

Monday, November 26, 12



# The CouchDB Implementation

**Jan Lehnardt**

**jan@a.o**

**@janl**

**CouchDB PMC Chair**

**Committer #2**

Monday, November 26, 12

Thanks for the invite

Glad to be here

JSConf EU / BBuzz / JSFAB

# Any Database

Monday, November 26, 12

- FS integration / raw storage
- core data structures
- core features
- API

**API**

**CORE FEATURES**

**CORE DATA STRUCTURES**

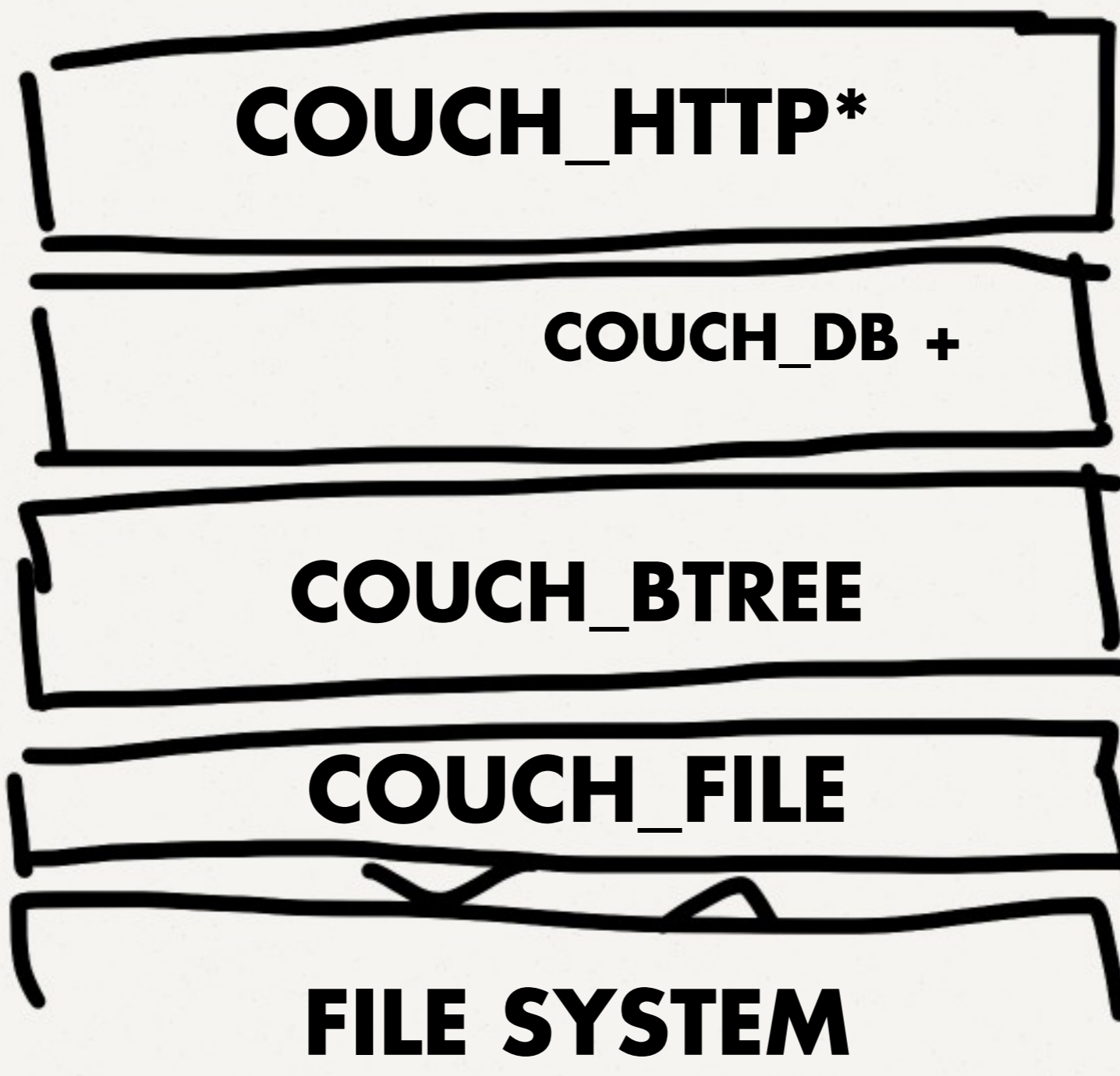
**FILE SYSTEM ACCESS**

**FILE SYSTEM**

# CouchDB is no different

Monday, November 26, 12

- couch\_file
- couch\_btree
- couch\_db / couch\_doc / couch\_mr / couch\_replicator / etcpp
- couch\_httpd\*



**COUCH\_DOC**  
**COUCH\_MR**  
**COUCH\_REPL...**

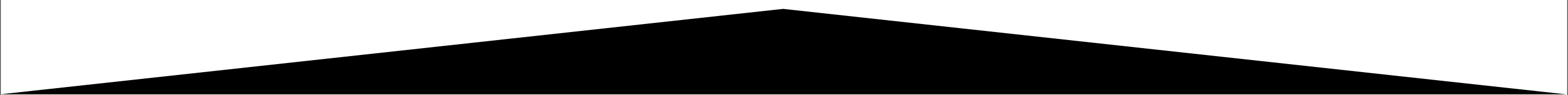
# Core Datastructures

Monday, November 26, 12

- 2 x b+tree & data interleaved
- append only, mvcc
- full fsync control
  
- Can answer:
  - Data for \$key
  - What happened \$since
  
- Used for core data storage
- As well as indexes
  
- Everything else is built on top



# Behold the b+tree



Monday, November 26, 12

Monday, November 26, 12

**by-id**

<b>A</b>													
<b>B</b>													
<b>C</b>													
<b>D</b>													
<b>E</b>													
<b>F</b>													
<b>G</b>													
<b>H</b>													
<b>I</b>													
<b>J</b>													
<b>K</b>													
<b>L</b>													
<b>...</b>													

