

Rethinking Services with Stateful Streams

Ben Stopford
@benstopford



THIS IS ME

- ENGINEER
AT CONFLUENT
- Ex THOUGHTWORKS
+ UK FINANCE

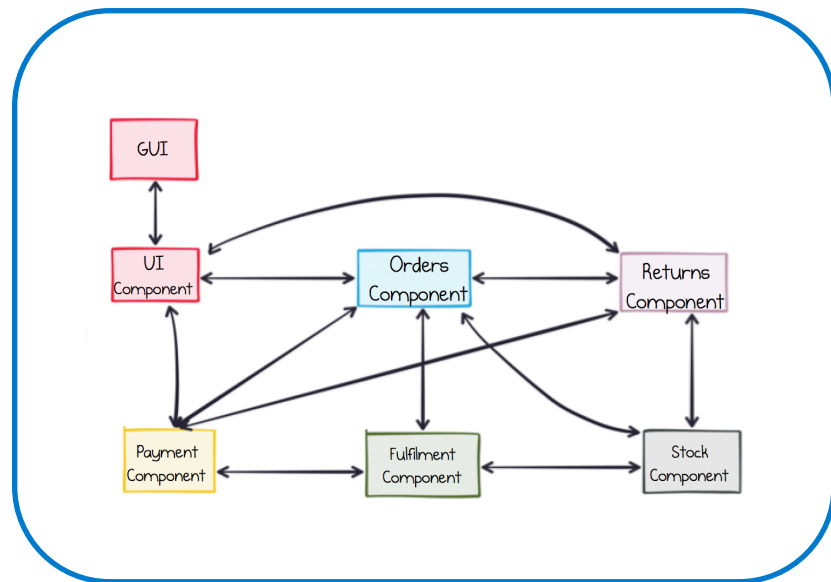


What are
microservices
really
about?

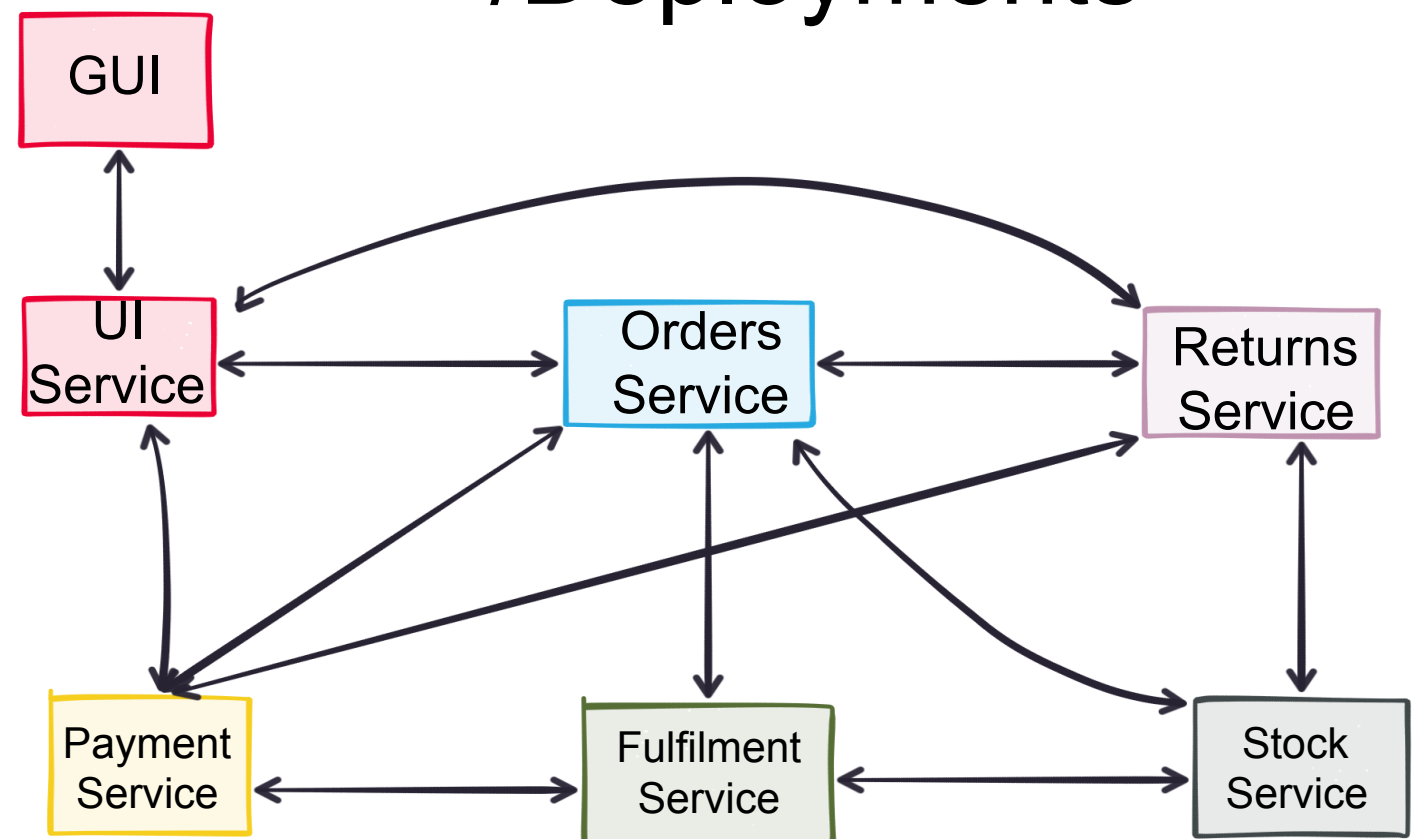


Splitting the Monolith

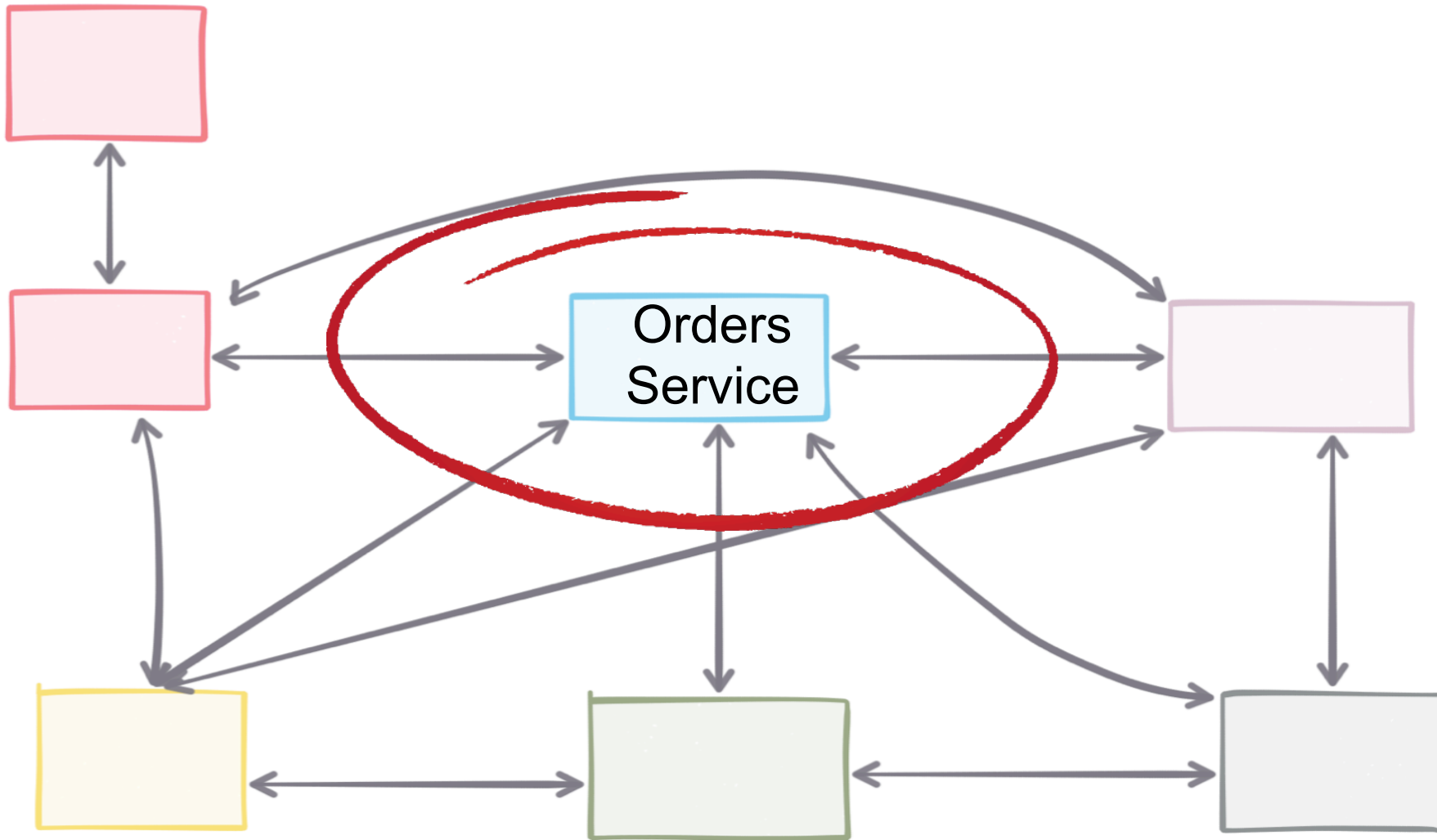
Single Process
/Code base
/Deployment

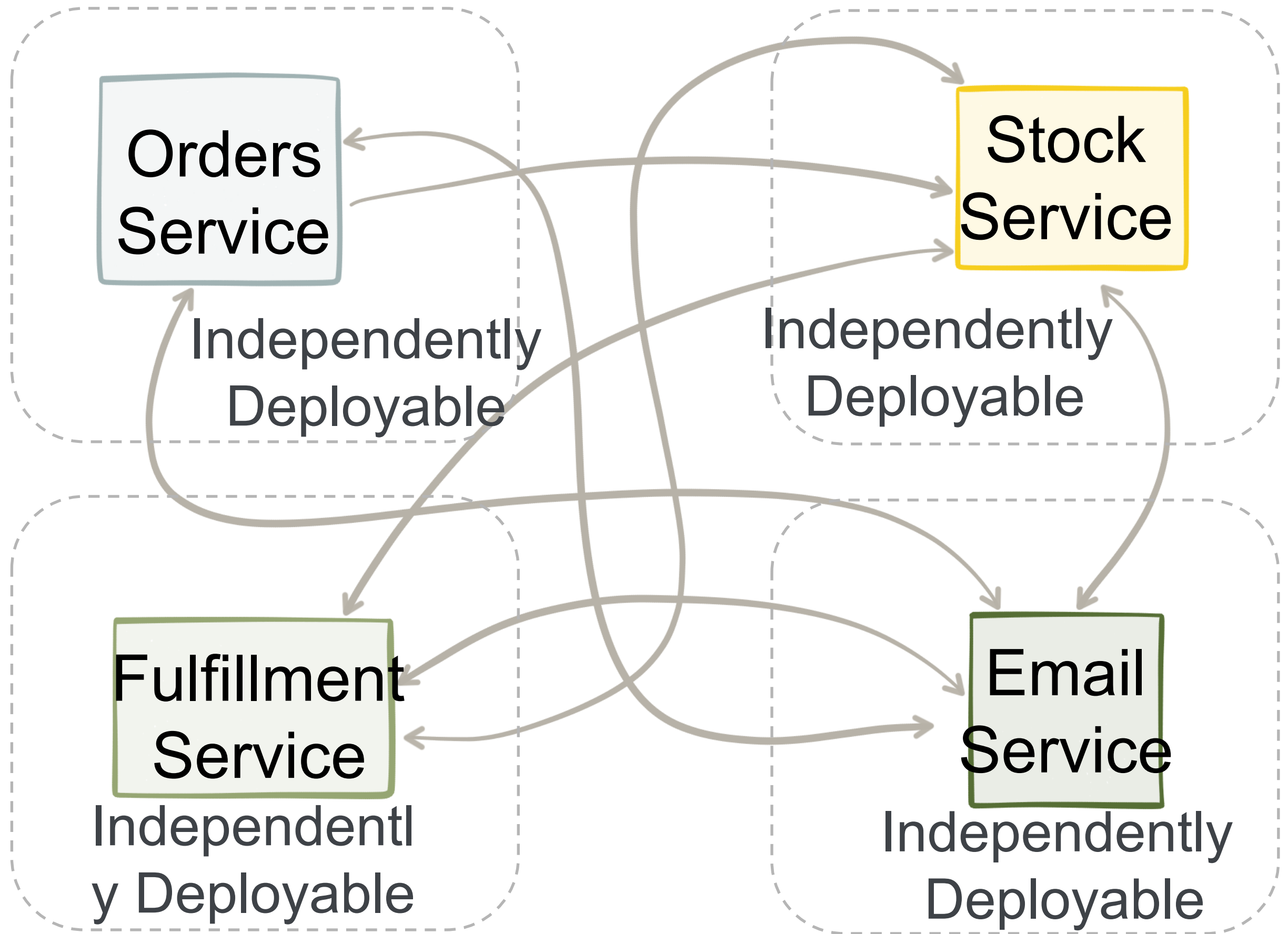


Many Processes
/Code bases
/Deployments



Autonomy?

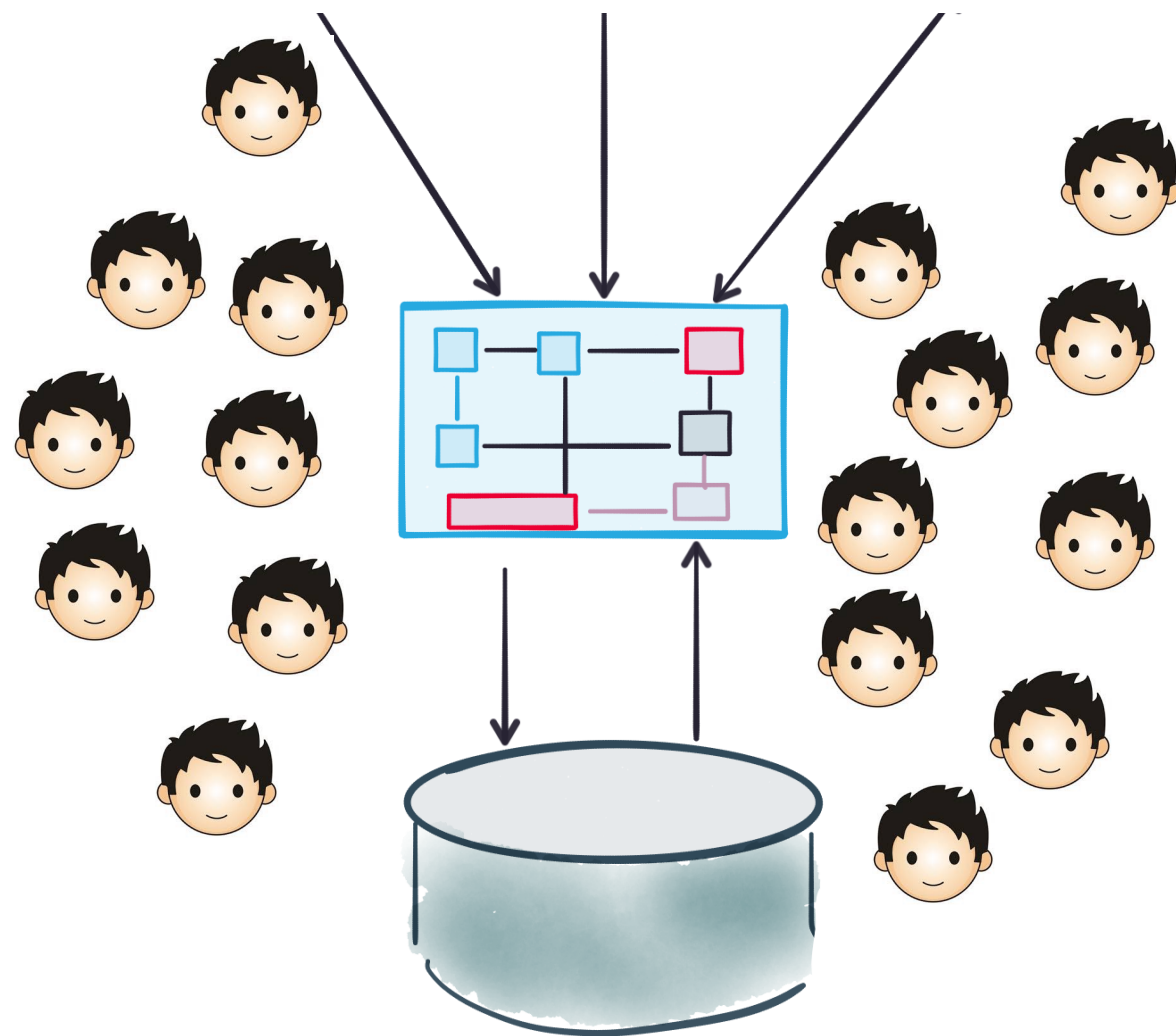




Independence is
where services
get their value

Allows Scaling

Scaling in terms of people

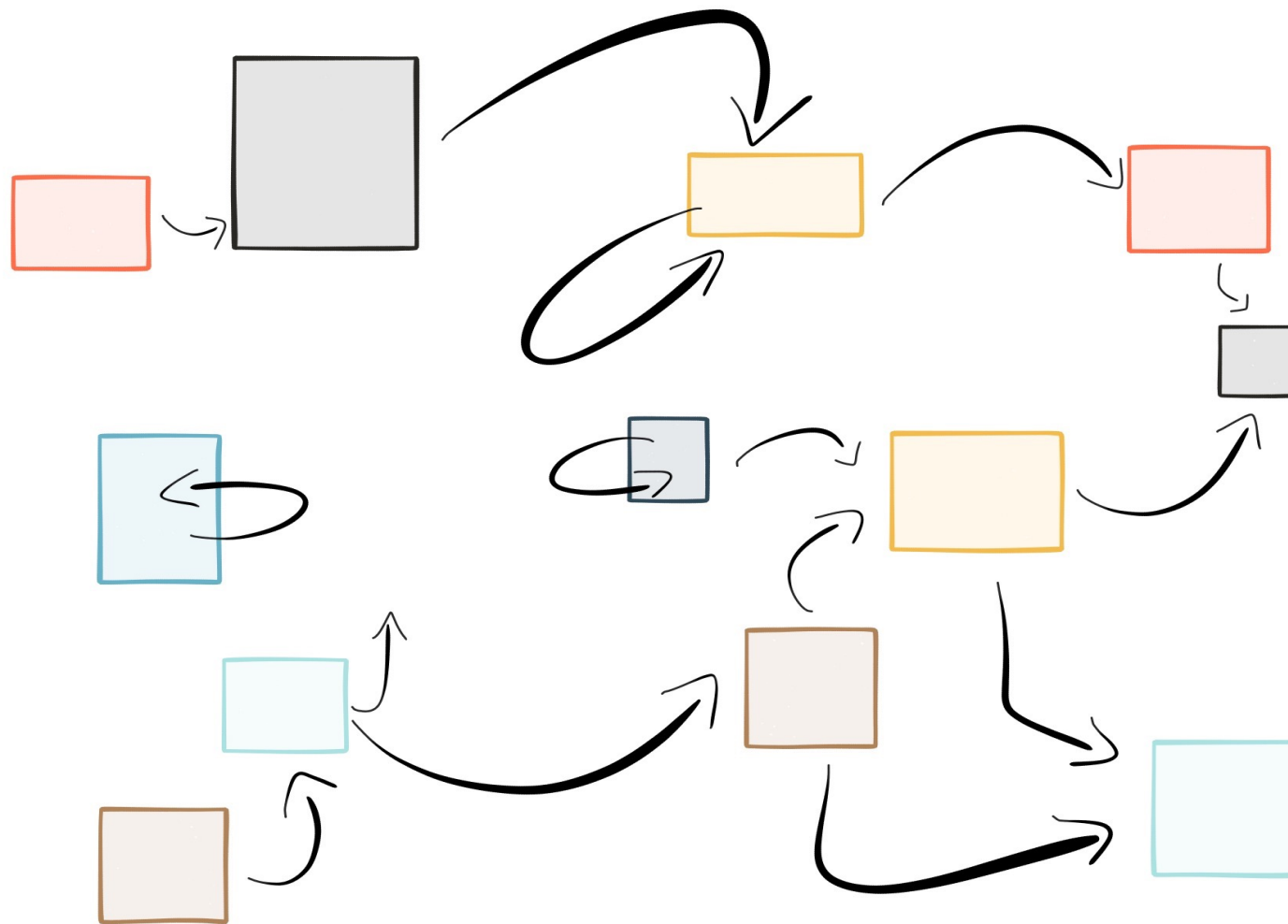


What happens when we grow?

