

The Streaming Mindset

... what, why, how?

Marta Paes (@morsapaes)

Developer Advocate

About Ververica



Original Creators of
Apache Flink®



Enterprise Stream Processing
With Ververica Platform



Part of
Alibaba Group



Working in DevRel



J. Doe • 00:00

I have 1.5 years of experience in writing pyspark batch jobs, now I wanted to get my hands dirty in real time processing

Can you please guide me how should I proceed



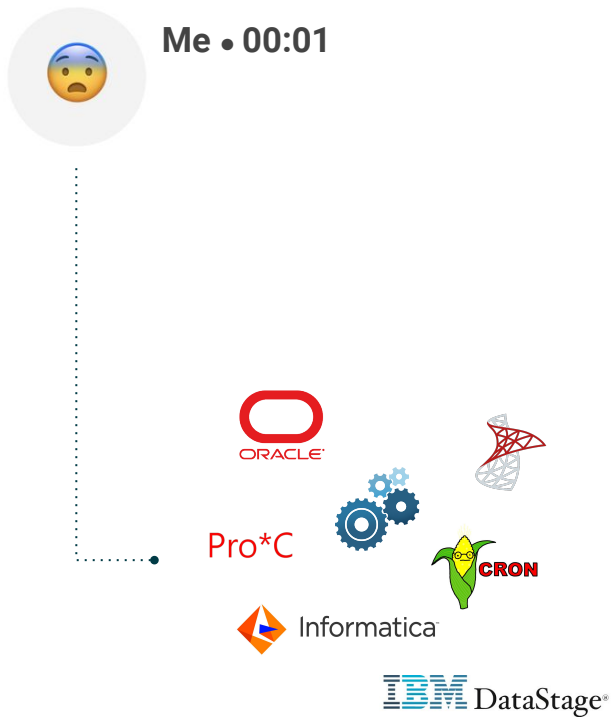
Working in DevRel



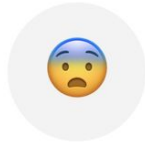
Me • 00:01



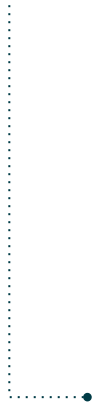
Working in DevRel



Working in DevRel



Me • 00:01



ORACLE
Pro*C
Informatica
IBM DataStage®
CRON



hadoop MapReduce
STORM
APACHE Spark



Working in DevRel



Me • 00:01



Pro*C



Informatica



CRON



IBM DataStage®

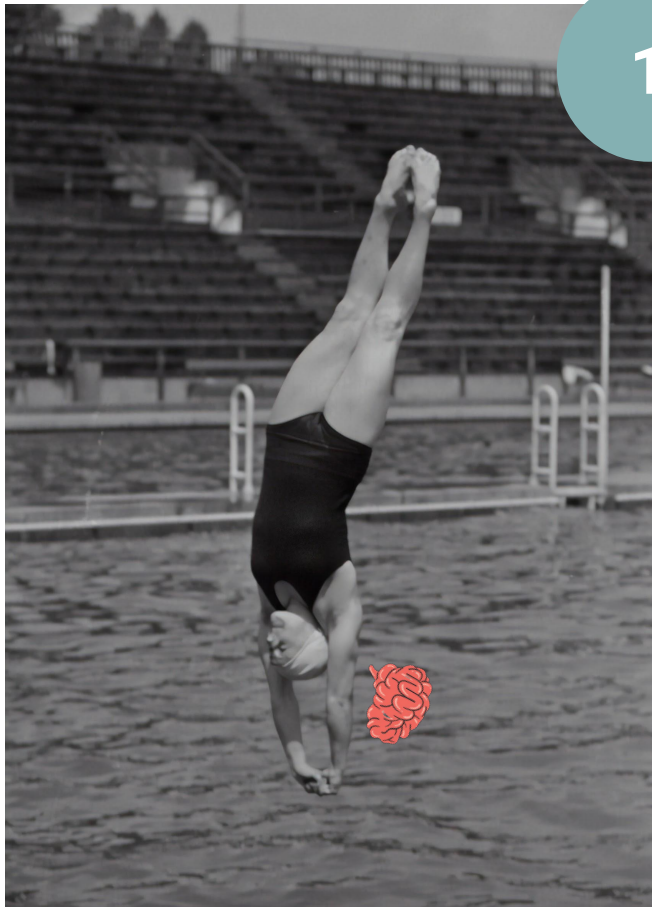


Where do you start?