<b>A</b>														
<b>B</b>														
<b>C</b>														
<b>D</b>														
<b>E</b>														
<b>F</b>														
<b>DOC_G</b>														
<b>H</b>														
<b>I</b>														
<b>J</b>														
<b>K</b>														
<b>L</b>														
<b>...</b>														

<b>A</b>														
<b>B</b>														
<b>C</b>														
<b>DOC_D</b>														
<b>E</b>														
<b>F</b>														
<b>DOC_G</b>														
<b>H</b>														
<b>I</b>														
<b>J</b>														
<b>K</b>														
<b>L</b>														
<b>...</b>														

<b>A</b>
<b>B</b>
<b>C</b>
<b>DOC_D</b>
<b>E</b>
<b>F</b>
<b>DOC_G</b>
<b>H</b>
<b>I</b>
<b>J</b>
<b>DOC_K</b>
<b>L</b>
<b>...</b>



**DOC\_D**

**DOC\_G**

**DOC\_K**

**by-req**  
**or**  
**“what happened  
since?”**

**1. DOC\_G**

**2. DOC\_D**

**3. DOC\_K**

# The CouchDB File Format

Monday, November 26, 12

- 2 x b+tree & data interleaved
- append only, mvcc
- full fsync control
  
- Can answer:
  - Data for \$key
  - What happened \$since
  
- Used for core data storage
- As well as indexes
  
- Everything else is built on top



Monday, November 26, 12

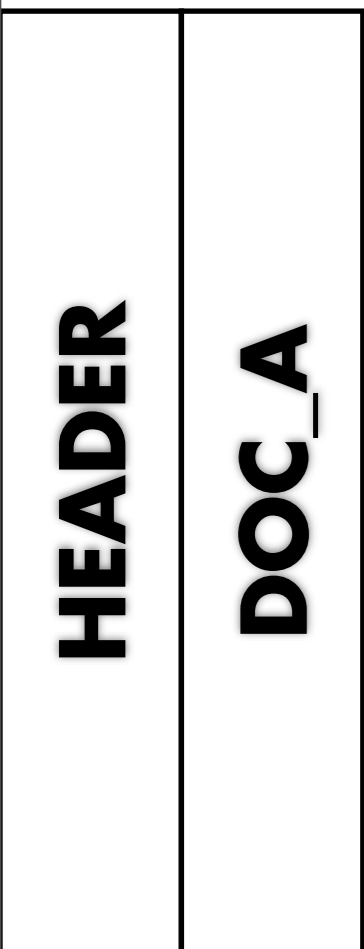
- 2 x b+tree & data interleaved
- append only, mvcc
- full fsync control
  