Companies are inevitably a collection

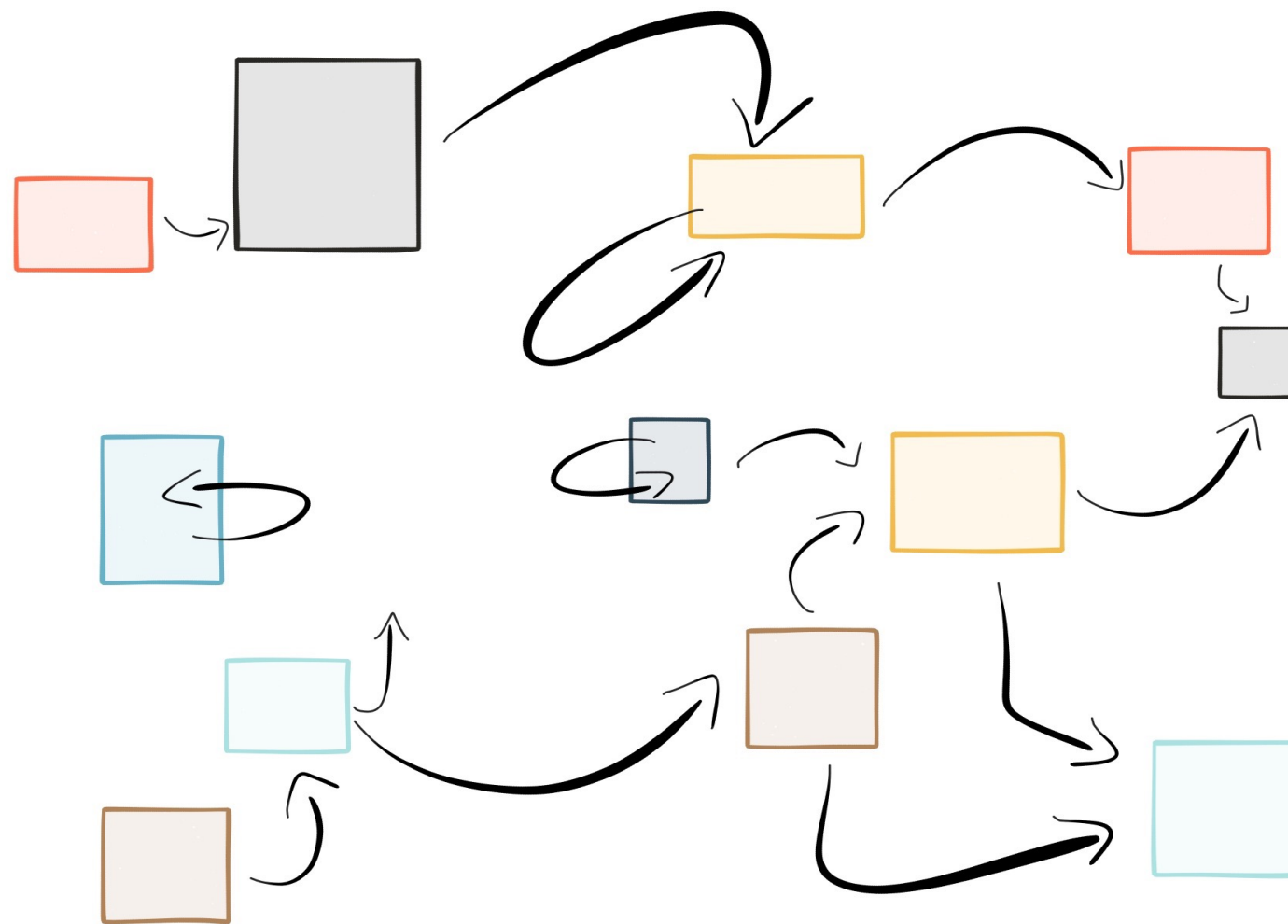
of applications

They must work together in some degree

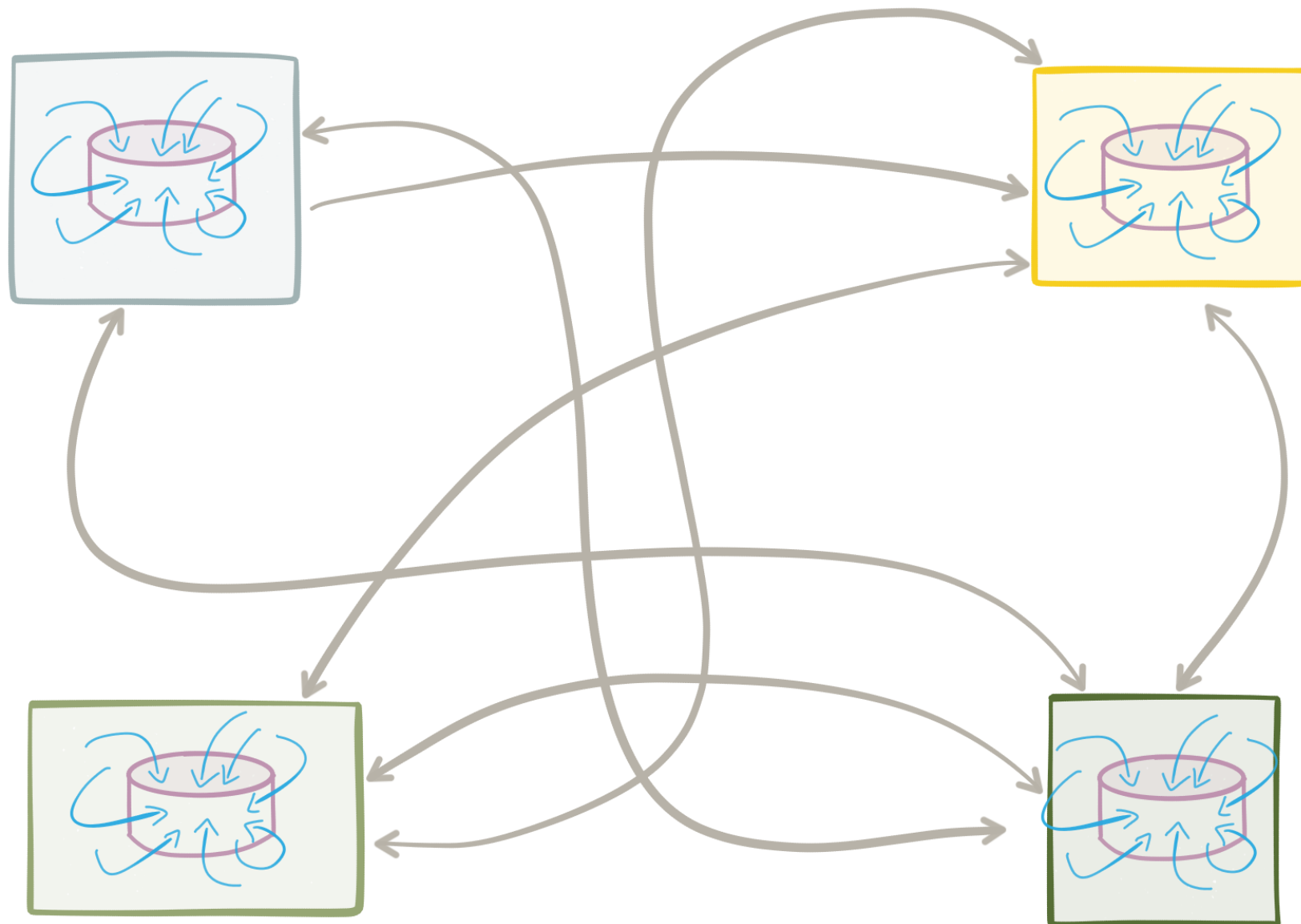


Interconnection is an afterthought

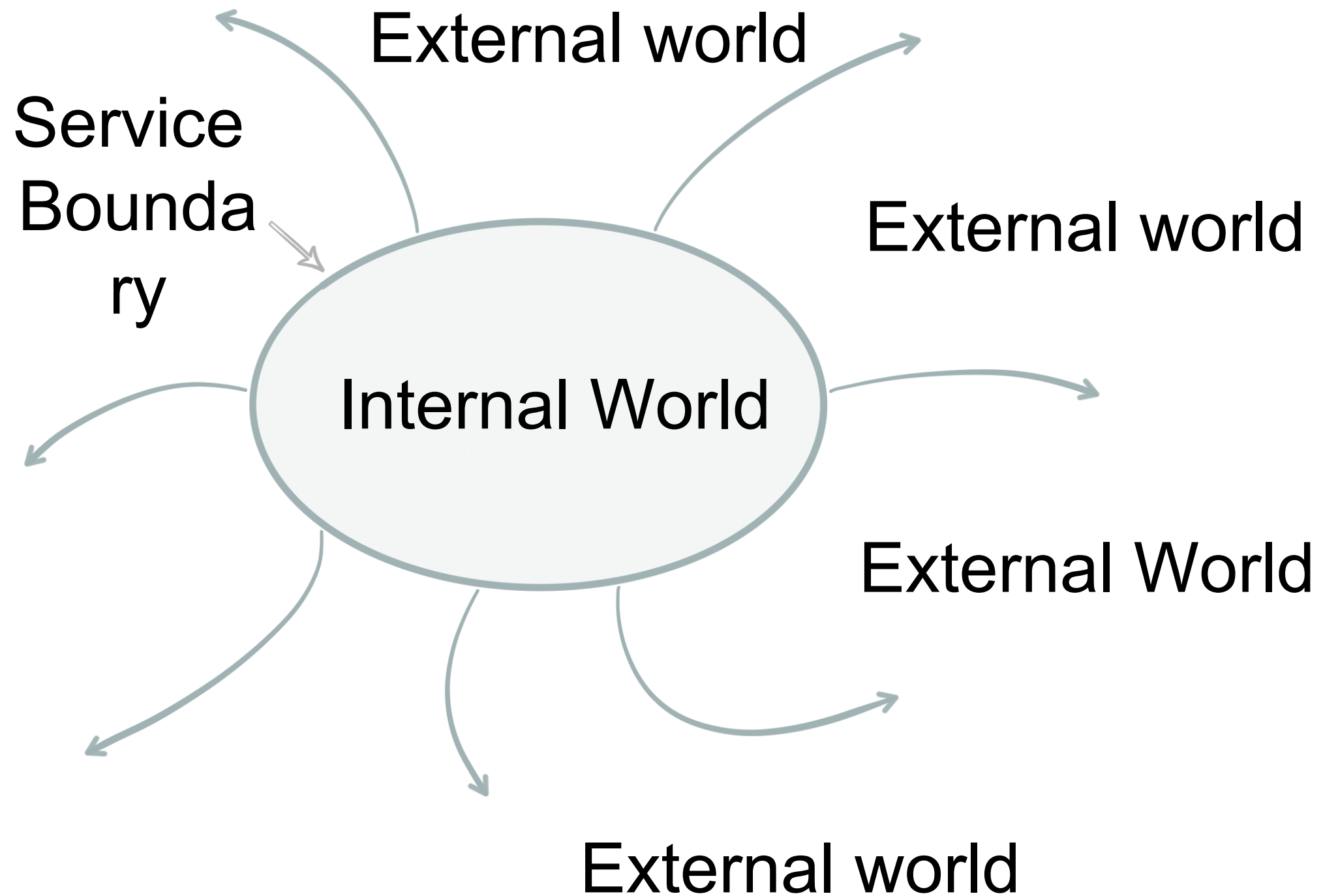
FTP / Enterprise Messaging etc



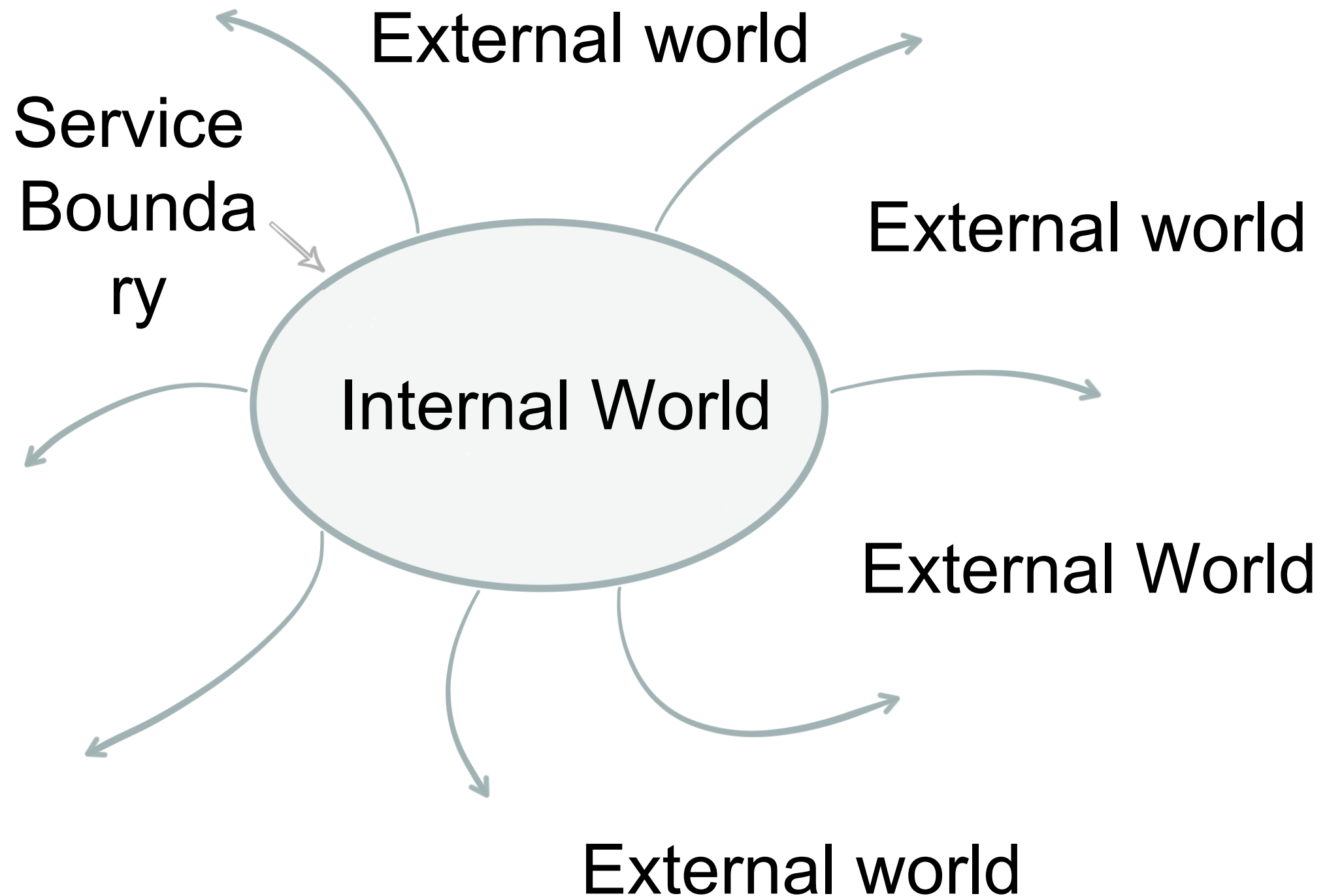
Microservices / SOA are patterns for multi-team architectures



Services Force Us To Consider The External World

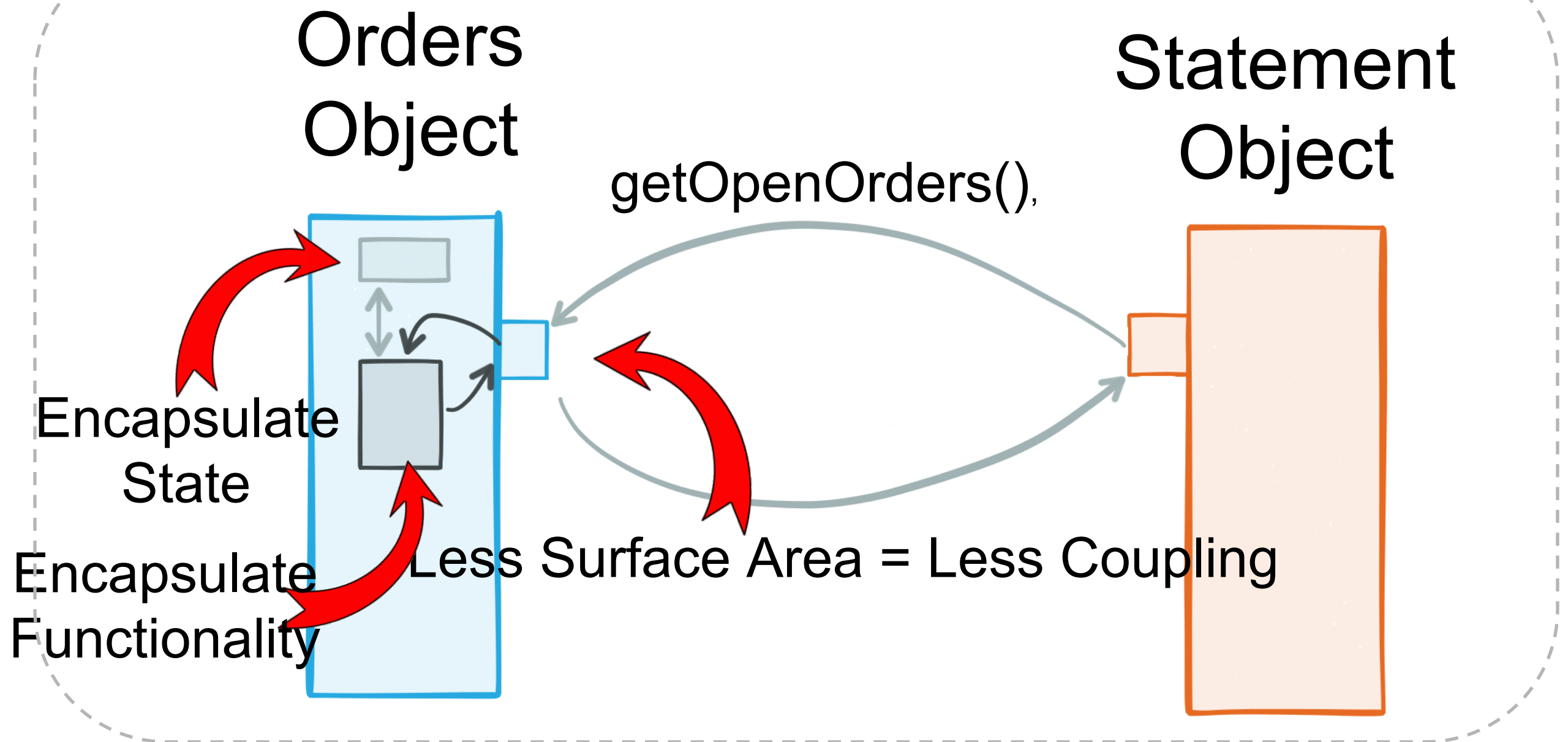


External World is something we should Design For



Independence
comes at a
cost
\$\$\$\$

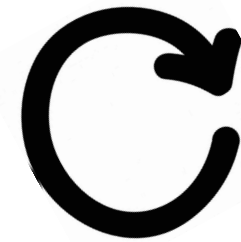
Consider Two Objects in one address space



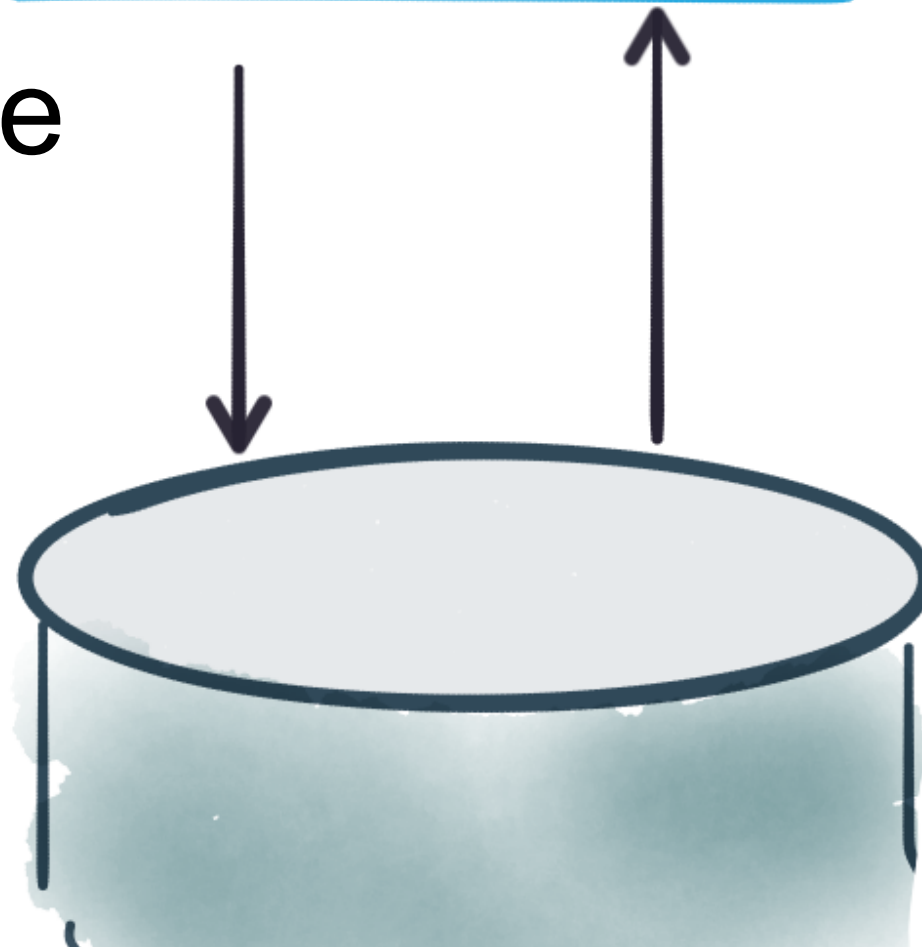
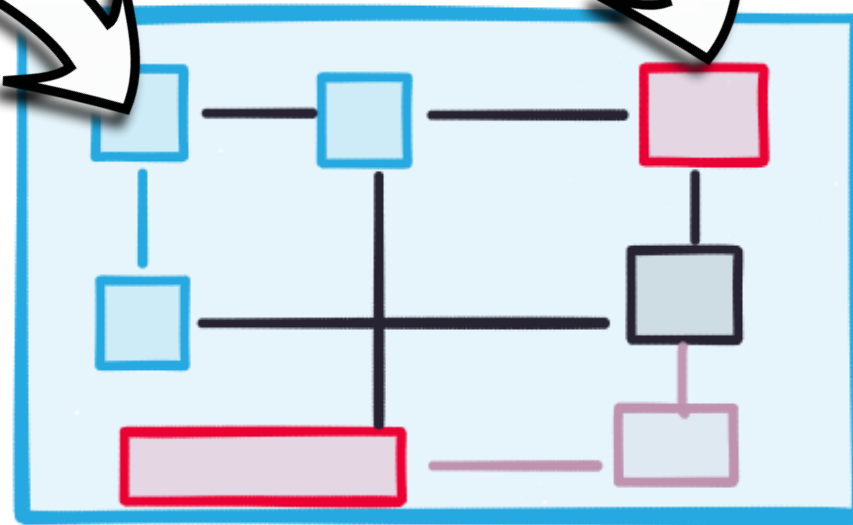
Encapsulation => Loose Coupling

Change 1 Change 2

Redeploy



Singly-deployable
apps are easy

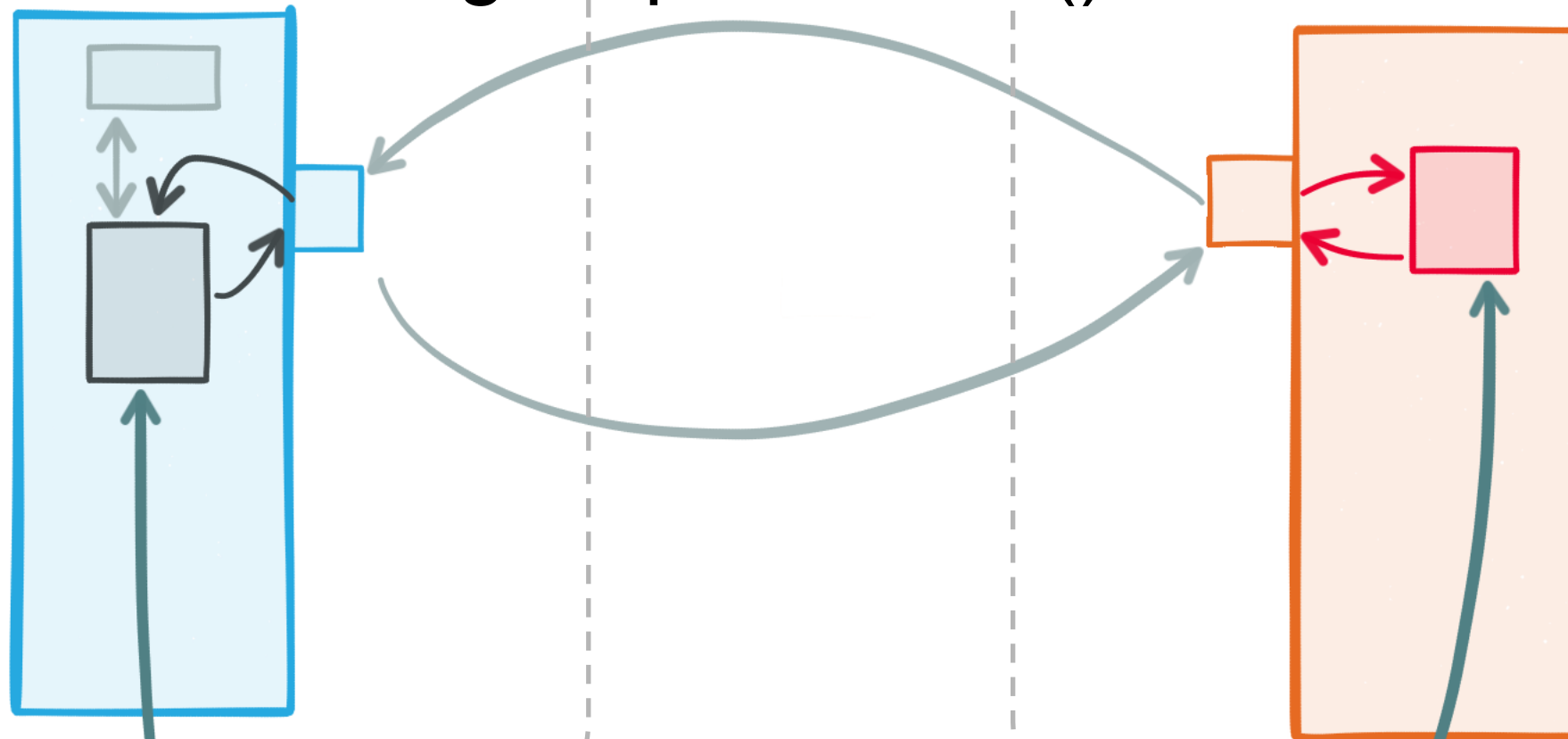


Independently Deployable

Independently Deployable

Orders Service Statement Service

getOpenOrders()