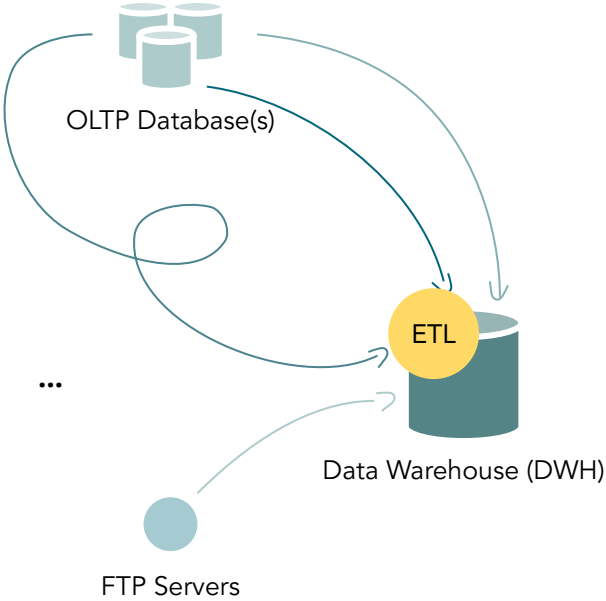
1

Go Headfirst

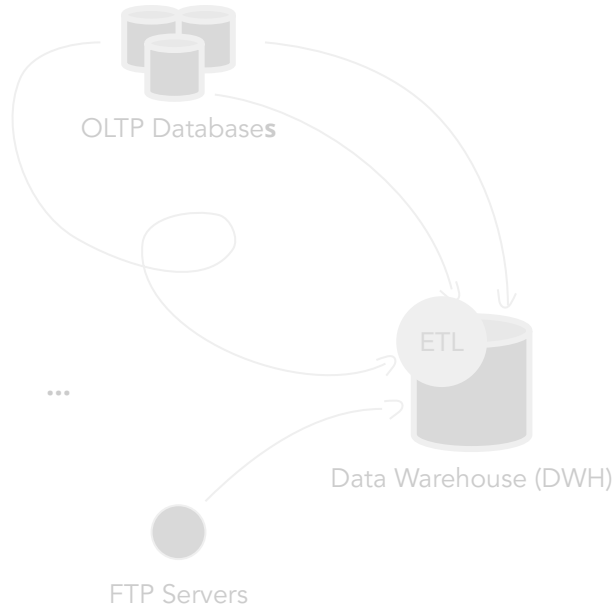
- Stream Processing 101



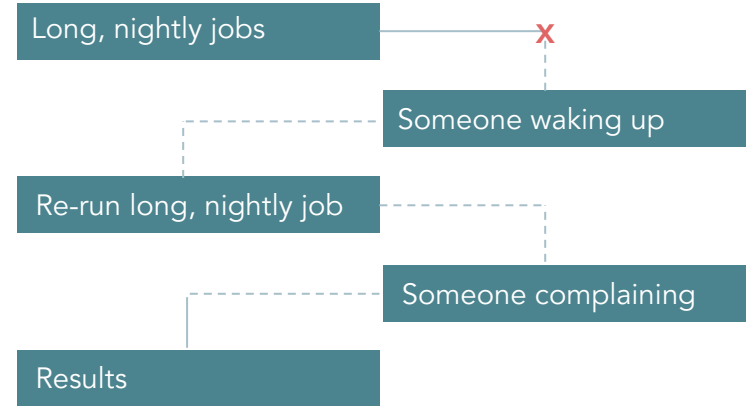
Analytics...Not that Long Ago



Analytics...Not that Long Ago



The quest for data...



But in the end...

