

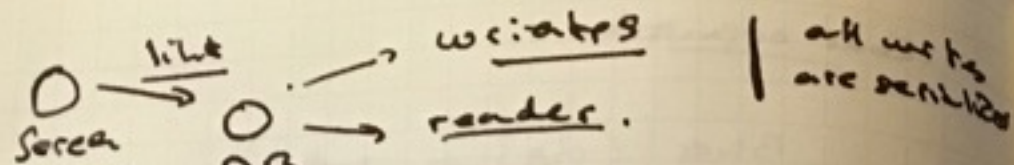


BUILD A P2P DOCUMENT ORIENTED DATABASE

BARREL

<https://barrel-db.org>

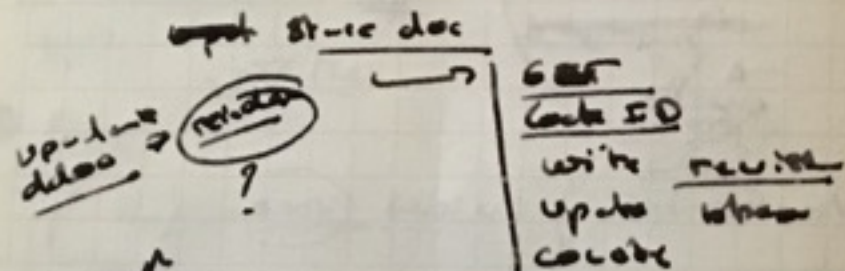




name:
name use &
Path:
type (doc or word)

open (DB New.)
DB resource
or create.

constructor.
read_header
write_header?



view DB.

create a
p-trin du doc

view of
log doc ID.
view-by-obj.

DB
view
NULL

DOC DB

Local
Doc ID
Doc.

REQUIREMENTS

- writer/reader
- read are
concurrent
are low - cost
- compact.

use DB ...
more - u-file
creation + DB
associate 2 1
path as can
1 - plus view
ajuta de view doc
of collabor

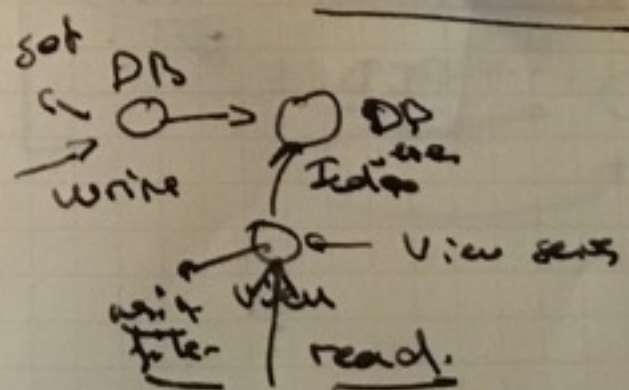
write doc.

Read ops / DB
update file.

for all readers with
diff.
write.

