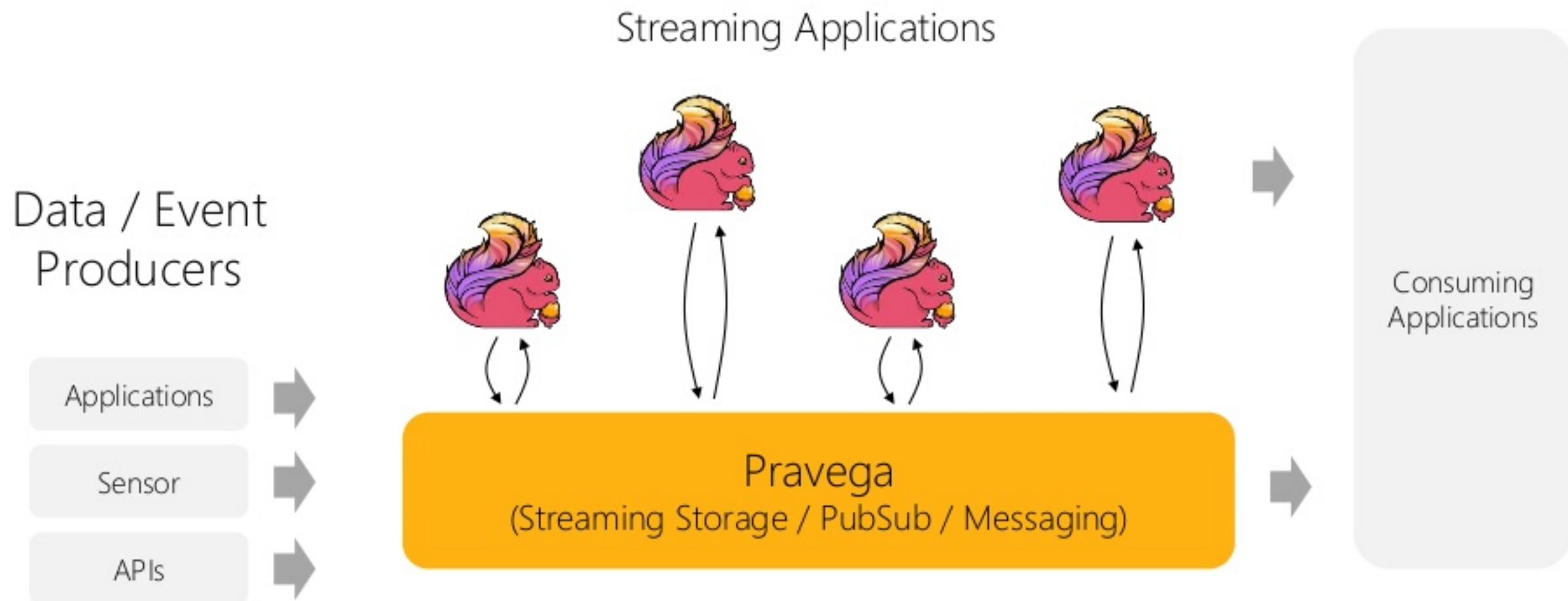
- Security

- ... and much more

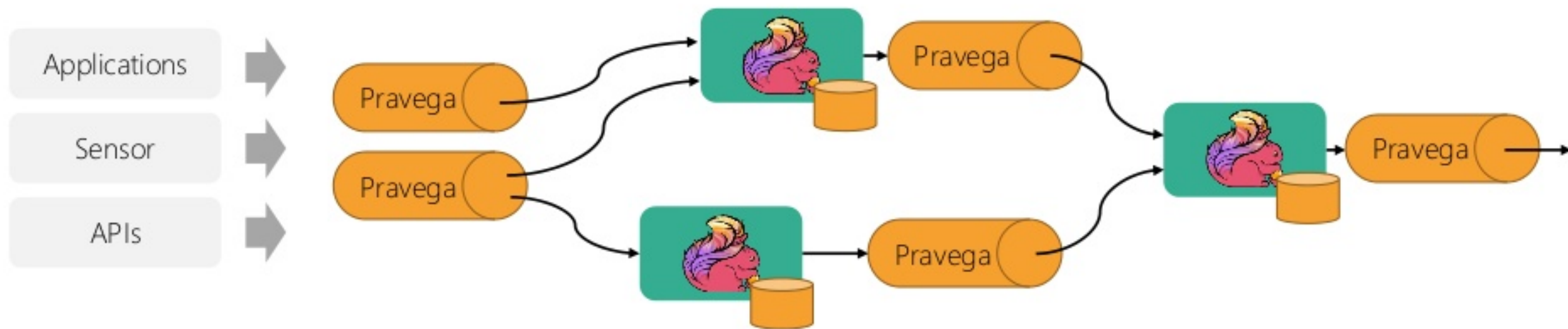http://pravega.io

E-mail: fpj@apache.org
Twitter: @fpjunqueira

Join the community!