Synchronized
changes are painful

Services work best
where
requirements are
isolated in a single
bounded context

Single Sign On

Single Sign On `authorise()`, Business Service

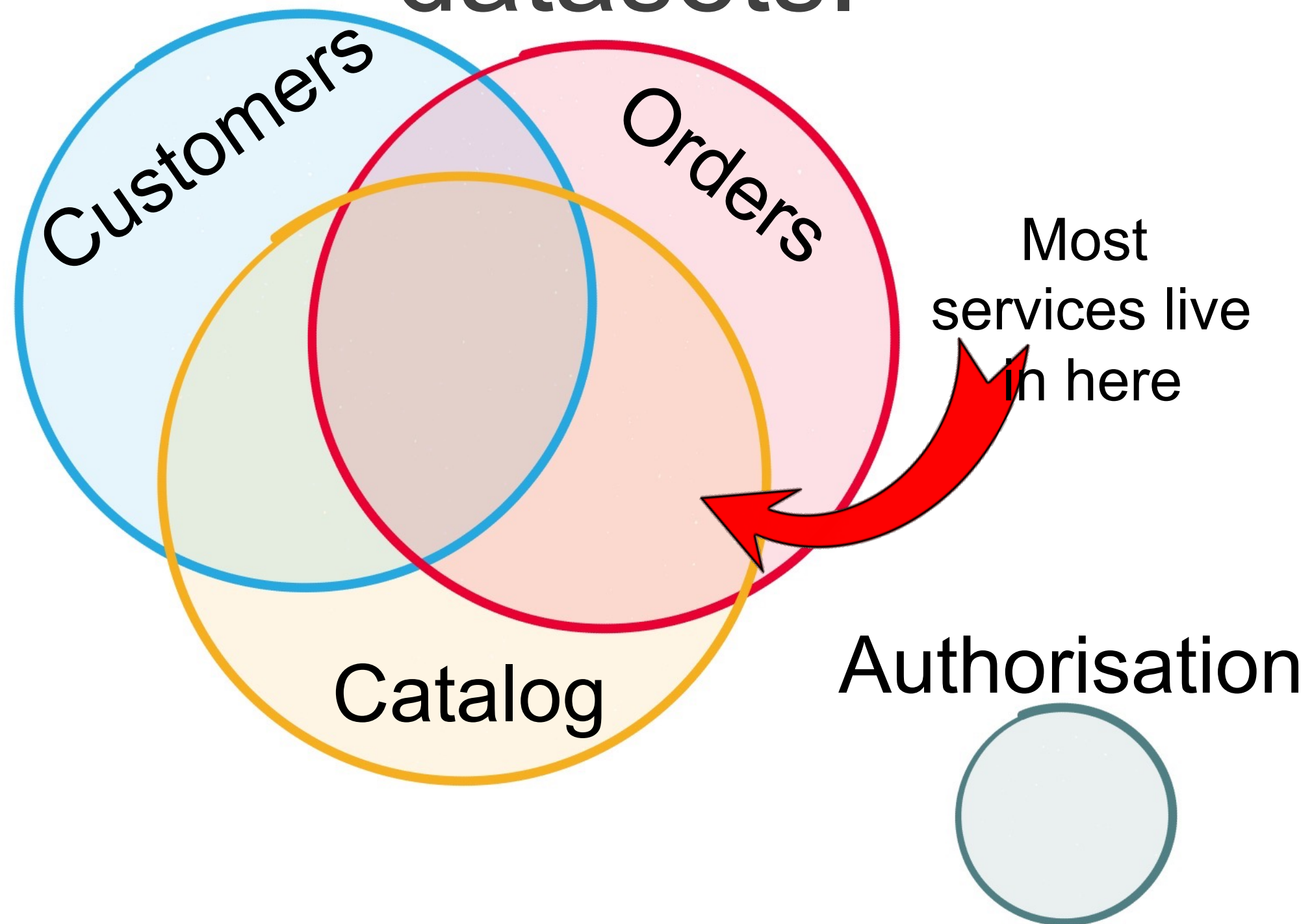


It's unlikely that a Business Service would need the internal SSO state / function to change

SSO has a
tightly
bounded
context

But business
services are
different

Most business services
share the same core
datasets.



The futures of
business
services are
far more
tightly



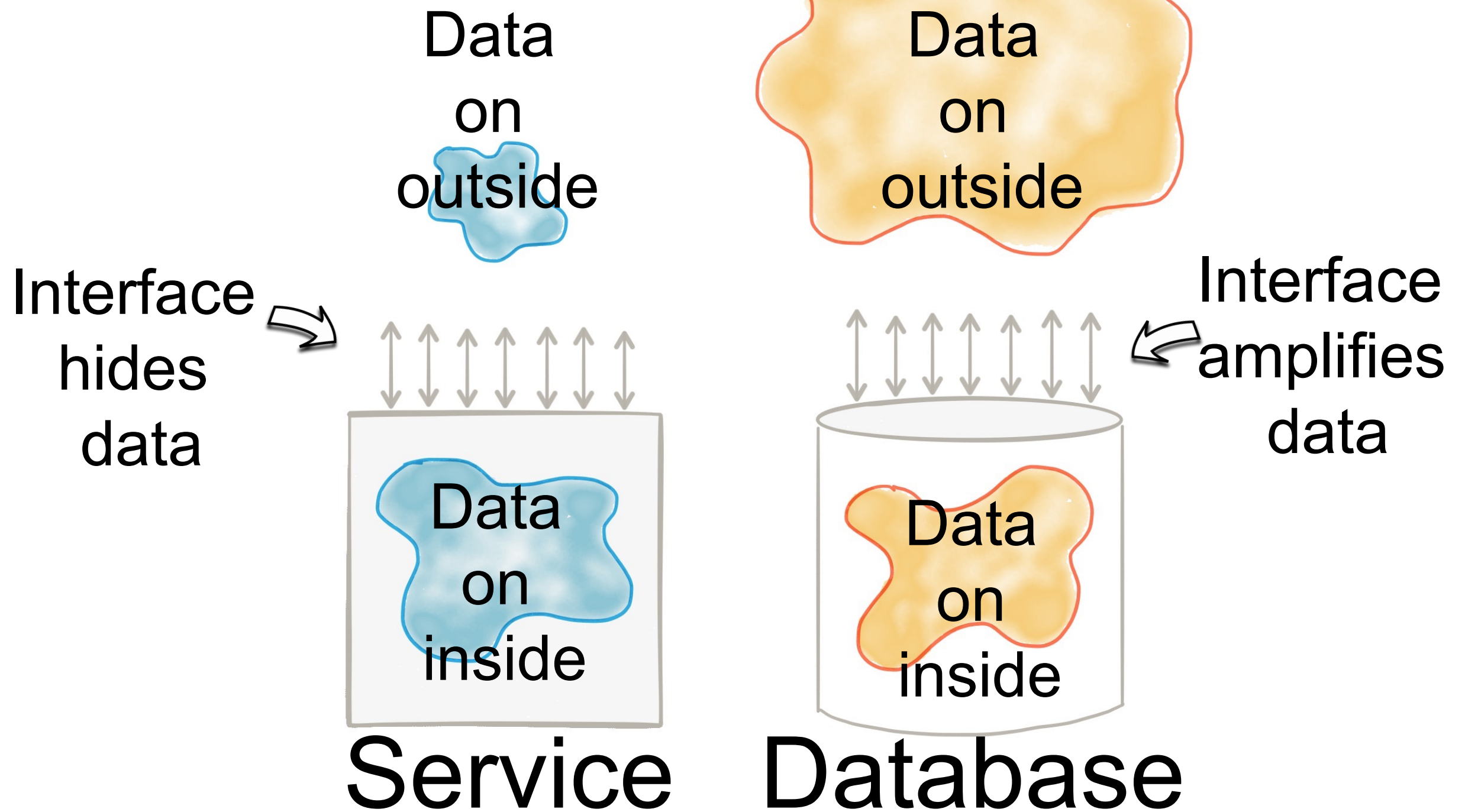
intertwined

We need encapsulation to
hide internal state. Be
loosely coupled.

But we need the freedom to
slice & dice shared data
like any other dataset

But data
systems have
little to do with
encapsulation

Databases amplify the data they hold



The data dichotomy

Data systems are about exposing data

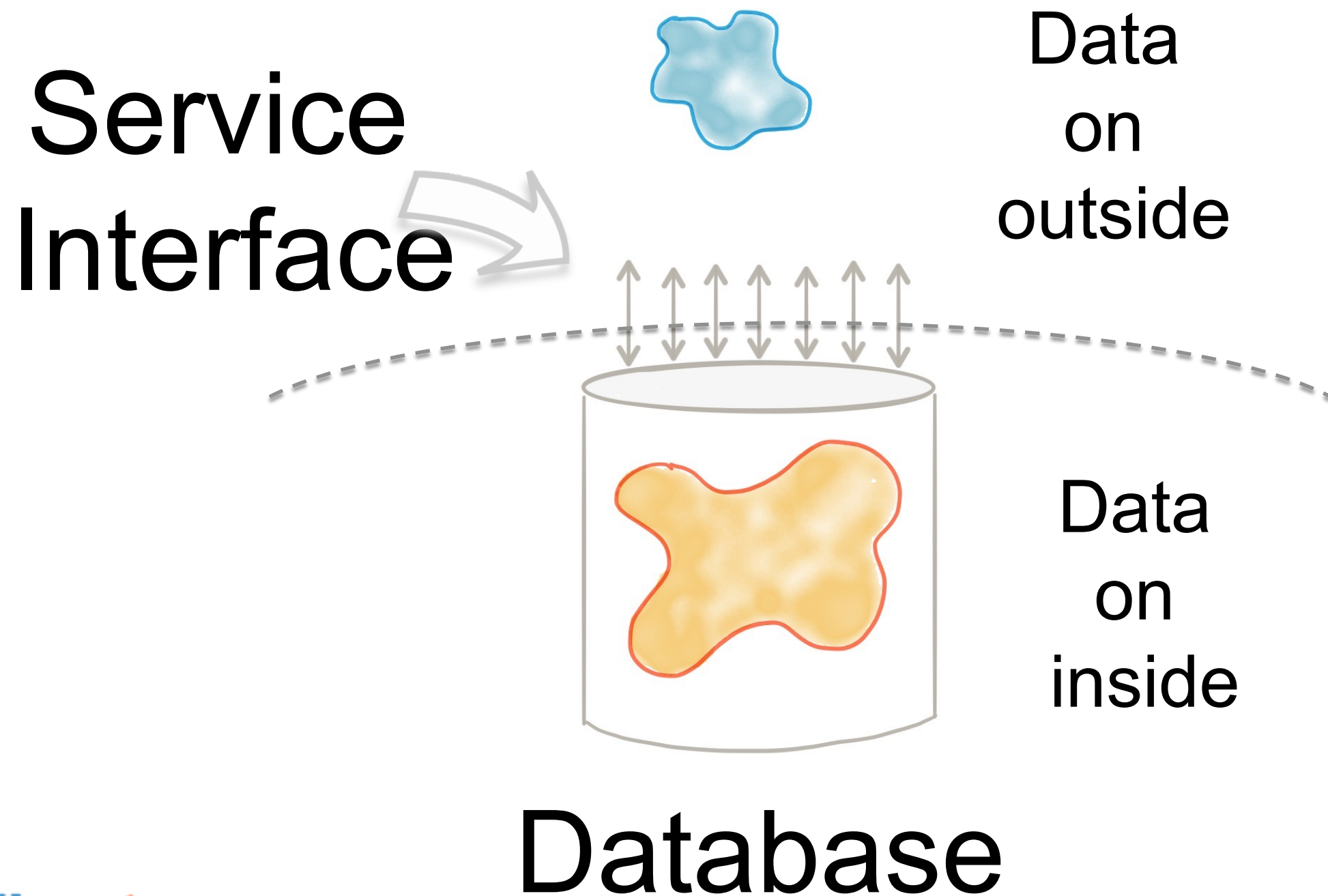
Services are about hiding it.

Microservices shouldn't share a database

Good Advice!

So what do we
do instead?

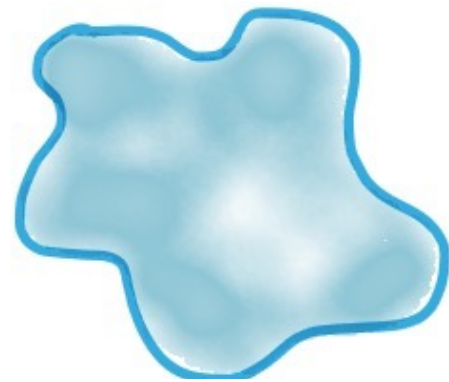
We wrap a database in a service interface



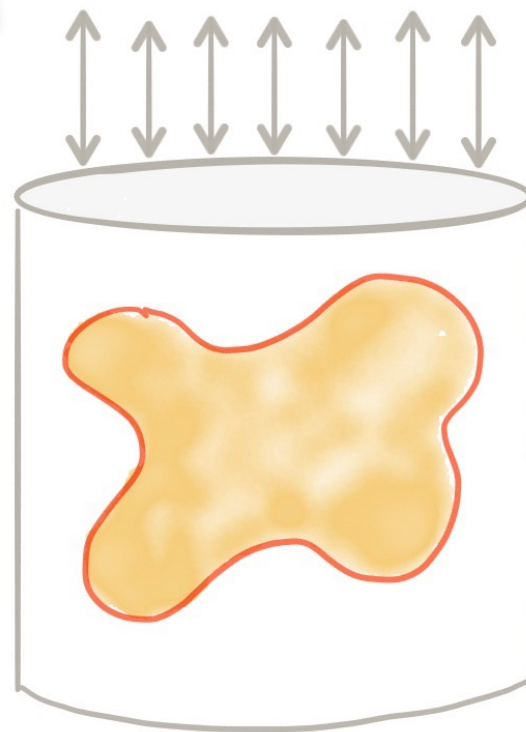
One of two
things
happens next

Either (1) we constantly add to the
interface,
as datasets grow

Service
Interface

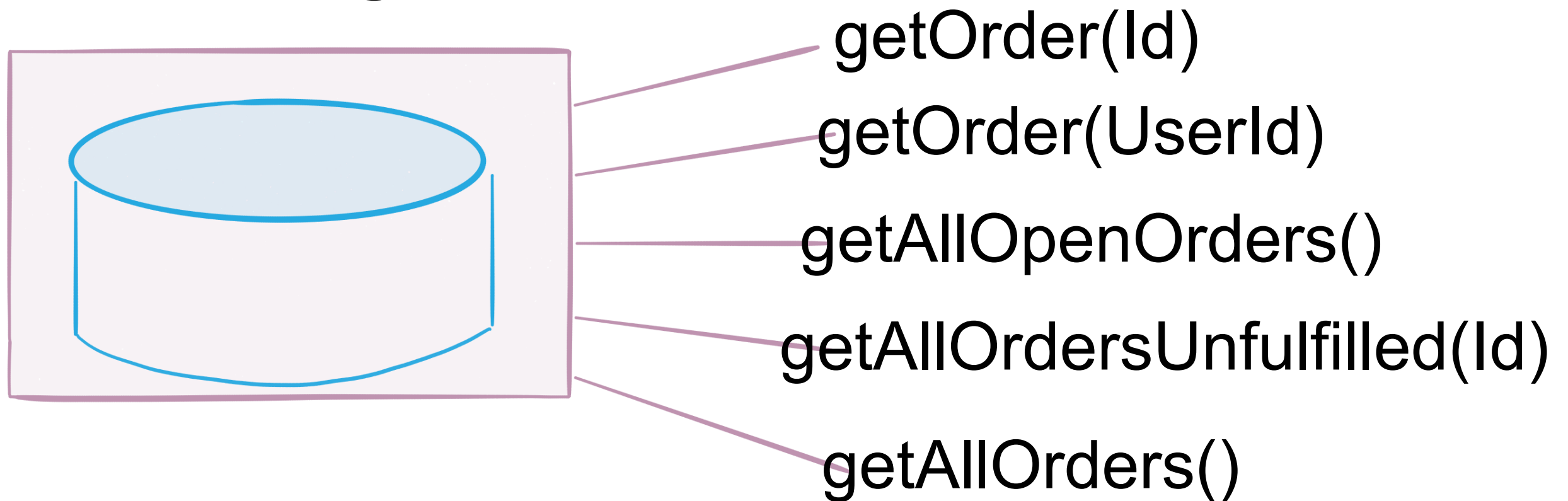


```
getOpenOrders(  
  fulfilled=false,  
  deliveryLocation=CA,  
  orderValue=100,  
  operator=GreatherThan)
```