write

span. block
write



open -> DB

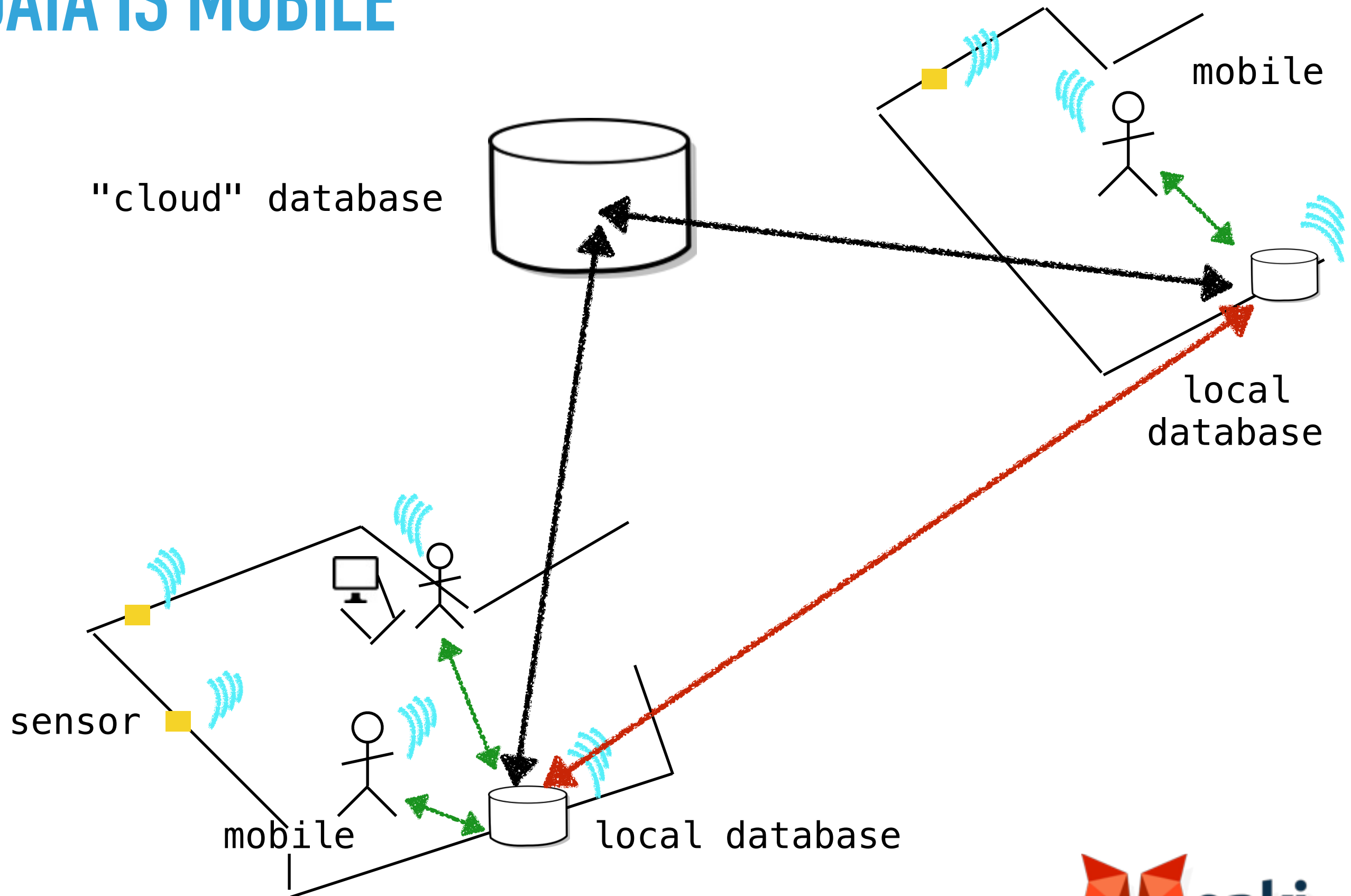




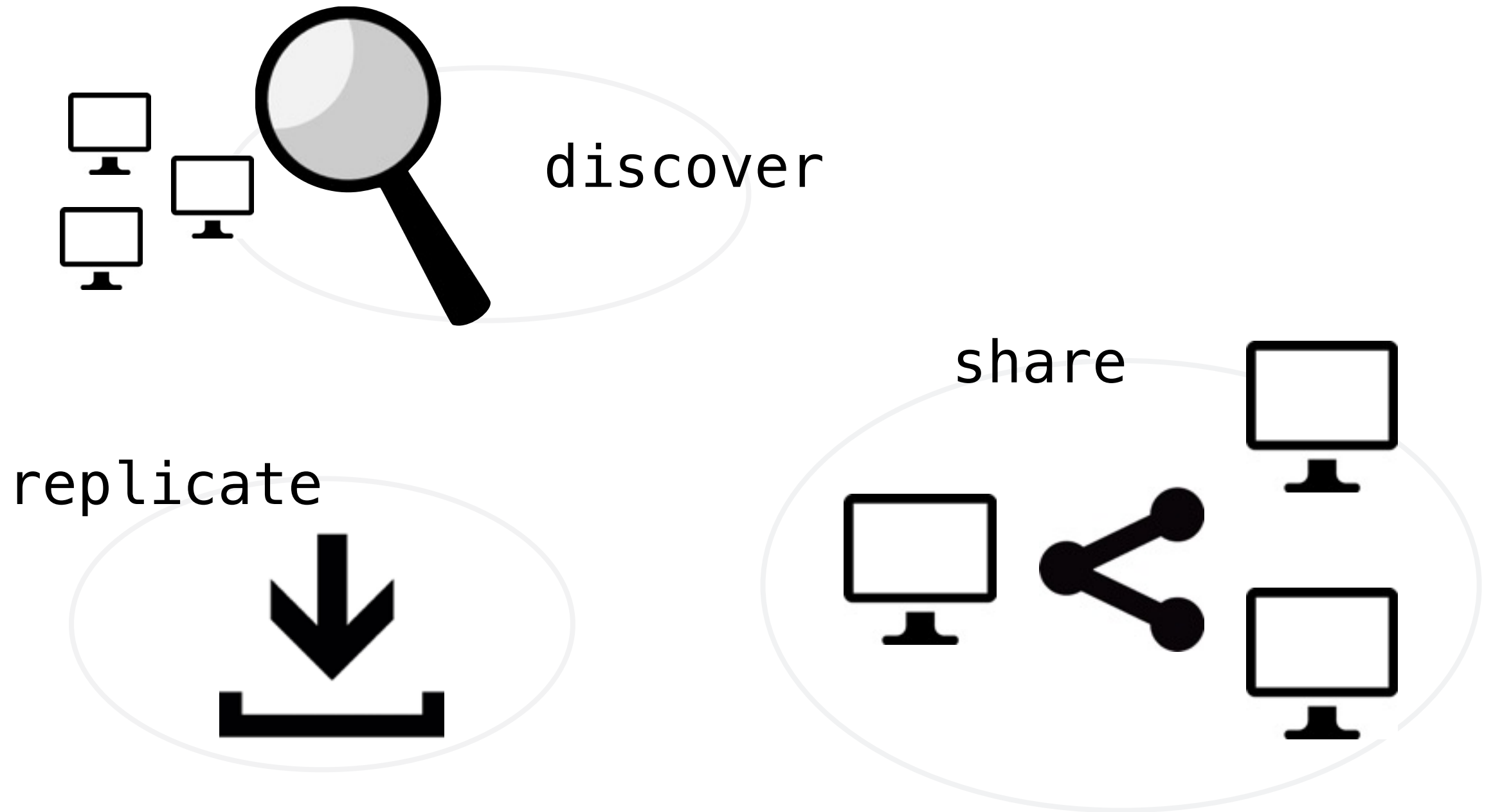


VISION AND CONCEPT

DATA IS MOBILE



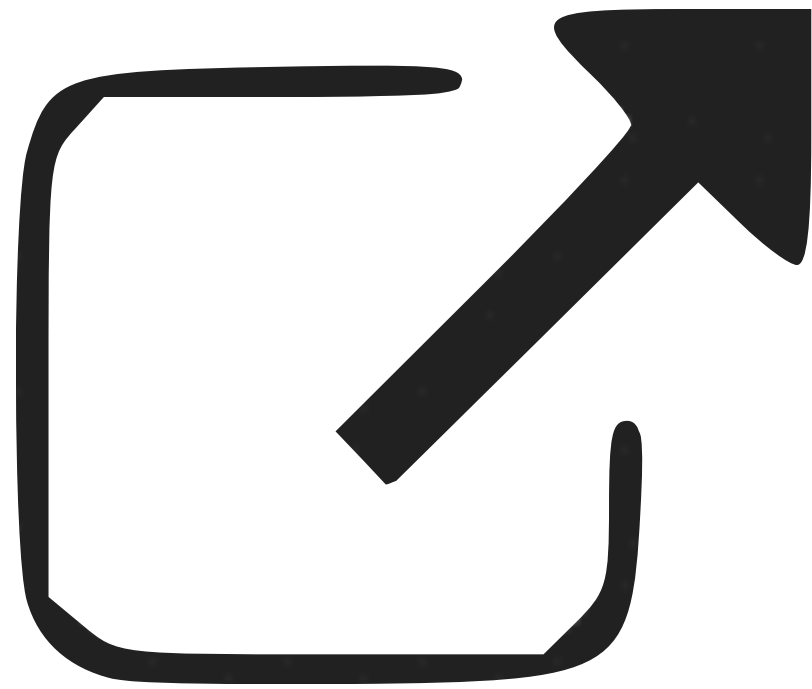
PEER TO PEER (P2P)



- ▶ Local first
- ▶ Put/Match the data next to you
- ▶ Query Locally
- ▶ Replicate a view of the data you need



WHAT IS BARREL



- ▶ a document database
- ▶ documents are JSON with attachments and links
- ▶ changes feed for documents and indexes
