# Databases and Stream Processing:

## A Future of Consolidation

Ben Stopford

Office of the CTO, Confluent

# Marc Andreessen:
# Software is Eating the World

# Weak Form
**Companies are USING MORE SOFTWARE**

# Strong Form
**Companies are BECOMING SOFTWARE**

# Loan Application Using Software

**1** BORROWER

**2** APPLICATION FORM

**3** CREDIT OFFICER

**4** RISK OFFICER

**5** LOAN OFFICER
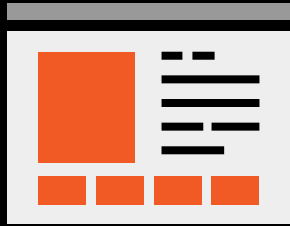
**6** APPROVE

DENY

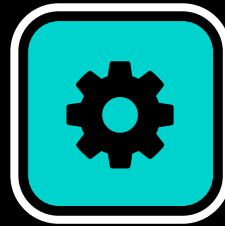# Loan Application in Software
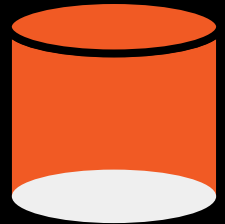
**Using** Software:
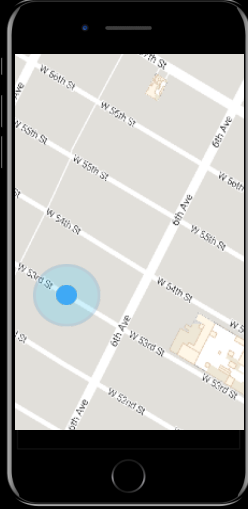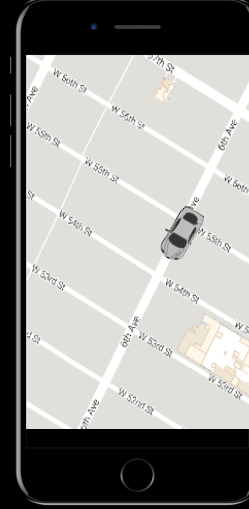Classic Three-Tier Architecture

USER → UI → SERVICE → DATABASE

# Becoming Software:
## Services Talking To Each Other With APIs



**SERVICE**　　　**SERVICE**　　　**SERVICE**　　　**SERVICE**

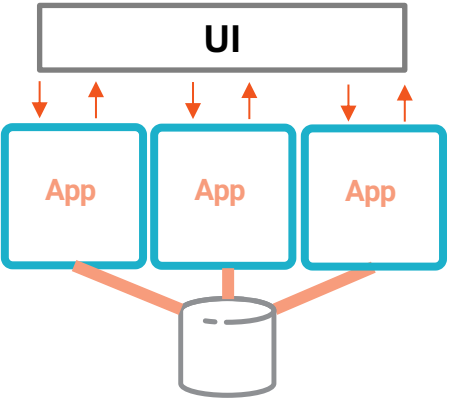# Evolution of software systems
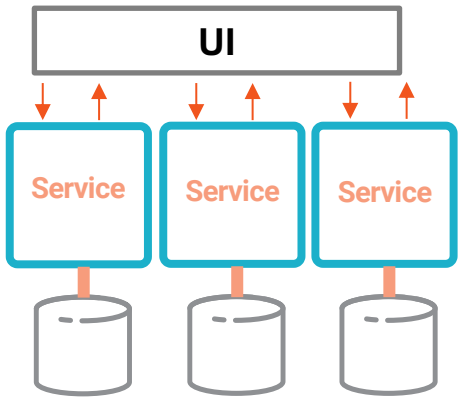
Monolith
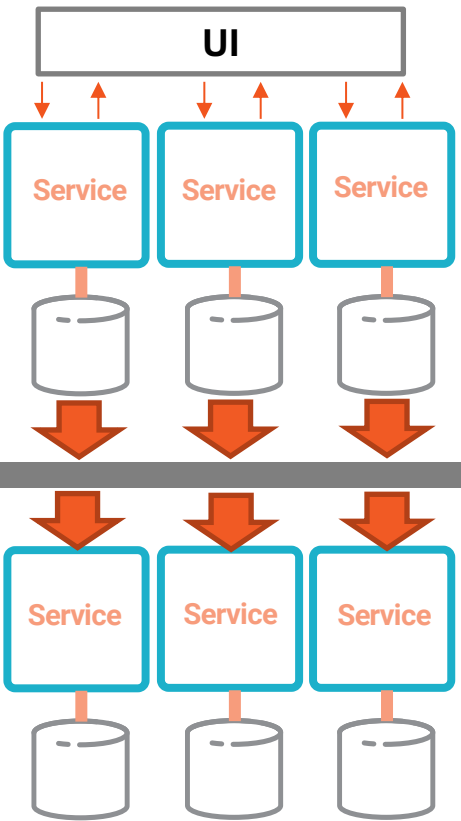
Distributed Monolith

Microservices

Event-Driven Microservices

THE USER OF THE SOFTWARE

IS MORE
SOFTWARE

What does this mean for databases?

