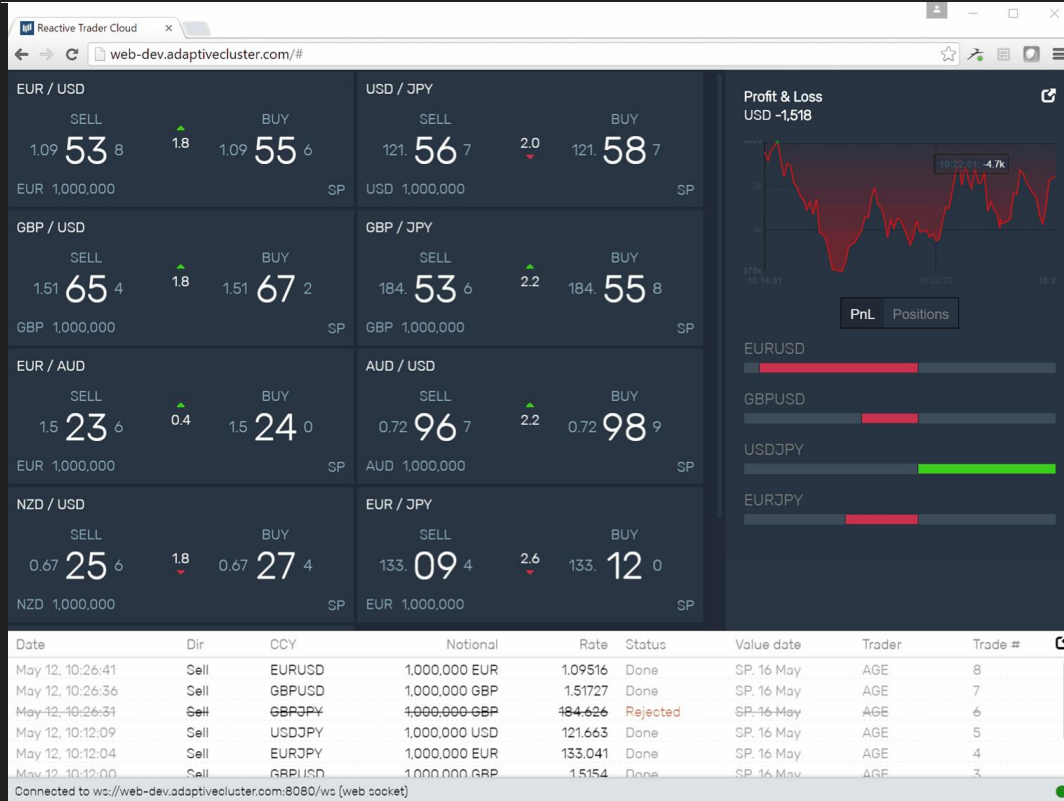


Clustered Event-Driven Services

Simple model in a complex event-driven world

YOW! 2017 - Sydney

Yow 2014 - Reactive user interfaces

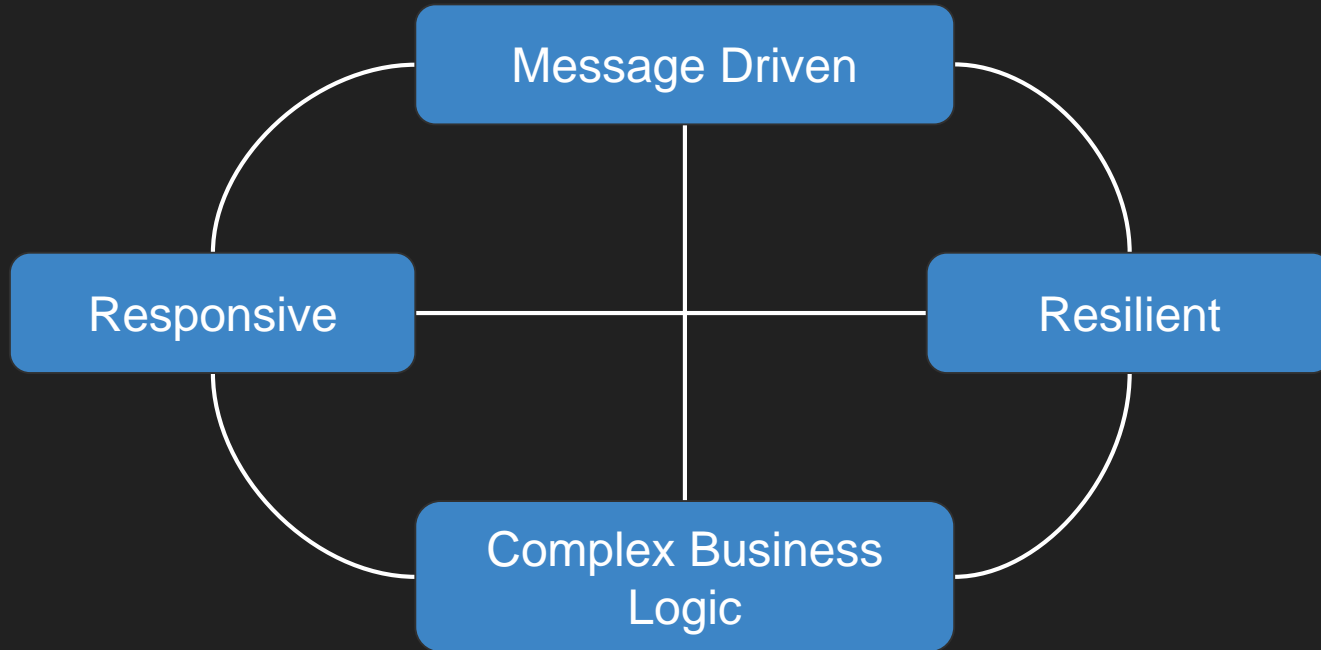


Lee Campbell



Matt Barrett

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



Martin Thompson

Our journey to the LMAX Architecture

Stage 1

Simple Model
Easy To Debug

Stage 2

Fault Tolerance
State Replication

Stage 3

Durability without a
Database



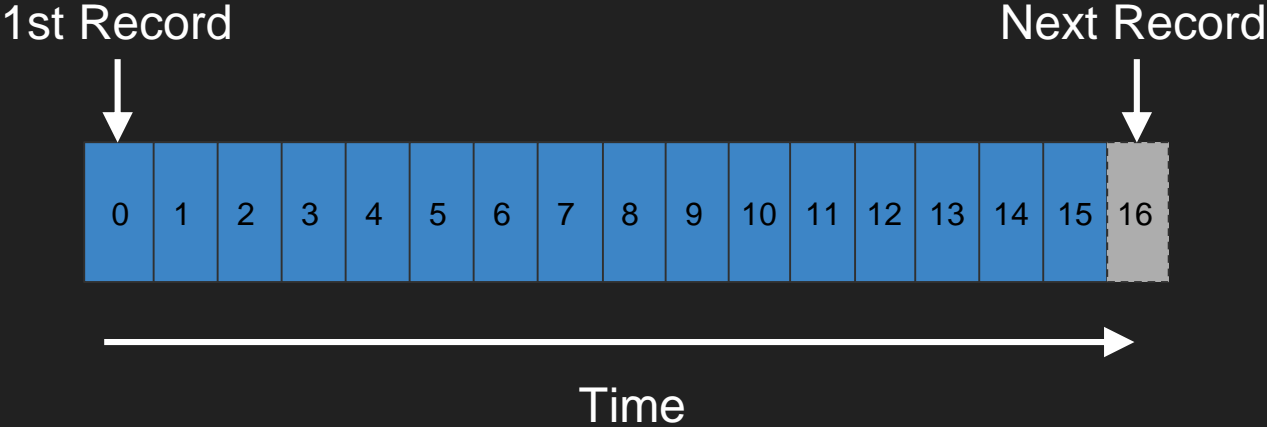
Deterministic execution in a distributed event system

Determinism

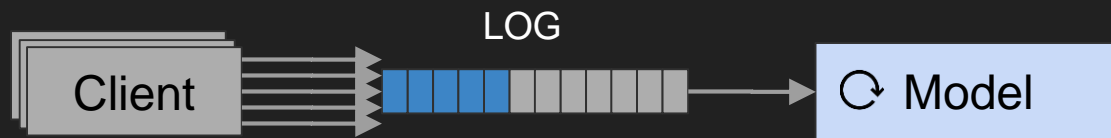
Given an **initial state** and a **command**, a deterministic model will always produce the same **output state** and **side-effect(s)**



Log as Sequencing Mechanism



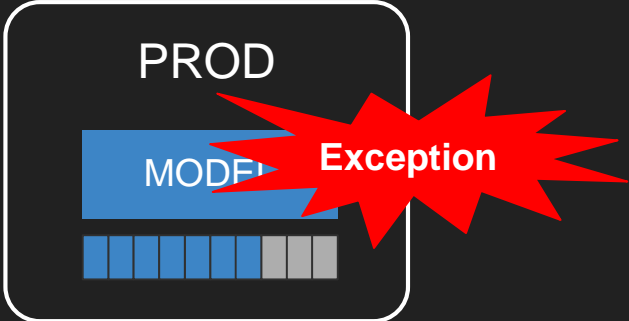
Single Threaded Log Consumer



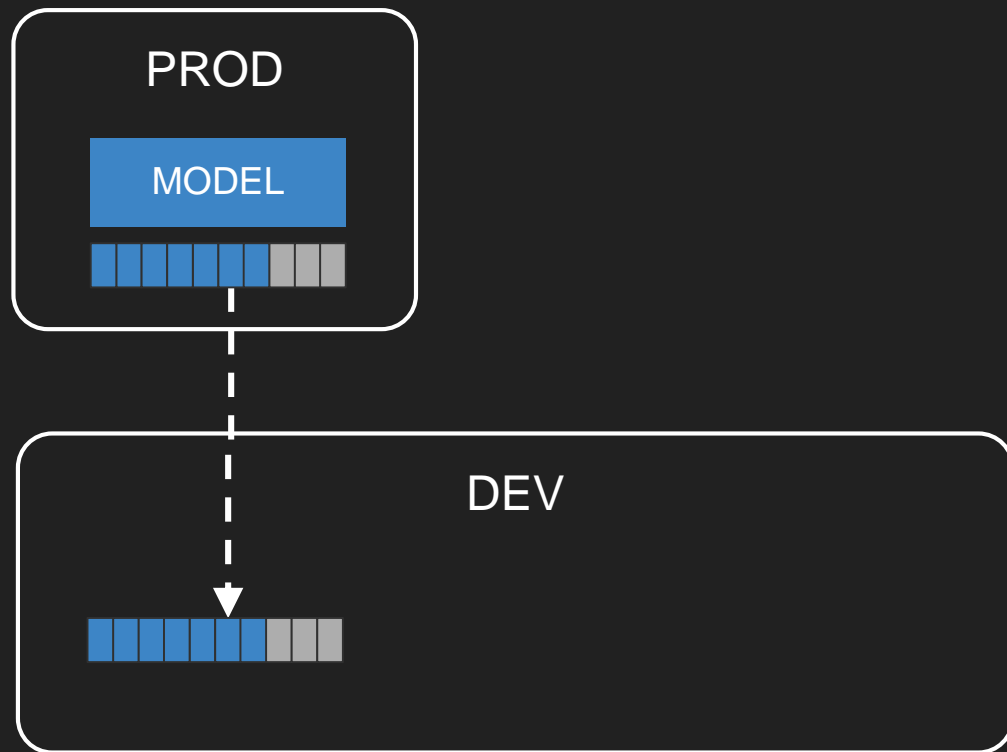
Easy Debugging by Replicating Exact State



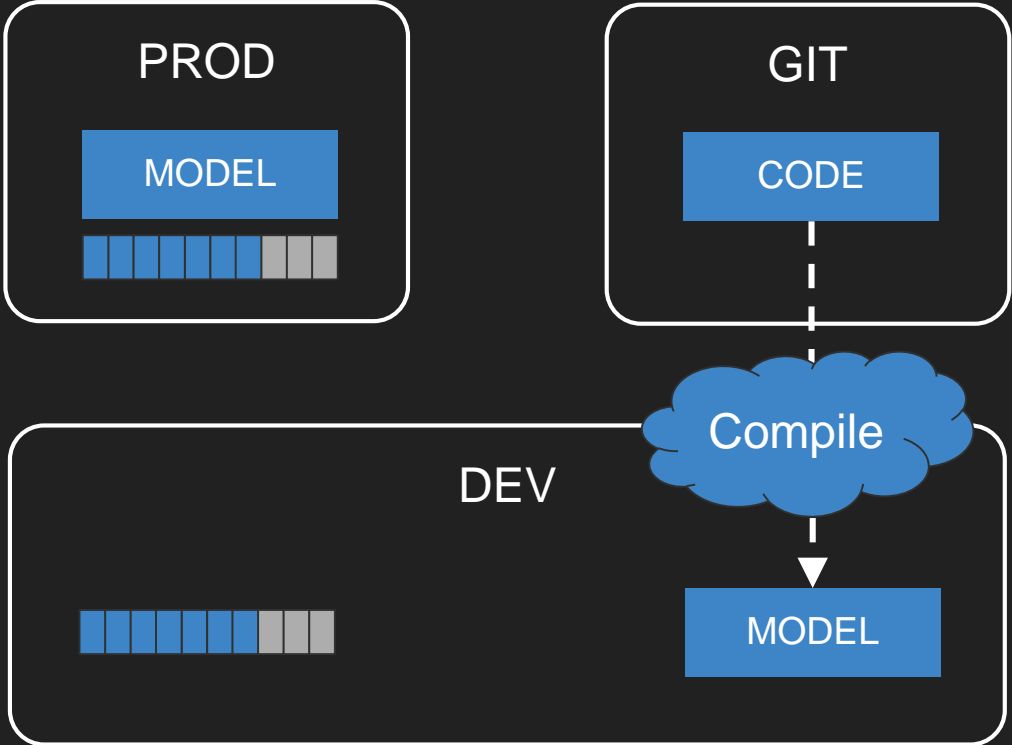
Production issue



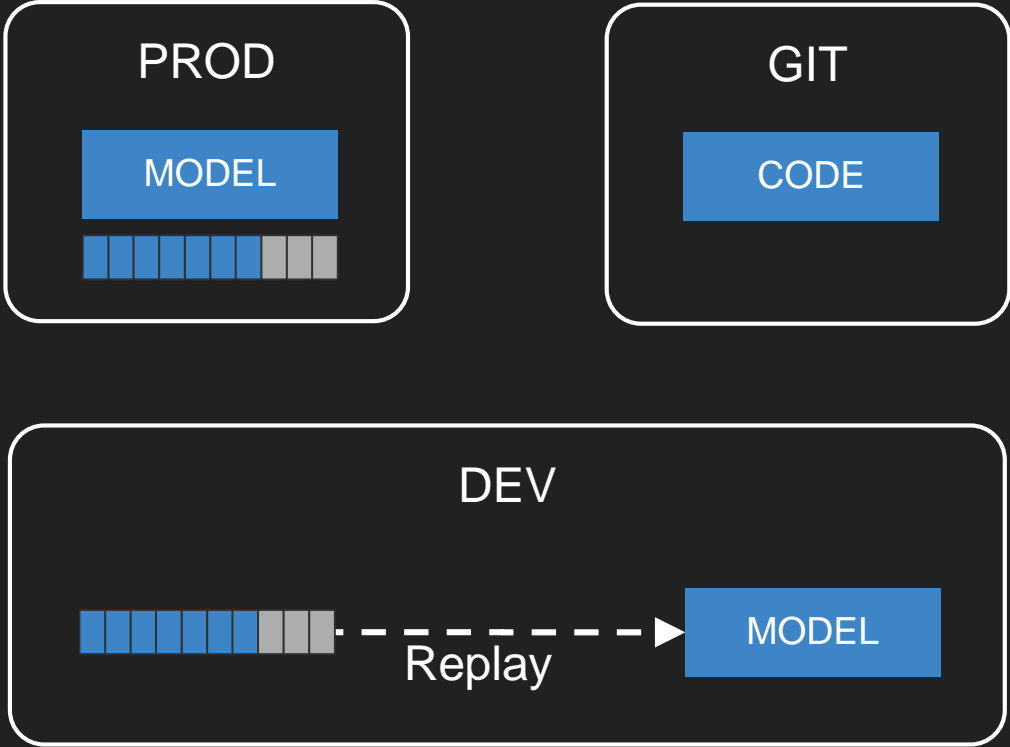
Get the log



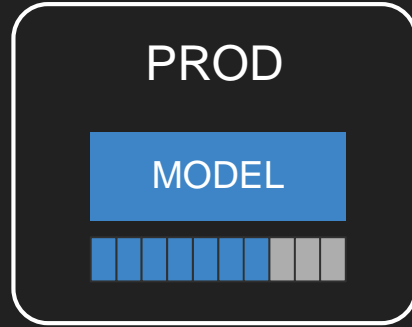
Get the same version of the code, compile



