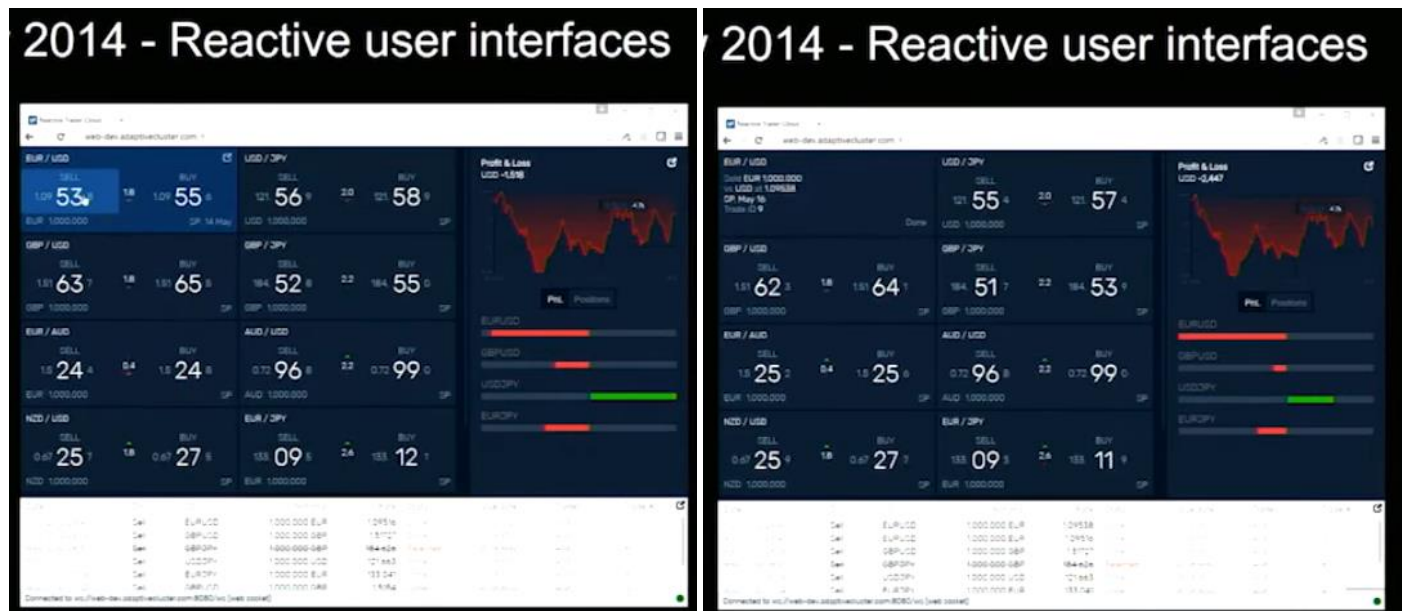


OLIVIER DEHEURLES / JAMES WATSON

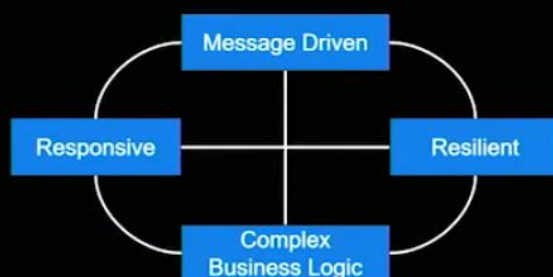
CLUSTERED EVENT-DRIVEN SERVICES ARCHITECTURE + DESIGN

In 2010 I came across a strange, new (to me!) architecture that the LMAX team used for their Foreign Exchange system. You might have heard about the Disruptor; it came out of this project. The core of our system is a clustered service which uses the Raft consensus algorithm to reliably replicate state between the different nodes and hosts our application logic. We will take a quick look at Raft and then at the benefits of this design compared to more “mainstream” architectures.



This architecture offers a clean separation of concerns between the infrastructure – which takes care of the concurrency, I/O and high availability aspects – and the application logic. The clean architecture is a great fit for domain-driven design. If you fancy building fast, resilient services without a database you should come to this talk. Our UIs and the backend systems are constantly sending messages between each other back and forth repeatedly in milliseconds.

Attributes of a trading system



The problem

Complex multi-threaded programs

⇒ Hard to reason about

⇒ Hard to debug

The Solution - LMAX Architecture

Blueprint for reactive services

Clean separation of concerns

- technical requirements
- Business logic

Keep the business logic as simple as possible

Very easy to debug



Dave Farley



The business logic runs on a single thread in this architecture.

Our journey to the LMAX Architecture

- Stage 1
 - Simple Model
 - Easy To Debug
- Stage 2
 - Transparent fault tolerance
 - Clustering with consistent replication
- Stage 3
 - Durability without a database

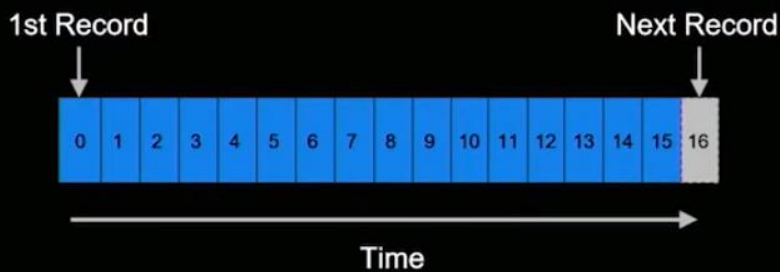
Deterministic execution in a distributed event system

Determinism

Given an **initial state** and a **command**, a deterministic model will always produce the same **output state** and **side-effect(s)**



Log as Sequencing Mechanism



Single Threaded Log Consumer



Single threaded, isn't that slow? No!

