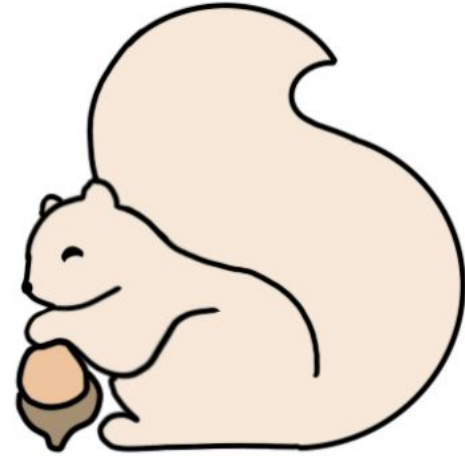


Stream Processing With Apache Flink



Big Data Demystified 2024

About the presenter



Dunitth Dhanushka

Senior Developer Advocate, Redpanda Data

- Event streaming, real-time analytics, and stream processing enthusiast
- Frequent blogger, speaker, and an educator



twitter.com/dunitthd



medium.com/event-driven-utopia



linkedin.com/in/dunitthd/

Agenda

1. An introduction to stream processing and Apache Flink
2. Deploying a Flink cluster
3. Generating a Flink application project for Java
4. Stateless operations
5. Q & A

Code is in GitHub!

<https://github.com/dunithd/kafka-meetup-frankfurt>

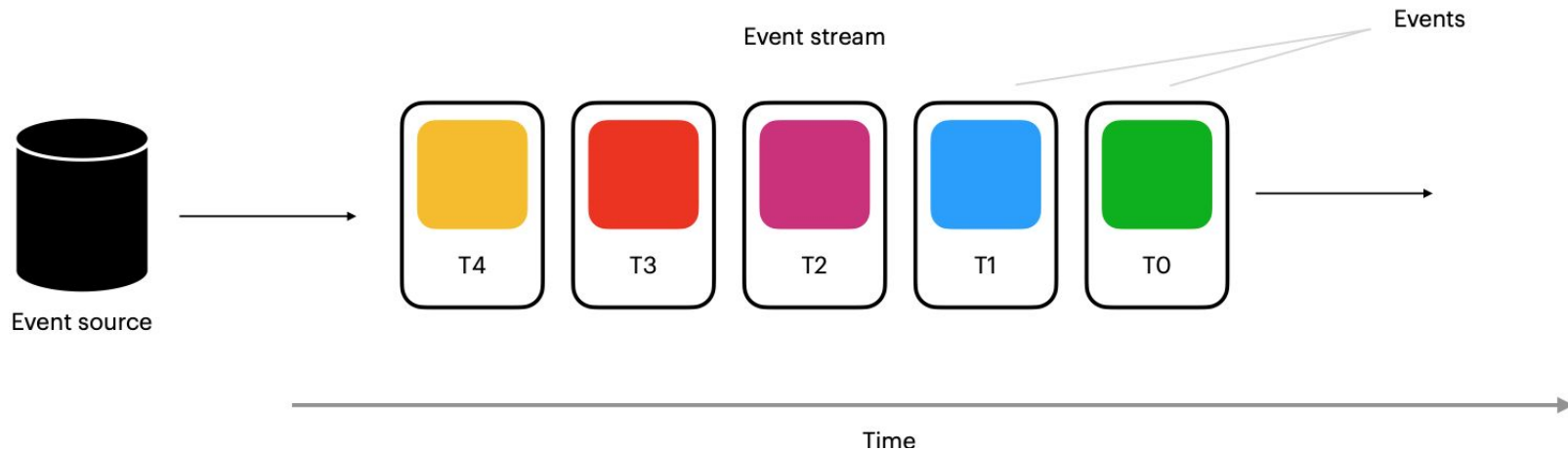


What is Stream Processing?