Replay the log



Reproduce the same behaviour, debugger attached



Performance

Locality of state and behaviour

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

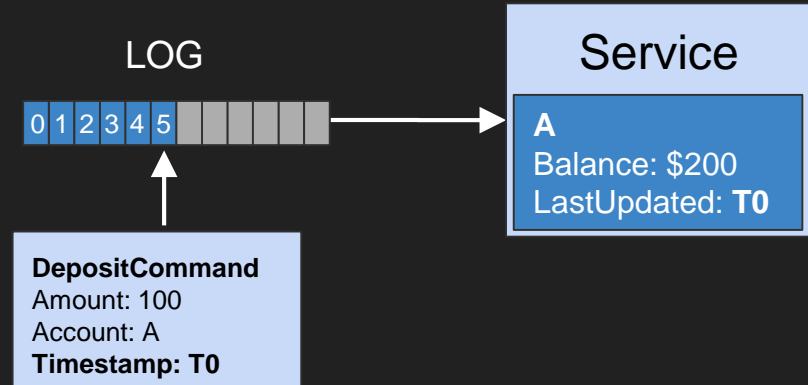
100K ops/sec - profiling required

1M+ ops/sec

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

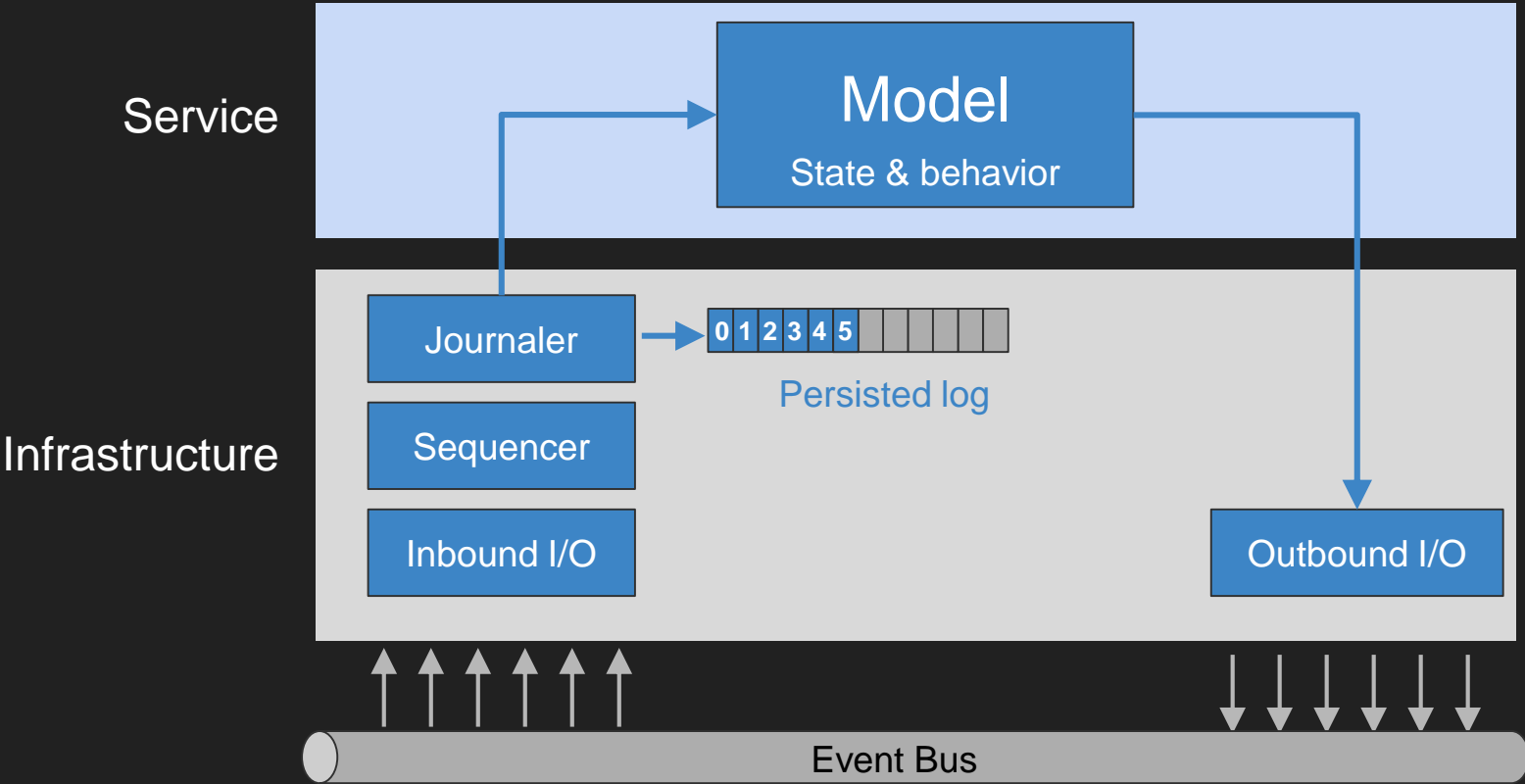
Handling Time Deterministically

- System time is not deterministic
- Inject time through the log
- Great for testability

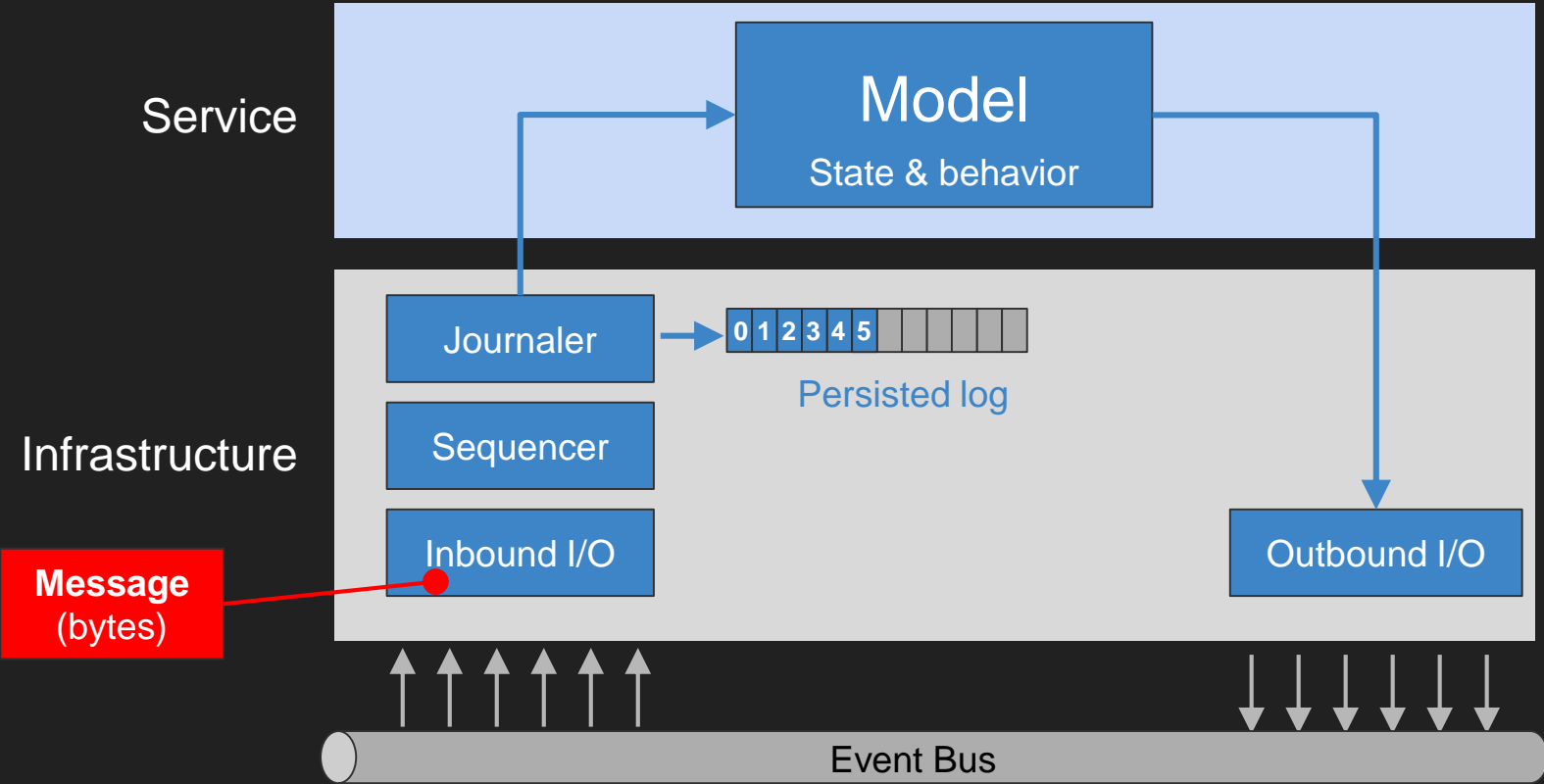


Service Anatomy

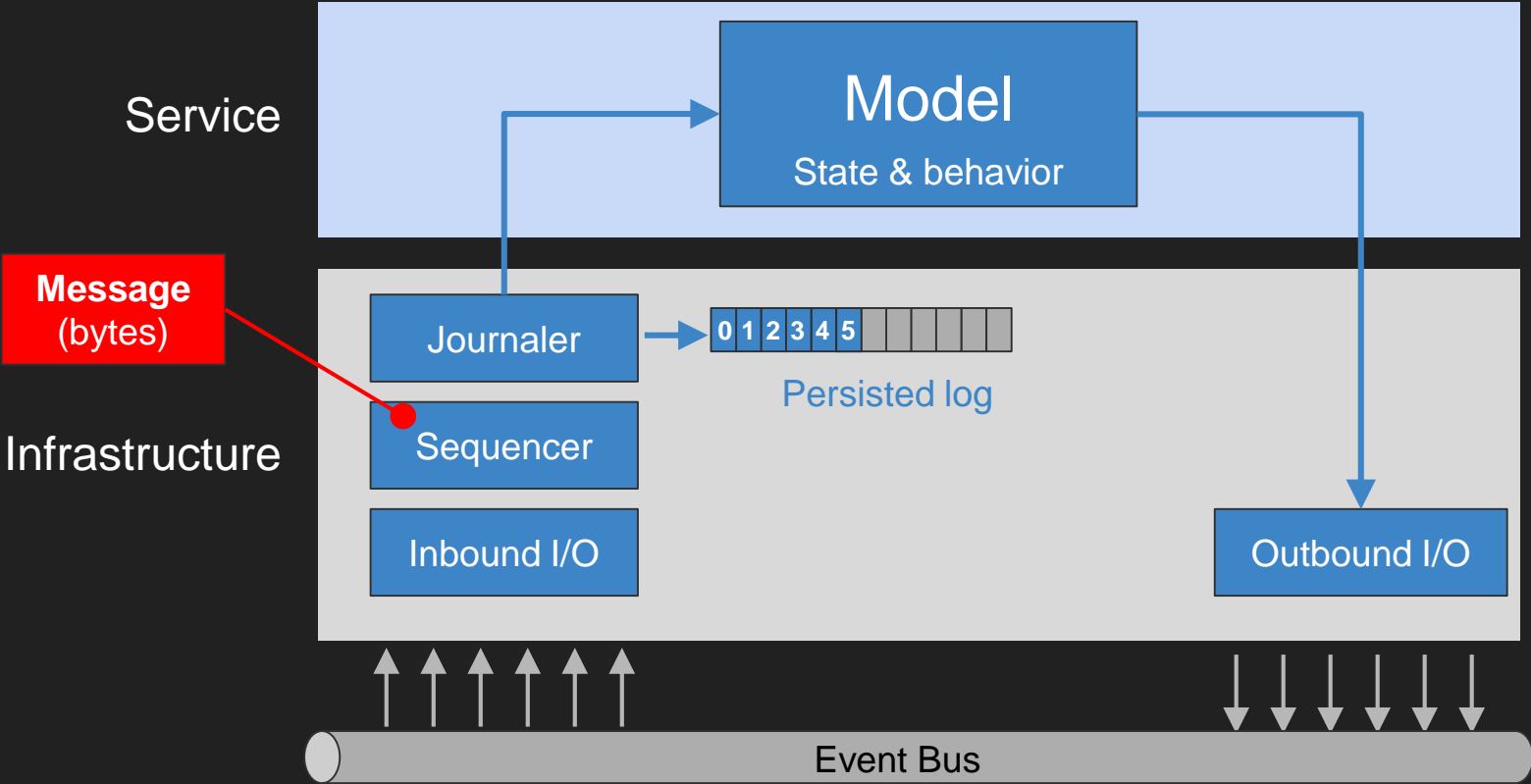
Service Anatomy: Separate Infrastructure from Model



