

Turning  
the DB

~inside-out~

with

Samza

*dataintensive.net*

O'REILLY®

# Designing Data-Intensive Applications



Martin Kleppmann

@martinkl

Samza.incubator.  
apache.org

Samza

Typical web app.

Browser / client app

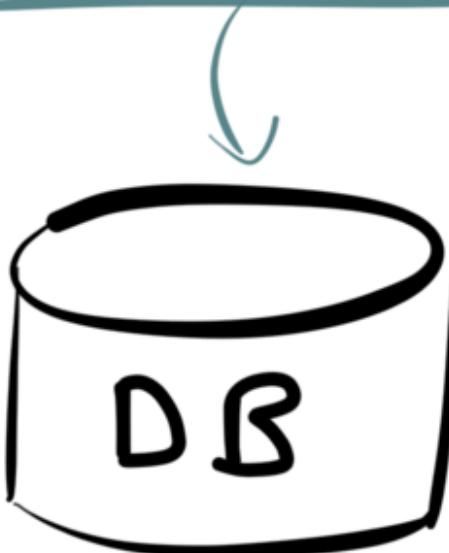


Typical web app.

Browser / client app

HTTP (stateless)

"Backend" (stateless)

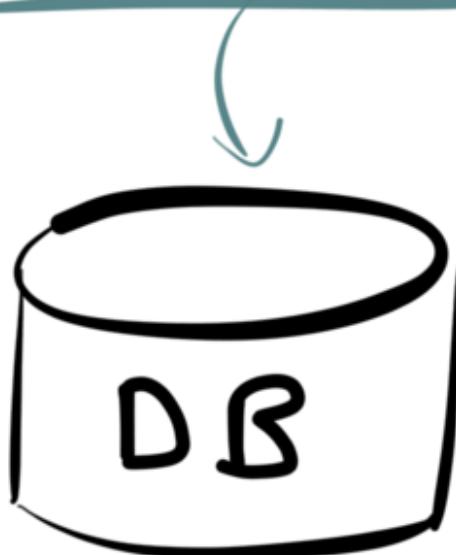


Typical web app.

Browser / client app

HTTP (stateless)

"Backend" (stateless)



(mutable  
global state)

"It's always been  
that way" ?



# 1. Replication

# Shopping cart

customer_id	product_id	quantity
123	888	1
123	999	1
234	444	2
234	555	3
345	666	1

update cart set quantity = 3  
where customer\_id=123 and  
product\_id=999

customer_id	product_id	quantity
123	888	1
123	999	1
234	444	2
234	555	3
345	666	1

update cart set quantity = 3  
where customer\_id=123 and  
product\_id=999

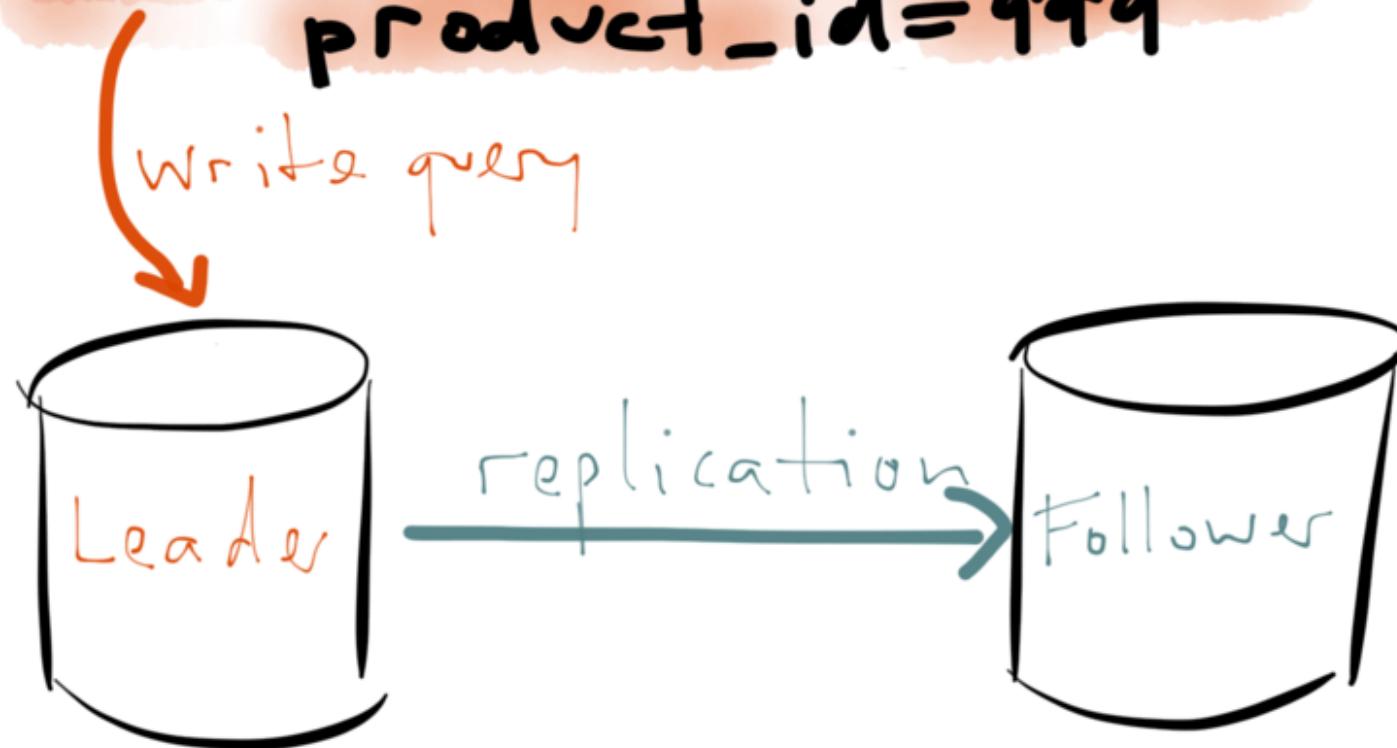
customer_id	product_id	quantity
123	888	1
123	999	3
234	444	2
234	555	3
345	666	1

update cart set quantity = 3  
where customer\_id=123 and  
product\_id=999

write query



update cart set quantity = 3  
where customer\_id=123 and  
product\_id=999



update cart set quantity = 3  
where customer\_id=123 and  
product\_id=999

write query



Change row 8765

old=[123, 999, 1]

new=[123, 999, 3]

update cart set quantity = 3  
where customer\_id=123 and  
product\_id=999

---

VS.

---

"at 2014-09-18 14:10:59,  
customer 123 changed quantity of  
product 999 in their cart  
from 1 to 3"

update cart set quantity = 3  
where customer\_id=123 and  
product\_id=999

(imperative state mutation)

VS.

"at 2014-09-18 14:10:59,  
customer 123 changed quantity of  
product 999 in their cart  
from 1 to 3"

(immutable event / fact)

2.

Secondary  
indexes



customer	product	qty
123	888	1
123	999	3
234	444	2
234	555	3

customer	product	qty
123	888	1
123	999	3
234	444	2
234	555	3

CREATE INDEX ON  
cart (customer\_id);  
CREATE INDEX ON  
cart (product\_id);

customer	product	qty
123	888	1
123	999	3
234	444	2
234	555	3



123 → [123 | 888 | 1]

123 → [123 | 999 | 3]

234 → [234 | 444 | 2]

234 → [234 | 555 | 3]

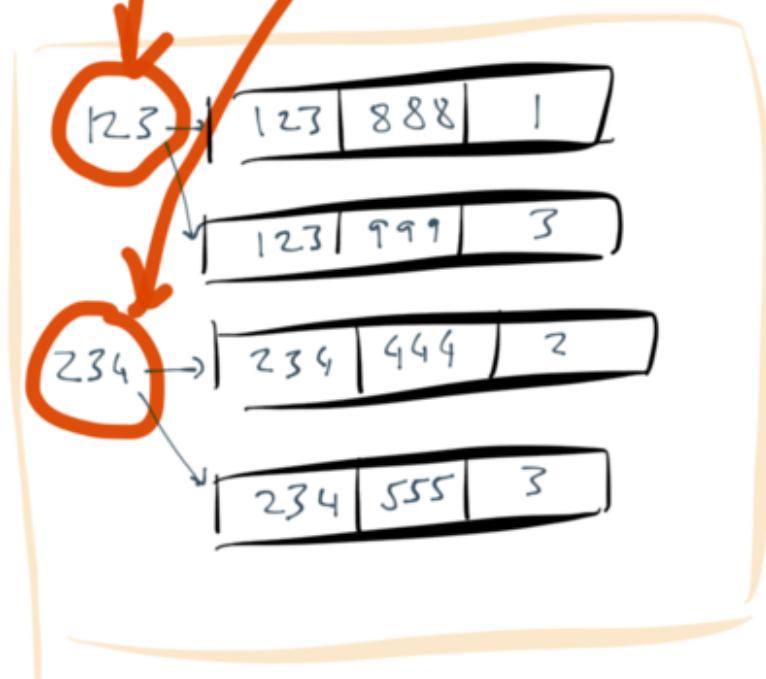
444 → [234 | 444 | 2]

555 → [234 | 555 | 3]

888 → [123 | 888 | 1]

