

On-the-fly Reconfiguration of Query Plans for Stateful Stream Processing Engines

BTW 2019 - 7. March 2019

Adrian Bartnik

Bonaventura Del Monte

Dr. Tilmann Rabl

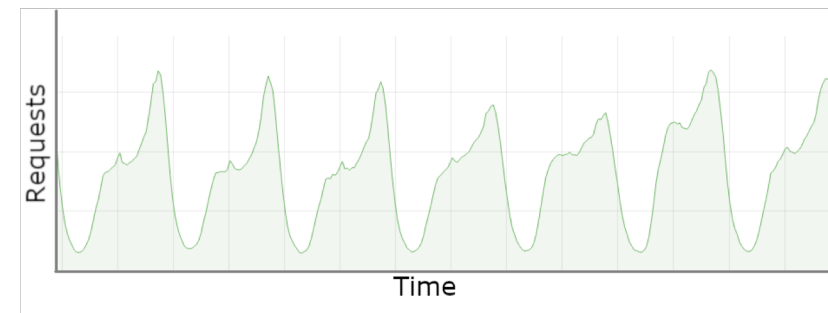
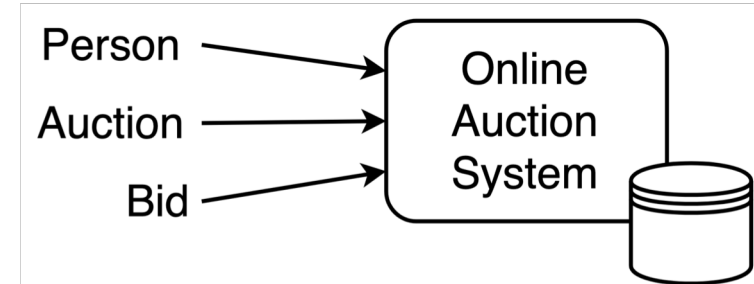
Dr. Volker Markl



Fachgebiet Datenbanksysteme und Informationsmanagement
Technische Universität Berlin

<http://www.dima.tu-berlin.de/>

- Stream processing engines emphasize **processing velocity**
- SPEs can handle **unbounded data sets**
- Many streaming workloads **vary over time**



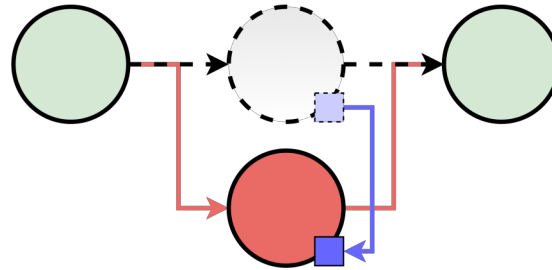
➡ Adapt streaming job to incoming workload

- Restarting job only way to modify running streaming job
 - May involve distribution of large state
 - Require external system to guarantee correct processing semantics
 - Not possible if other systems rely on streaming output

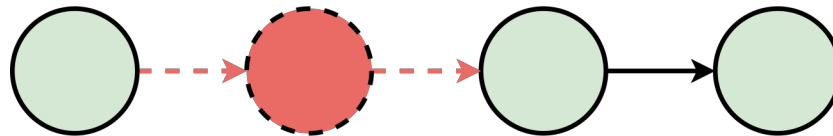


- Modify a running streaming job in Apache Flink

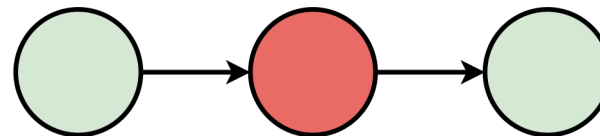
- Stateful operator migration



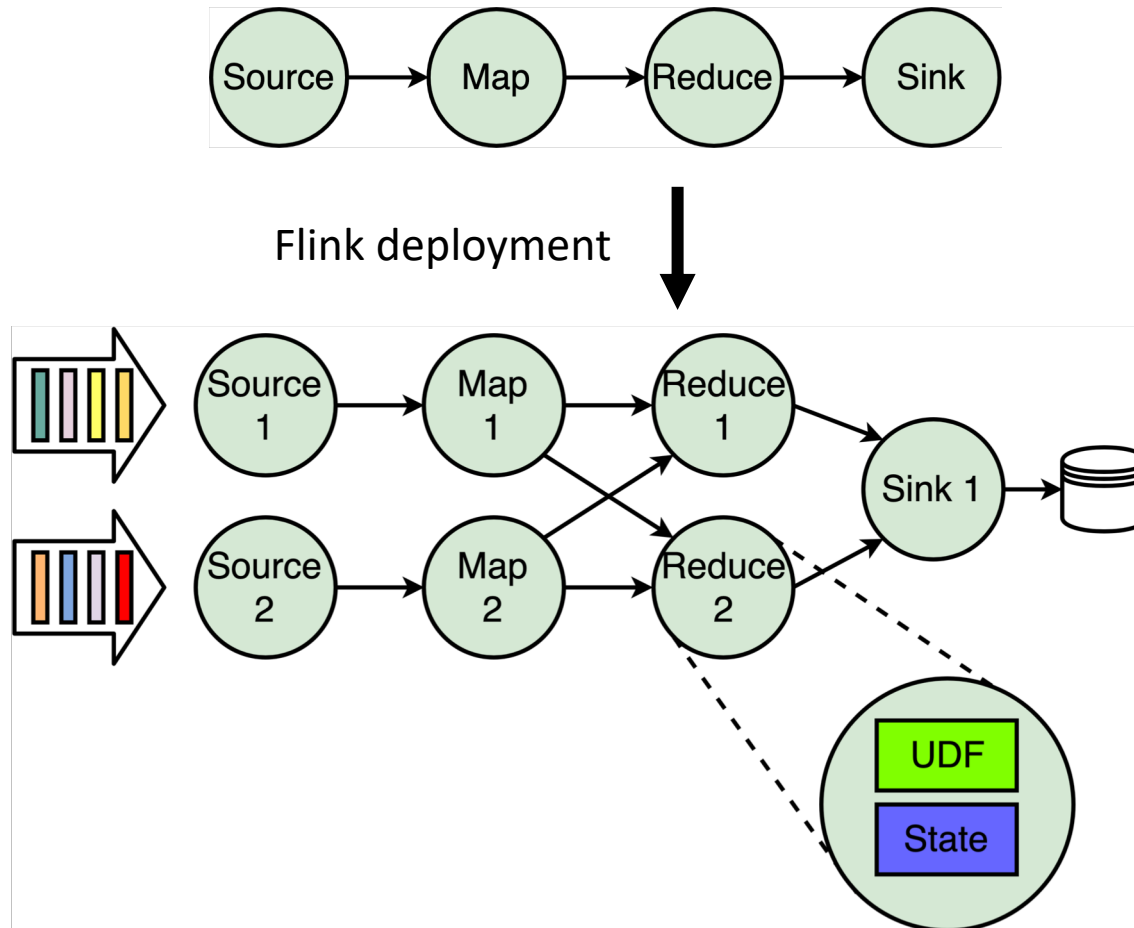
- Introduction of new operators



- Optimize or replace existing operators

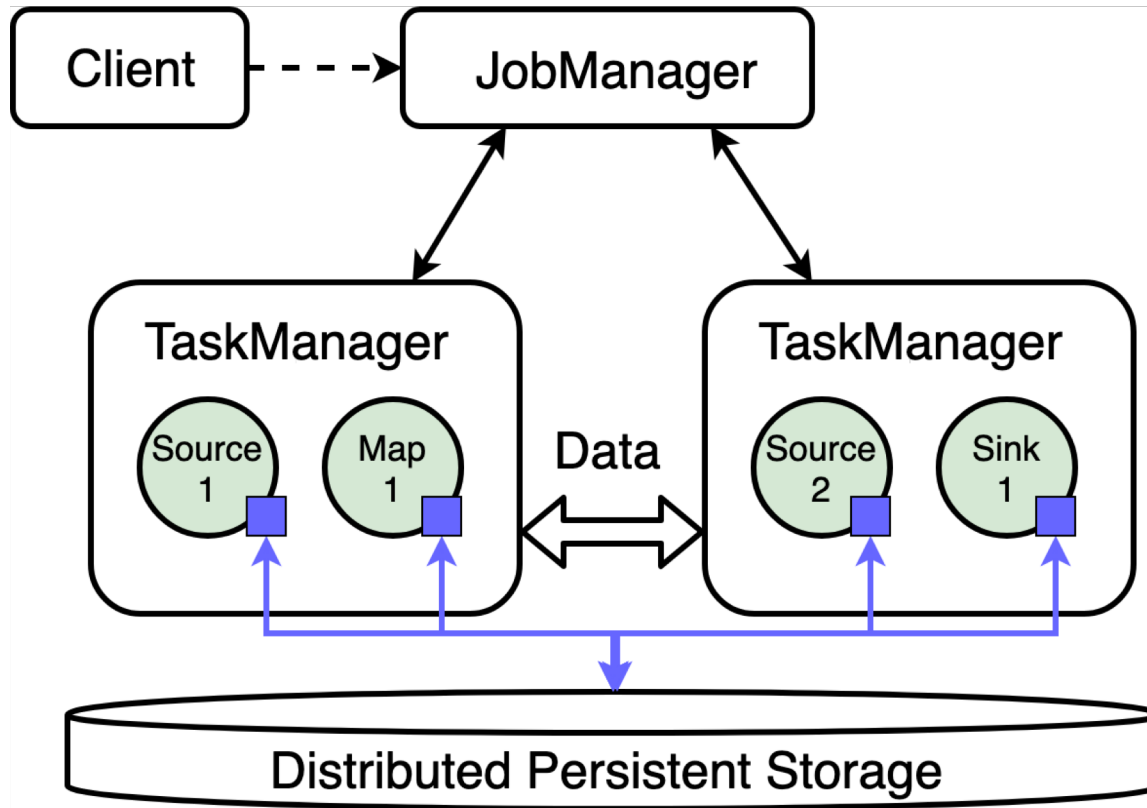


- Dataflow of streams and operators form DAG

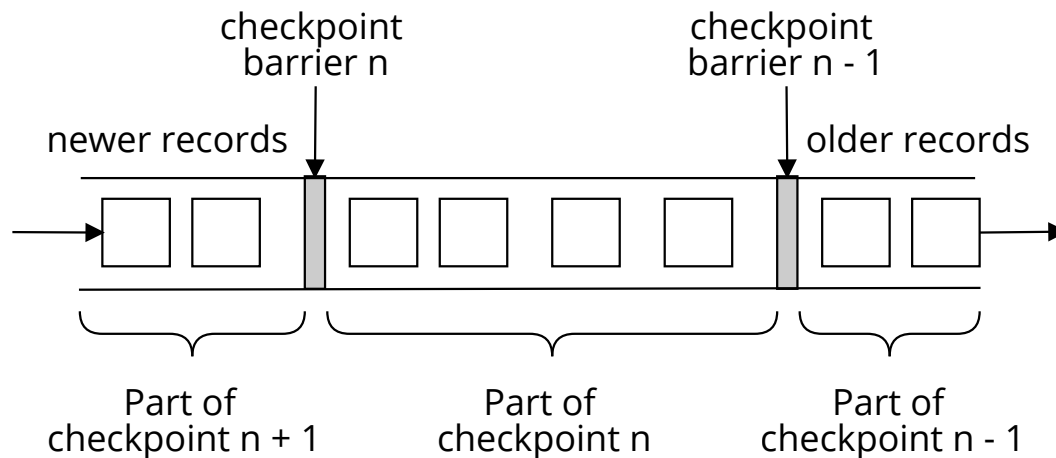




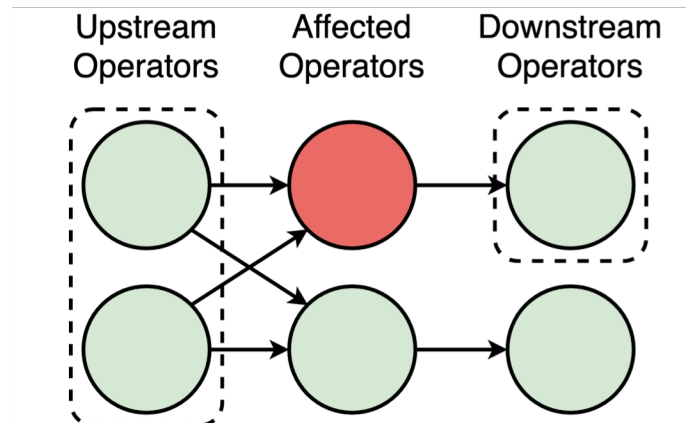
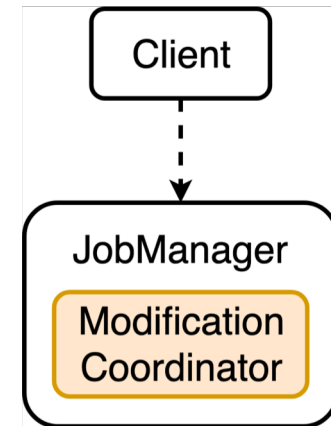
- Dataflow of streams and transformations form DAG