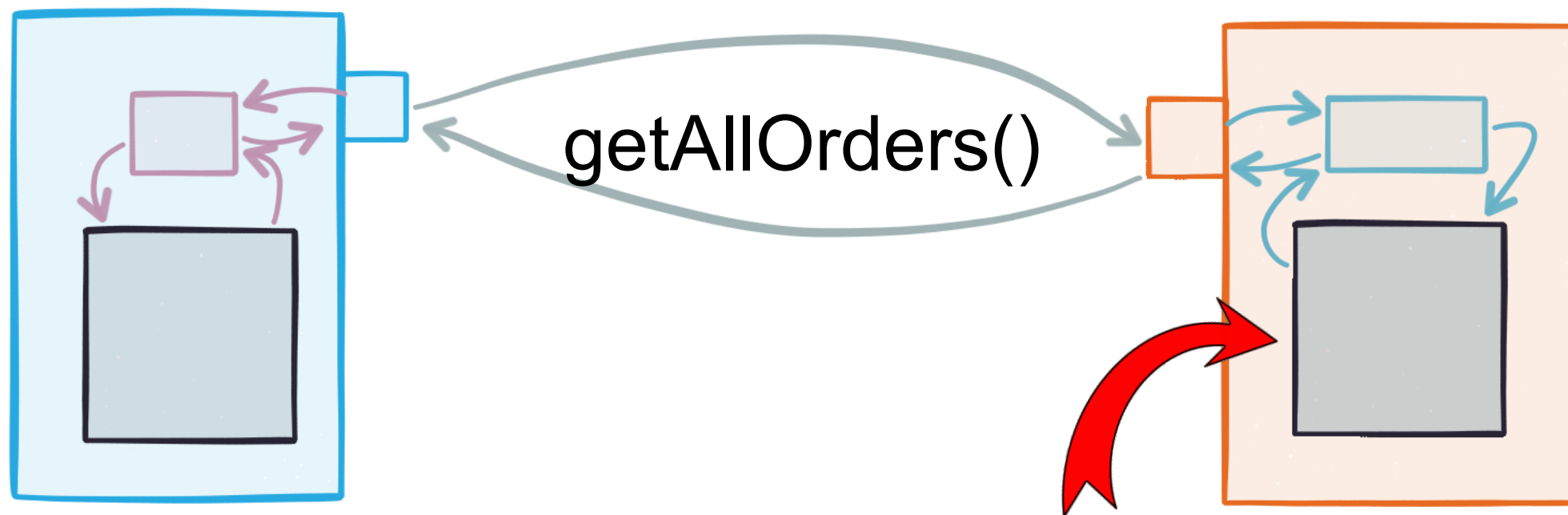
Database

(1) Services can end up looking like kookie home-grown databases



...and **DATA**
amplifies this
“Data-Service”
problem

(2) Give up and move whole datasets en masse

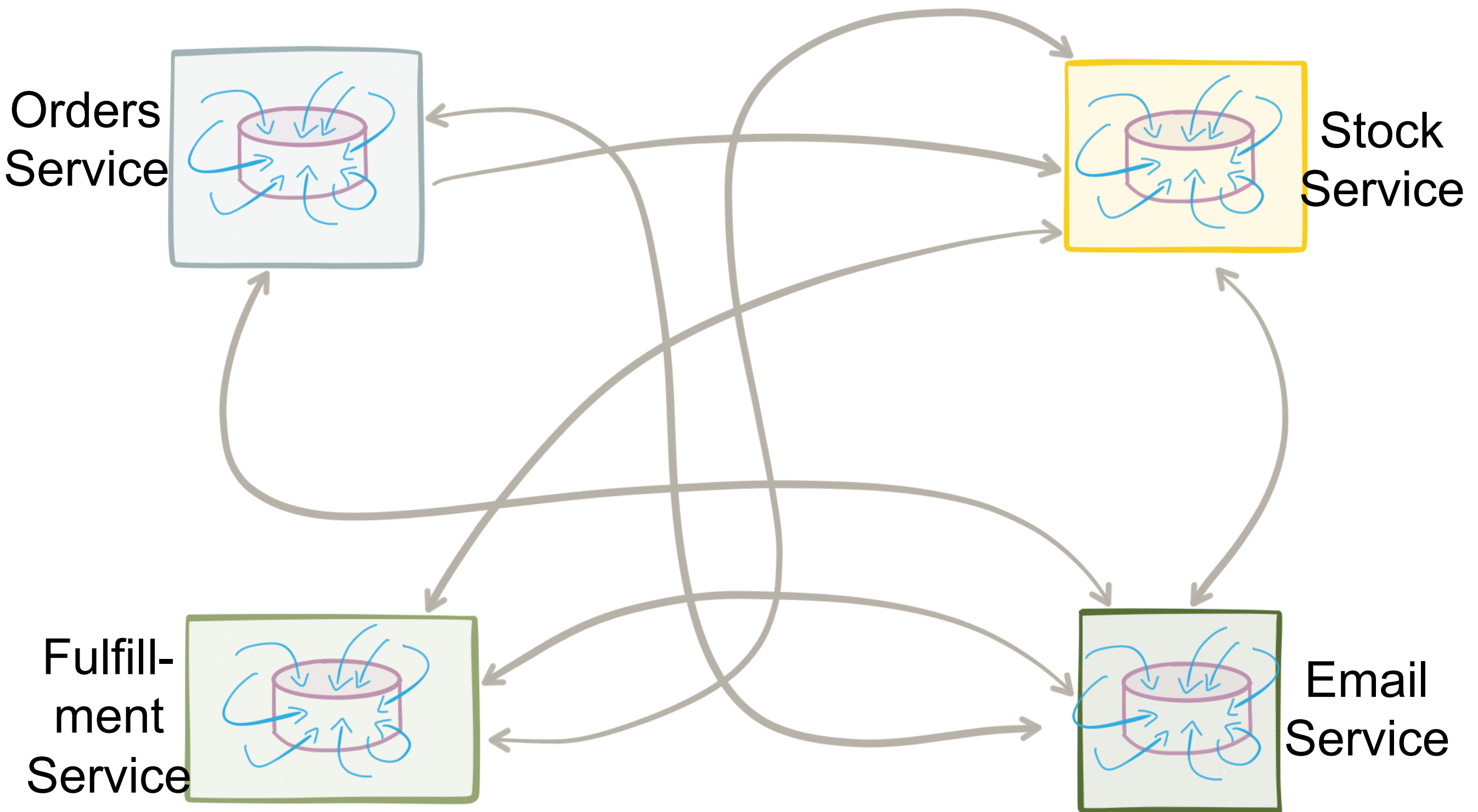


Data Copied
Internally

- a) Too slow to change
- b) To perform joins

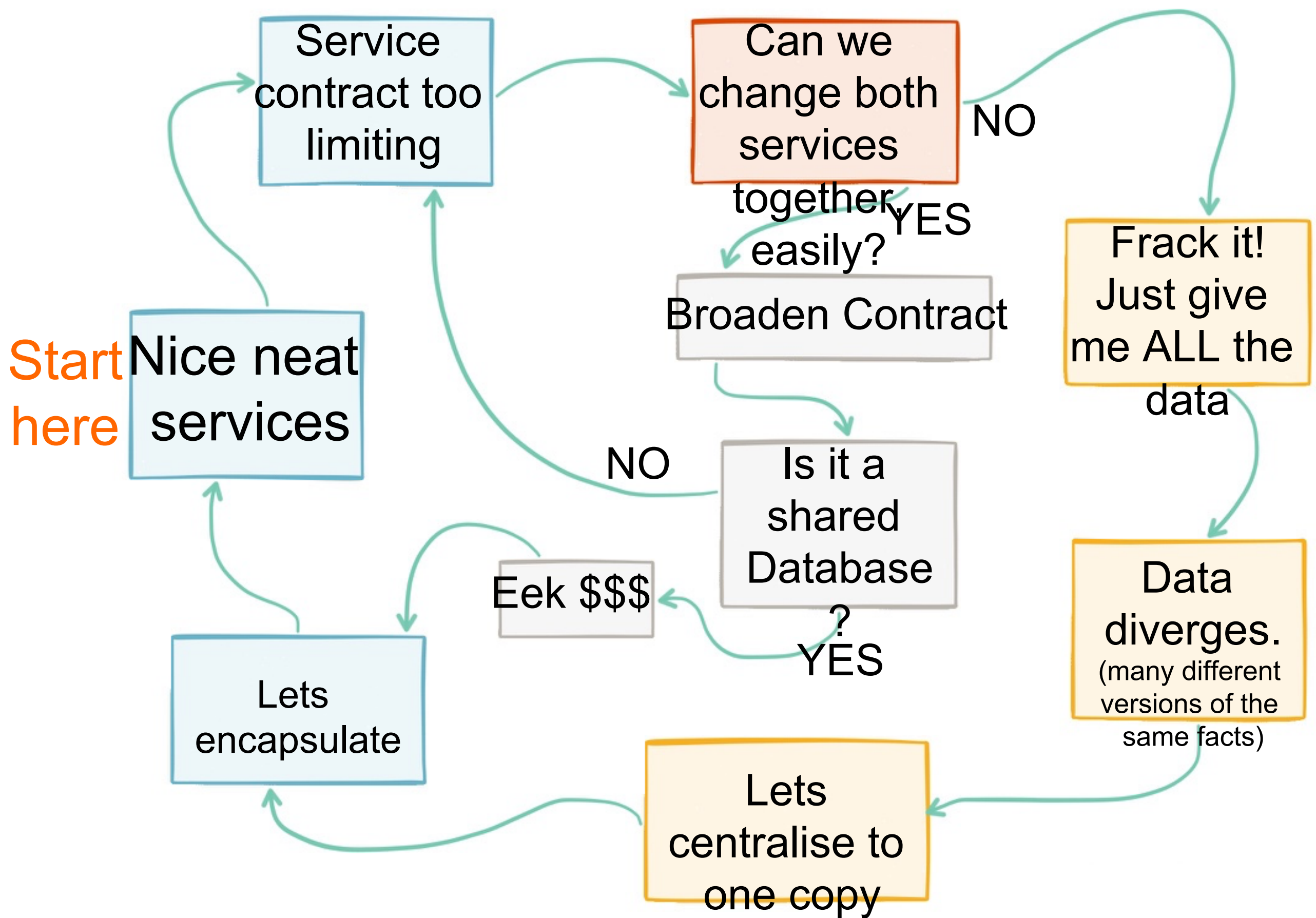
This leads to a
different
problem

Data diverges over time



The more
mutable
copies, the
more data will
diverge over
time

cycle of inadequacy:



These forces compete in the systems we build



Accessibility

VS



Coupling

Divergence

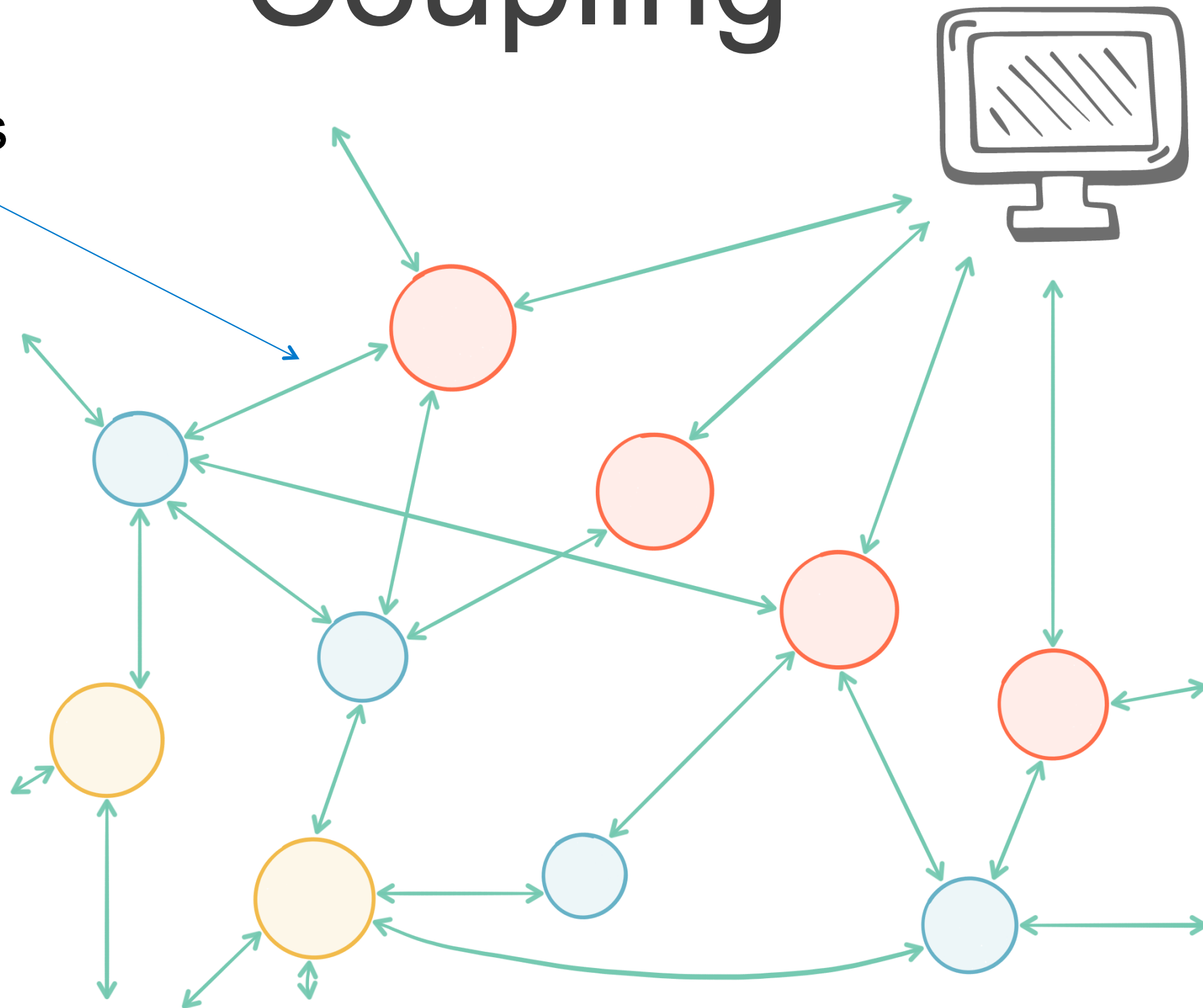


Is there a
better way?

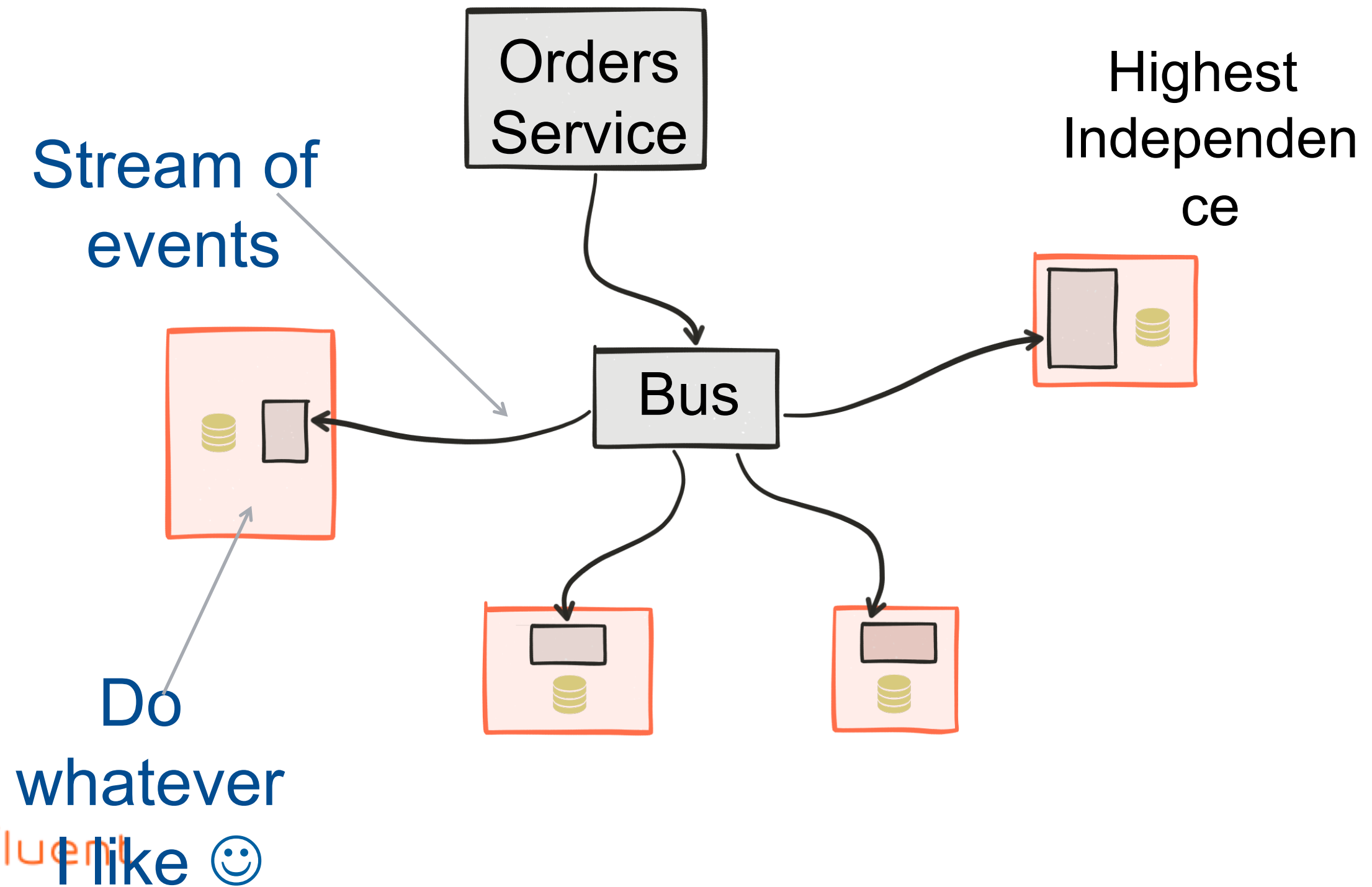
Step 1:
Build on a backbone of
events.
Layer in requests where
necessary.

Request Driven -> highest Coupling

Commands
& Queries

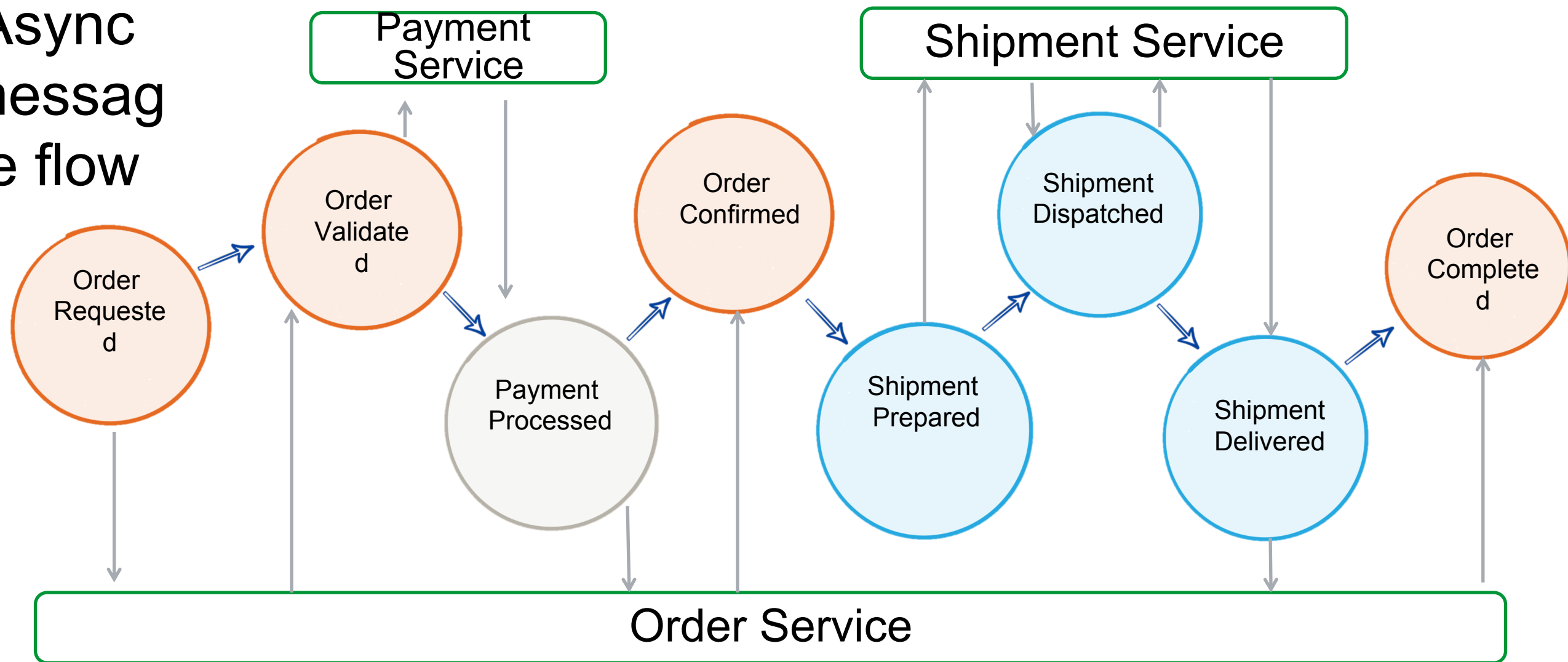


Event Broadcast => lowest coupling

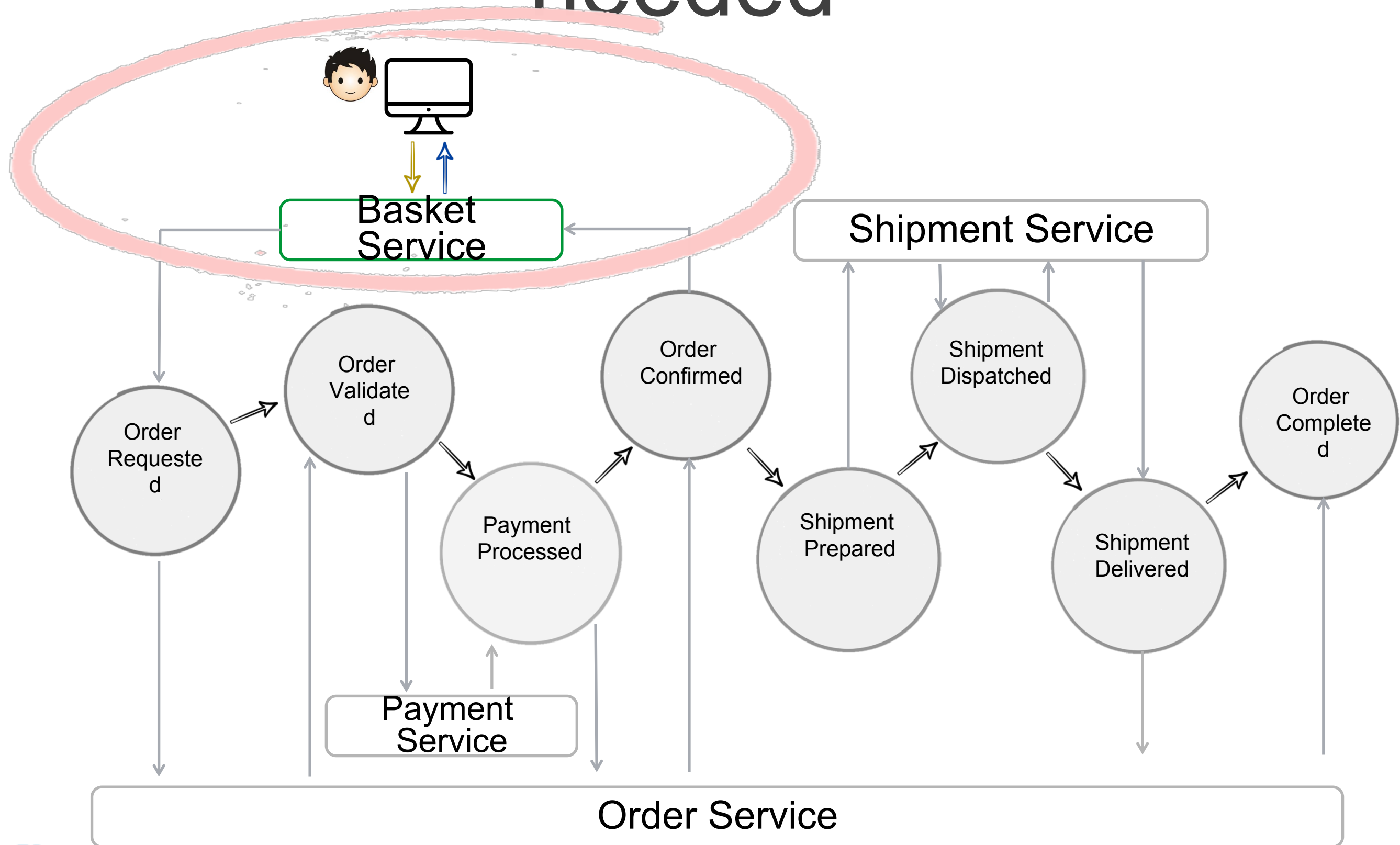


Event Driven Services Chain (Event Collaboration)