999 → [123 | 999 | 3]

customer	product	qty
123	888	1
123	999	3
234	444	2
234	555	3



444	234	444	2
555	234	555	3
888	123	888	1
999	123	999	3

customer	product	qty
123	888	1
123	999	3
234	444	2
234	555	3

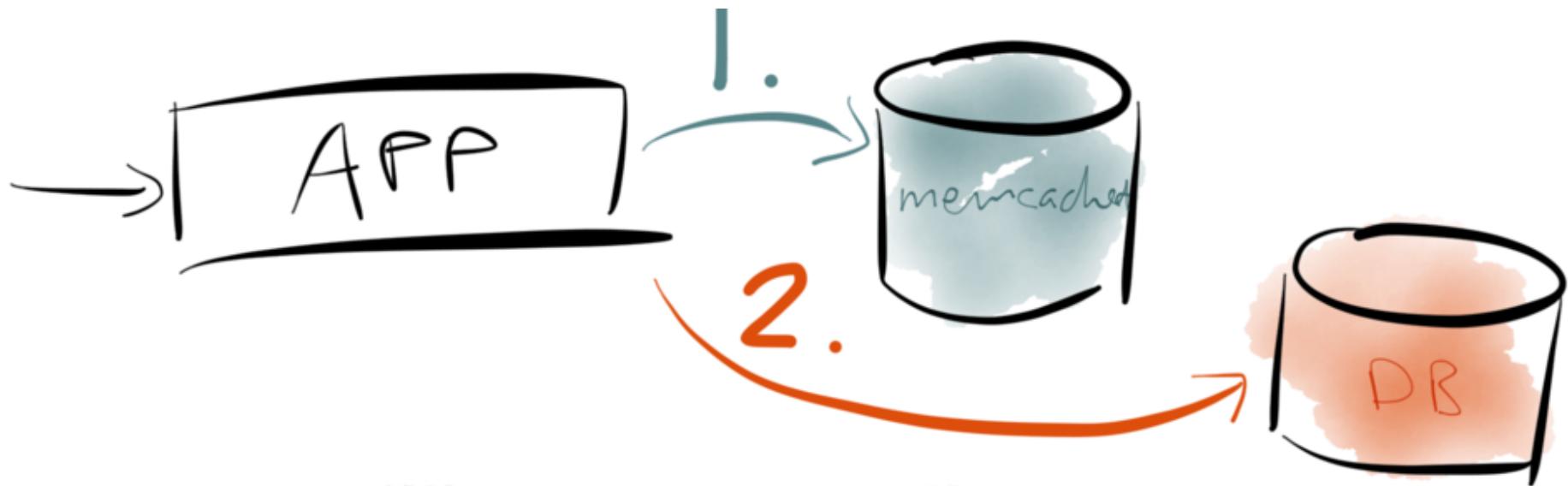
CREATE INDEX  
CONCURRENTLY...

3.

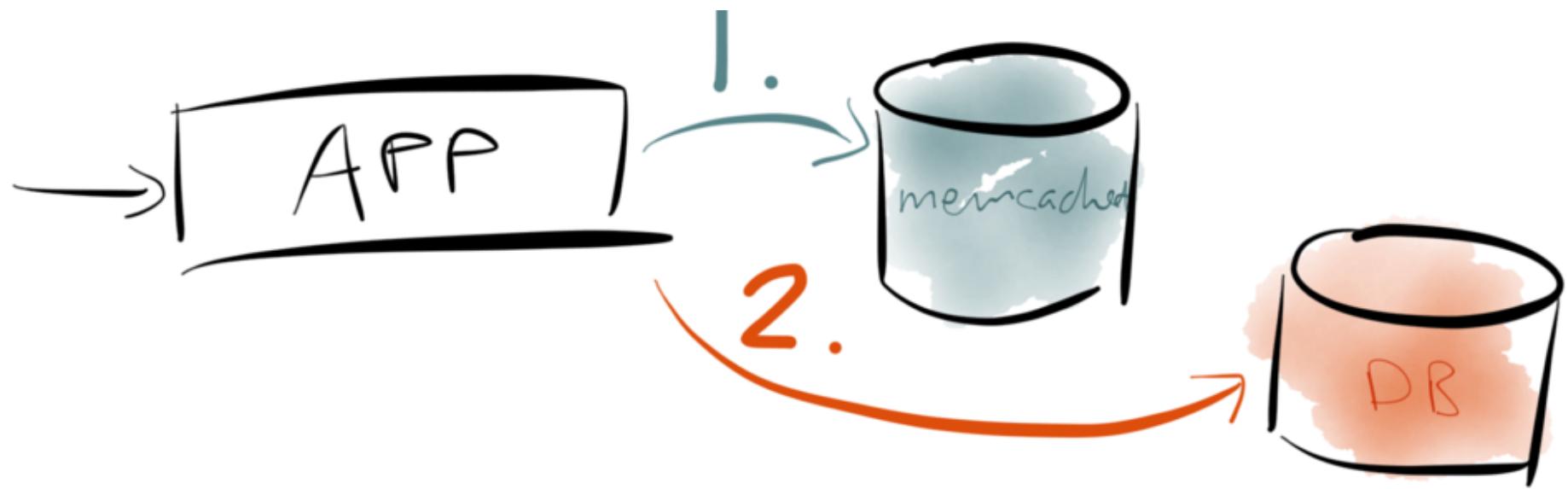
Caching



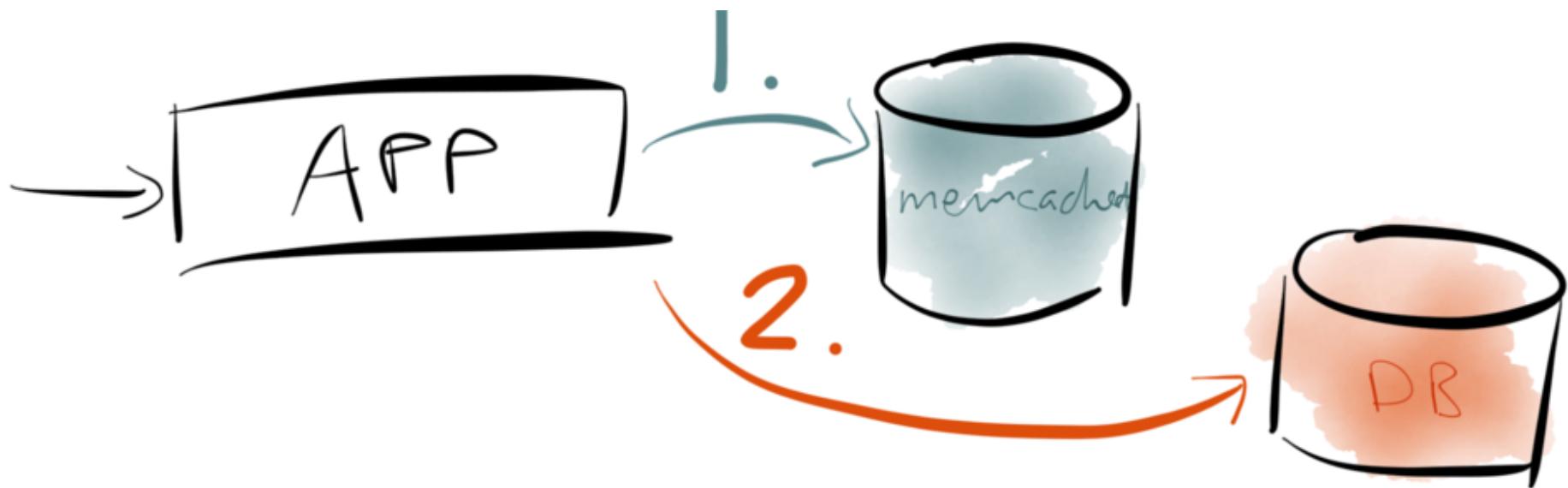
r = cache.get(key)



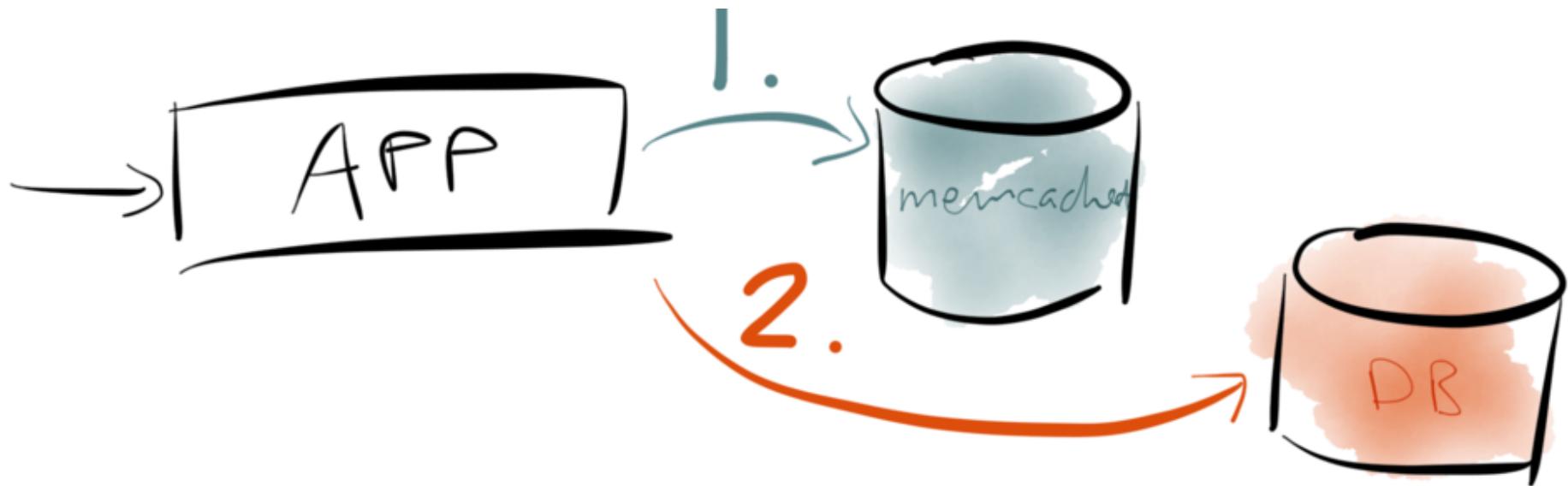
```
r = cache.get(key)
if (!r) {
    r = db.get(key)
    cache.put(key, r)
}
return r;
```



- invalidation?