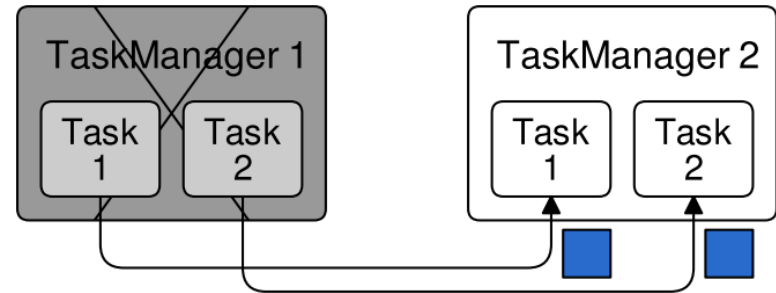
- Checkpoint mechanism ensures **fault tolerance** and allows jobs to be restarted with previous state
- Checkpoint barriers divide stream into independent and **disjoint subsets**



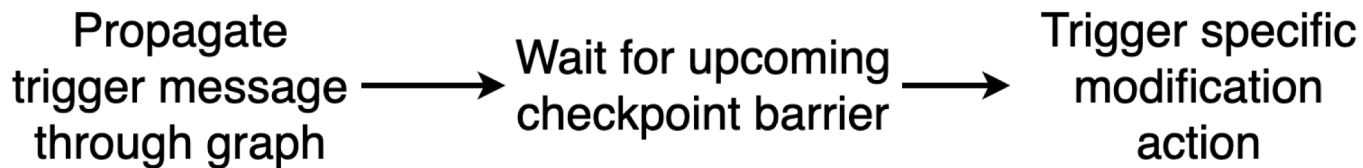
- Integration with checkpoint mechanism
 - Ensures exactly-once processing semantics
 - Reuse of Flink state backend
- ModificationCoordinator
 - checks validity of modification
 - Ingests trigger messages
 - starts and supervises modification execution

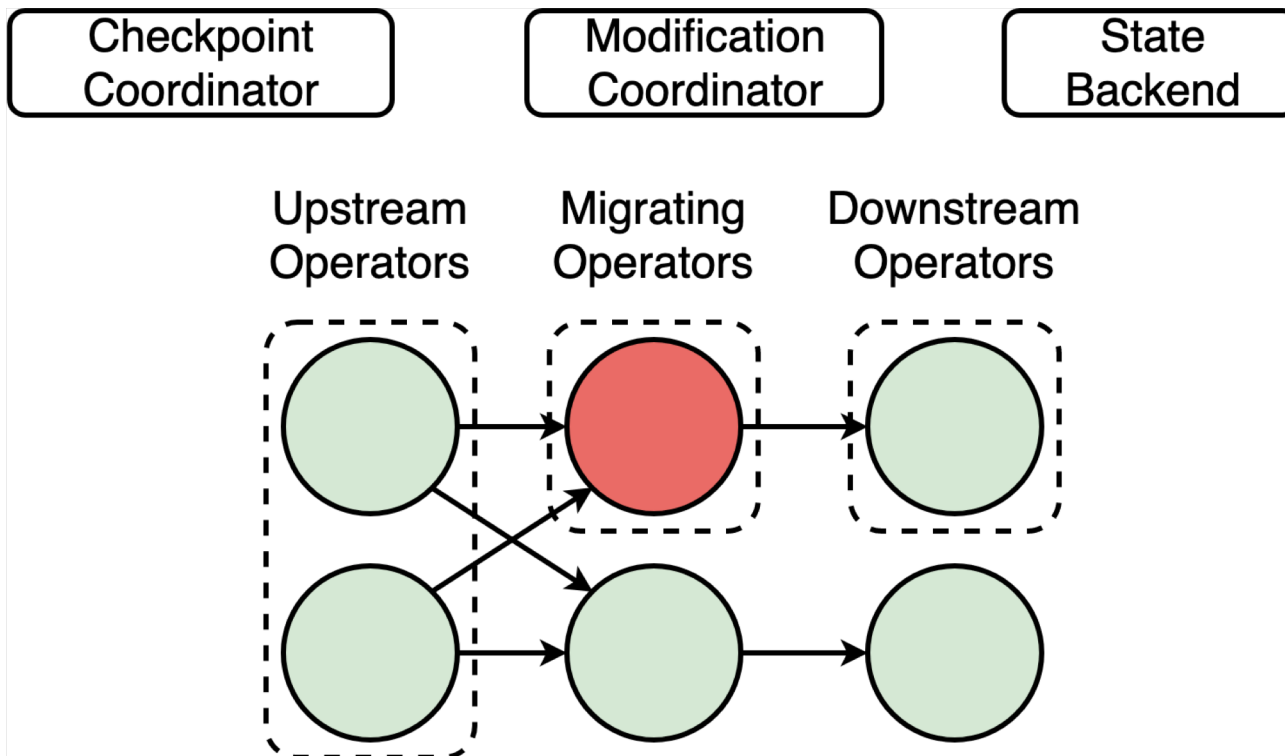


- Motivation: TaskManager metrics indicate an upcoming hardware failure but job should not stop