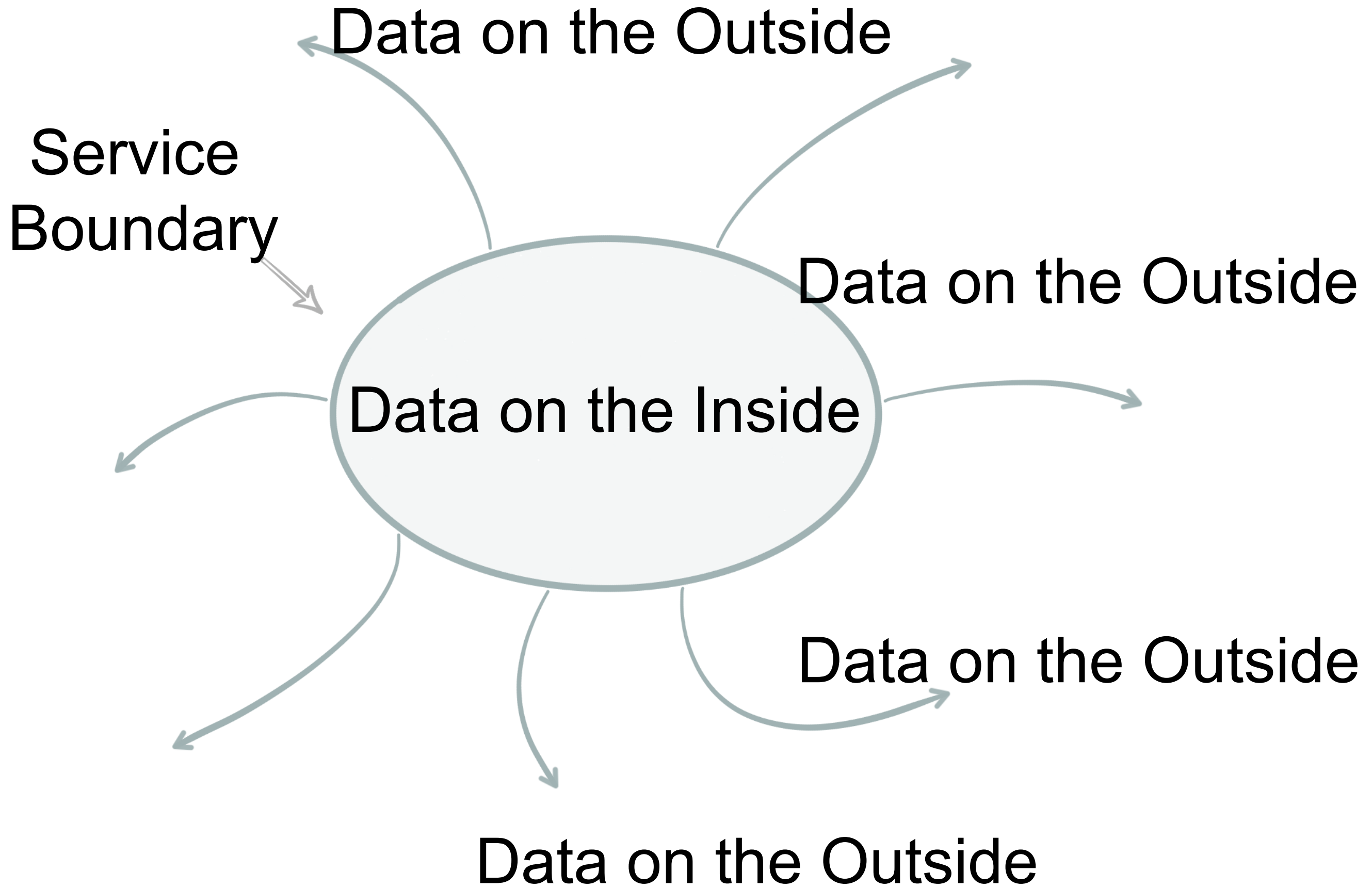
Async
messag
e flow



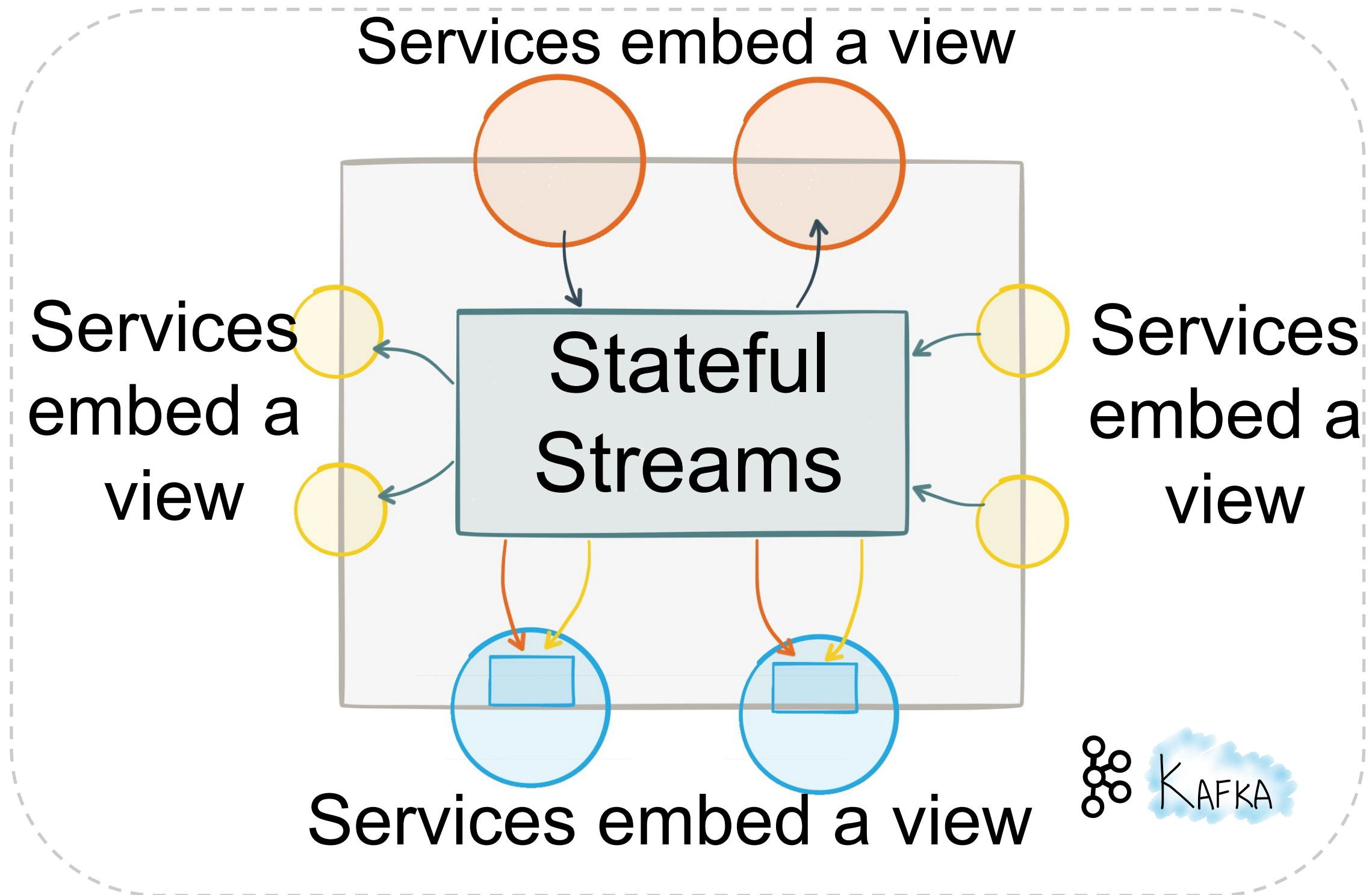
Layer queries/commands where needed

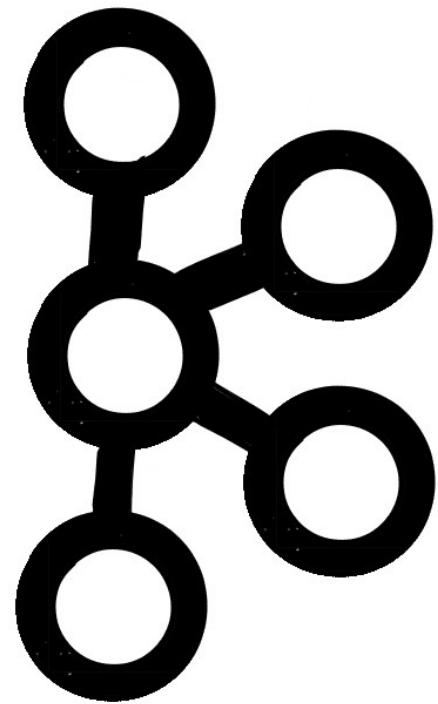


Step 2:
Make data-on-the-
outside
a 1st class citizen



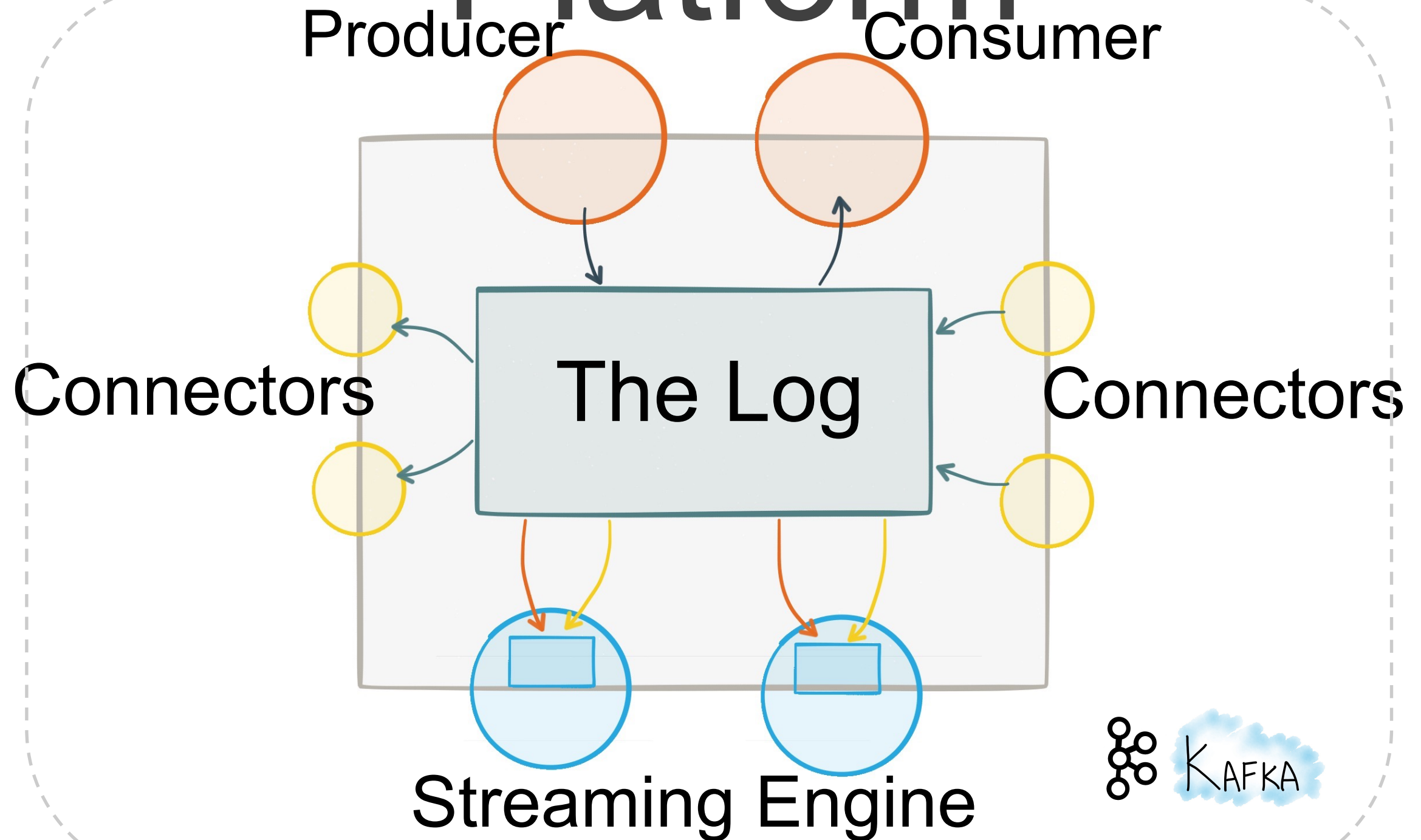
Stateful Streams



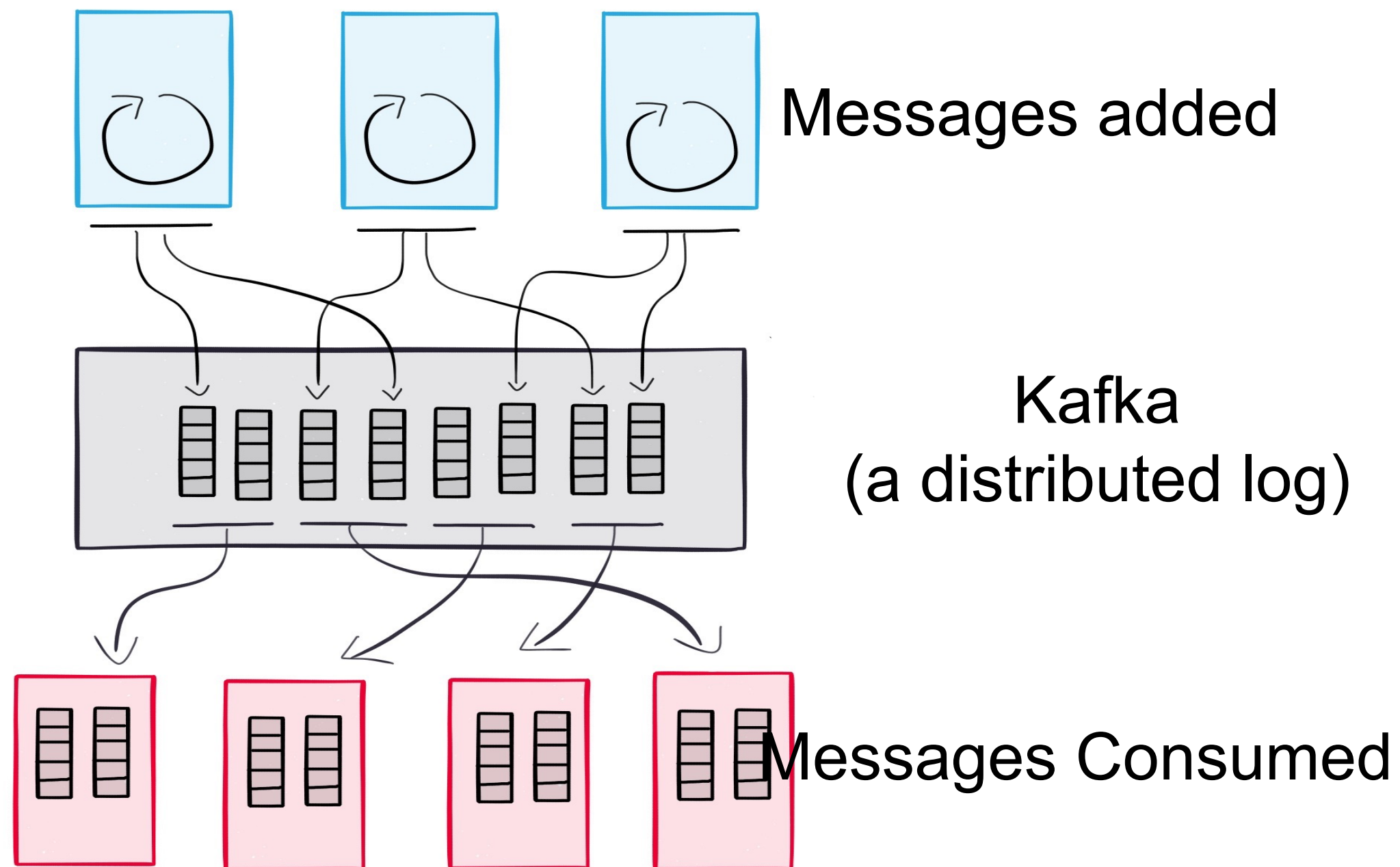


Kafka helps with this

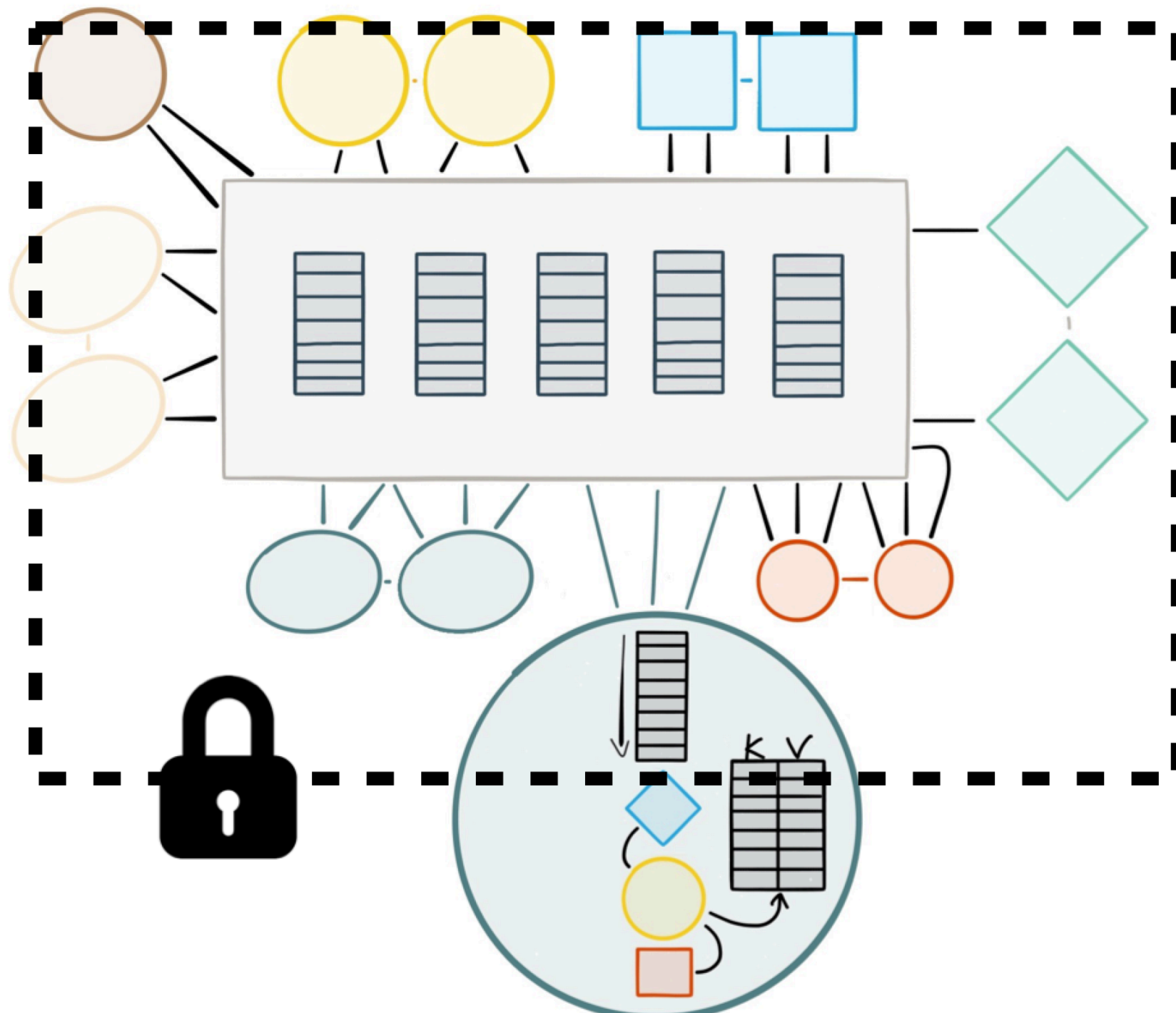
Katka: a Streaming Platform



The Log



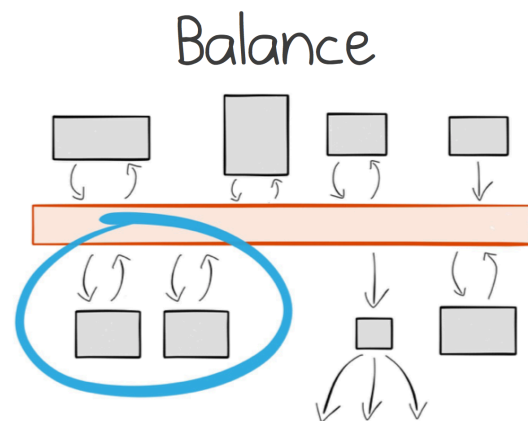
Transactions



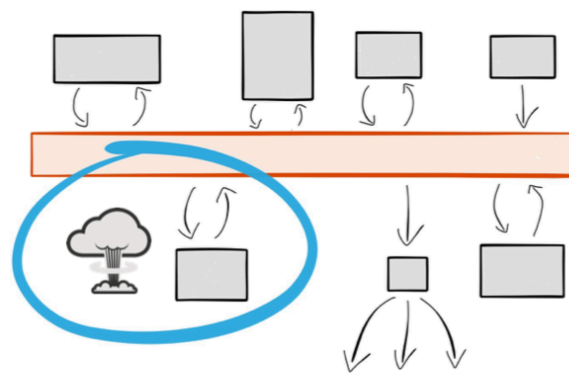
Service Backbone

Scalable, Fault Tolerant, Concurrent, Strongly Ordered,
Transactional*, Retentive

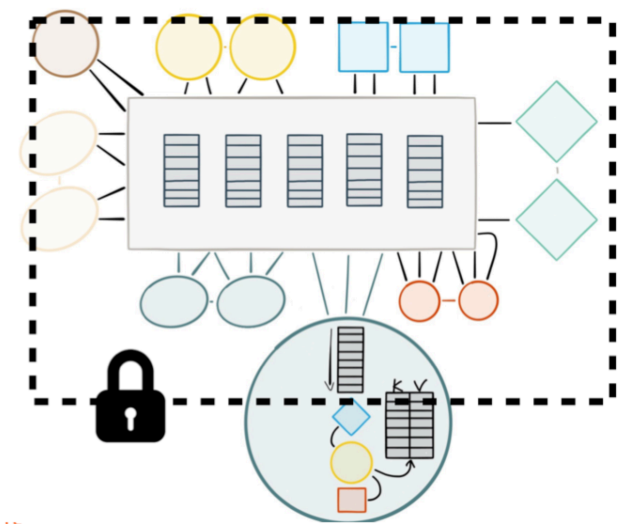
Services Naturally Load



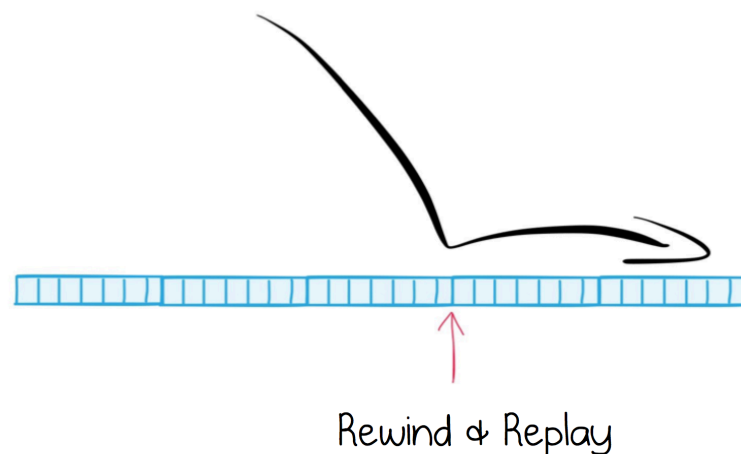
This provides Fault Tolerance



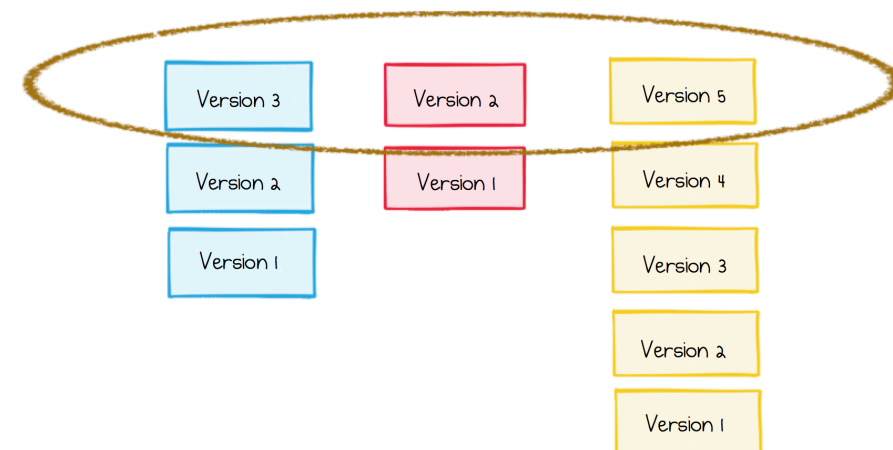
Transactions



Roll back in "time"



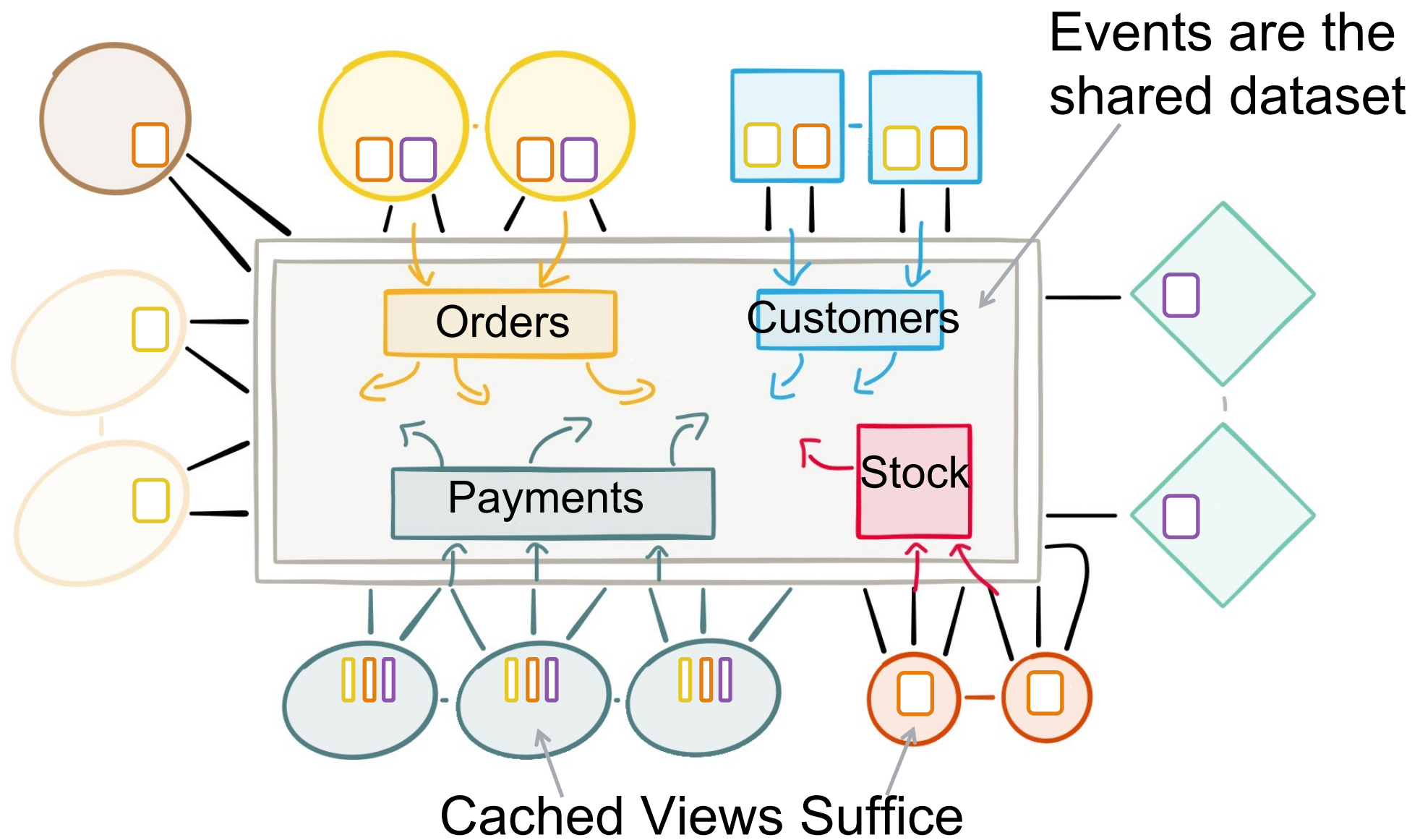
Compacted Log
(tables & streams)



A place to keep
the data-on-the-
outside

Leave data in Kafka

-> Services only need caches



Now add
Stream
Processing

What is Stream Processing?