- invalidation?
- race conditions / consistency issues?

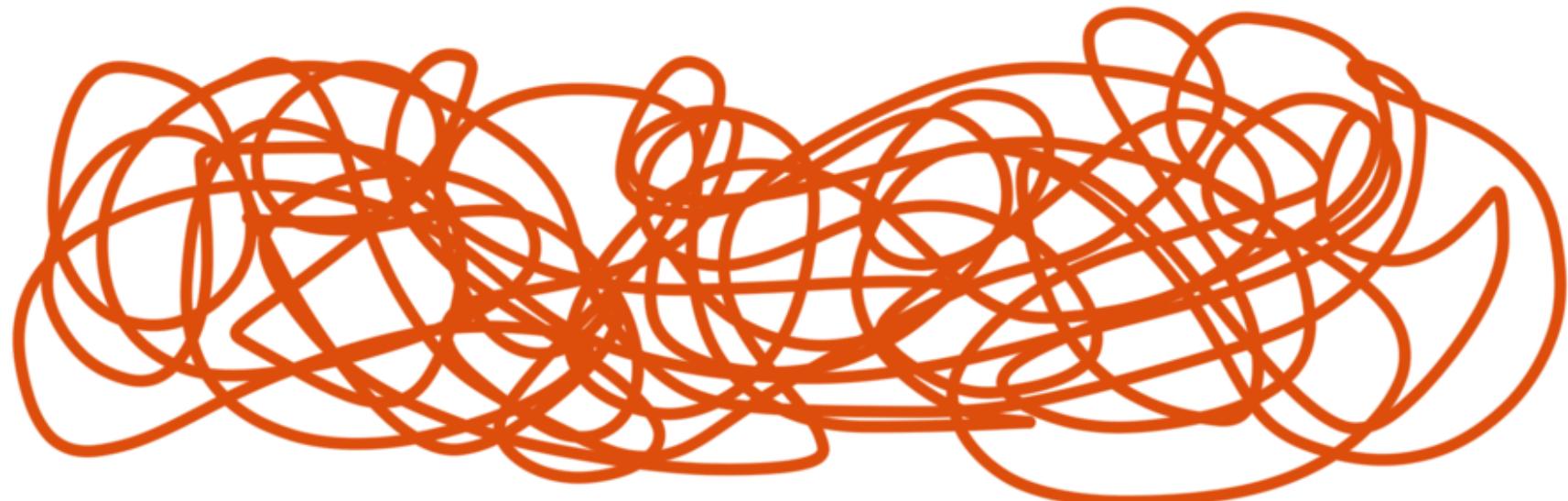


- invalidation?
- race conditions / consistency issues?
- cold start / bootstrapping

**CREATE INDEX ...**

(1 line of SQL!)

vs.





# 4. Materialized views



CREATE VIEW example (foo)

AS SELECT foo  
FROM bar  
WHERE ...

CREATE VIEW example (foo)

AS SELECT foo  
FROM bar  
WHERE ...

query:  
SELECT \* FROM example

CREATE VIEW example (foo)

AS SELECT foo  
FROM bar  
WHERE ...

query:

SELECT x FROM example

(Rewrite at query time)

SELECT foo FROM bar  
WHERE ...

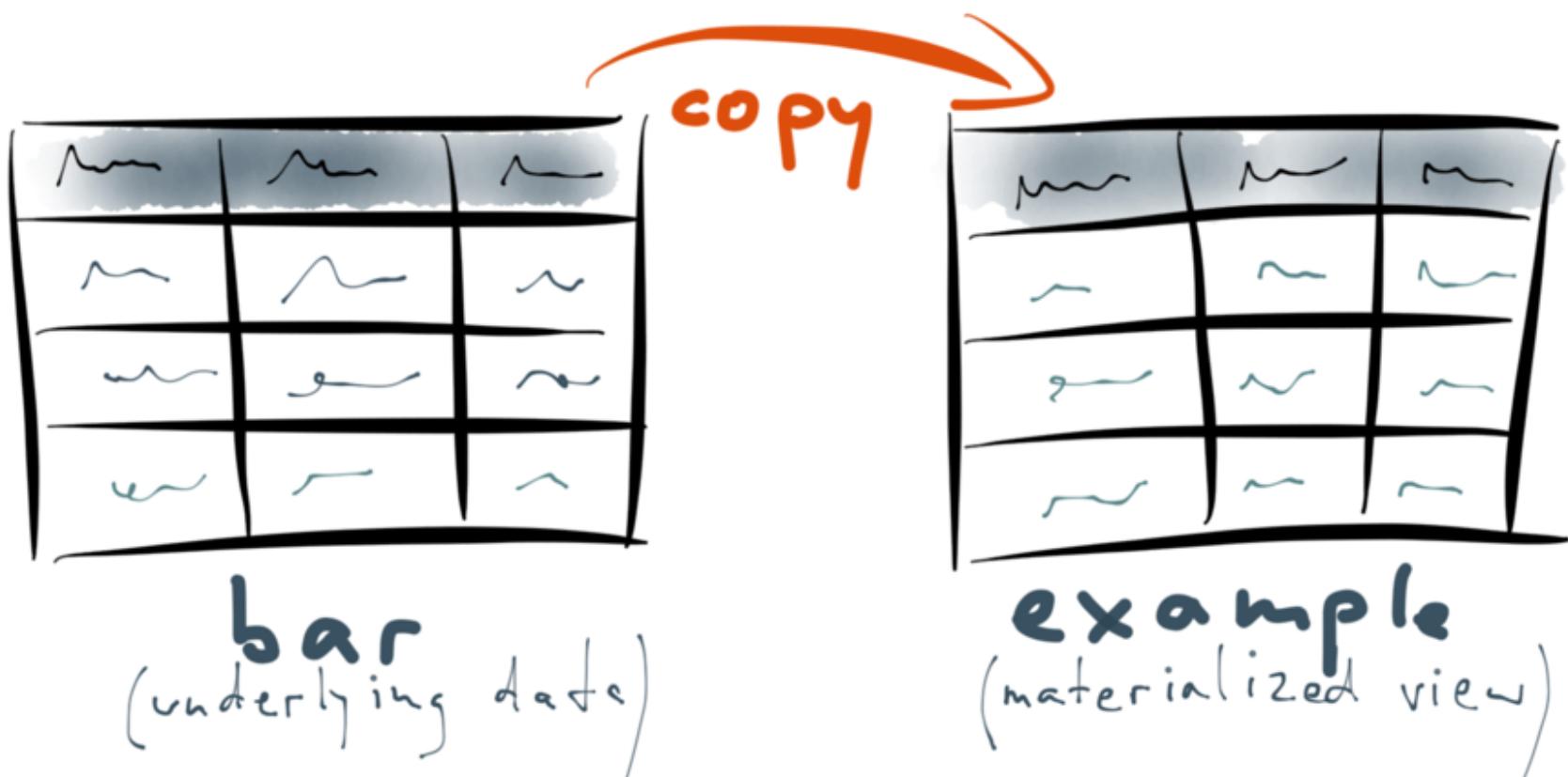
CREATE MATERIALIZED  
VIEW example (foo)  
AS SELECT foo FROM bar ...

CREATE MATERIALIZED  
VIEW example (foo)  
AS SELECT foo FROM bar ...

m	m	m
m	m	m
m	m	m
m	m	m

bar

CREATE MATERIALIZED  
VIEW example (foo)  
AS SELECT foo FROM bar ...