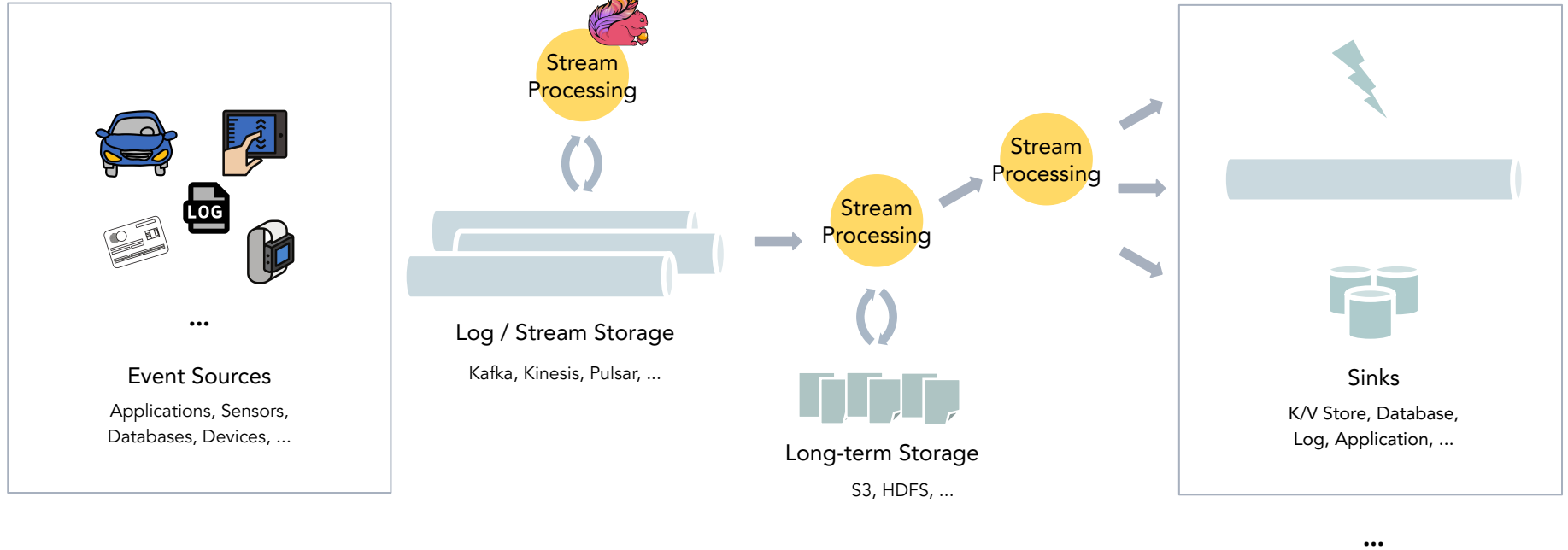
- Most source data is continuously produced
- Not everyone can wait for yesterday's data
- Most logic is not changing that frequently



Everything is a Stream

Everything is a Stream

Your static data records become **events** that are **continuously produced** and should be **continuously processed**.



