



# Deconstructing the Database

Rich Hickey

Most programs are outside the bounds of and single process we write  
in any single FP language

# What is Datomic?

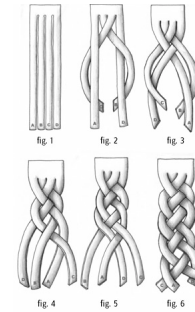
- A new database
- A sound model of **information**, with time
- Provides **database as a value** to applications
- Bring **declarative programming** to applications
- Focus on reducing complexity

Not going to do a full rationalization or overview

Focus on information important for Datomic, not necessarily for all use of dbs/stores

# DB Complexity

- Stateful
- Same query, different results
  - no basis
- Over there
- 'Update' poorly defined
  - Places



Inherent complexity in state

# Update

- What does update mean?
- Does the new replace the old?
- Granularity? new \_\_\_\_ replace the old \_\_\_\_
- Visibility?

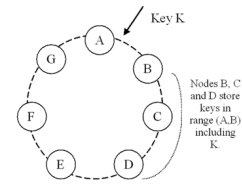
Granularity and replacement the big issues

# Manifestations

- Wrong programs
- Scaling problems
- Round-trip fears
- Fear of overloading server
- Coupling, e.g. questions with reporting

read committed vs repeatable, serializable

# Consistency and Scale



- What's possible?
- Distributed redundancy and consistency?
- Elasticity
- Inconsistency huge source of complexity

dynamo and bigtable

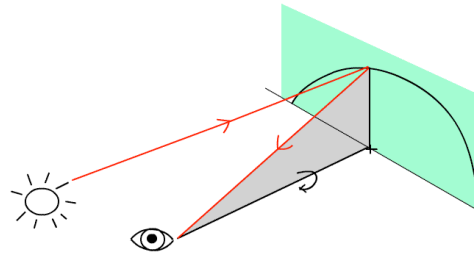
# Information and Time

- Old-school memory and records
- The kind you remember  
... and keep
- Auditing and more



# Perception and Reaction

- No polling
- Consistent





# Coming to Terms

## Value

- An *immutable* magnitude, quantity, number... or immutable composite thereof

## Identity

- A putative entity we associate with a series of causally related values (states) over time