- ▶ replication between any nodes in both ways
- ▶ views (~ map)
- ▶ HTTP 1.1/2 API

WHAT IS BARREL

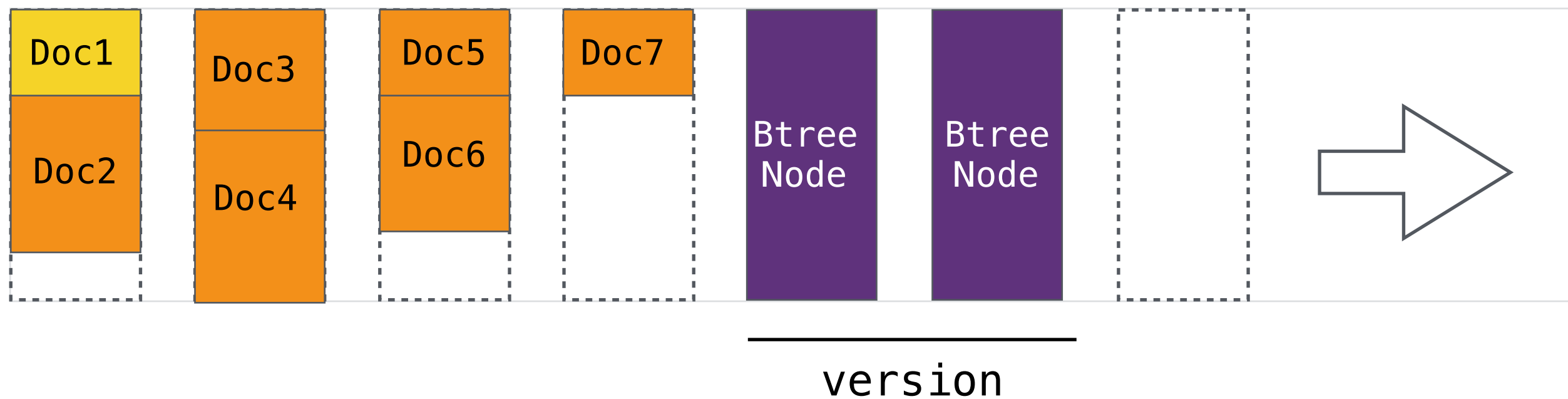
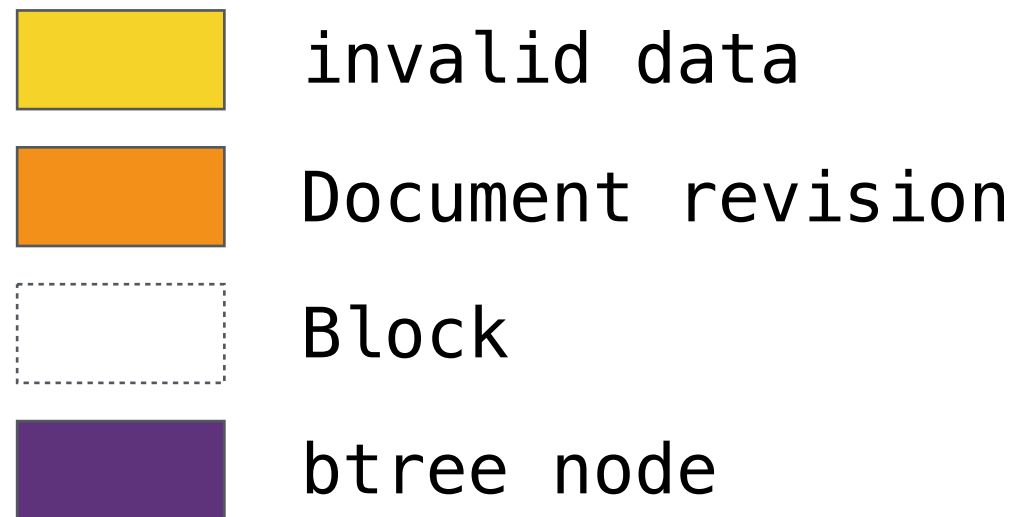


- ▶ DATA: not just blobs
- ▶ Replicated APPs
- ▶ Couchapps but extended and revisited

REPLICATED APPS



DECONSTRUCT



APPEND ONLY & MVCC



- ▶ Create a new file to remove the fragmentation
- ▶ A race between copy and the addition of new data
- ▶ Require at least twice of the storage