Pravega +

# Streaming Storage and Compute

Streaming Applications

Data / Event
Producers

Consuming
Applications

Applications

Sensor

APIs

**Pravega**
(Streaming Storage / PubSub / Messaging)
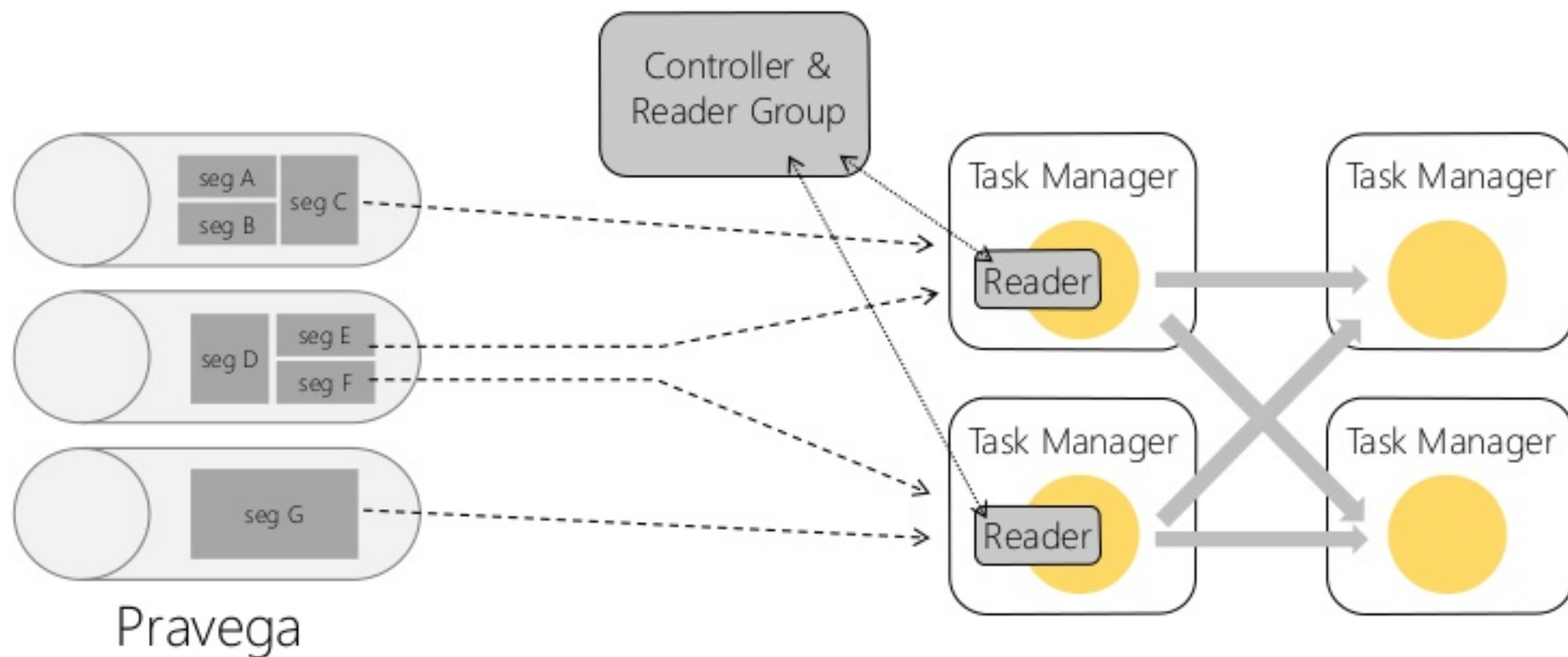
# Streaming Pipelines
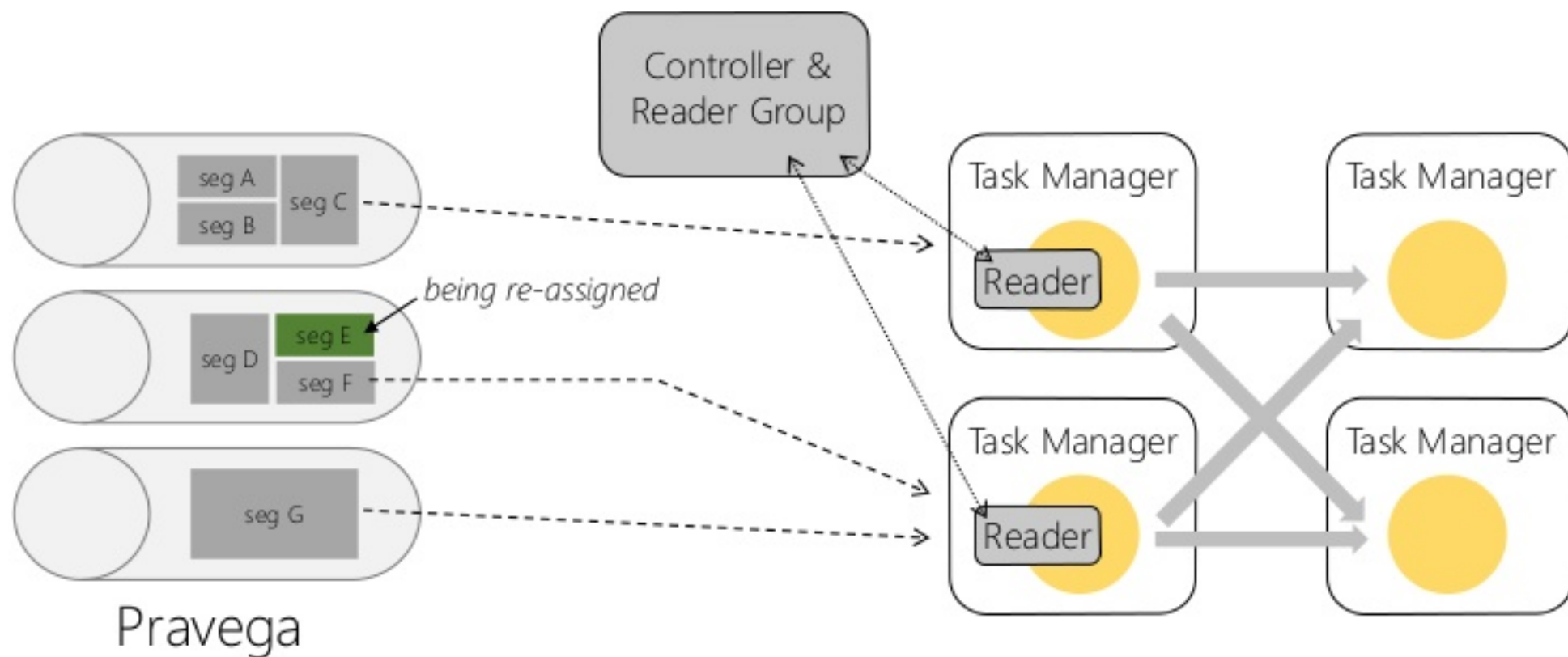
# Flink reading from Pravega