- ① Replication
- ② Secondary indexing
- ③ Caching
- ④ Materialized views

- ① Replication
  - ② Secondary indexing
  - ③ Caching
  - ④ Materialized views
- DERIVED DATA

- ① Replication
  - ② Secondary indexing
  - ③ Caching
  - ④ Materialized views
- DERIVED DATA



- ① Replication
  - ② Secondary indexing
  - ③ Caching
  - ④ Materialized views
- DERIVED DATA

- ① Replication
  - ② Secondary indexing
  - ③ Caching
  - ④ Materialized views
- DERIVED DATA



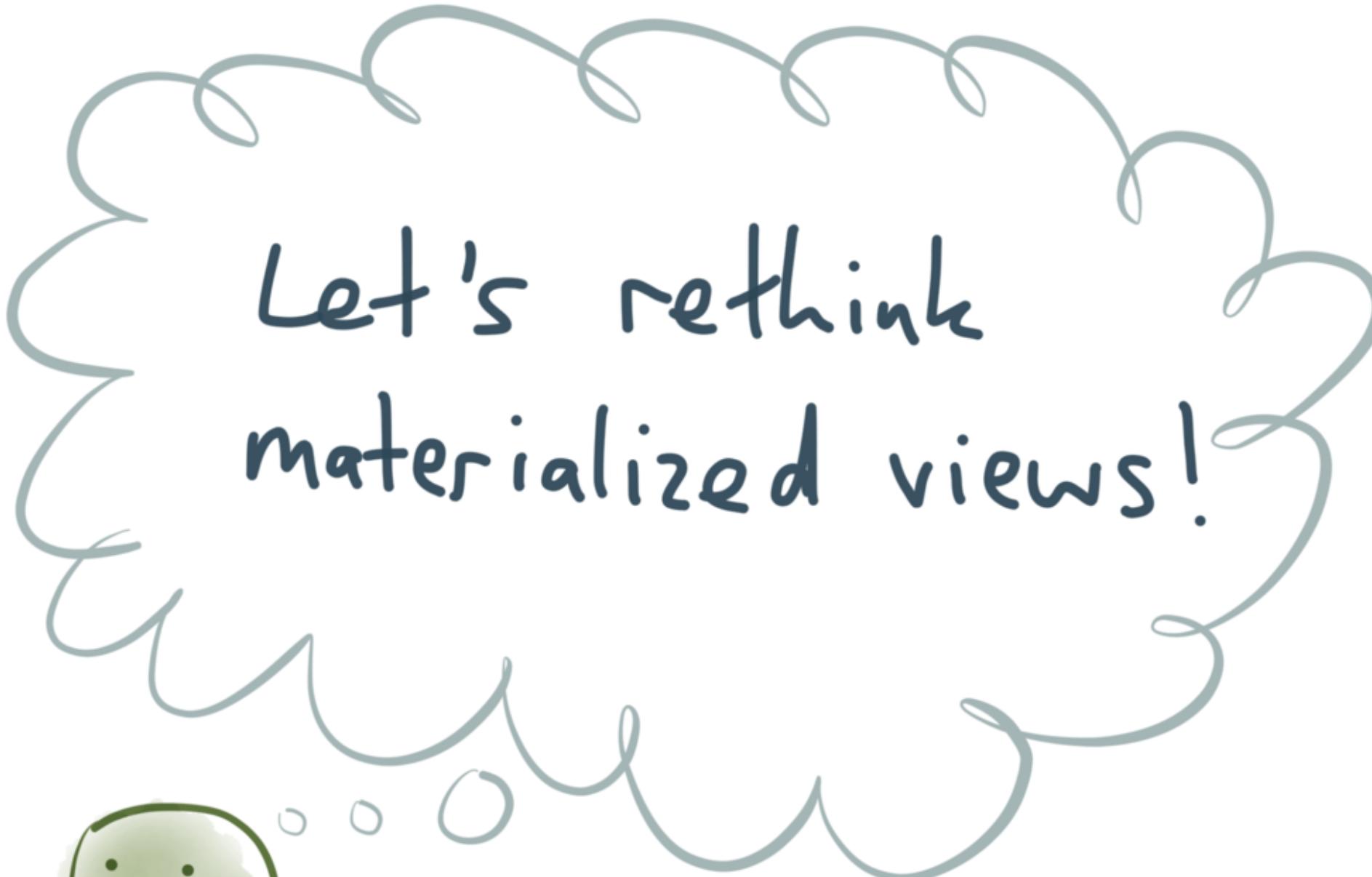
- ① Replication
  - ② Secondary indexing
  - ③ Caching
  - ④ Materialized views
- DERIVED DATA





Magically  
self-updating  
cache ...



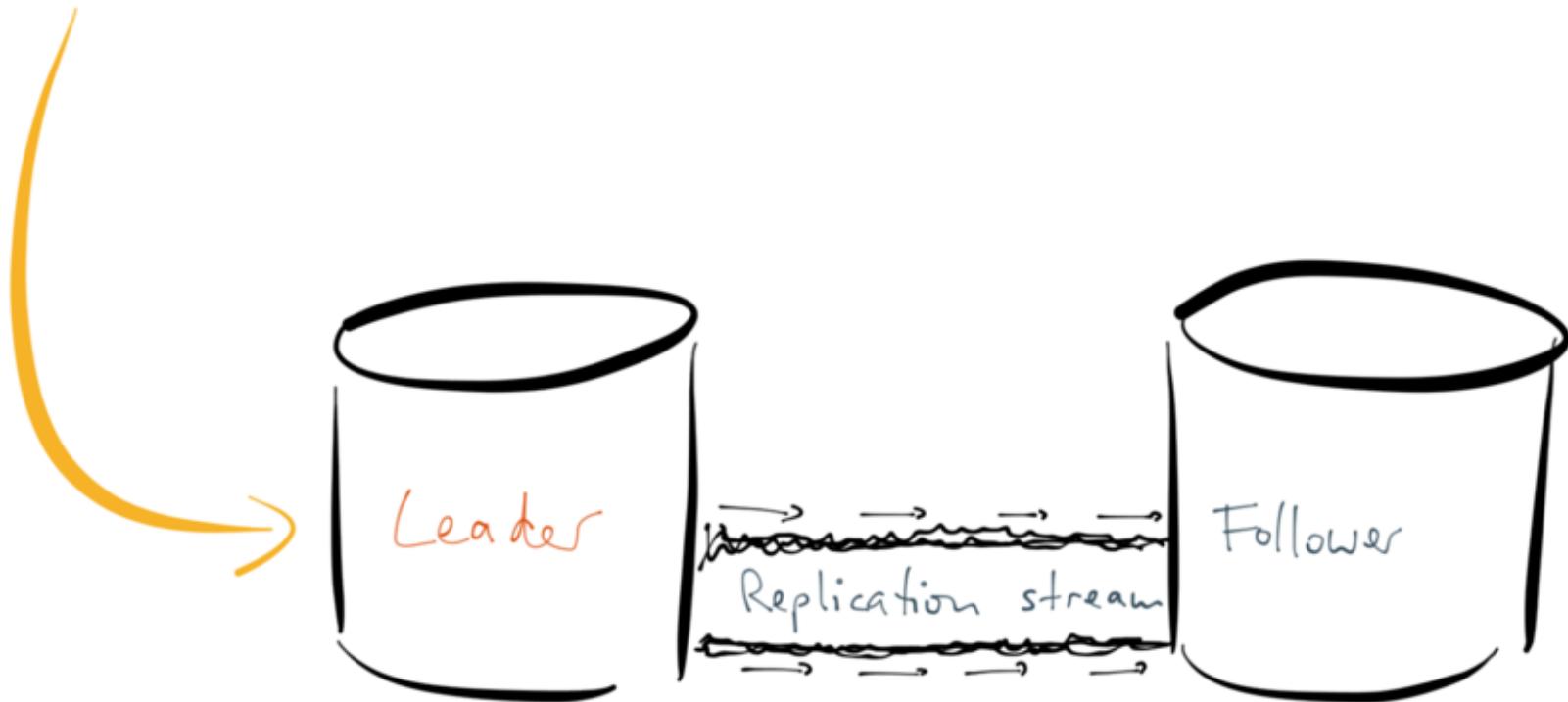


Let's rethink  
materialized views!