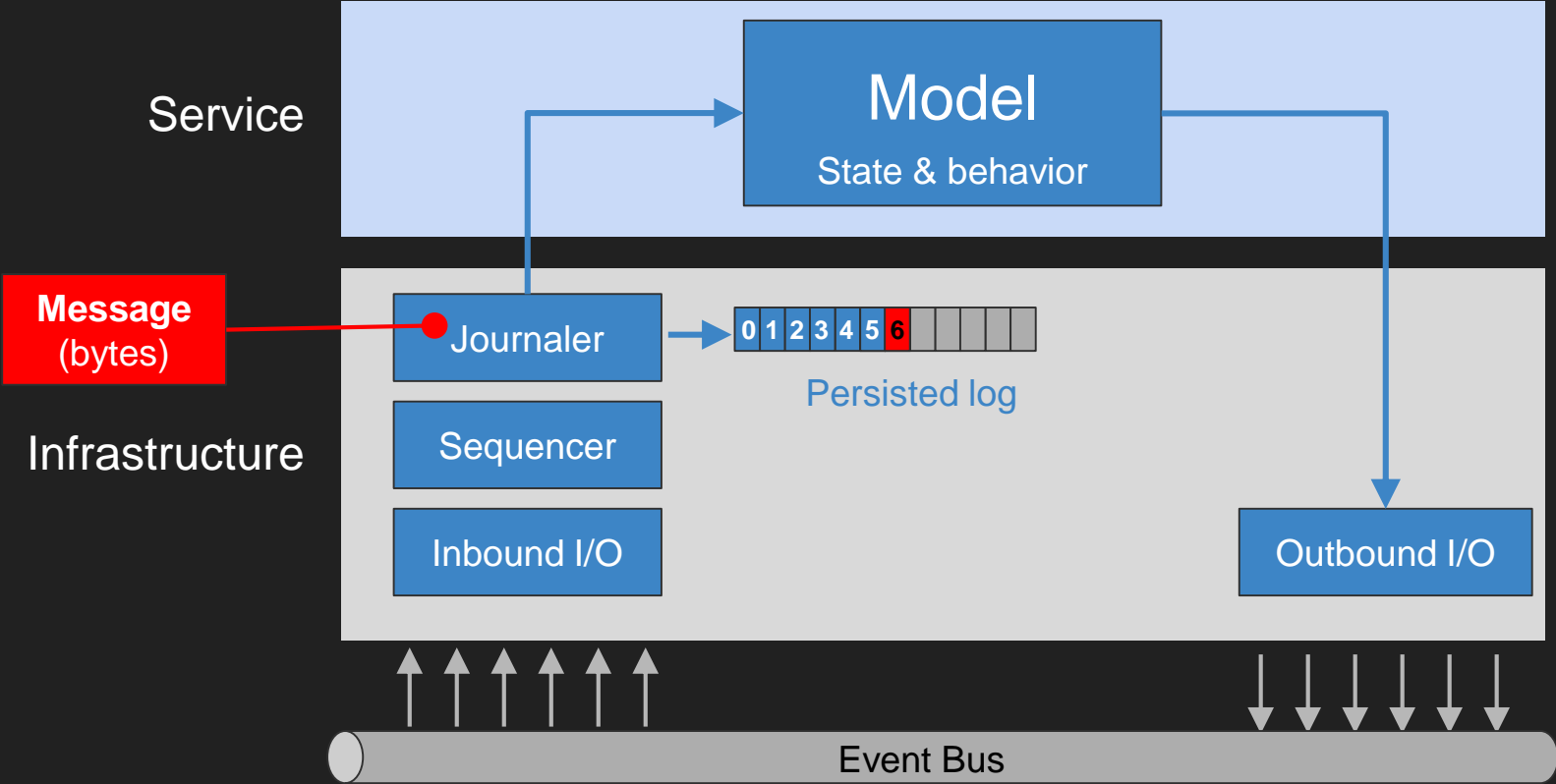
Incoming message



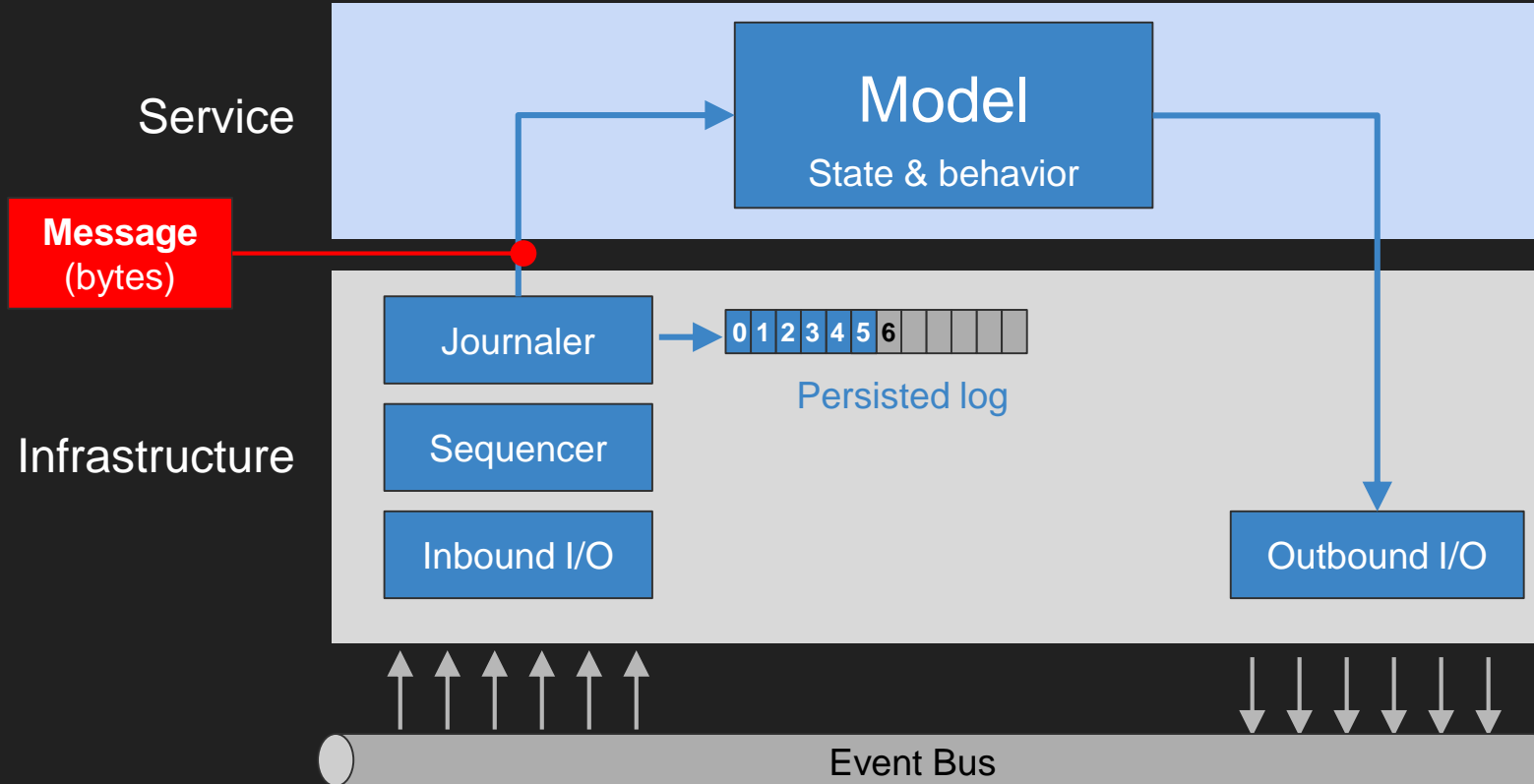
Message is sequenced



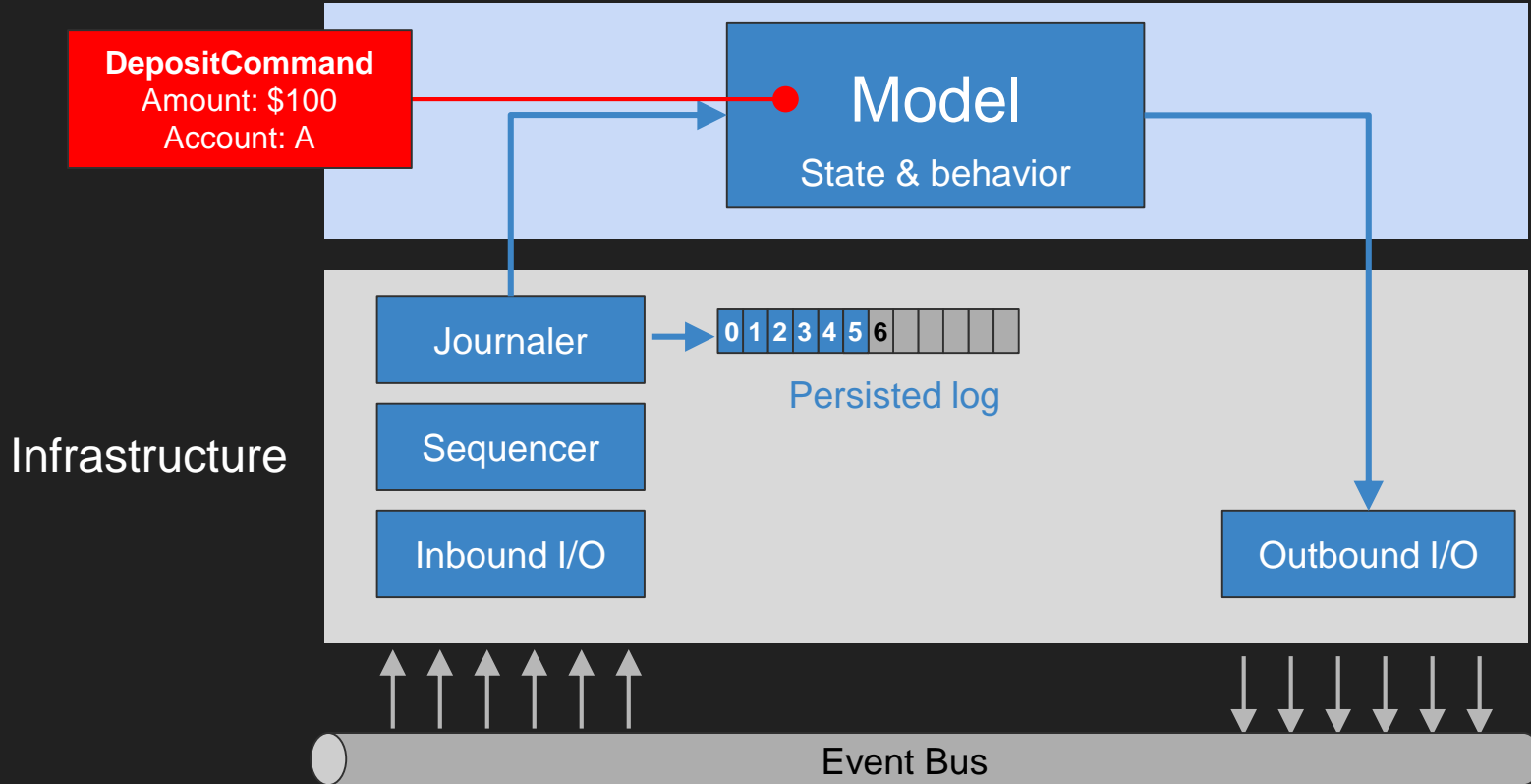
Message is written to disk



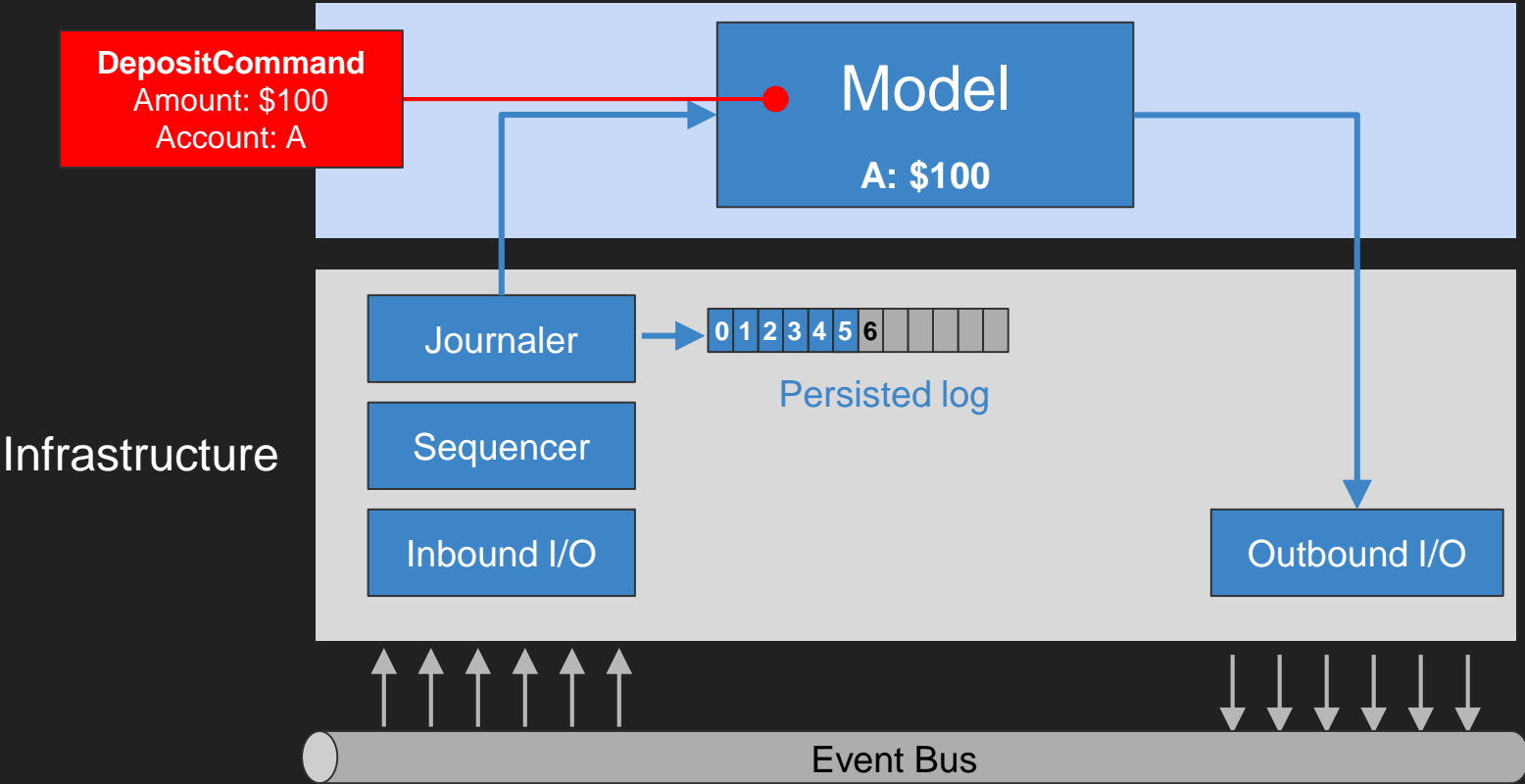
Message passed to the model



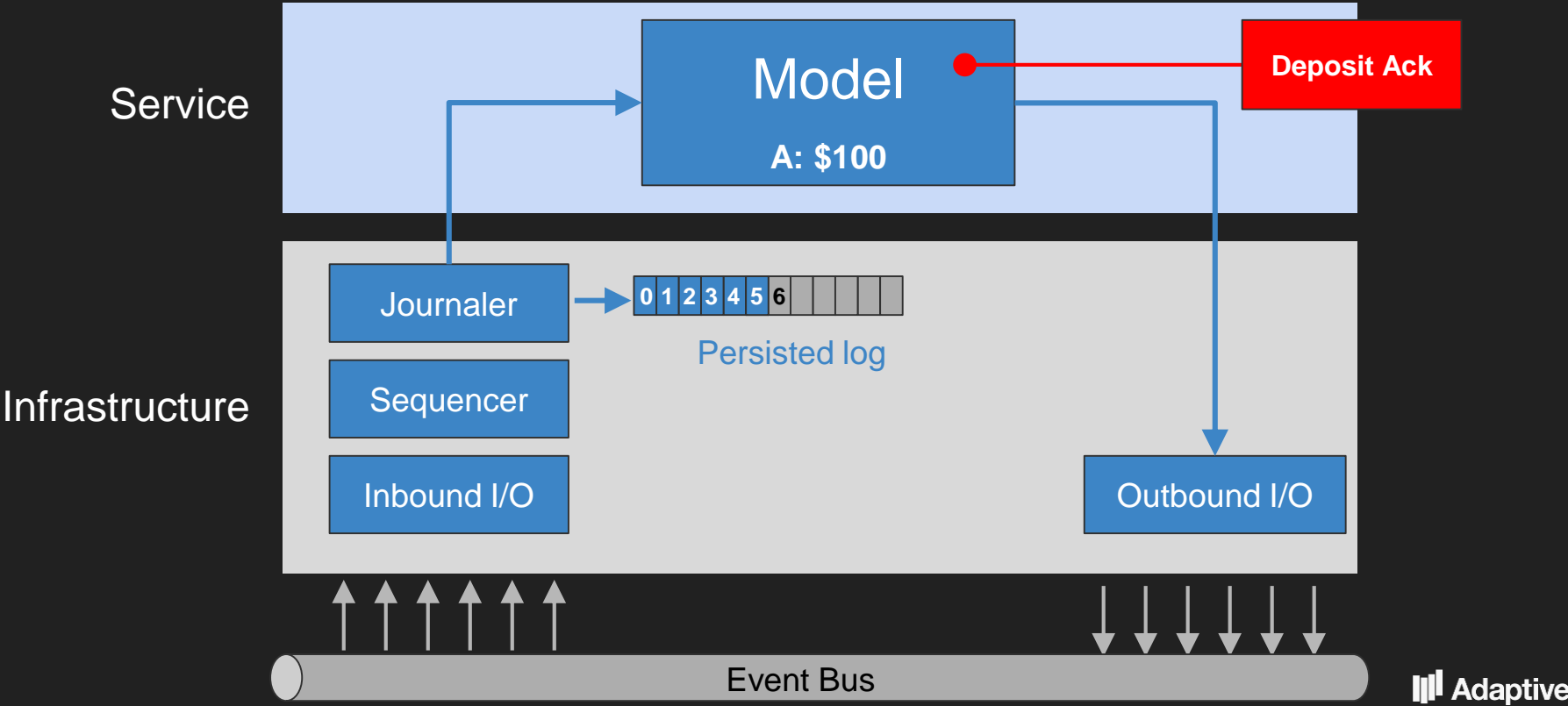
Message is decoded



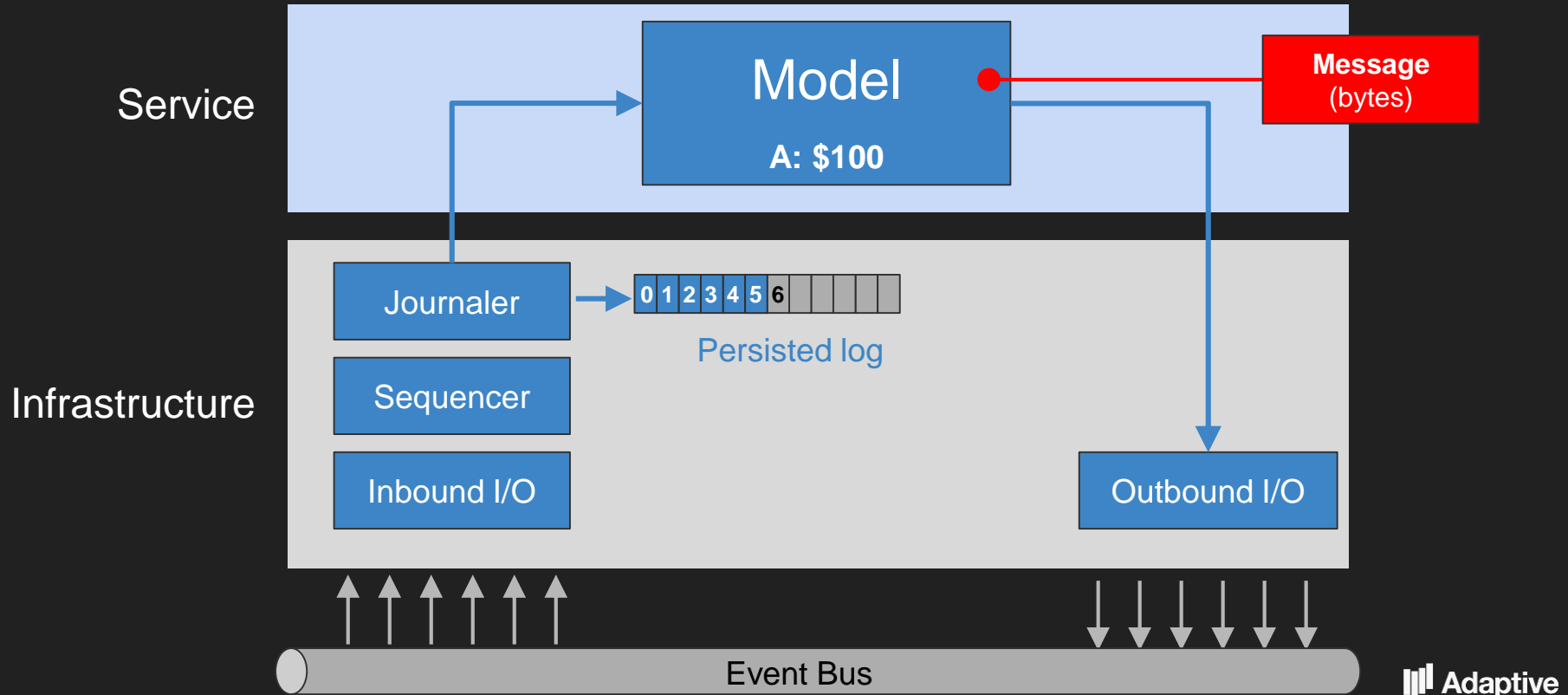
Command is applied to the model



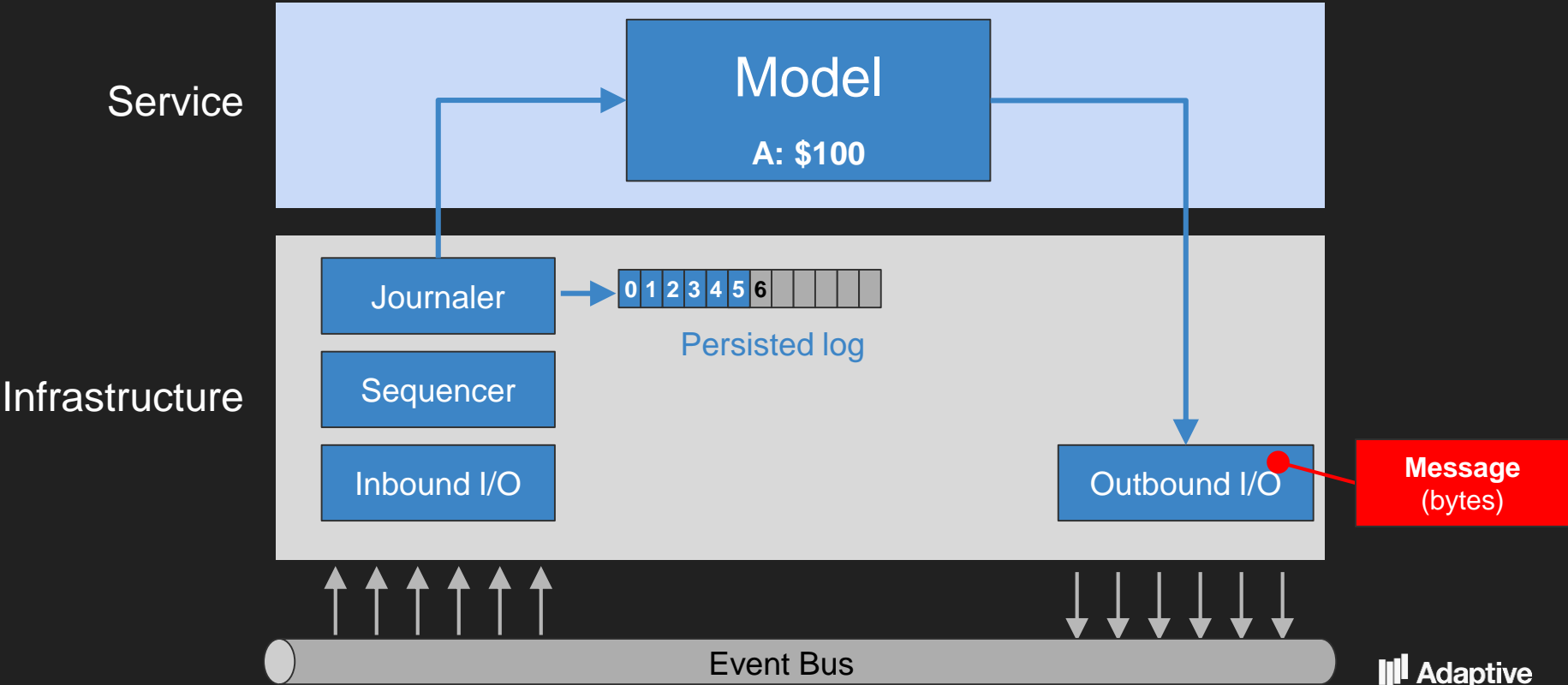
Response message is created



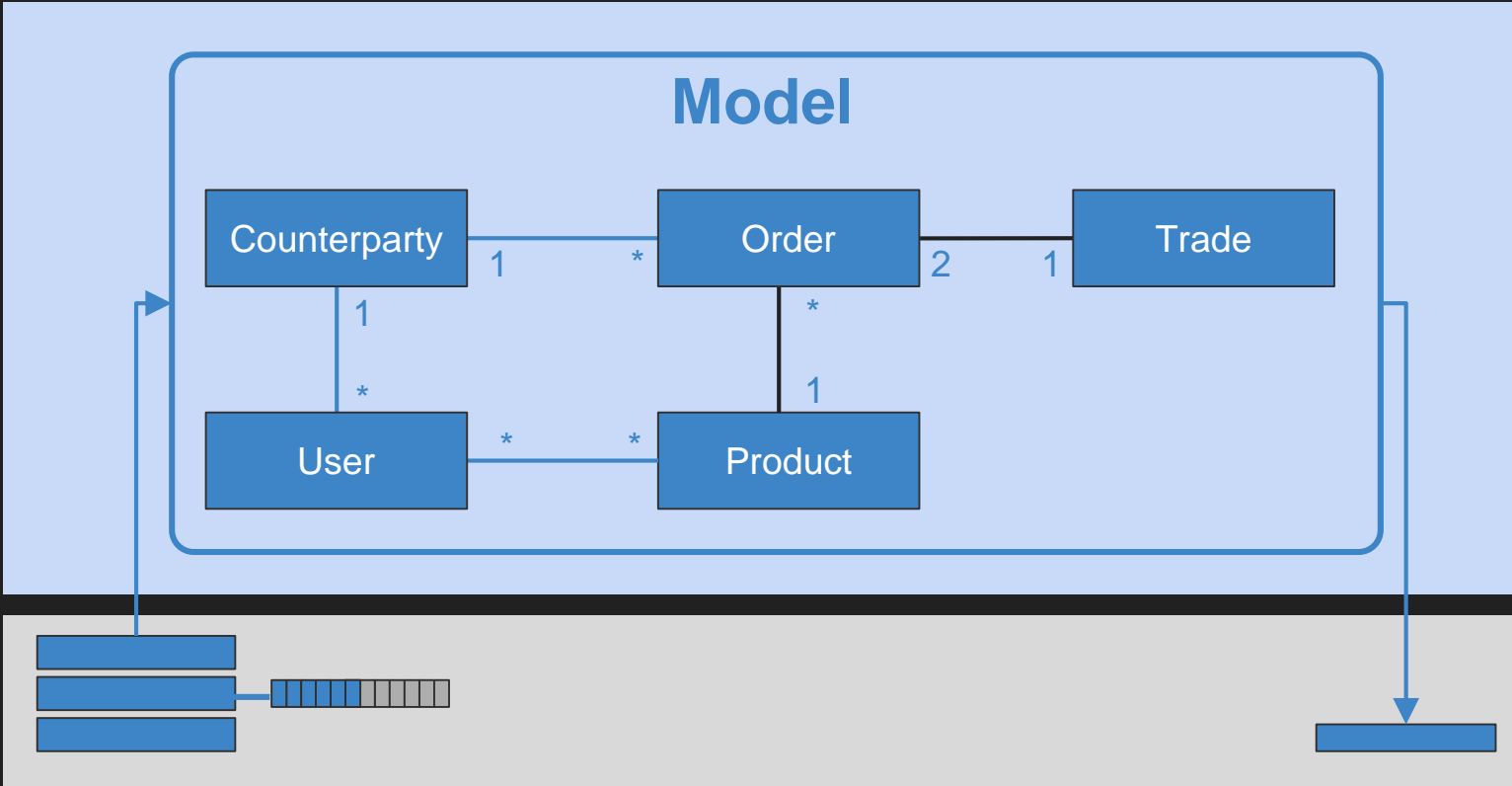
Response message is encoded