- Can answer:
  - Data for \$key
  - What happened \$since
  
- Used for core data storage
- As well as indexes
  
- Everything else is built on top

**HEADER**



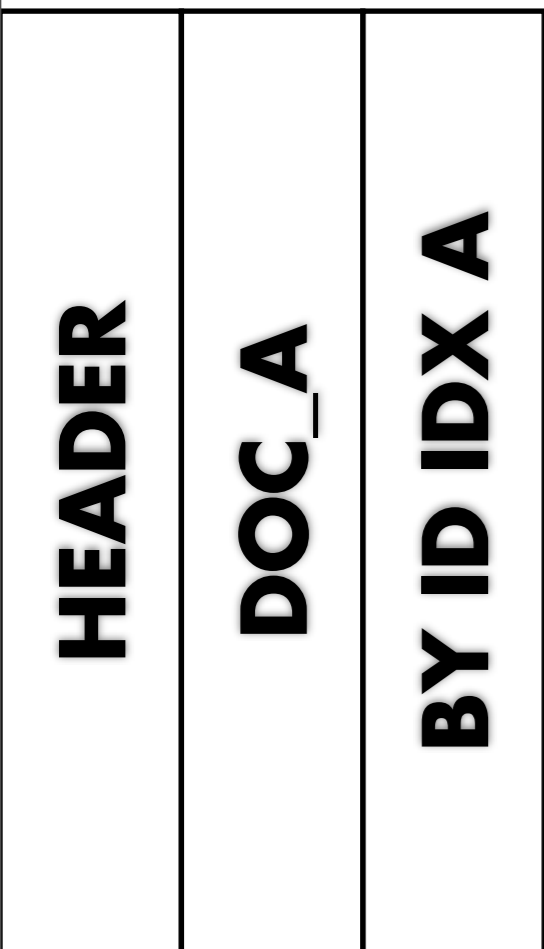
Monday, November 26, 12

- 2 x b+tree & data interleaved
- append only, mvcc
- full fsync control
  
- Can answer:
  - Data for \$key
  - What happened \$since
  
- Used for core data storage
- As well as indexes
  
- Everything else is built on top



Monday, November 26, 12

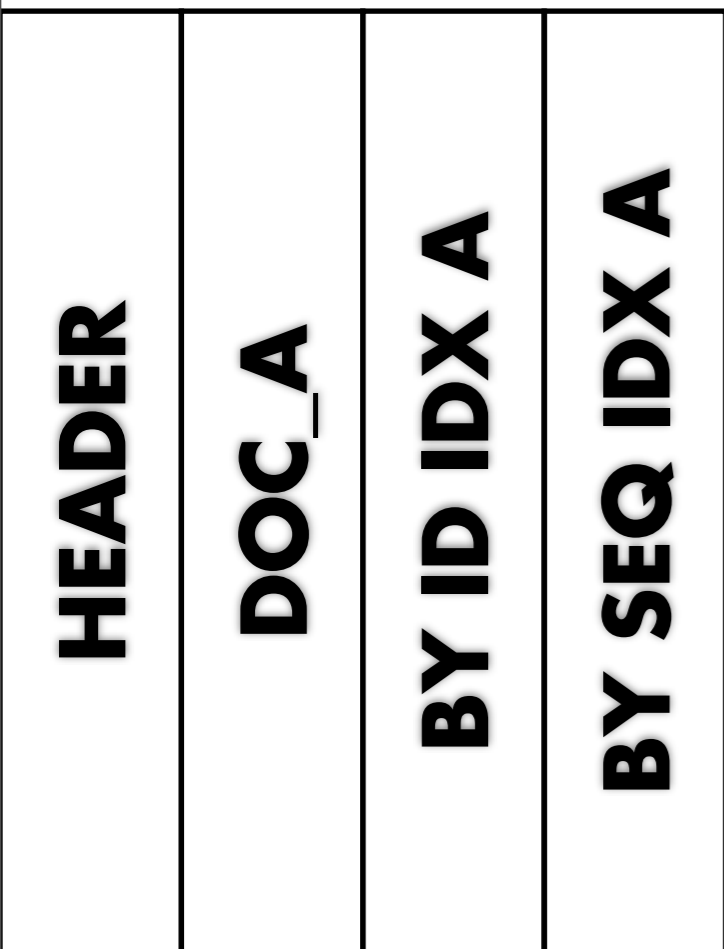
- 2 x b+tree & data interleaved
- append only, mvcc
- full fsync control
  
- Can answer:
  - Data for \$key
  - What happened \$since
  
- Used for core data storage
- As well as indexes
  
- Everything else is built on top



Monday, November 26, 12

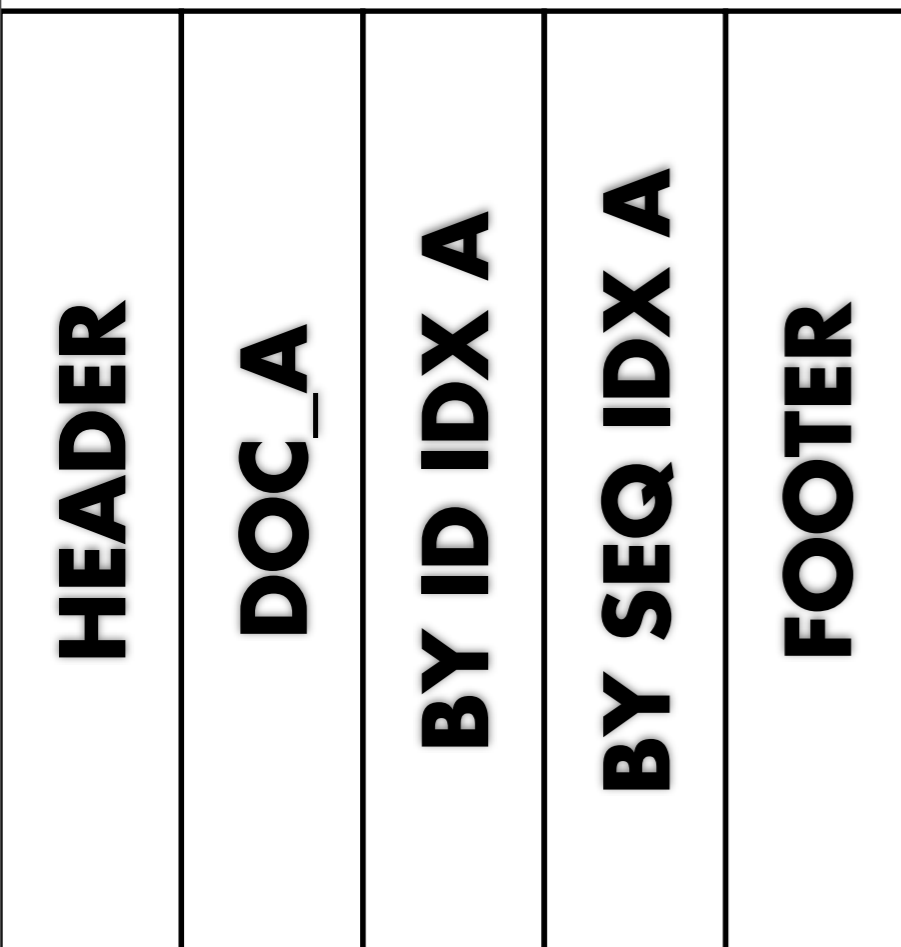
- 2 x b+tree & data interleaved
- append only, mvcc
- full fsync control
  