Event streams

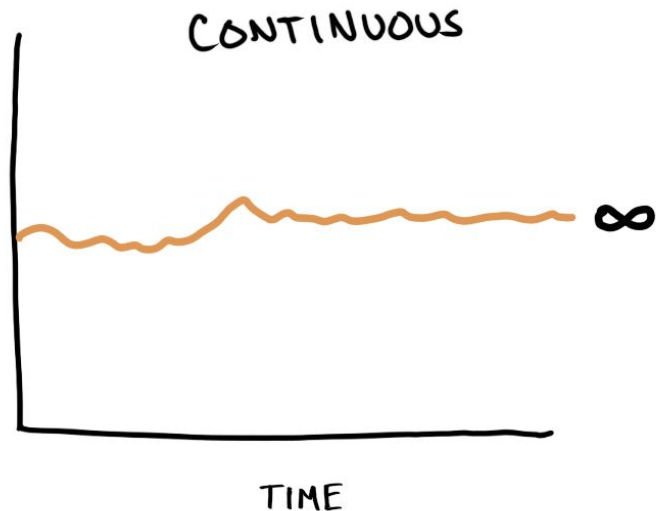
A data stream consists of a series of data points **ordered in time**.

Each data point represents an “**event**” or a change in the state of the business.



Stream processing

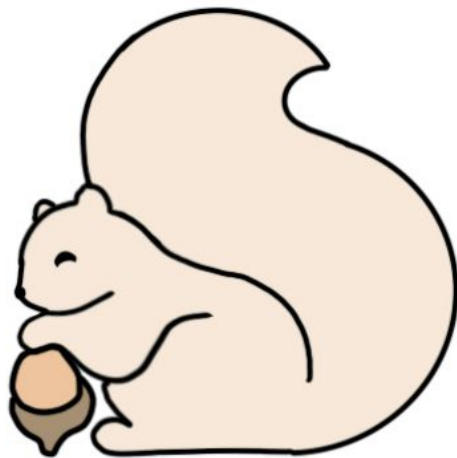
- Enriching, transforming, and reacting to data **continuously**
- Great for time-sensitive business use cases
- Jobs run indefinitely



What is Flink?

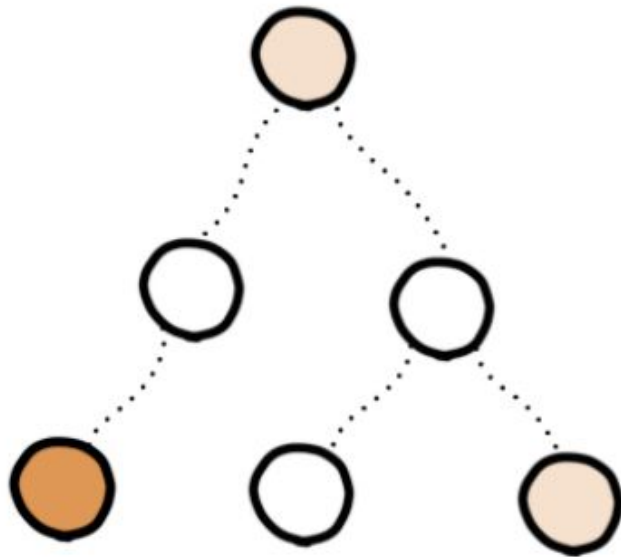
What is Apache Flink?

- Data processing framework
- Open-source project
- Very active community
- Battle-tested



