# Kafka Streams Testing

## A Deep Dive

Ivan Ponomarev, John Roesler

# Who Are We



Ivan Ponomarev:

- Software Engineer at KURS, tutor at MIPT
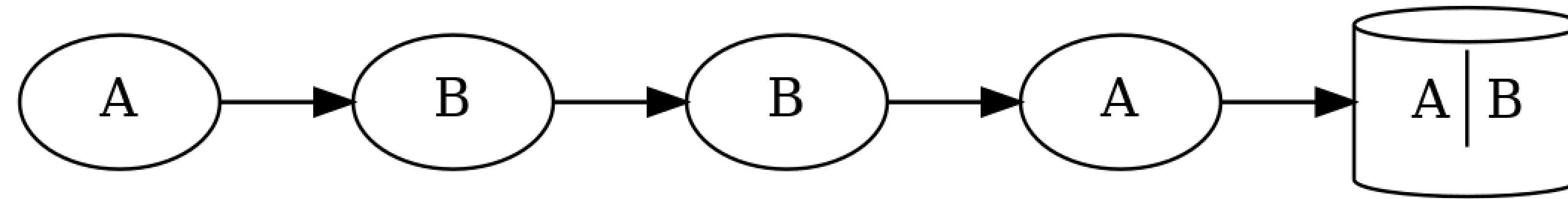- Apache Kafka Contributor

# Who Are We



John Roesler:

- Software Engineer at Confluent
- Apache Kafka Committer and PMC member

# Kafka Streams Testing: A Deep Dive

1. Purpose: cover testing methodologies for Kafka Streams
   - "Unit" Testing: TopologyTestDriver
   - Integration Testing: KafkaStreams
2. Start with motivating example (from Ivan's production)
3. A flawed testing approach: unit testing doesn't work for this example
4. Deep-dive into the testing framework
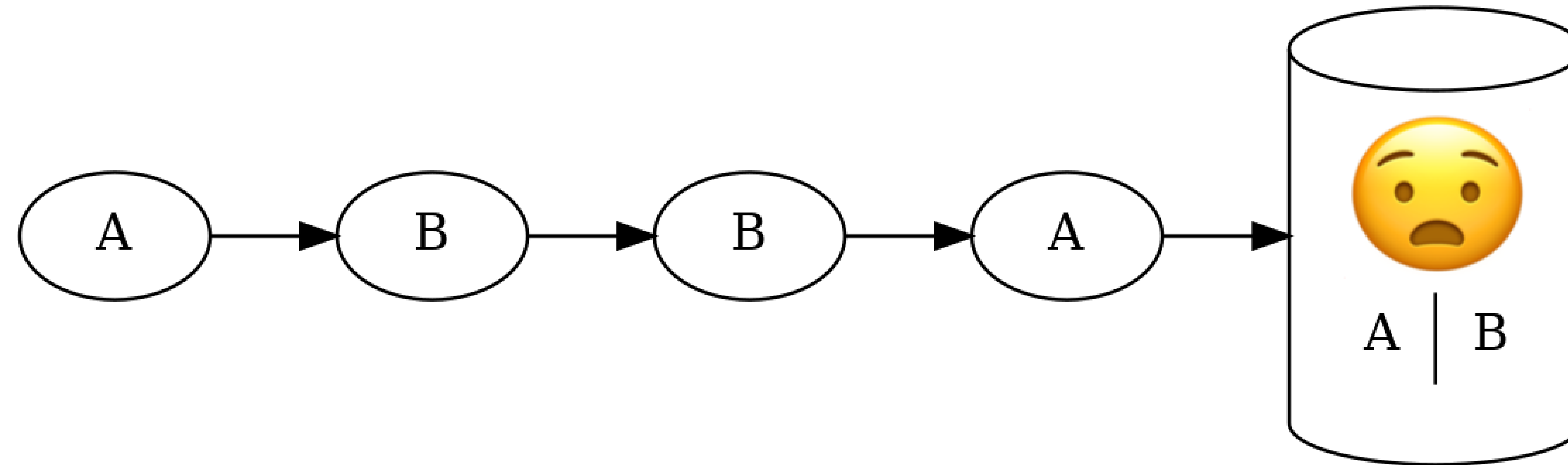5. Correctly testing the example with integration tests

