# Unified
# Batch & Stream Processing
# in Apache Flink

Apache Flink Meetup Berlin #6
*April 29, 2015*

**Ufuk Celebi**
uce@apache.org

# Batch vs. Stream Processing

**Batch** Processing

**Stream** Processing

| **High** Latency |
|---|

| **Low** Latency |
|---|

| **Batch Processors** |
|---|

| **Stream Processors** |
|---|

| **Static Files** |
|---|

| **Event Streams** |
|---|

What's the *difference* between a **Batch** and **Streaming** *Runtime*?

# Batch Processing

HDFS

Read *complete* file →

**map**

**reduce**

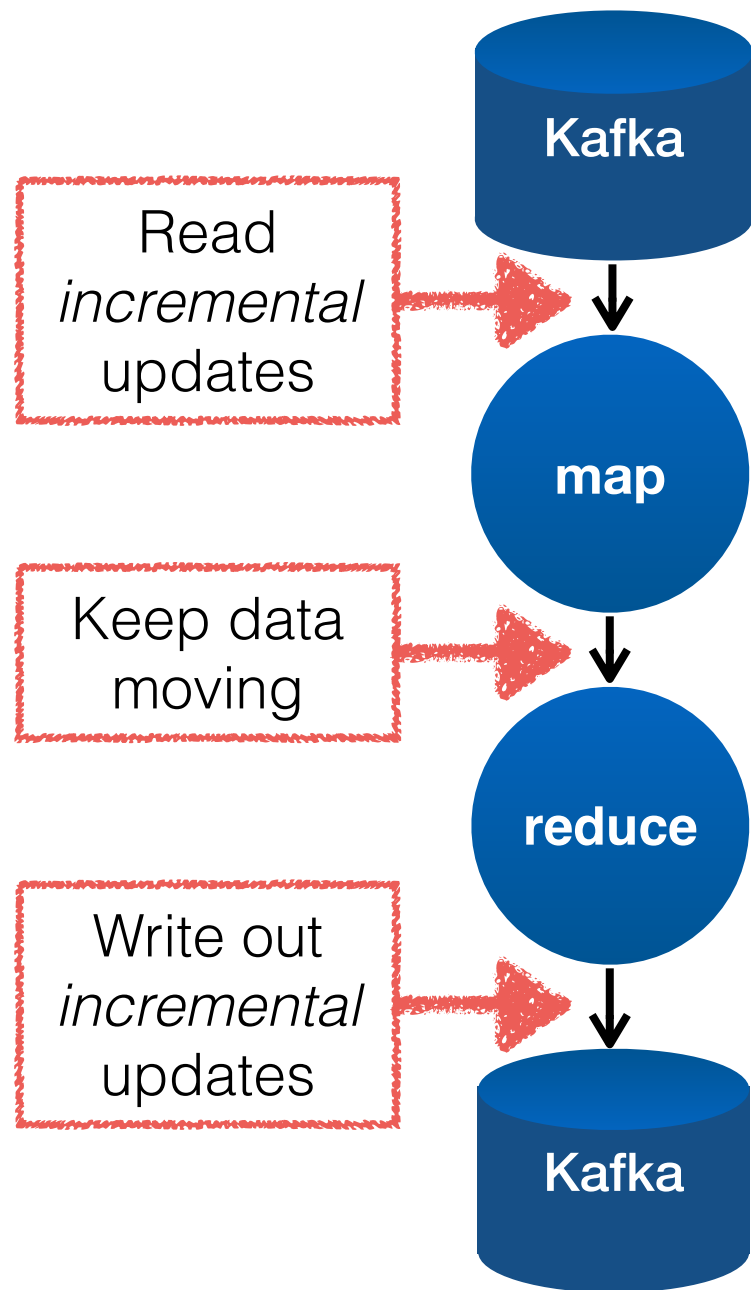Write out *complete* result →

HDFS

O Romeo, Romeo,
wherefore art thou Romeo?

Romeo, 1
Romeo, 1
wherefore, 1
art, 1
thou, 1
Romeo, 1

← Write out intermediate data set

Romeo, 3
wherefore, 1
art, 1
thou, 1

# Stream Processing

**Kafka**

Read *incremental* updates →

**map**

Keep data moving →

**reduce**

Write out *incremental* updates →

**Kafka**

| O | Romeo | Romeo | |
|---|---|---|---|
| wherefore | art | thou | Romeo |

Romeo, 1
Romeo, 1
wherefore, 1
art, 1
thou, 1
Romeo, 1

Romeo, 3
wherefore, 1
art, 1
thou, 1

# Batch ~~vs.~~ *and* Stream
# Processing

# Do we *really* need two **separate systems** for this?

# Apache Flink

| **Batch** Processing | **Stream** Processing |
| --- | --- |
| **DataSet** (Java/Scala) | **DataStream** (Java/Scala) |

**Flink Runtime**

# Flink Runtime

**Operator**

**Intermediate Result**

# Operators

**User**

Map  Reduce  Join  CoCroup  ...

**Internal**

Sort  Hash  ...

Operators represent computations over data.

# Operators

# Intermediate Results



Logical handle to data produced by operators.

# Intermediate Results



Logical handle to data **produced by operators**.

# Map-Reduce Example

```
DataSet<Tuple2<String, Integer>>
   counts;

counts = input
   .flatMap(new LineSplitter())
   .groupBy(0)
   .sum(1);
```

**Map**

**Intermediate Result**

**Reduce**

**Logical handle to data** produced by operators.

# Result Characteristics

**Pipelined** vs. **Blocking**

**Ephemeral** vs. **Checkpointed**

# Result Characteristics

## Pipelined vs. Blocking

Ephemeral vs. Checkpointed

**How and when to do data exchange?**

# Blocking Results

1101

0101

0100

**Map** → **Blocking Result**

# Blocking Results

0101

0100

1101 → **Blocking Result**

# Blocking Results

# Blocking Results

0100

0101

1101
**Blocking
Result**

# Blocking Results



0100

Map

1101

0101

Blocking Result
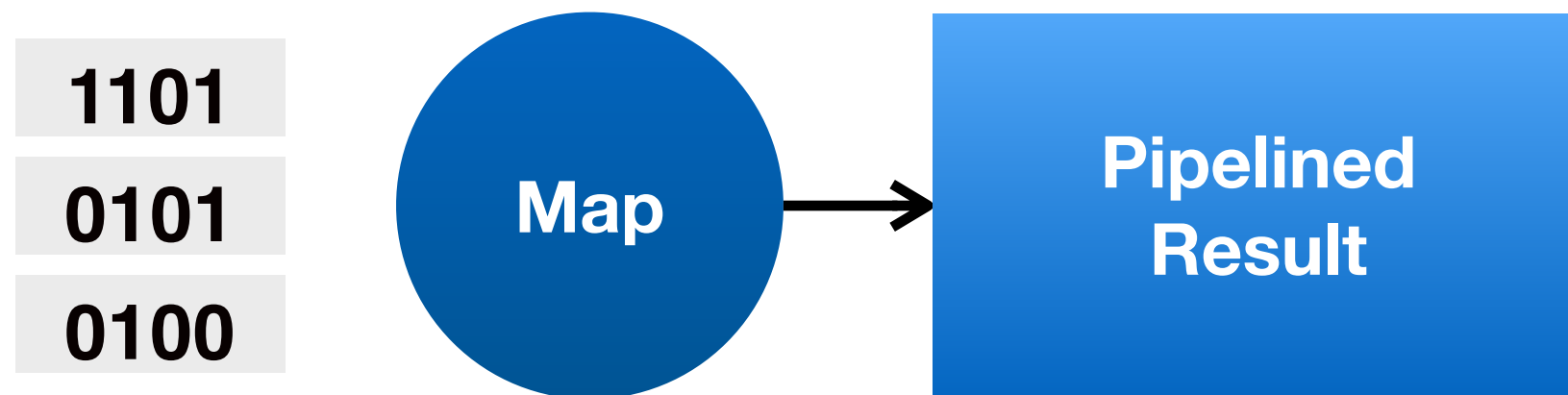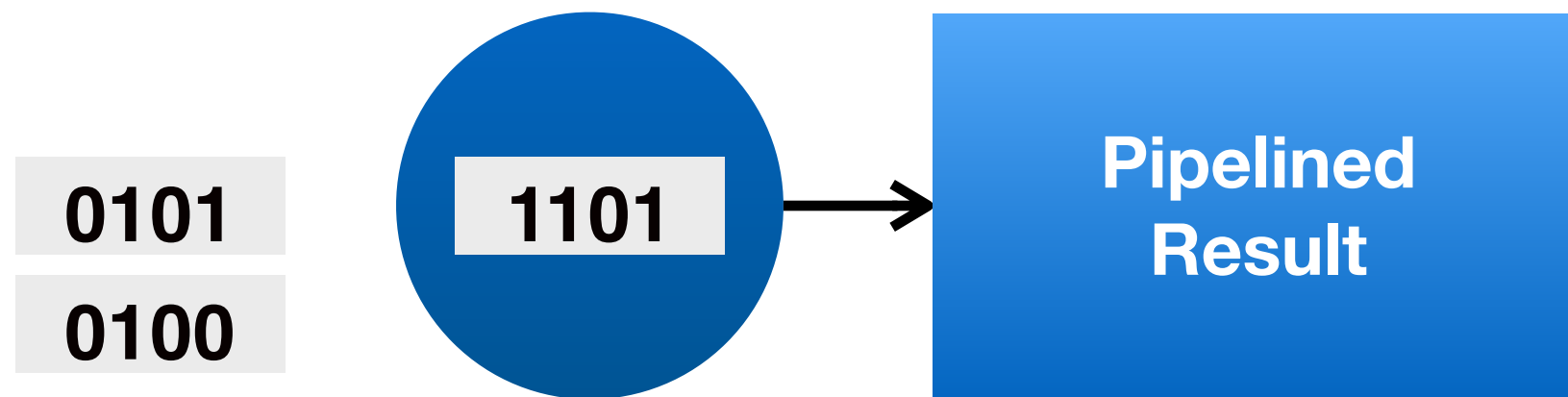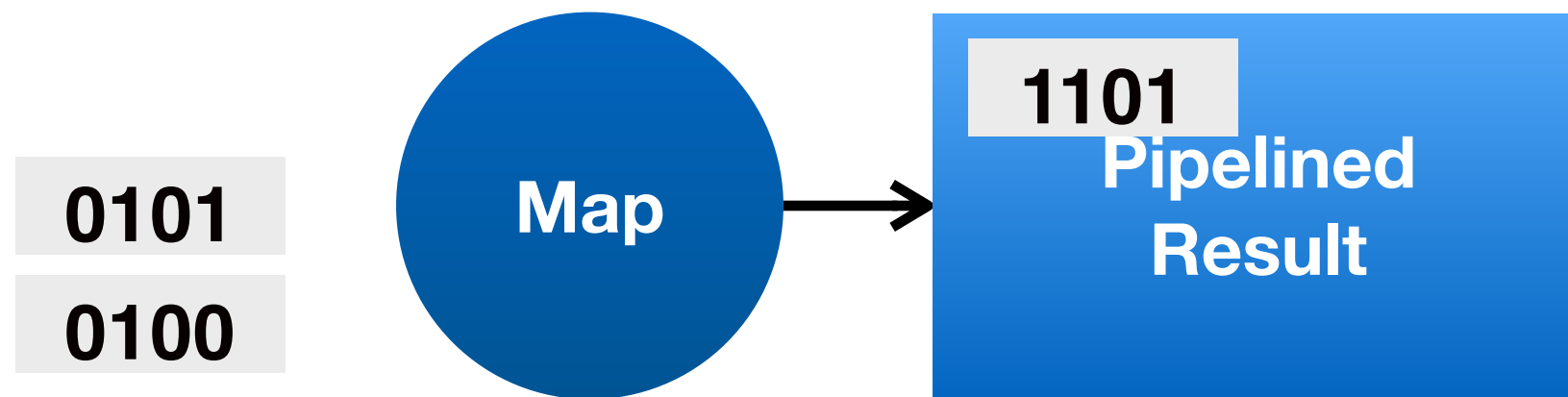
# Blocking Results