Stream Processing 101

Batch Processing

query/logic changes fast

data changes slowly

E.g: Ad-hoc queries, data exploration, ML model training

Continuous Streaming

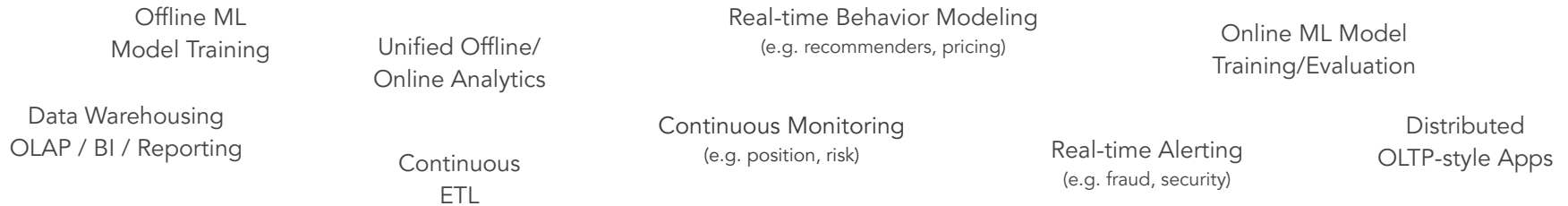
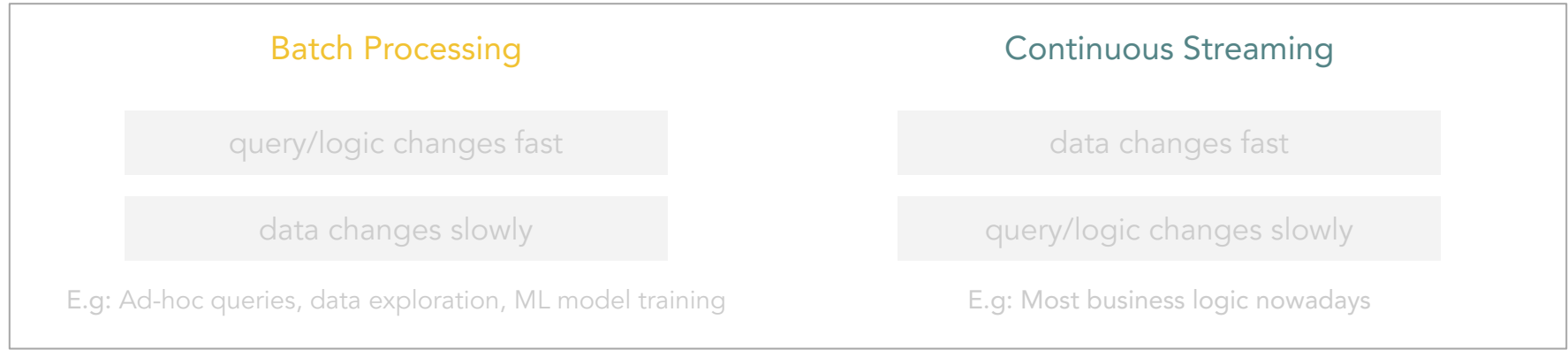
data changes fast