10K ops/sec: quite easy to achieve with sensible code

100K ops/sec: some profiling required

1M+ ops/sec:

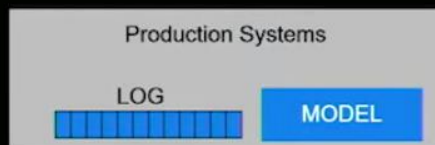
- Optimised data structures
- Low complexity algorithms - $O(1)$, $O(\log(N))$
- Minimise allocation

Determinism - Handling Time

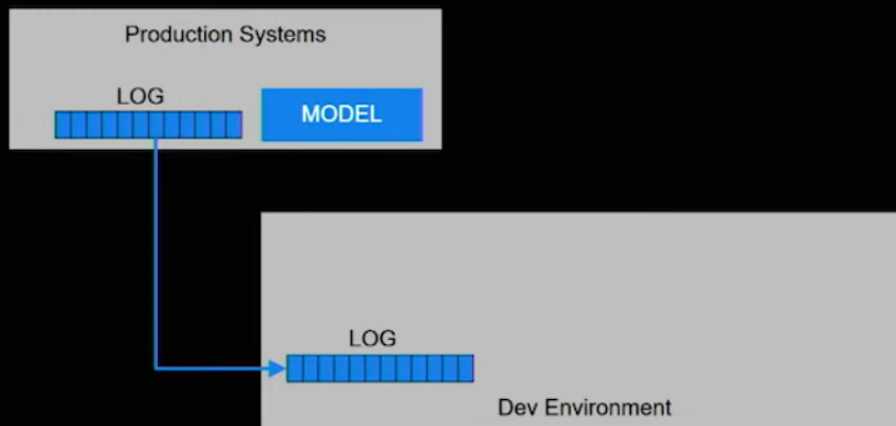


We add the timestamp of the pickup time to the message

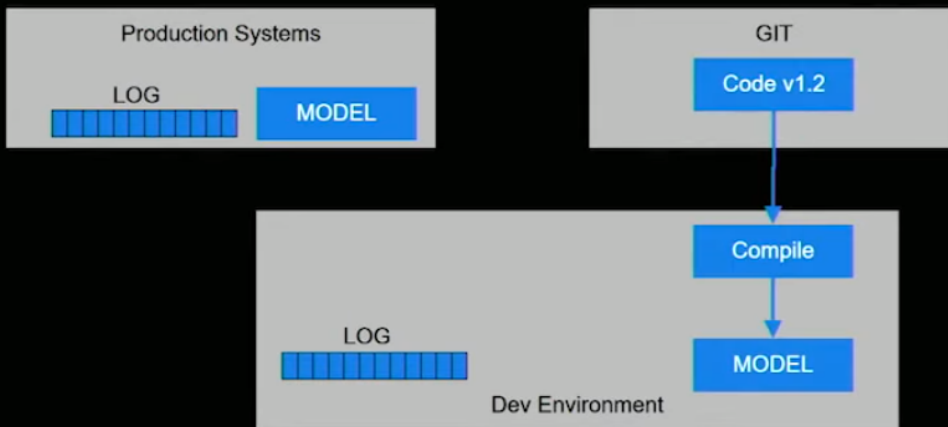
Easy Debugging by Replicating Exact State



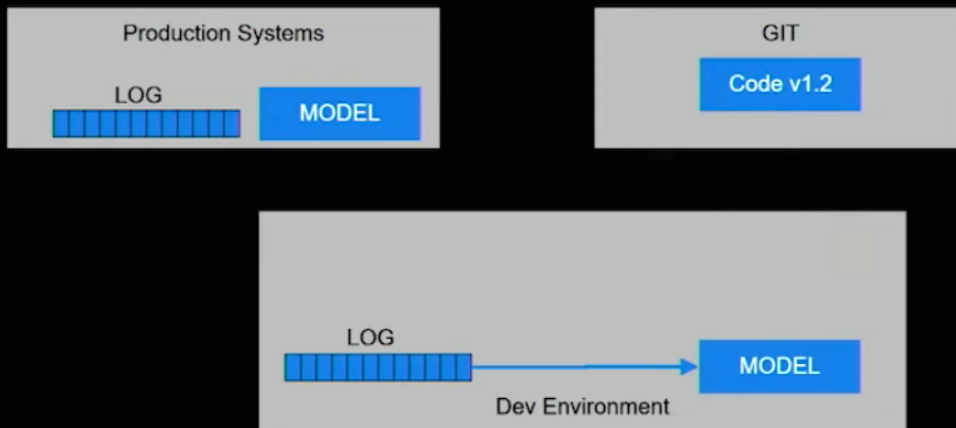
Get the log



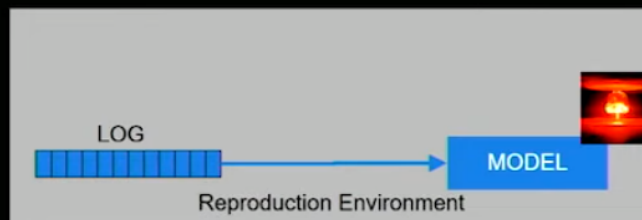
Get the same version of the code, compile



Replay the log

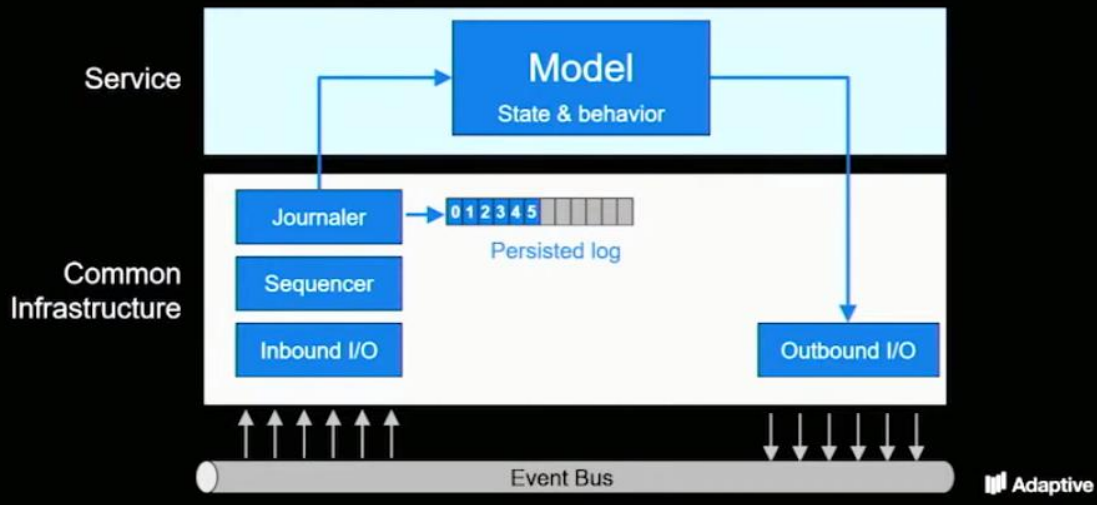


Reproduce the issue, with debugger attached!