Message is sent over the network

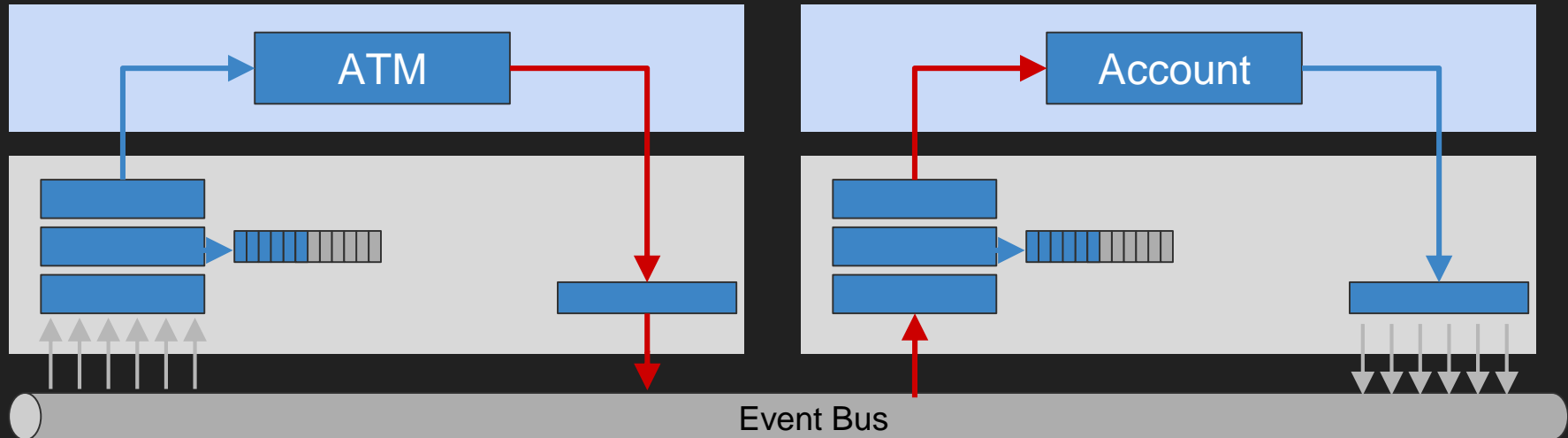


Stateful Model

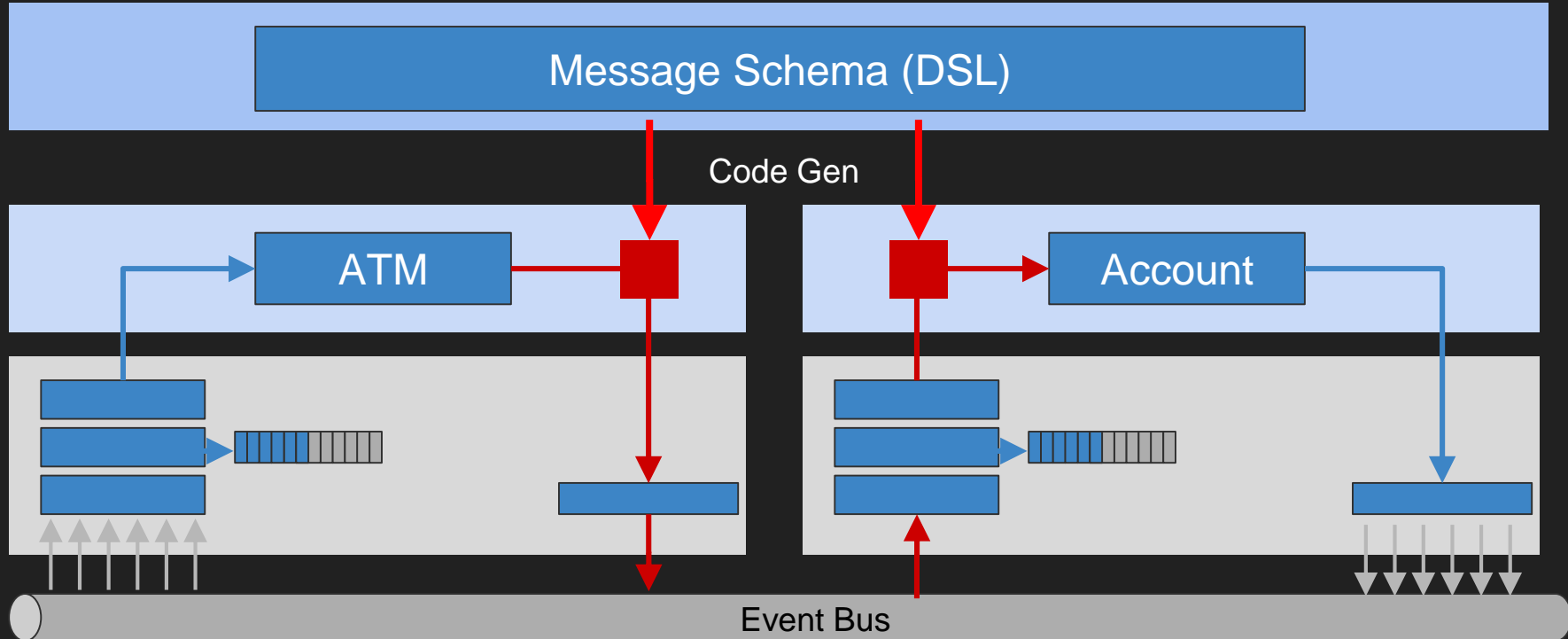


Service Contracts & Message Domain Model

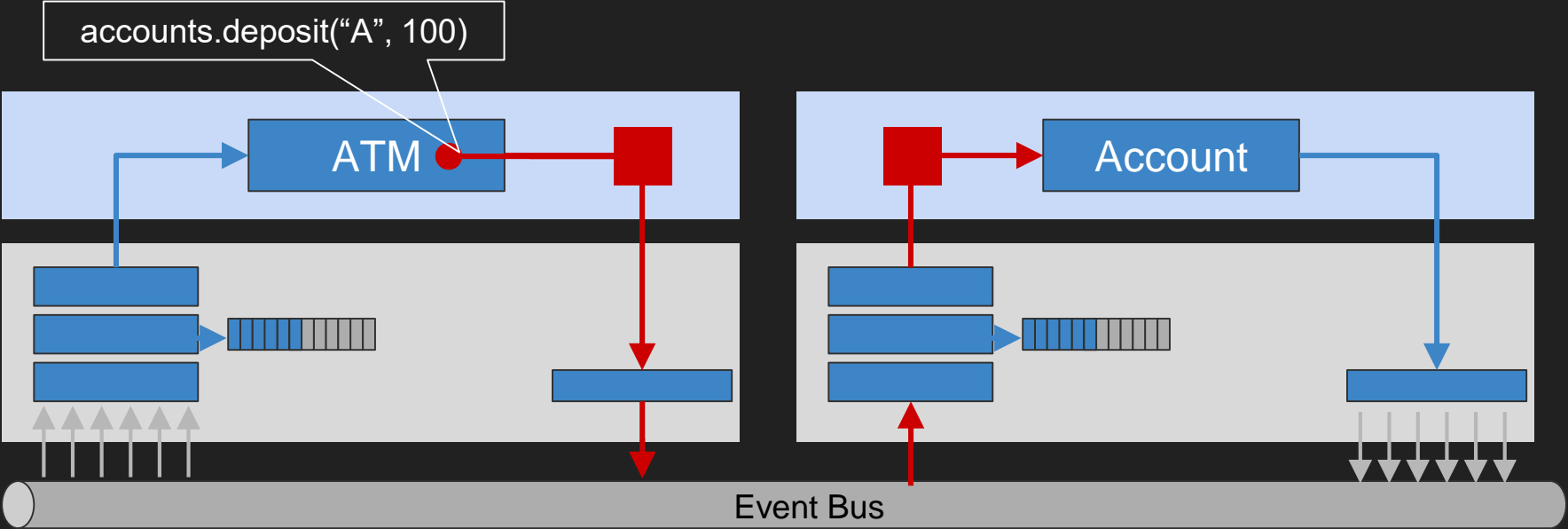
Service Binding



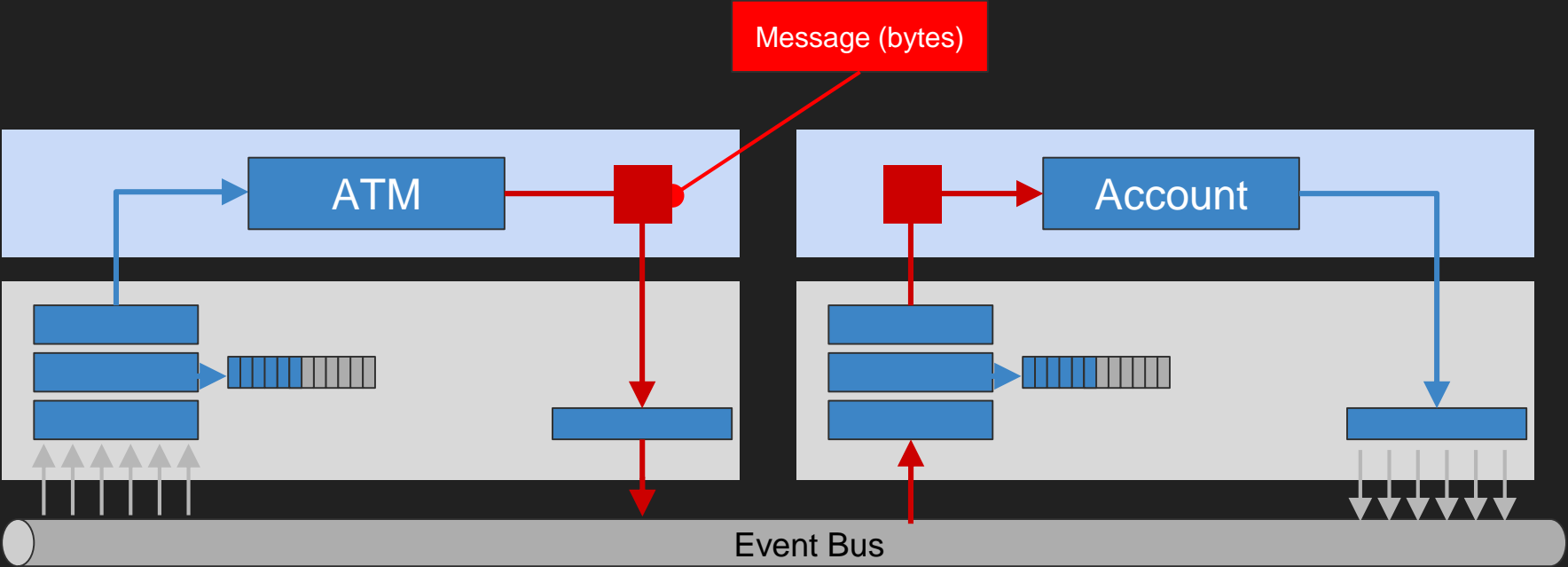
Code Generation of Marshallers



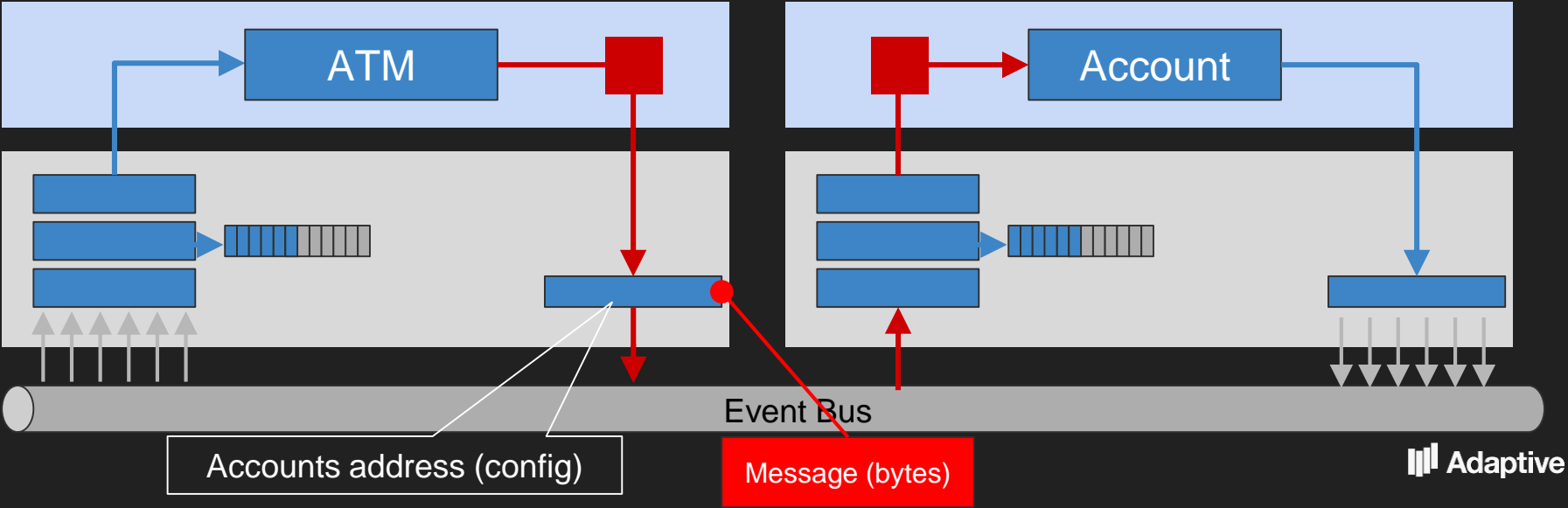
Service Interface injected in ATM Model



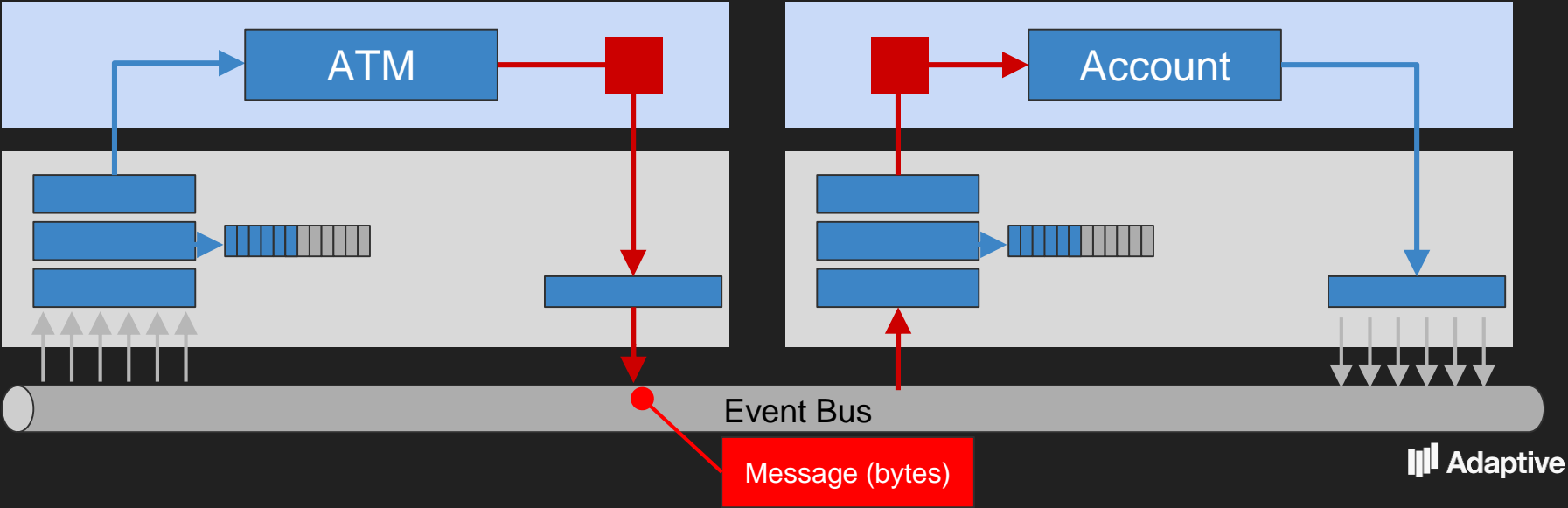
Message is encoded



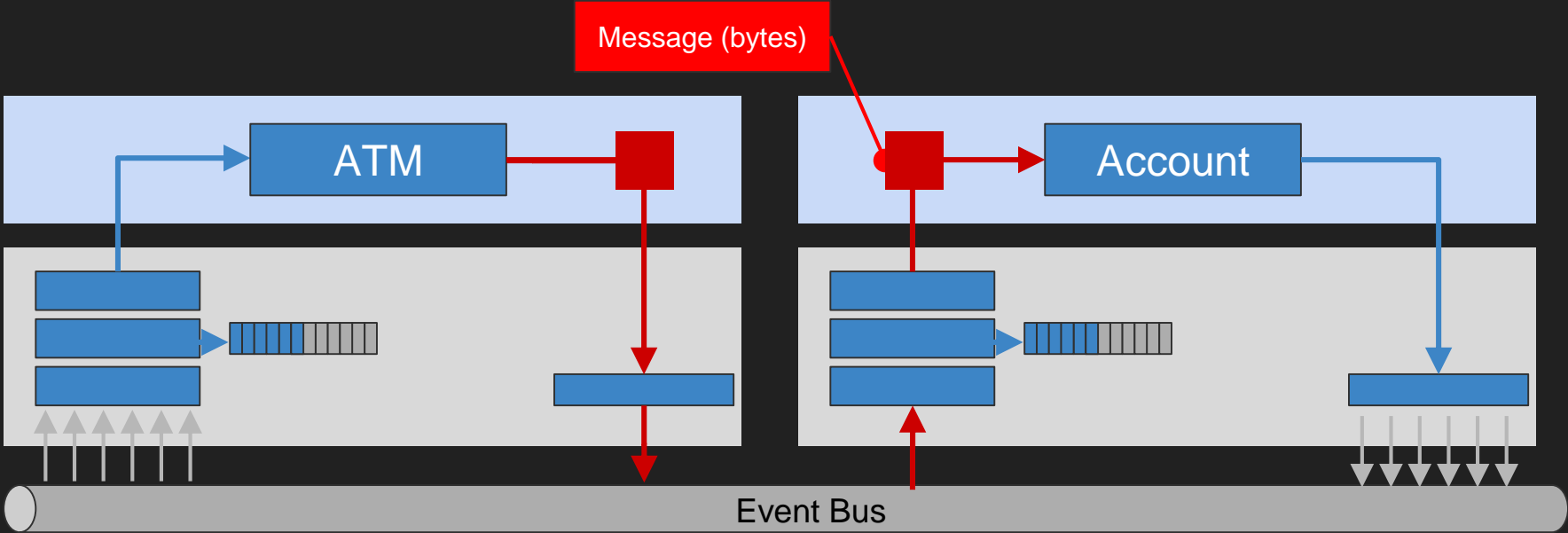
Message is routed