# Reading via ReaderGroup

- Readers do not choose their own segments
- ReaderGroup automatically assigns and re-balances segments
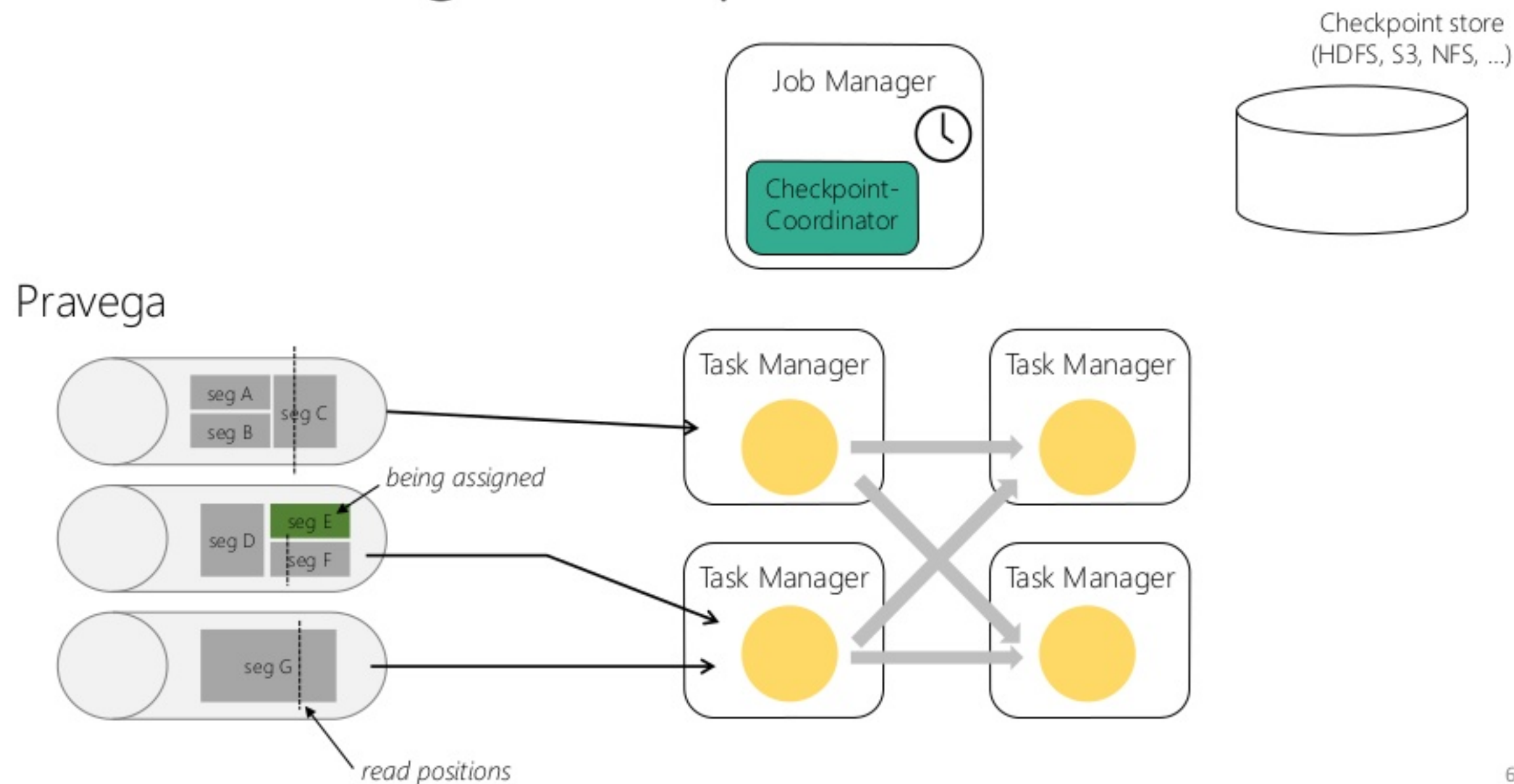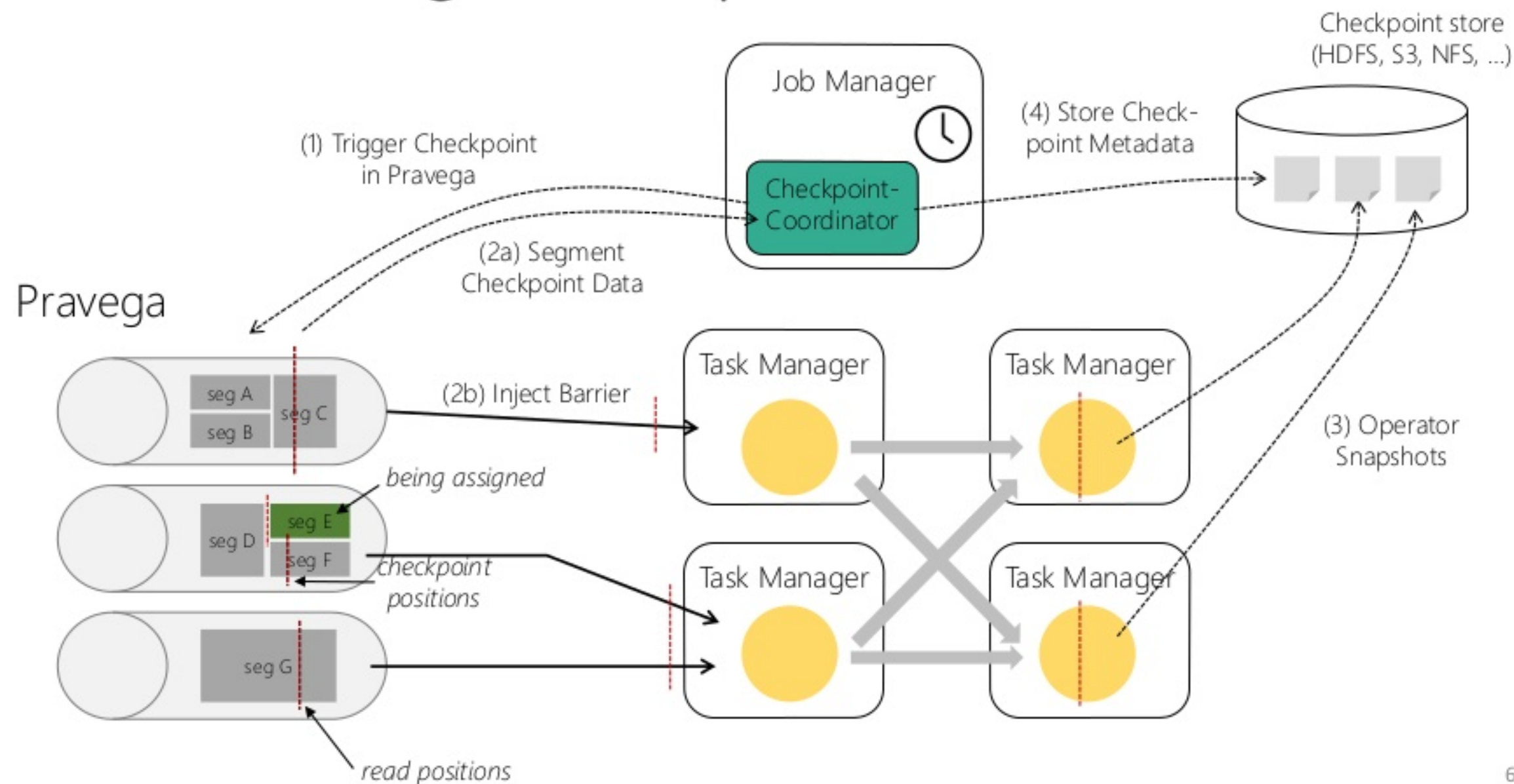- Leaving the ReaderGroup in charge is key to automatic scaling