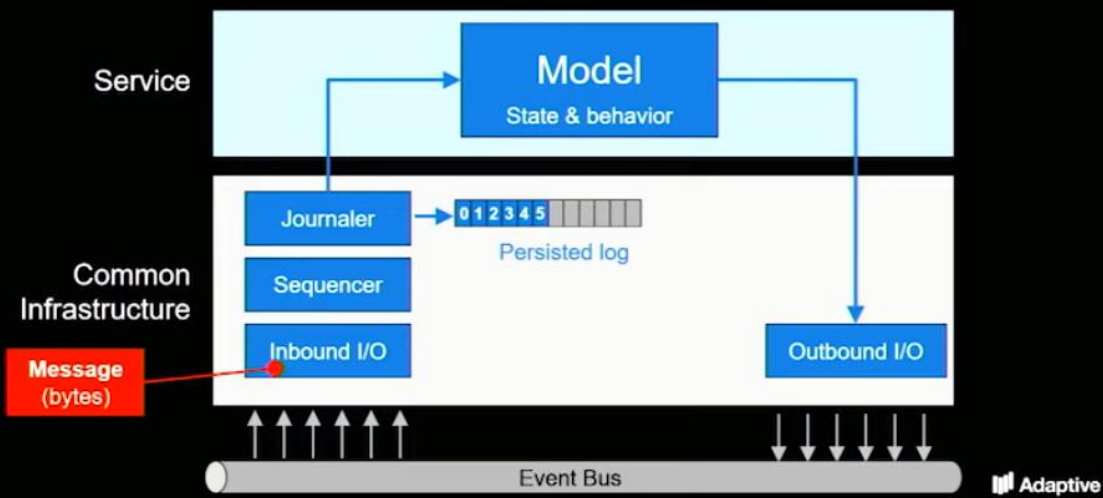


Service Anatomy

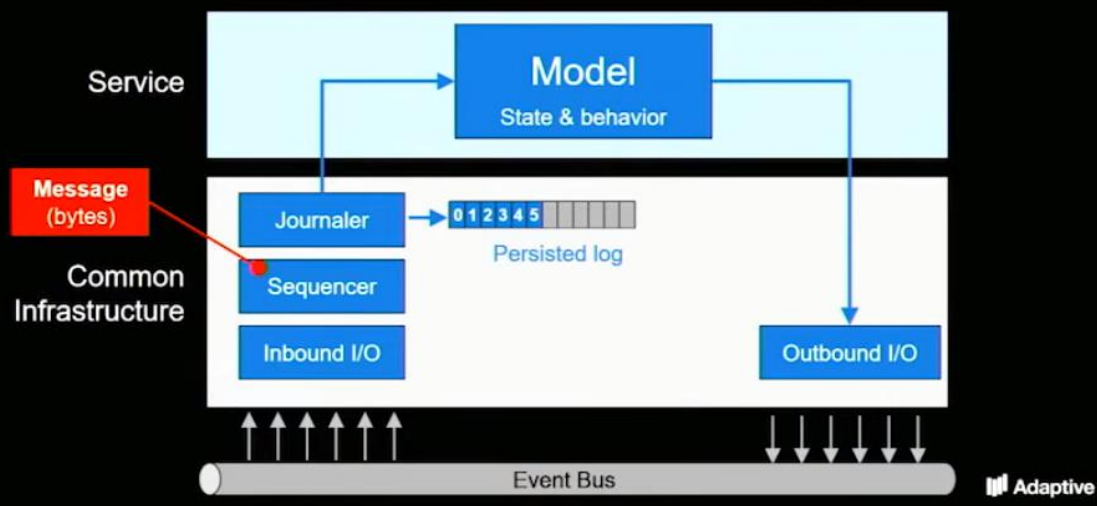
Service Anatomy: Separate Infrastructure from Model