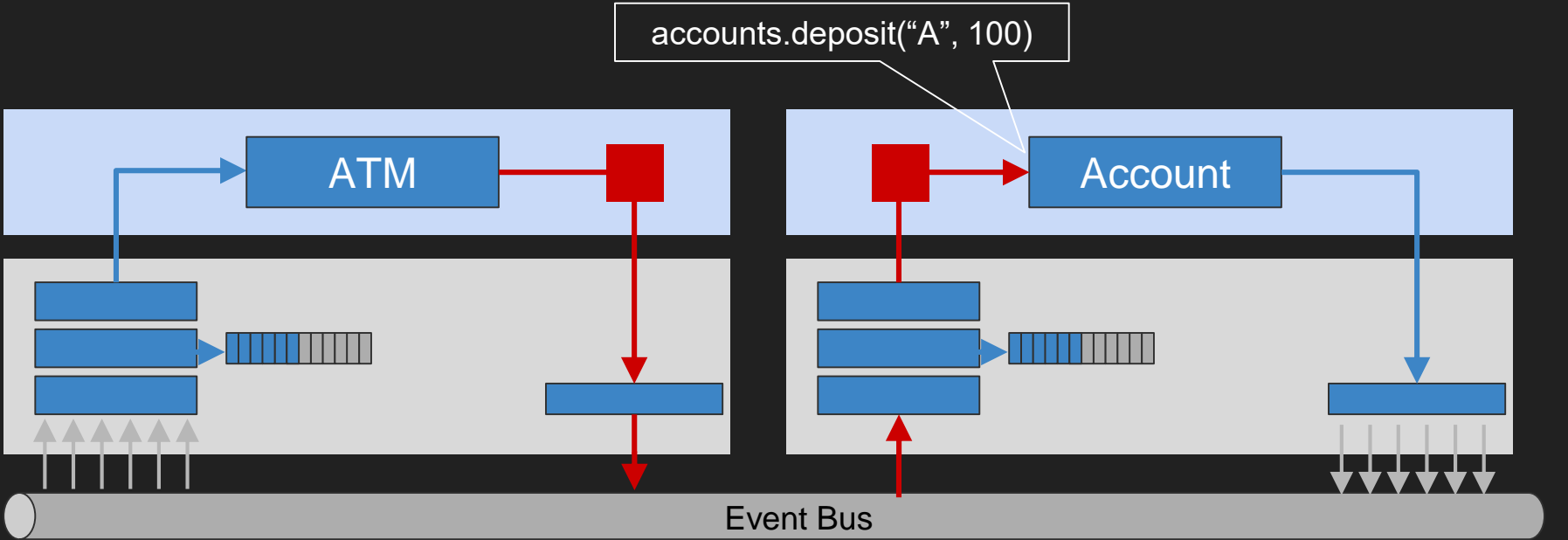
Message is sent on the wire



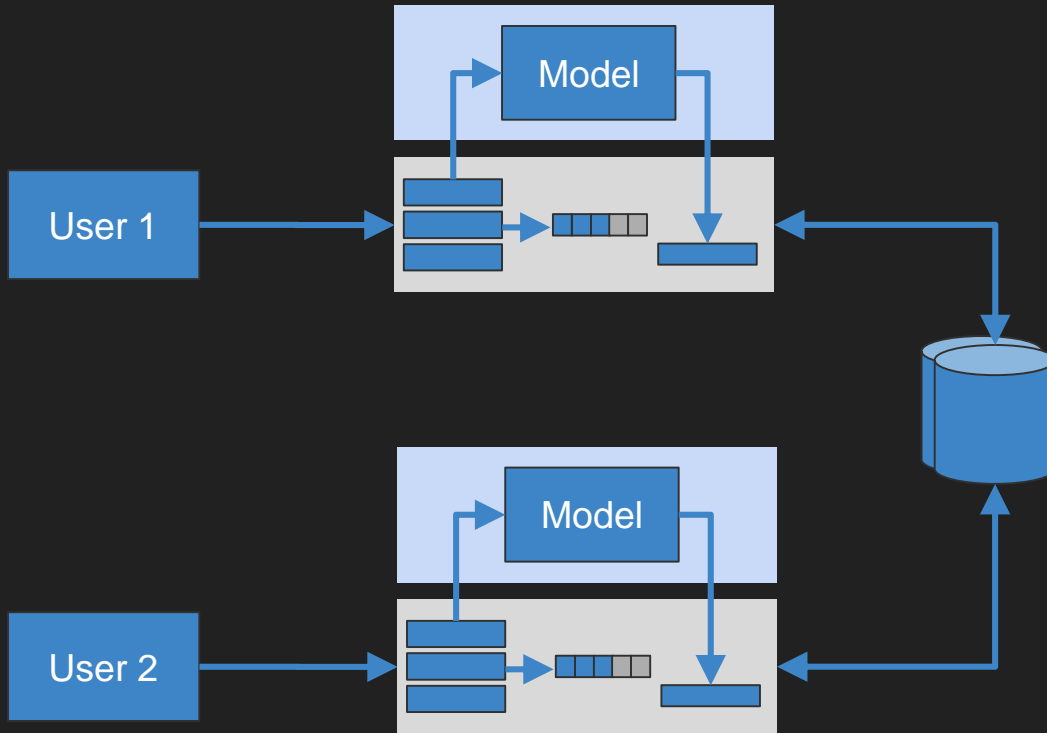
Message header contains message type



Marshaller decodes and a method on account



1st Iteration - State Synchronized via Database



Benefits

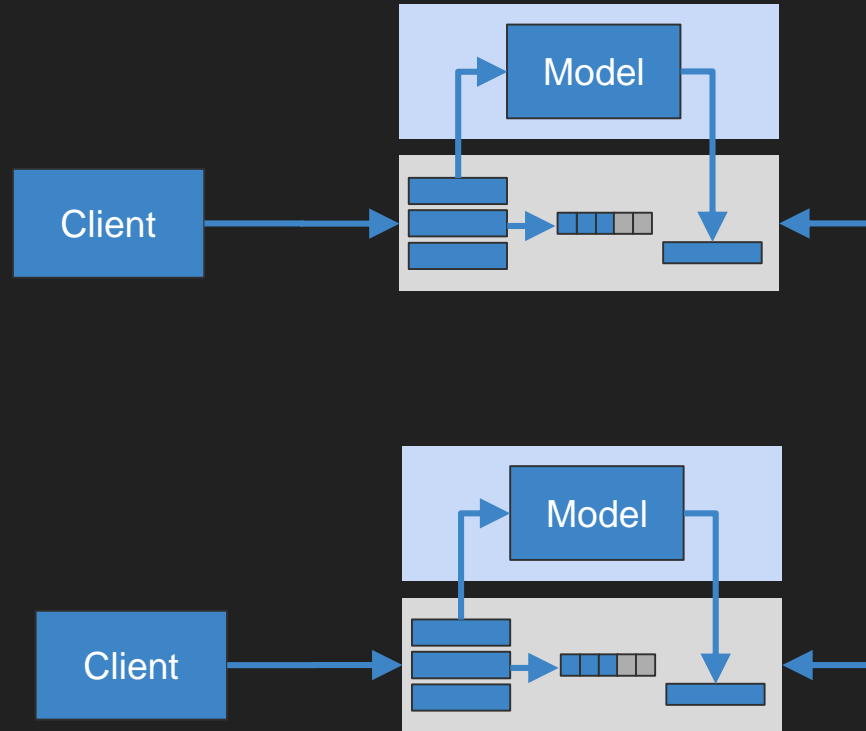
- Simple model
- Debugability

Pain points

- Synchronisation through database

Fault Tolerance & State Replication

2nd Iteration - Clustering



Fault tolerant data replication

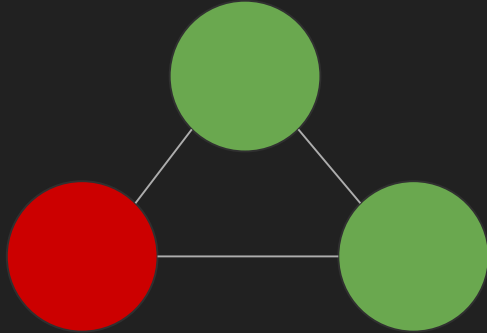
Difficult problem:

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

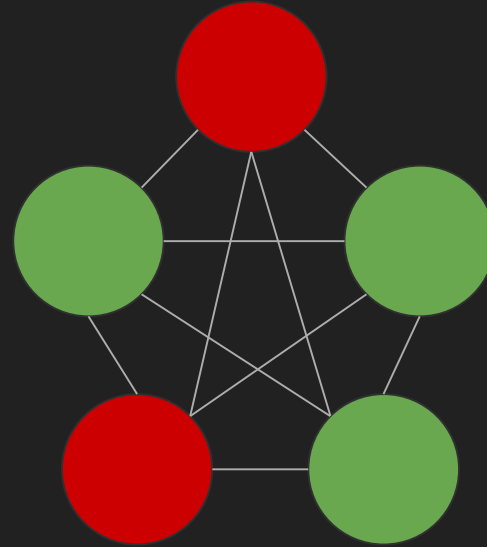


Distributed consensus problem

RAFT

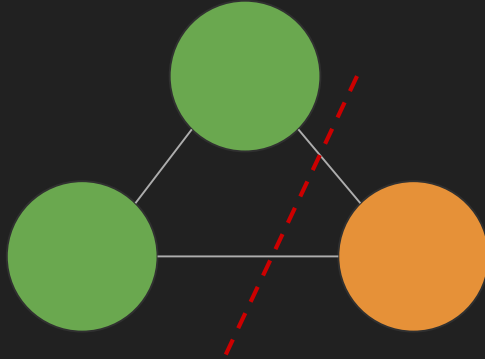


Tolerates fault of 1 node

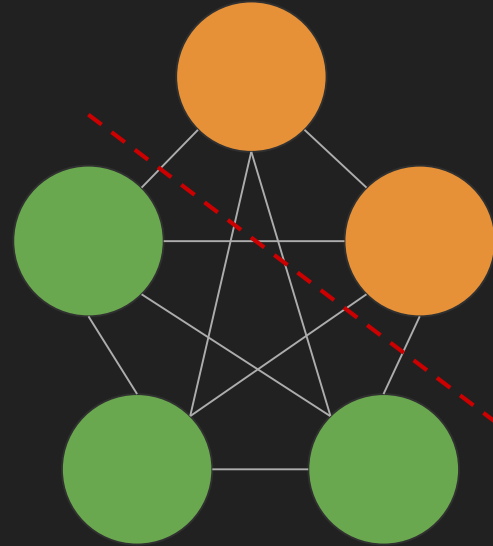


Tolerates fault of 2 nodes

RAFT - Network failure



1 node partitioned from the other nodes



2 nodes partitioned from the other nodes

Leader election



Follower



Follower



Follower

Leader election



Follower



Leader



Follower

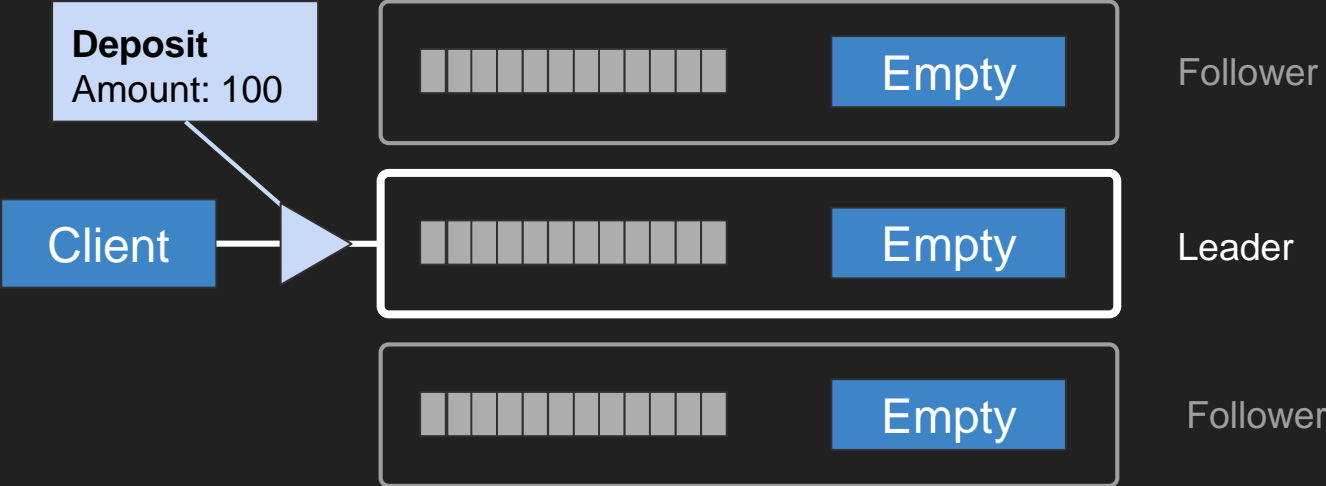
Leader discovery



Client connects



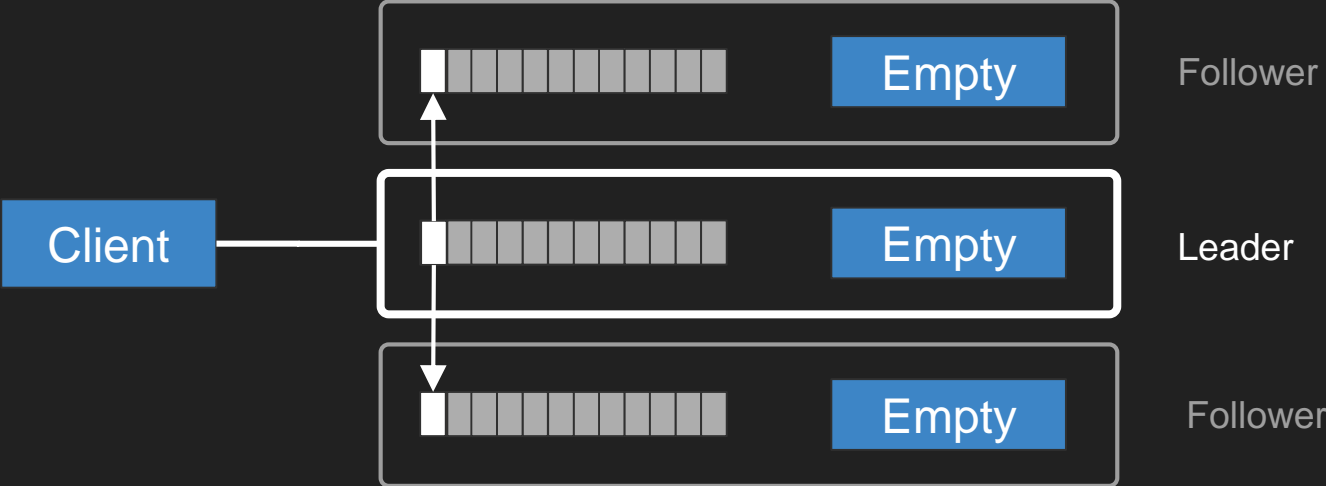
Client sends a command



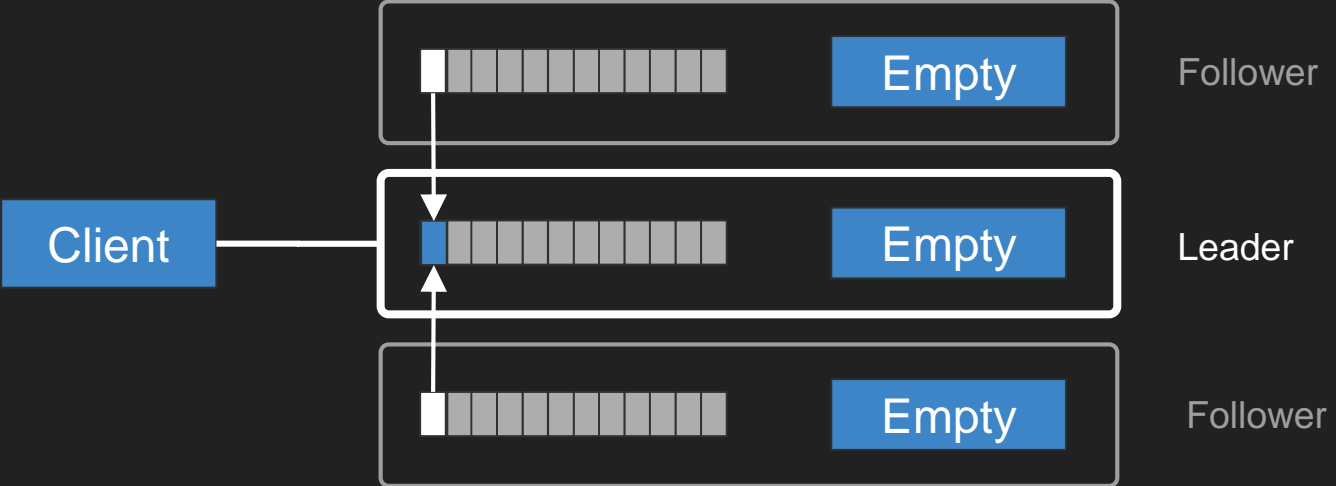
Leader writes command to its log



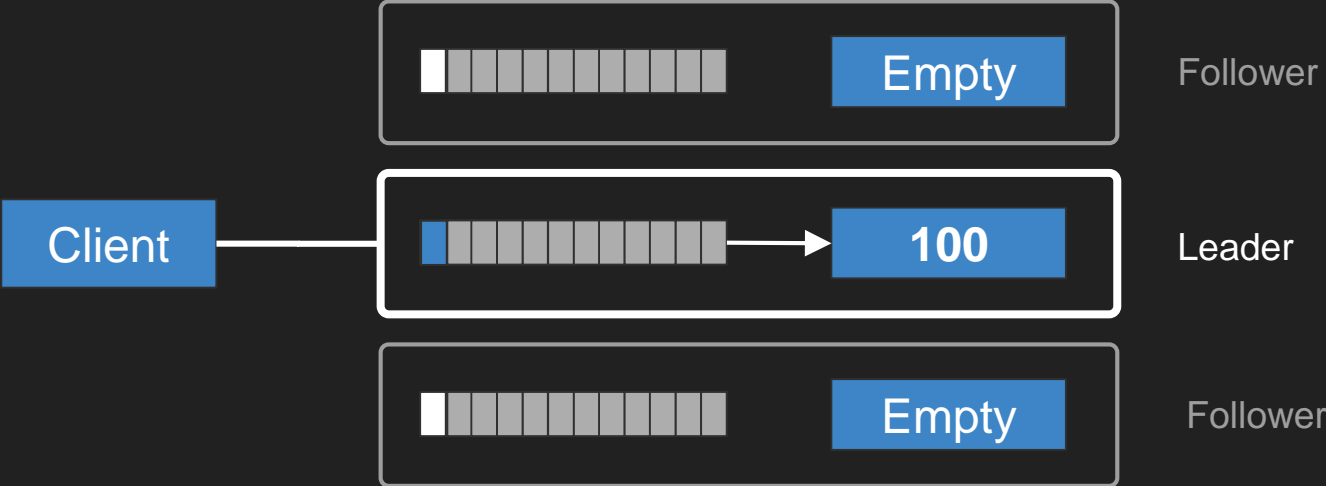
Command is replicated



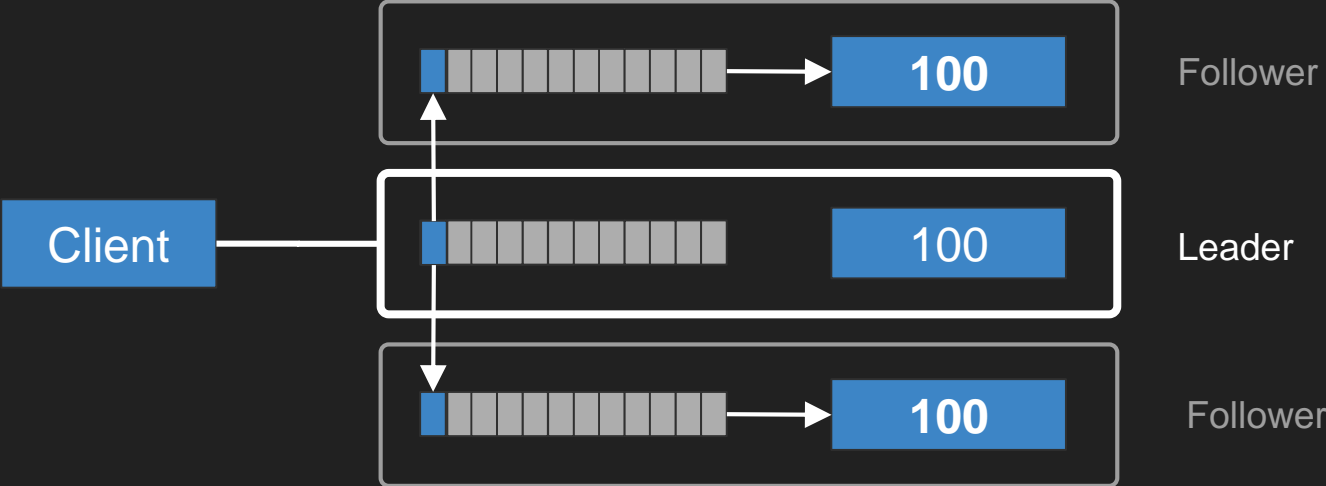
Followers Acknowledge



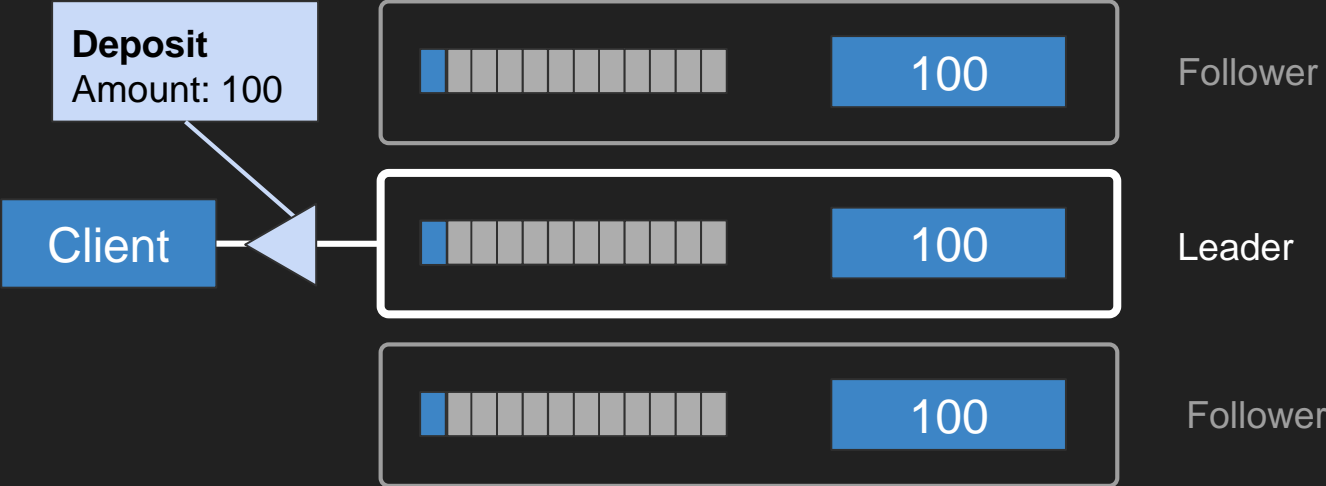
Leader commits command and applies to model