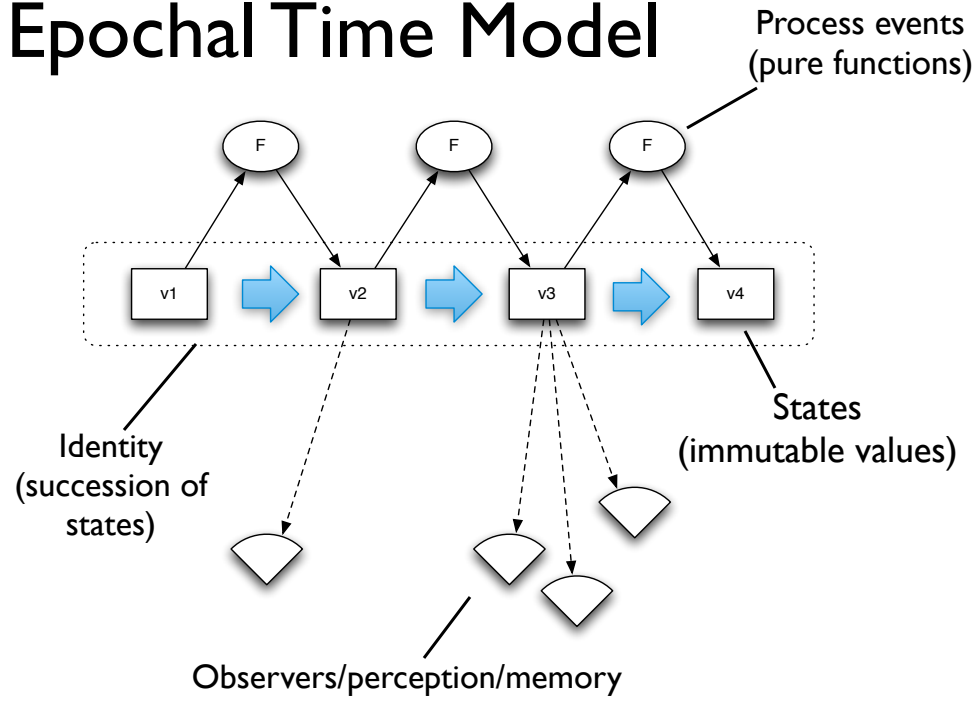
## State

- Value of an identity at a moment in time

## Time

- Relative before/after ordering of causal values

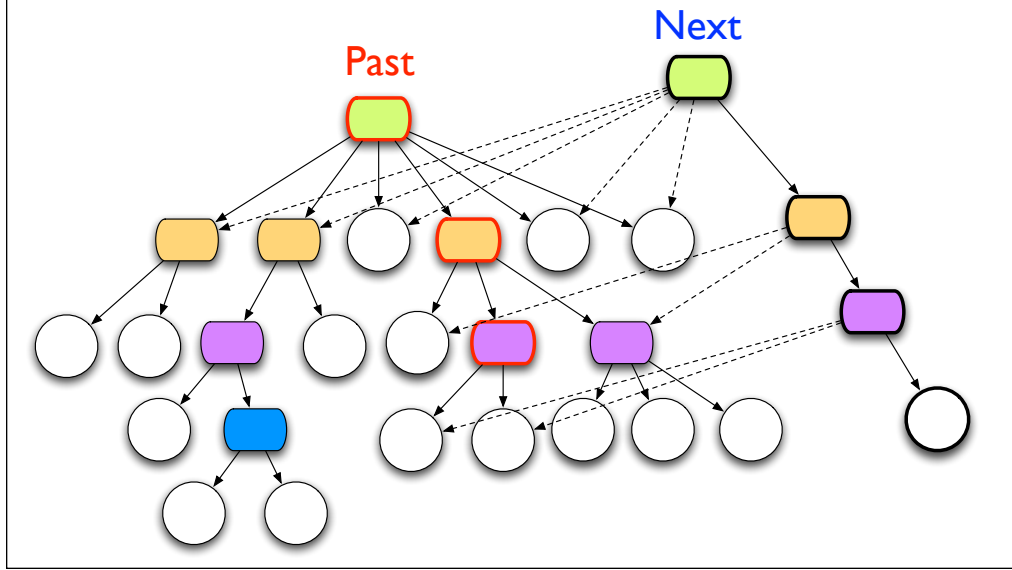