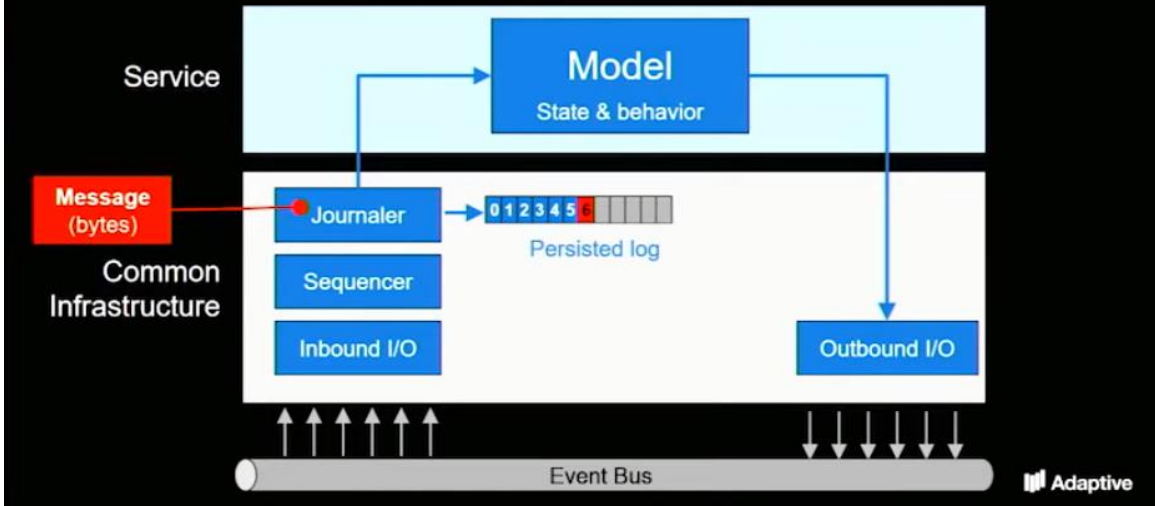
Incoming message



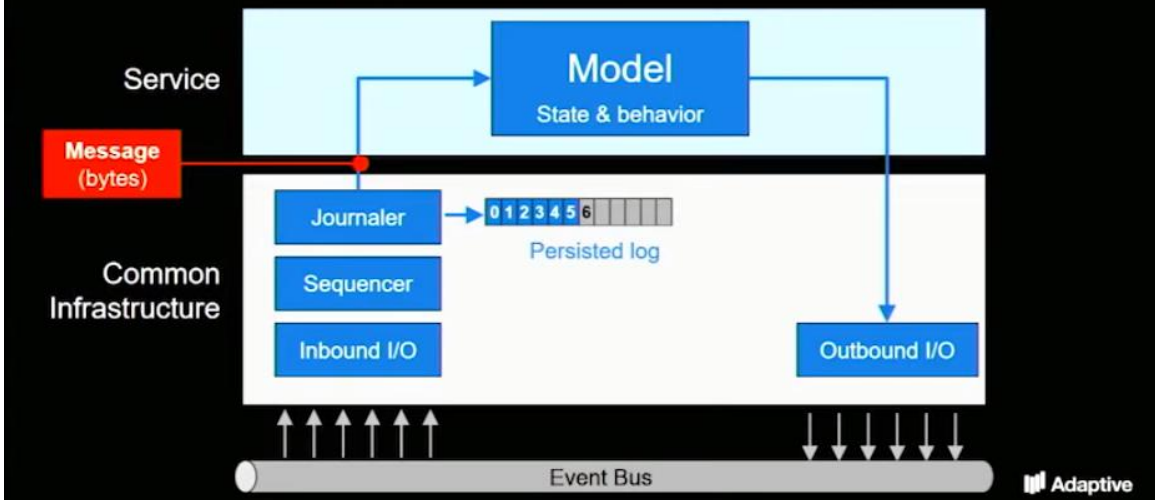
Message is sequenced



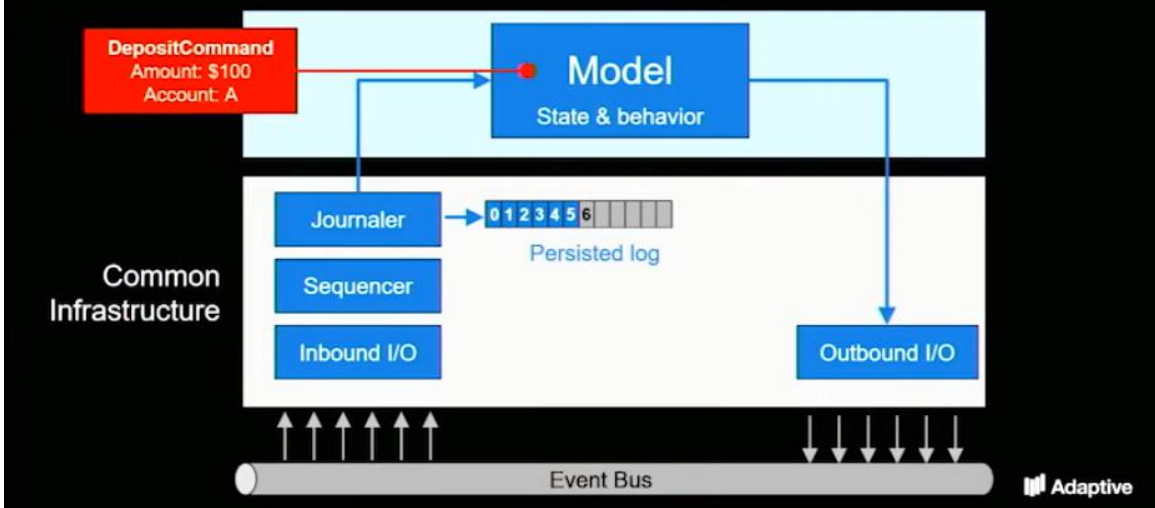
Message is written to disk



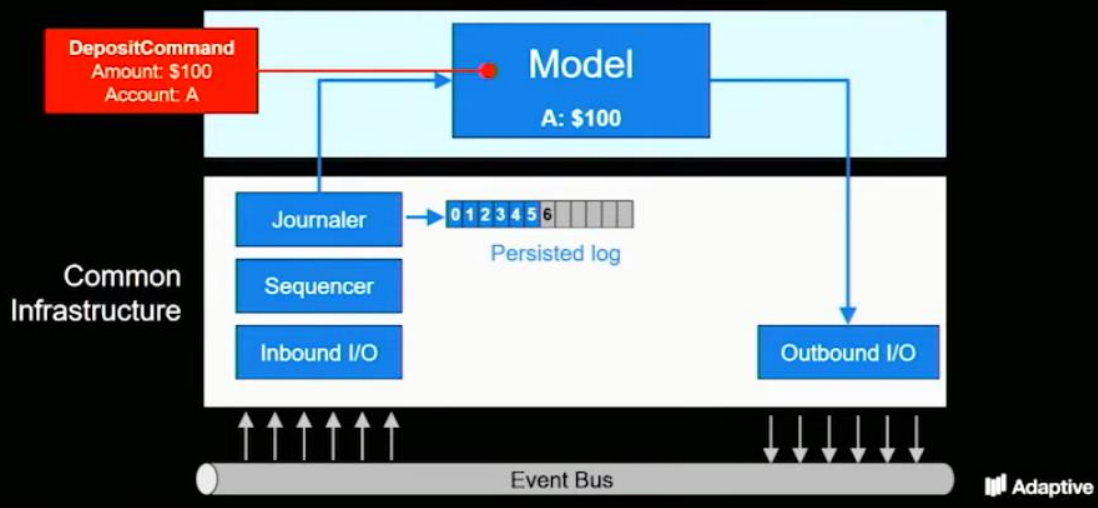
Message passed to the model



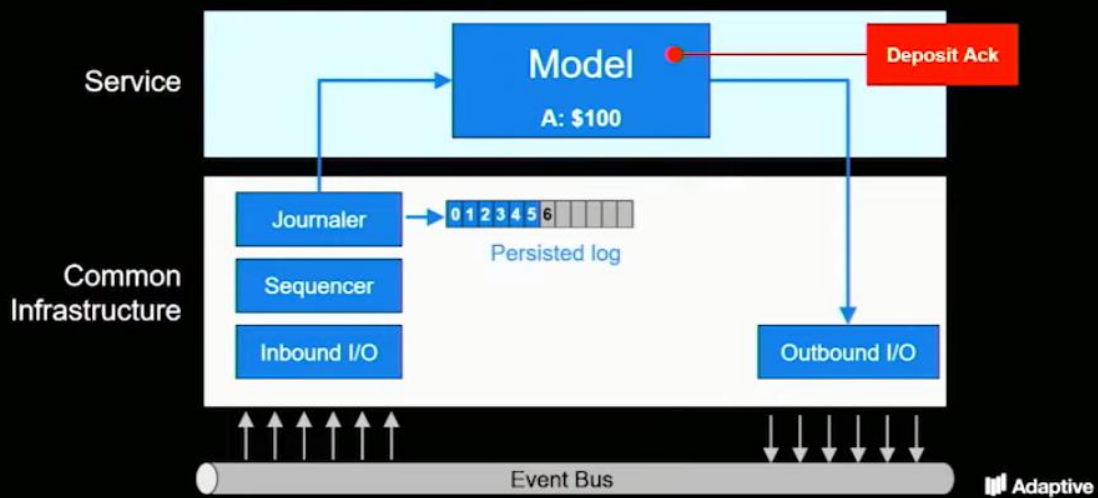
Message is decoded



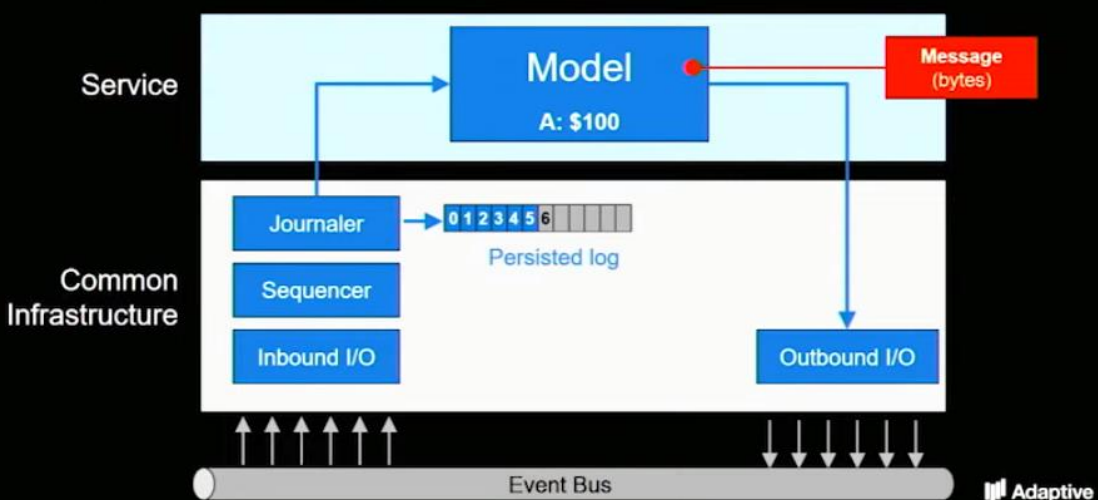
Command is applied to the model



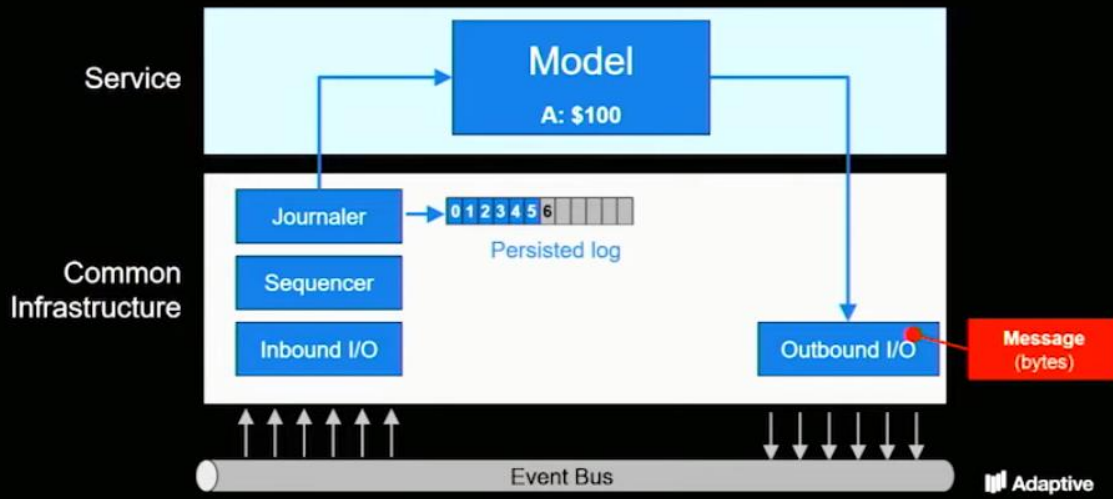
Response message is created



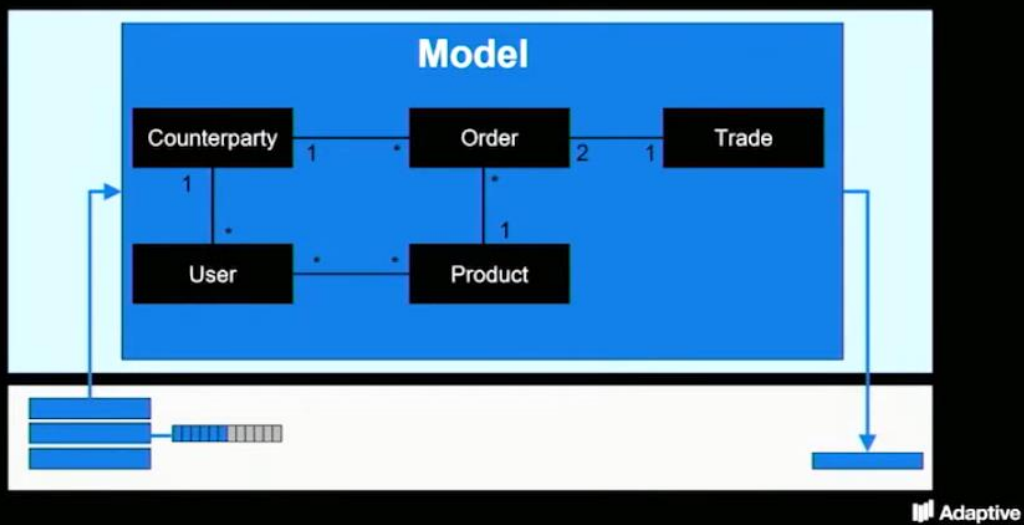
Response message is encoded