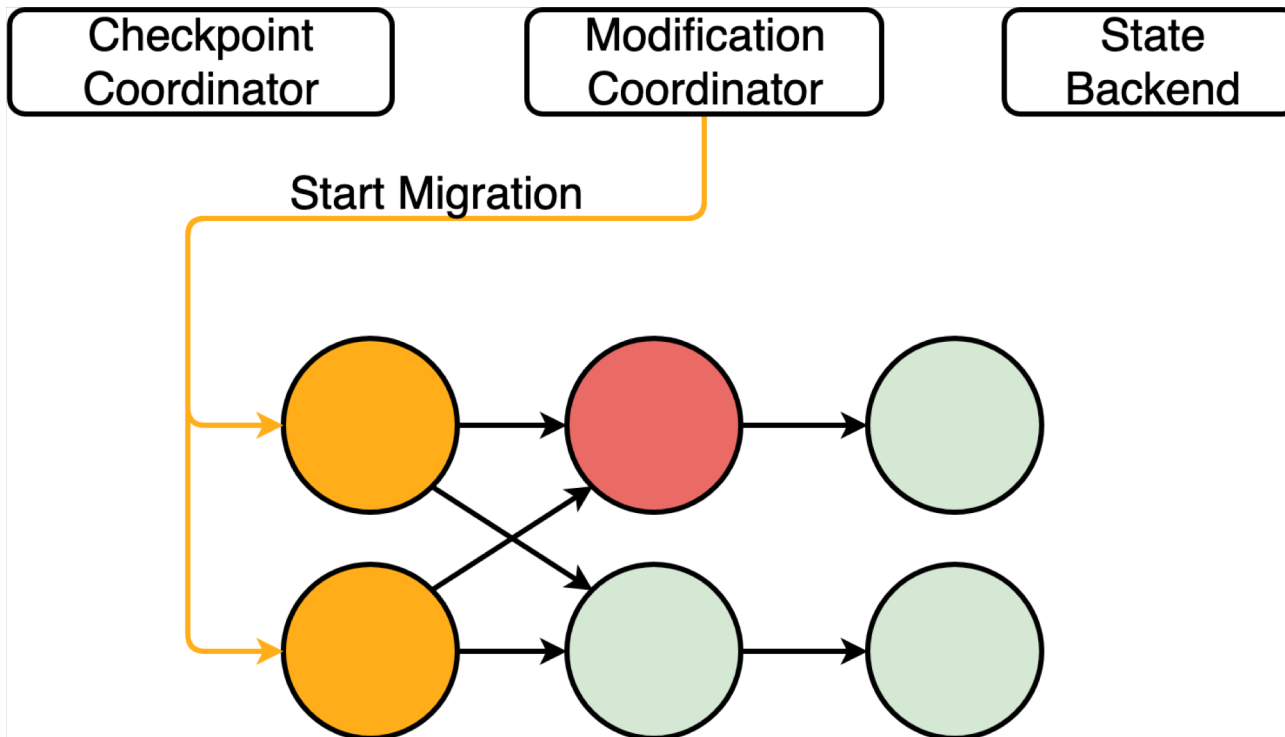


1. Preparing Migration – 2. Waiting – 3. Performing Migration

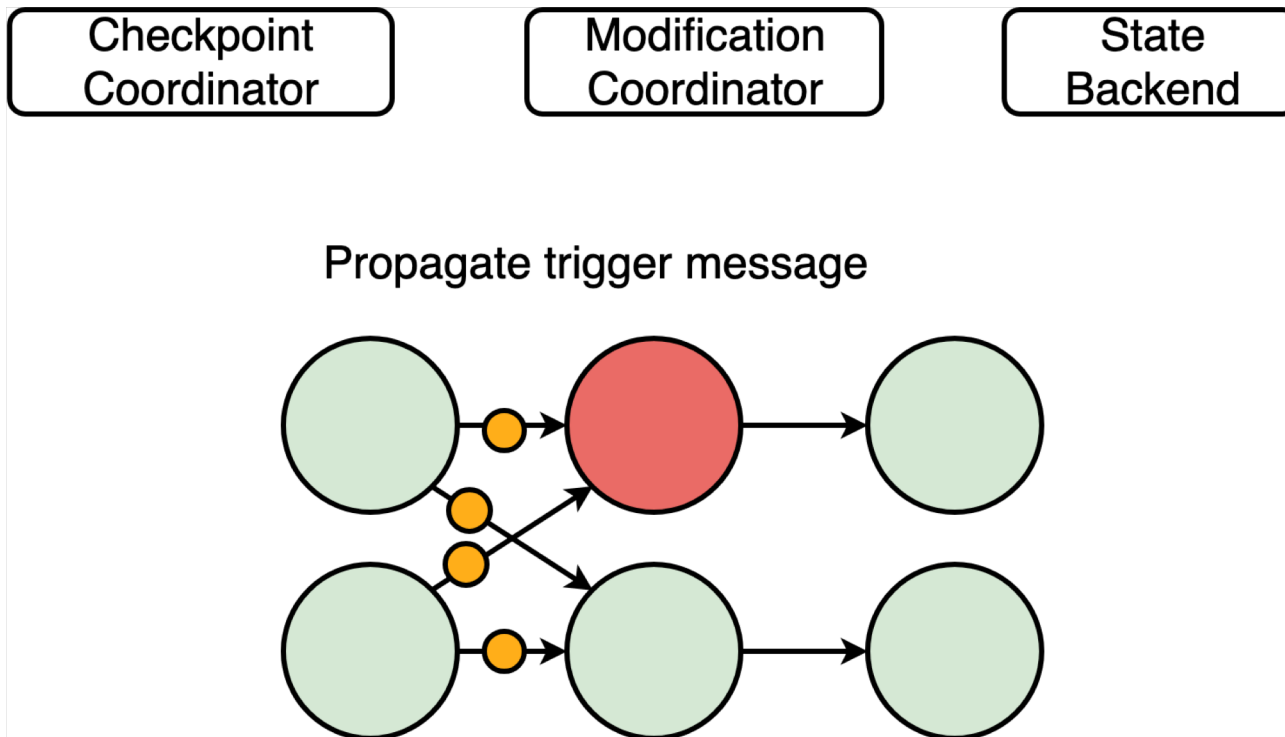




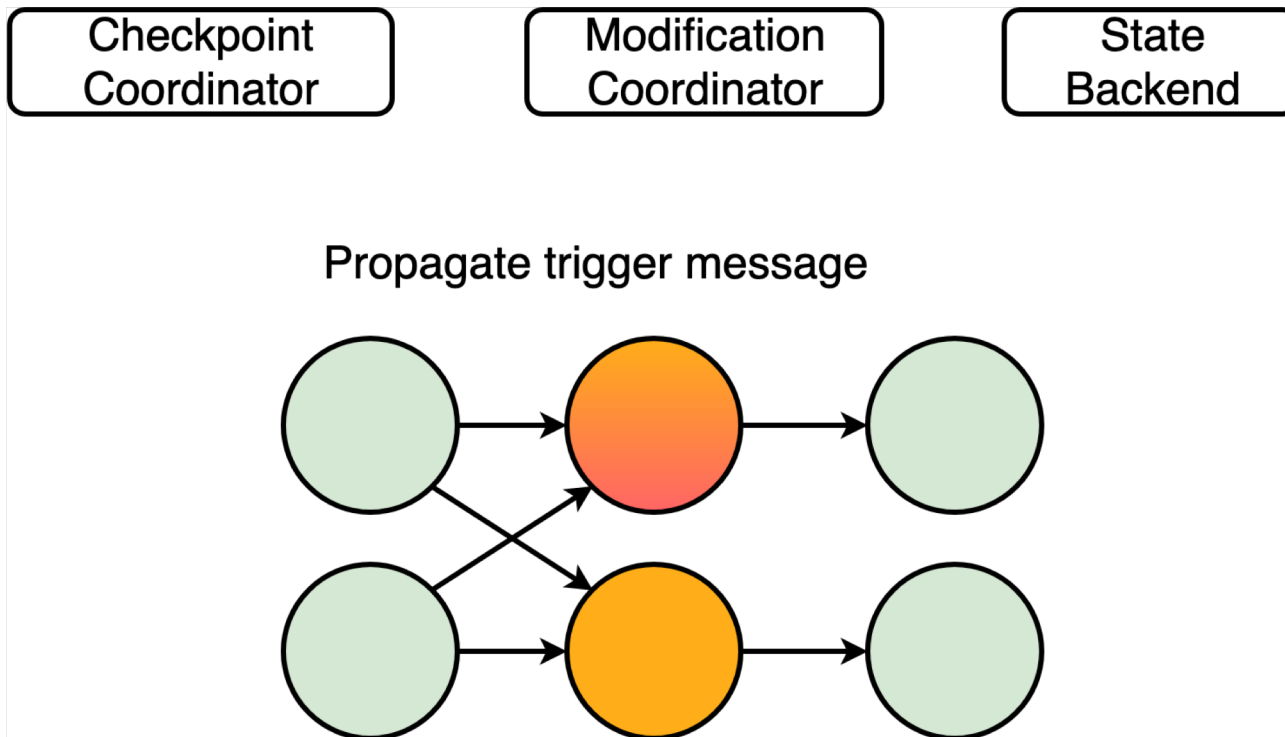
Preparing Migration – Waiting – Performing Migration



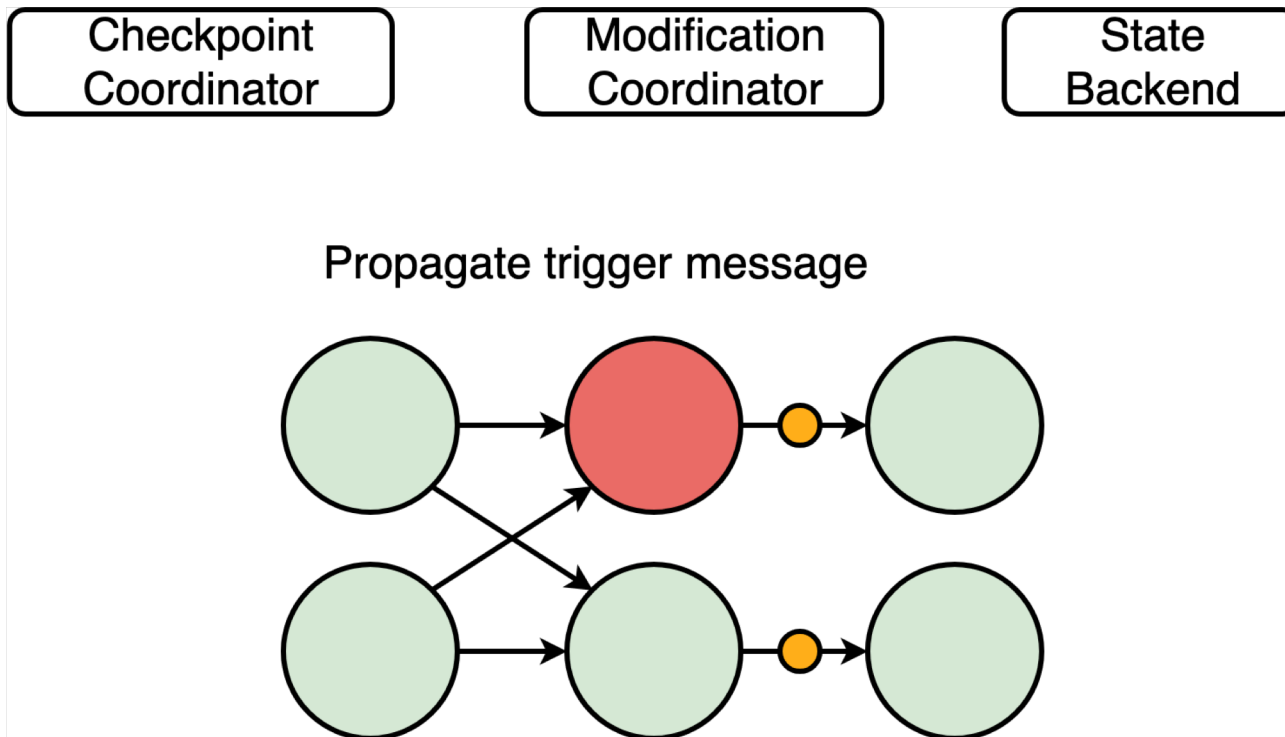
Preparing Migration – Waiting – Performing Migration



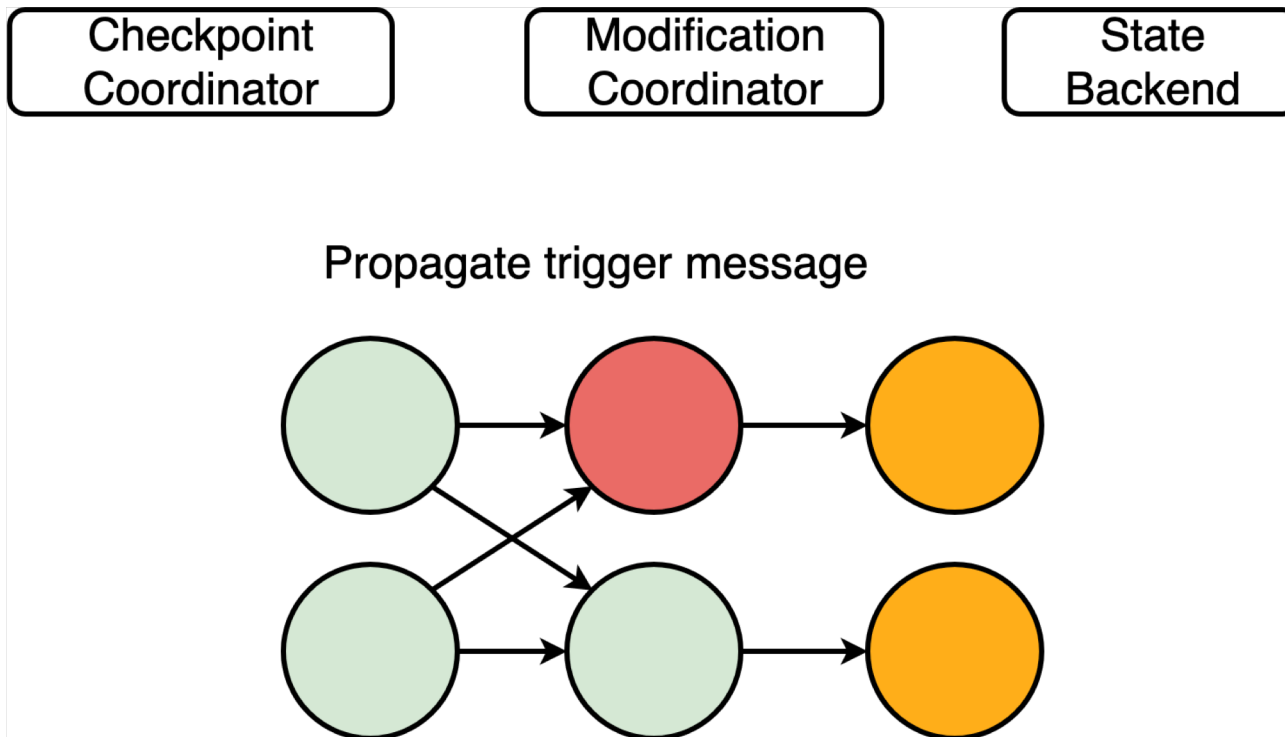
Preparing Migration – Waiting – Performing Migration