query/logic changes slowly

E.g: Most business logic nowadays



Stream Processing 101



Stream Processing Use Cases

Examples

NETFLIX

[Large-scale Data Pipelines](#)

ING 

[ML-Based Fraud Detection](#)

aws 

[Service Monitoring & Anomaly Detection](#)



Stream Processing Use Cases

Examples

NETFLIX

Large-scale Data Pipelines

ING 

ML-Based Fraud Detection

aws 

Service Monitoring & Anomaly Detection



Unified Online/Offline Model Training

Uber

E2E Streaming Analytics Pipelines

criteo. 

ML Feature Generation



2

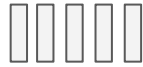
Bridge Concepts

- Bounded vs. Unbounded data
- Event time vs. Processing time
- Fault tolerance



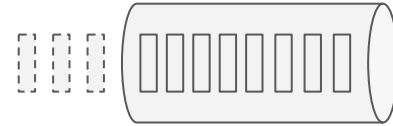
Bounded vs. Unbounded Data

Batch Processing



- Data "at rest"
- Hard boundaries (e.g. process 1 day of data)

Continuous Streaming



- Data "on the fly"
- Ever-growing, infinite data set



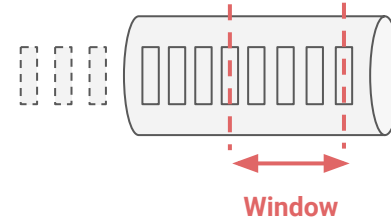
Bounded vs. Unbounded Data

Batch Processing



- Data "at rest"
- Hard boundaries (e.g. process 1 day of data)

Continuous Streaming



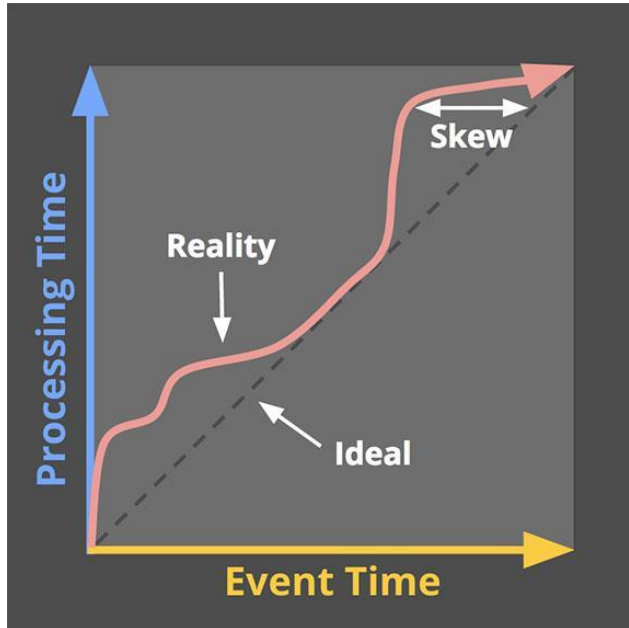
- Data "on the fly"
- Ever-growing, infinite data set



Windows split the stream into buckets of finite size, over which you can apply computations