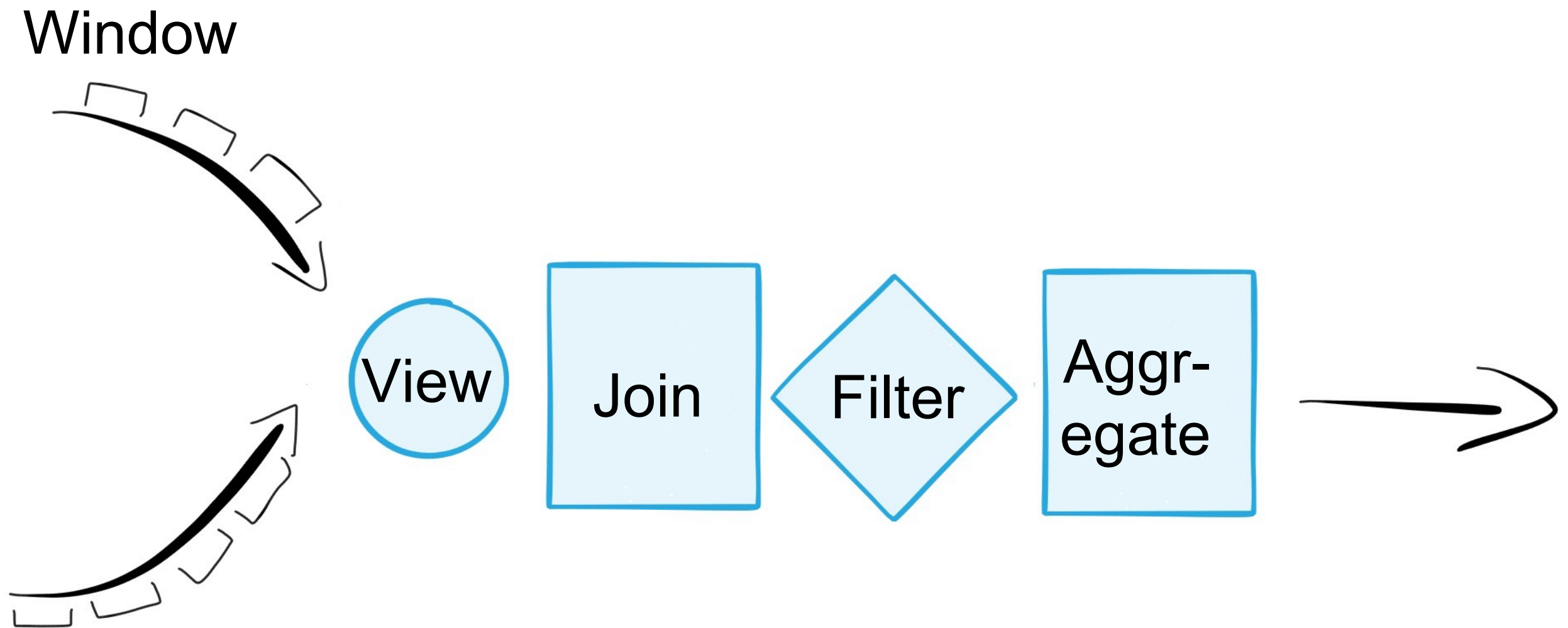


A machine for combining
and processing streams of
events

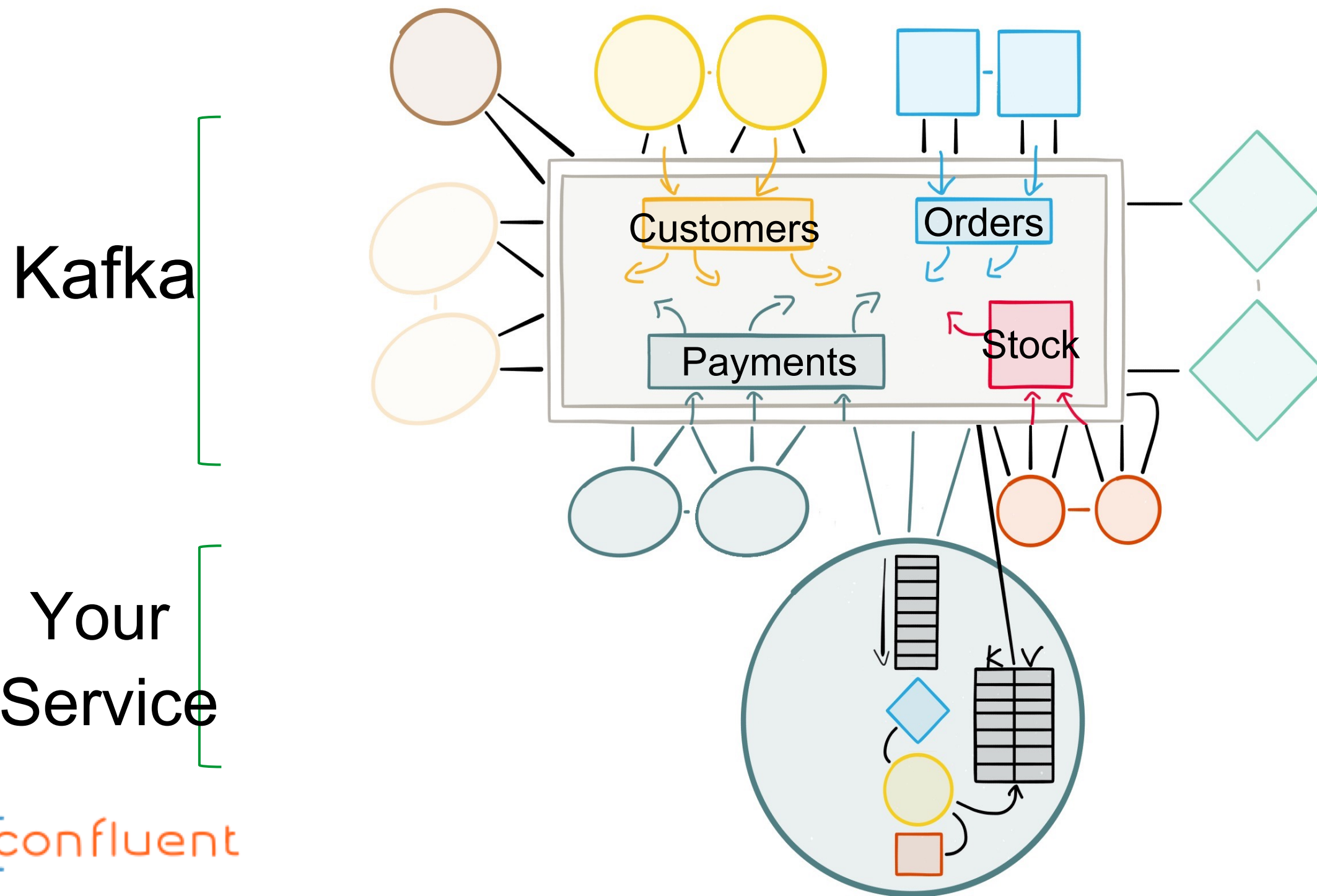
Stateful stream processing

- Data constantly updated
- Some data accumulated in each service
 - Accumulated via a window
 - Accumulated via a key

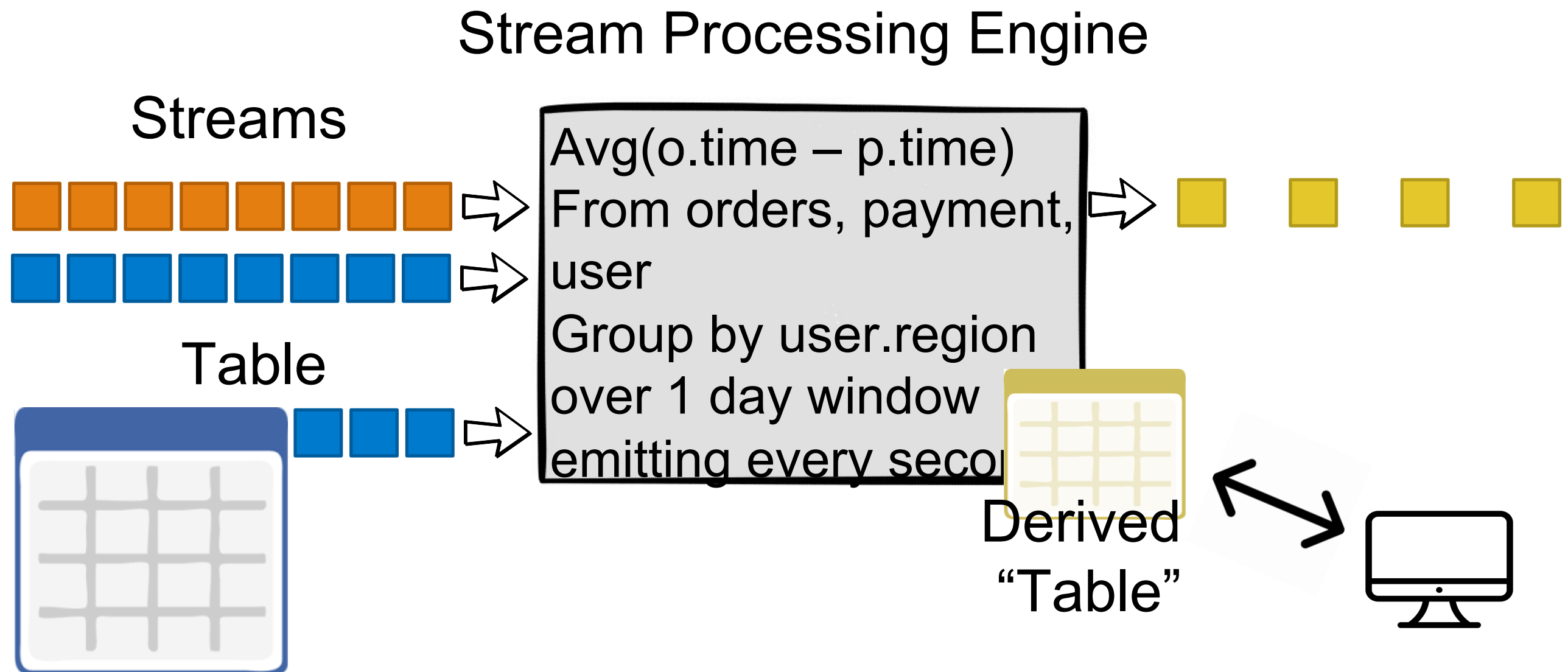
A Query Engine inside your service



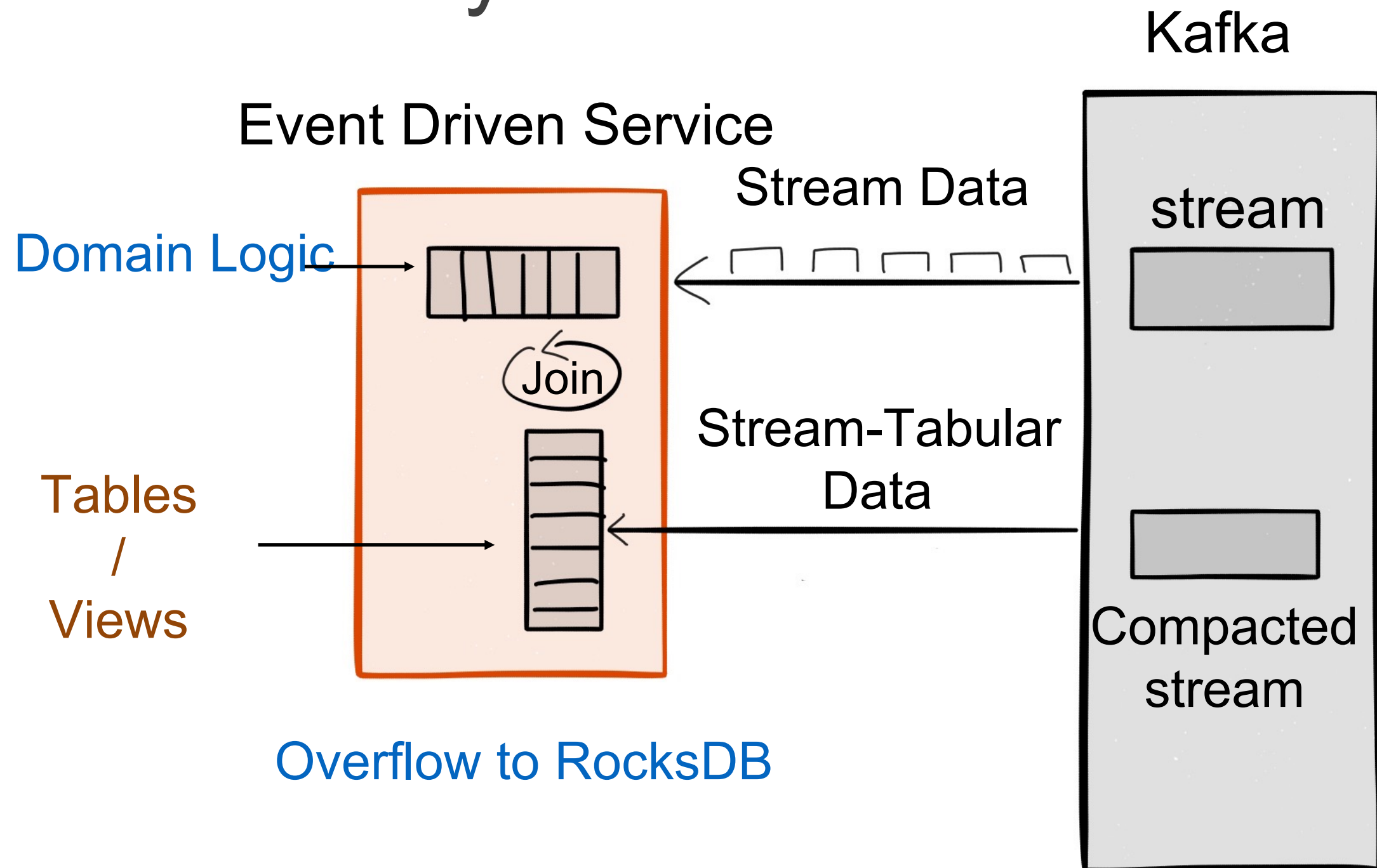
A database embedded in your service. One designed for handling streams.



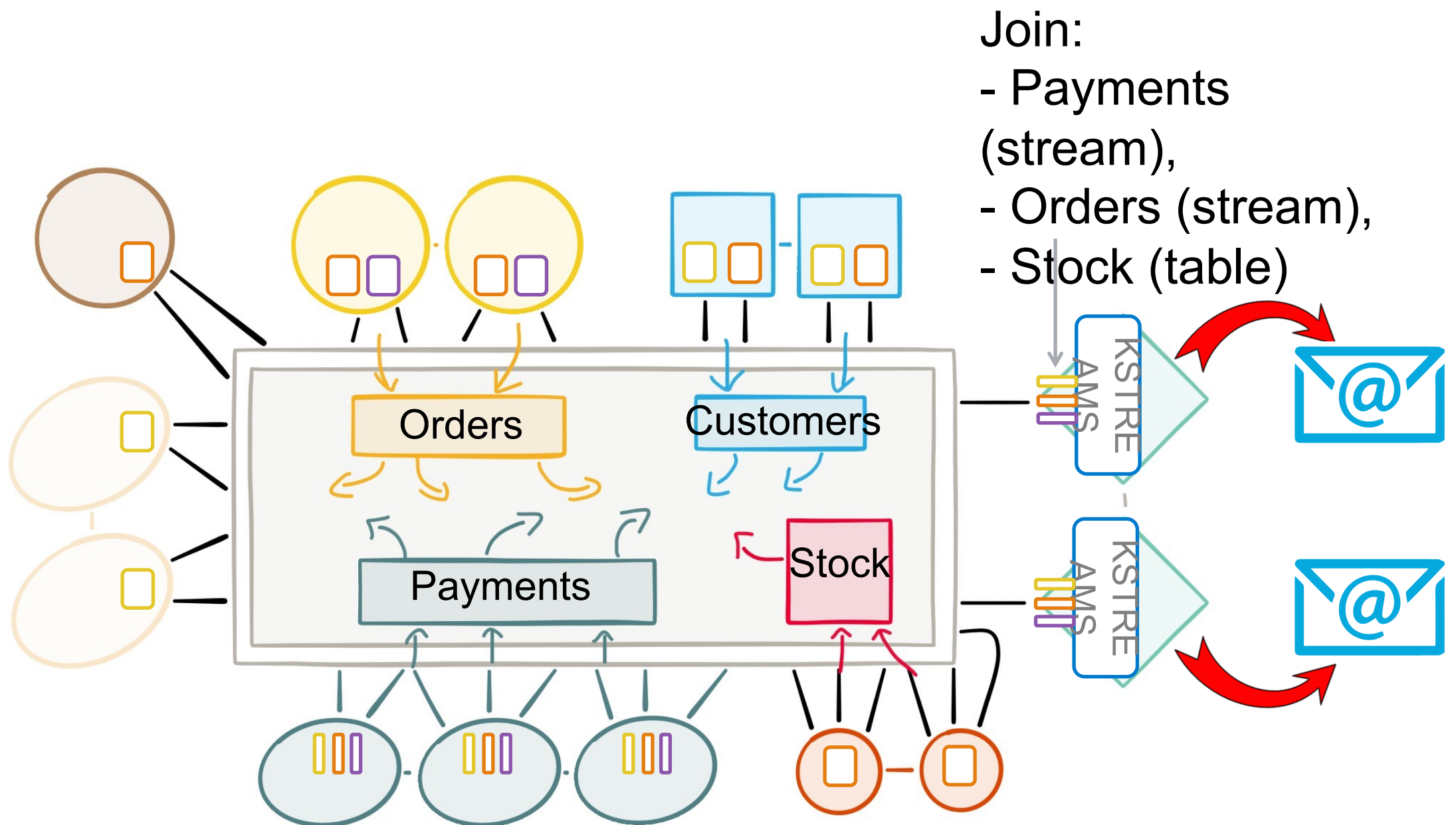
View as a Stream or a Table (in & out)



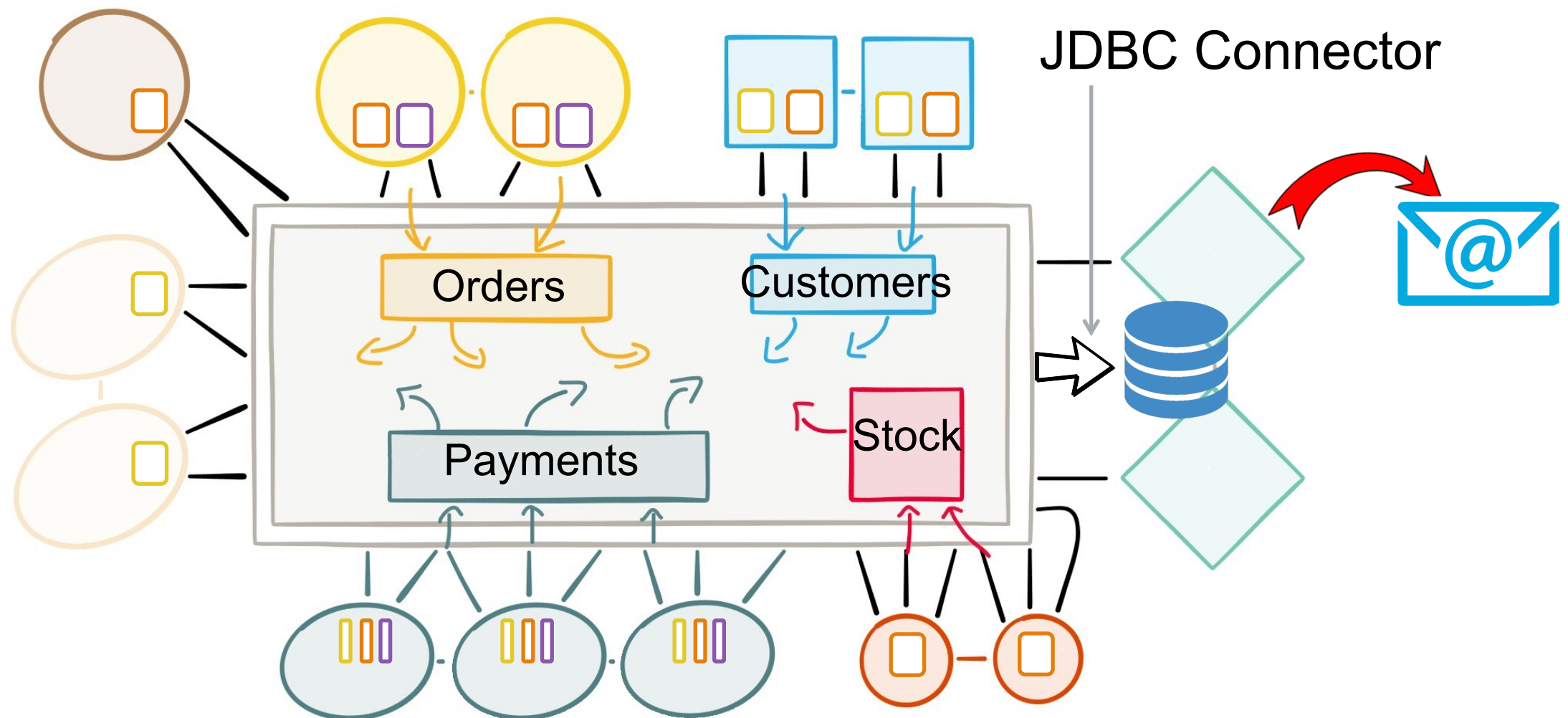
Window / Tables Cached on disk in your Service