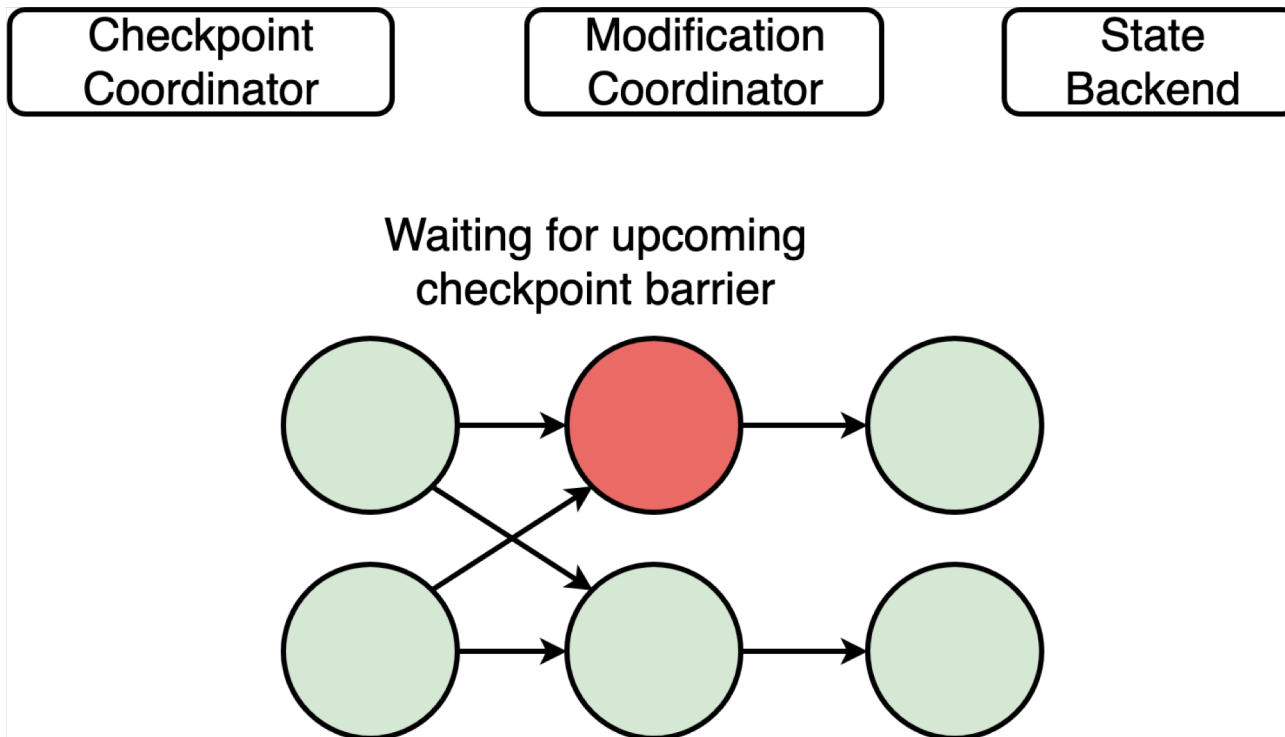
Preparing Migration – Waiting – Performing Migration



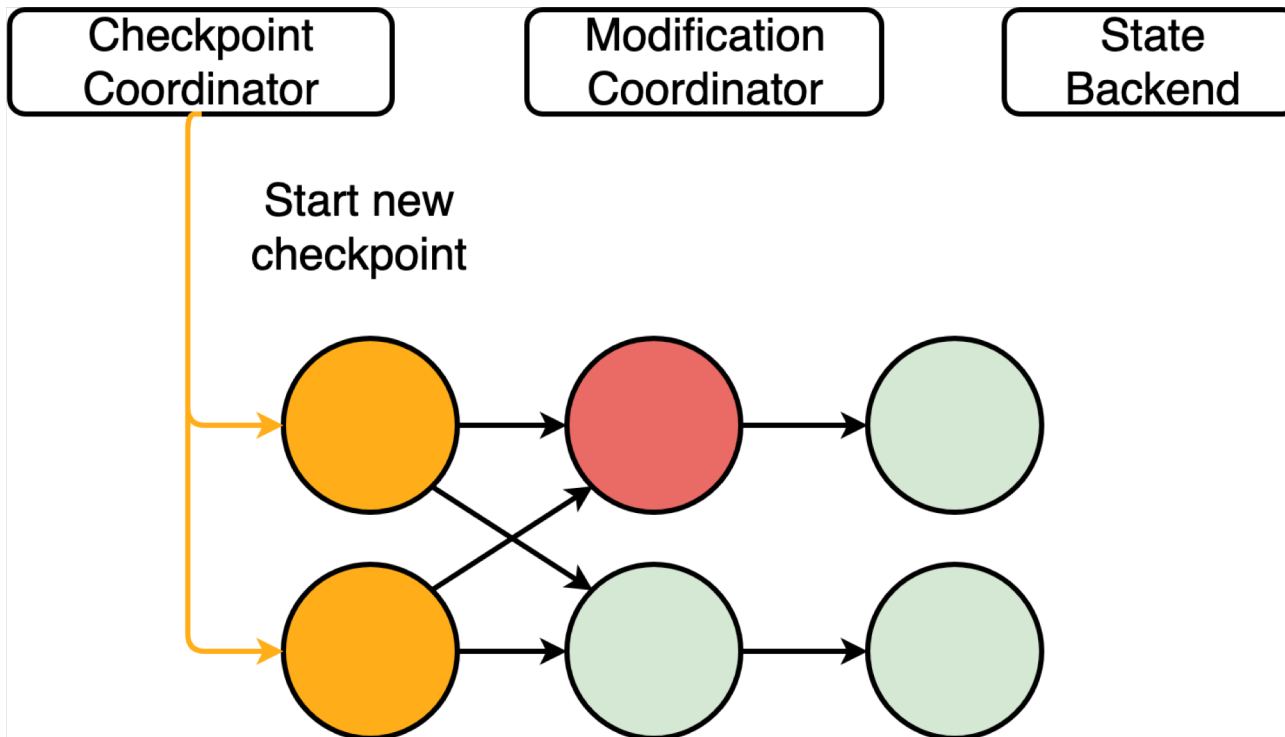
Preparing Migration – Waiting – Performing Migration



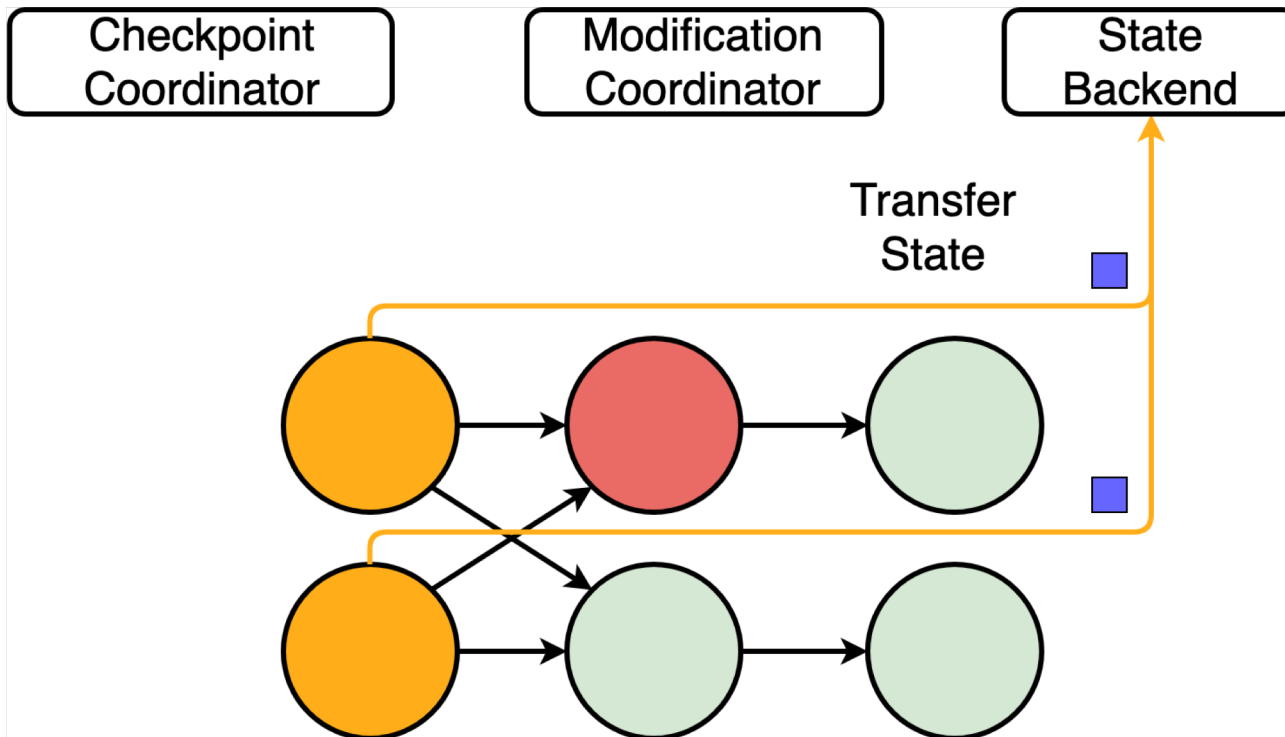
Preparing Migration – Waiting – Performing Migration



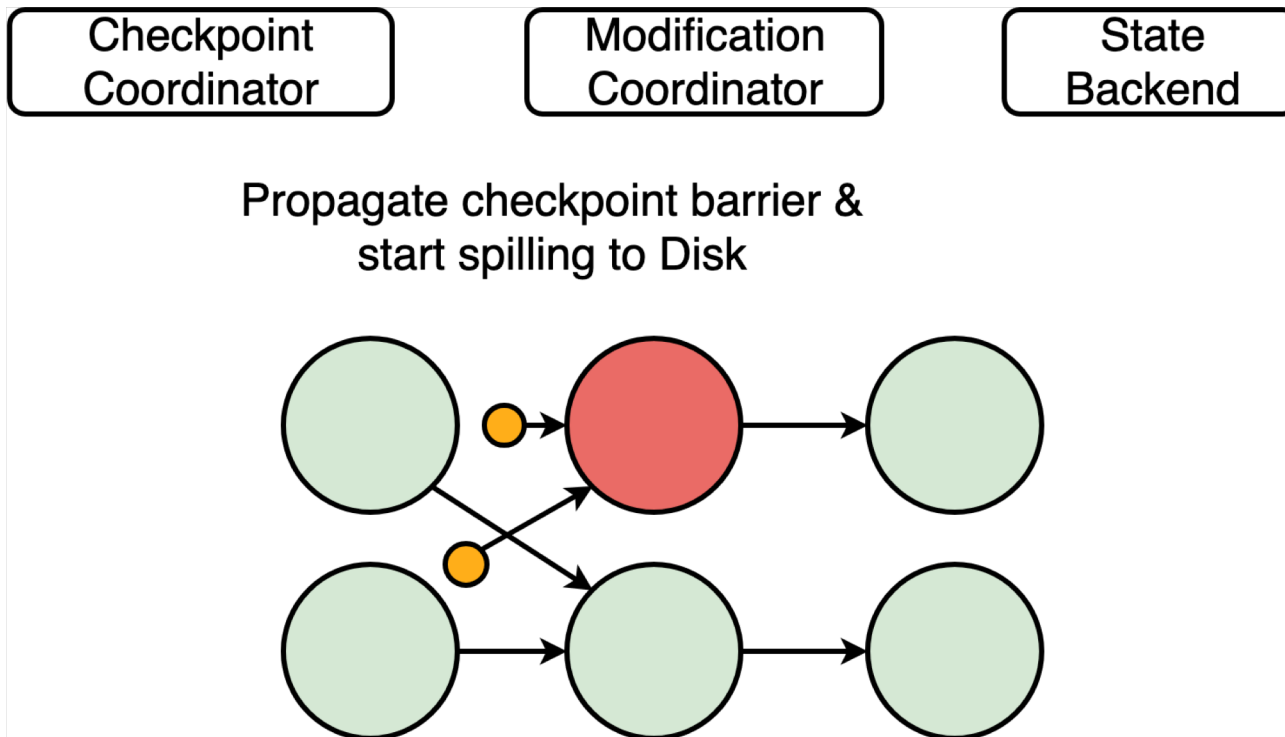
Preparing Migration – Waiting – Performing Migration