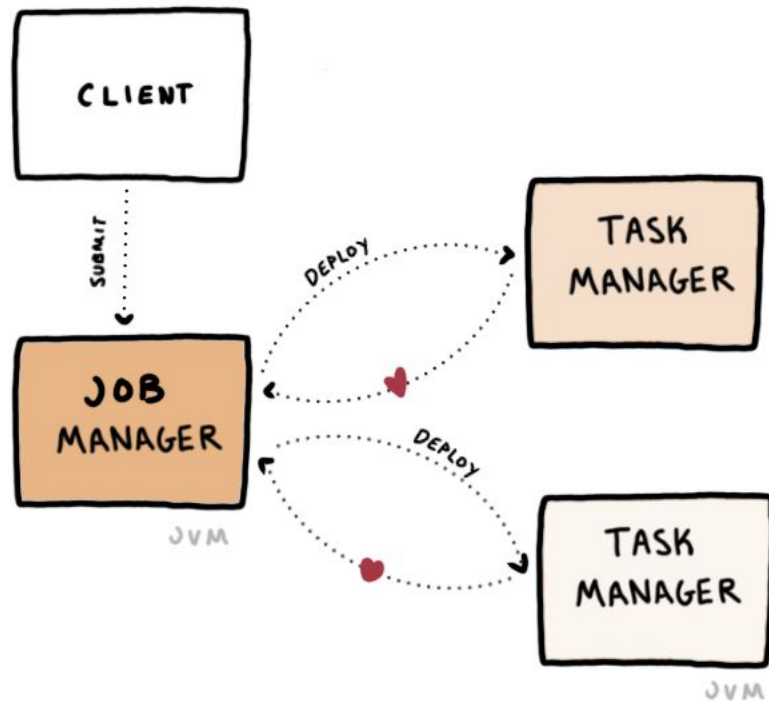
Dataflow programming

- Dataflow graphs, or *Job graphs*
- Nodes represent **operators** (functions)
- Edges represent the flow of data between operators
- Logical representation
- Can be **parallelized**