Pravega

# Reading via ReaderGroup

- At points in time, some segments may be not assigned to any reader
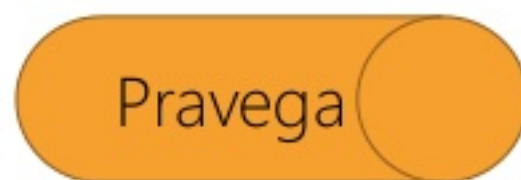- Example: New segments, re-balancing segments, …



Pravega

# Flink + Pravega Checkpoints

Checkpoint store
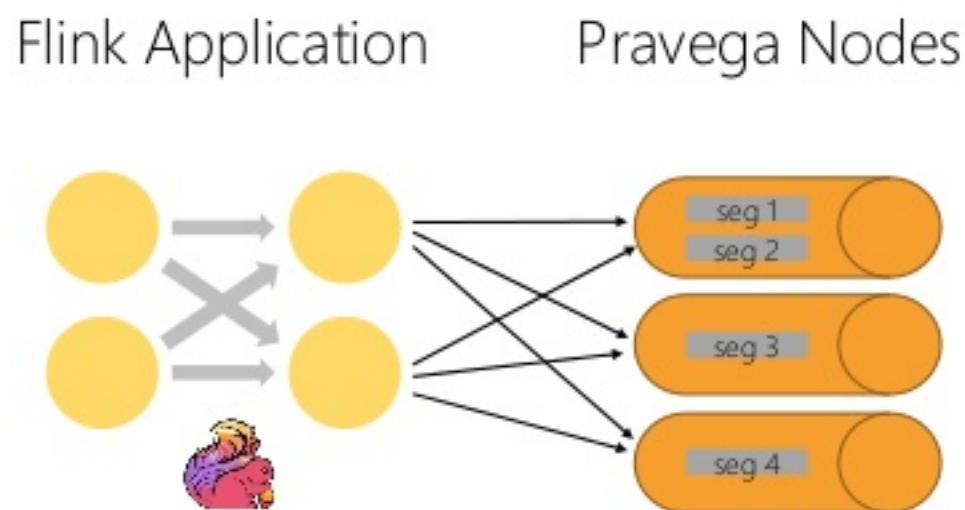(HDFS, S3, NFS, ...)

Job Manager

Checkpoint-
Coordinator

Pravega

seg A

seg B

seg C

*being assigned*

seg D

seg E

seg F

seg G

*read positions*

Task Manager

Task Manager

Task Manager

Task Manager

# Flink + Pravega Checkpoints



Checkpoint store
(HDFS, S3, NFS, ...)