- Can answer:
  - Data for \$key
  - What happened \$since
  
- Used for core data storage
- As well as indexes
  
- Everything else is built on top





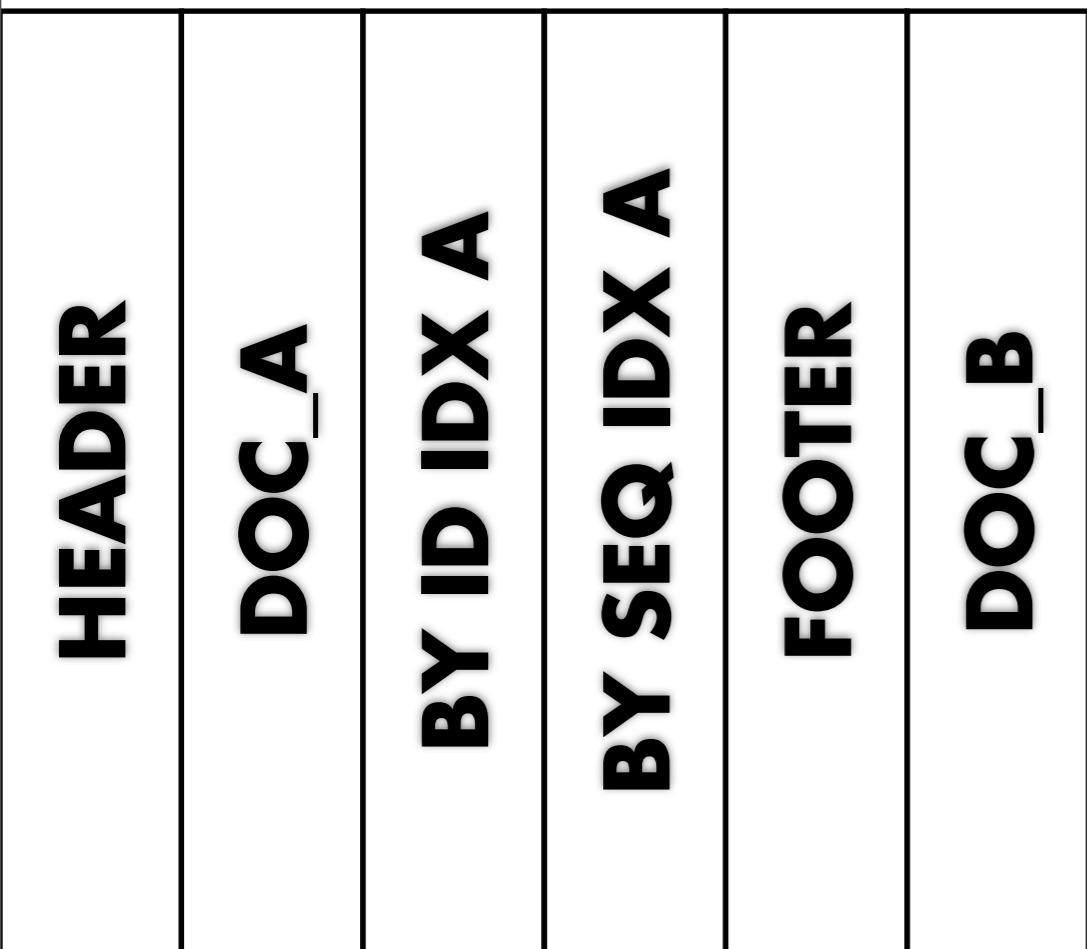
Monday, November 26, 12

- 2 x b+tree & data interleaved
- append only, mvcc
- full fsync control
  
- Can answer:
  - Data for \$key
  - What happened \$since
  
- Used for core data storage
- As well as indexes
  
- Everything else is built on top



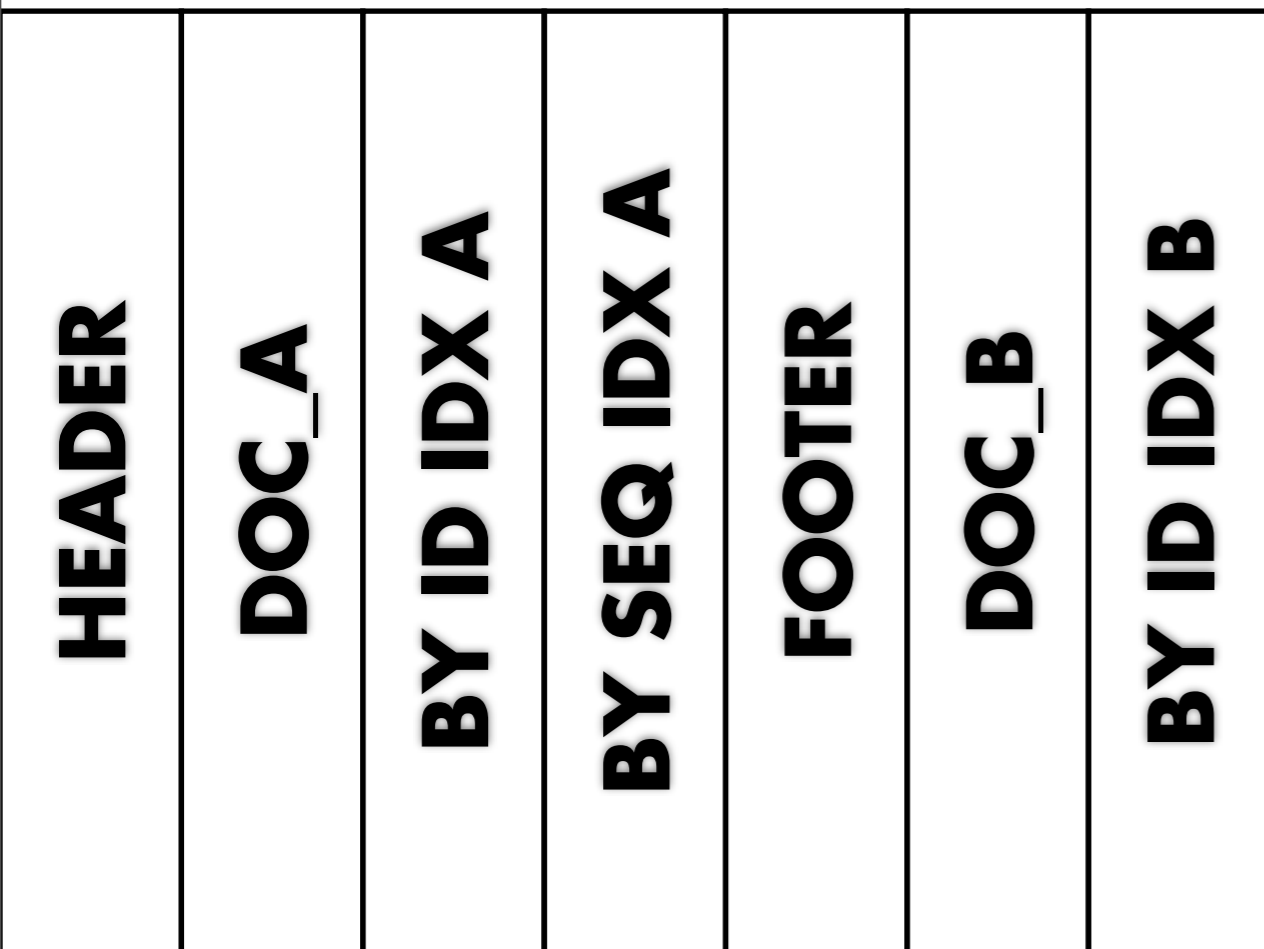
Monday, November 26, 12

- 2 x b+tree & data interleaved
- append only, mvcc
- full fsync control
  
- Can answer:
  - Data for \$key
  - What happened \$since
  
- Used for core data storage
- As well as indexes
  
- Everything else is built on top



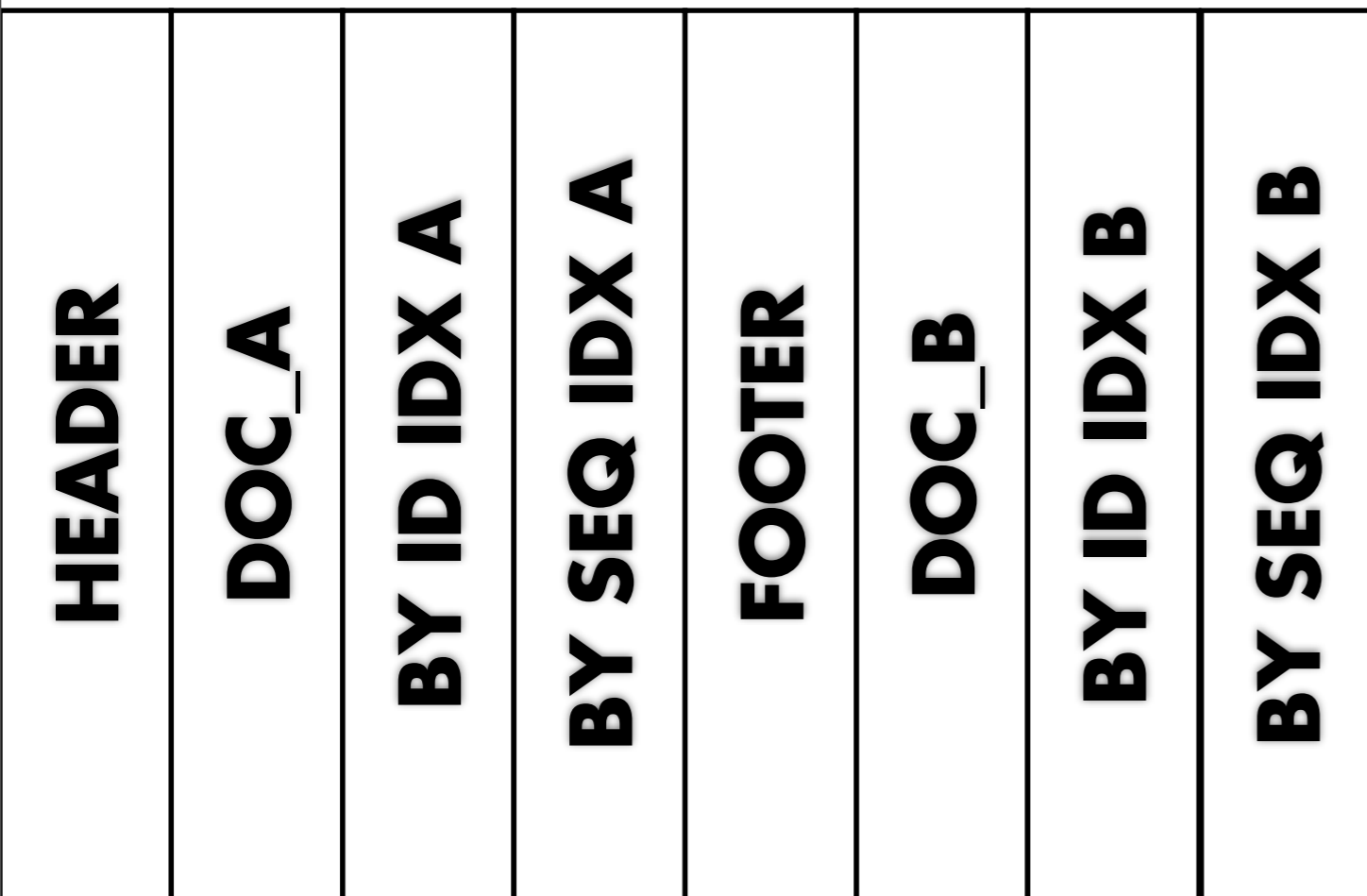
Monday, November 26, 12

- 2 x b+tree & data interleaved
- append only, mvcc
- full fsync control
  
- Can answer:
  - Data for \$key
  - What happened \$since
  
- Used for core data storage
- As well as indexes
  
- Everything else is built on top



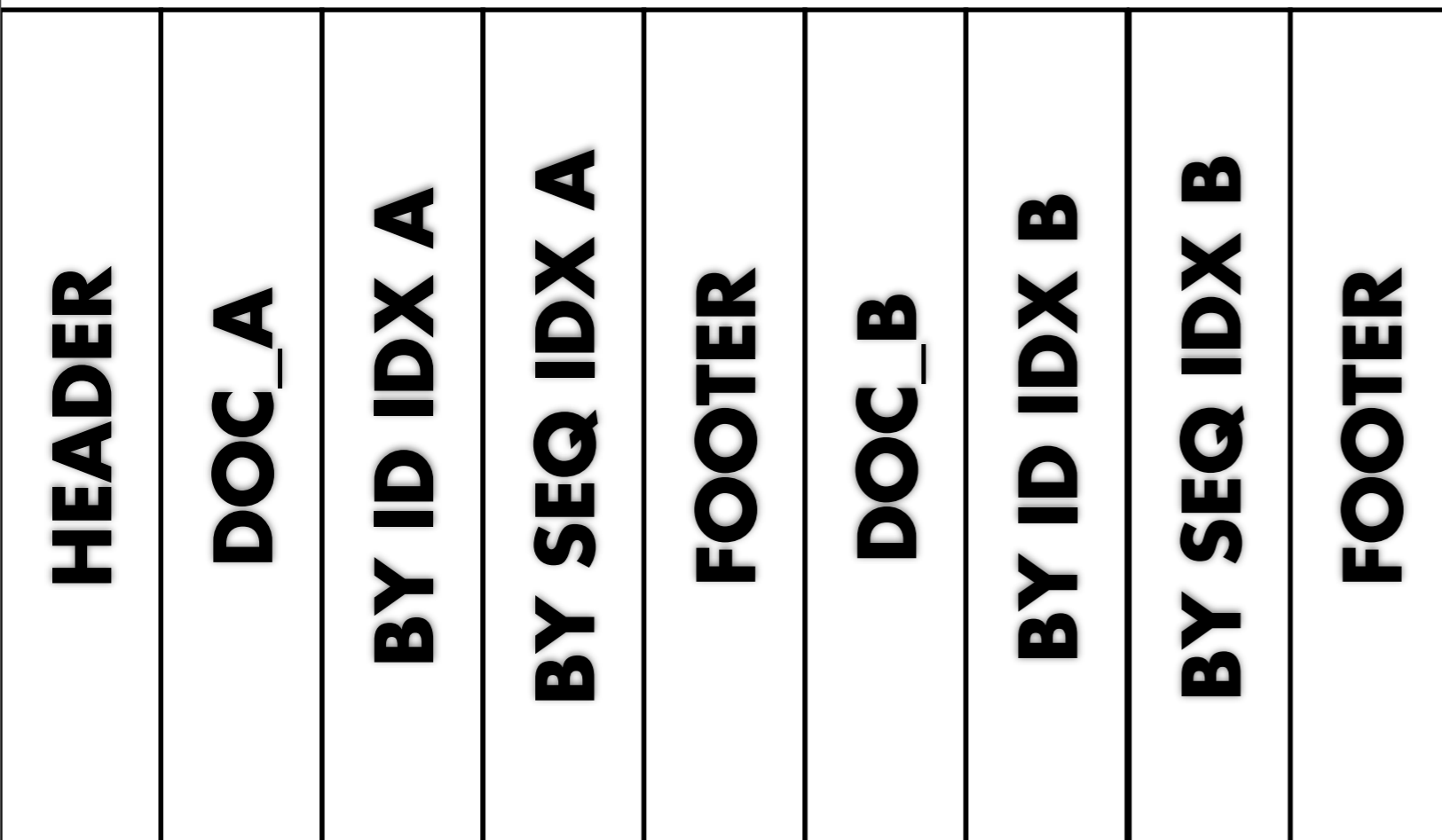
Monday, November 26, 12

- 2 x b+tree & data interleaved
- append only, mvcc
- full fsync control
  
- Can answer:
  - Data for \$key
  - What happened \$since
  
- Used for core data storage
- As well as indexes
  
- Everything else is built on top



Monday, November 26, 12

- 2 x b+tree & data interleaved
- append only, mvcc
- full fsync control
  
- Can answer:
  - Data for \$key
  - What happened \$since
  
- Used for core data storage
- As well as indexes
  
- Everything else is built on top



Monday, November 26, 12

- 2 x b+tree & data interleaved
- append only, mvcc
- full fsync control
  
- Can answer:
  - Data for \$key
  - What happened \$since
  
- Used for core data storage
- As well as indexes
  
- Everything else is built on top



Monday, November 26, 12

**Bulk add + Delete**

**DOC\_A**



Monday, November 26, 12

Bulk add + Delete



**DOC\_A**

**DOC\_B**



Monday, November 26, 12

Bulk add + Delete

**DOC\_A**

**DOC\_B**

**BY ID IDX A**



**DOC\_A**

**DOC\_B**

**BY ID IDX A**

**BY ID IDX B**



**DOC\_A**

**DOC\_B**

**BY ID IDX A**

**BY ID IDX B**

**BY SEQ IDX A**



**DOC\_A**

**DOC\_B**

**BY ID IDX A**

**BY ID IDX B**

**BY SEQ IDX A**

**BY SEQ IDX B**



Monday, November 26, 12

Bulk add + Delete