# Epochal Time Model

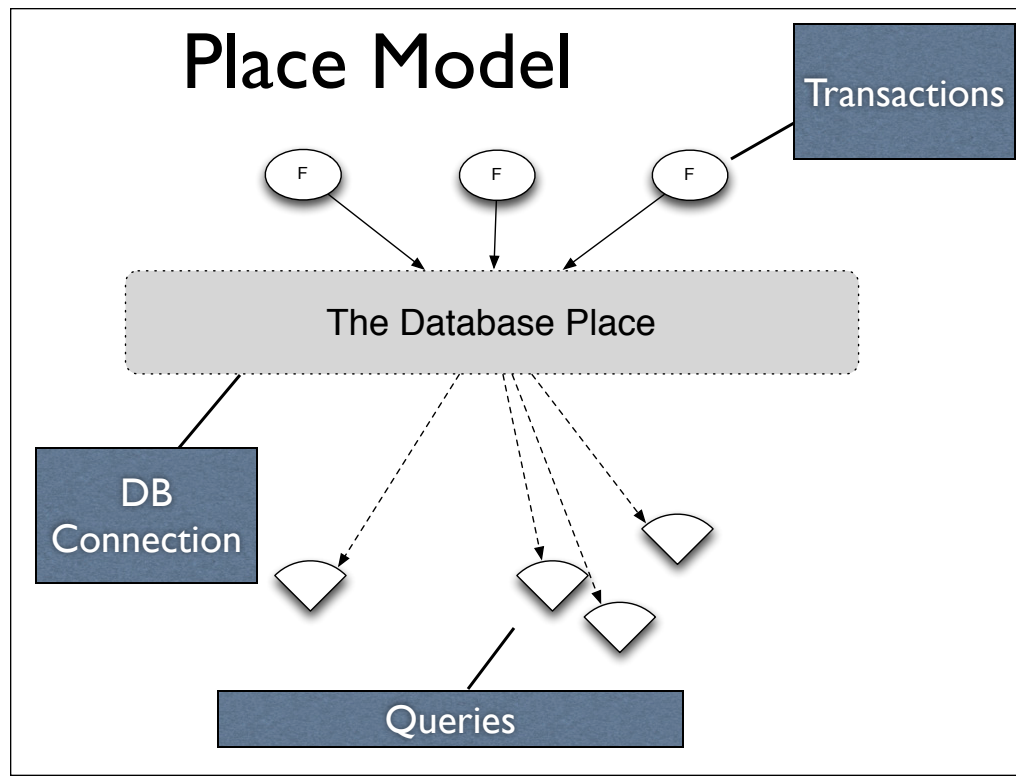


# Implementing Values

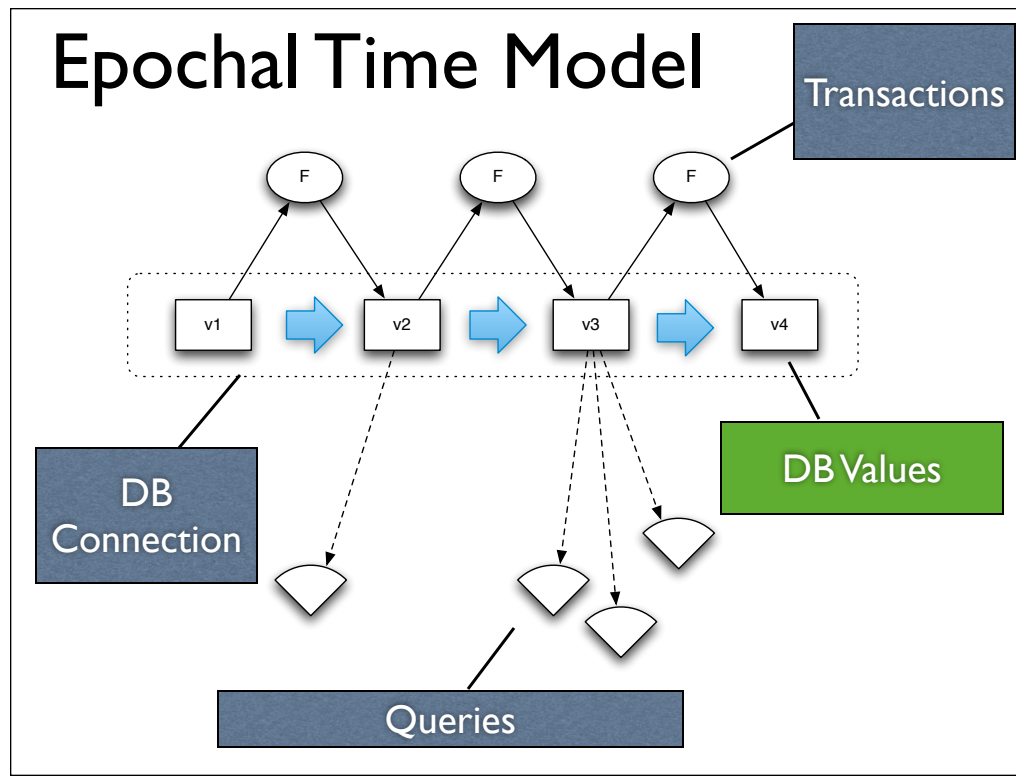
- Persistent data structures
- Trees
- Structural sharing