# The task
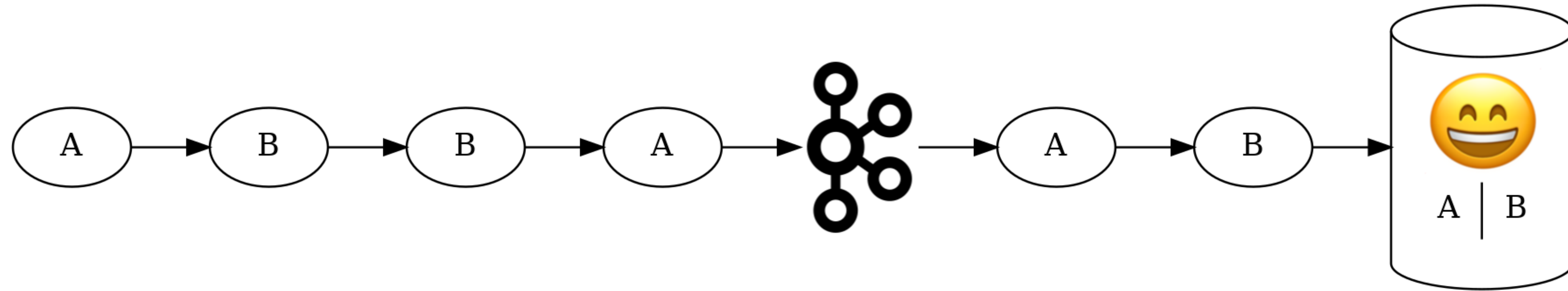
## Save different source IDs in the database

# The problem
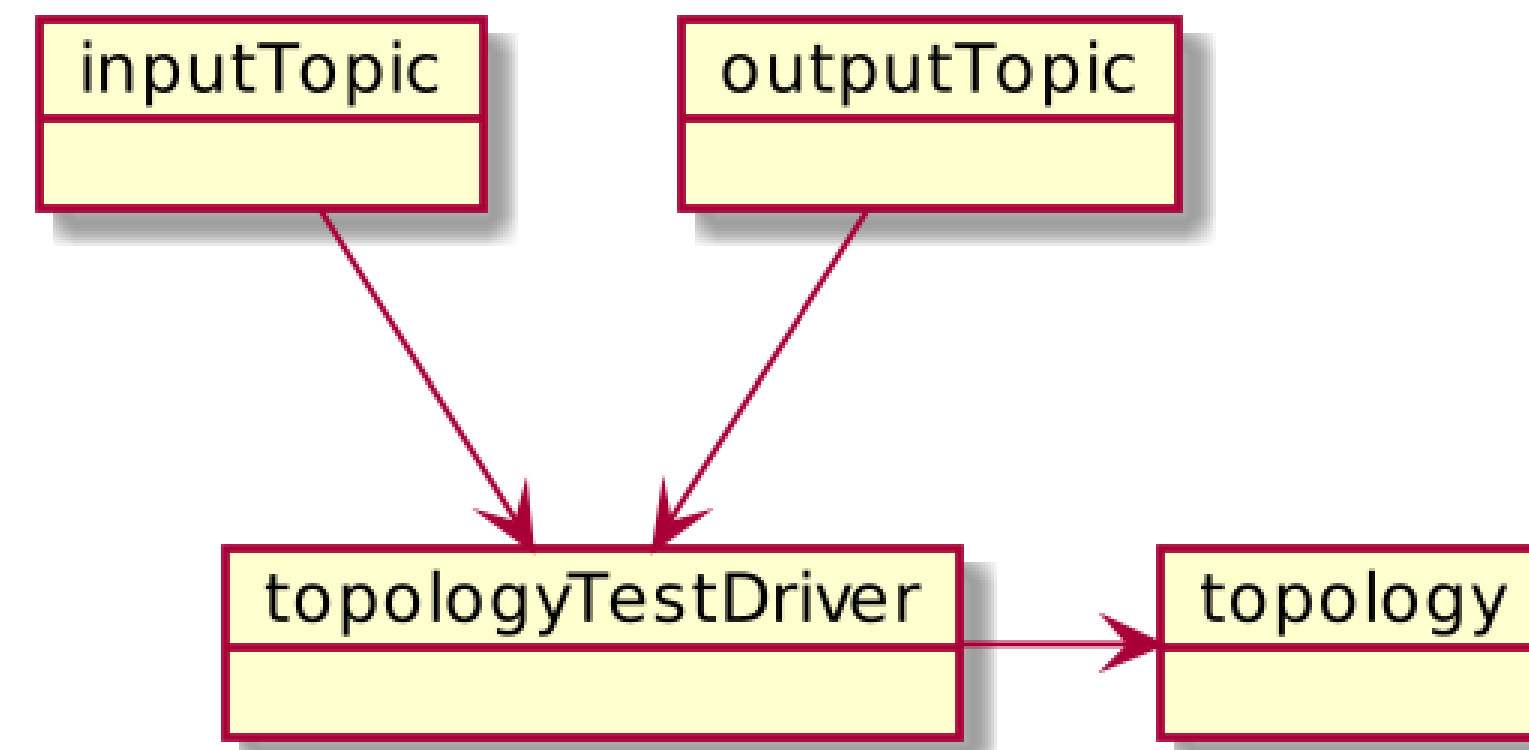
## Too many writes to the database

# The solution

Let's deduplicate using Kafka Streams!

# TopologyTestDriver

# `TopologyTestDriver` capabilities

```
                    pipeInput (V)                   readValue ()
┌───┬───┐           pipeInput (K, V)   ╭──────────╮ readKeyValue()    ┌───┬───┐
│ K │ V │ ─────────────────────────▶  │ Topology │ ───────────────▶  │ K │ V │
└───┴───┘                             ╰──────────╯                    └───┴───┘
```

| What is being sent/received | TestInputTopic methods | TestOutputTopic methods |
|---|---|---|
| A single value | `pipeInput (V)` | `V readValue ()` |
| A key/value pair | `pipeInput (K, V)` | `KeyValue<K,V> readKeyValue()` |

# `TopologyTestDriver` capabilities

```
                pipeValueList(List<V>)          readValuesToList()
 K | V          pipeKeyValueList                readKeyValuesToList()      K | V
 K | V             (List<KeyValue<K,V>>)   Topology    readKeyValuesToMap()   K | V
 K | V                                                                        K | V
```

| What is being sent/received | TestInputTopic methods | TestOutputTopic methods |
| --- | --- | --- |
| A list of values | `pipeValueList (List<V>)` | `List<V> readValuesToList()` |
| A list of key/value pairs | `pipeKeyValueList (List<KeyValue<K,V>>)` | `List<KeyValue<K,V>> readKeyValuesToList()` <br> `Map<K,V> readKeyValuesToMap()` |

# `TopologyTestDriver` capabilities



| What is being sent/received | TestInputTopic methods | TestOutputTopic methods |
|---|---|---|
| A list of Records | `pipeRecordList (List<? extends TestRecord<K, V>>)` | `List<TestRecord<K, V>> readRecordsToList()` |