**DOC\_A**

**DOC\_B**

**BY ID IDX A**

**BY ID IDX B**

**BY SEQ IDX A**

**BY SEQ IDX B**

**FOOTER**



Monday, November 26, 12

Bulk add + Delete

**DOC\_A**

**DOC\_B**

**BY ID IDX A**

**BY ID IDX B**

**BY SEQ IDX A**

**BY SEQ IDX B**

**FOOTER**

**DEL DOC\_A**



**DOC\_A**

**DOC\_B**

**BY ID IDX A**

**BY ID IDX B**

**BY SEQ IDX A**

**BY SEQ IDX B**

**FOOTER**

**DEL DOC\_A**

**BY ID IDX A**





**DOC\_A**

**DOC\_B**

**BY ID IDX A**

**BY ID IDX B**

**BY SEQ IDX A**

**BY SEQ IDX B**

**FOOTER**

**DEL DOC\_A**

**BY ID IDX A**

**BY SEQ IDX A**





<b>DOC_A</b>
<b>DOC_B</b>
<b>BY ID IDX A</b>
<b>BY ID IDX B</b>
<b>BY SEQ IDX A</b>
<b>BY SEQ IDX B</b>
<b>FOOTER</b>
<b>DEL DOC_A</b>
<b>BY ID IDX A</b>
<b>BY SEQ IDX A</b>
<b>FOOTER</b>

# Operational Consequences

Monday, November 26, 12

- efficient on spinning disk, “tape”
- btree = wide, upper layers in disk cache
- backup with `cp $a $b`
  
- compaction hurts

# Core Features (using by-req)

Monday, November 26, 12

- Replication
- Indexing / Views / GeoCouch / Lucene / ES etc.
- /\_changes
- Compaction

Monday, November 26, 12

## Replication

**DATABASE A**

Monday, November 26, 12

Replication



Monday, November 26, 12

Replication