Components

- Clients
- Job manager
- Task manager



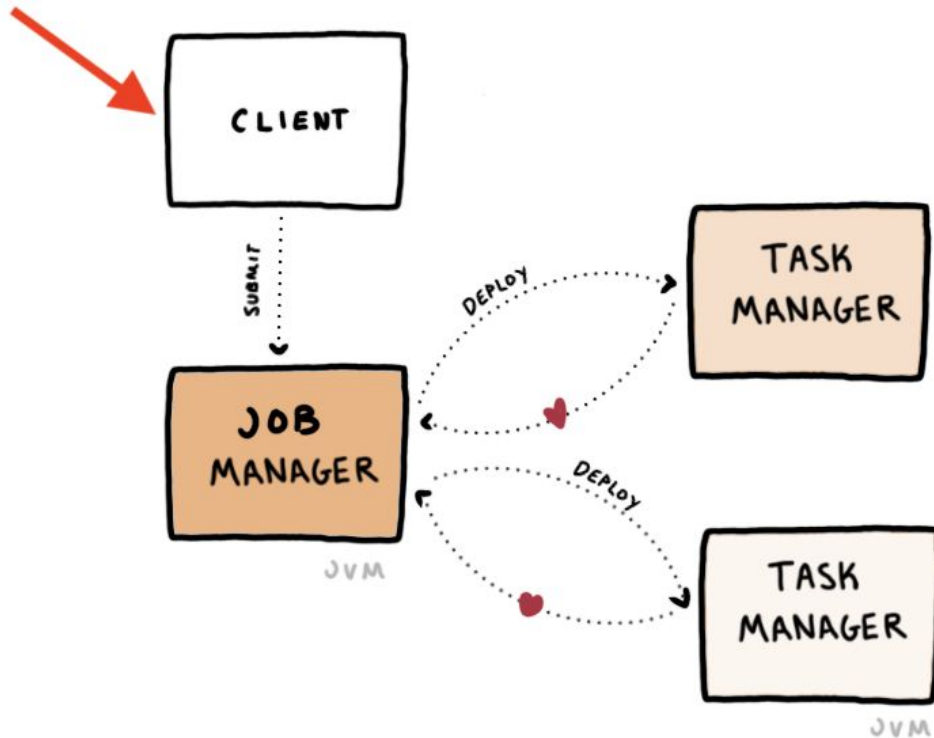
Clients

- SQL client
- PyFlink
- REST client

`POST jobs/`

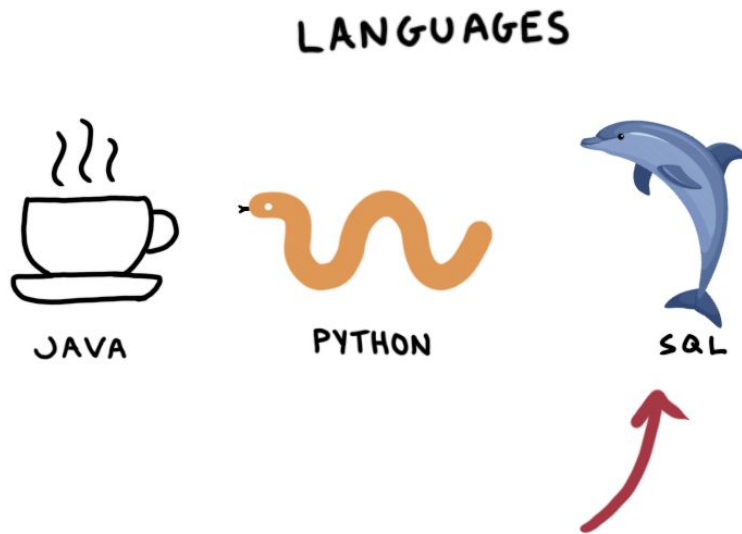
- CLI

`./bin/flink run`



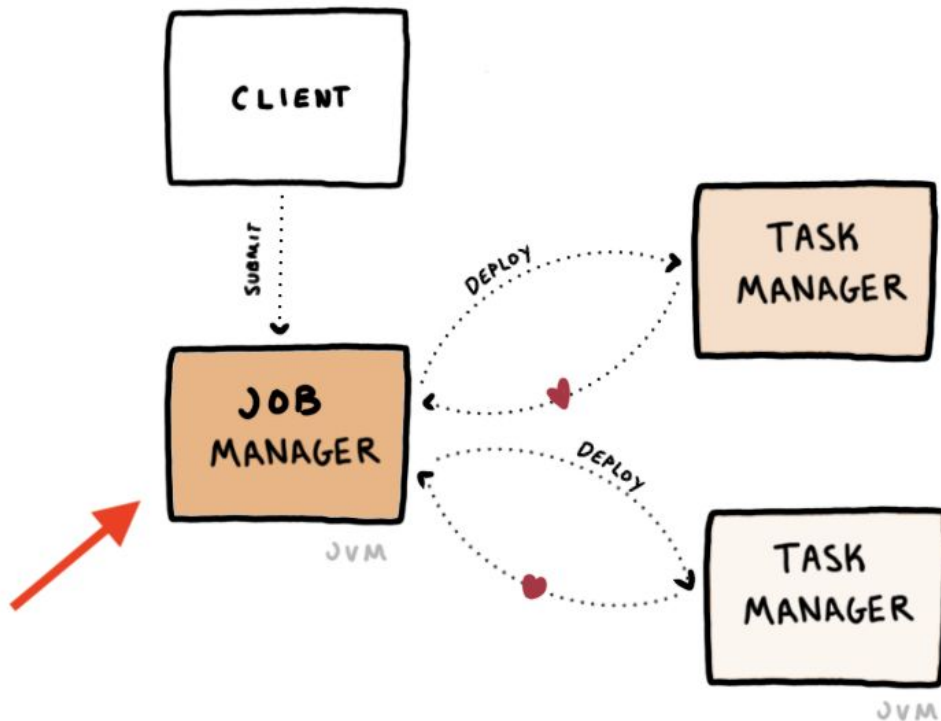
Languages

- JVM
- PyFlink
- Flink SQL



Job manager

- Coordinator
- Initiates checkpoints
- UI
- Different deployment strategies and cluster modes
- HA optional



Demo

Redpanda

orders-raw

orders



Flink Application

order-processing-job

What is Redpanda?