THE COMPACTION ISSUE



ID-Index

ID 1	METADATA 1
ID 2	METADATA 2
ID 3	METADATA 3

Seq-Index

SEQ 1	METADATA 1
SEQ 2	METADATA 2
SEQ 3	METADATA 3

DB file



Indexed document

DOCUMENT STORAGE



- ▶ 2 indexes (btree): by sequence, by id
- ▶ 1 index for local documents without conflict handling
- ▶ A revision tree is stored in indexes pointed to the revision offset
- ▶ The revision is stored in the file separately

HOW ARE STORED DOCUMENTS



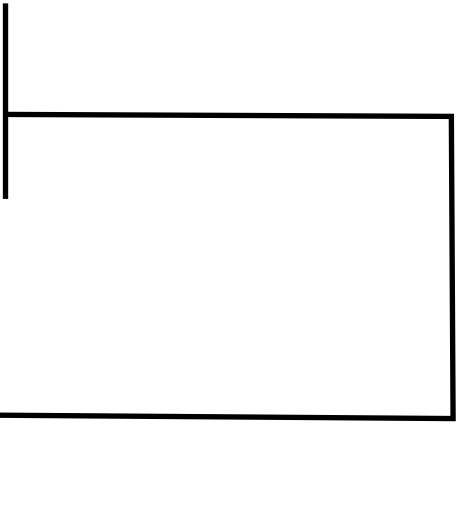
- ▶ Reverse index (map)
- ▶ Index using a function
- ▶ Function in javascripts, erlang, ..
- ▶ Incremental index
- ▶ Retrieves changes (aka view changes)
- ▶ View are regrouped by groups (1 db file/group)

VIEWS



Log-Index

DOCID	View 1	KEY 1	SEQ 1	ADD
		KEY 2	SEQ 2	DEL
	View 2	KEY 1	SEQ 1	ADD



Key-Index

[KEY 1, DOCID]	[VALUE, DOCREV, SEQ]
[KEY 2, DOCID]	[del, DOCREV, SEQ]
[KEY 3, DOCID 2]	[VALUE, DOCREV, SEQ]

SEQ-Index

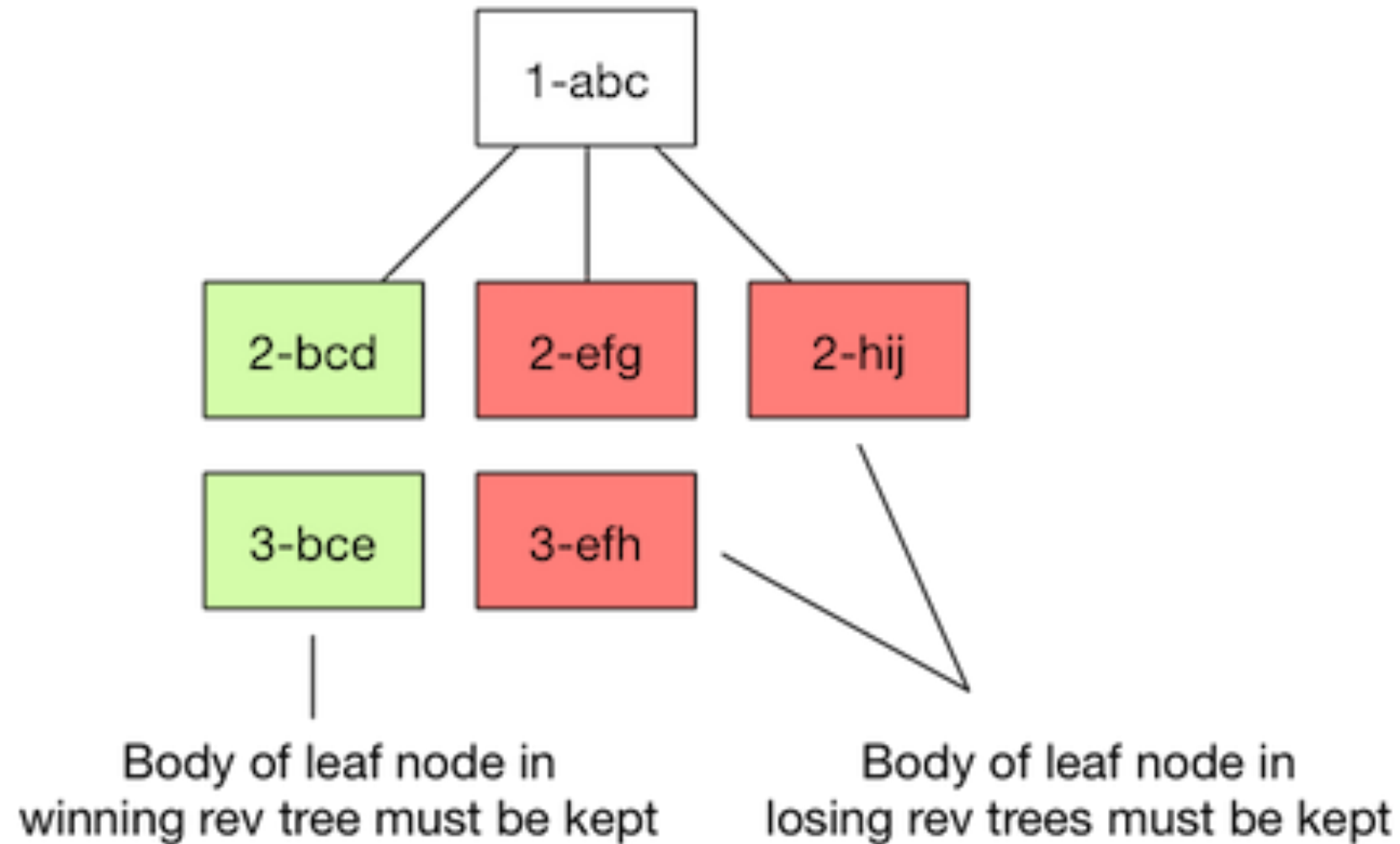
[SEQ 1, KEY]	[VALUE, DOCREV, SEQ]
[SEQ 2, KEY 2]	[del, DOCREV, SEQ]
[SEQ 3, KEY]3	[VALUE, DOCREV, SEQ]