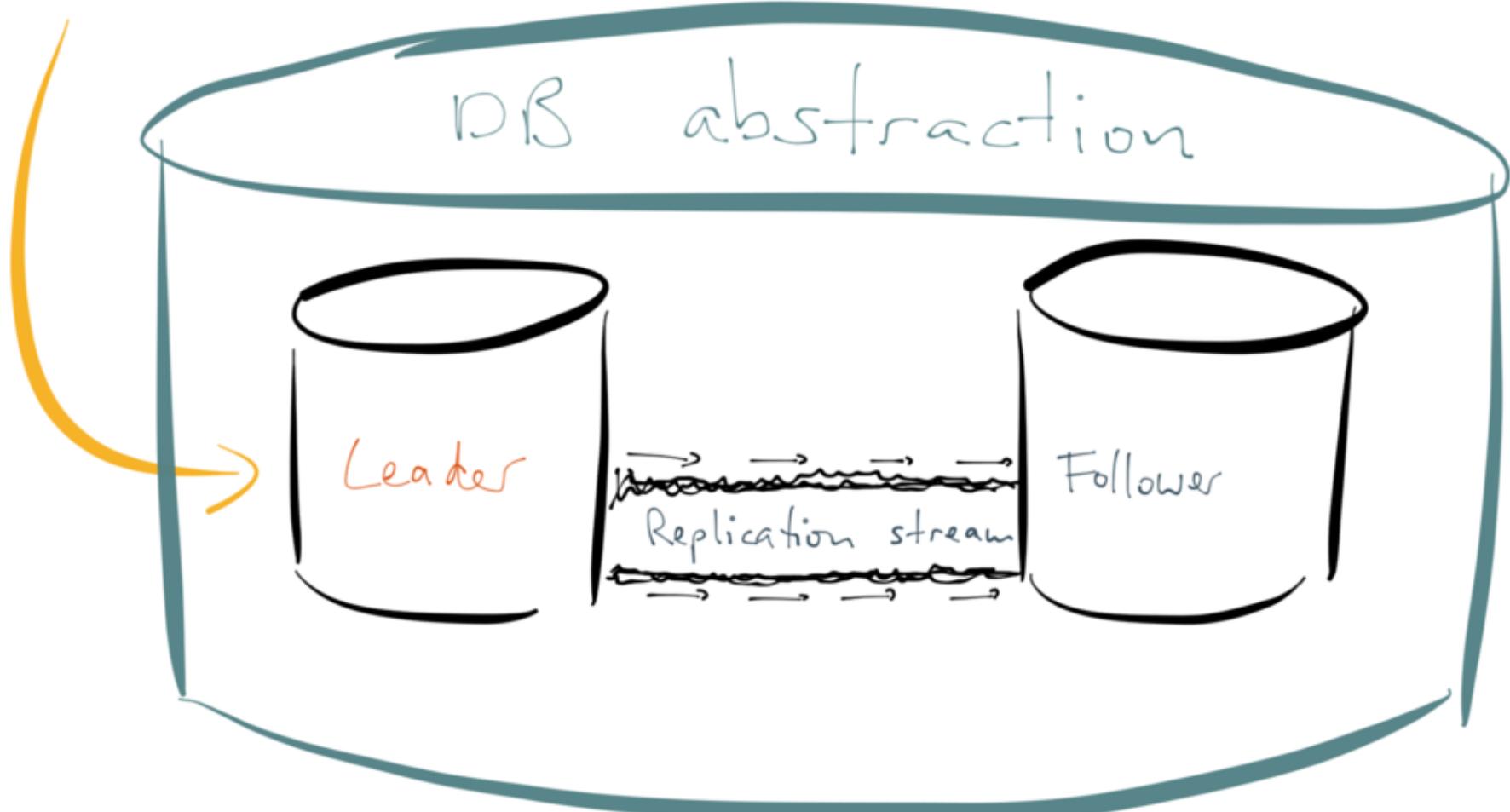
writes

traditional DB



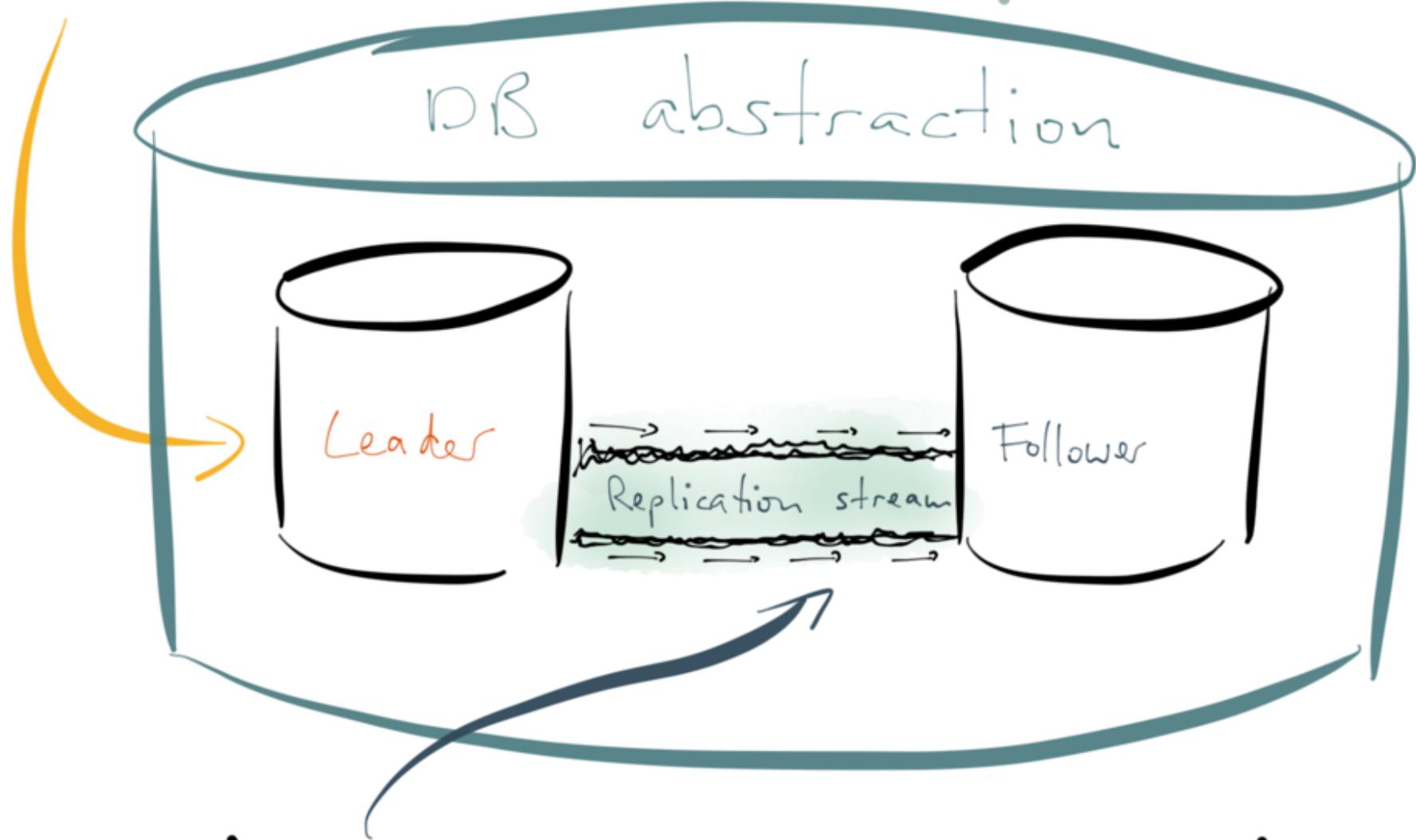
writes

traditional DB



writes

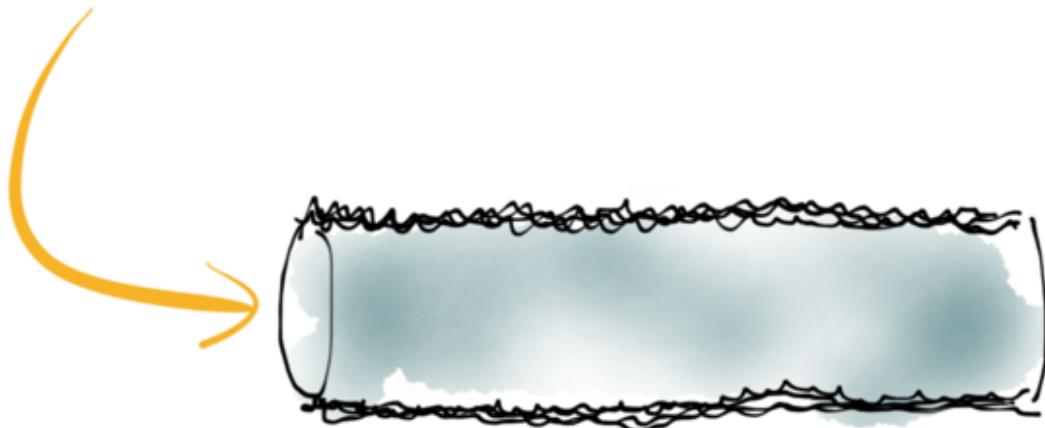
traditional DB



implementation detail.