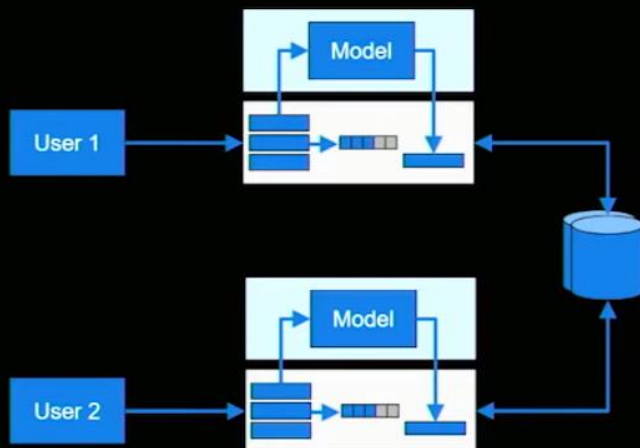
Message is sent over the network



Domain model unaware of infrastructure



1st Iteration - State Synchronized via Database



Benefits

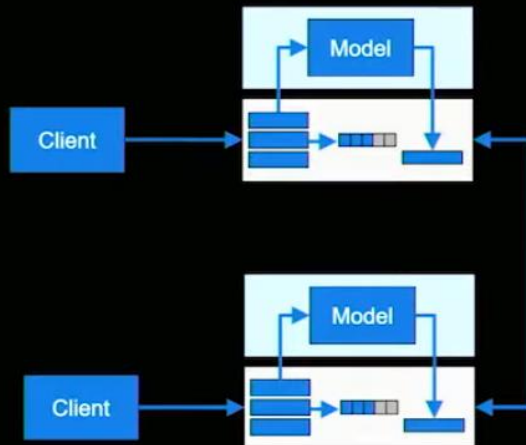
- Simple model
- Debugability

Pain points

- Synchronisation through database

Fault Tolerance & State Replication

2nd Iteration - Clustering



Fault tolerant data replication

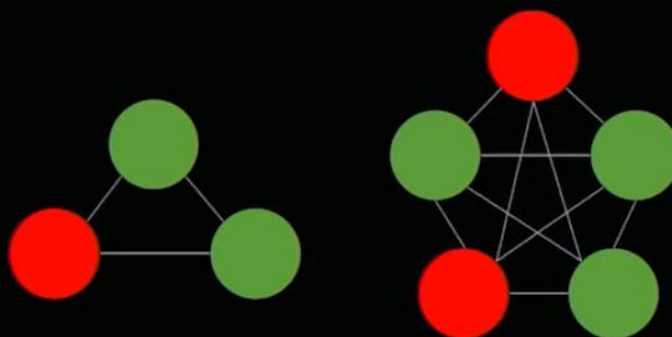
Difficult problem:

- The log must stay **consistent** between servers
- Even if a **server fails** or the **network fails**