Redpanda is a Kafka API compatible streaming data platform

- Written in C++
- Thread-per-core architecture
- Designed for modern hardware

Simple to deploy, use and manage



Single binary



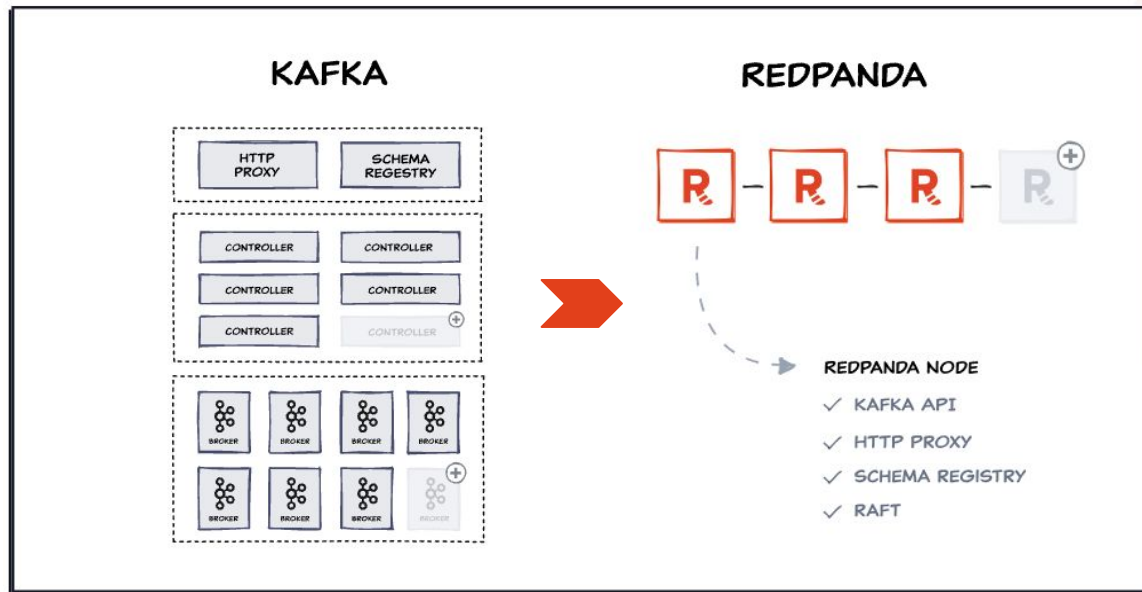
Kafka-compatible APIs



Easy Day 2 Ops



Dev-friendly interface



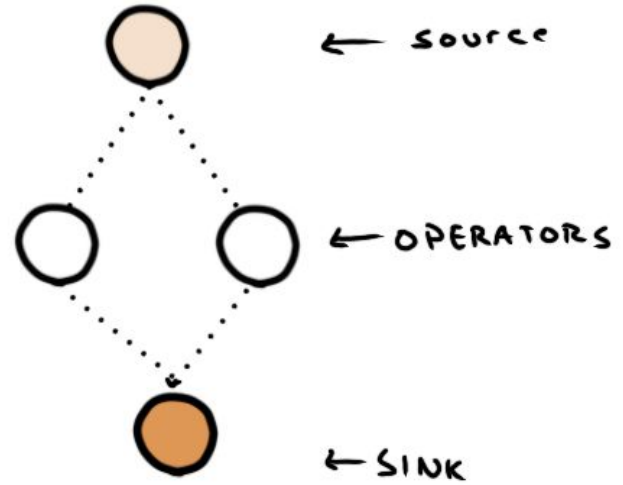
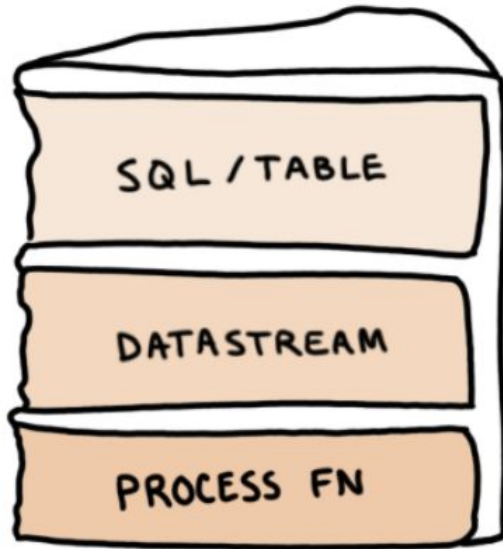
What you need?