# Blocking Results

# Blocking Results



Map → Blocking Result

1101
0101
0100

Reduce

17

# Blocking Results

# Blocking Results

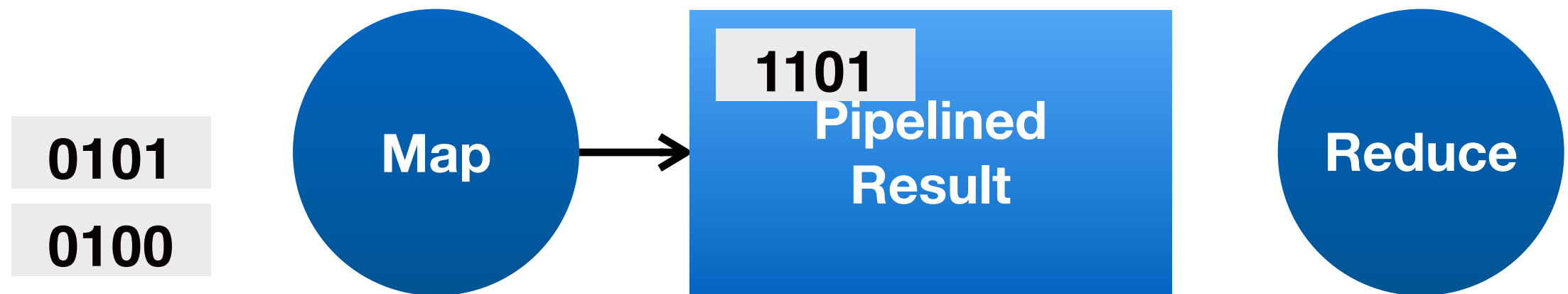# Blocking Results

# Blocking Results

# Pipelined Results

1101

0101

0100

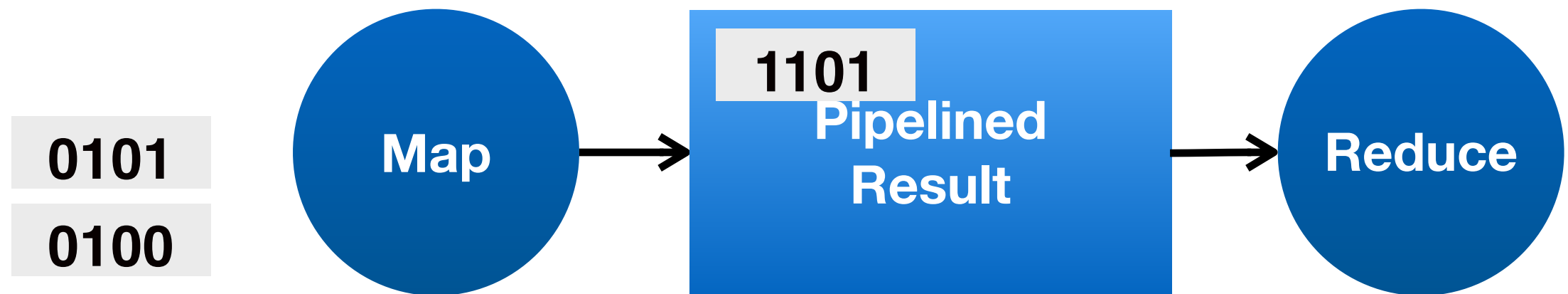**Map** → **Pipelined Result**

# Pipelined Results

0101
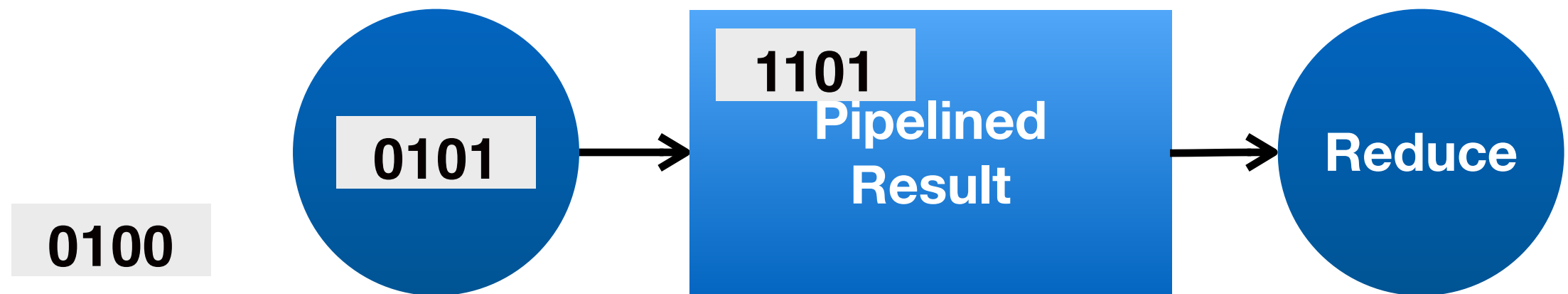
0100

1101

**Pipelined Result**

# Pipelined Results
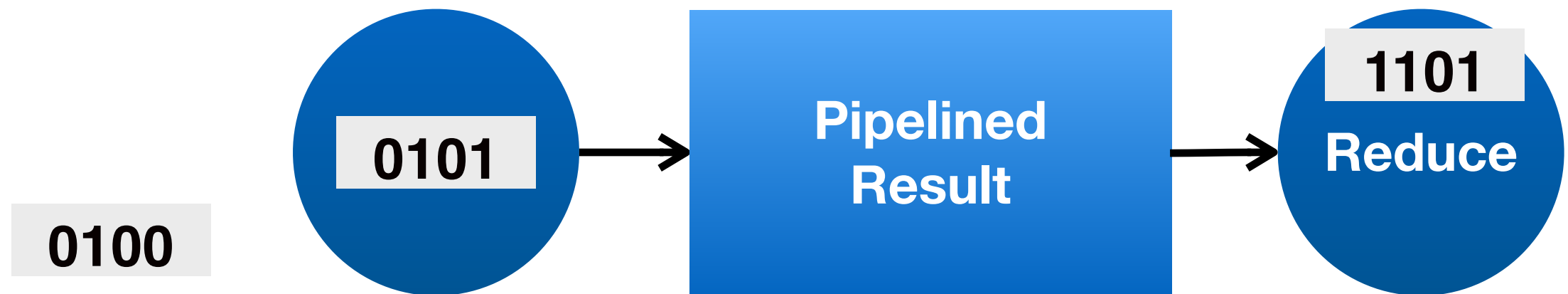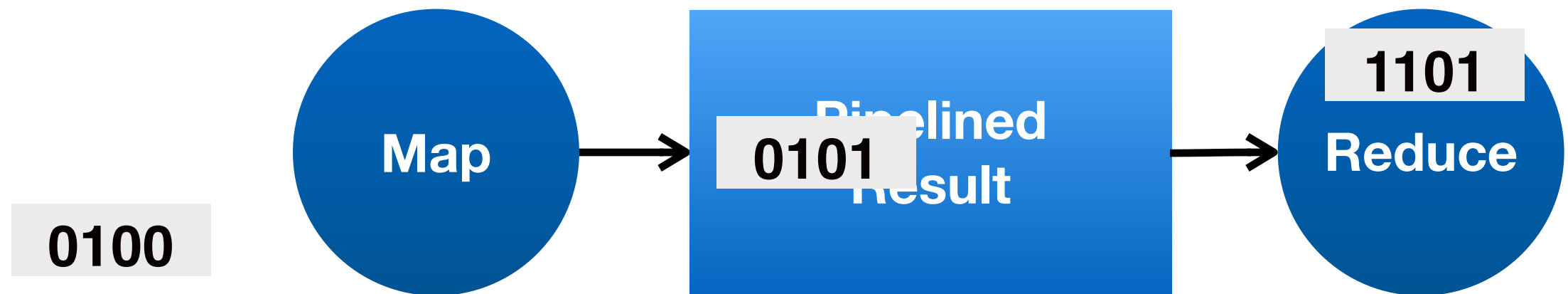
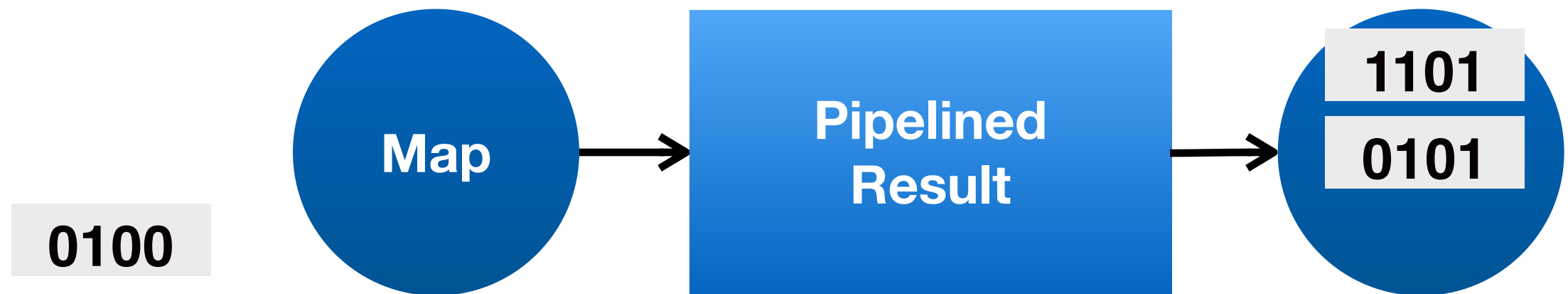# Pipelined Results