# Demo

1. Spring Boot app
2. Let's do some test-driven development and first write a test
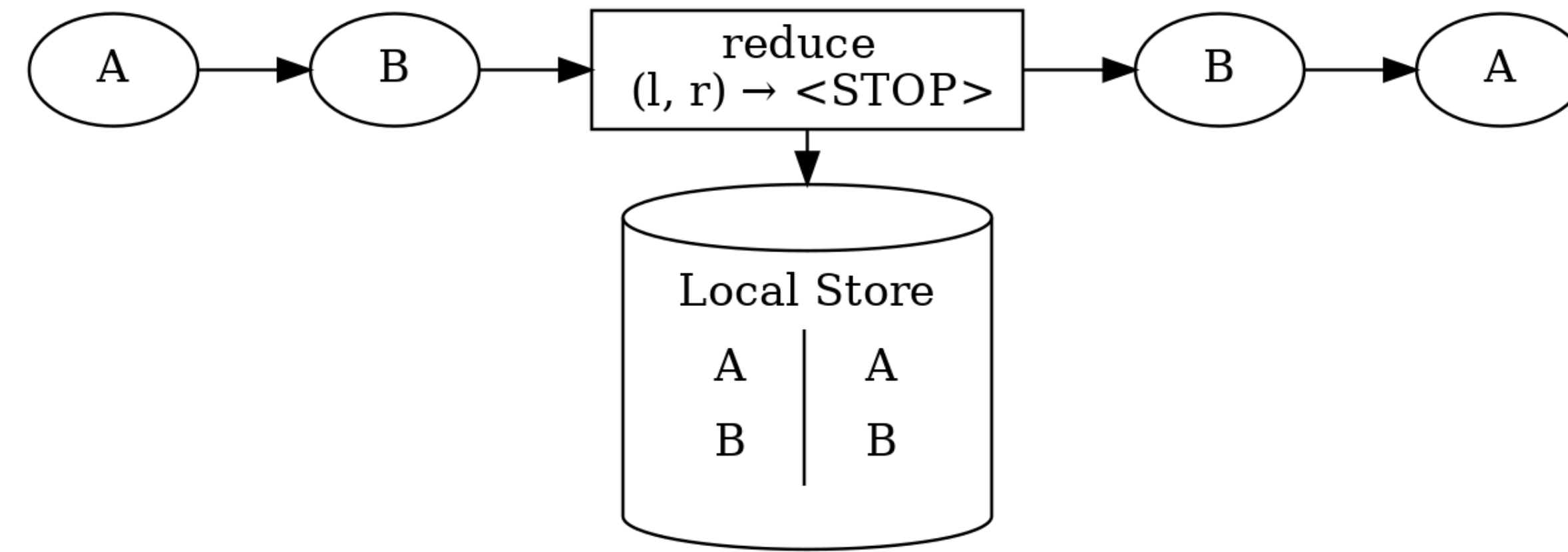3. Writing a test with TTDriver
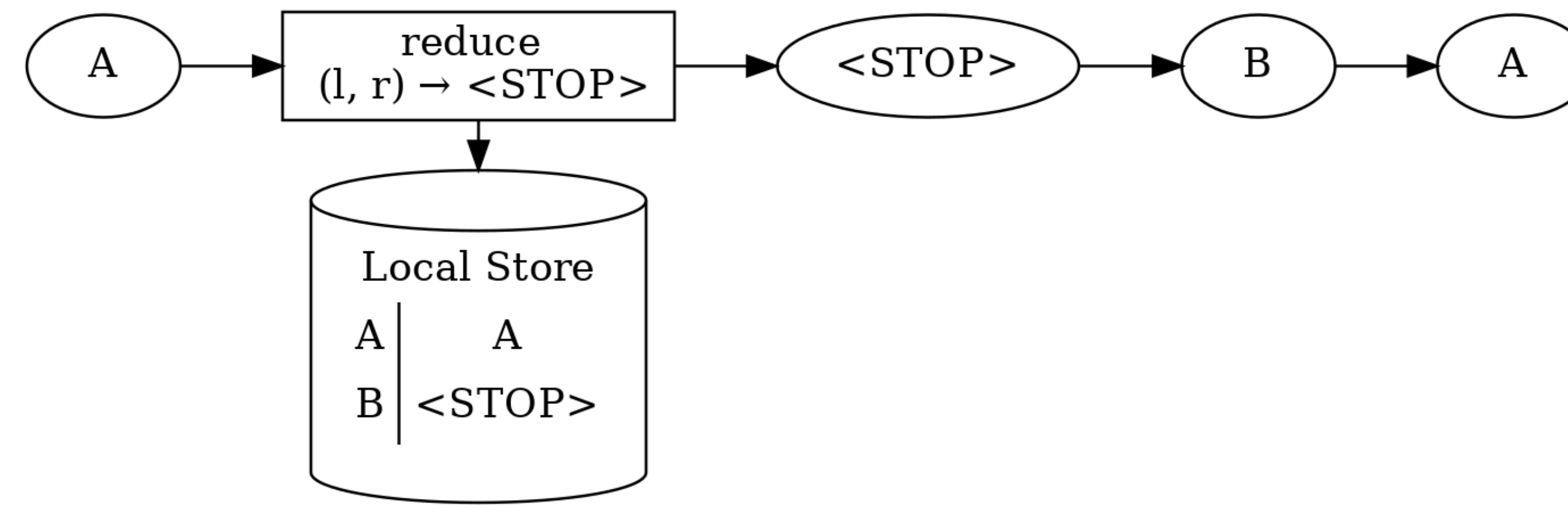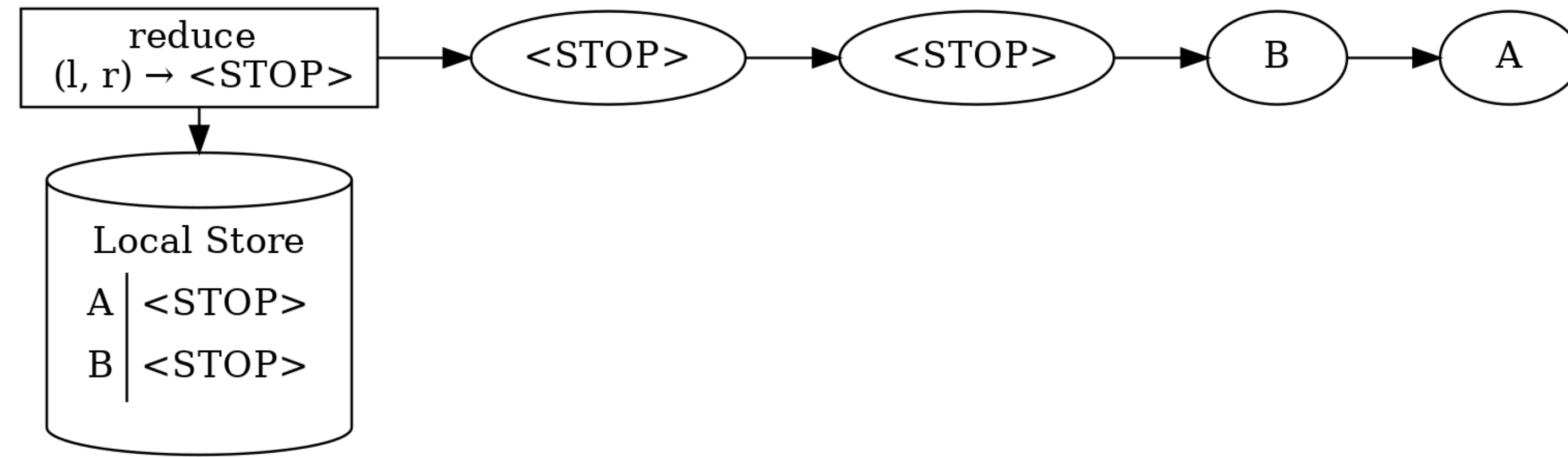
# A "Simple Solution"