"unbundled" DB

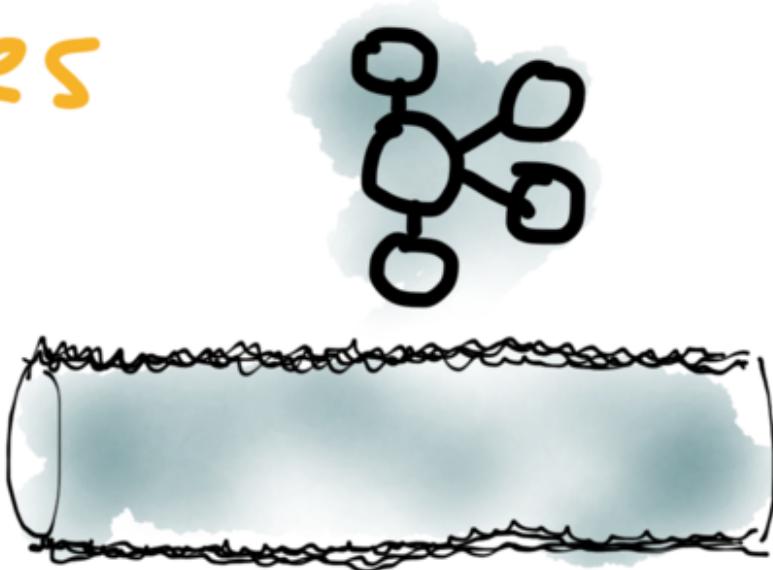
writes



Event  
Stream/  
transaction  
log

"unbundled" DB

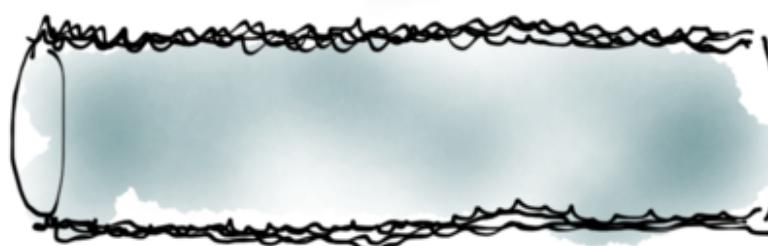
writes



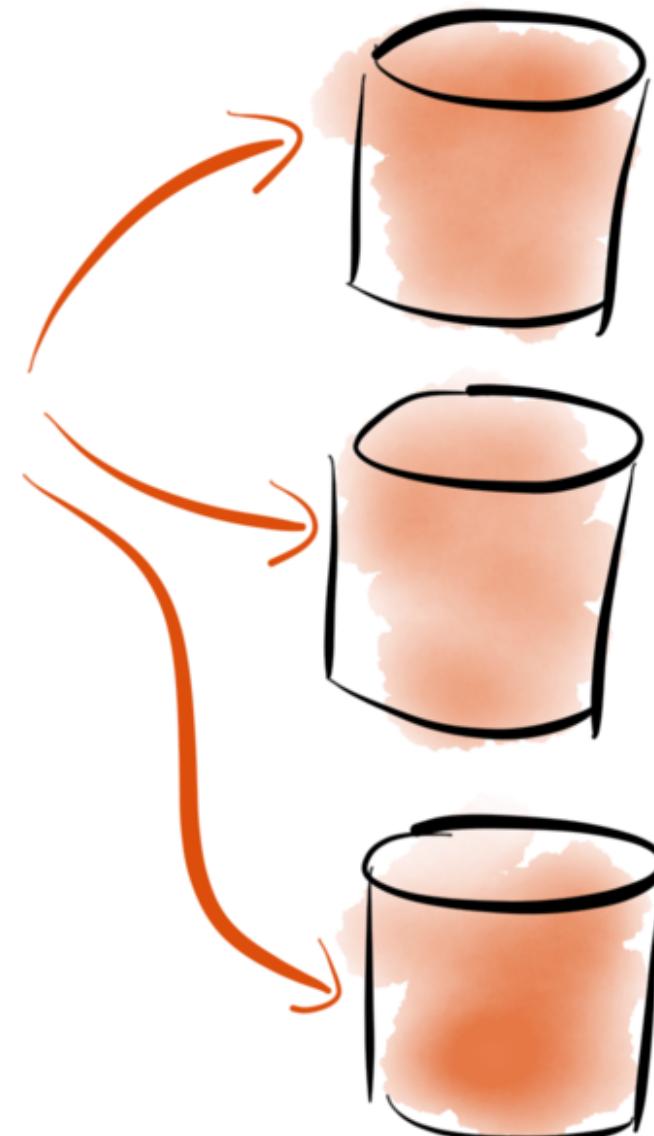
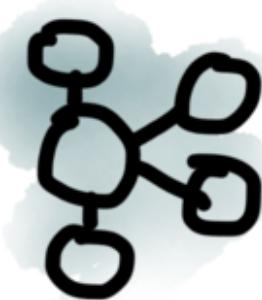
Event  
Stream/  
transaction  
log