# Structural Sharing



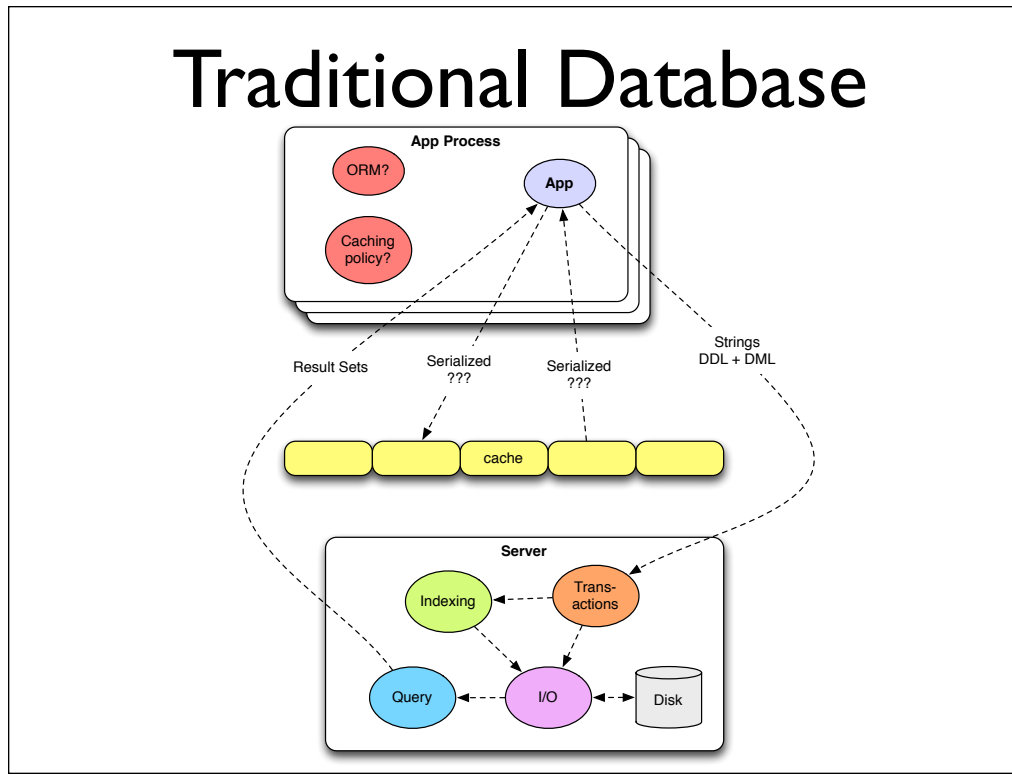


We should recognize – same problems as OO  
Conflates identity and value, collapses time



This is what we want – transactions are functions of db values queries as well

# Traditional Database



what's makes up a database?

cache over DB access, disk locality was important

# The Choices

- Coordination
  - how much, and where?
  - process requires it
  - perception shouldn't
- Immutability
  - sine qua non

coordinate up front or later

immutability advantages with eventually consistent systems



# Approach

- Move to information model
- Split process and perception
- Immutable basis in storage
- Novelty in memory

# Information

- Inform
  - 'to convey knowledge via facts'
  - 'give shape to (the mind)'
- Information
  - the facts

# Facts

- **Fact** - 'an event or thing known to have happened or existed'
  - From: factum - 'something done'
  - Must include time
  - Remove structure (a la RDF)
  - Atomic **Datom**
    - Entity/Attribute/Value/Transaction(time)

don't get more flexibility by trading tables for documents  
factum – “something done”

# Database State

- The database as an expanding **value**
- An accretion of **facts**
- The past doesn't change - immutable
- Process requires new space
- Fundamental move away from **places**

What is the state – snapshot of a referential model?

# Accretion

- Root per transaction doesn't work
- Latest values include past as well
  - The past is sub-range
- Important for information model

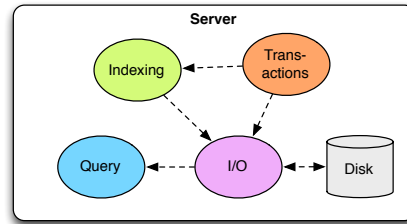
can we just do shared structure from before on disk?

# Process

- Reified
- Primitive representation of novelty
  - Assertions and retractions of **facts**
  - **Minimal**
- Other transformations expand into those

Important to accretion that novelty representation is minimal

# Deconstruction



- Process

- Transactions
- Indexing
- ○

- Perception/Reaction

- Query
- Indexes
- I

# State

- Must be organized to support query
- Sorted set of facts
- Maintaining sort live in storage - bad
  - BigTable - mem + storage merge
  - occasional merge into storage
  - persistent trees