- Docker Compose (Docker Desktop installed locally)

Developing Flink Applications

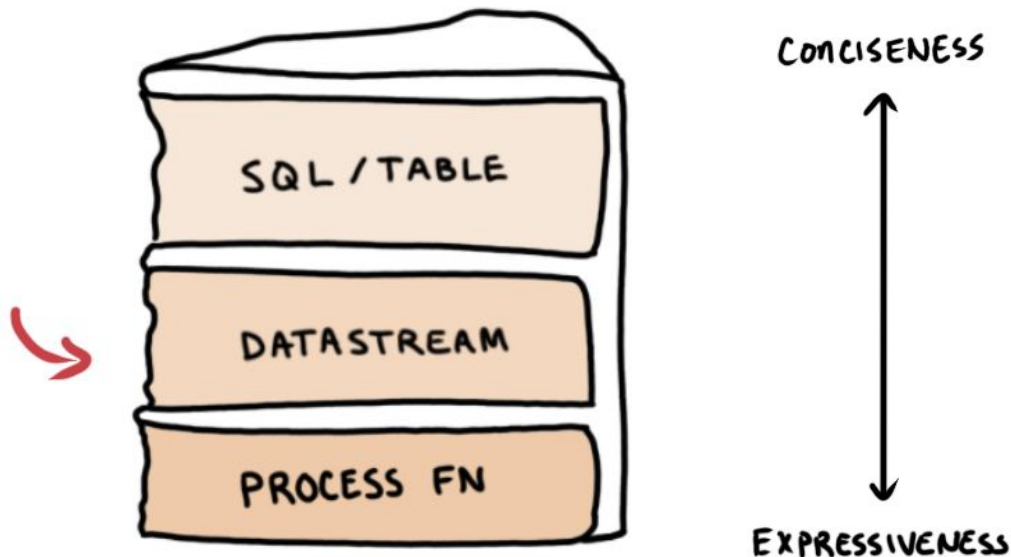
With Java

APIS



Lower levels

- More expressive
- Imperative approach
- Use for more complex use cases



DataStream API

```
StreamExecutionEnvironment env=  
    StreamExecutionEnvironment.getExecutionEnvironment();  
  
DataStream<Tuple2<String, Integer>> dataStream = env  
    .socketTextStream("localhost", 9999)  
    .flatMap(new Splitter())  
    .keyBy(value -> value.f0)  
    .window(TumblingProcessingTimeWindows.of(Time.seconds(5)))  
    .sum(1);  
  
dataStream.print();  
  
env.execute("Window WordCount");
```


Higher levels

- More concise
- Declarative approach
- Use when lower-level state / time access isn't needed