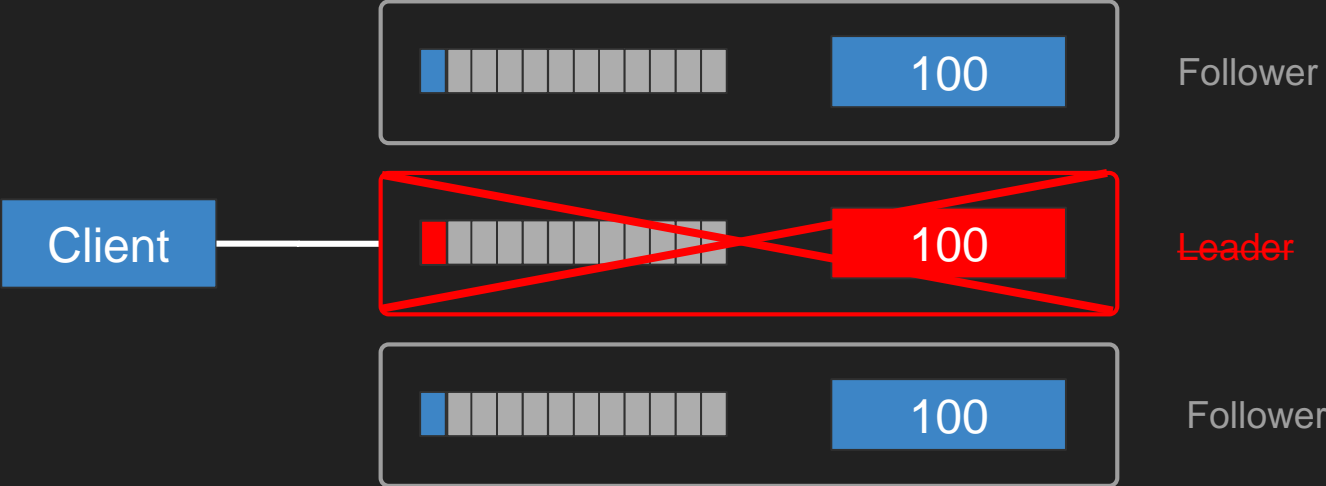
Leader notifies followers to commit



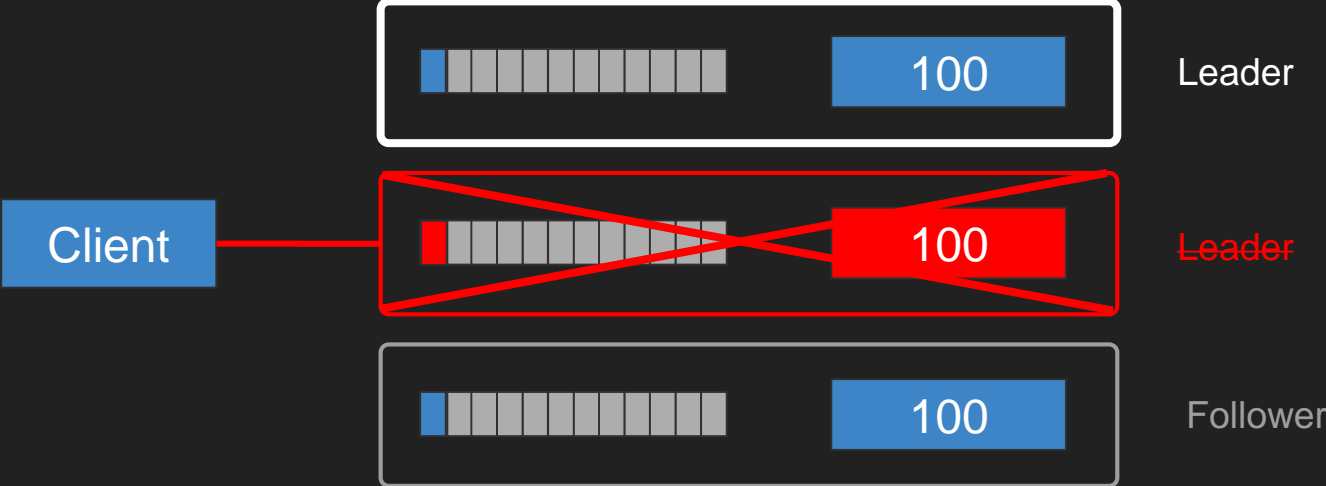
Leader notifies client



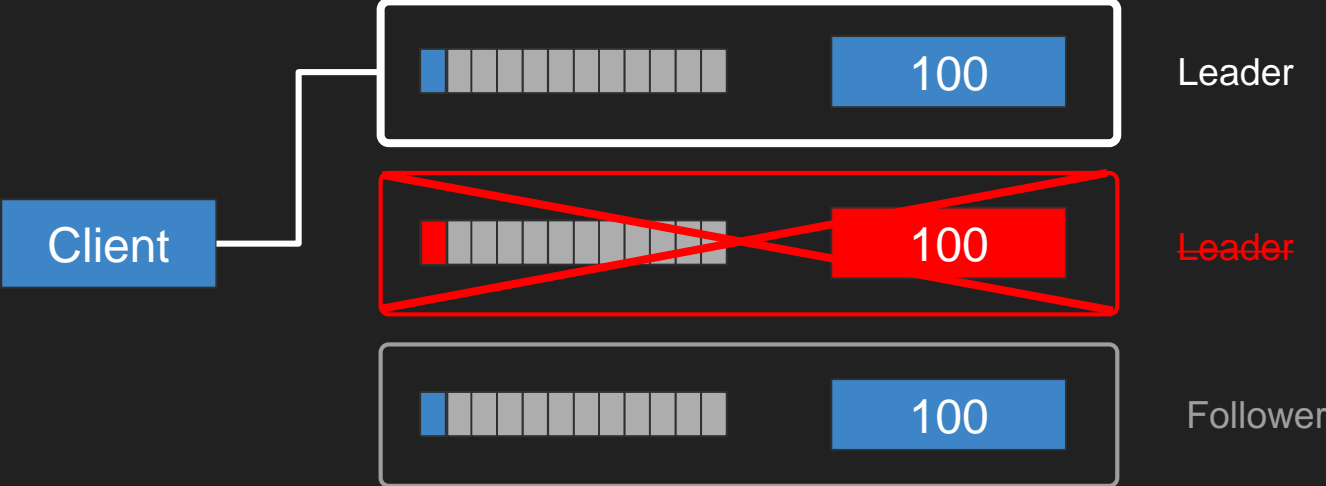
Leader dies



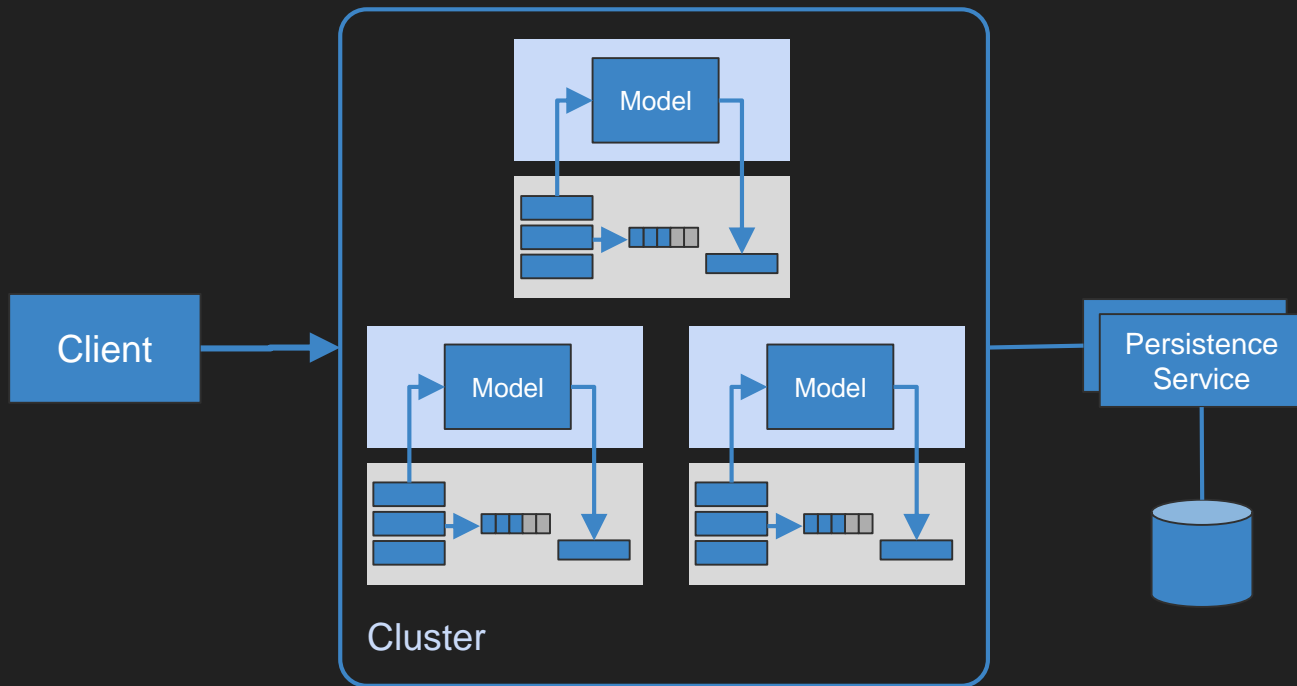
New leader election



Client times out and connects to new leader



2nd Iteration - Clustering



Benefits

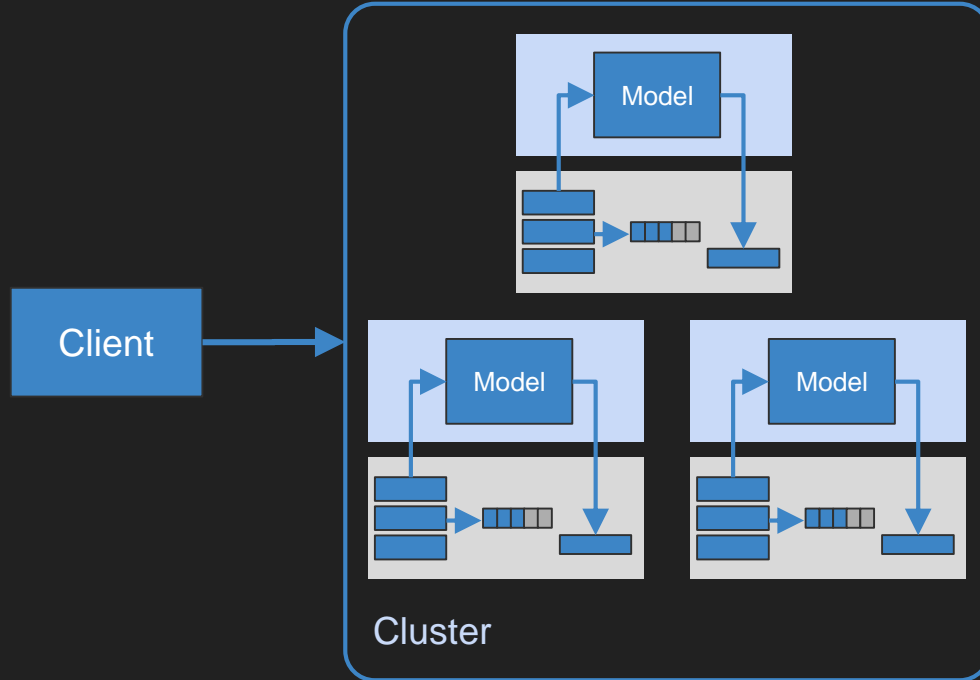
- Fault tolerant
- Consistent model

Pain points

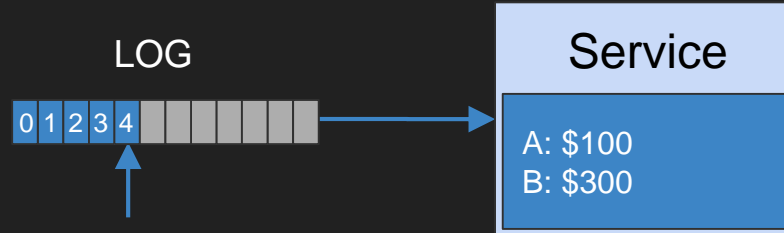
- Database as golden source

Durability

3rd Iteration - Durable Services



System restart

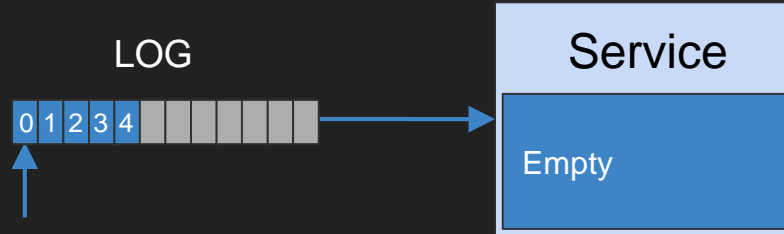


Stop the service

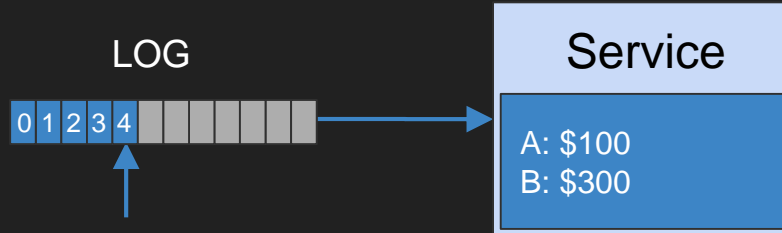
LOG



Start the service



Replay commands



Model is back in the same state because it's deterministic.

Managing the log

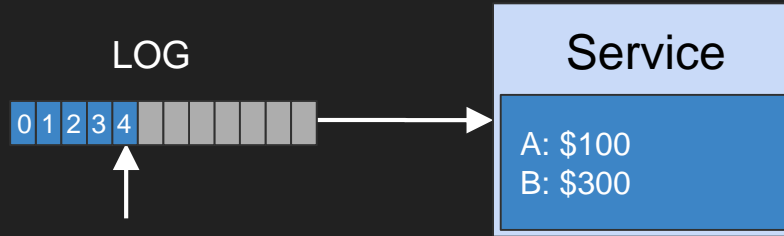
OK if

- Log grows slowly
- Your model has short lifetime

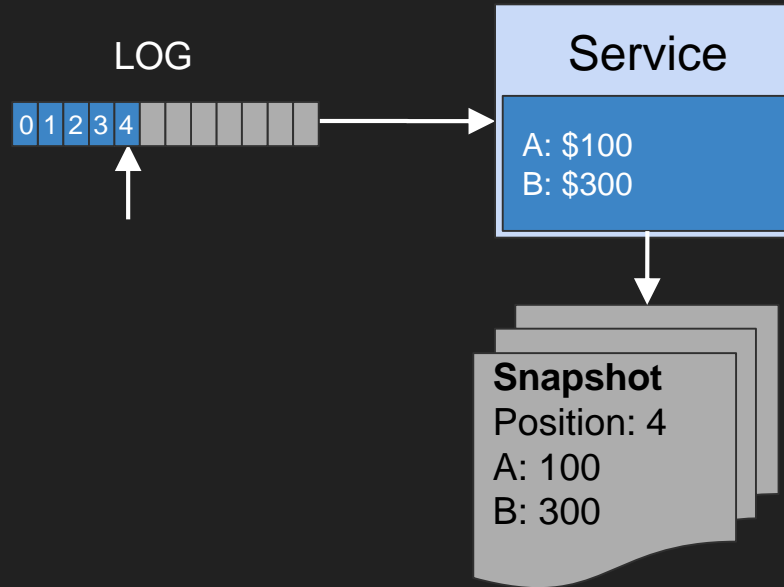
Problem: the size of the log affects recovery time

Need log compaction.