Distributed consensus problem

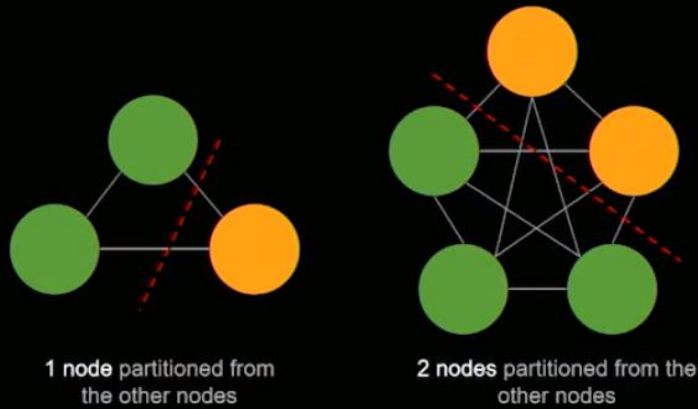
RAFT



Tolerates fault of 1 node

Tolerates fault of 2 nodes

RAFT - Network failure



Log replication with RAFT

Leader election



Log replication with RAFT

Clients connect to the **leader**



Log replication with RAFT

Client sends a DepositCommand to the leader



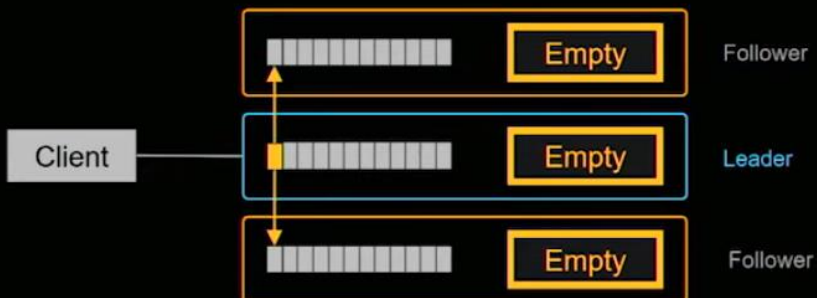
Log replication with RAFT

Leader writes the command to its log



Log replication with RAFT

Leader sends the message to followers (replication)