view

VIEW STORAGE



REVISION TREE



BUILT



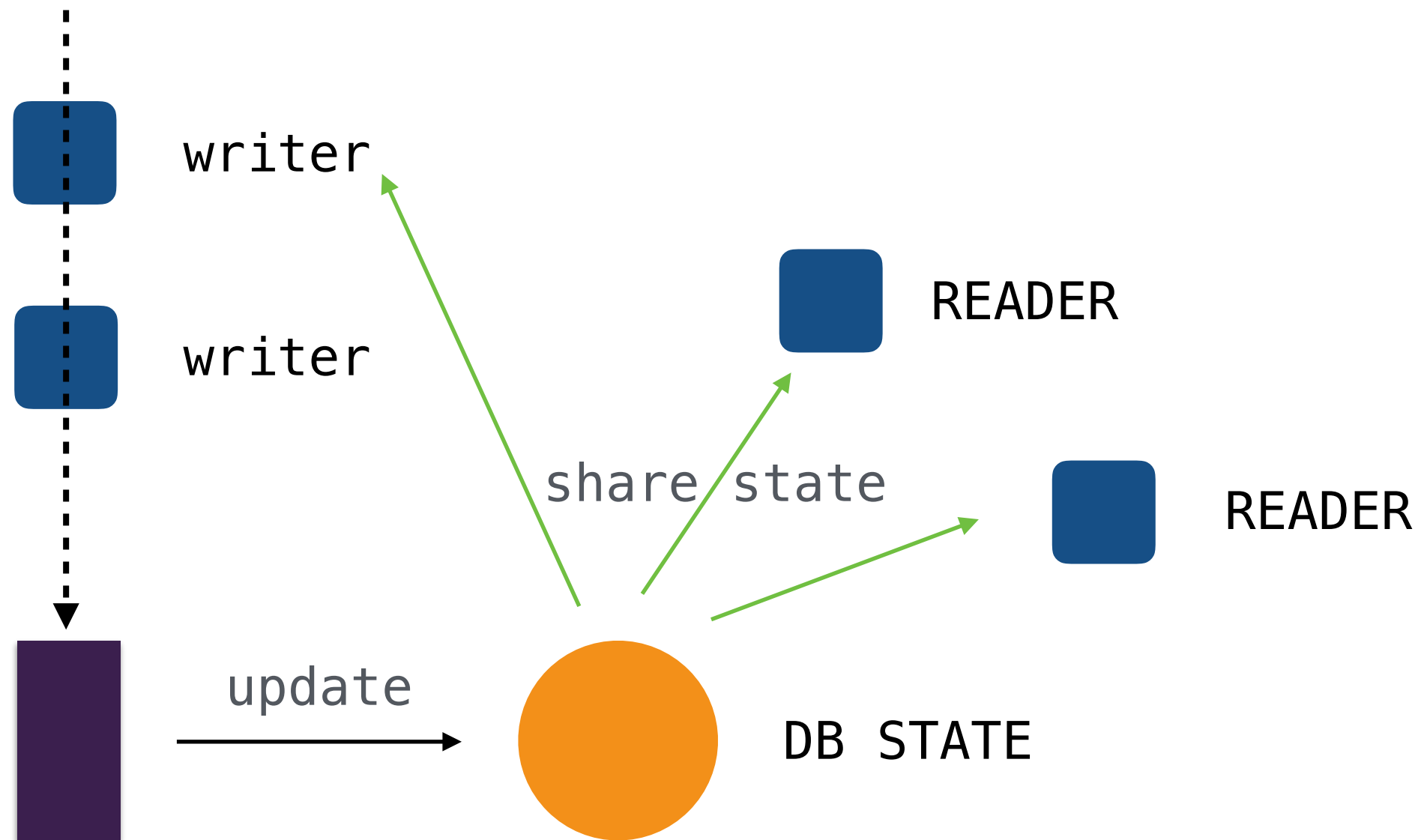
IN ERLANG



- ▶ Write is slow
- ▶ Read should not being blocked by writes
- ▶ No shared memory
- ▶ No atomic integer trick
- ▶ Only actors and message passing
- ▶ Operations on a doc are atomic

CHALLENGES





READ/WRITE OPERATIONS

- ▶ LRU to cache blocks
<https://github.com/barrel-db/erlang-lru>
- ▶ 1 File process, Operations are limited
- ▶ DB users are linked to the database process
- ▶ Optional Write buffer to reduce the latency
- ▶ Optional wal