<b>2</b>
<b>1</b>
<b>DATABASE A</b>

Monday, November 26, 12

Replication



<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

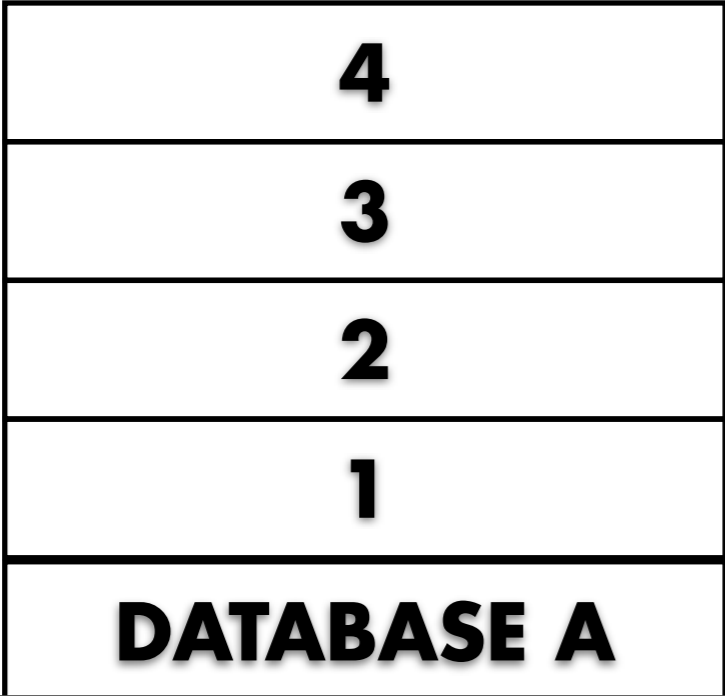
Monday, November 26, 12

Replication

<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

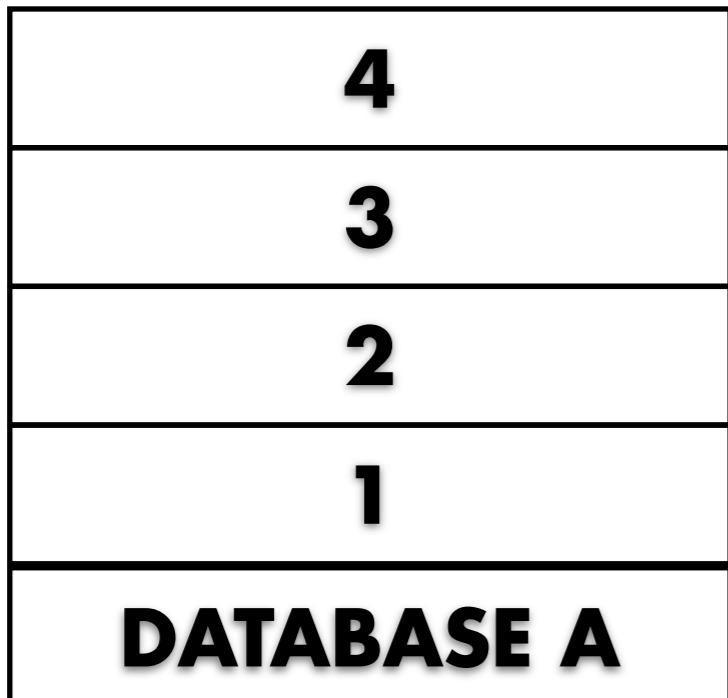
Monday, November 26, 12

Replication



Monday, November 26, 12

Replication



Monday, November 26, 12

Replication

<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

<b>2</b>
<b>1</b>
<b>DATABASE B</b>

Monday, November 26, 12

Replication

<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE B</b>

Monday, November 26, 12

Replication

<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE B</b>

Monday, November 26, 12

Replication

<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE B</b>

Monday, November 26, 12

Replication



<b>6</b>
<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE B</b>

Monday, November 26, 12

Replication

<b>7</b>
<b>6</b>
<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE B</b>

Monday, November 26, 12

Replication

<b>8</b>
<b>7</b>
<b>6</b>
<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE B</b>