Log replication with RAFT

Followers store the message in their log and acknowledge back to the leader



Log replication with RAFT

Leader marks the message as committed, process the message and notifies followers to commit the message



Log replication with RAFT

Followers commit the message in their log and process it



Log replication with RAFT

Leader sends a deposit acknowledgement to the client



Log replication with RAFT

What if the leader dies?



Log replication with RAFT

Followers hold an election, one of them becomes leader



Log replication with RAFT

Client detects disconnection from leader and connects to new leader



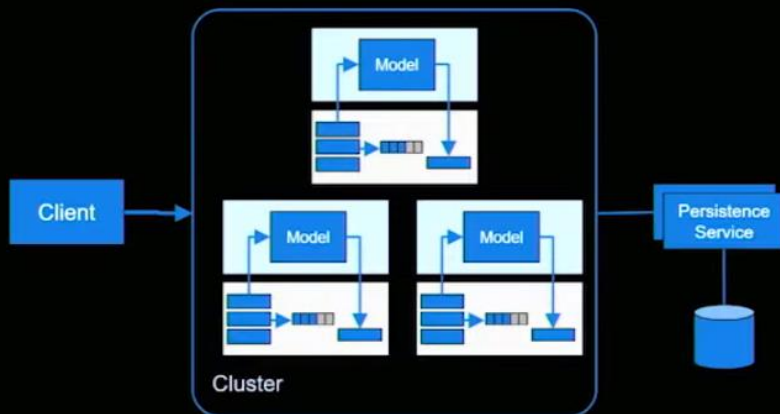
Log replication with RAFT

Client is connected to another instance of the service, in the same state

⇒ System is tolerant to server and network failures



2nd Iteration - Clustering



Benefits

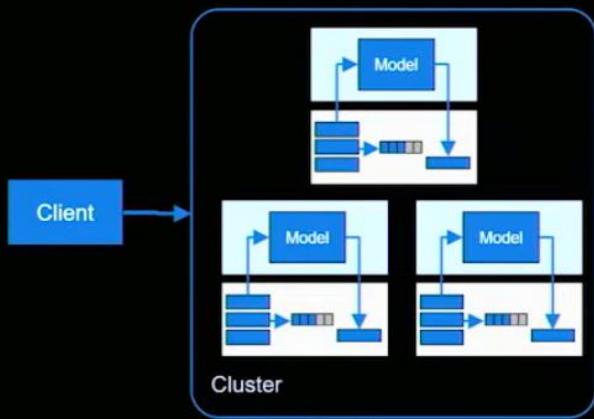
- Fault tolerant
- Consistent model & it doesn't change

Pain points

- Database as golden source

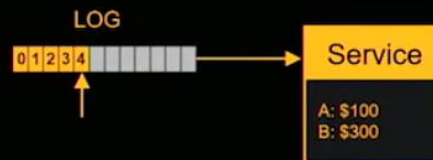
Durability

3rd Iteration - Durable Services



System restart

Given the current state...



System restart

Stop the service



System restart

Service restarts, reads from beginning of the log



System restart

Replay all commands up to latest

Services is guaranteed to get back in the same state due to **deterministic** logic



Managing the log

Problem: the size of the **log** affects **recovery time**

Not a problem if

- Log grows slowly
- Your model has short lifetime

Otherwise consider **snapshotting** or other forms of **compaction**

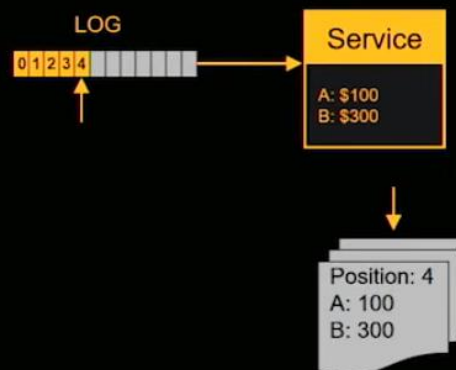
Snapshotting

Considering the system in this state...



Snapshotting

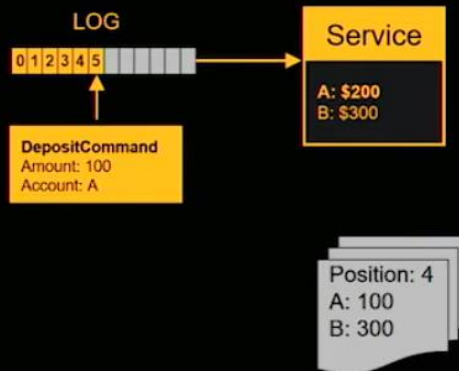
Service writes current state to disk with position in the log