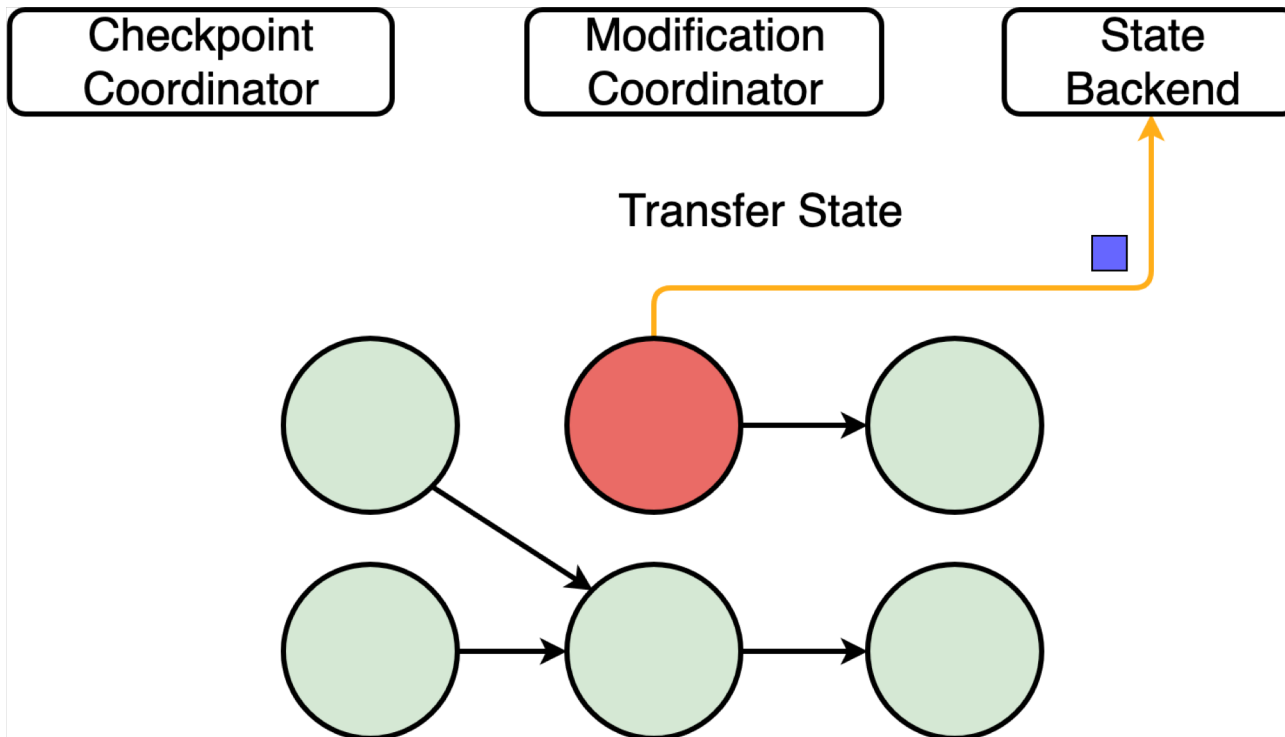
Preparing Migration – Waiting – Performing Migration



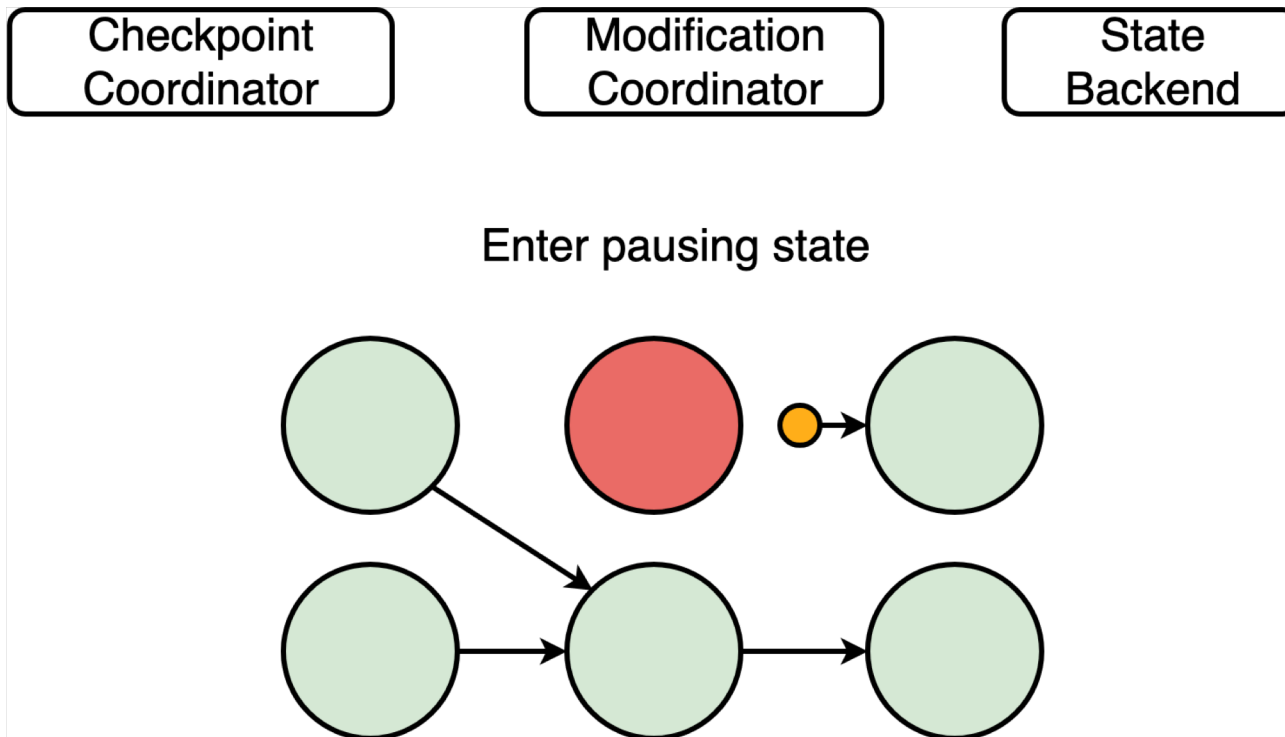
Preparing Migration – Waiting – Performing Migration



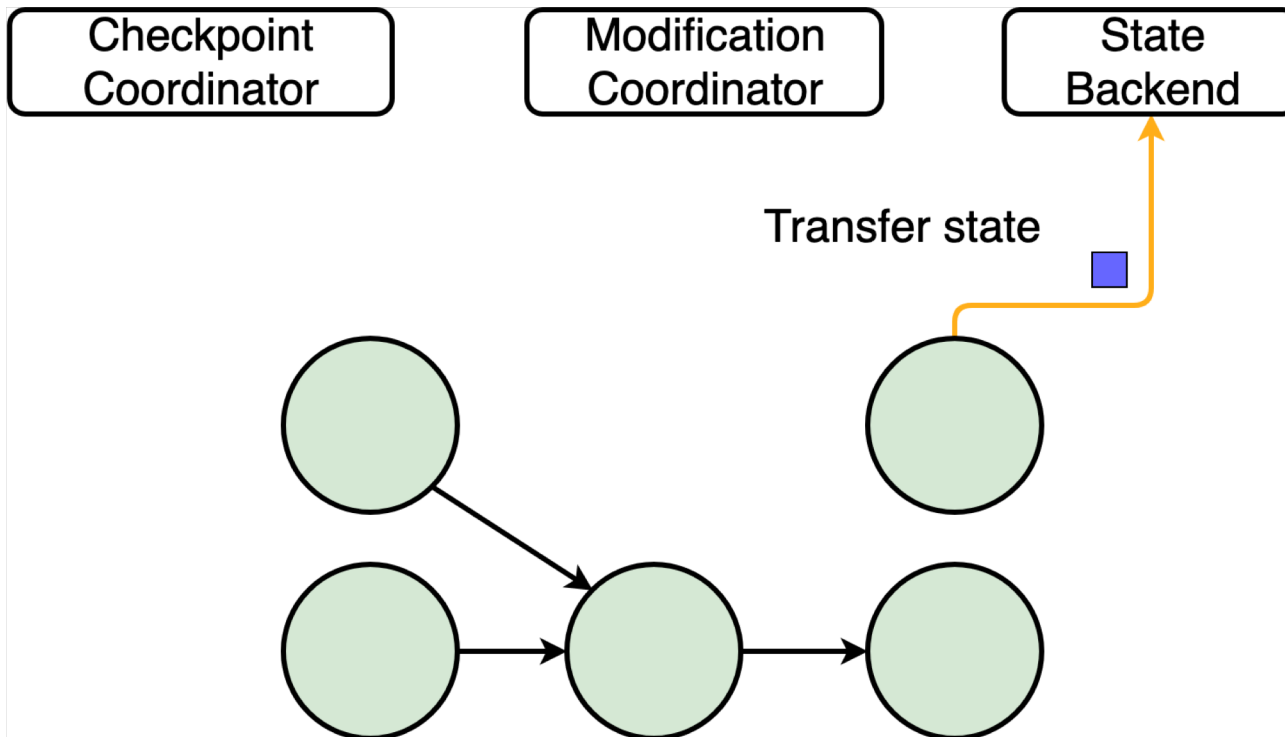
Preparing Migration – Waiting – Performing Migration



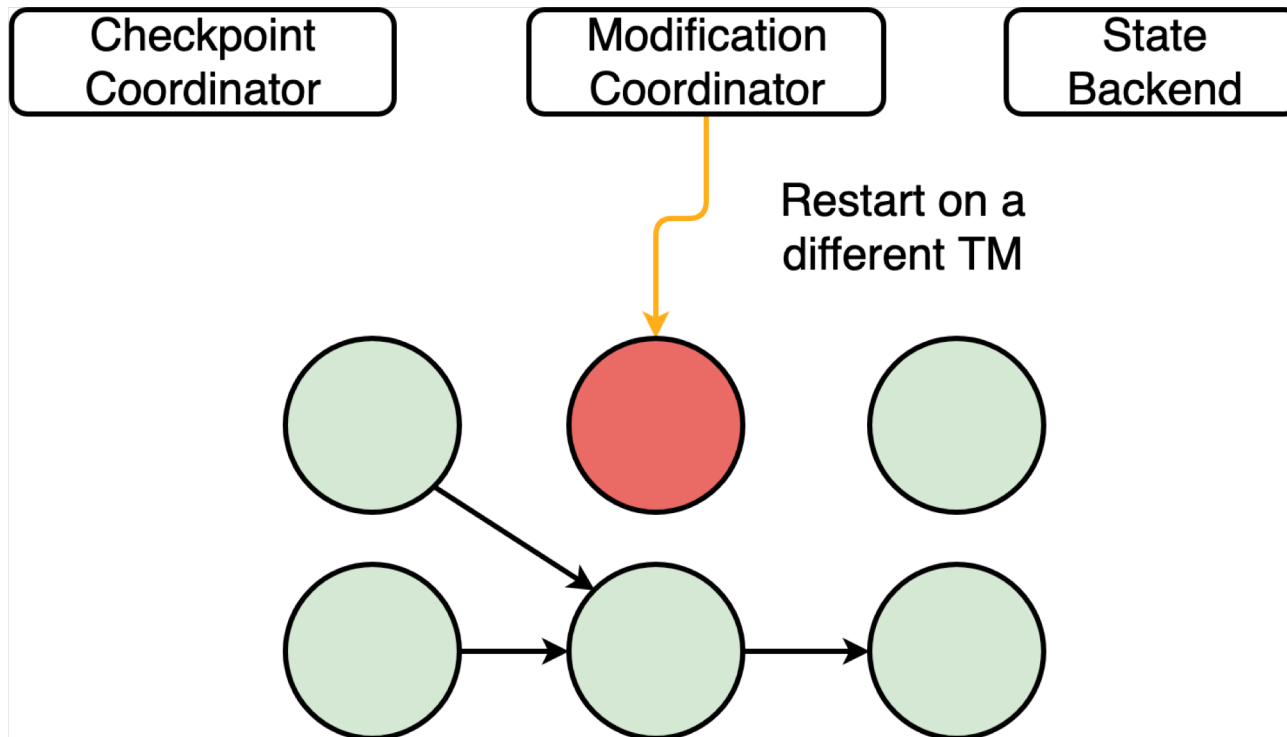
Preparing Migration – Waiting – Performing Migration