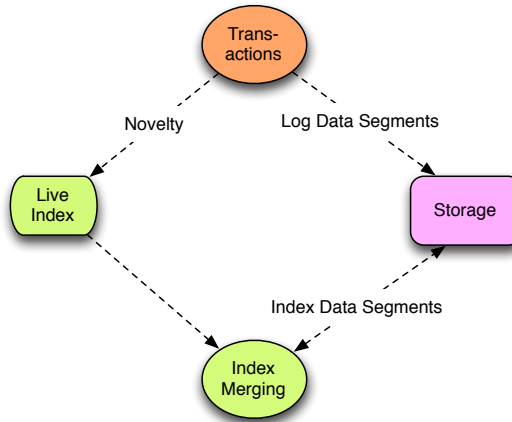
Databases are about leverage



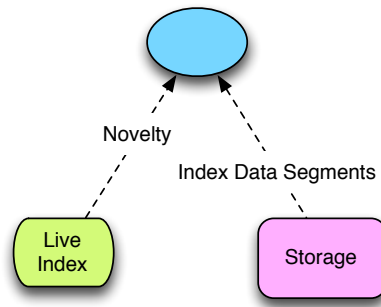
# Indexing

- Maintaining sort live in storage - bad
- BigTable et al:
  - Accumulate novelty in memory
  - Current view: mem + storage merge
  - Occasional integrate mem into storage  
Releases memory

# Transactions and Indexing



# Perception



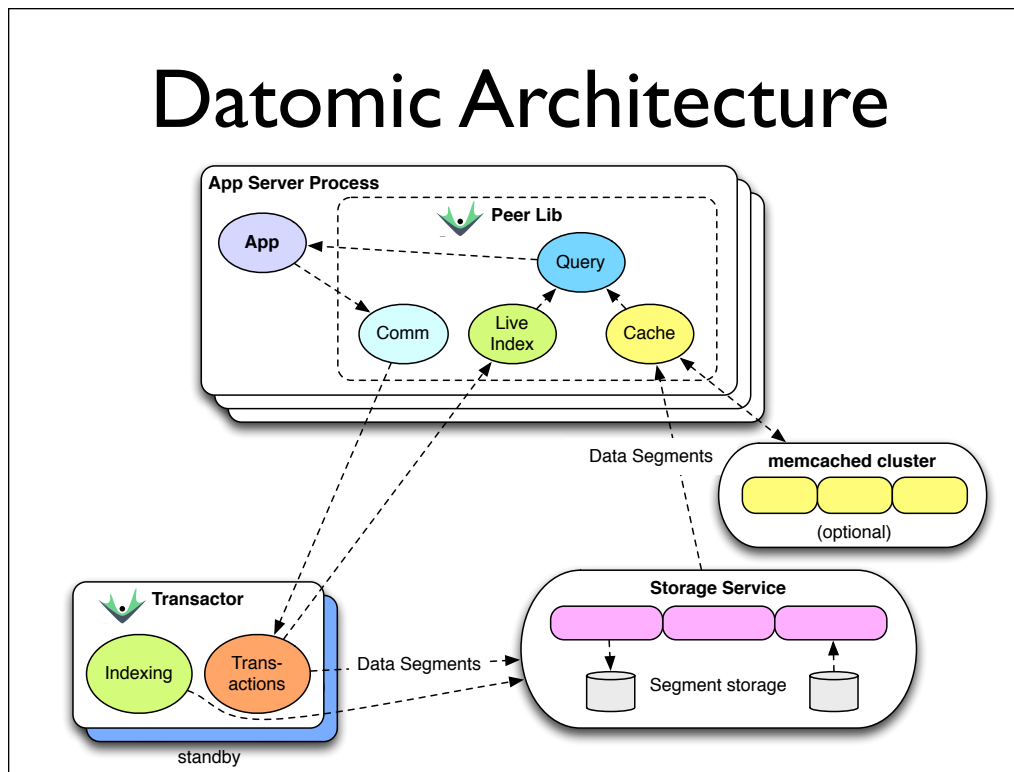
Just a merge join. Any coordination required? Contention?  
How does live index get updated?  
Live vs periodic now separate decision

# Components

- Transactor
- Peers
  - Your app servers, analytics machines etc
- Redundant storage service

Physical architecture

# Datomic Architecture



But, storage now remote, slow? No – cache everywhere

# Transactor

- Accepts transactions
  - Expands, applies, logs, broadcasts
- Periodic indexing, in background
- Indexing creates garbage
  - Storage GC

# Peer Servers

- Peers directly access storage service
- Have own query engine
- Have live mem index and merging
- Two-tier cache
  - Datoms w/object values (on heap)
  - Segments (memcached)

# Consistency and Scale

- Process/writes go through transactor
  - traditional server scaling/availability
- Immutability supports consistent reads
  - without transactions
- Query scales with peers
  - Elastic/dynamic e.g. auto-scaling

consistency as in db satisfies consistency predicate  
loopholes 7 and 10?