# Pipelined Results

# Pipelined Results

# Pipelined Results

# Pipelined Results

**0100**

**Map** → **Pipelined Result** **0101** → **Reduce** **1101**

18

# Pipelined Results

**0100**

**Map** → **Pipelined Result** → **1101** **0101**

# Pipelined Results



0100 → Pipelined Result → 1101 / 0101

# Pipelined Results

# Pipelined Results

Map → Pipelined Result → 1101 / 0101 / 0100

# Result Characteristics

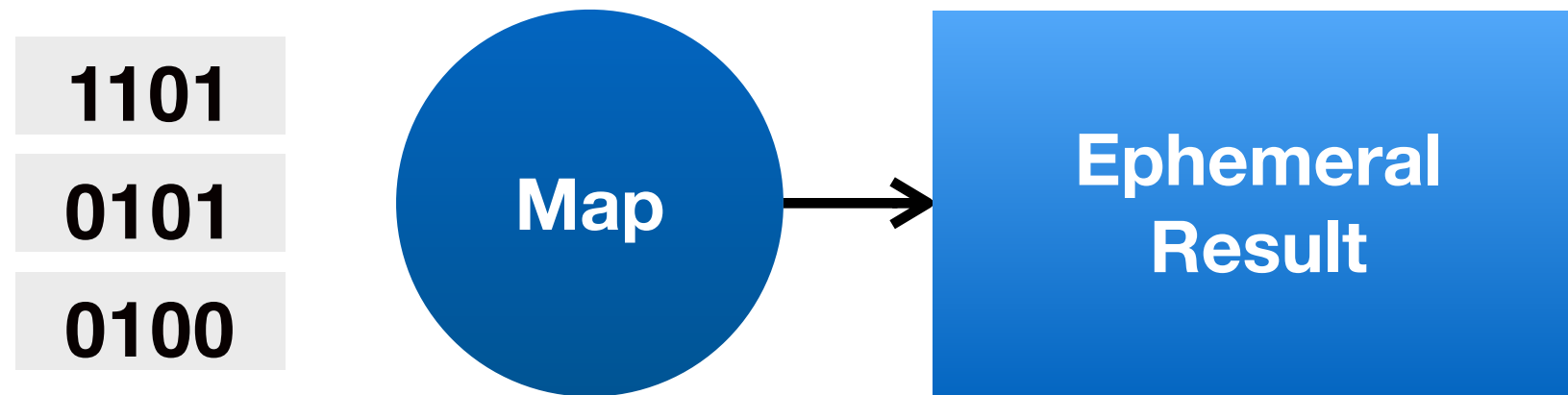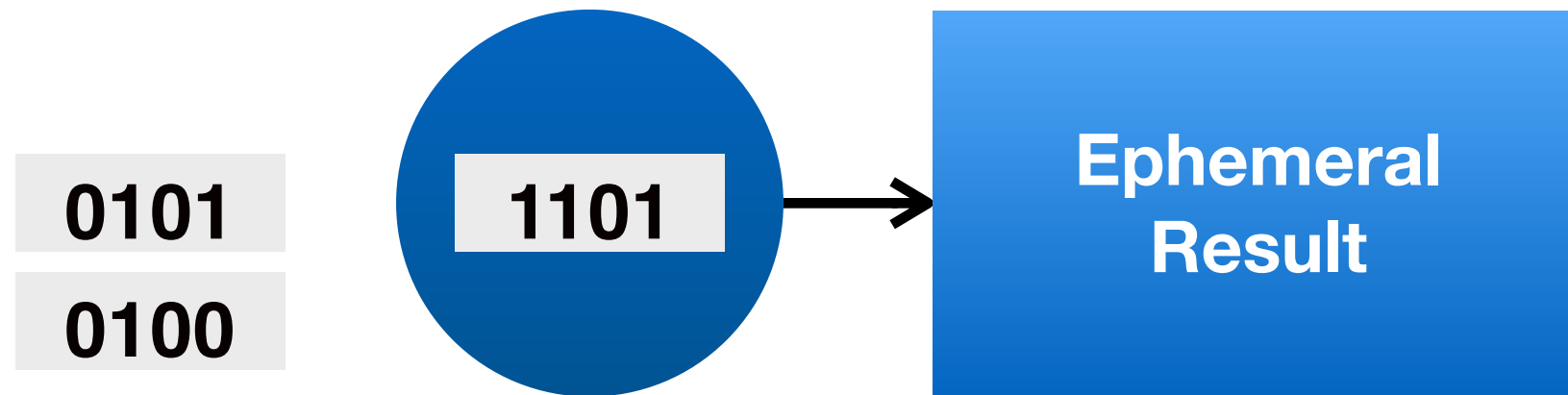|  | Ephemeral | Checkpointed |
|---|---|---|
| **Pipelined** | Low-latency | Low-latency |
| **Blocking** | Easy to reason about resource consumption | Easy to reason about resource consumption |