Event Time vs. Processing Time



Event time

- Deterministic results
- Handle out-of-order or late events
- Trade-off result completeness/correctness and latency

Processing time

- Non-deterministic results
- Best performance and lowest latency
- Speed > completeness/correctness



Fault Tolerance

Batch Processing

pipelines run on a fixed schedule

Continuous Streaming

Long-running pipelines



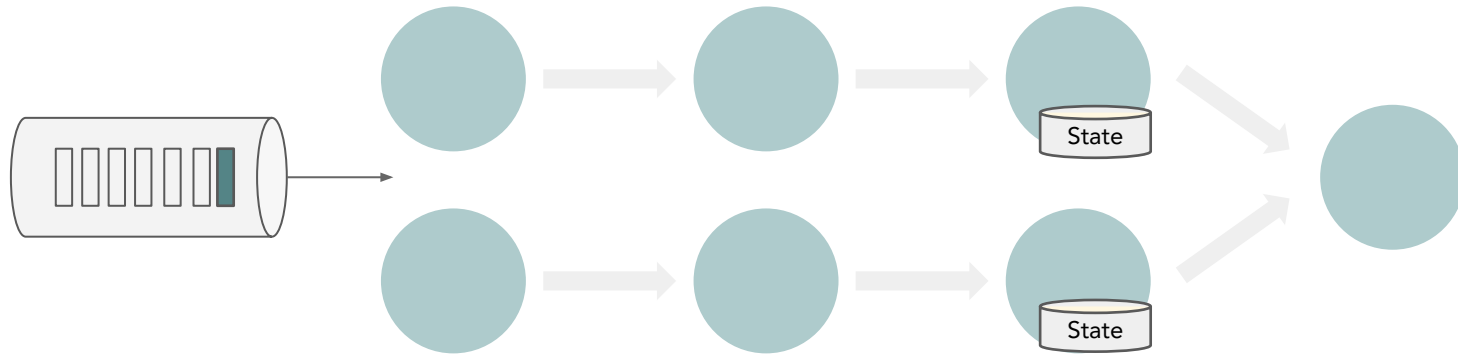
Fault Tolerance

Batch Processing

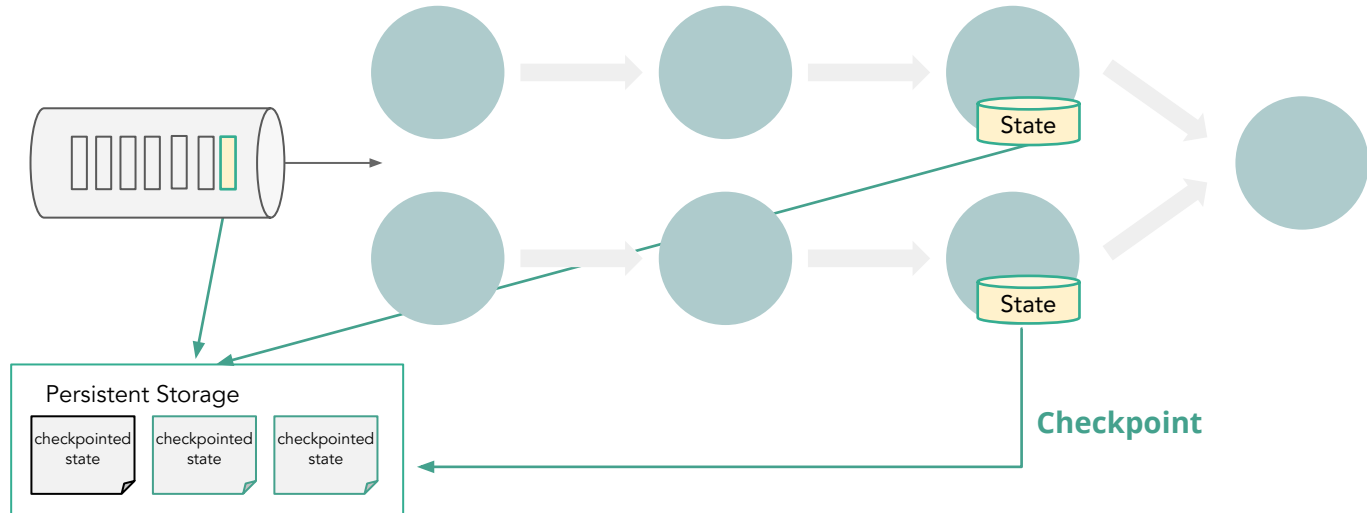
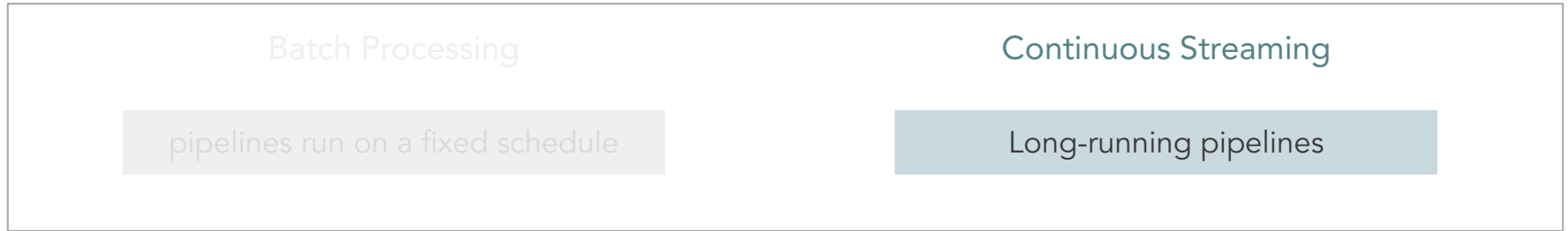
pipelines run on a fixed schedule