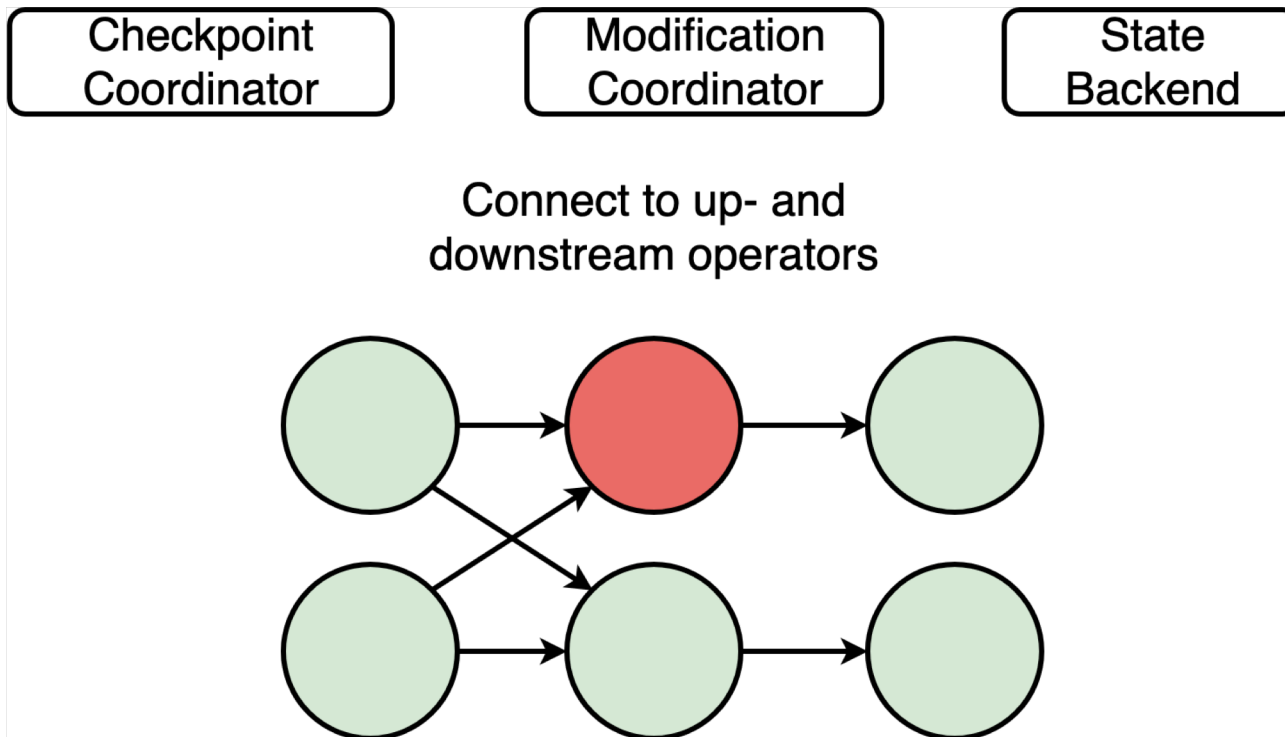
Preparing Migration – Waiting – Performing Migration



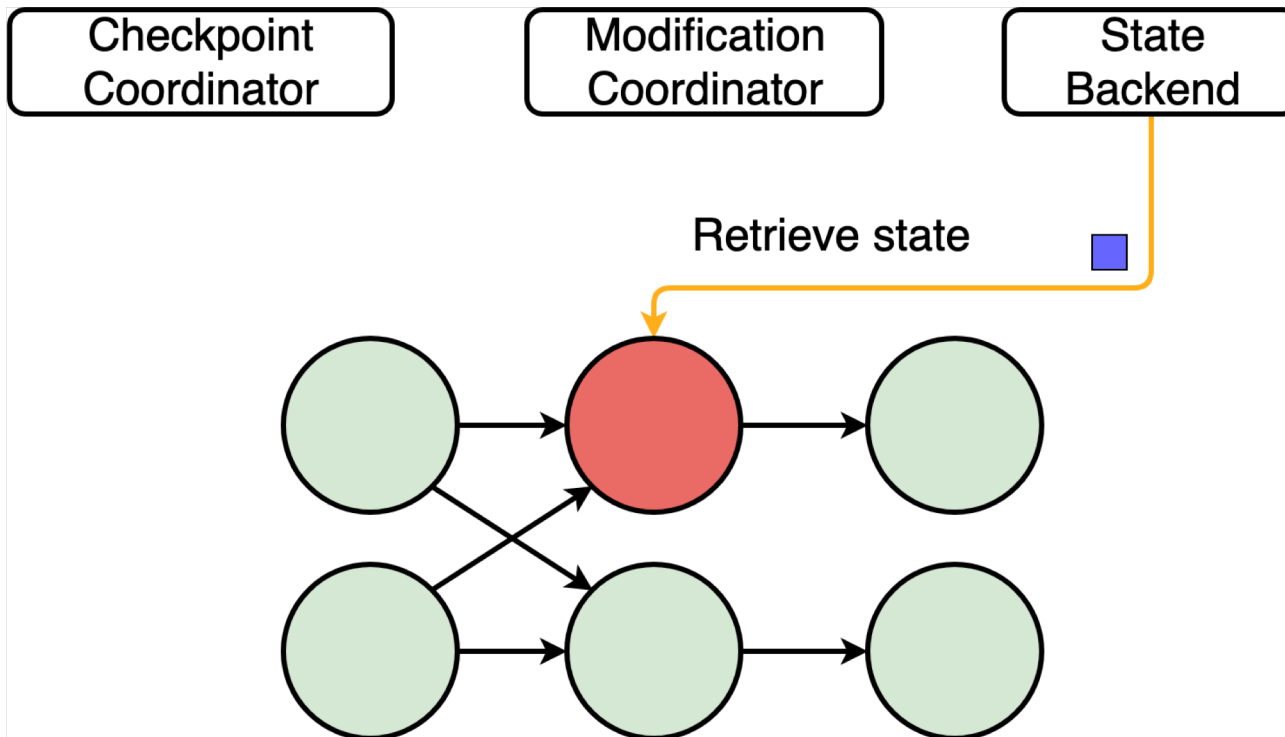
Preparing Migration – Waiting – Performing Migration



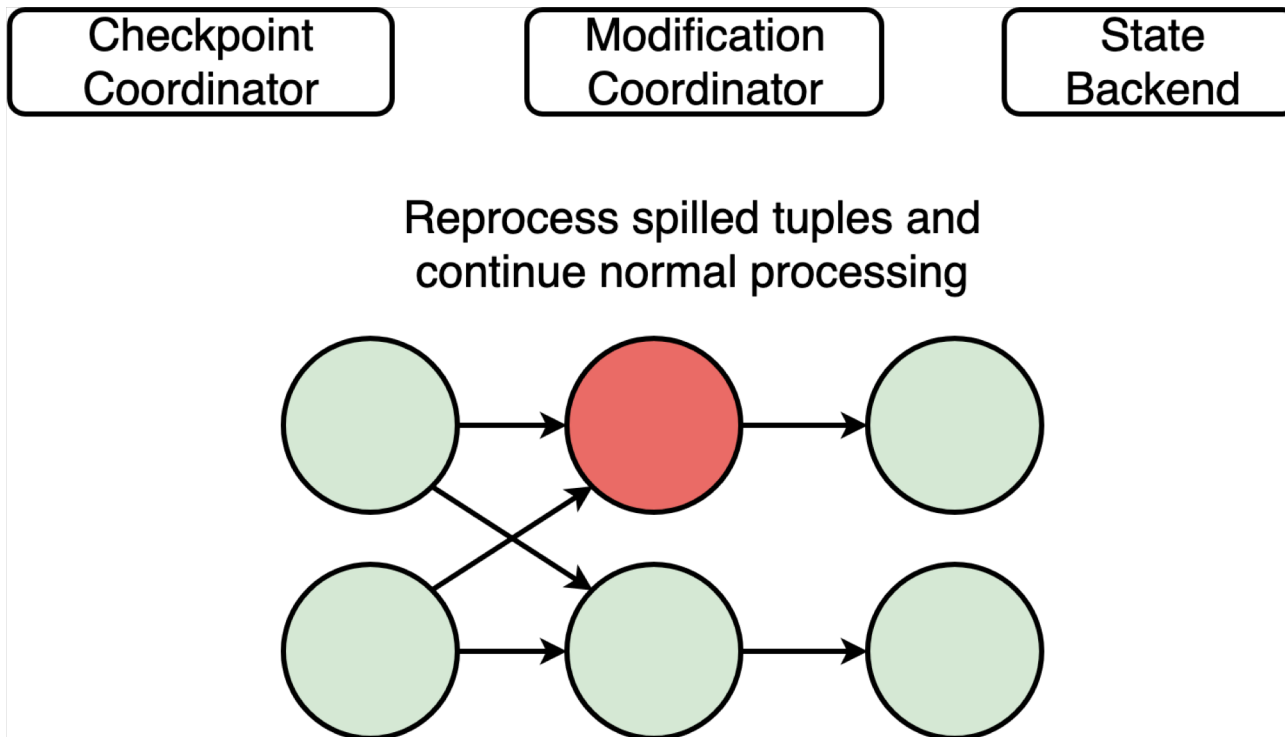
Preparing Migration – Waiting – Performing Migration



Preparing Migration – Waiting – Performing Migration

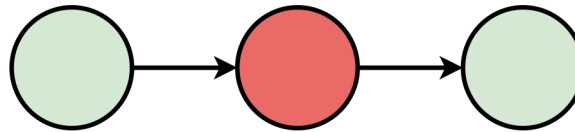


Preparing Migration – Waiting – Performing Migration

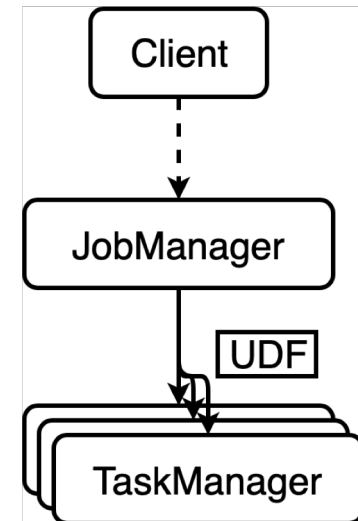


Preparing Migration – Waiting – Performing Migration

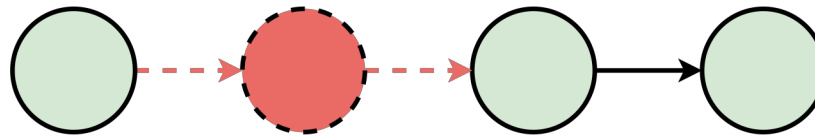
- Motivation: Optimize UDF at runtime, e.g. based on runtime metrics