# Result Characteristics

Pipelined vs. Blocking

**Ephemeral** vs. **Checkpointed**
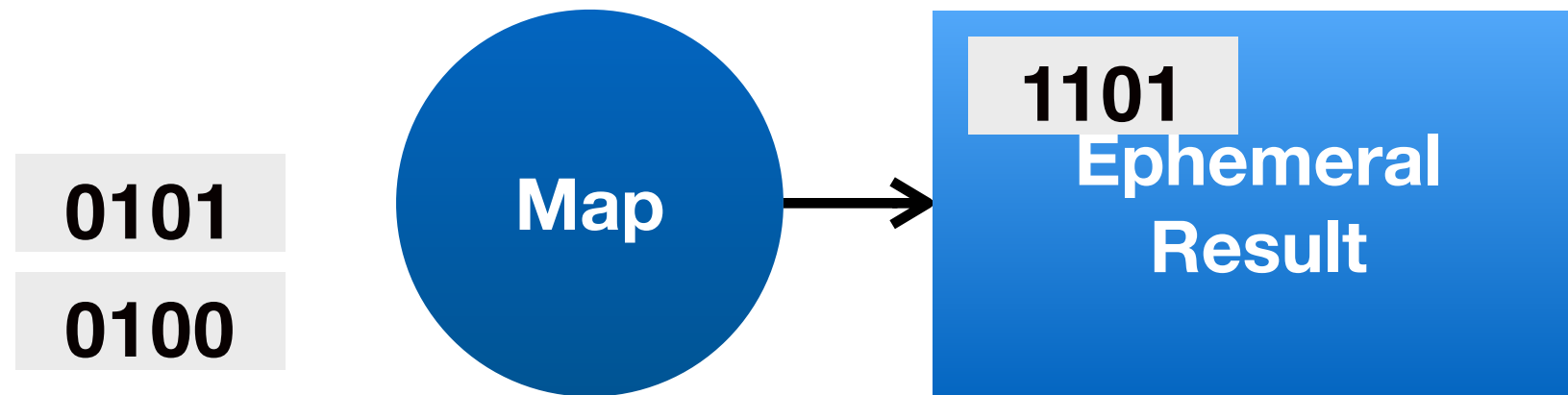
**How long to keep results around?**

# Ephemeral Results

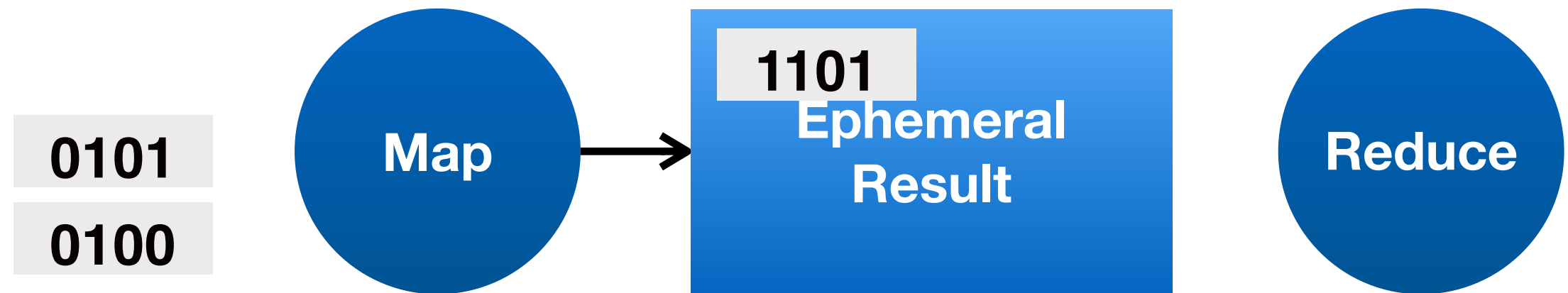# Ephemeral Results



0101
0100
1101
Ephemeral Result

# Ephemeral Results

**0101**

**0100**

**Map**

**1101**

**Ephemeral Result**

# Ephemeral Results

**0101**

**0100**

Map