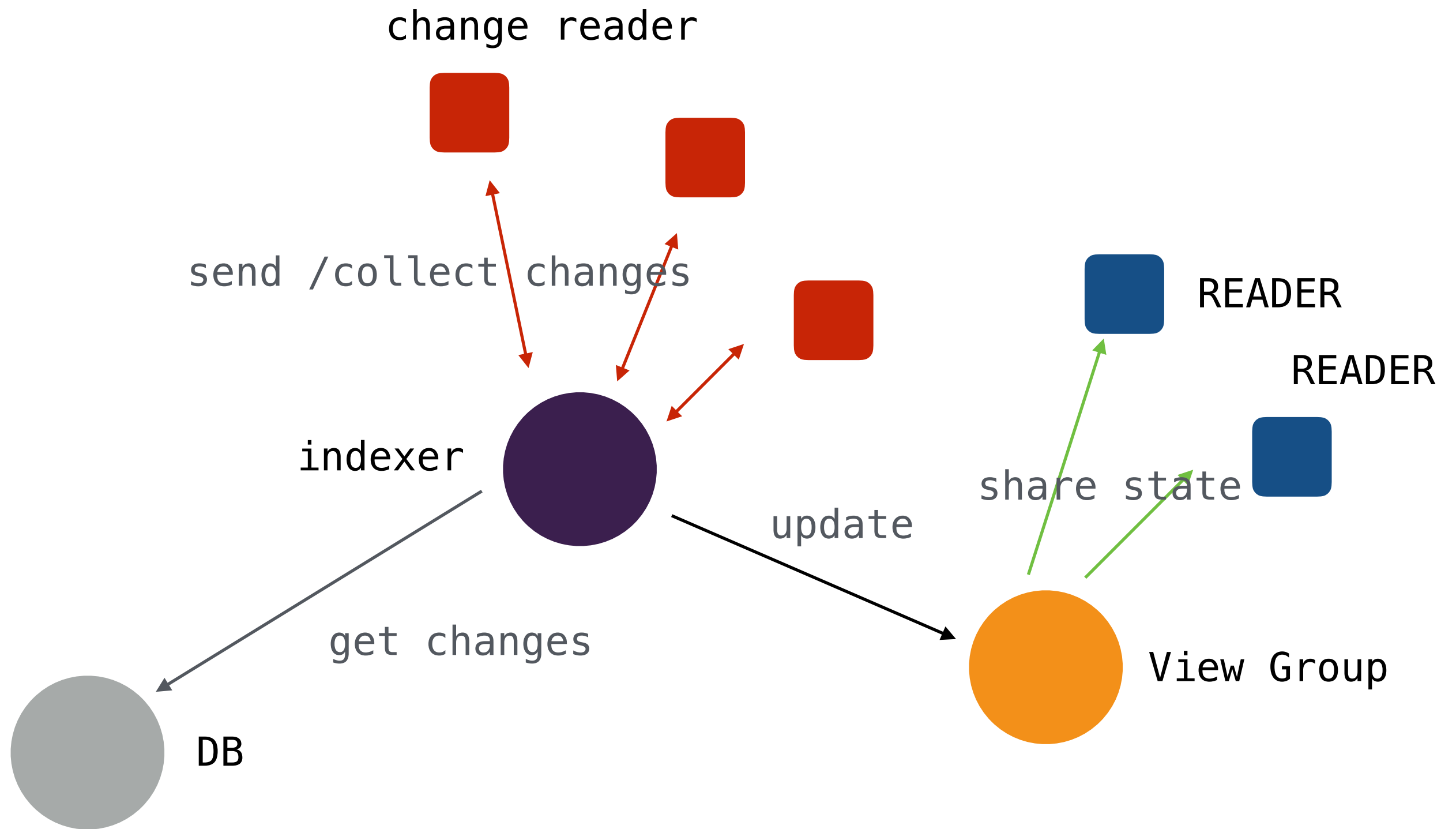
READ/WRITE OPERATIONS



- ▶ STORE SEGMENTS of data for compaction
- ▶ IO is "relatively" slow in erlang
- ▶ USE a "native KV store" as a nif.

SPEEDUPS





INDEX OPERATIONS



- ▶ Credit Flow Based
- ▶ The View group keep the state
- ▶ View group is created on demande
- ▶ kept open until it has readers
- ▶ Indexer ask for updates
- ▶ Read functions (Map functions) are processed in //