Example: Trigger an email when a payment is confirmed. Include Order & Stock description

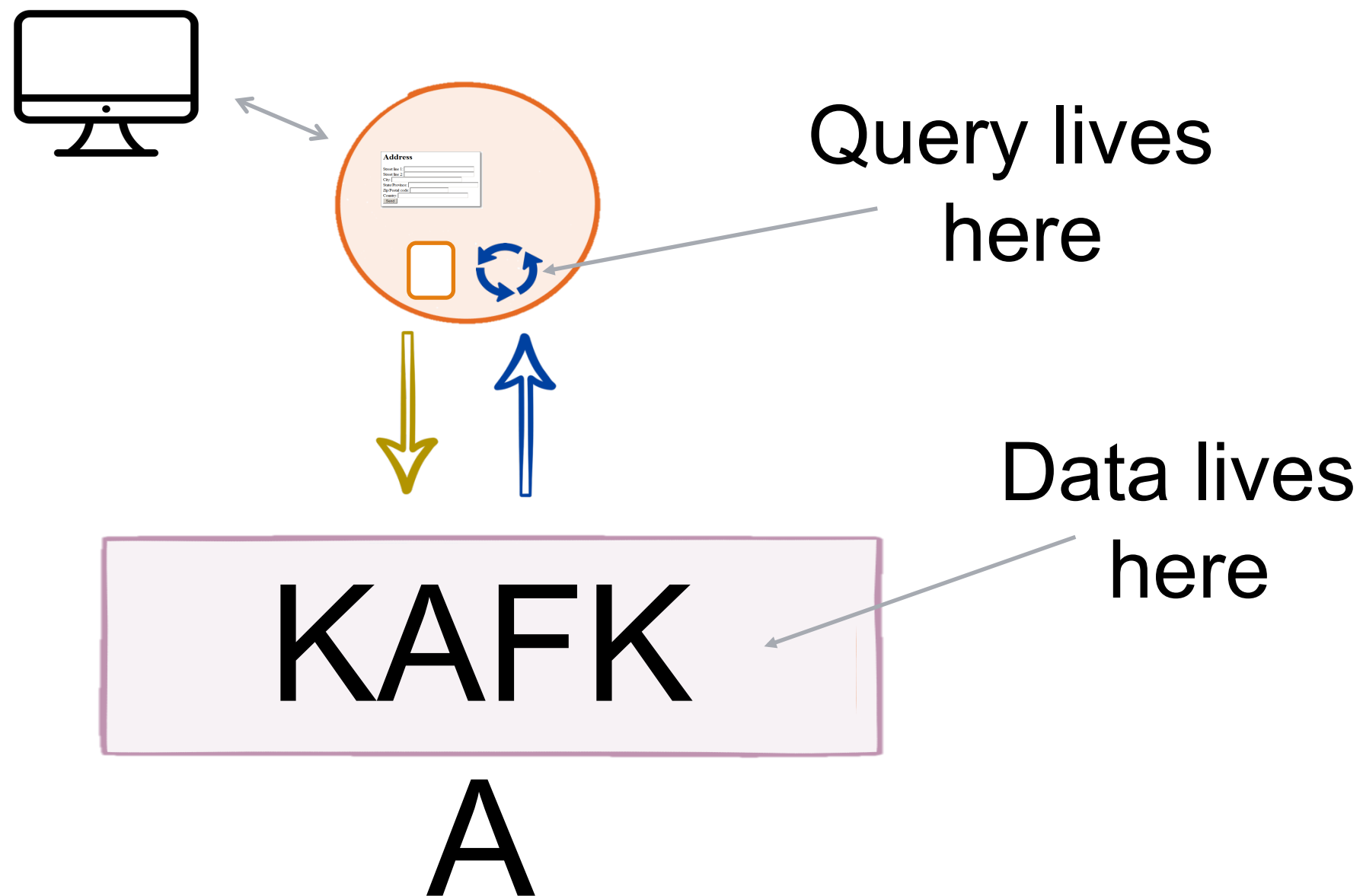


If a streaming engine doesn't cut it,
branch out to database but keep it
ephemeral

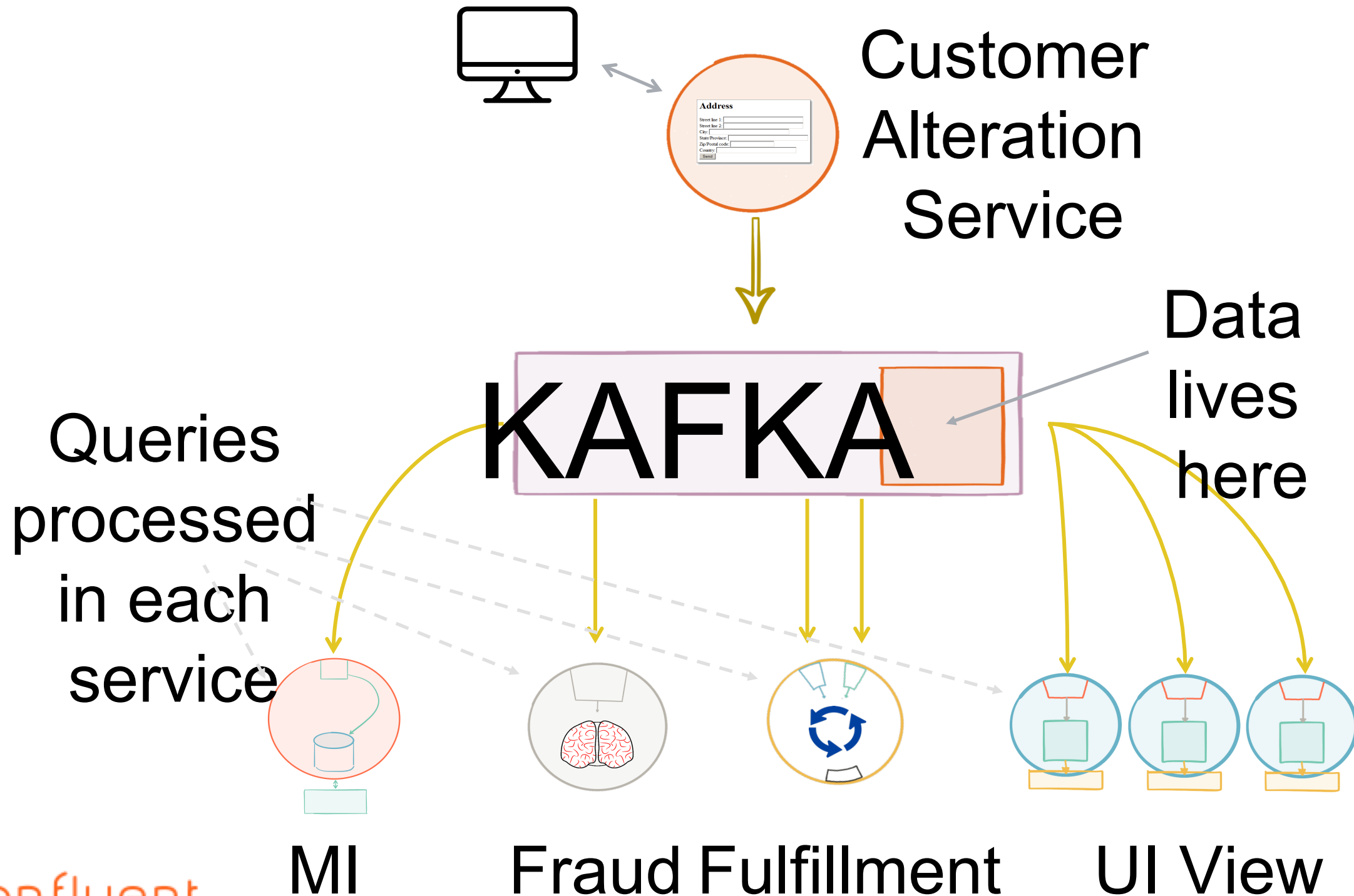


So we have
shared
storage in the
Log, and a
query engine
layered on top

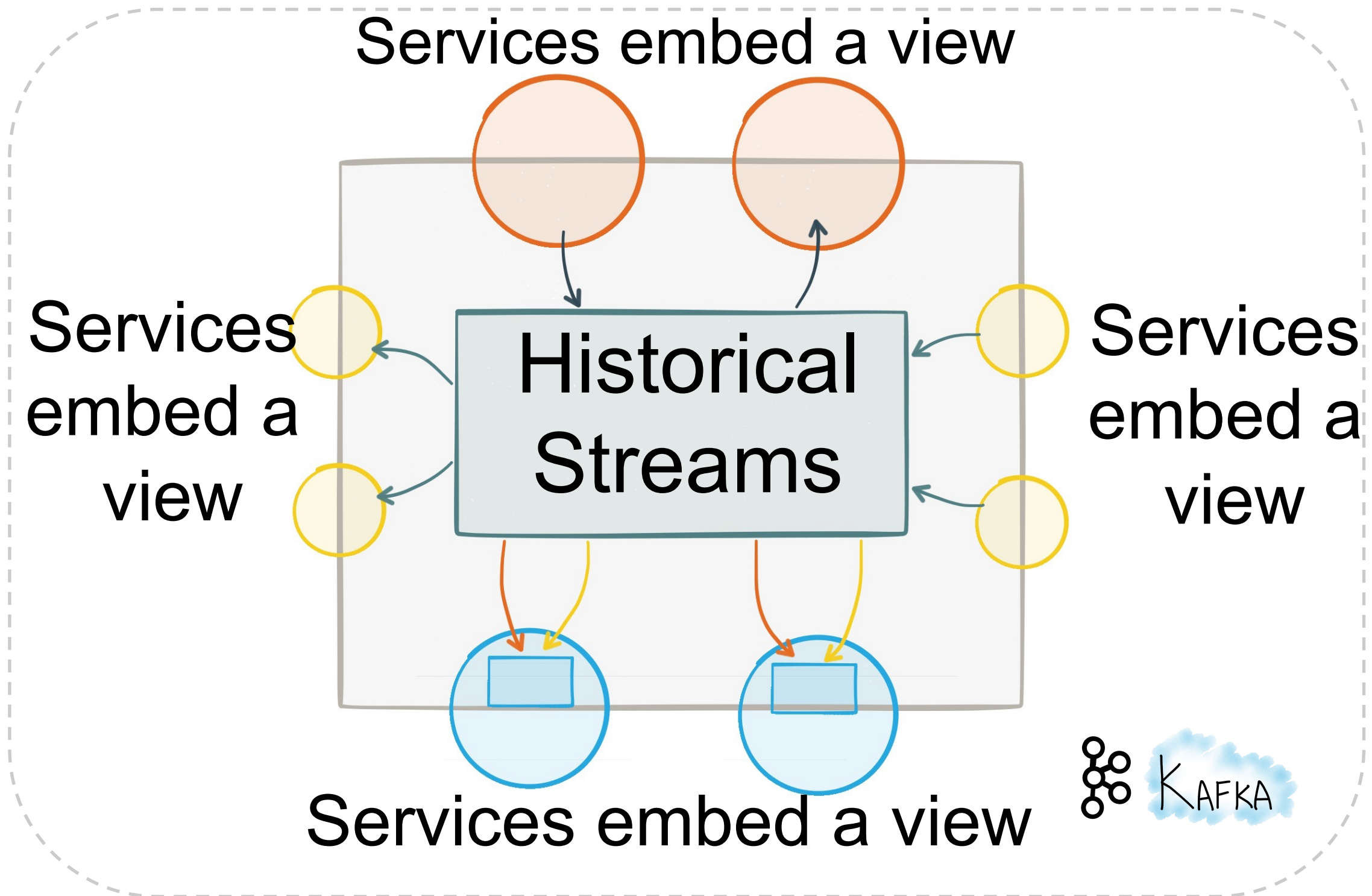
Data Storage + Query Engine == Database?



1 x Data Storage + n x Query == Shared Database?



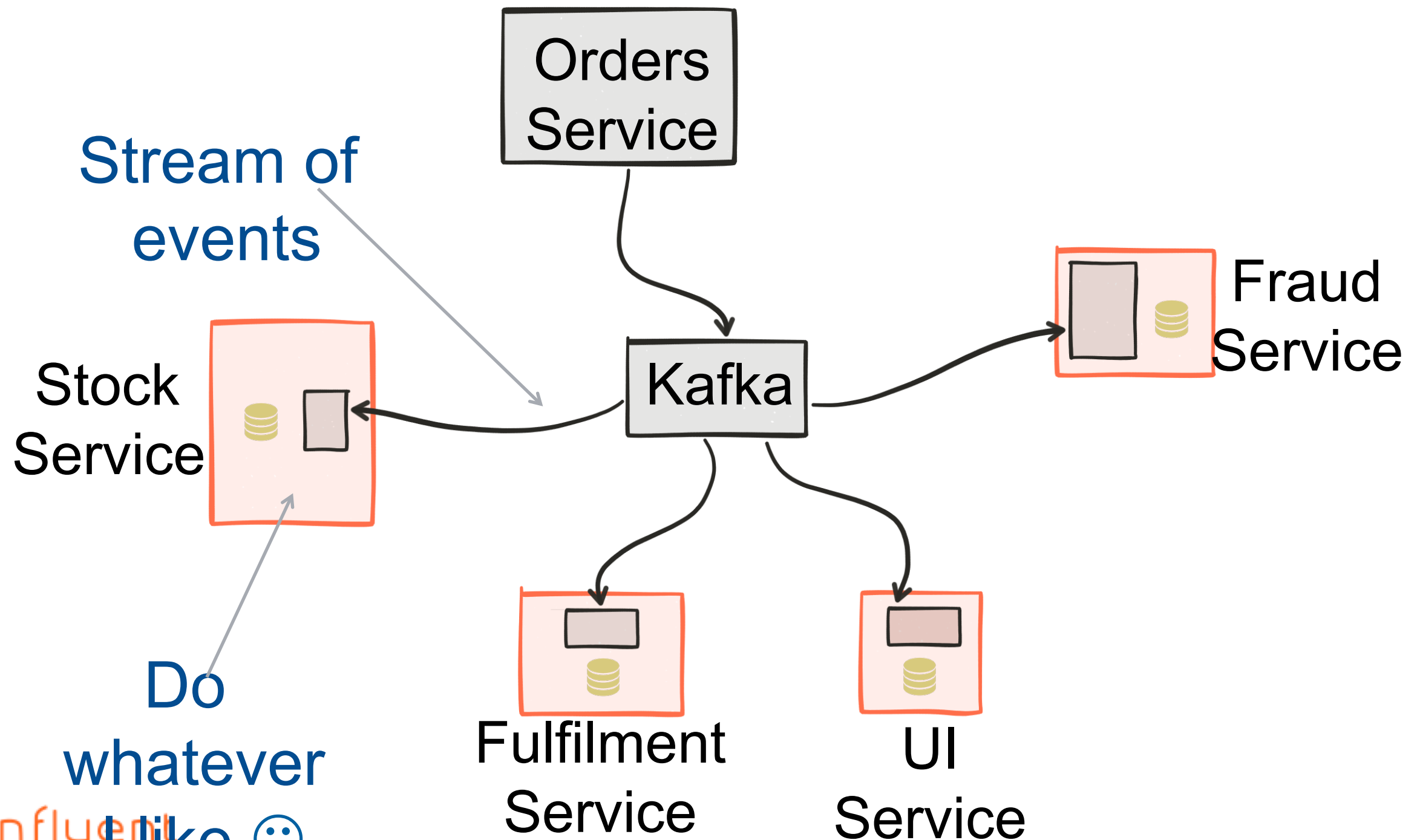
A Database, Inside Out



Microservices
shouldn't
share a
database

But this isn't a
normal
database

Event Broadcast has the lowest coupling



Stream of
events

Stock
Service

Orders
Service

Kafka

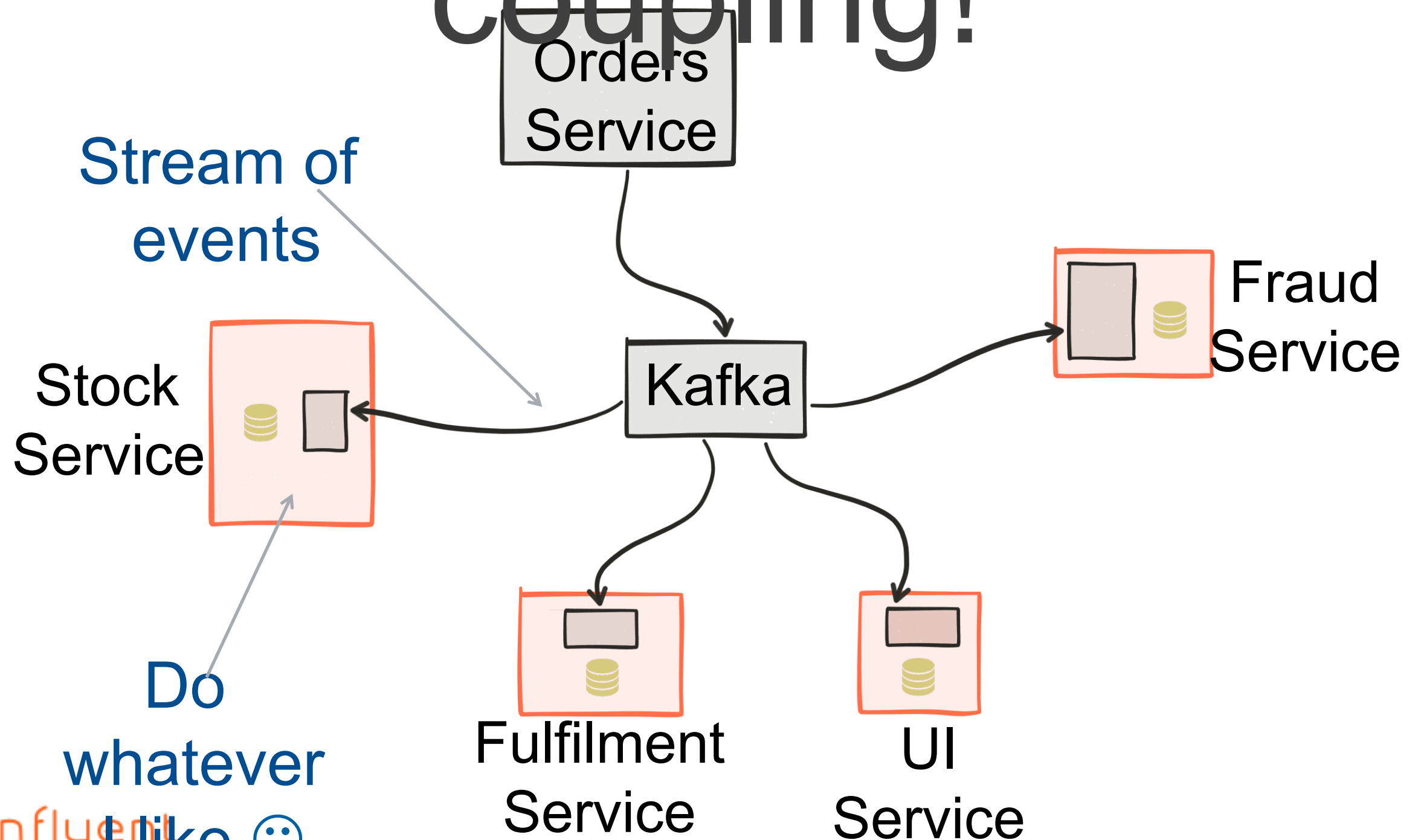
Fraud
Service

Do
whatever

Fulfilment
Service

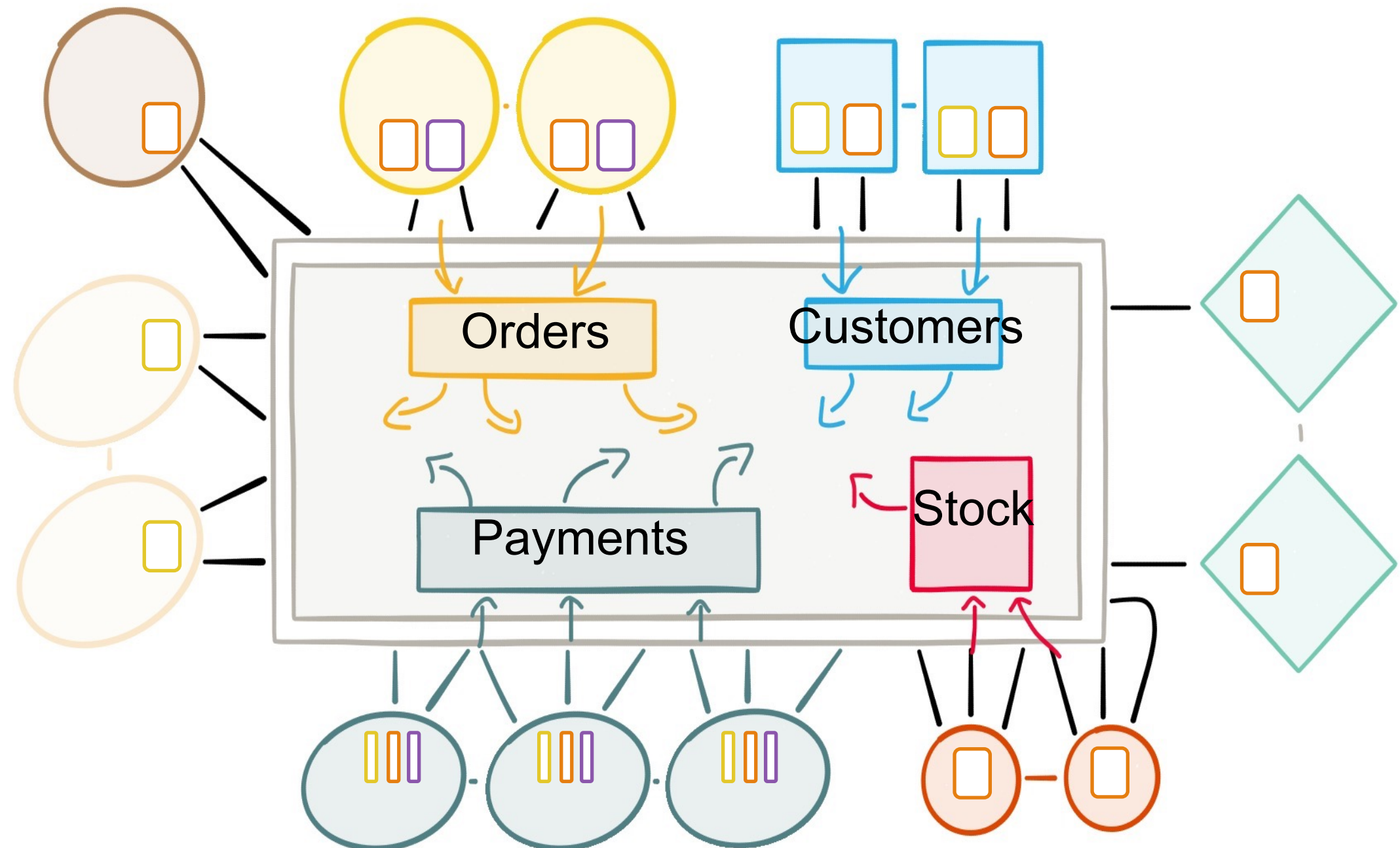
UI
Service

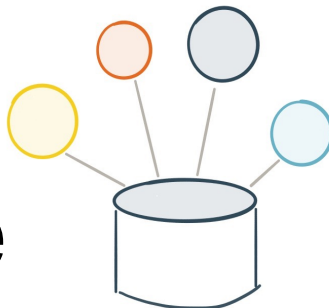



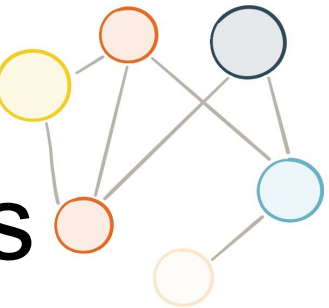



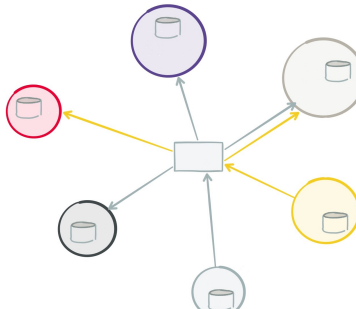



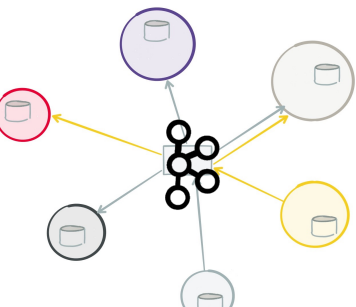



Centralizing Immutable data doesn't affect coupling!