1101

**Ephemeral Result**

Reduce

# Ephemeral Results

**0101**

**0100**

**Map**
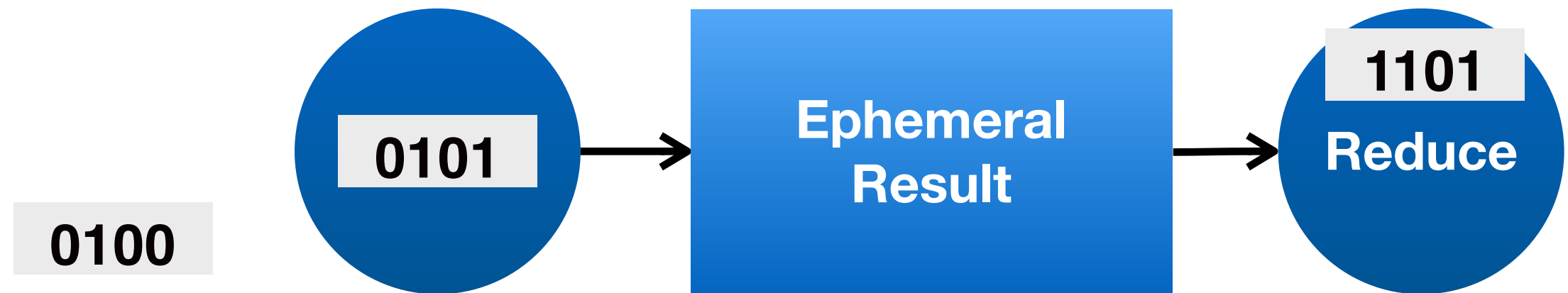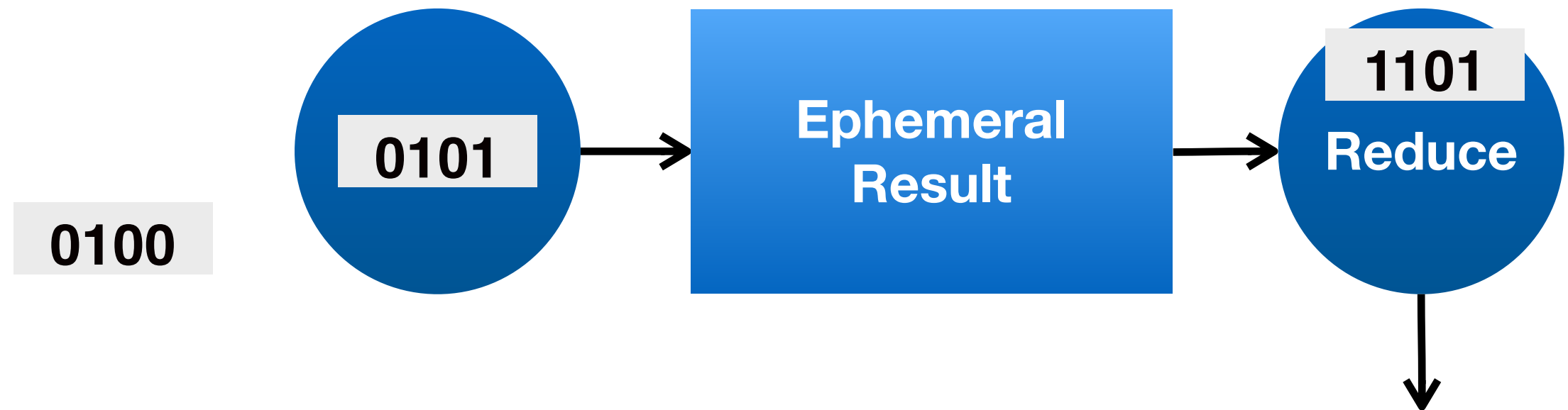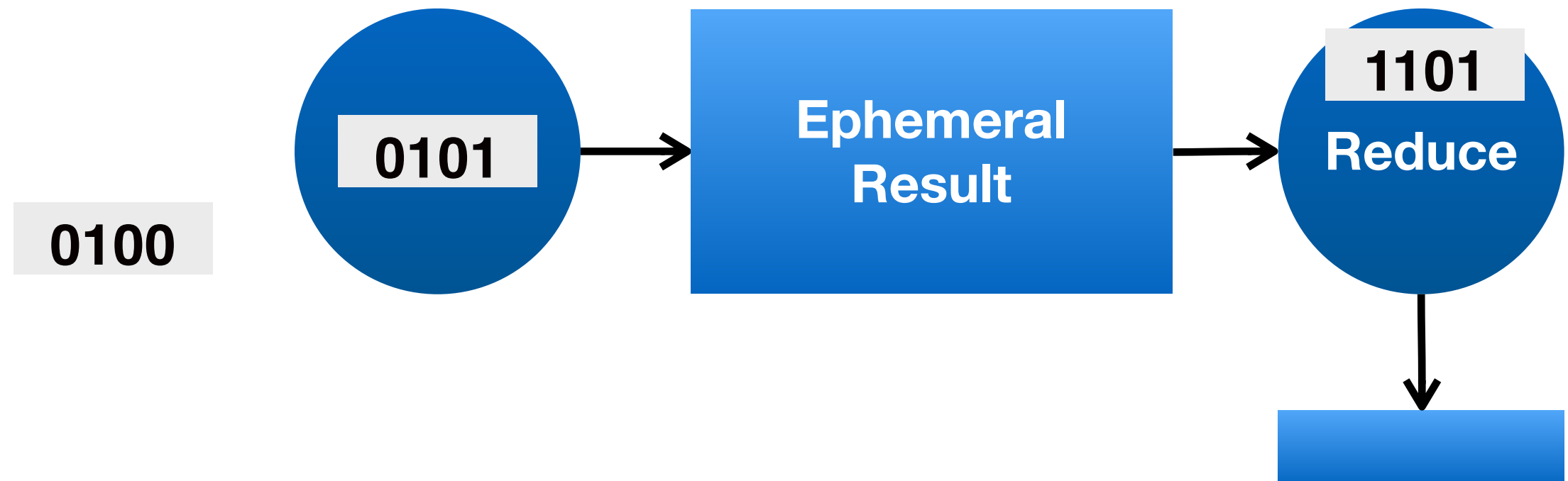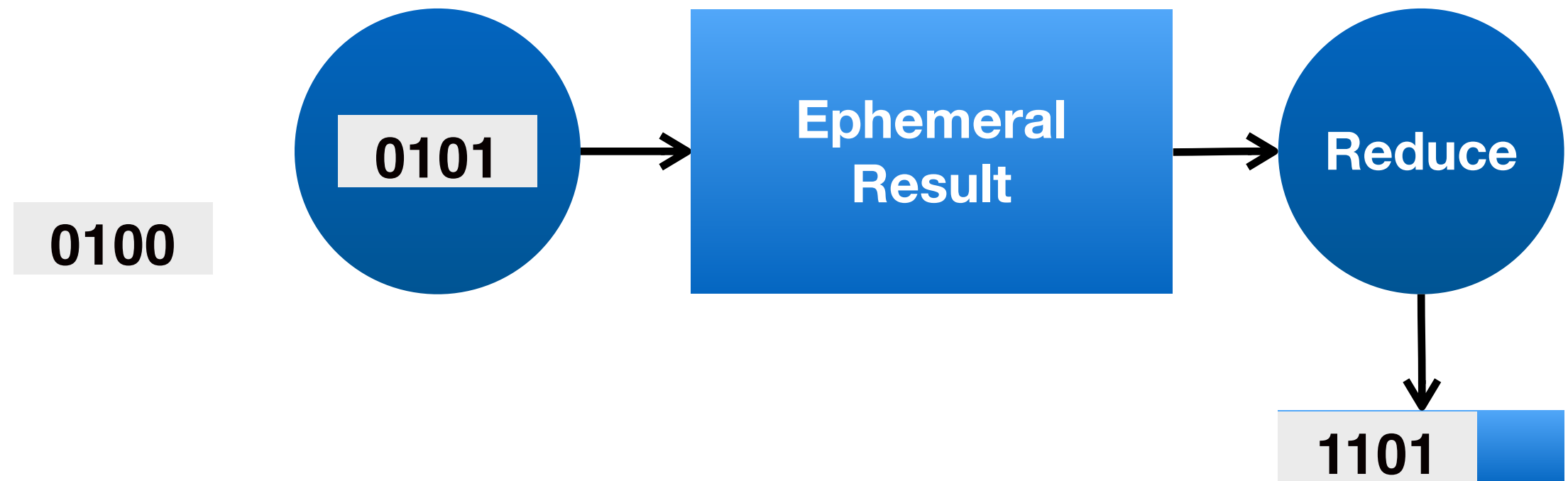
**1101**
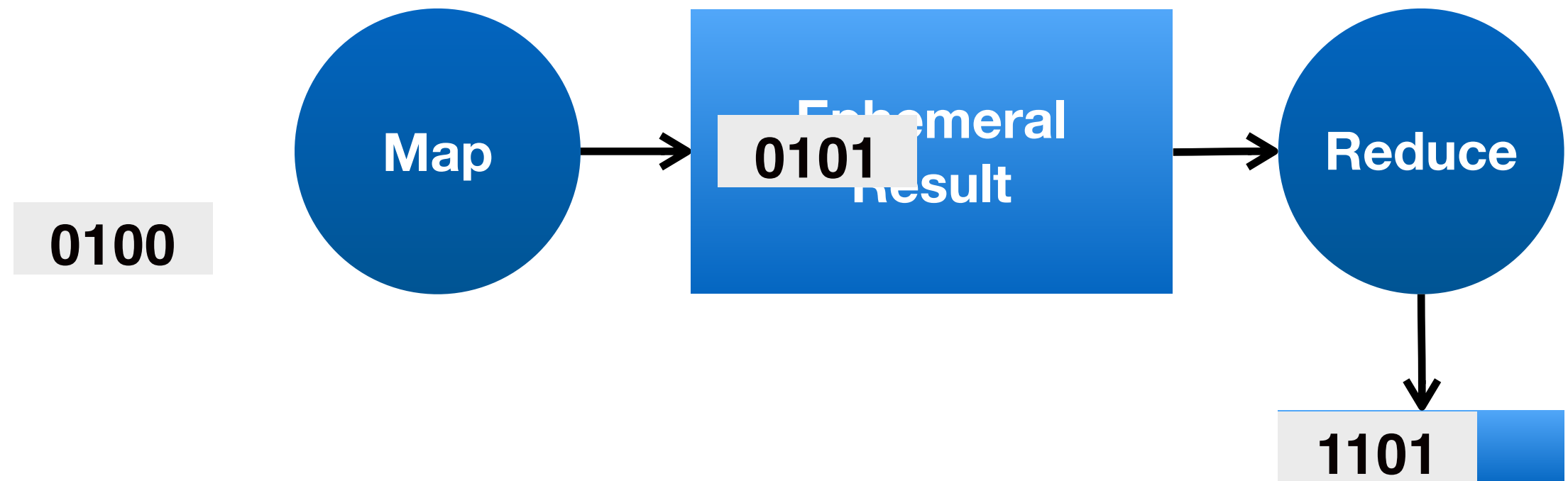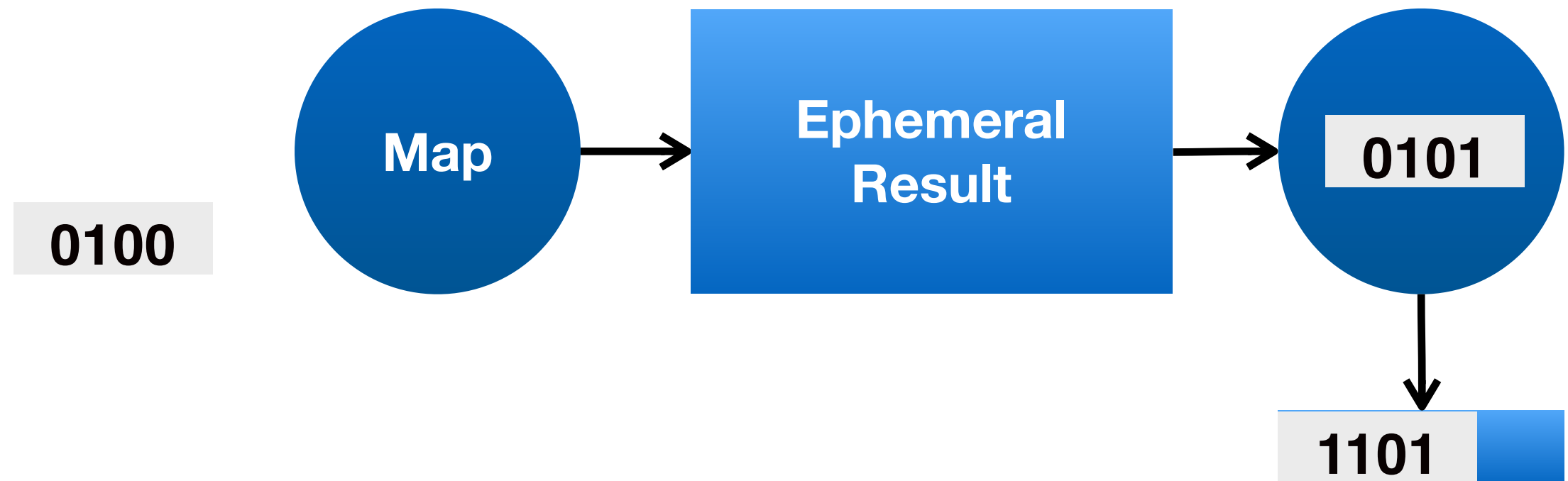
**Ephemeral Result**