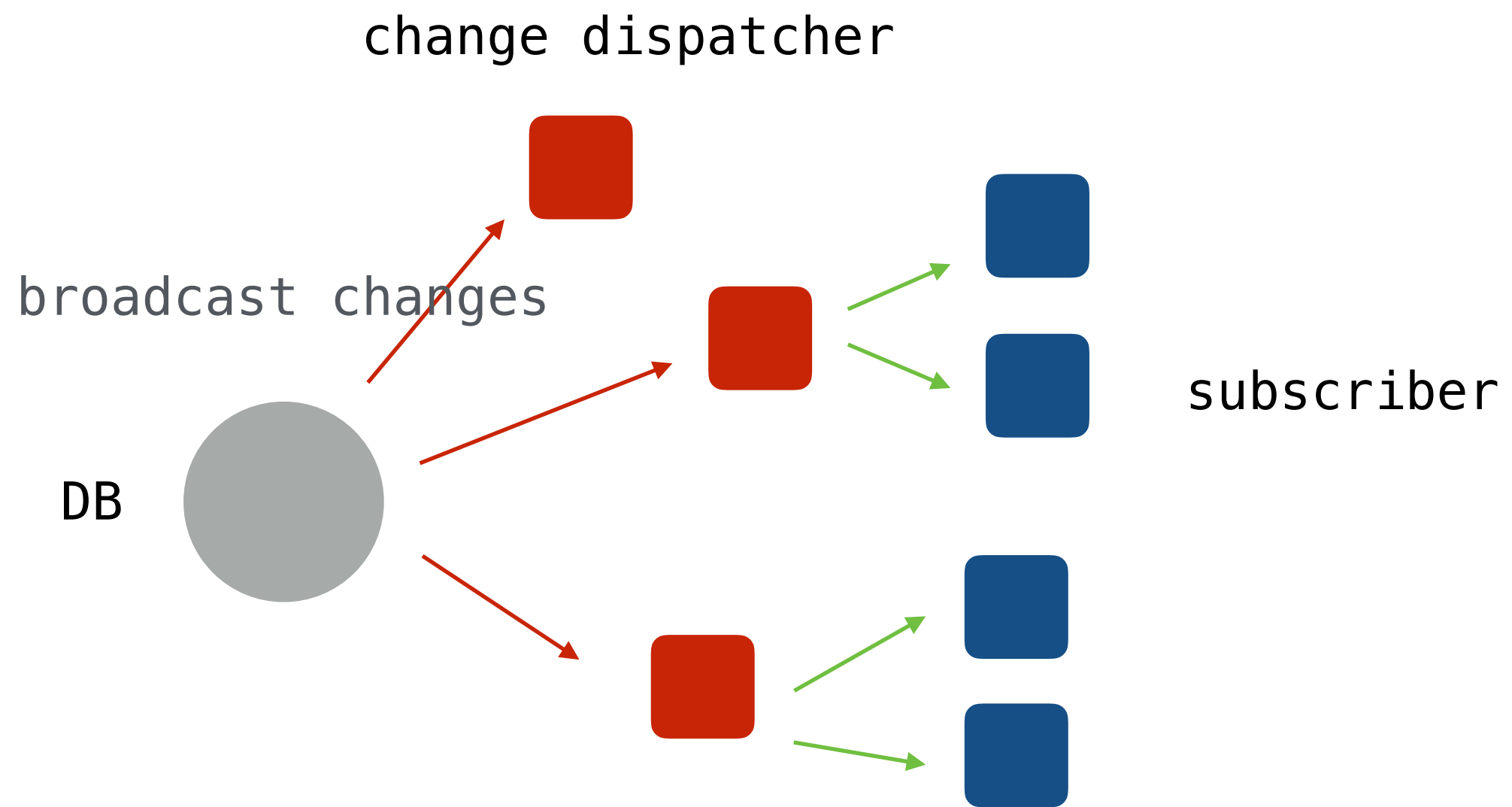
INDEX OPERATIONS



- ▶ Added 2 features:
 - ▶ MOVE: move doc(s) to another node or database (like copy but with delete)
 - ▶ User hooks functions (run in background) using hooks:
<https://github.com/barrel-db/hooks>
- ▶ Partition on demand
- ▶ Decision depends on the application needs

NEW FUNCTIONS





CHANGES HANDLER



- ▶ Use the sequence index
- ▶ changes load balancing
- ▶ consumer subscribe on patterns (delete, update, ...)
- ▶ Create changes Load Balancer on demand
- ▶ Allows remote nodes to subscribe to a queue
- ▶ Based on primer (release on March 2016)

CHANGES EVENTS



- ▶ Use the sequence index
- ▶ changes load balancing
- ▶ consumer subscribe on patterns (delete, update, ...)
- ▶ Create changes Load Balancer on demand
- ▶ Allows remote nodes to subscribe to a queue
- ▶ Based on primer (release on March 2016)