Continuous Streaming

Long-running pipelines



Fault Tolerance



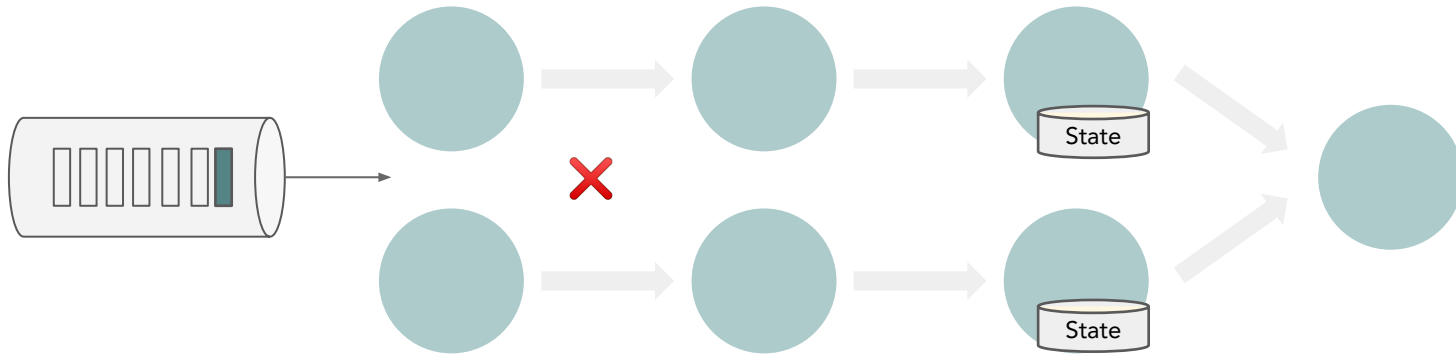
Fault Tolerance

Batch Processing

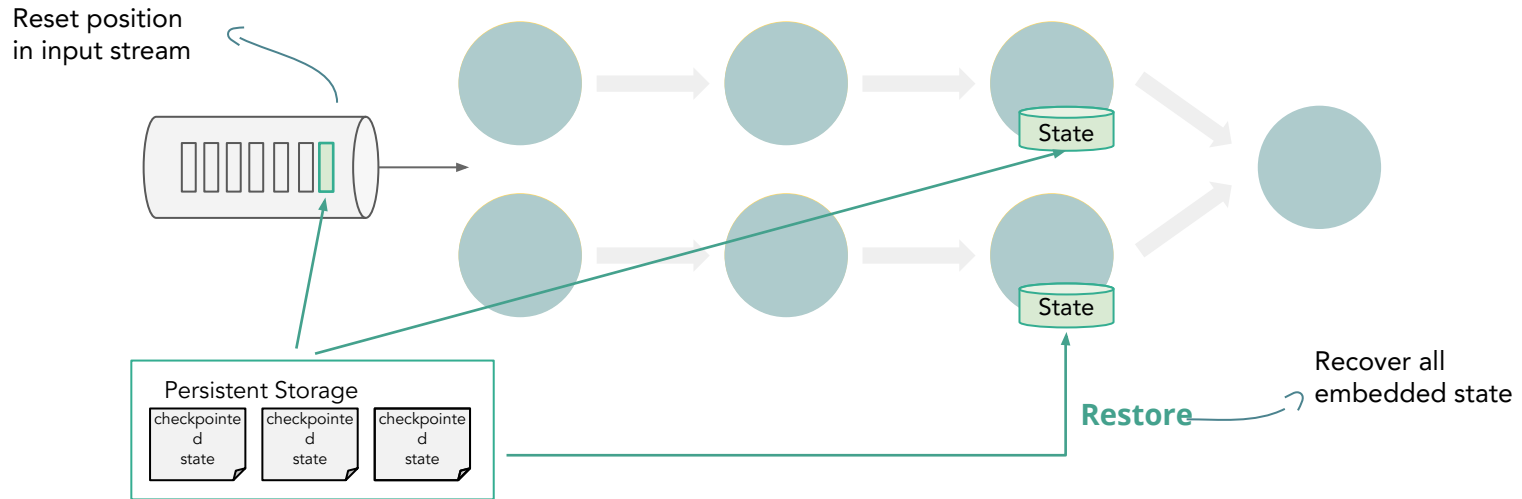
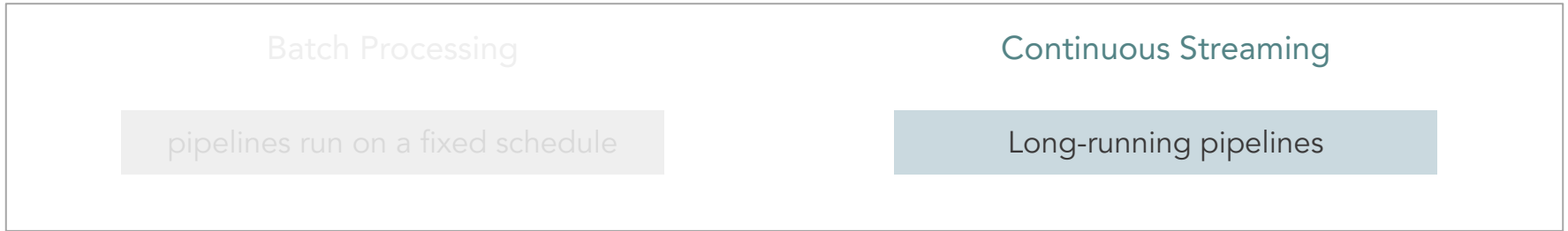
pipelines run on a fixed schedule