**Reduce**

# Ephemeral Results

# Ephemeral Results

# Ephemeral Results

# Ephemeral Results

0100

0101 → **Ephemeral Result** → 1101 **Reduce**

# Ephemeral Results
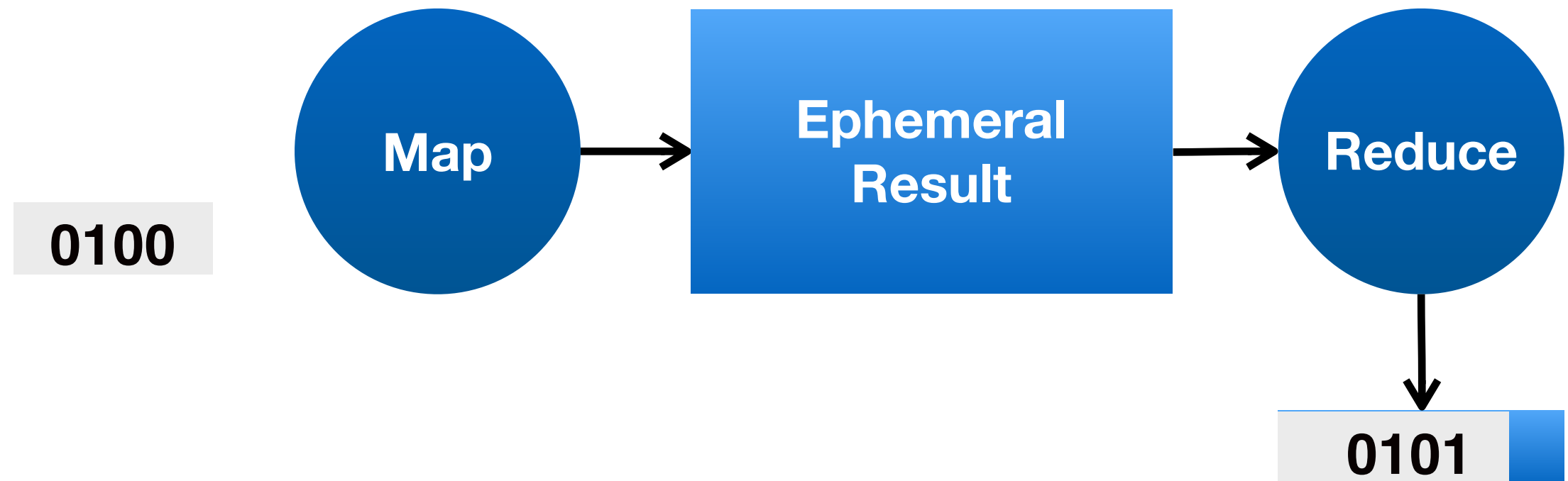
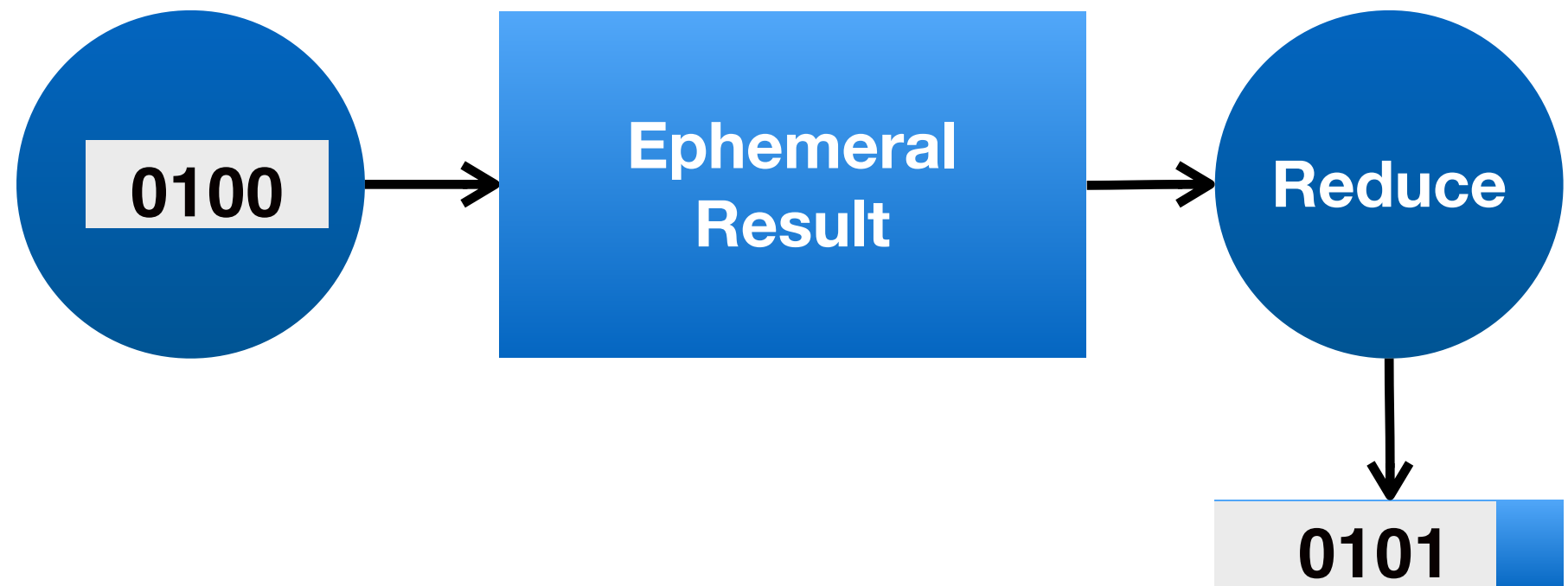# Ephemeral Results

# Ephemeral Results

# Ephemeral Results

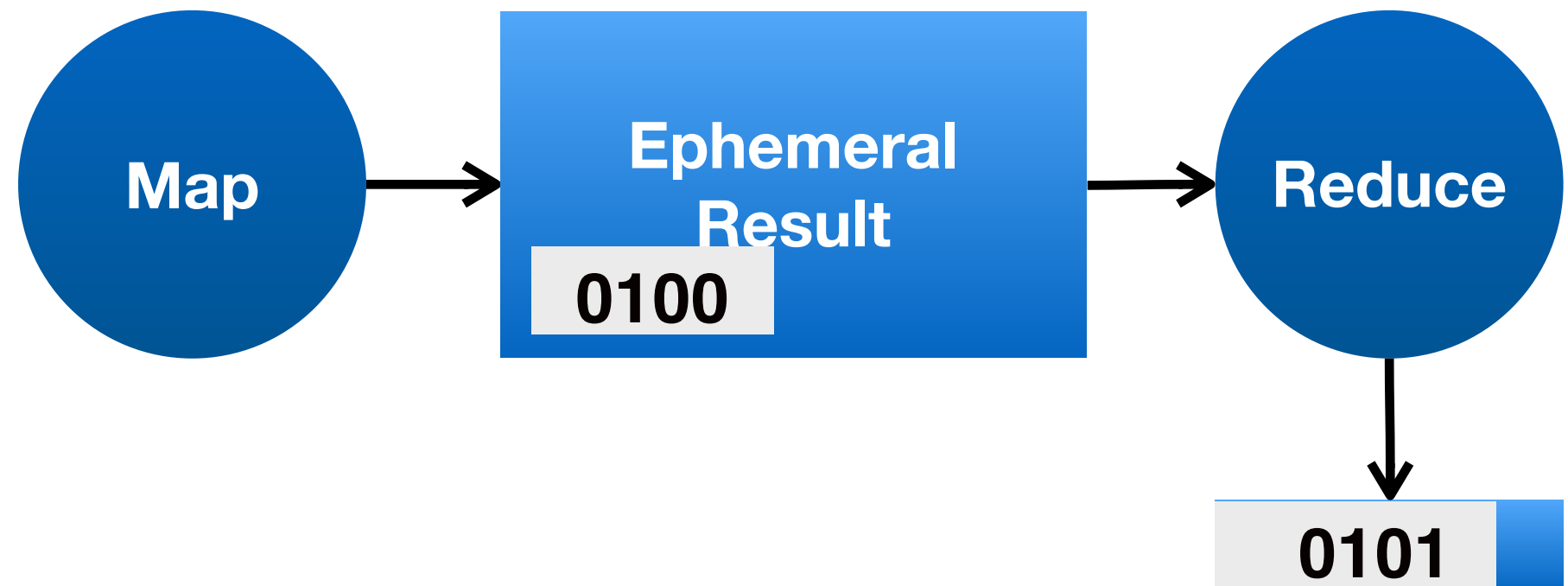# Ephemeral Results

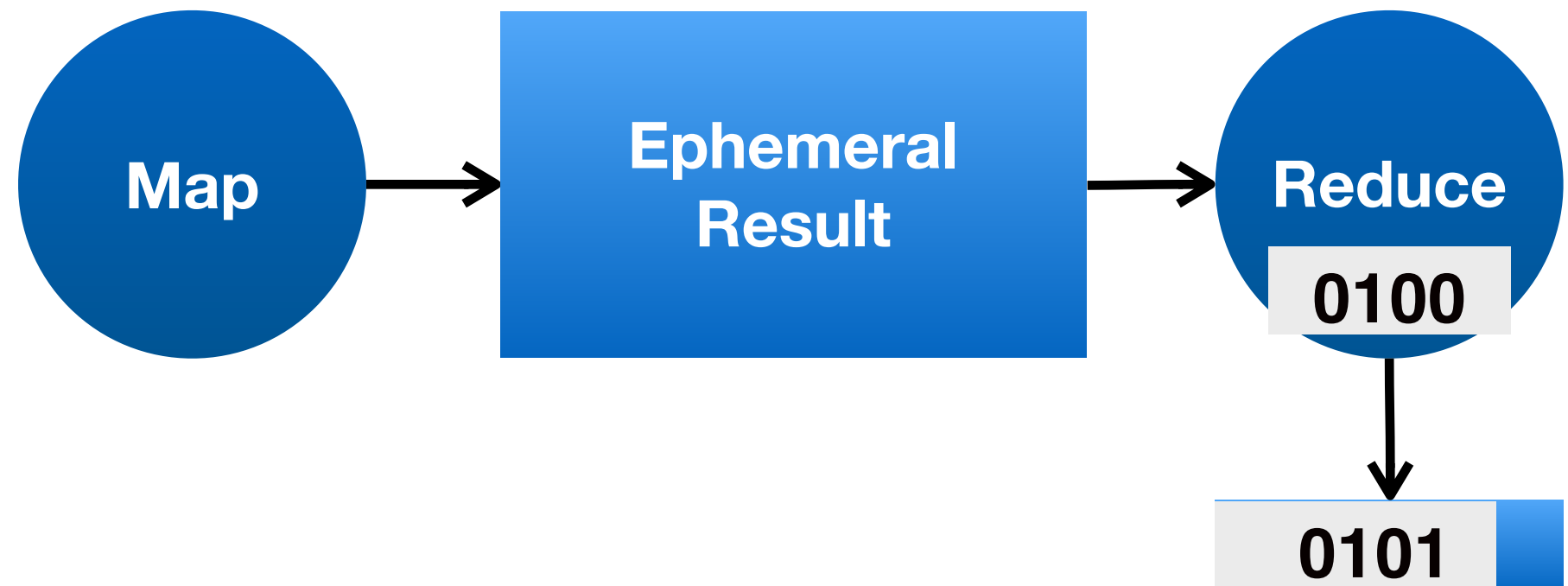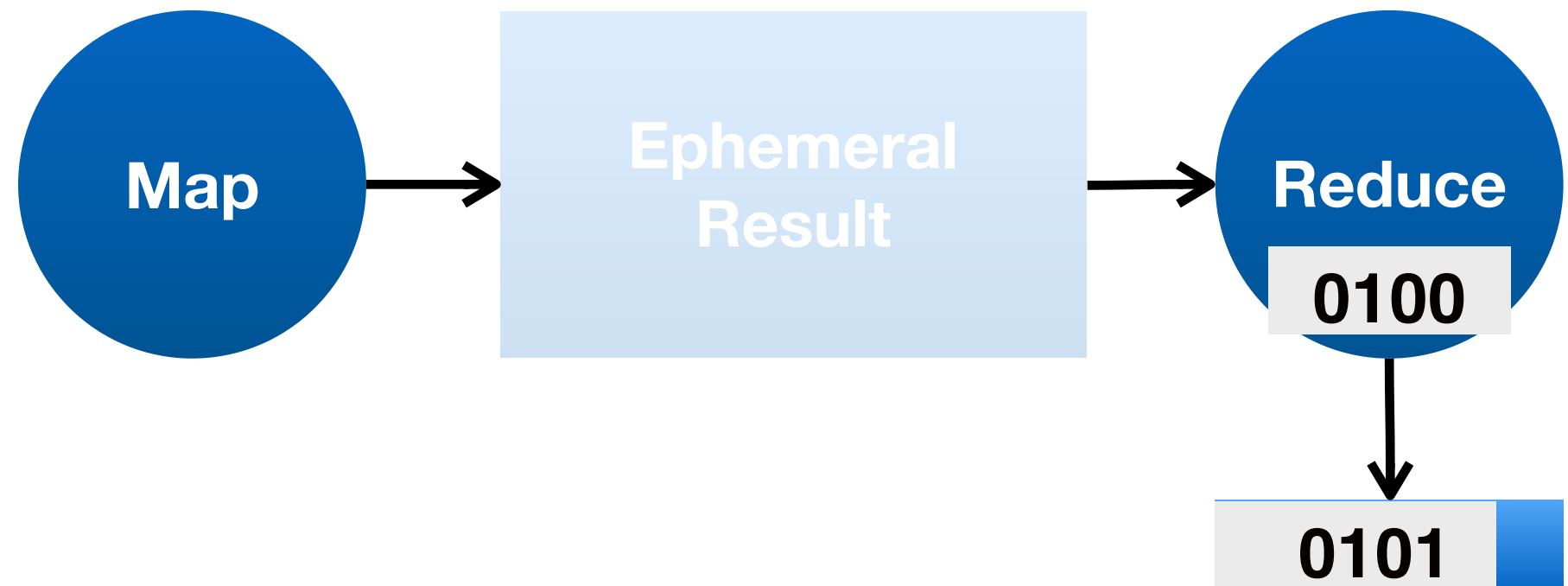# Ephemeral Results