CHANGES EVENTS

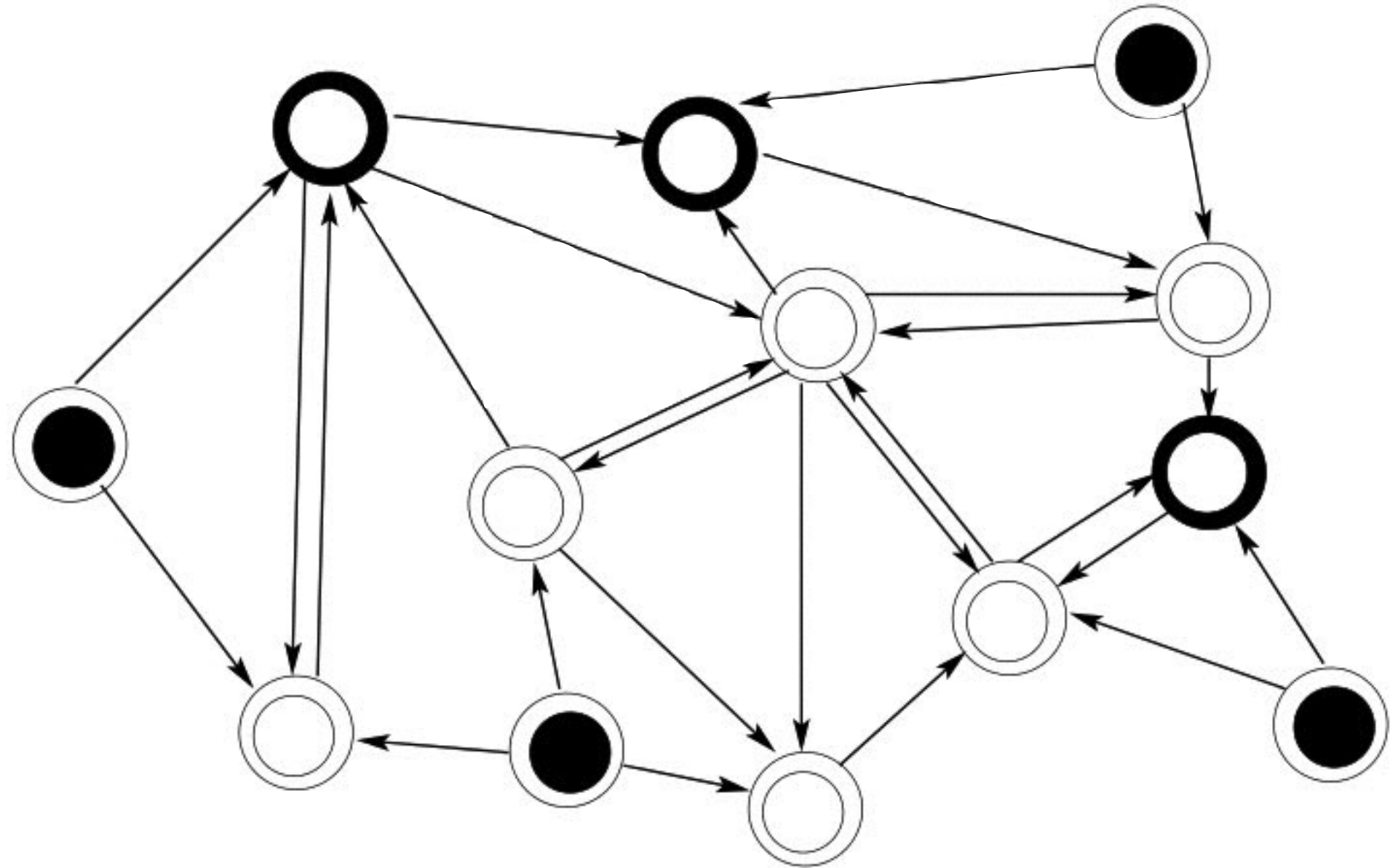


- ▶ inherited the HTTP api in mochiweb
- ▶ small changes to makes the server more resilient
- ▶ chatterbox
- ▶ wip in cowboy.
- ▶ yaws ?

HTTP API



P2P



- ▶ Over HTTP
- ▶ Replication is the core
- ▶ Each nodes can replicate each others
- ▶ PUSH/PULL
- ▶ Chained replication

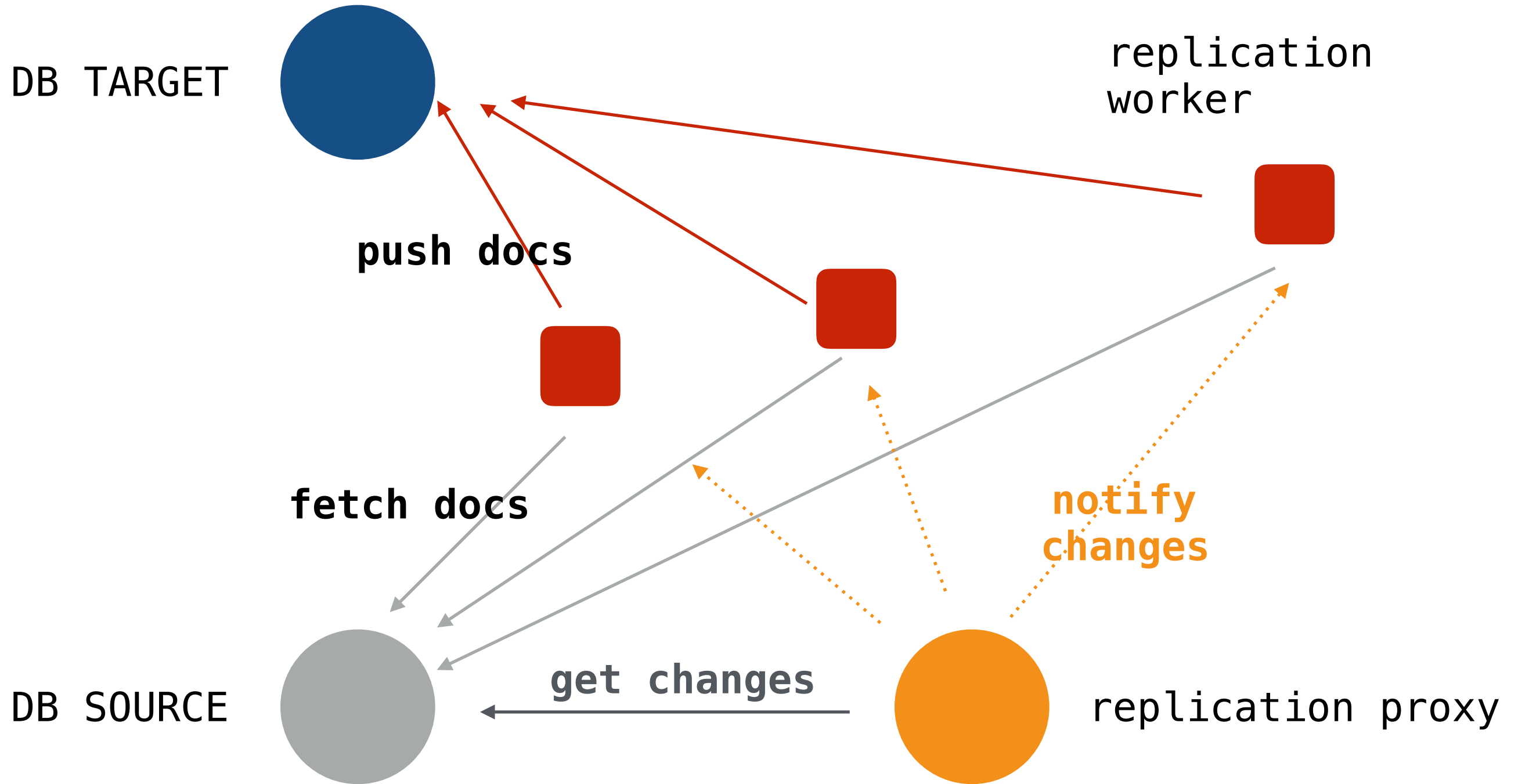
P2P



- ▶ Based on the change feed
- ▶ fetch the revisions and their attachments not present on the node
- ▶ continuous or not
- ▶ try to collect multiple docs at once
- ▶ use hackney:
<http://github.com/benoitc/hackney>
- ▶ Use a Flow-based pattern instead of a classic pool

REPLICATION





REPLICATION OPERATIONS



- ▶ Replication state is stored at least on one node
- ▶ checkpoints
- ▶ get the revisions not actually stored on the nodes ("_rev_diffs")
- ▶ the replication proxy maintains routes
- ▶ build replication chains, by replicating status

REPLICATION





Barrel

[HTTPS://BARREL-DB.ORG](https://barrel-db.org)

Enki Multimedia

[HTTP://ENKIM.EU](http://enkim.eu)