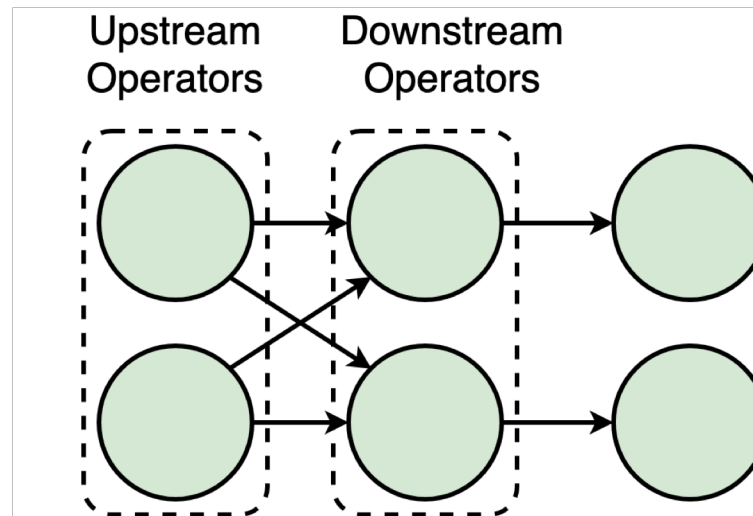
- New UDF distributed via JobManager to all TaskManagers
- Synchronization upcoming checkpoint barrier
 - All operators will replace their UDF



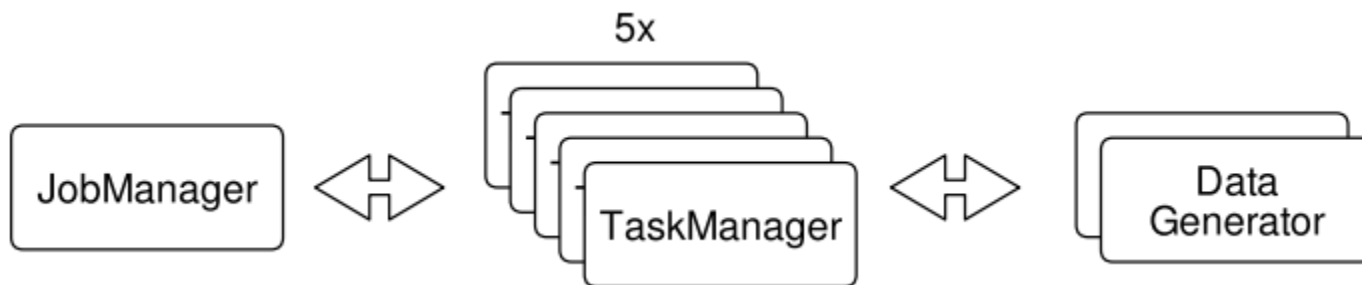
- Motivation: Modify data flow at runtime



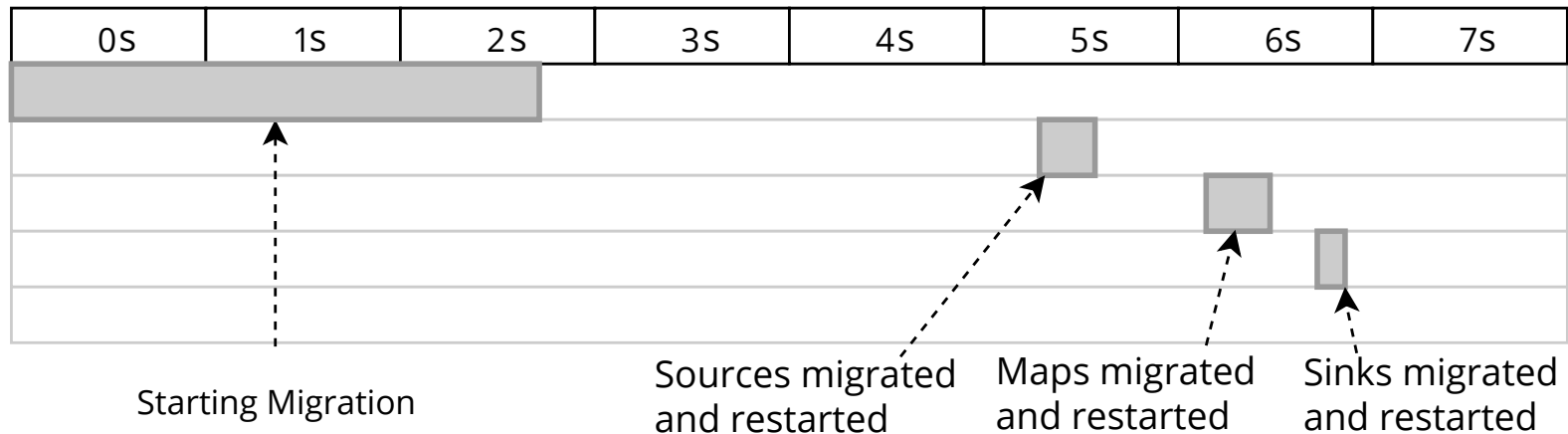
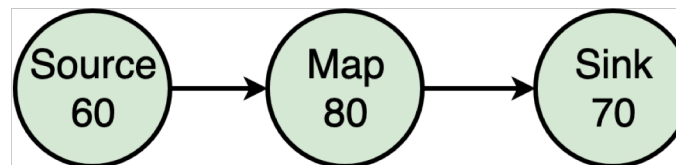
- Separation in up- and downstream operators



- All benchmarks performed on IBM cluster
 - 8 machines with 48 cores and 48 GB of memory each
 - 1 GB/s network between all machines
- Implemented custom data generator with constant load
- Warmup of 10 minutes for each benchmark



- Stateful Map Query covers operator migration with **small state**
 - 1 second checkpointing interval
 - Migrated 34 operator instances in less than 7.1 seconds with 17kB state size



-
- ```

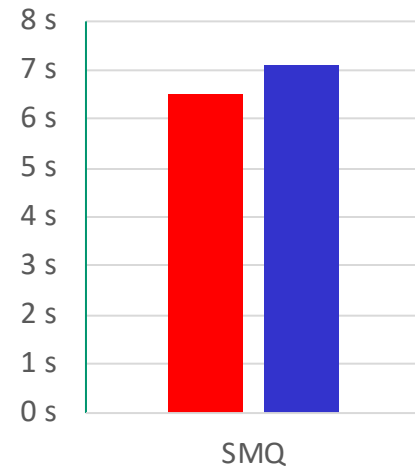
graph LR
 PS((Person Source 40)) --> WS((Window/Sink 40))
 AS((Auction Source 40)) --> WS

```



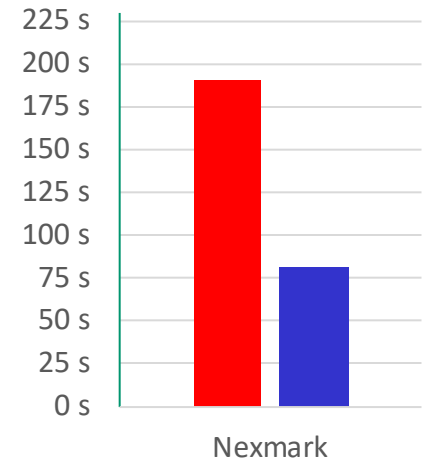
- Big performance improvements when migrating large state by a factor of more than 2.3
- No data loss between restarts with migration
- Most time is not spent on actual migration
  - submitting migration
  - waiting for checkpoint barriers

Stateful Map Query



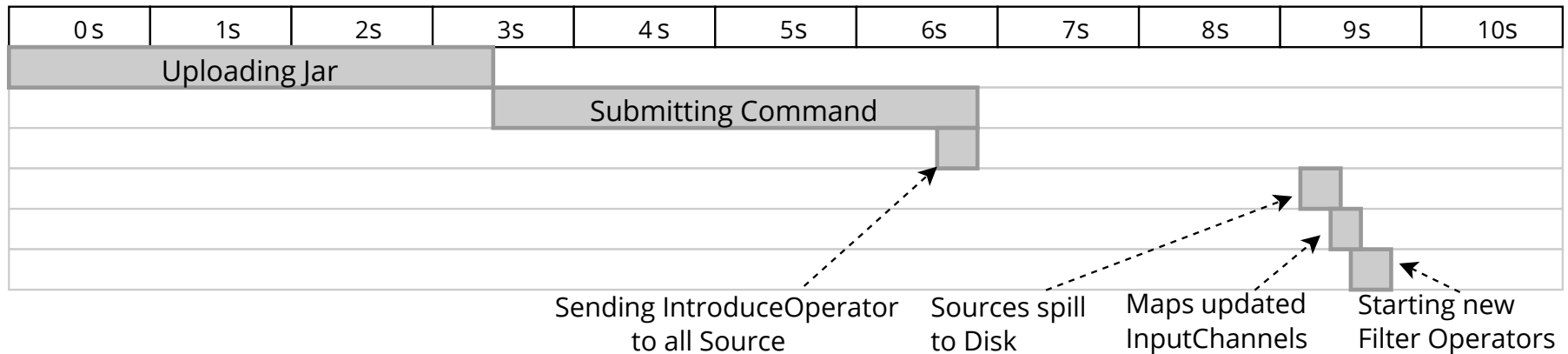
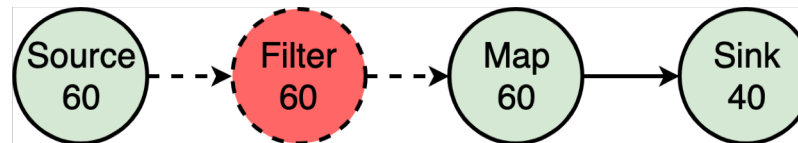
■ Flink Savepoint ■ Migration

Nexmark Query

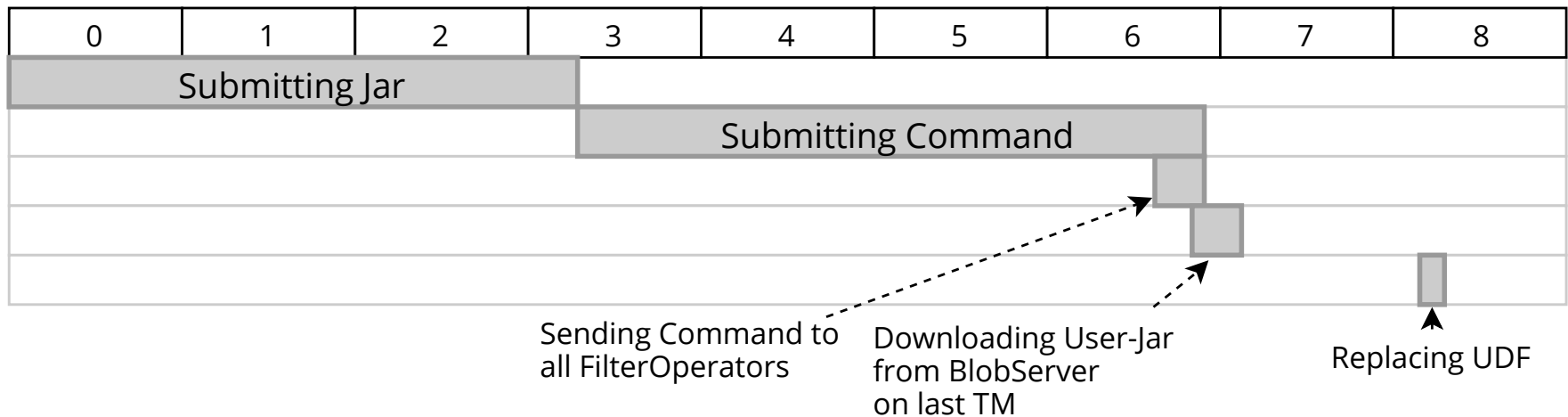
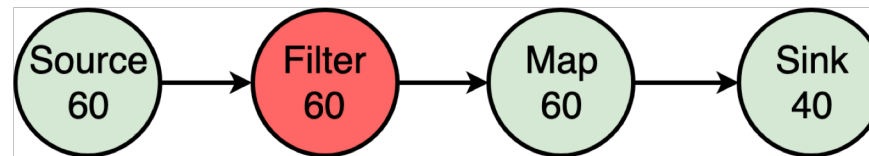