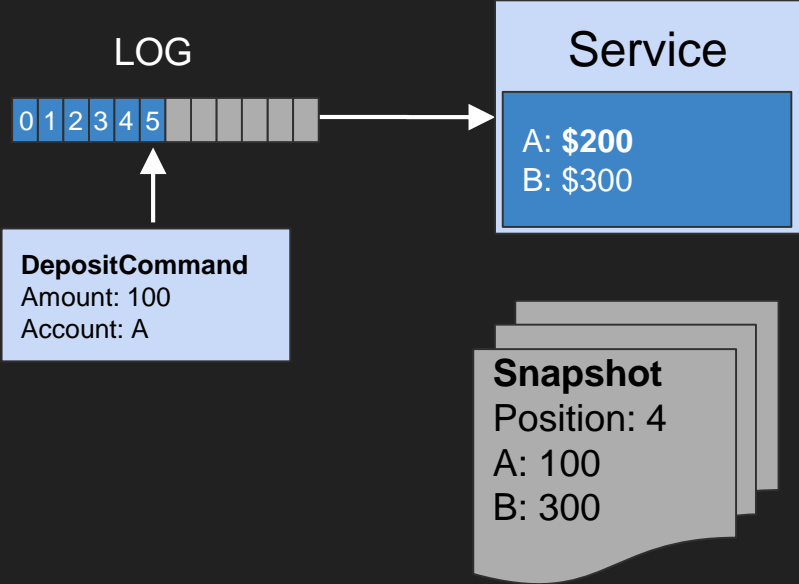
Snapshotting



Write state to disk

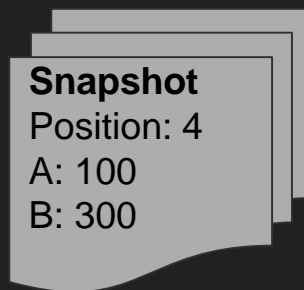


New Commands Are Processed

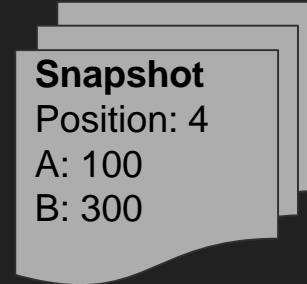
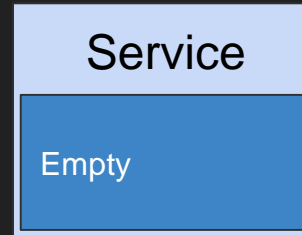
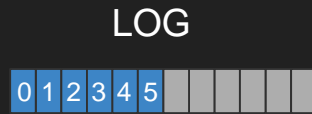


Stop the Service

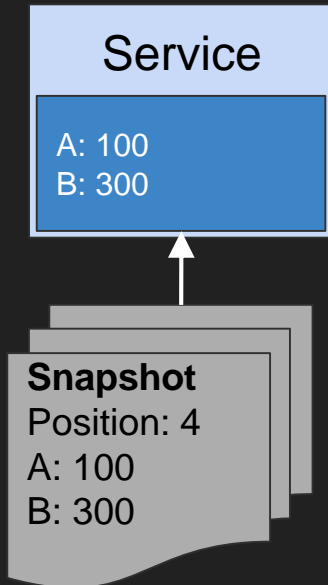
LOG



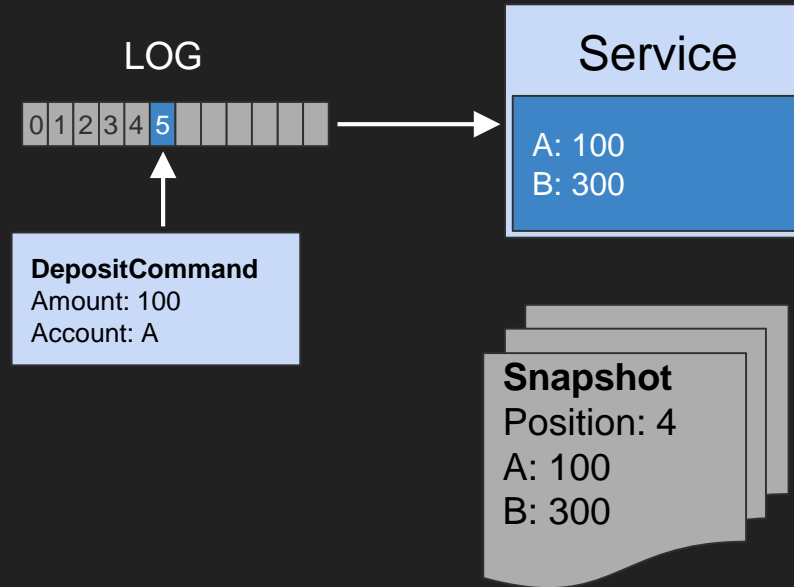
Restart the service



Restore the snapshot



Replay only commands from after the snapshot

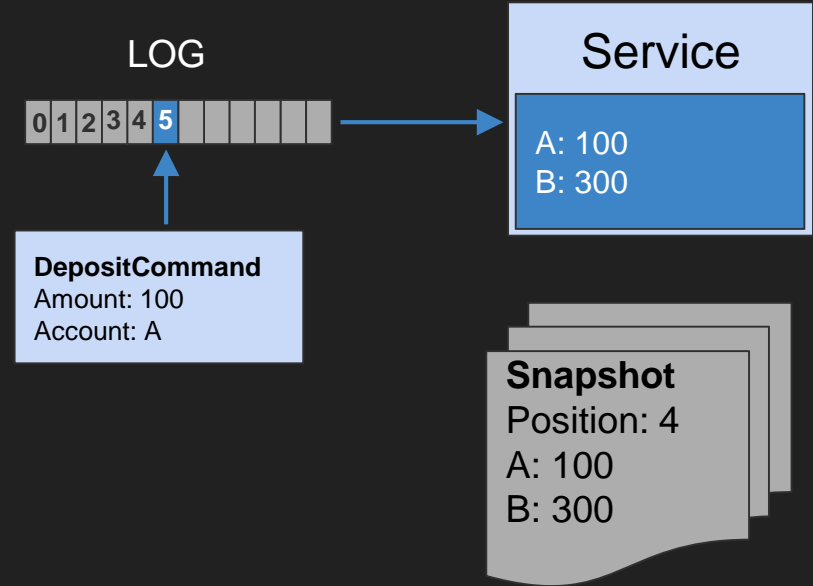


Snapshotting

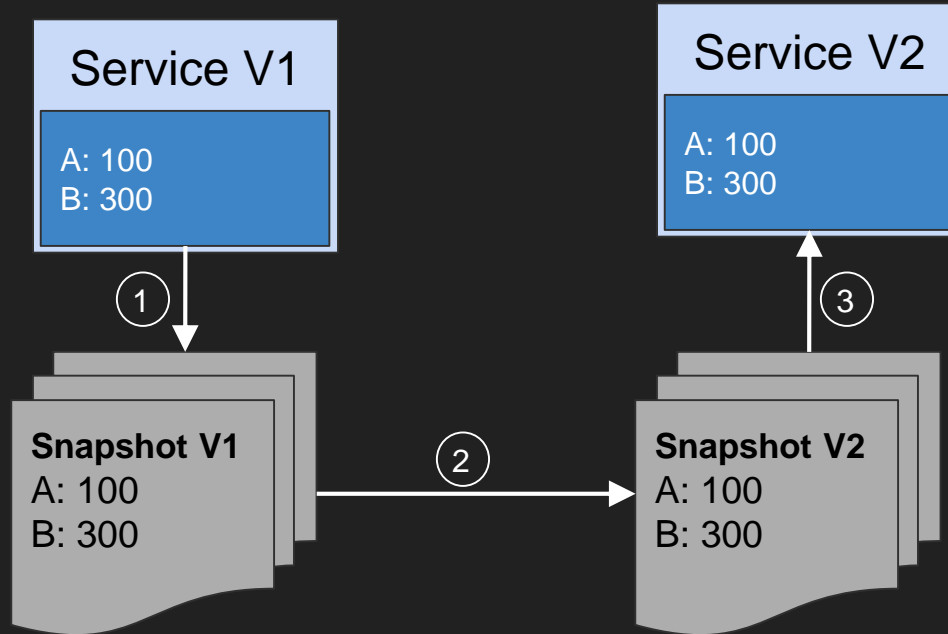
System is back in the same state.

Snapshotting makes it possible to trim the log.

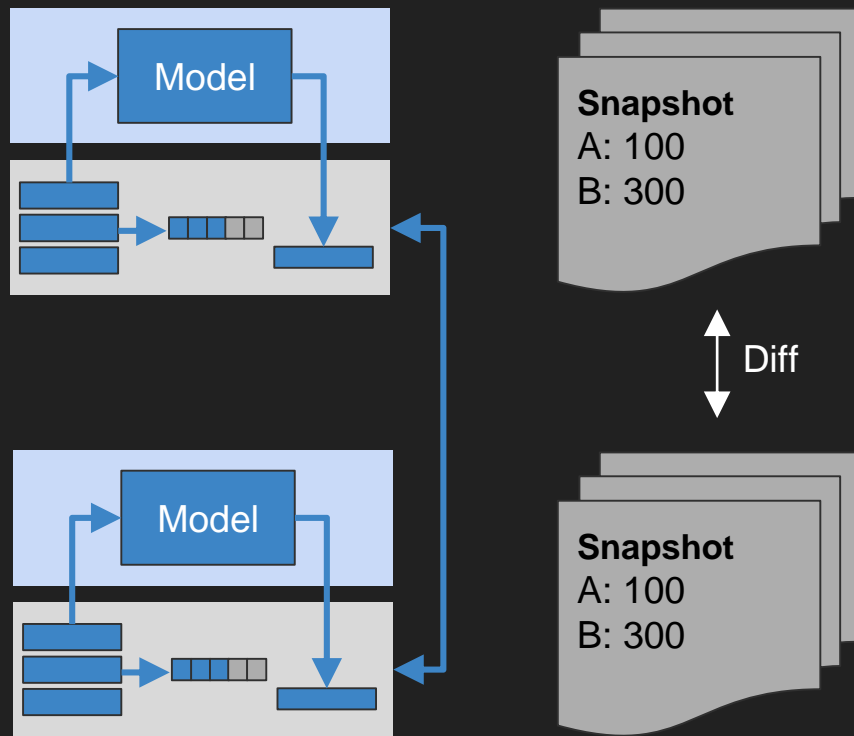
With a snapshot at position N, we can archive the log up to the same position.



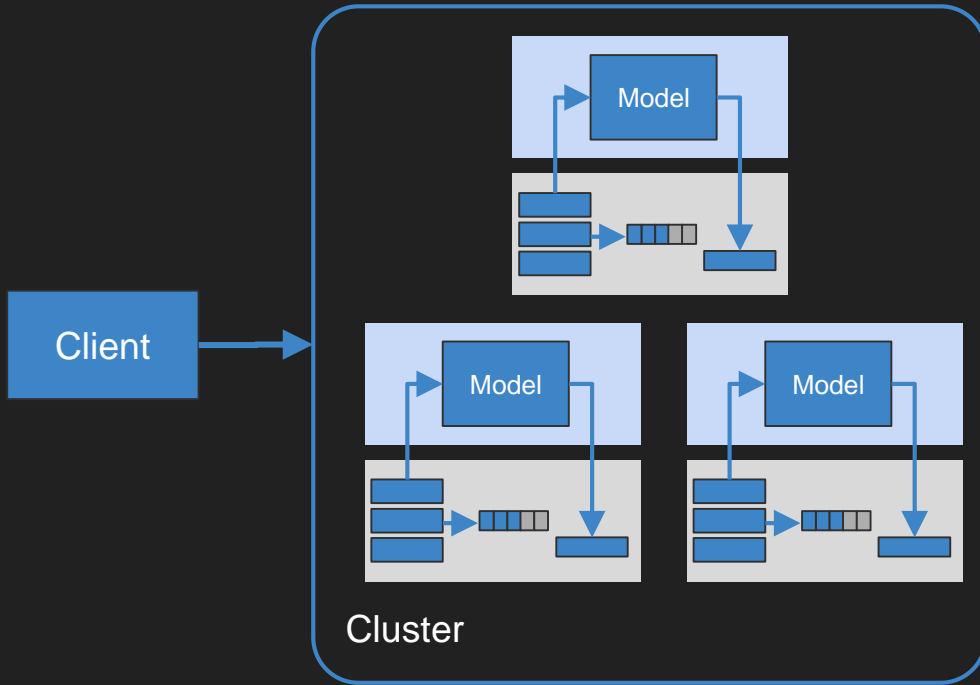
Upgrades



Consistency Check



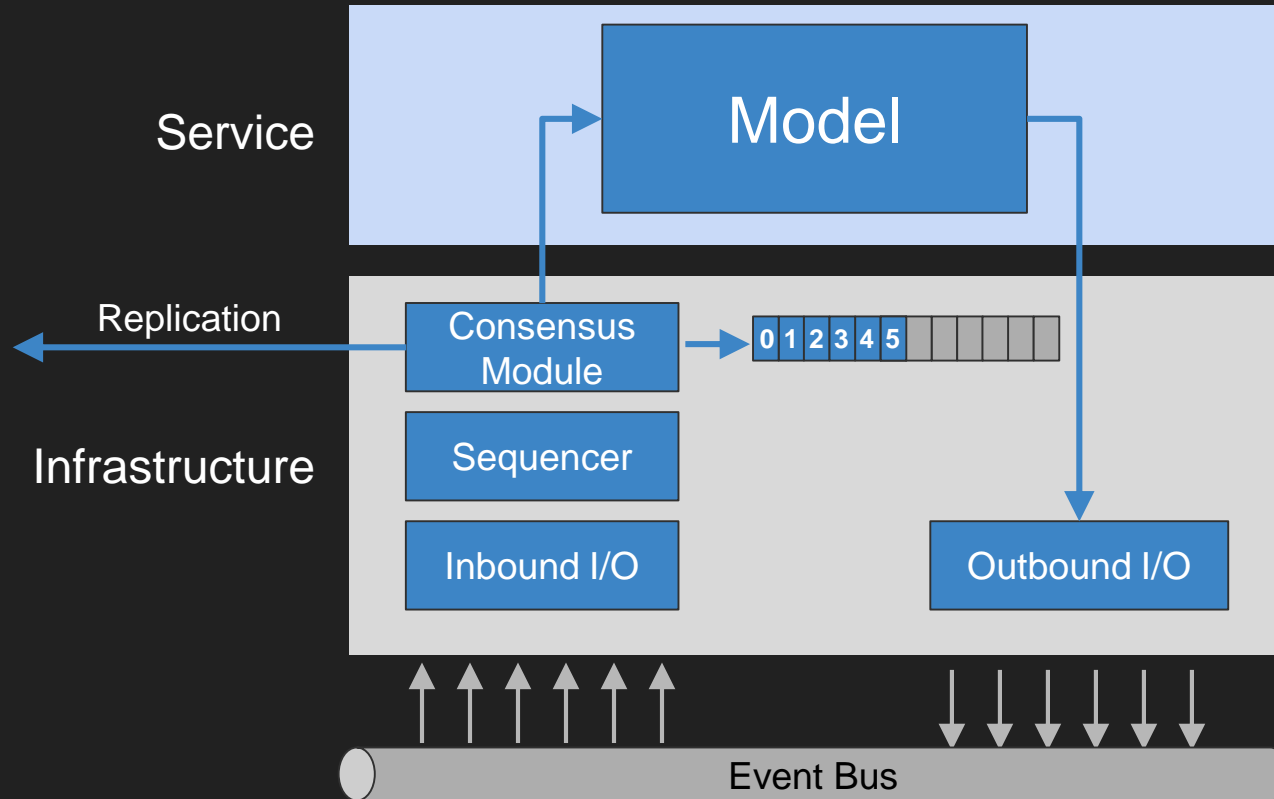
Iteration 3 - Durable Services



Benefit

State is managed in one place

Reactive Application Server



Simple Model
Easy To Debug
Fast

Fault Tolerance
State Replication
Durability
Concurrency
Messaging

Don't build the infrastructure yourself!

Open source

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



Martin Thompson



Todd Montgomery

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>

Thank you!

Any questions?

[http://weareadaptive.com/careers/
careers@weareadaptive.com](http://weareadaptive.com/careers/careers@weareadaptive.com)

Possible topics

Out of scope

- Timers & scheduling
- Testing with time
 - 23-25 hour day exemple
- Messaging
 - Code gen stuff
 - Error scenario / state machine