Continuous Streaming

Long-running pipelines



Fault Tolerance



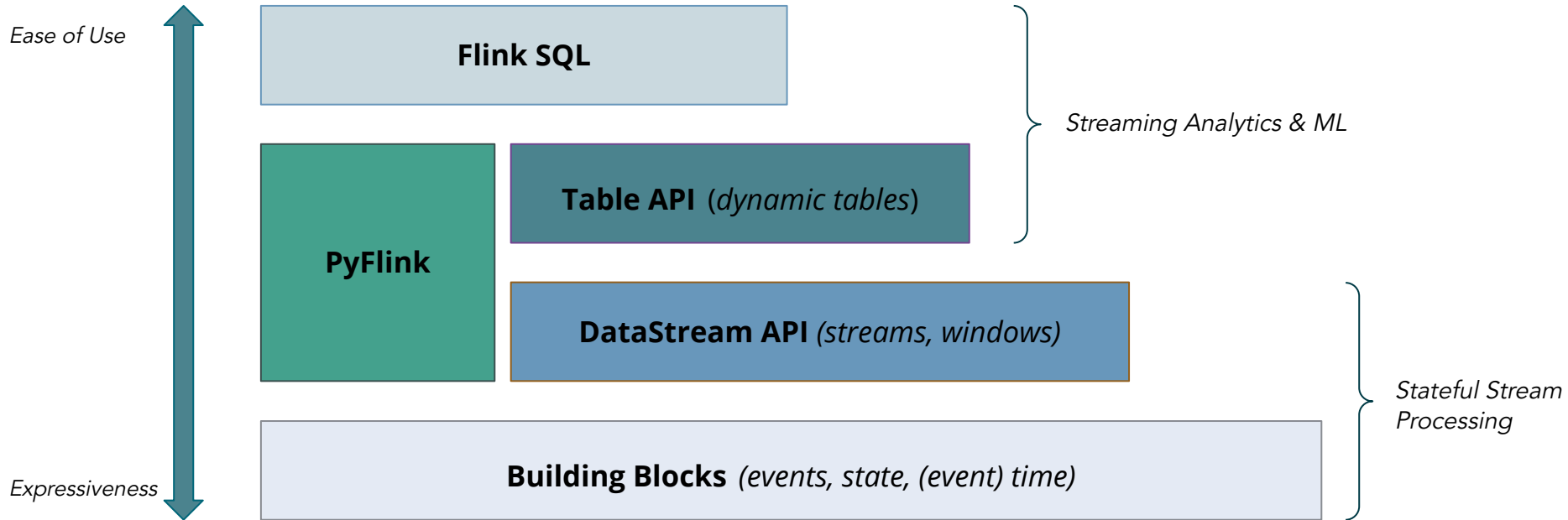


3

Pick a Flavour & Build

The Flink API Stack

Layered, with different tradeoffs for **expressiveness** and **ease of use**. You can mix and match all the APIs!



How to Get Hands-On?

Start with whatever language and/or abstractions are more familiar to you!

Java/Scala



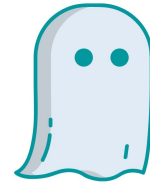
- [Self-paced Training Course](#)
- [DataStream API Walkthrough](#)

SQL



- [Flink SQL Cookbook](#)
- [Table API Walkthrough](#)

Python

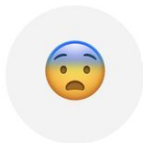


- [PyFlink Walkthrough](#)
- [Zeppelin Notebooks](#)



Starting from the beginning

From being dumbfounded...



J. Doe • 00:00

I have 1.5 years of experience in writing pyspark batch jobs, now I wanted to get my hands dirty in real time processing

Can you please guide me how should I proceed



Me • 00:01



...to actually having a plan!



J. Doe • 00:00

I have 1.5 years of experience in writing pyspark batch jobs, now I wanted to get my hands dirty in real time processing

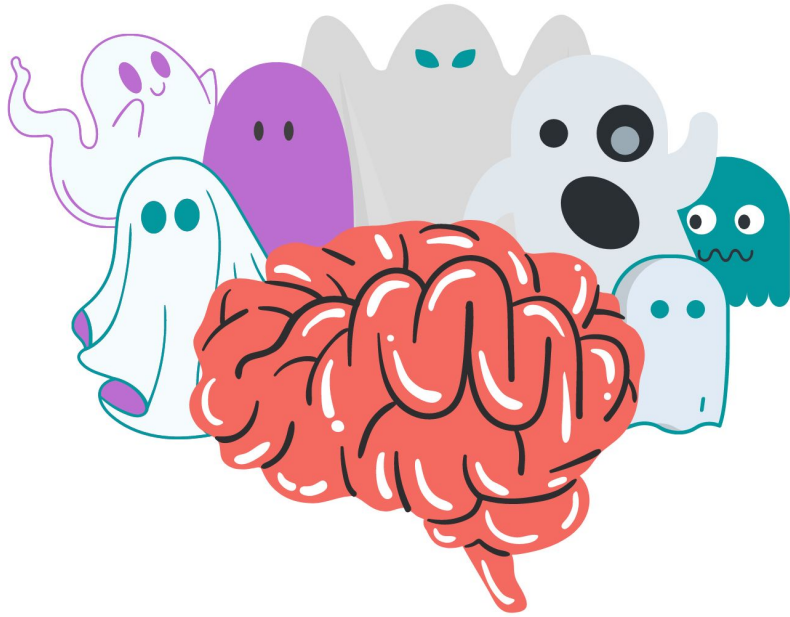
Can you please guide me how should I proceed



Me • 00:01

- ✓ Invest in learning the Stream Processing 101
- ✓ Take the time to understand how it differs from Batch Processing
- ✓ Start with something familiar and increase complexity gradually
- ✓ Ask questions!





Thank you, Bristech!

Follow me on Twitter: @morsapaes

Learn more about Flink: <https://flink.apache.org/>