# Ephemeral Results

# Ephemeral Results

# Checkpointed Results

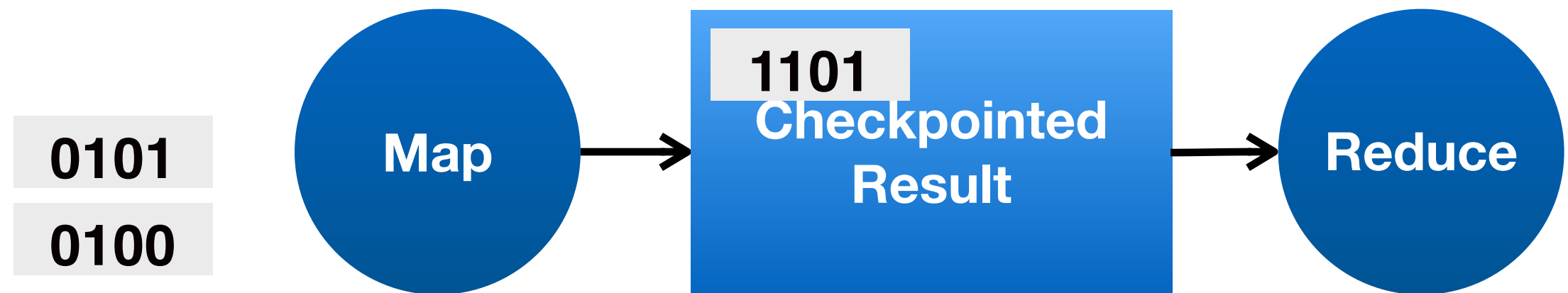# Checkpointed Results



0101
0100
1101
Checkpointed Result
Reduce

# Checkpointed Results

# Checkpointed Results

# Checkpointed Results

0100

0101 → 1101 Checkpointed Result → Reduce

# Checkpointed Results

**0100**

**Map** → **1101** **0101** Checkpointed Result → **Reduce**

22

# Checkpointed Results



0100

Map

1101

0101

Checkpointed Result

Reduce

# Checkpointed Results



22

# Checkpointed Results

# Checkpointed Results

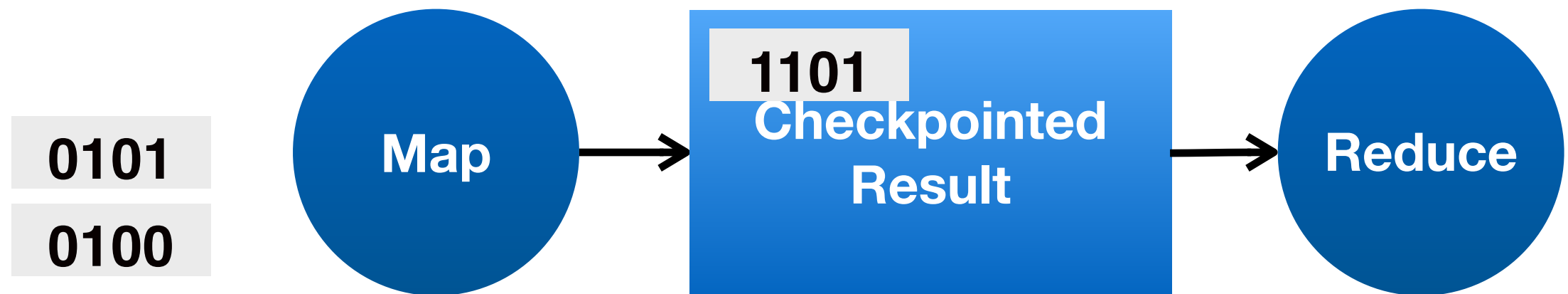Map ➝ **1101** **Checkpointed** **0101** **Result** **0100** ➝ Reduce

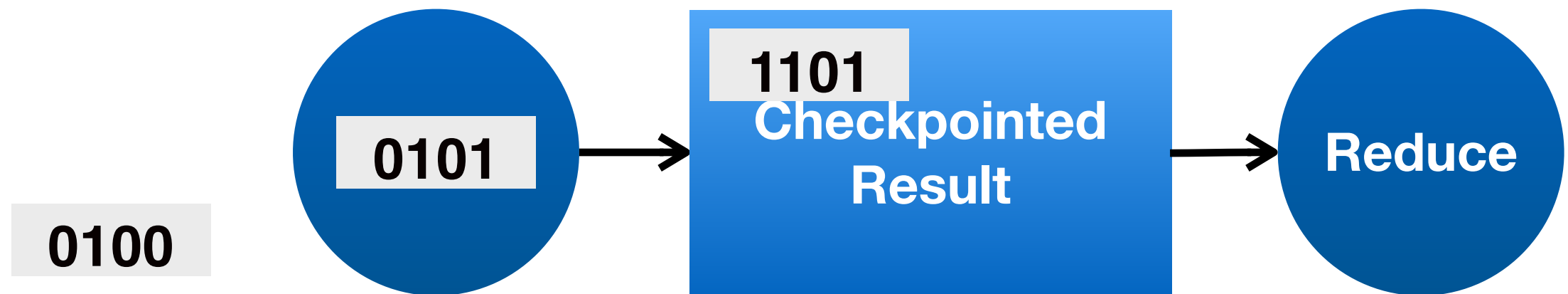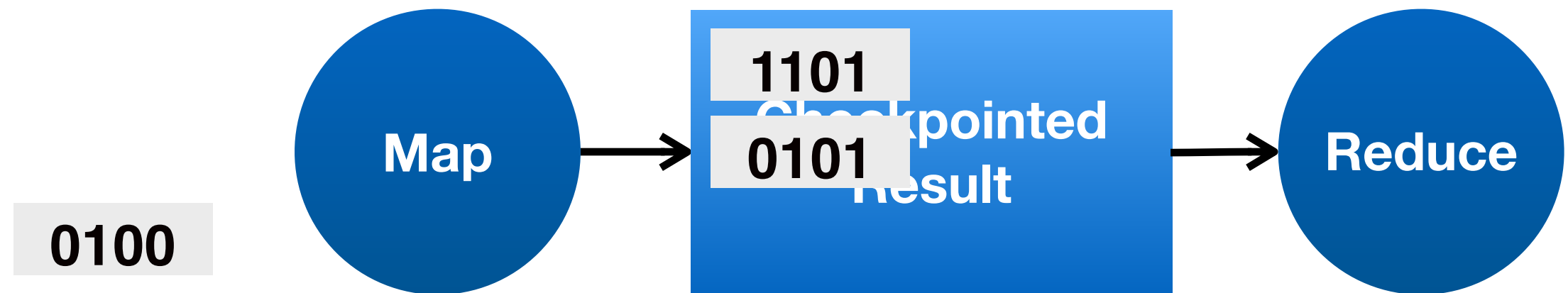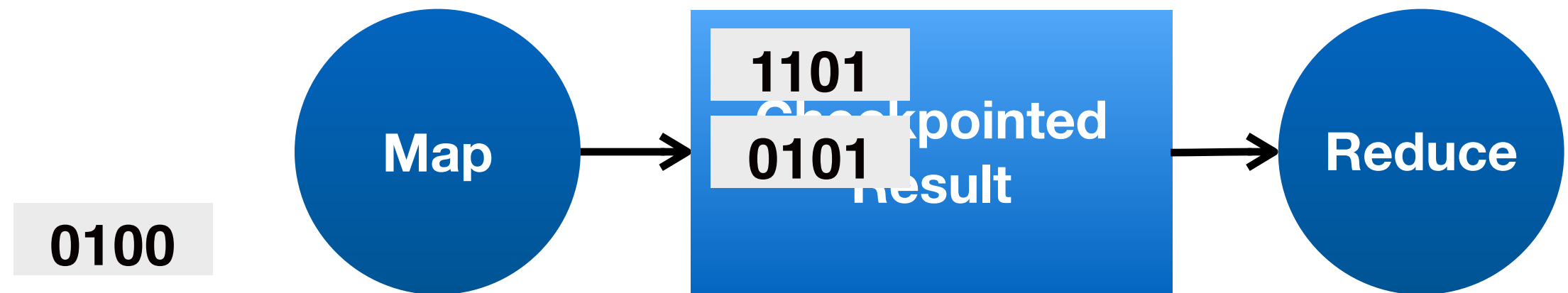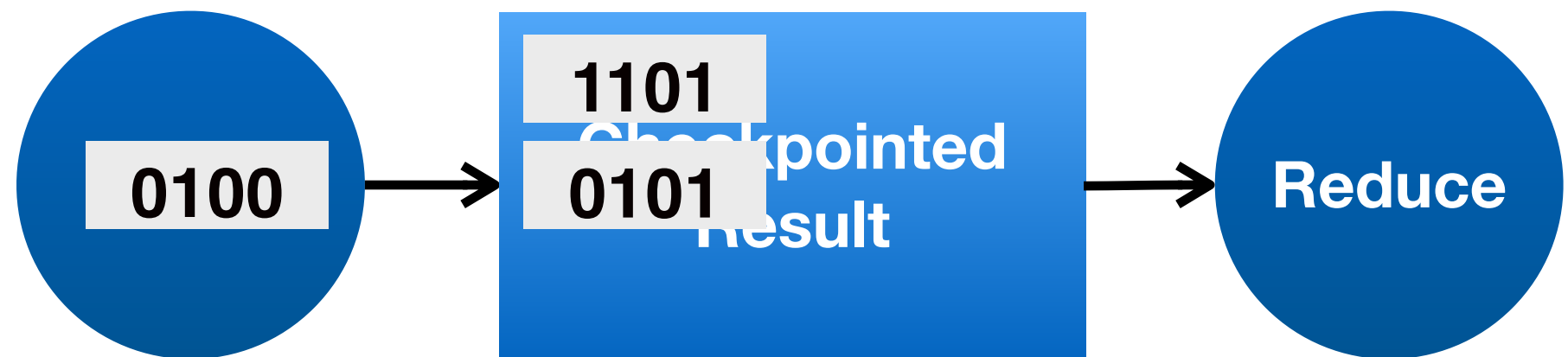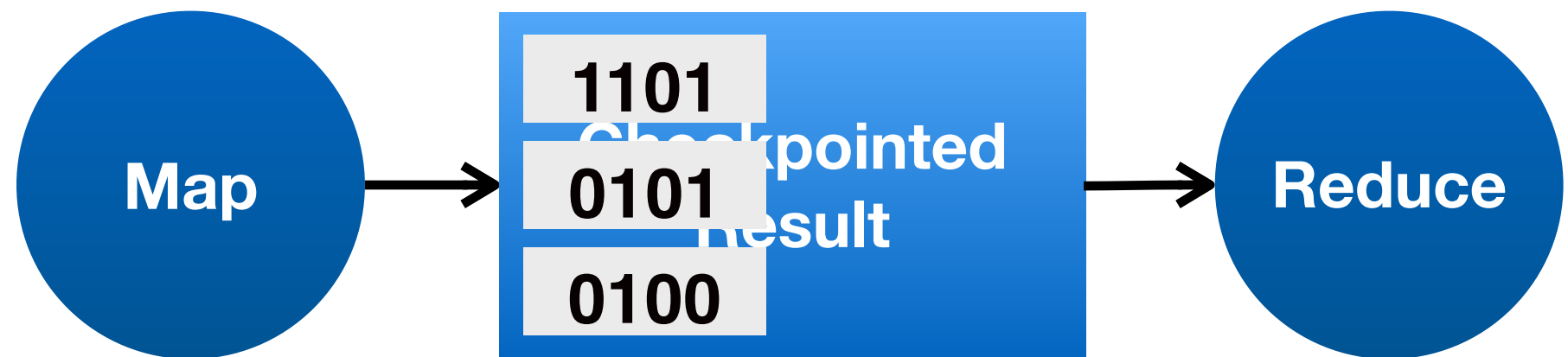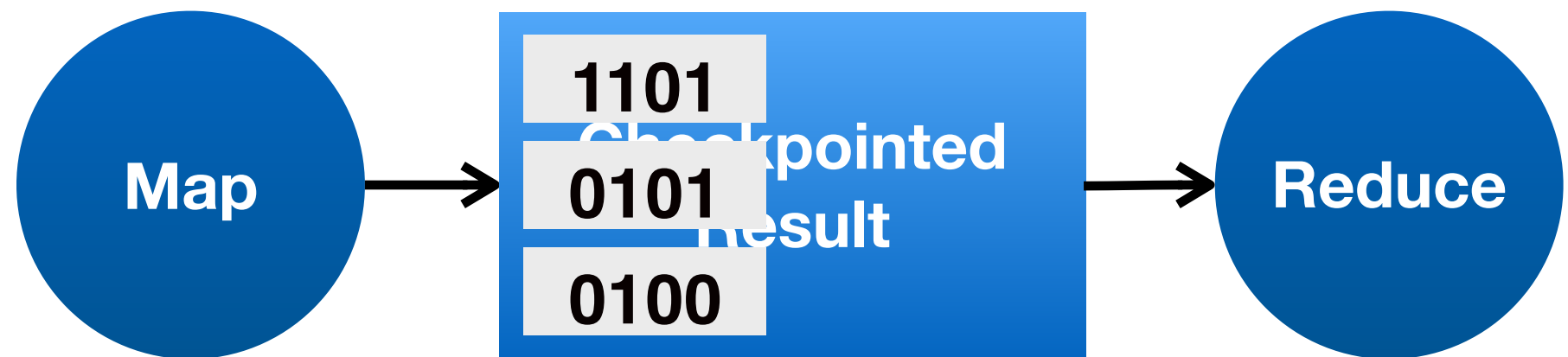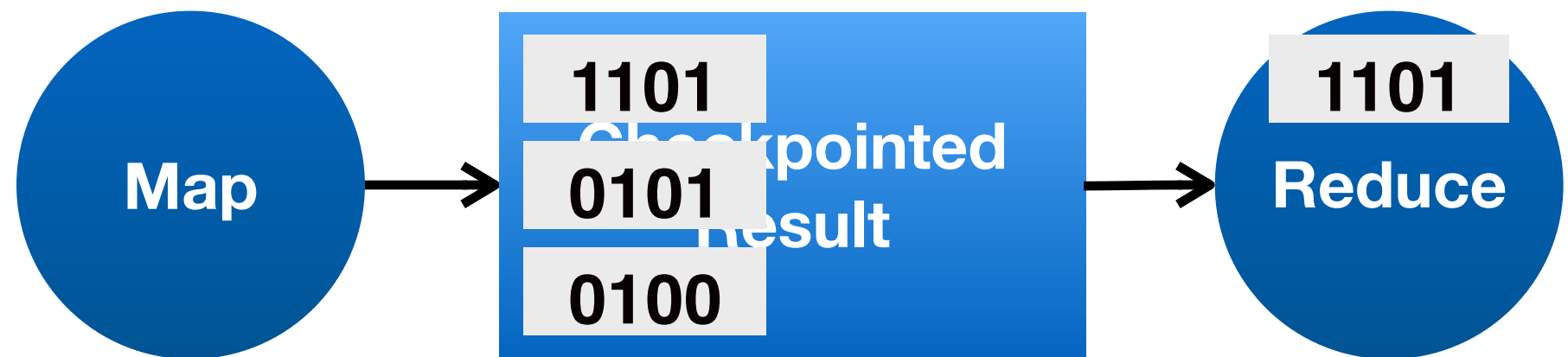# Checkpointed Results