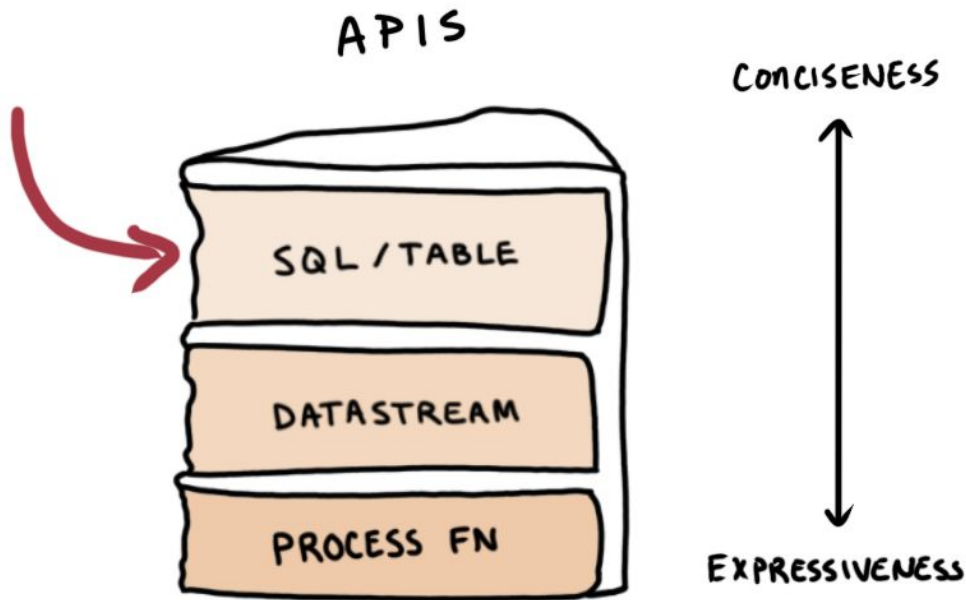


Table API

```
EnvironmentSettings settings = ...

TableEnvironment tEnv = TableEnvironment.create(settings);

// specify table program

Table orders = tEnv.from("Orders"); // schema (a, b, c, rowtime)

Table counts = orders

    .groupBy($"a")

    .select($"a", $"b".count().as("cnt"));

counts.execute().print();
```

Choosing an API

Table API

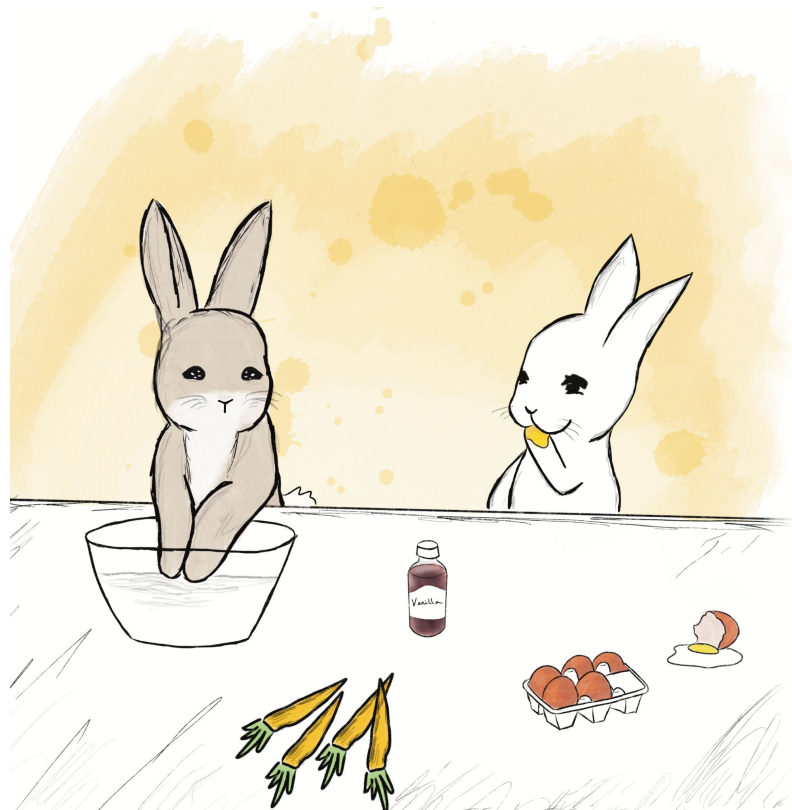
- If you want a **declarative** interface or closer alignment with **SQL**
- You don't need low-level control of **state** or time-based operations

DataStream API

- If you need more fine-grained control of streaming primitives (state, time)
- Imperative interface
- You want to work with a more **raw** stream of events
- Custom business logic

Mixing APIs

- Mixing is possible
- Example:
 - Table API for cleansing, filtering, or re-shaping input data
 - DataStream API for the primary application logic



Demo

What you need?