We have hundreds
of databases…

**Databases are designed to help *you*!**

Unless there is a **user and UI waiting,** why should it be **synchronous**?

# The Alternative: Event Streams

Stream Processors are built for Asynchronicity

# Stream Processors have a different interaction model

## TRADITIONAL DATABASE

**Active Query**

**Passive Data**

SELECT *
FROM
DB_TABLE

DB Table

## EVENT STREAM PROCESSING

**Active Data**

**Passive Query**

CREATE TABLE AS
SELECT * FROM
EVENT_STREAM

Event Stream

# Streams or Tables?

# An Event
# records the fact that something happened

A good
was sold

An invoice
was issued

A payment
was made

A new customer
registered

**Events are state changes, they carry intent**

State:

Bob works at Google

Event:

Bob moved from Google to Amazon

# Streams
**record exactly what happened**

# Tables
**current state**



Where you have been   vs.   Where you are now

Payments you made   vs.   Your account balance

# Streams
## A sequence of moves

1. e4        e5
2. Nf3       Nc6
3. Bc4       Bc5
4. d3        Nf6
5. Nbd2

# Tables
## Position of each piece

**Streams = INSERT only**
**Immutable, append-only**

**Tables = INSERT, UPDATE, DELETE**
**Mutable, Primary Key**

**A stream can be considered as an immutable, append-only table**

Stream Processors Communicate Through Streams

INPUT STREAMS

SP

OUTPUT STREAMS

# But internally they use tables



| USER | PAYMENTS |
|------|----------|
| JAY | 42 |
| SUE | 18 |
| FRED | 65 |
| … | … |

Payments Stream

```
CREATE TABLE credit_scores AS
    SELECT user, updateScore(p.amount)…
```

| USER | CREDIT_SCORE |
|------|--------------|
| JAY | 695 |
| SUE | 430 |
| FRED | 710 |

Credit Score Table

Credit Score Stream

20

# Streams
## record history

# Tables
## represent state

projection
(Group By Key, SUM, COUNT)

Duality

table changes

*See Streams and Tables: Two Sides of the Same Coin, M. Sax et al., BIRTE '18

# Similar to a materialized view in a database



STREAM PROCESSOR

Payments Stream

Credit Score Table

| USER | CREDIT_SCORE |
|------|--------------|
| JAY | 695 |
| SUE | 430 |
| FRED | 710 |

Credit Score Stream

APP

- Asynchronous
- Push query

ACTIVE DATABASE

Payments Table

| USER | PAYMENTS |
|------|----------|
| JAY | 42 |
| SUE | 18 |
| FRED | 65 |
| ... | ... |

Credit Score Table

| USER | CREDIT_SCORE |
|------|--------------|
| JAY | 695 |
| SUE | 430 |
| FRED | 710 |

- Synchronous
- Pull  query

20

# Joins

# Joining a stream with a table

Orders

Lookup Customer

Customers → Table of Customers (with Primary Key)

# Joining two streams

Bob's
Order

Jill's
Order

Orders

Payments

Jill's
Payment

Bob's
Payment

`orders.join(payments)`

# Joining two streams

Bob's
Order

Jill's
Order

Jill's
Payment

Bob's
Payment

`orders.join(payments)`

# Joining two streams



Bob's Order

Jill's Order

Jill's Payment

Bob's Payment

`orders.join(payments)`

# Joining two streams

Bob's
Order

Jill's
Order

Jill's
Payment

Bob's
Payment

`orders.join(payments)`

# Joining two streams



Key-value store

Bob's Order

Jill's Order

Jill's Payment

Bob's Payment

# Joining two streams

Jill's Order

Key-value store

Bob's Order

Jill's Payment

Bob's Payment

# Joining two streams

Jill's
Order

Key-value store

Bob's
Order

Jill's
Payment

Bob's
Payment

# Joining two streams

Jill's Order

Bob's Order

Jill's Payment

Bob's Payment

K-V

K-V

# Joining two streams



Bob's Order

Jill's Order

Jill's Payment

Bob's Payment

# Joining two streams

# Joining two streams

Bob's Order

K-V

Jill's Order

Jill's Payment

K-V

Bob's Payment

# Joining two streams



Bob's Order

Jill's Order

Bob's Payment

Jill's Payment

# Streams represent history –> Cartesian Product

**Orders Stream**

| | |
|---|---|
| **101** | Boots |
| 200 | Hat |
| **101** | Boots2 |
| 105 | Pants |
| 200 | Hat2 |