# Checkpointed Results

# Checkpointed Results

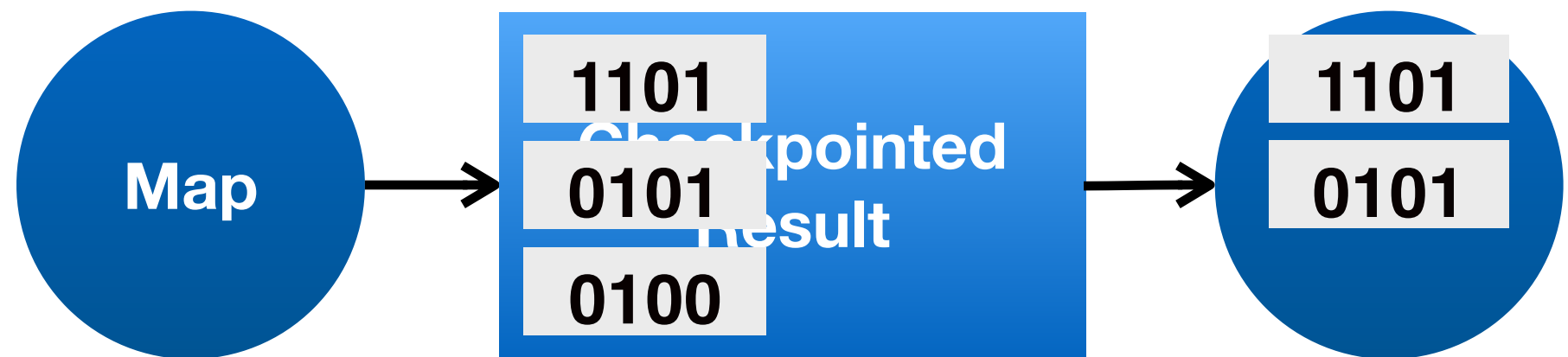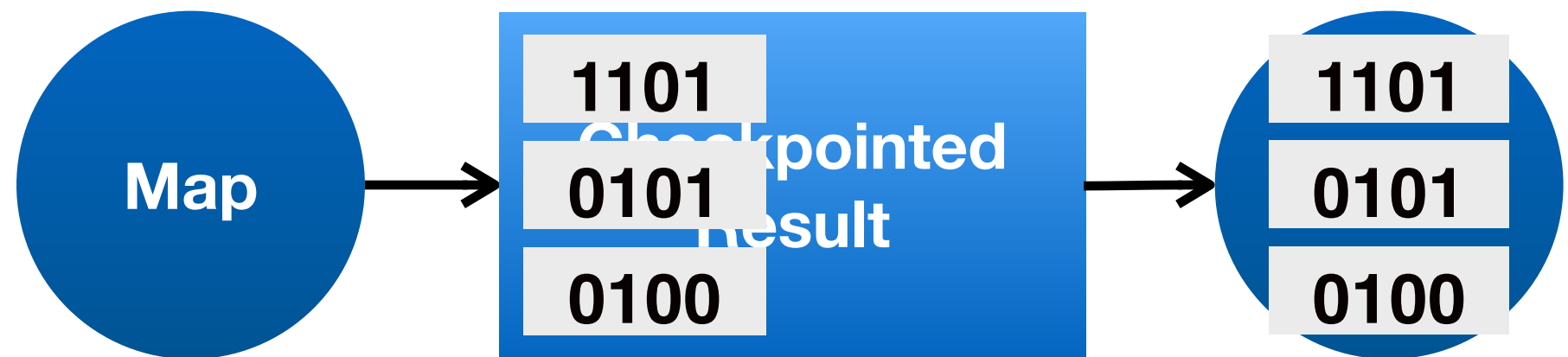# Result Characteristics

|  | Ephemeral | Checkpointed |
|---|---|---|
| **Pipelined** | Low-latency | Low-latency<br><br>**Fine-grained fault tolerance** |
| **Blocking** | Easy to reason about resource consumption | **Fine-grained fault tolerance**<br><br>Easy to reason about resource consumption |

# Benefits

- Very flexible design

- Decouples **high-level requirements** from runtime
  - Fault tolerance for batch vs. streaming
  - Different program optimization paths
  - Iterative programs
  - Interactive queries

# Use cases



Heavy data pipelines

Graph analytics

Large-scale Machine Learning

Real-time stream processing

# Flink Stack

# Current Implementations

**Pipelined** vs. Blocking

**Ephemeral** vs. Checkpointed

**Backpressure** vs. No Backpressure

# Current Implementations

Pipelined vs. **Blocking**

Ephemeral vs. **Checkpointed**

Backpressure vs. **No Backpressure**

# Pipelined vs. Blocking in Flink

- Default type for both batch and streaming programs: **Pipelined**

- In batch mode only: use blocking exchange if necessary (e.g. to avoid deadlocks or break up long pipelines)

- More details: https://cwiki.apache.org/confluence/display/FLINK/Data+exchange+between+tasks

# ExecutionConfig

```java
// Set up the execution environment
ExecutionEnvironment env = ExecutionEnvironment
  .getExecutionEnvironment();

ExecutionConfig conf = env.getConfig();

// Remote data exchange is blocking (local is pipelined)
conf.setExecutionMode(ExecutionMode.BATCH);

// Remote and local data exchange is blocking
conf.setExecutionMode(ExecutionMode.BATCH_FORCED);

// Remote and local data exchange is pipelined except
// when necessary to avoid deadlocks etc. [DEFAULT]
conf.setExecutionMode(ExecutionMode.PIPELINED);

// Remote and local data exchange is always pipelined
conf.setExecutionMode(ExecutionMode.PIPELINED_FORCED);
```

**flink.apache.org**

@ApacheFlink