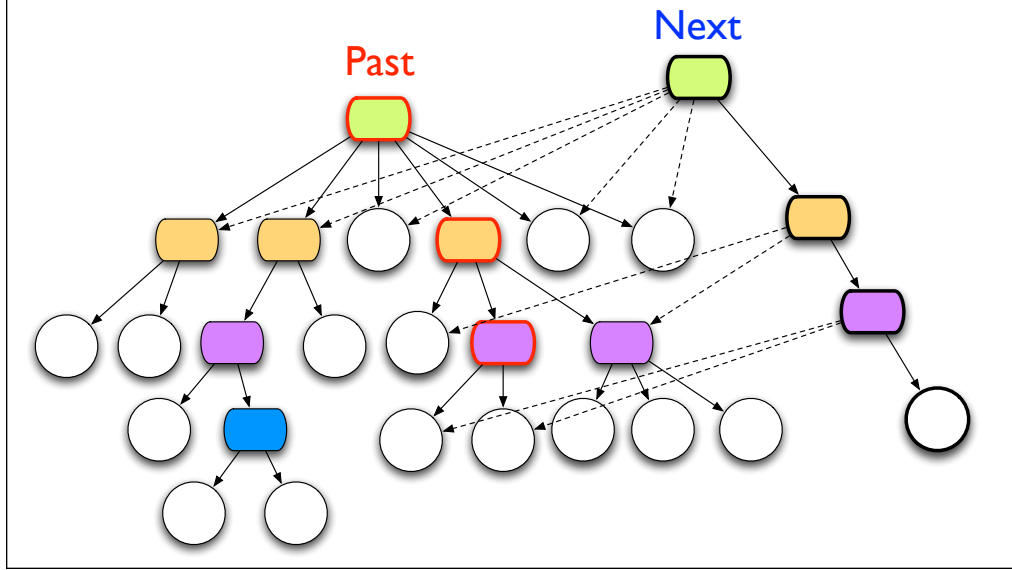
# Memory Index

- Persistent sorted set
- Large internal nodes
- Pluggable comparators
- 2 sorts always maintained
  - EAVT,AEVT
- plus AVET,VAET