**Payments Stream**

| | |
|---|---|
| 101 | $50 |
| 200 | $10 |
| 105 | $3 |
| 200 | $12 |
| 101 | $60 |

**Join Output (Stream)**

# Joining Streams to Streams

**Orders Stream**

| | |
|---|---|
| **101** | Boots |
| 200 | Hat |
| **101** | Boots2 |
| 105 | Pants |
| 200 | Hat2 |

**Payments Stream**

| | |
|---|---|
| 101 | $50 |
| 200 | $10 |
| 105 | $3 |
| 200 | $12 |
| 101 | $60 |

Use time window

**Join Output (Stream)**

# Tools for correlating recent events in time

# More advanced temporal functions



Orders

Page Visits

Session

Join Output
(Stream)

# Late and out-of-order data

Orders

Page Visits

Window 1

Window 2

**Join Output
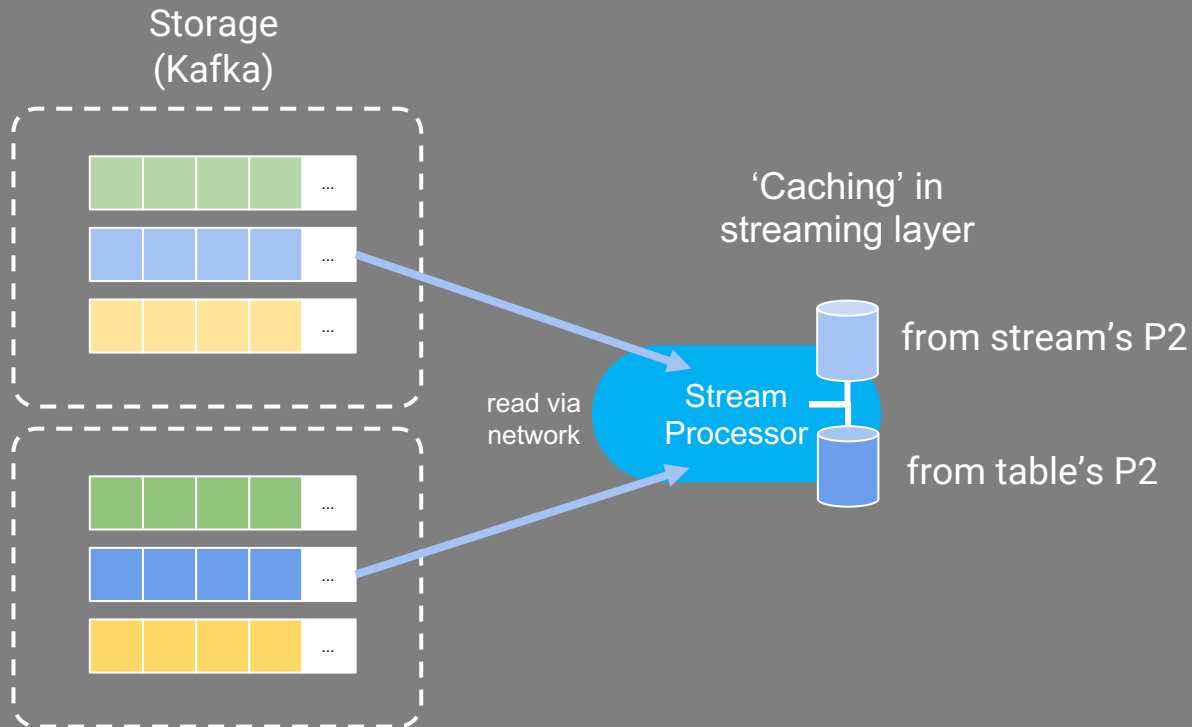(Stream)**

**Stream processors provide tools that handle asynchronicity, leverage time and focus on 'now'**

# Data Placement

# Layered storage model



Storage
(Kafka)