Job Manager

(1) Trigger Checkpoint
in Pravega

(4) Store Check-
point Metadata

Checkpoint-
Coordinator

(2a) Segment
Checkpoint Data

Pravega

seg A
seg B
seg C

(2b) Inject Barrier

Task Manager

Task Manager

being assigned

seg E

seg D
seg F

checkpoint
positions

Task Manager

Task Manager

(3) Operator
Snapshots

seg G

read positions

64

# Flink writing to Pravega

# The FlinkPravegaWriter
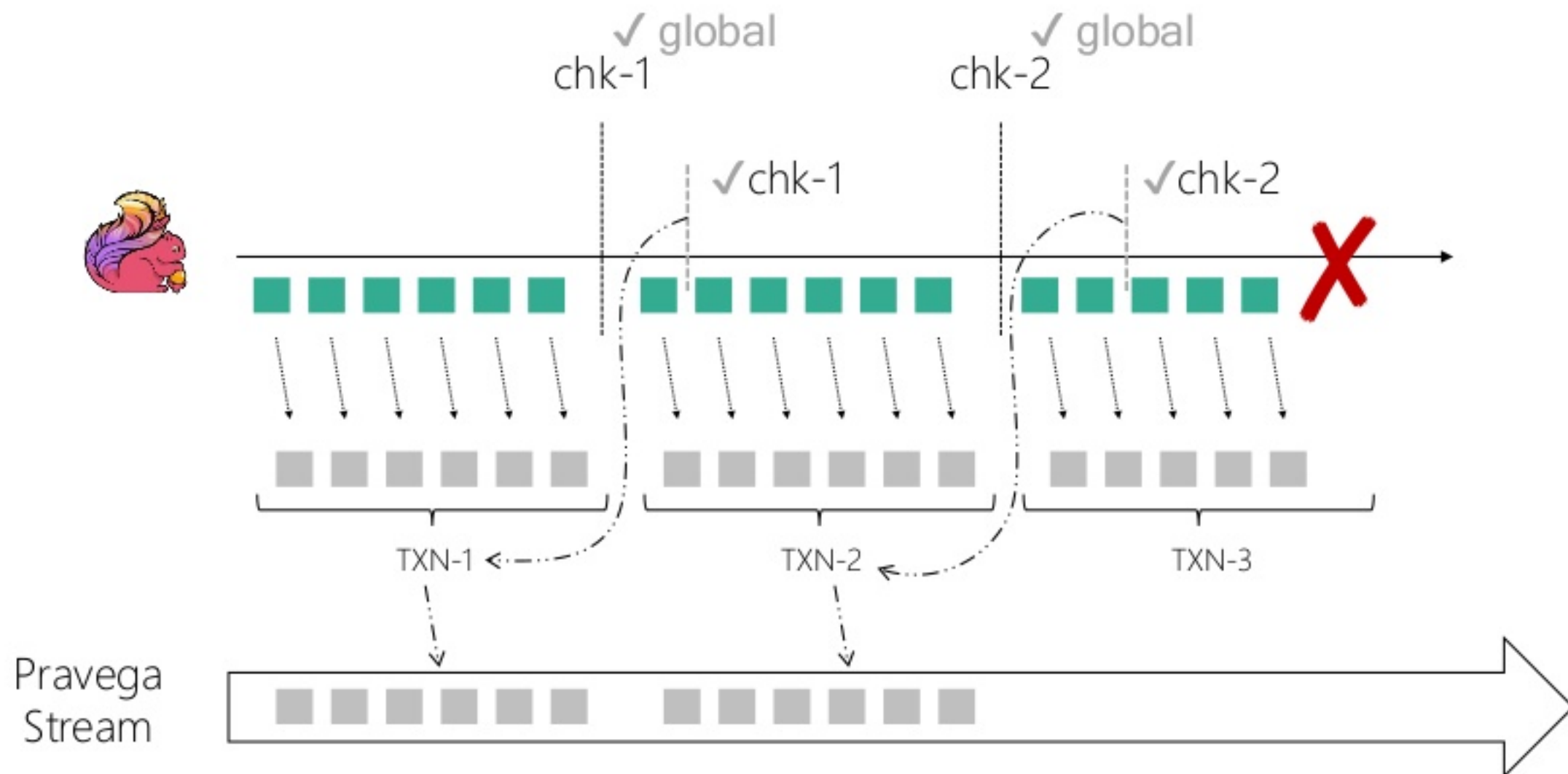
Flink Application        Pravega Nodes

- **Regular Flink SinkFunction**

- No partitioner, but a "routing key"

- Remember: No partitions in Pravega
  - Just dynamically created segments

- Same key always goes to the same segment

- Order of elements guaranteed per key!
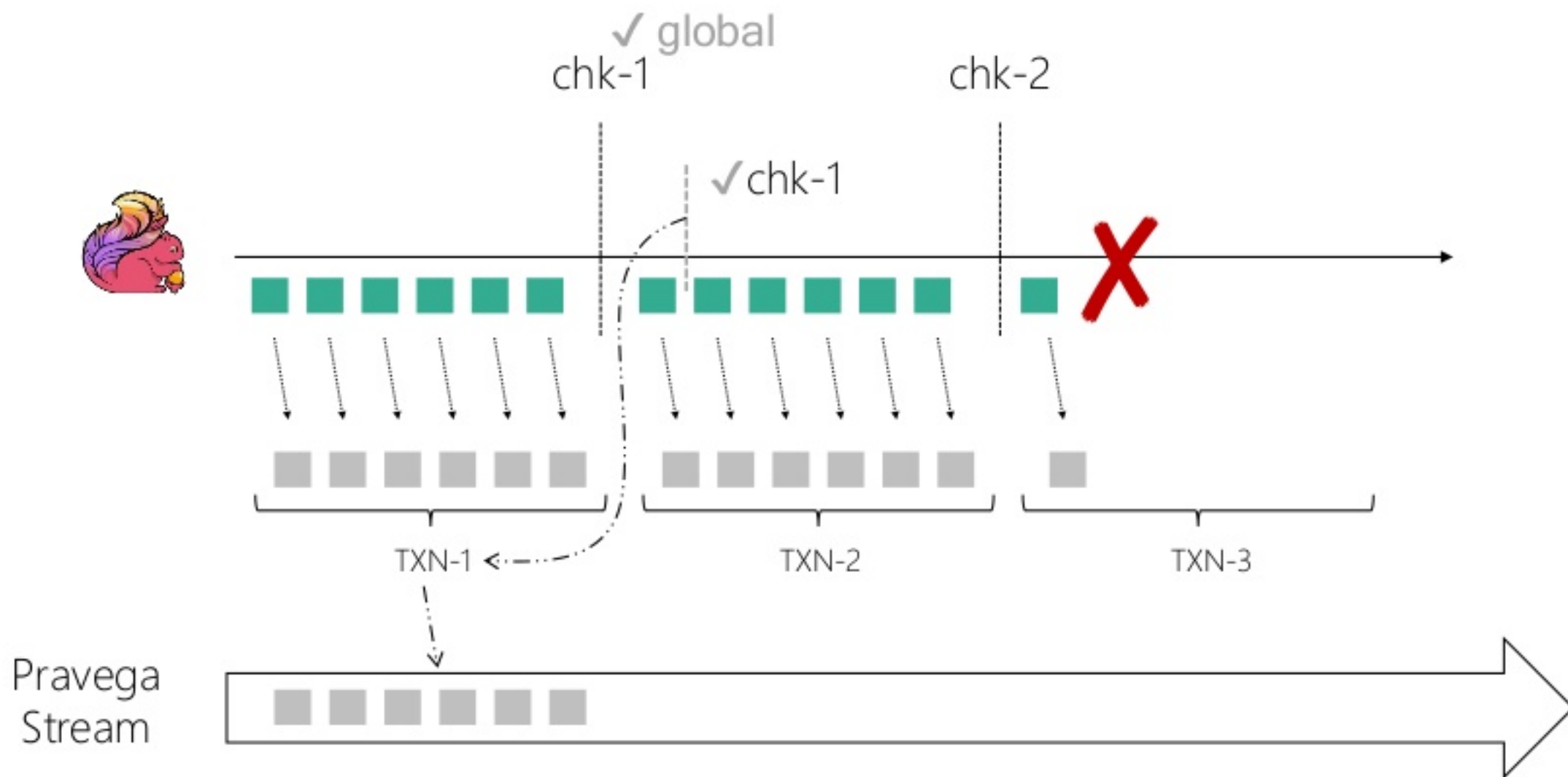
seg 1
seg 2
seg 3
seg 4