# A "Simple Solution"

# A "Simple Solution"



A → B → reduce (l, r) → <STOP> → B → A

Local Store
A | A
B | B

# A "Simple Solution"

# A "Simple Solution"

# Demo

- writing the topology
- TopologyTestDriver test is green

# Tests are green

# Tests are green

build passing

Should we run this in production?

# What we saw in production:

# Why it's not working

| Kafka Streams | TopologyTestDriver |
| --- | --- |
| is a big data streaming framework | is a fast, deterministic testing framework |

# Why it's not working

| Kafka Streams | TopologyTestDriver |
| --- | --- |
| is a big data streaming framework | is a fast, deterministic testing framework |

- designed for high throughput
- throughput demands batching, buffering, caching, etc.
- caching is the culprit in this example

# Why it's not working

| Kafka Streams | TopologyTestDriver |
|---|---|
| is a big data streaming framework | is a fast, deterministic testing framework |
| • designed for high throughput | • designed for synchronous, immediate results |
| • throughput demands batching, buffering, caching, etc. | • flush cache after every update |
| • caching is the culprit in this example | |

# Why it's not working

## Caching in Kafka Streams

- don't immediately emit every aggregation result
- "soak up" repeated updates to the same key's aggregation
- configure cache size: `max.bytes.buffering` (10MB)
- configure cache flush interval: `commit.interval.ms` (30s)
- emit *latest* result on flush or eviction