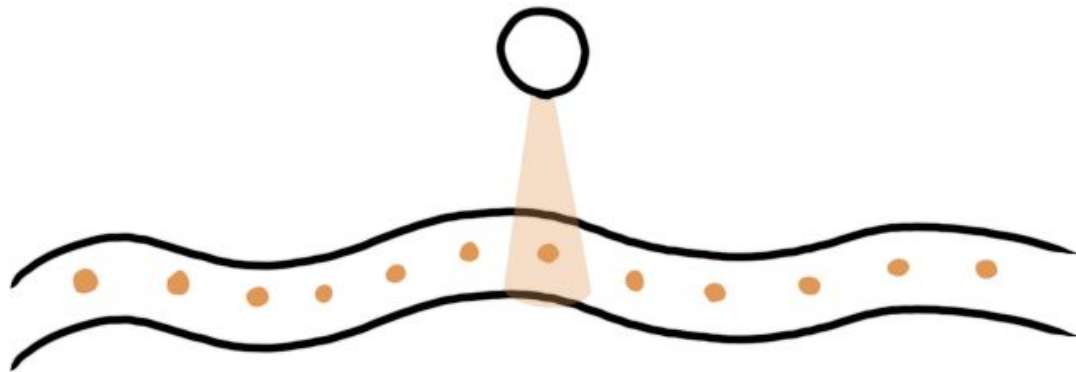
- Java 11 or higher
- Maven 3+
- Internet connection
- Terminal
- IDE - IntelliJ IDEA or Eclipse

```
mvn archetype:generate \  
-DarchetypeGroupId=org.apache.flink \  
-DarchetypeArtifactId=flink-quickstart-java \  
-DarchetypeVersion=1.16.2 \  
-DgroupId=kafka.meetup.flink \  
-DartifactId=flink-quickstart-project \  
-DinteractiveMode=false
```

Stateless Operations

Stateless operators

observe events without historical context



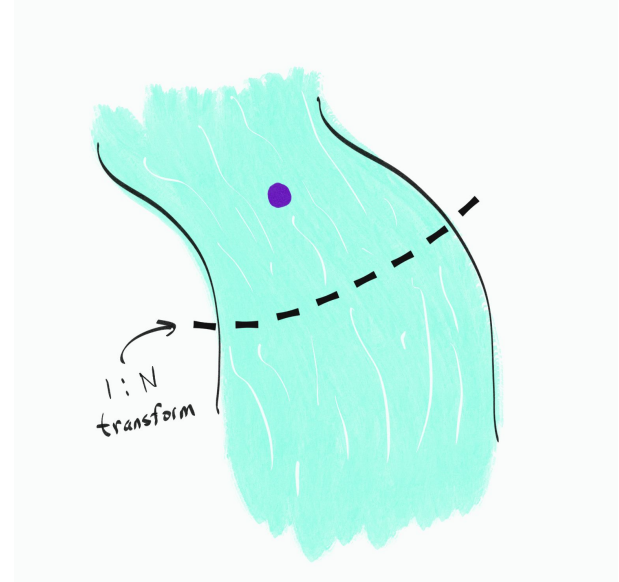
Stateless applications

- No memory
- Lightweight storage requirements
- Simple to deploy



Examples of stateless operations

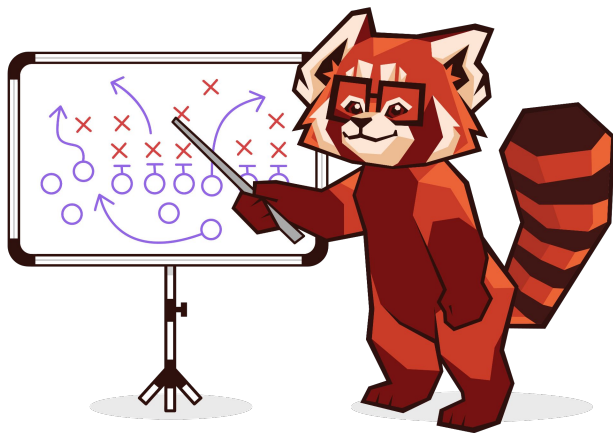
- Filtering data
- Single record data transformations
- Scrubbing records of sensitive information



Redpanda University

Free, self-paced online learning

<https://university.redpanda.com>



- Learn the fundamentals of data streaming and Redpanda
- Install Redpanda and use the rpk CLI to configure it
- Create producers and consumers in Java, Python and NodeJS
- Sign up today for free!

Thanks for joining!

Let's keep in touch

 @redpandadata

 redpanda-data

 redpanda-data

 hello@redpanda.com