# "unbundled" DB

writes

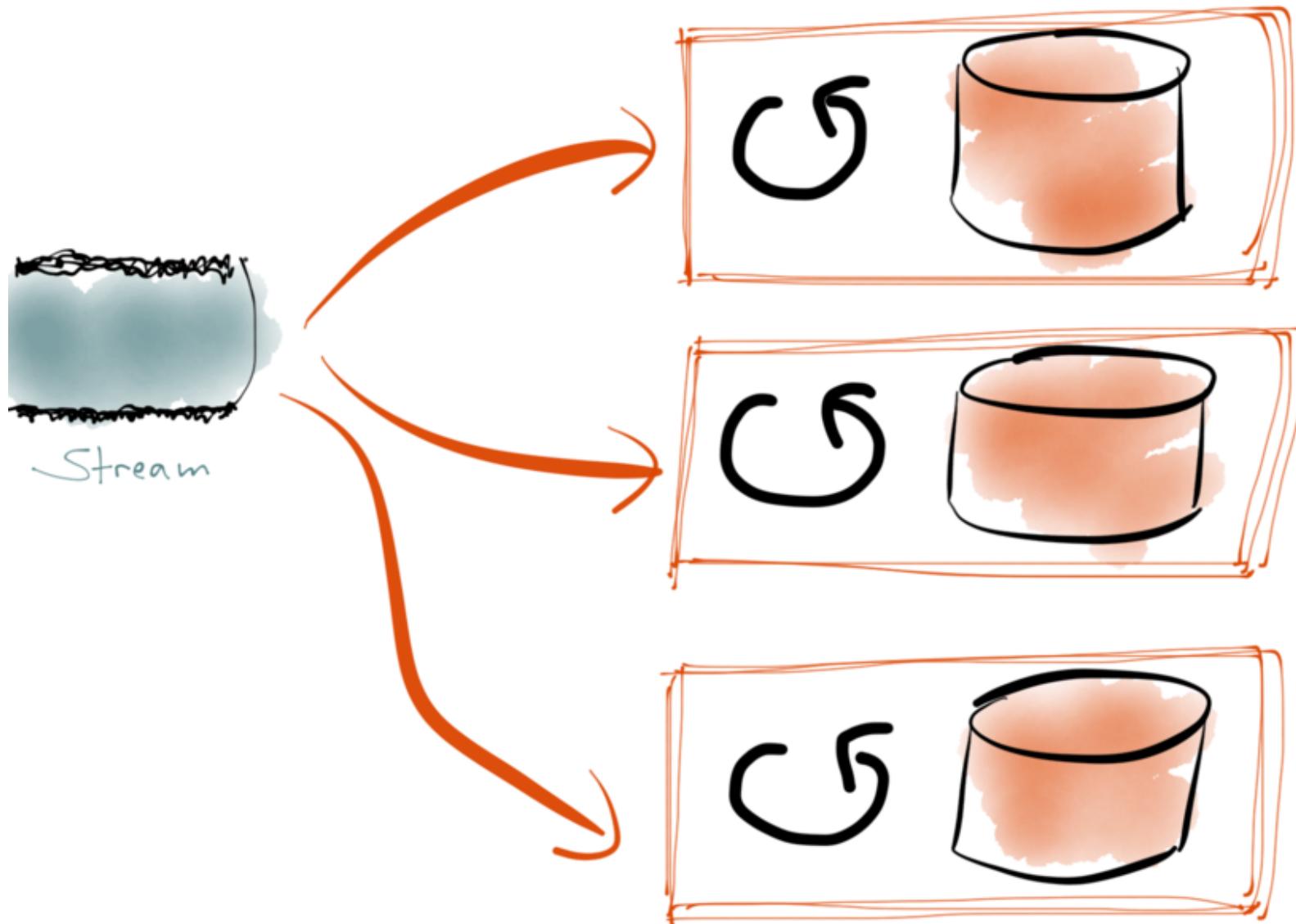


Event  
Stream/  
transaction  
log

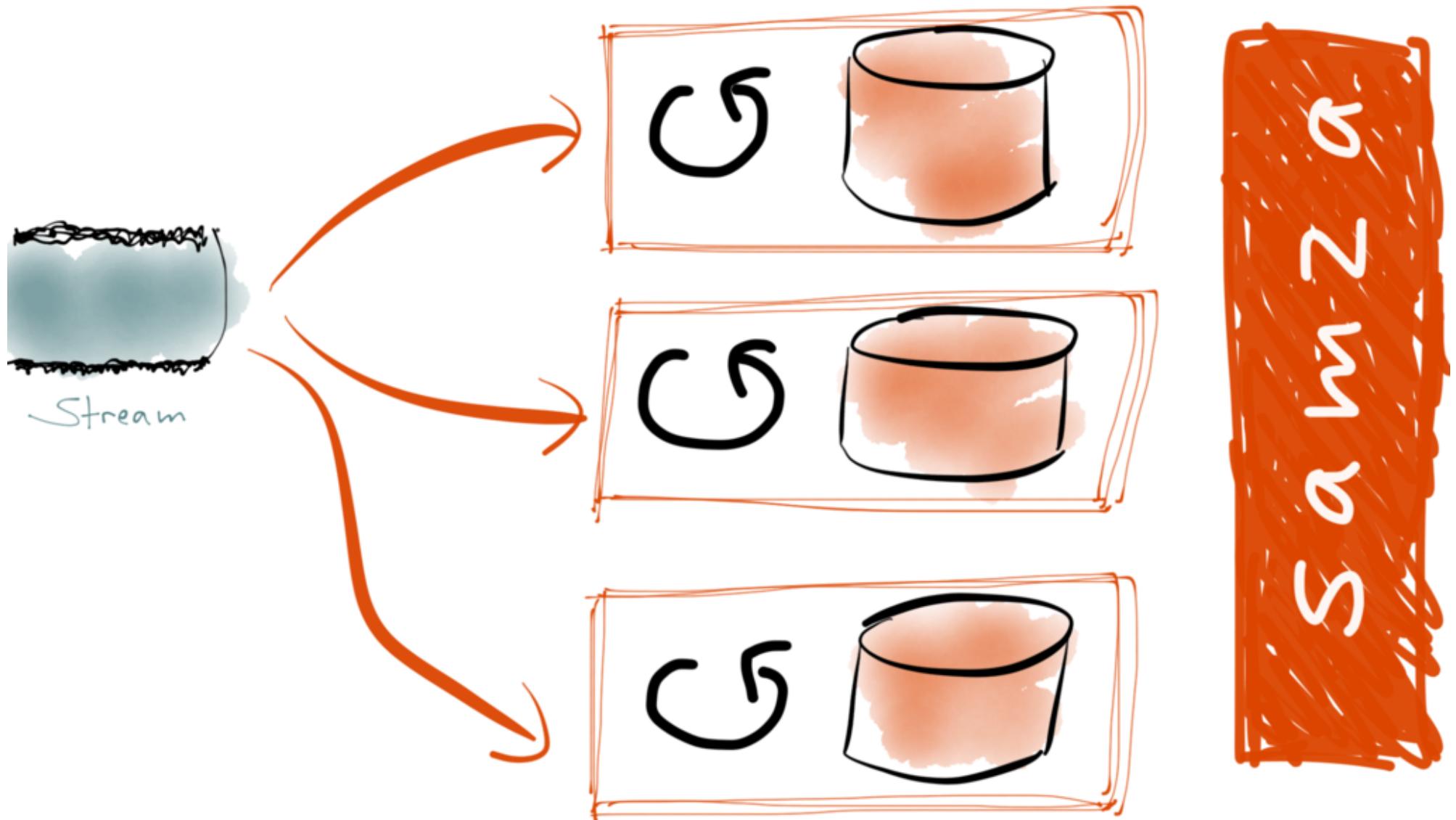


materialized views

# "unbundled" DB



# "unbundled" DB



Writes:

~~mutate state in place~~

## Writes:

~~mutate state in place~~

append-only stream of  
immutable facts

— Simple, Scalable, fault-tolerant, fast sequential I/O

(see also: Datalog, Datomic, Lambda Architecture)

## Writes:

~~mutate state in place~~

append-only stream of  
immutable facts

— Simple, Scalable, fault-tolerant, fast sequential I/O

(see also: Datalog, Datomic, Lambda Architecture)

→ Apache Kafka

Reads:

~~query the database~~

Reads:

~~query the database~~

- consume & join streams
- maintain materialized views

Reads:

~~query the database~~

- consume & join streams
- maintain materialized views

→ Apache Samza

## Cache vs. materialized view:

~~read-through cache~~  
~~updated by app~~

## Cache vs. materialized view:

~~read-through cache~~

~~updated by app~~

-materialized view = precomputed

# Cache vs. materialized view:

read-through cache

updated by app

-materialized view = precomputed

-to build new view, consume  
input since beginning of time

(just like building a new DB index)

κ

# Mechanics:

deployment monitoring

cluster operations integrations

fault-tolerance performance

security onboarding isolation

tooling documentation support

...being worked on

(contributions  
welcome!)



Why?

1.

Better data

# 1.

# Better data

- Good for analytics
- Separation of concerns between writing and reading
- Write once, read from many different views
- Historical point-in-time queries
- Recovery from human error (eg. bad code deployed)

②.

Fully  
precomputed  
caches

②.

Fully  
precomputed  
caches

No cold cache/warming, no hit/miss rate,  
no race conditions, no complex invalidation logic,  
better isolation/robustness, etc ...

③.

Streams  
everywhere!