# Why it's not working

# Why it's not working

# Why it's not working

# Why it's not working

# Why it's not working

# Why it's not working

# Demo

TopologyTestDriver vs. Kafka Streams execution loop

# Kafka Streams execution loop

# Kafka Streams execution loop

# Kafka Streams execution loop

# Kafka Streams execution loop

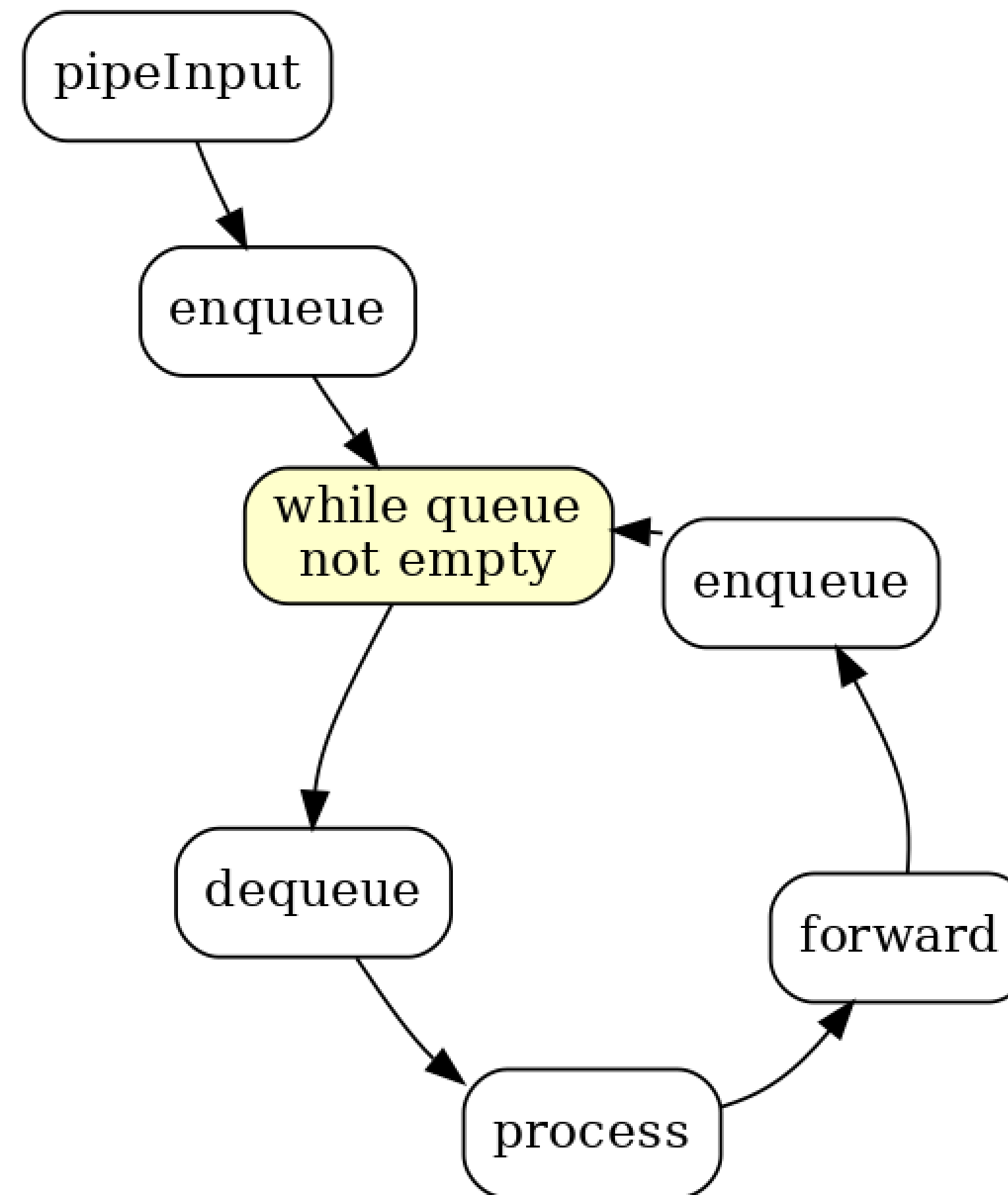# Kafka Streams execution loop

# Kafka Streams execution loop

# Kafka Streams execution loop

# Kafka Streams execution loop

# TopologyTestDriver