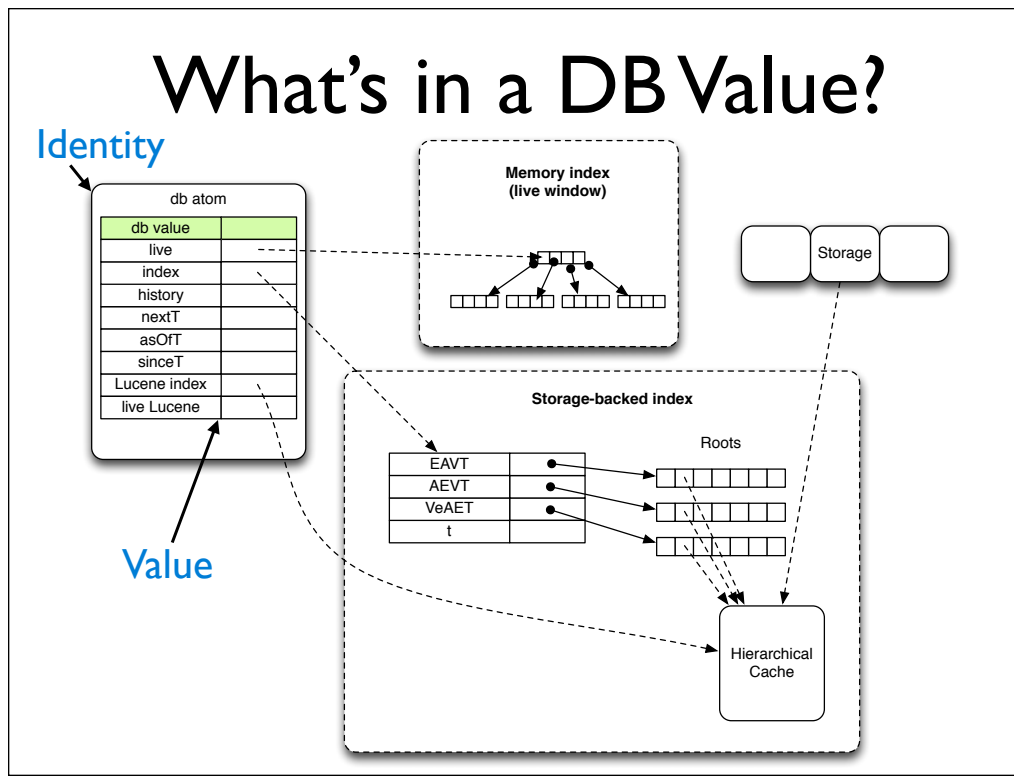
# Storage

- Log of tx asserts/retracts (in tree)
- Various covering indexes (trees)
- Storage service/server requirements
  - Data segment values (K->V)
  - atoms (consistent read)
  - pods (conditional put)

# Structural Sharing



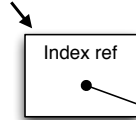
# What's in a DB Value?