DB



DB



raw data as written

Business logic

Cache



DB



raw data as written

Business logic

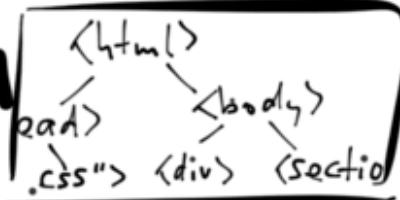
Cache



eg. JSON over REST

UI logic

HTML DOM



template rendering

DB



raw data as written

Business logic

Cache



eg. JSON over REST

UI logic

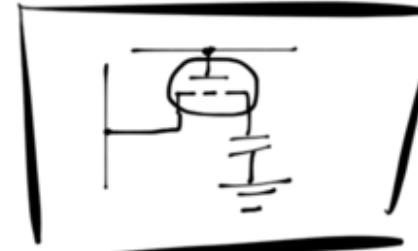
HTML DOM

`<html>`  
  `<head>`   `<body>`  
    `<div>`   `<section>`  
  `<css>`

template rendering

Browser rendering

Video mem



pixels

DB



raw data as written

Business logic

Cache



eg. JSON over REST

UI logic

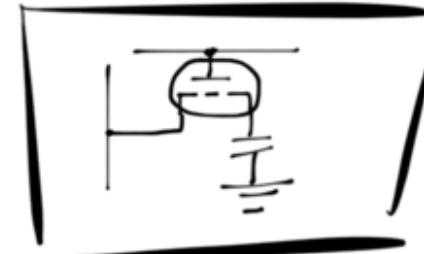
HTML DOM

`<html>  
 <head> <body>  
 <div> <section>  
 <css>`

template rendering

Browser rendering

Video mem



pixels



DB



raw data as written

Business logic

Cache

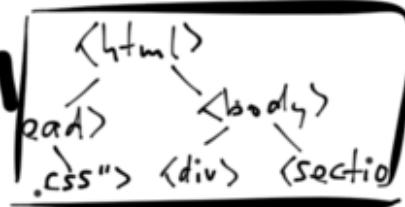


eg. JSON over REST

UI logic



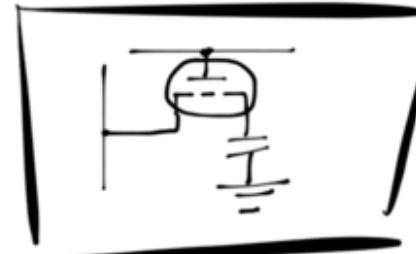
HTML DOM



template rendering

Browser rendering

Video mem



pixels



DB



raw data as written

Business logic



Cache

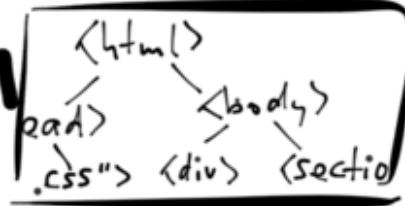


eg. JSON over REST

UI logic



HTML DOM

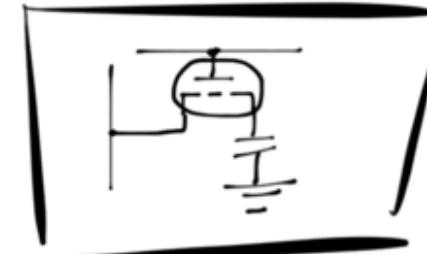


template rendering

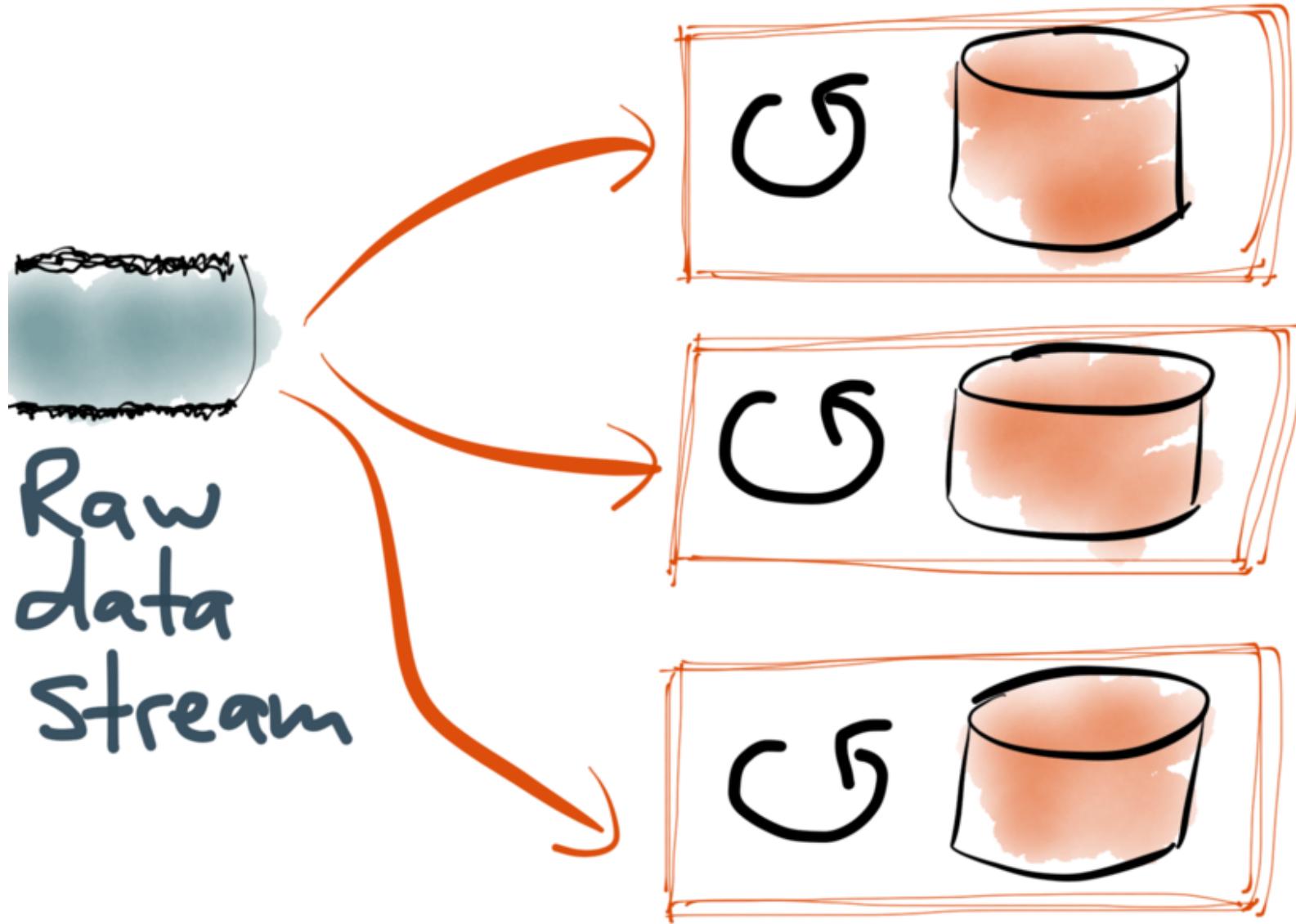
Browser rendering



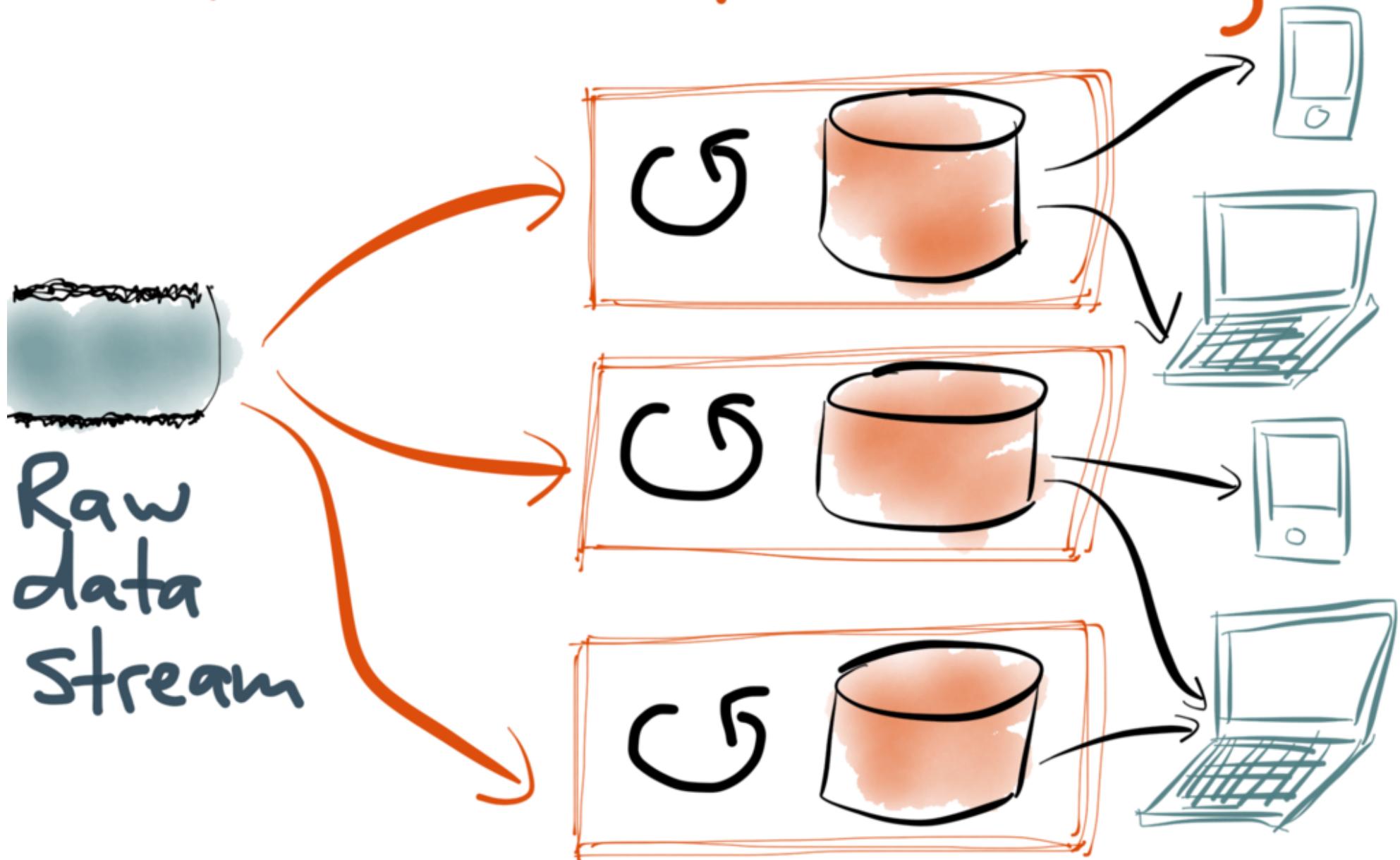
Video mem



pixels

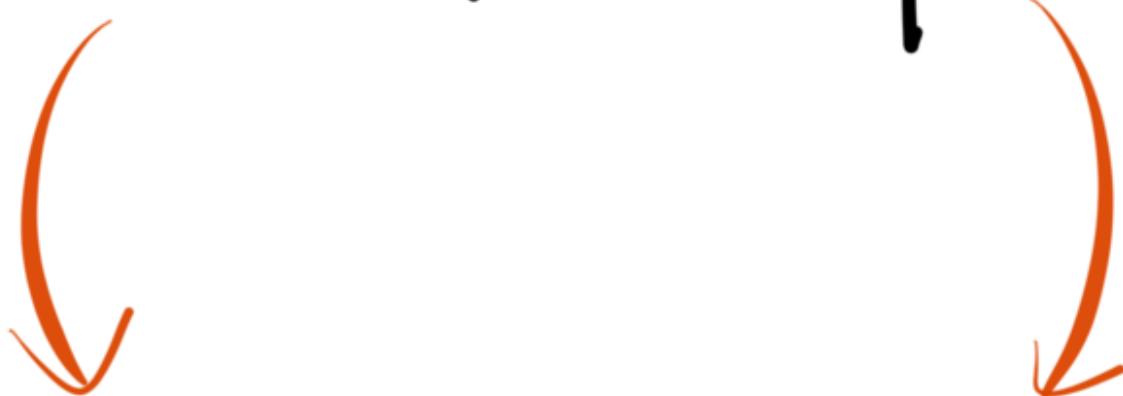


# Clients subscribe to MV changes



Request / response

~~Request / response~~



Subscribe / notify

(e.g. Meteor, Firebase — need much more)

# Streams

everywhere!

Including updates to materialized views.

# References / further reading

- Martin Kleppmann: “Rethinking caching in web apps.” 1 October 2012. <http://martin.kleppmann.com/2012/10/01/rethinking-caching-in-web-apps.html>
- Martin Kleppmann: “Designing data-intensive applications.” O’Reilly, to appear in 2015. <http://dataintensive.net/>
- Jay Kreps: “The Log: What every software engineer should know about real-time data’s unifying abstraction.” 16 December 2013. <http://engineering.linkedin.com/distributed-systems/log-what-every-software-engineer-should-know-about-real-time-datas-unifying>
- Jay Kreps: “Questioning the Lambda Architecture.” 2 July 2014. <http://radar.oreilly.com/2014/07/questioning-the-lambda-architecture.html>
- Jay Kreps: “Why local state is a fundamental primitive in stream processing.” 31 July 2014. <http://radar.oreilly.com/2014/07/why-local-state-is-a-fundamental-primitive-in-stream-processing.html>
- Nathan Marz and James Warren: “Big Data: Principles and best practices of scalable realtime data systems.” Manning MEAP, to appear January 2015. <http://manning.com/marz/>
- Apache Samza documentation. <http://samza.incubator.apache.org/>
- Alexandros Labrinidis, Qiong Luo, Jie Xu, and Wenwei Xue: “Caching and Materialization for Web Databases,” *Foundations and Trends in Databases*, volume 2, number 3, pages 169–266, March 2010.
- Stefano Ceri, Georg Gottlob, and Letizia Tanca: “What You Always Wanted to Know About Datalog (And Never Dared to Ask),” *IEEE Transactions on Knowledge and Data Engineering*, volume 1, number 1, pages 146–166, March 1989.

*dataintensive.net*

O'REILLY®

# Designing Data-Intensive Applications

THE BIG IDEAS BEHIND RELIABLE, SCALABLE,  
AND MAINTAINABLE SYSTEMS



Martin Kleppmann

@martinkl

Samza.incubator.  
apache.org

Samza