execution loop

# TopologyTestDriver execution loop

# TopologyTestDriver execution loop

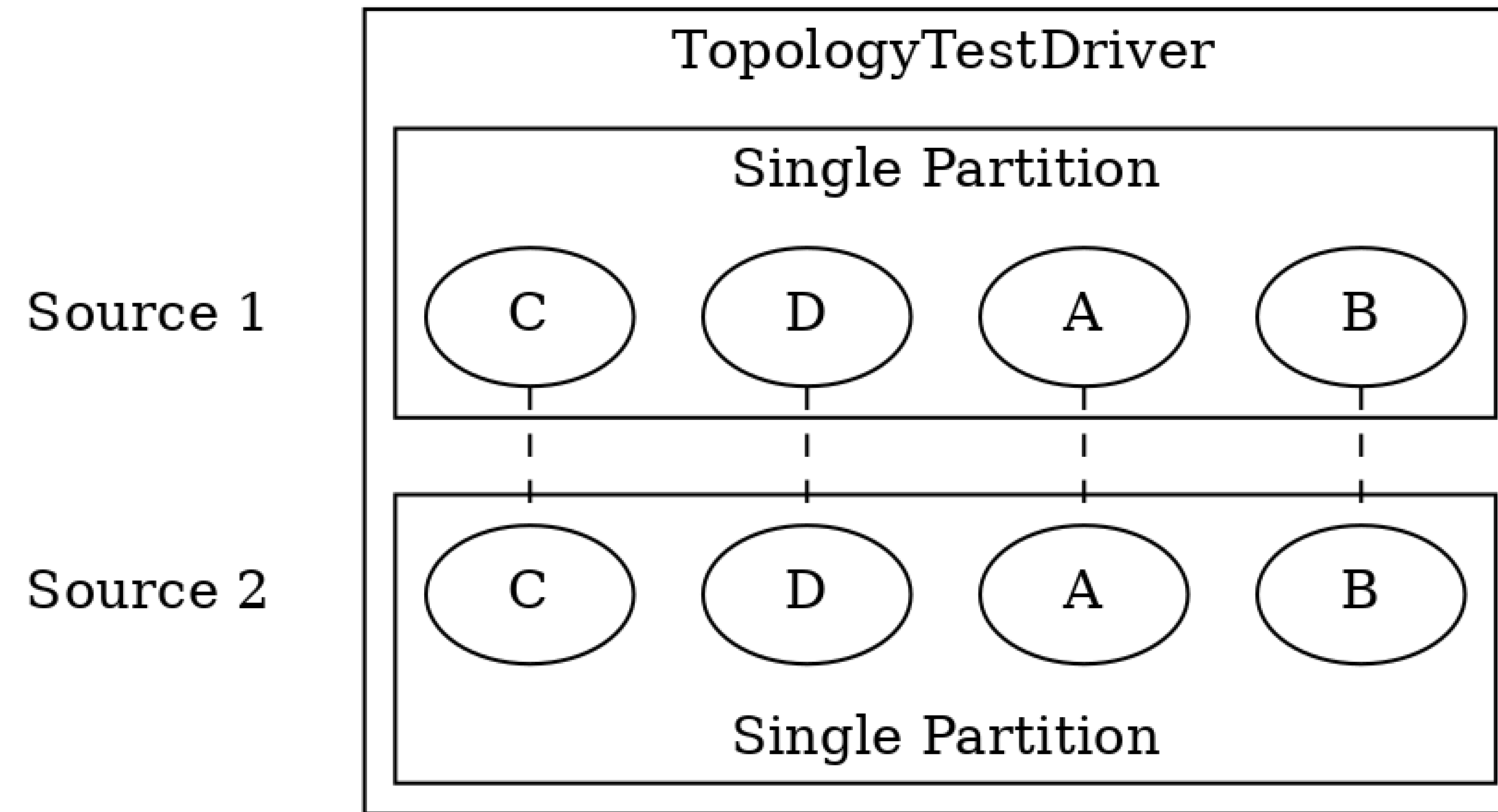# TopologyTestDriver execution loop

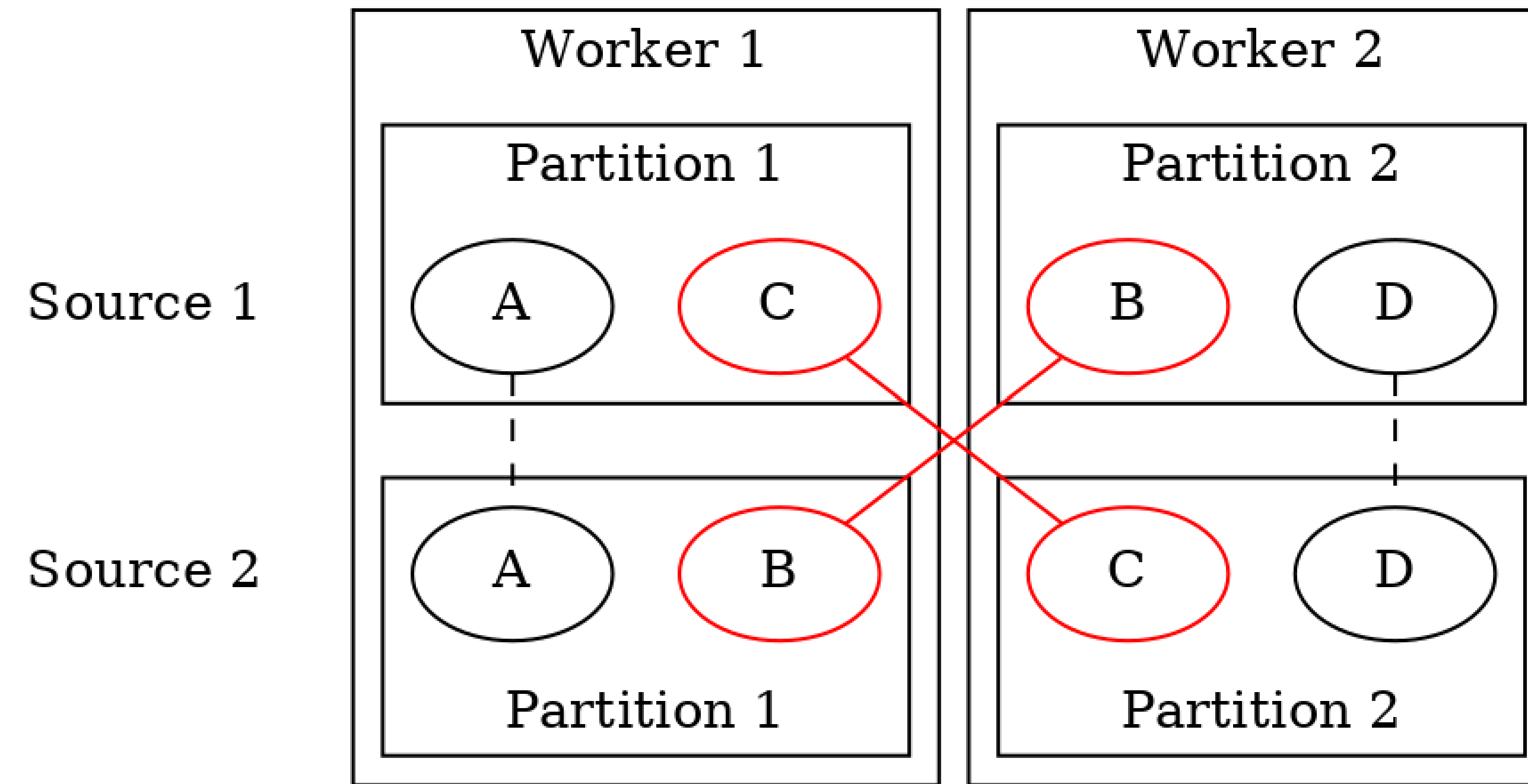# TopologyTestDriver execution loop

# TopologyTestDriver execution loop

# TopologyTestDriver execution loop

# What else?

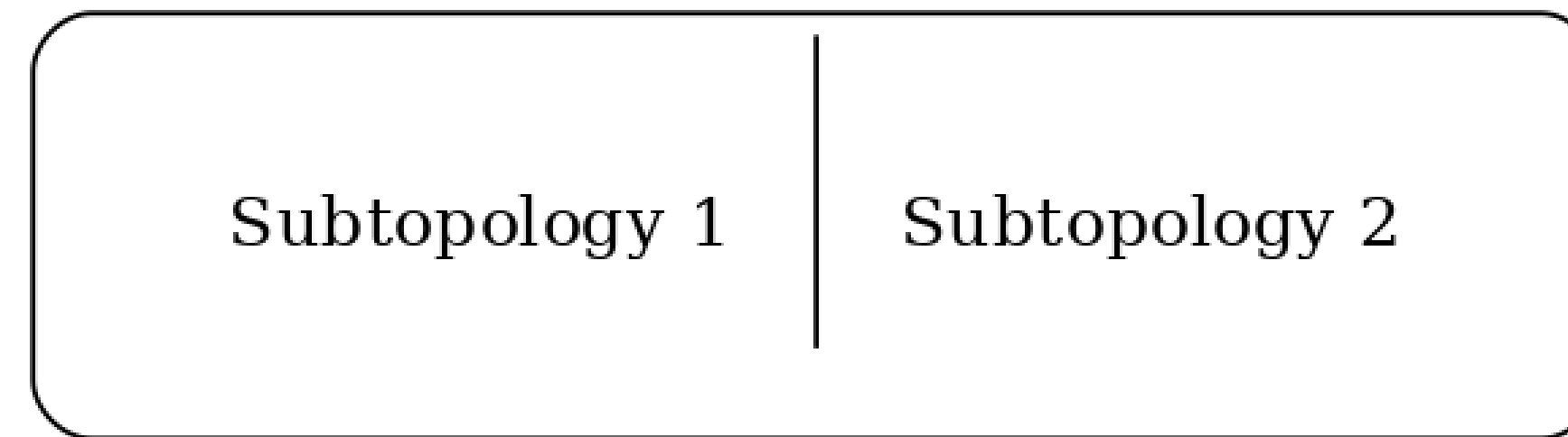What are other problems that can't be surfaced with `TopologyTestDriver`?
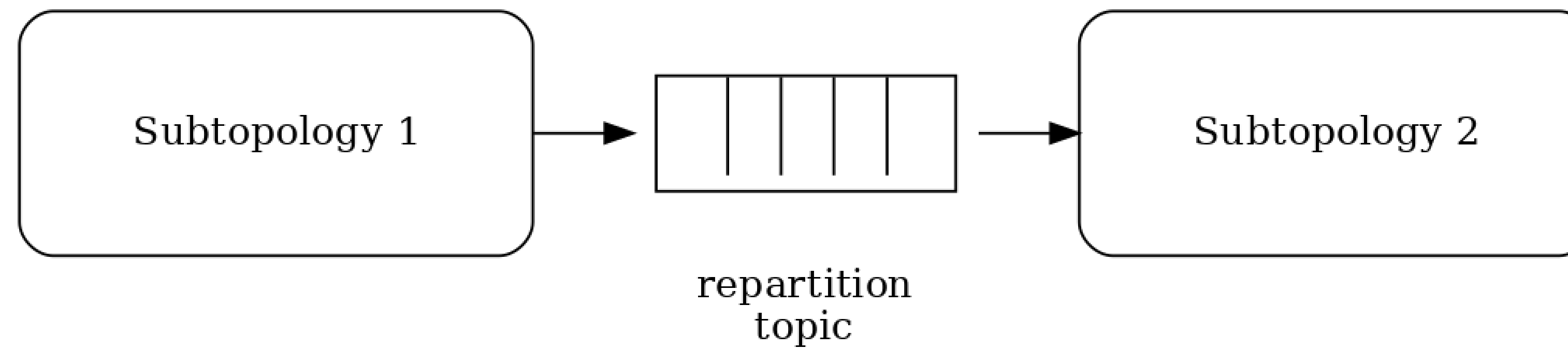
# Kafka Streams: co-partitioning problems