■ Flink Savepoint ■ Migration

- Introduce new filter operator between source and map operators
  - 1 second checkpointing interval



- Change filter condition for an existing filter operator
  - 1 second checkpointing interval



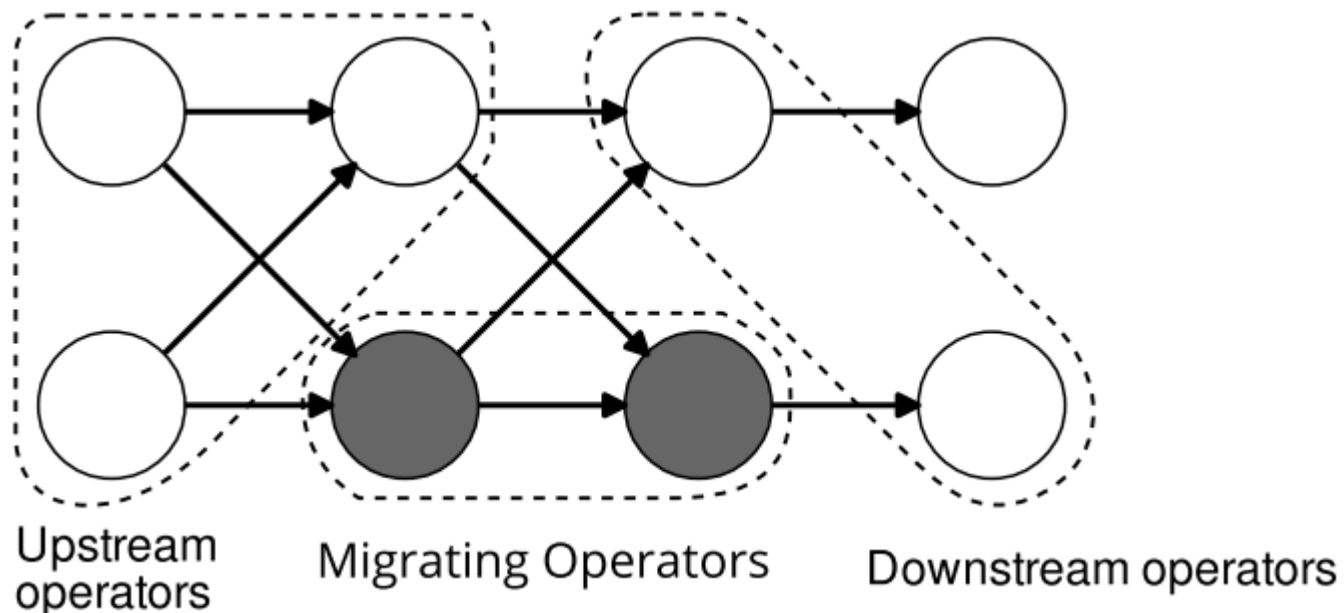
- Mechanism to modify jobs at runtime
  - Modifications independent of actual operator
  - Prevents data loss during restart without persistent message queue
  - Implemented three modification use cases
    - Operator migration
    - Introduction of new operators
    - Replace UDF
  
- Benchmark results
  - Migrating operators is almost as fast as Flink's savepoints
  - For large state 2.3x faster
  - Actual Modifications take less than a few seconds
  
- Solid foundation for dynamic runtime optimizations

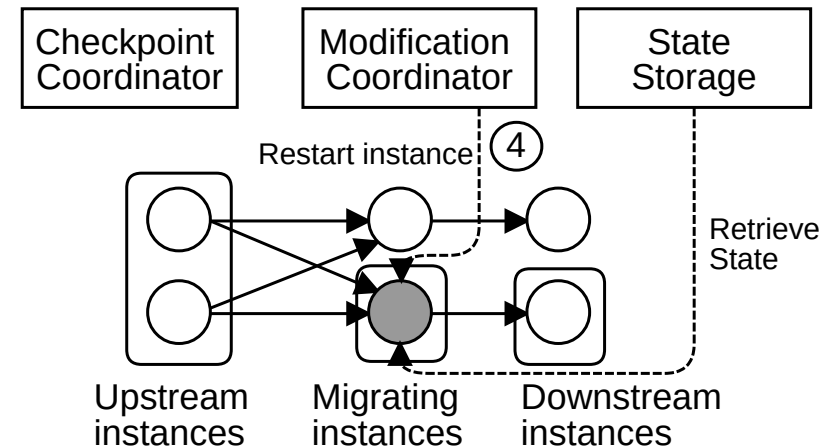
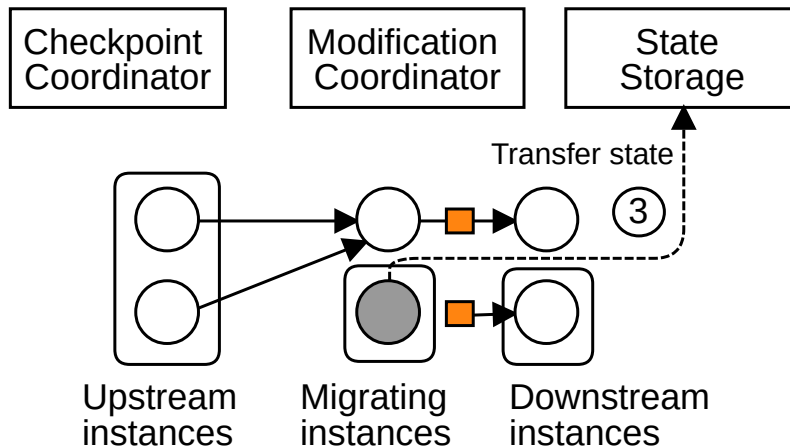
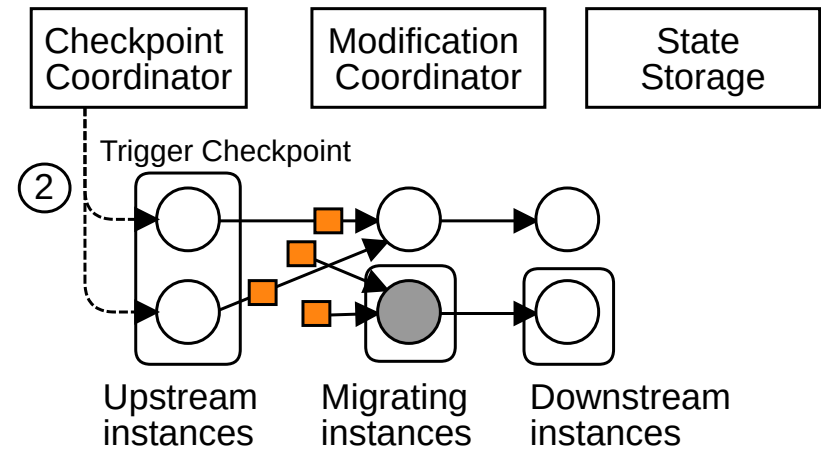
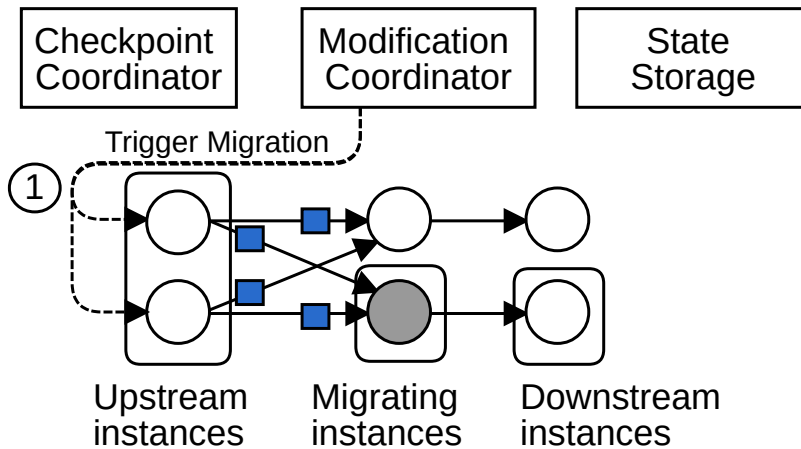


## Backup Slides



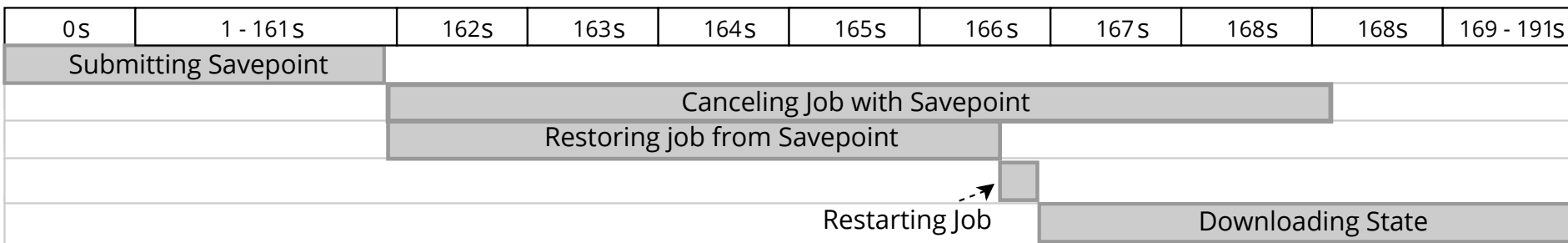
- Real-life jobs require migration of multiple, successive operator instances
- Main insights
  - Upstream operator must still guarantee data safety
  - It does not matter if upstream operator is spilling to disk or migrating itself





■ Trigger Migration Marker   
 ■ Checkpoint Marker   
 → TCP Data Exchange   
 - - -> RPC Exchange

- Big data: [https://commons.wikimedia.org/wiki/File:BigData\\_2267x1146\\_white.png](https://commons.wikimedia.org/wiki/File:BigData_2267x1146_white.png)
- Flink & logos: <https://flink.apache.org/poweredby.html>
- Graph icon made by Smashicons from [www.flaticon.com](http://www.flaticon.com) is licensed by CC 3.0 BY



| Benchmark          | Migrated state size | Waiting for Synchronization | Migration duration | Overall duration | Migrated Operators |
|--------------------|---------------------|-----------------------------|--------------------|------------------|--------------------|
| Stateful Map Query | 17 kB               | 2 s                         | 2.4 s              | 7.1 s            | 34                 |
| Nexmark            | 2,7 GB              | 54.9 s                      | 23.7 s             | 81.8 s           | 24                 |

- Big performance improvements when migrating large state by a factor of more than 2.3
- No data loss between restarts with migration
- Most time is not spent on actual migration
  - submitting migration
  - waiting for checkpoint barriers