'Caching' in
streaming layer

read via
network

Stream
Processor

from stream's P2

from table's P2

# Partitioned Data (Fact-Fact joins)

Storage (Kafka)

**Partitioned
KTable / TABLE**



P1 ........ → SP 1 — 2 GB

P2 ........ → SP 2 — 3 GB

P3 ........ → SP 3 — 5 GB

P4 ........ → SP 4 — 2 GB

# Broadcast Data (Fact-Dimension Joins)

# Architecturally there are parallels e.g. Data Warehousing



FACTS  DIMS

ETL

REPORTING

# Interaction Model

# TRADITIONAL DATABASE

# EVENT STREAM PROCESSING

**Active Query**

**Passive Data**

**Active Data**

**Passive Query**

**SELECT * FROM DB_TABLE**

**CREATE TABLE AS SELECT * FROM EVENT_STREAM**

**DB Table**

**Event Stream**

# Databases are Pull Queries

# Stream Processors are Push Queries

Payments

Payments

What is Ben's credit score now?

**APP**

695

Ben's credit score is 670

Ben's credit score is 710

Ben's credit score is 695

...

**APP**

# Hybrid stream processors provide both interaction models



**APP**

**Payments Stream**

**ksqlDB**

| USER | CREDIT_SCORE |
|------|--------------|
| JAY | 695 |
| SUE | 430 |
| FRED | 710 |

**Summarize & Materialize Credit Scores**

**Query Credit Scores**

**Stream Credit Scores**

**APP**

# Unified Model For:

1. The Asynchronous and the Synchronous
2. Interaction with Active or Passive Data