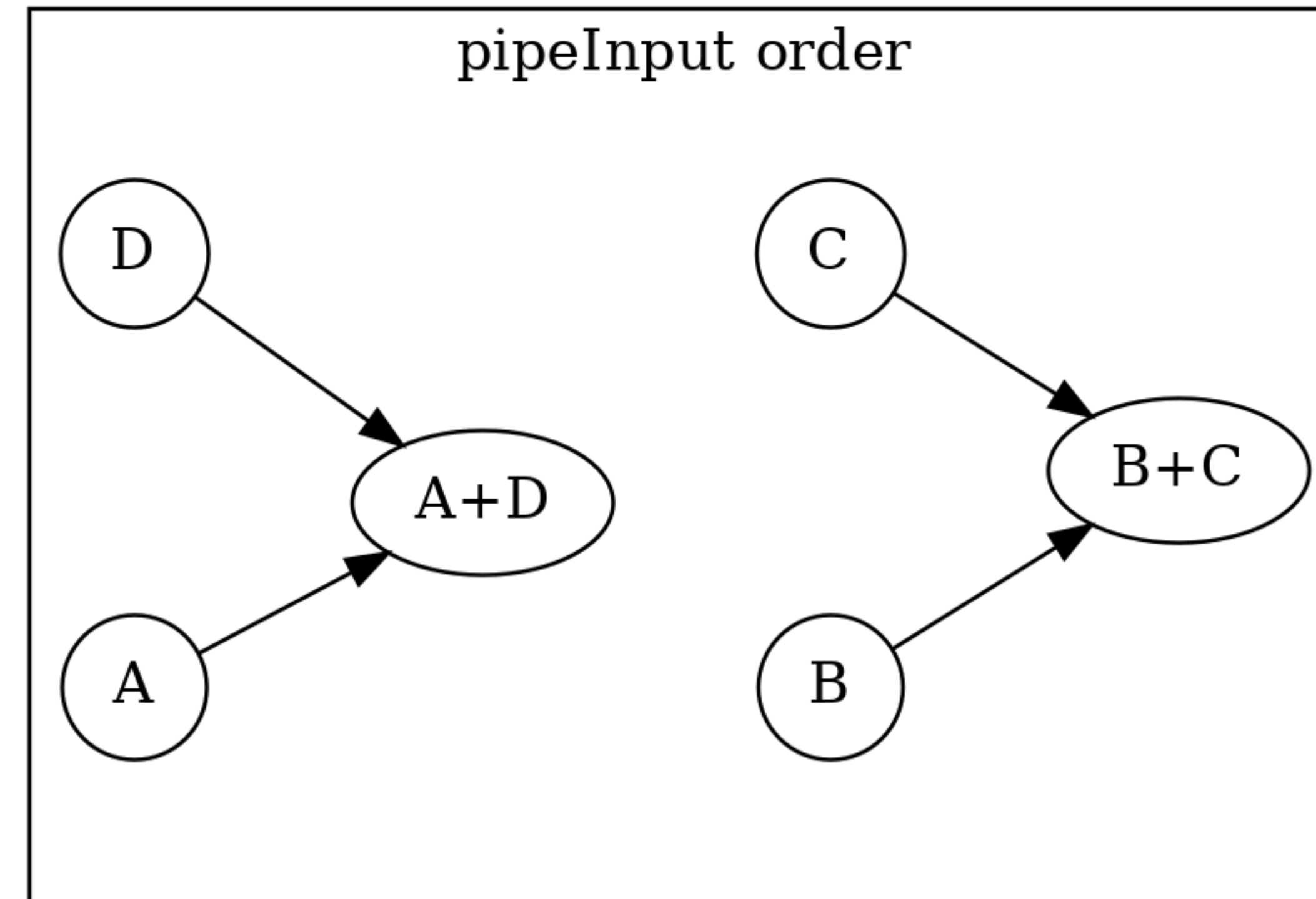# TopologyTestDriver: "Fused" subtopologies
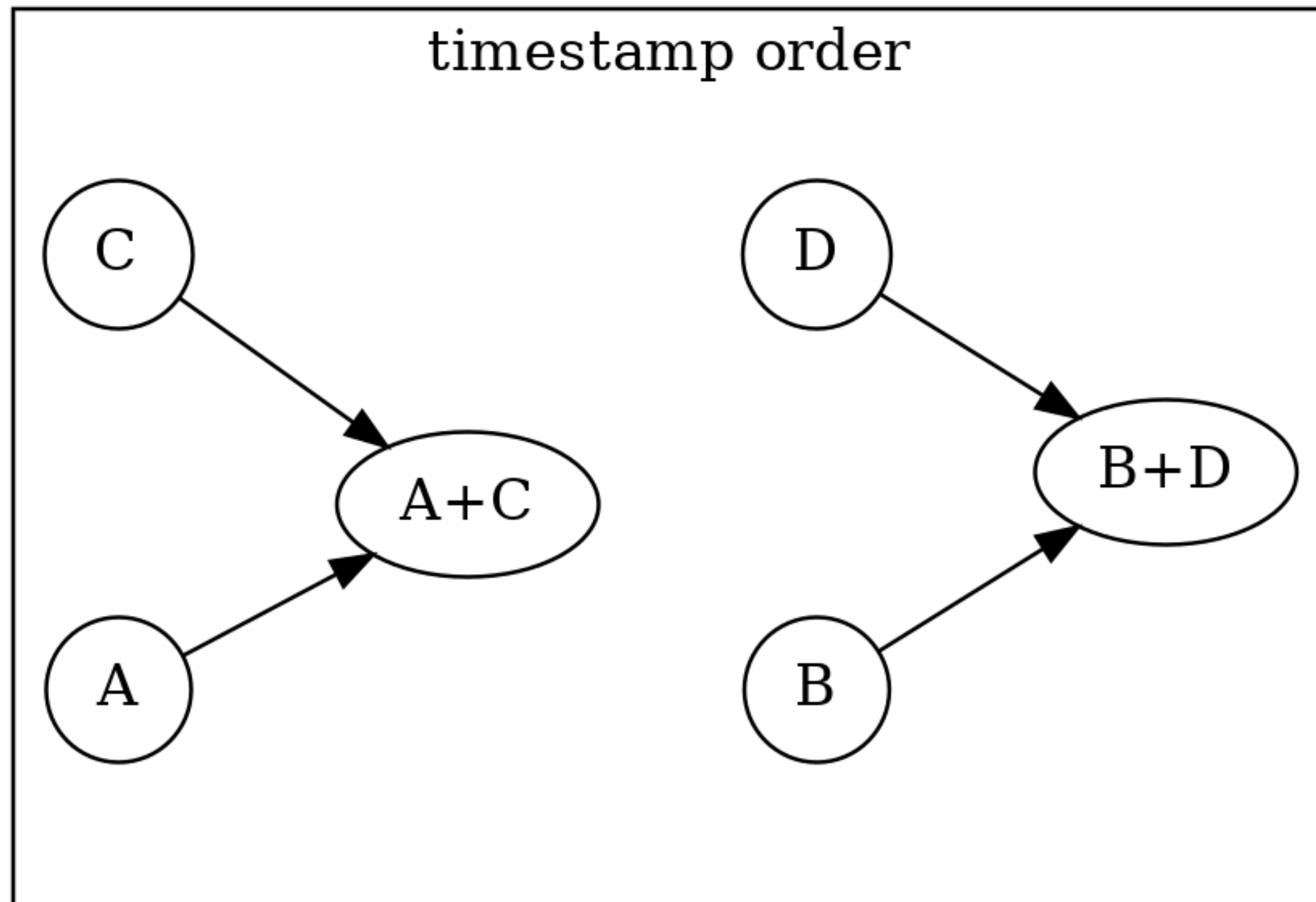
## ToplogyTestDriver