# Exactly-once via Transactions

- Similar to a distributed 2-phase commit

- Coordinated by asynchronous checkpoints, no voting delays

- Basic algorithm:
  - Between checkpoints: Produce into transaction
  - On operator snapshot: Flush local transaction *(vote-to-commit)*
  - On checkpoint complete: Commit transactions
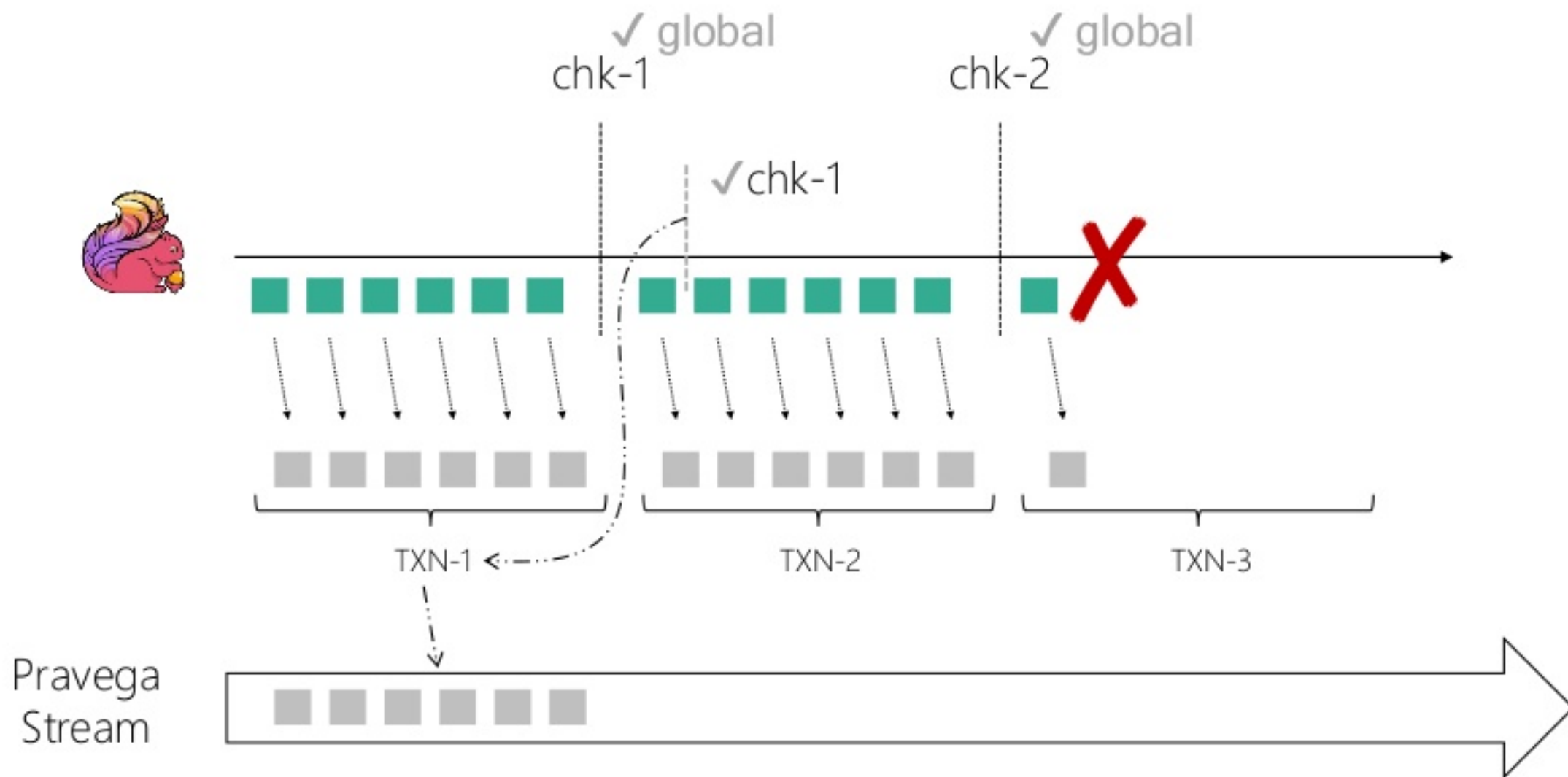  - On recovery: check and commit any pending transactions
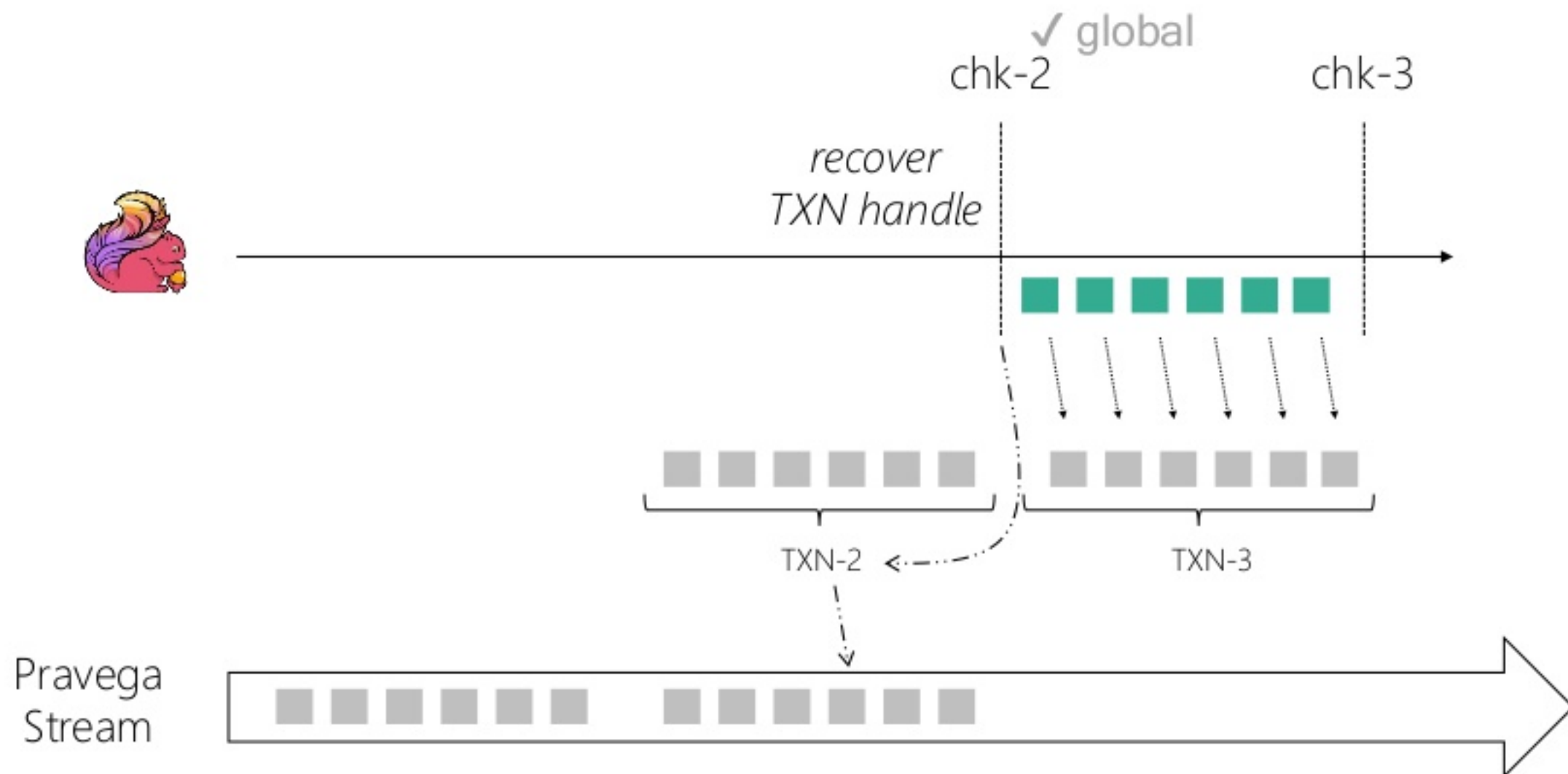
# Exactly-once via Transactions

# Transaction fails after local snapshot

# Transaction fails before commit...

# ... commit on recovery

# Looking ahead...

# Looking ahead...

Automatic Scaling

Flink follows Pravega's scaling
*(at least the first stage)*

High Availability through Pravega

Use synchronizers
instead of ZooKeeper
*(leader election,
distributed atomic counters, ...)*

# Questions?

http://pravega.io
http://github.com/pravega/pravega
http://github.com/pravega/flink-connectors

E-mail: sewen@apache.org, fpj@apache.org
Twitter: @StephanEwen, @fpjunqueira