# Unified interaction model



Earliest to now

Standard Database Query

The Past

Now

The Future

# Unified interaction model

# Unified interaction model
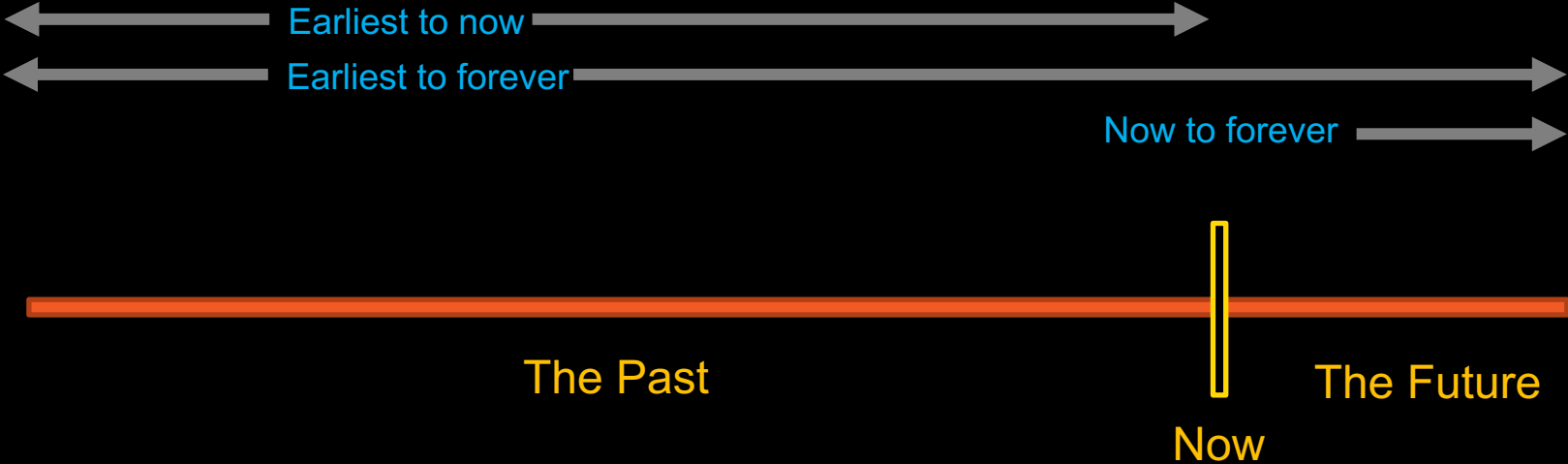
# Unified Interaction Model
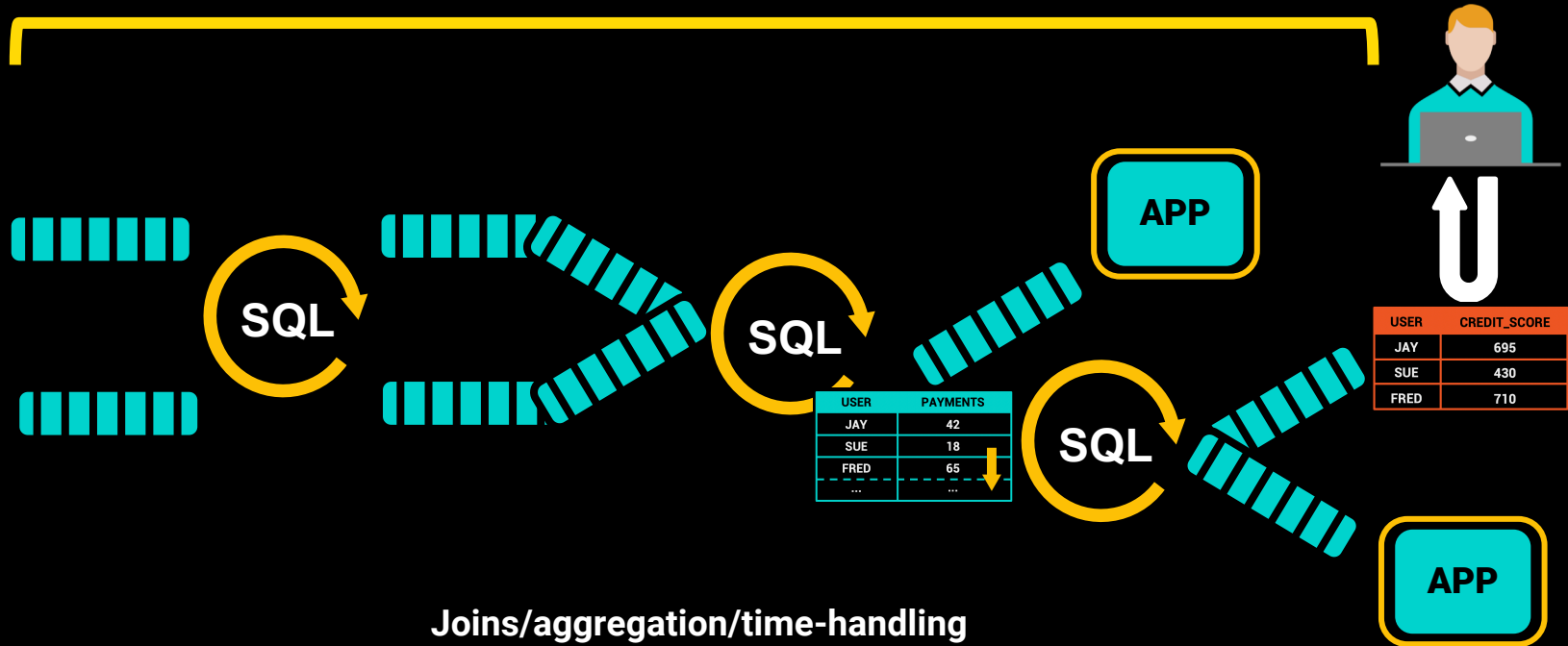
# PUSH

```
SELECT user, credit_score
  FROM orders
  WHERE ROWKEY = 'bob'
  EMIT CHANGES;
```

# PULL

```
SELECT user, credit_score
  FROM orders
  WHERE ROWKEY = 'bob';
```

# Asynchronous => Pipelines

**Transactions**

| USER | PAYMENTS |
|------|----------|
| JAY | 42 |
| SUE | 18 |
| FRED | 65 |
| ... | ... |

| USER | CREDIT_SCORE |
|------|--------------|
| JAY | 695 |
| SUE | 430 |
| FRED | 710 |

SQL

SQL

SQL

APP

APP

**Joins/aggregation/time-handling**

# Other important variants

- Stream processors are often programming frameworks today
  - Storm
  - Flink
  - Kafka Streams
- Today we have active databases that include change streams:
  - Mongo
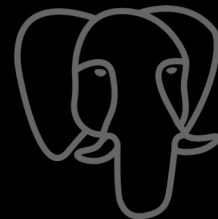  - Couchbase
  - RethinkDB

# As Software Eats the World

THE USER OF THE SOFTWARE

IS MORE
SOFTWARE

# We need
## Asynchronous + Synchronous
## Active + Passive

We still need
all of these

# So is the traditional perception of "a database" enough?

# Ben Stopford
**Confluent**
**@benstopford**
ben@confluent.io