## Kafka Streams

# Timing

- stream-stream joins can behave differently (`pipeInput` order vs. timestamp order)
- logic that depends on stream time (such as `suppress`) can behave differently

# Should we trust StackOverflow?

# Using Transformer

# Using Transformer

A → B → B → ValueTransformerWithKey → A

ValueTransformerWithKey →

Cache

A | A

→ Local Store

# Using Transformer

A → B → **ValueTransformerWithKey** → B → A

ValueTransformerWithKey → **Cache**

| Cache | |
|---|---|
| A | A |
| B | B |

Cache → **Local Store**

# Using Transformer

# Using Transformer

# Let's run tests on real Kafka!

- EmbeddedKafka
- TestContainers

# EmbeddedKafka vs TestContainers

## EmbeddedKafka

- Pro:
  - Just pull in a dependency
- Contra:
  - Pulls in Scala
  - Runs in the same JVM

## TestContainers

- Pro
  - Runs Kafka isolated in Docker
  - Not only for Kafka testing
- Contra
  - Needs Docker
  - Requires some time for the first start

## Demo

- Writing TestContainers test
  - An easy part: pushing messages to Kafka
  - A not so easy part: how do we check the output?

# Demo

- Deduplication: the correct implementation
- Now the test is green, but takes 5 seconds!

# Does it have to be so slow?

```java
List actual = new ArrayList<>();

while (true) {
  ConsumerRecords<String, String> records =
    KafkaTestUtils.getRecords(consumer, 5000 /* timeout in ms */);
  if (records.isEmpty()) break;
  for (ConsumerRecord<String, String> rec : records) {
    actual.add(rec.value());
  }
}

assertEquals(List.of("A", "B"), actual);
```

# Awaitility

```
Awaitility.await().atMost(10, SECONDS).until(
                () -> List.of("A", "B").equals(actual));
```

# Awaitility

```
Awaitility.await().atMost(10, SECONDS).until(
              () -> List.of("A", "B").equals(actual));
```

# Things we must keep in mind

- Cooperative termination
- Thread-safe data structure

# Demo

- Green test runs faster

# Will any extra messages appear?

- We can wait for extra 5 seconds (bad choice)
- We can put a 'marker record' at the end of the input and wait for it to appear in the output (not always possible)

# Summary

# Summary

- **Both** `TopologyTestDriver` and integration tests are needed

# Summary

- **Both** `TopologyTestDriver` and integration tests are needed
- Write unit tests with `TopologyTestDriver`. When it fails to surface the problem, use integration tests.

# Summary

- **Both** `TopologyTestDriver` and integration tests are needed
- Write unit tests with `TopologyTestDriver`. When it fails to surface the problem, use integration tests.
- Know the limitations of `TopologyTestDriver`.

# Summary

- **Both** `TopologyTestDriver` and integration tests are needed
- Write unit tests with `TopologyTestDriver`. When it fails to surface the problem, use integration tests.
- Know the limitations of `TopologyTestDriver`.
- Understand the difficulties and limitations of asynchronous testing.

Страницы / Index / Kafka Improvement Proposals

# KIP-655: Windowed Distinct Operation for Kafka Streams API

Создатель Ivan Ponomarev, отредактировано авг 24, 2020

- Status
- Motivation
- Public Interfaces
- Proposed Changes
- Compatibility, Deprecation, and Migration Plan
- Rejected Alternatives

## Status

**Current state**: Under Discussion

**Discussion thread**: *here*

**JIRA**: ↑ KAFKA-10369 - Introduce Distinct operation in KStream  **IN PROGRESS**

**Pull request**: PR-9210

Please keep the discussion on the mailing list rather than commenting on the wiki (wiki discussions get unwieldy fast).

# Useful links

- Confluent blog: Testing Kafka Streams – A Deep Dive
- **pro.kafka**: Russian Kafka chat in Telegram: https://t.me/proKafka
- Confluent community Slack: https://cnfl.io/slack

# Thank you!

**Ivan Ponomarev**

- ✉ iponomarev@curs.ru
- 🐦 @inponomarev
- ⭘ inponomarev



**John Roesler**

- ✉ john@confluent.io
- ✉ vvcephei@apache.org
- ⭘ vvcephei