To share a
database,
turn it inside out!

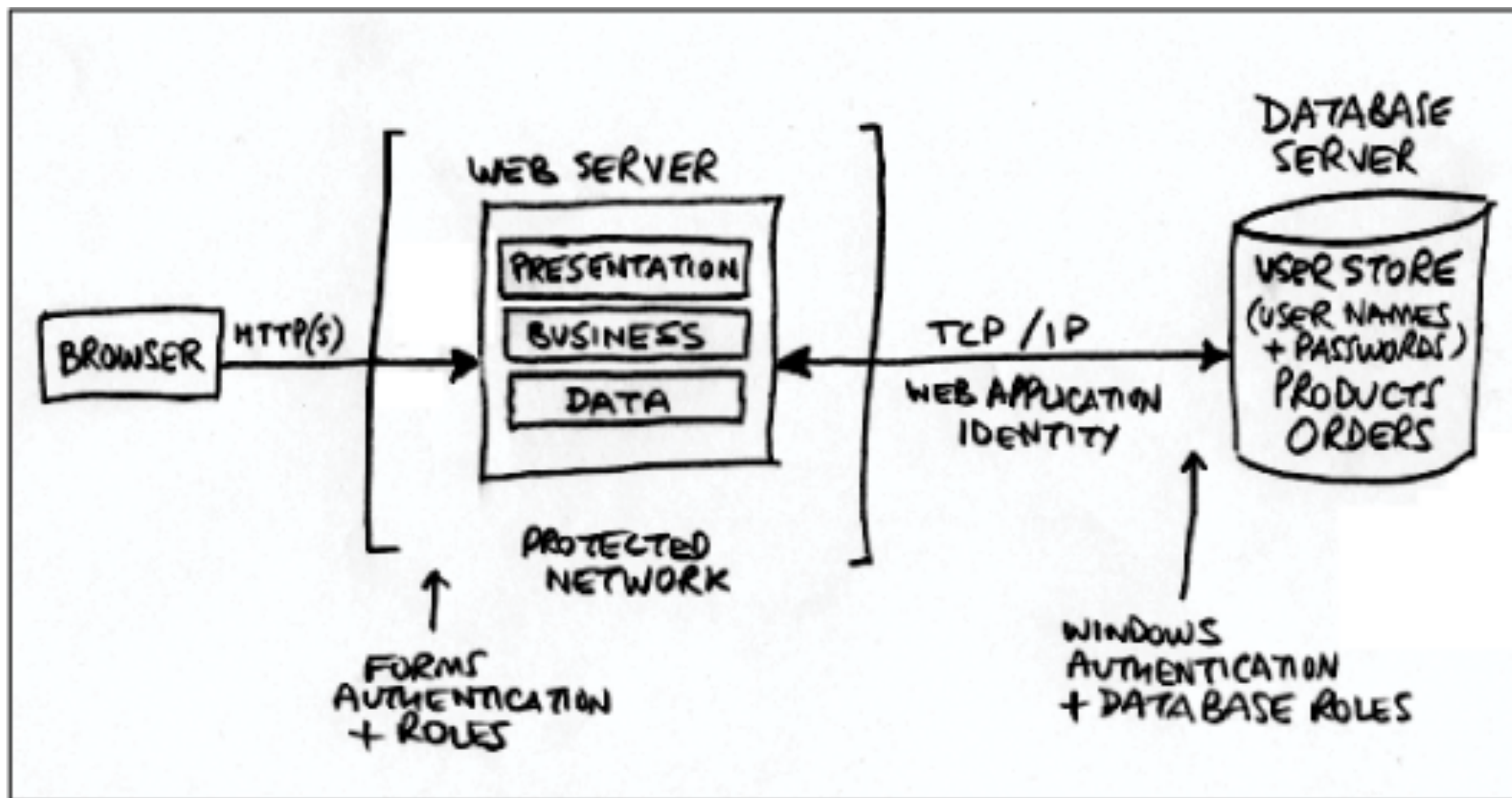
AKA: a machine for creating views



		Ease of change	Data Accessibility	Data Erosion
Shared database				
Service Interfaces				
Event Broadcast Stream				
Data Platform				

So...

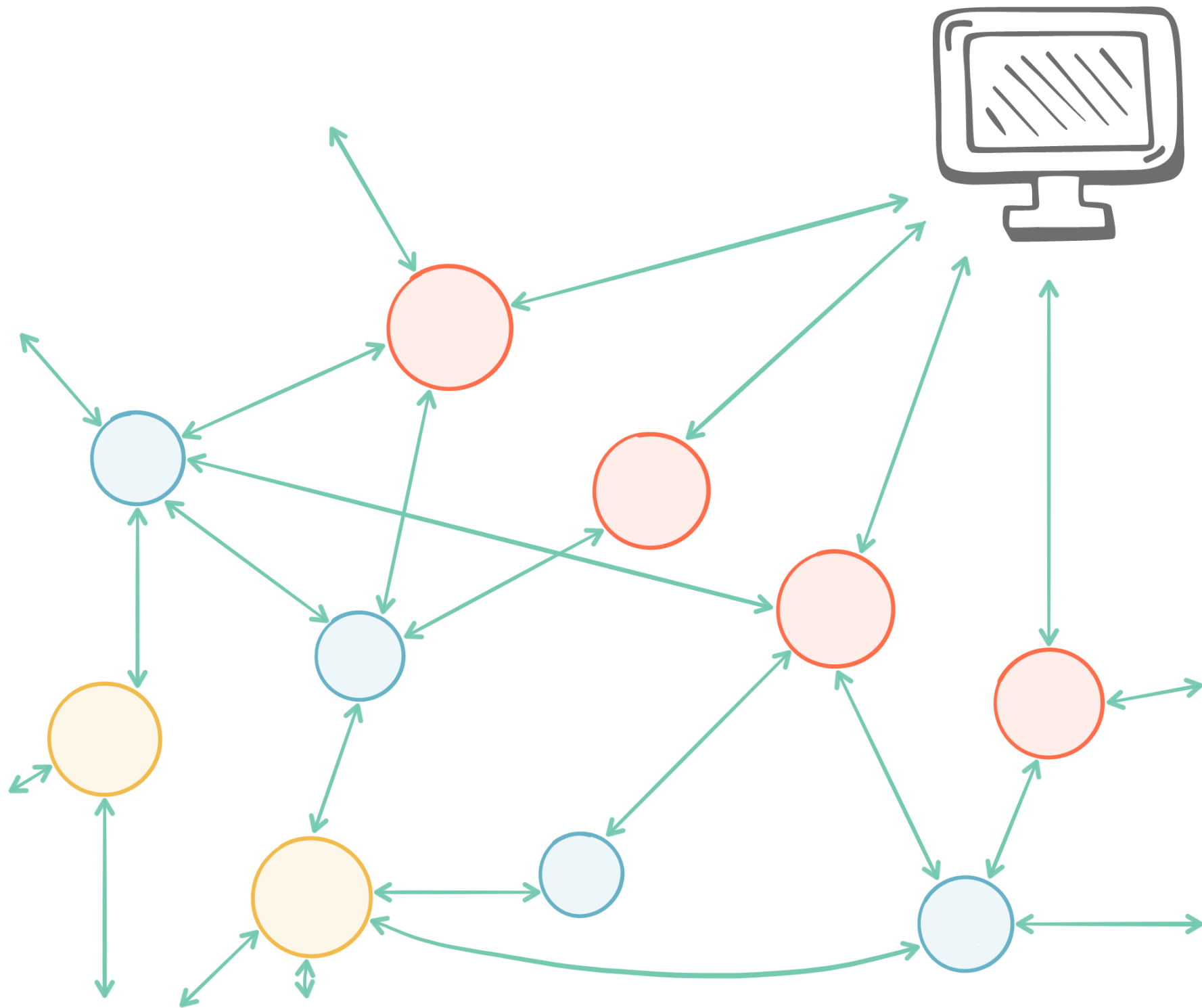
Good architectures have little to do with this:



It's about how systems evolves over time

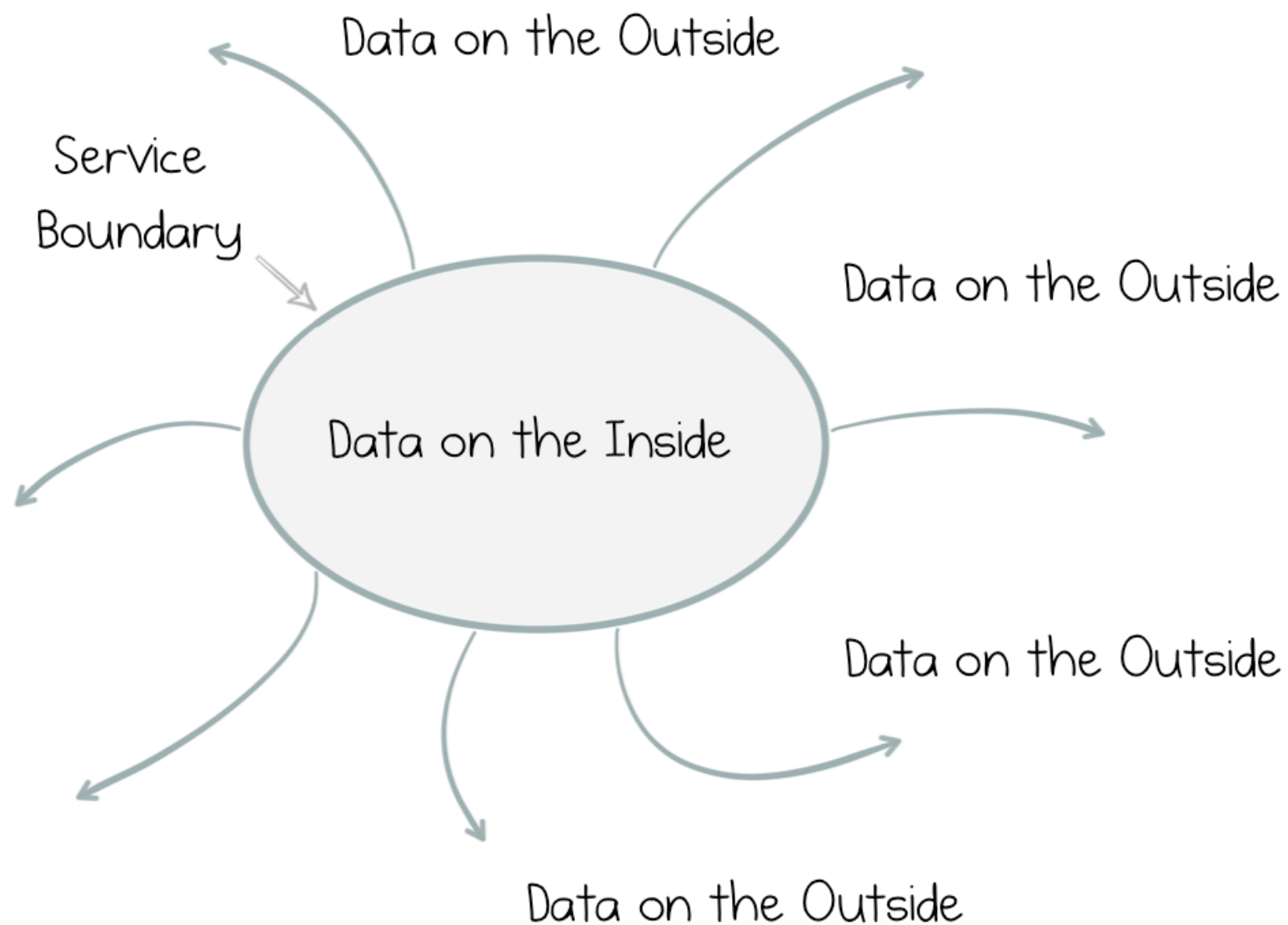


Request driven isn't enough

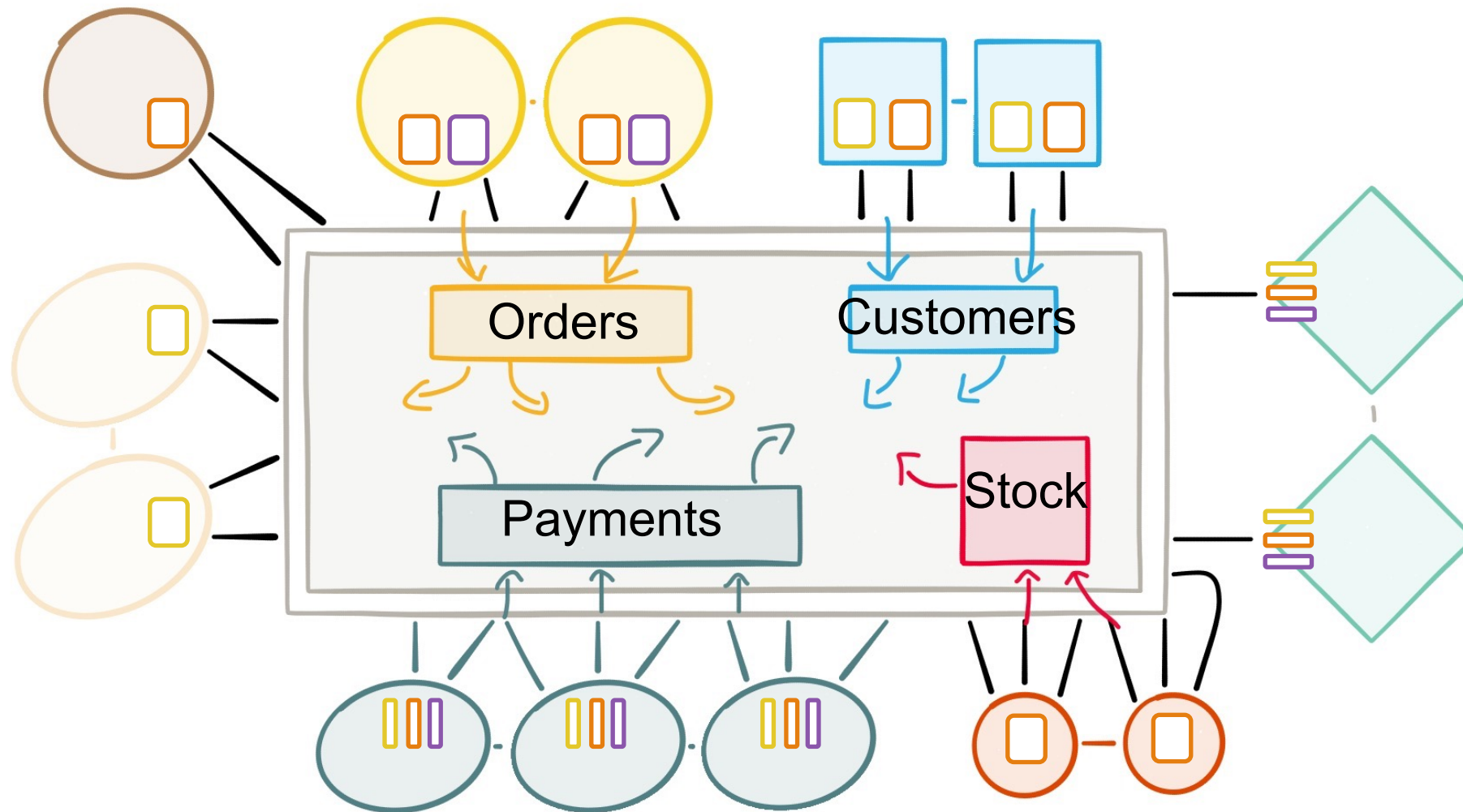


- High coupling
- Hard to handle async flows
- Hard to move and join datasets.

Embrace data that lives and flows between services

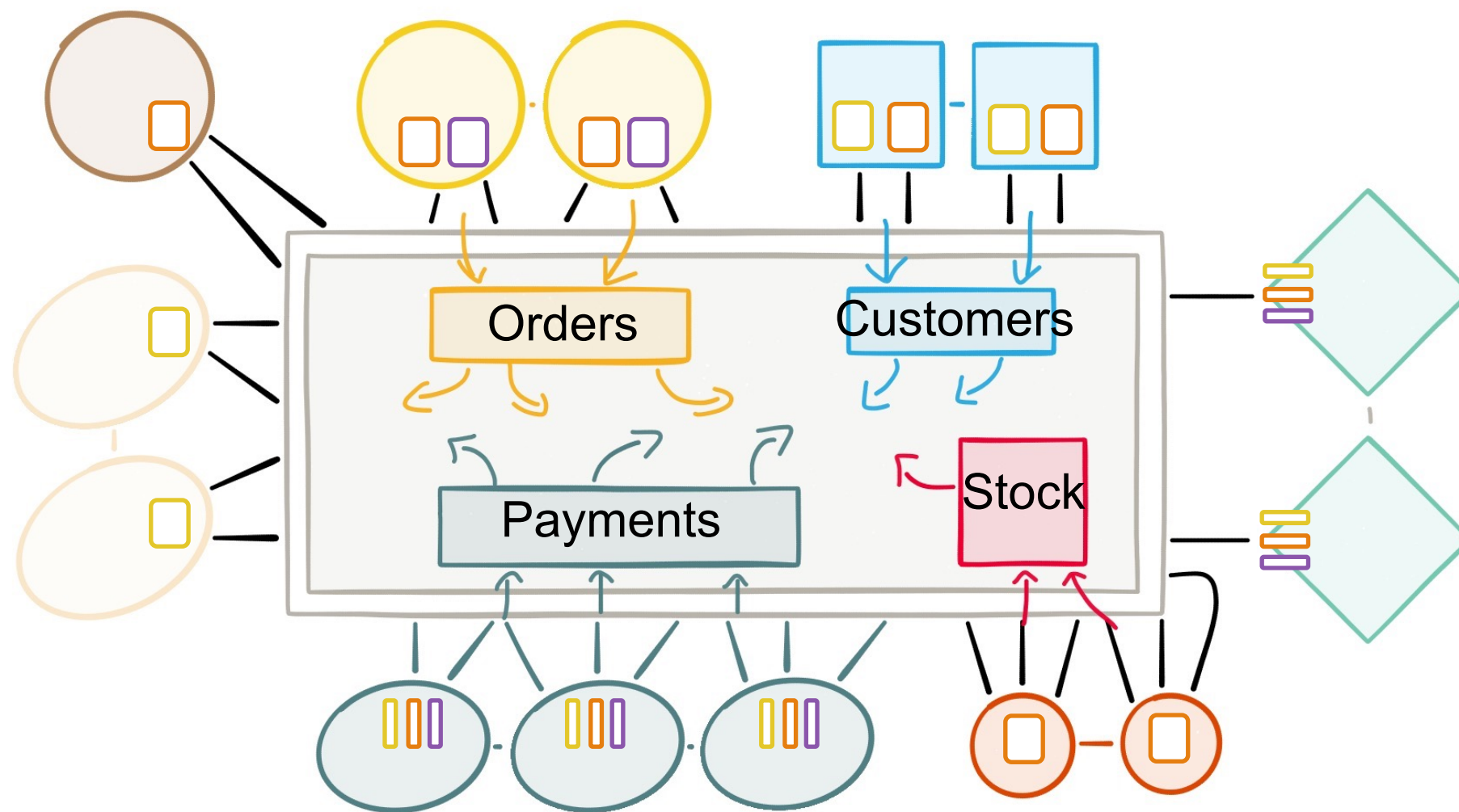


Give services independence



Freedom to tap into and manage shared c

But build on single stream data platform everyone can access and share



Mechanism for evolving an architecture efficiently over time

WIRED Principals

- **W**impy: Start simple, lightweight and fault tolerant.
- **I**mmutable: Build a retentive, shared narrative.
- **R**eactive: Leverage Asynchronicity. Focus on the now.
- **E**volutionary: Use only the data you need today.
- **D**ecentralized: Receiver driven.

References

- Stopford: The Data Dichotomy:
<https://www.confluent.io/blog/data-dichotomy-rethinking-the-way-we-treat-data-and-services/>
- Kleppmann: Turning the Database Inside Out:
<https://www.confluent.io/blog/turning-the-database-inside-out-with-apache-samza/>
- Helland: Immutability Changes Everything:
http://cidrdb.org/cidr2015/Papers/CIDR15_Paper16.pdf
- Helland: Data on the Inside vs Data on the Outside:
<http://cidrdb.org/cidr2005/papers/P12.pdf>
- Kreps: The Log:
<https://engineering.linkedin.com/distributed-systems/log-what-every-software-engineer-should-know-about-real-time-datas-unifying>
- Narkhede: Event Sourcing, CQRS & Stream Processing:
<https://www.confluent.io/blog/event-sourcing-cqrs-stream-processing-apache-kafka-whats-connection/>



Twitter:
@benstopford

Blog series at
<http://confluent.io/blog>

More coming very shortly