We take the model in memory and serialize it to save it

Snapshotting

Message 5 gets written to the log and processed

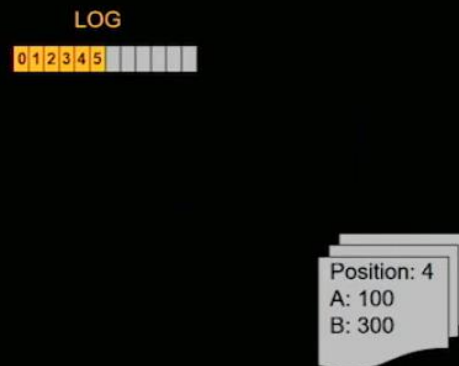


Snapshotting

Stop the service and restart

We still have

- The log
- The snapshot



Snapshotting

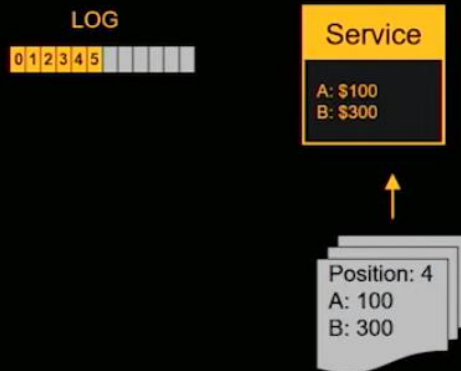
Service starts, initially with an empty model



Snapshotting

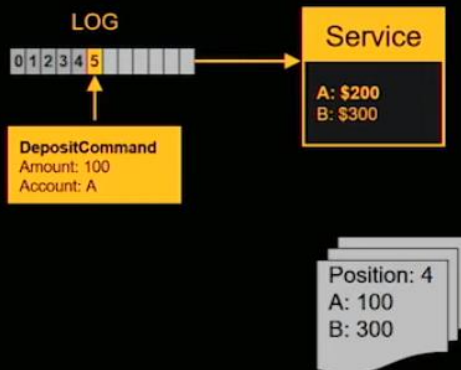
Service loads most recent snapshot

- Re-hydrate model



Snapshotting

Replay all messages after current position

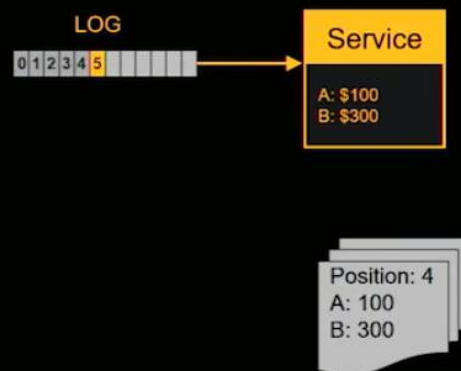


Snapshotting

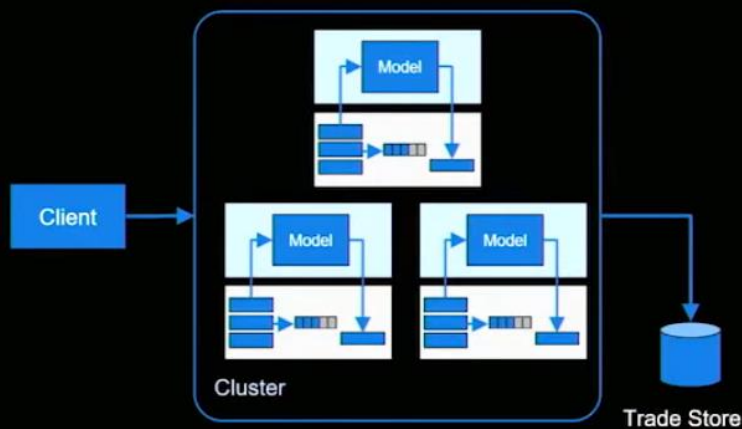
System is back in the same state

Snapshotting makes it possible to **trim** the log

With a snapshot at position N, we can **archive** the log up to the same position



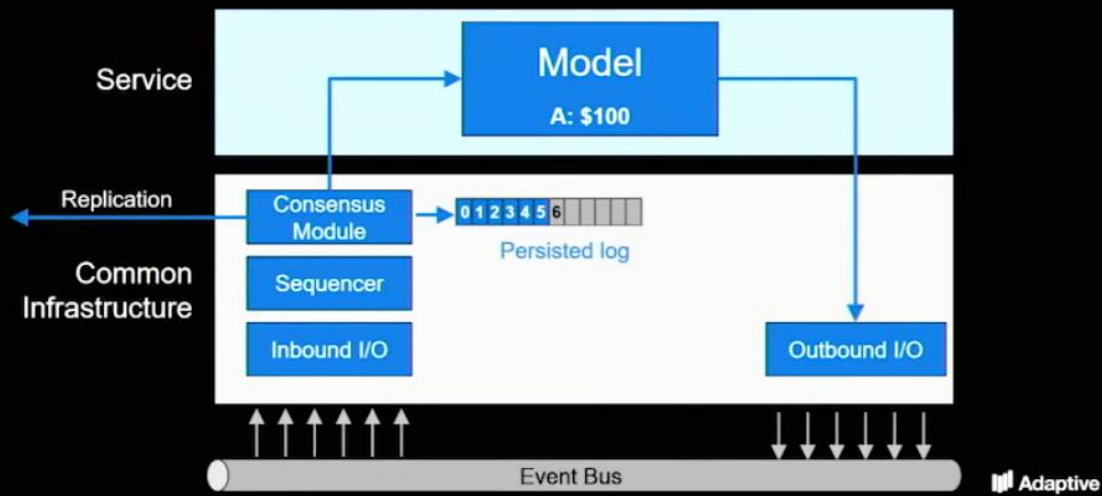
Iteration 3 - Durable Services



Benefits

- Model unaffected
- Fault tolerant
- Consistent model
- Durable without database

Reactive Application Server



Don't build the infrastructure yourself!

Open source

- Copycat - <http://atomix.io/>
- Aeron - <https://github.com/real-logic/aeron>



Todd Montgomery



Martin Thompson

Other use-cases

- Online games
 - MMO
 - Gambling
- Ticketing system
- Consistent Databases / Caches
- Many more..

Wrapping up

- Demonstrated the simplicity of this approach
- Widely applicable
- Deserves more attention!
- Learn more
 - LMAX Architecture <https://goo.gl/q1iSCB>
 - Raft paper and website - <https://raft.github.io/>
 - The Log - <https://goo.gl/m4iWqn>
 - White paper - <http://weareadaptive.com/blog>