Value can be lazily loaded, since source immutable

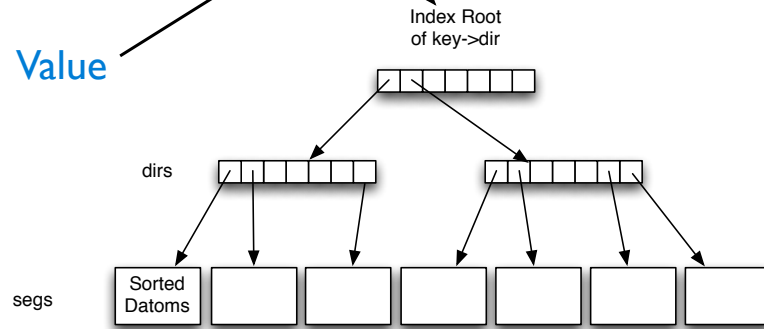
# Index in Storage

Identity



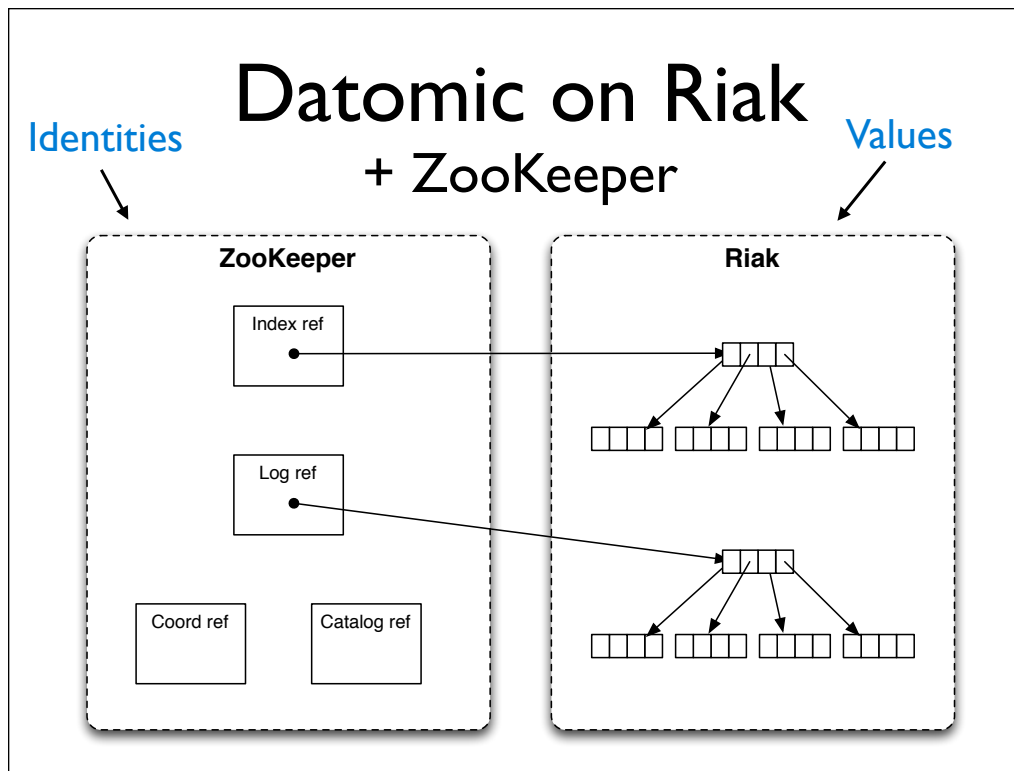
T	EAVT	AEVT	VeAET	AVET	Lucene
42					

Value



# Datomic on Riak + ZooKeeper

- Riak
  - redundant, distributed, highly available
  - durable
  - eventually consistent
- ZooKeeper
  - redundant, durable,
  - consistent (ordered ops + CAS)



pointer swap mentioned by Eric Brewer this morning

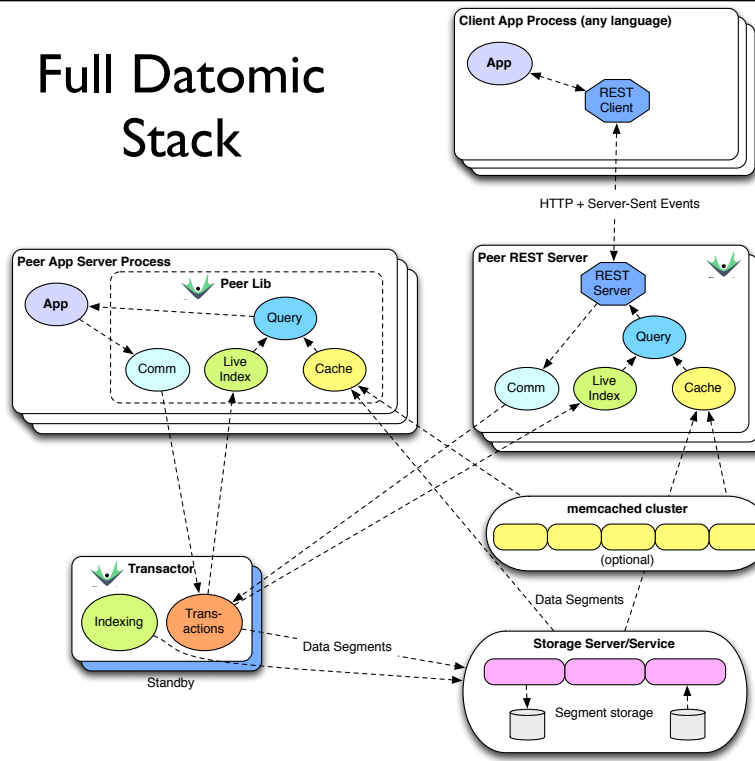
# Riak Usage

- Everything put into Riak is immutable
- $N=3, W=2, DW=2$
- $R=1$ , not-found-ok = false  
‘first found’ semantics
- There or not  
no vector clocks, siblings etc
- No speculative lookup

What notion of consistency?  
Application-level predicate



# Full Datomic Stack



# Stable Bases

```
//Peer  
Database db = connection.db().asOf(1000);  
Peer.q(aQuery, db);
```

```
//Client  
GET /data/mem/test/1000/datoms?index=aevt
```

basis



- Same query, same results
- db permalinks!
  - communicable, recoverable
- Multiple conversations about same value

Value of values

# DB Values

- Time travel
  - `db.asOf` - past
  - `db.since` - windowed
  - `db.with(tx)` - speculative
- dbs are arguments to query, not implicit
- mock with datom-shaped data:

```
[[:fred :likes "Pizza"]  
[:sally :likes "Ice cream"]]
```

# DB Simplicity Benefits

- Epochal state
  - Coordination only for process
- Transactions well defined
  - Functional accretion
- Freedom to relocate/scale storage, query
- Extensive caching
- Process events

# The Database as a Value

- Dramatically less complex
- More powerful
- More scalable
- Better information model



**Thanks for Listening!**