Monday, November 26, 12

Replication

<b>8</b>
<b>7</b>
<b>6</b>
<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE B</b>

Monday, November 26, 12

Replication

<b>8</b>
<b>7</b>
<b>6</b>
<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

<b>6</b>
<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE B</b>

Monday, November 26, 12

Replication

<b>8</b>
<b>7</b>
<b>6</b>
<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

<b>7</b>
<b>6</b>
<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE B</b>

Monday, November 26, 12

Replication

<b>8</b>
<b>7</b>
<b>6</b>
<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

<b>8</b>
<b>7</b>
<b>6</b>
<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE B</b>

Monday, November 26, 12

Replication

Monday, November 26, 12

**Indexing**



# **DATABASE A**

Monday, November 26, 12

Indexing



Monday, November 26, 12

Indexing

<b>2</b>
<b>1</b>
<b>DATABASE A</b>

Monday, November 26, 12

Indexing

<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

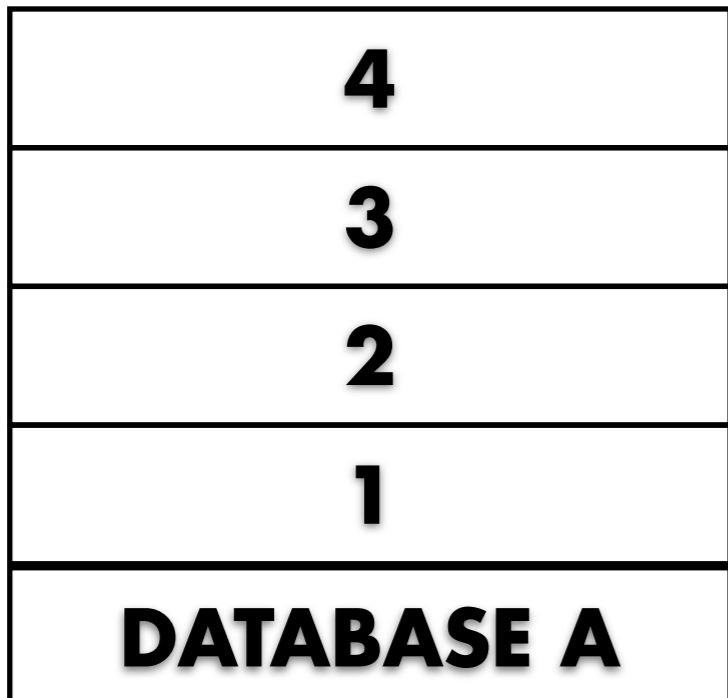
Monday, November 26, 12

Indexing

<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

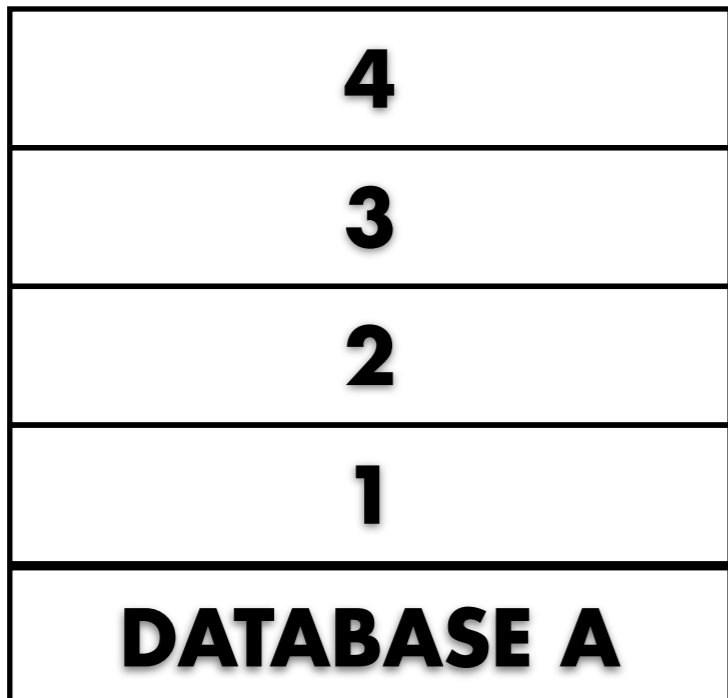
Monday, November 26, 12

Indexing



Monday, November 26, 12

Indexing



Monday, November 26, 12

Indexing

<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

<b>2</b>
<b>1</b>
<b>INDEX A</b>

Monday, November 26, 12

Indexing



<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

<b>3</b>
<b>2</b>
<b>1</b>
<b>INDEX A</b>

Monday, November 26, 12

Indexing

<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>INDEX A</b>

Monday, November 26, 12

Indexing

<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>INDEX A</b>

Monday, November 26, 12

Indexing

<b>6</b>
<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>INDEX A</b>

Monday, November 26, 12

Indexing

<b>7</b>
<b>6</b>
<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>INDEX A</b>

Monday, November 26, 12

Indexing

<b>8</b>
<b>7</b>
<b>6</b>
<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>INDEX A</b>

Monday, November 26, 12

Indexing

<b>8</b>
<b>7</b>
<b>6</b>
<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>INDEX A</b>

Monday, November 26, 12

Indexing

<b>8</b>
<b>7</b>
<b>6</b>
<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

<b>6</b>
<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>INDEX A</b>

Monday, November 26, 12

Indexing



<b>8</b>
<b>7</b>
<b>6</b>
<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

<b>7</b>
<b>6</b>
<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>INDEX A</b>

Monday, November 26, 12

Indexing

<b>8</b>
<b>7</b>
<b>6</b>
<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

<b>8</b>
<b>7</b>
<b>6</b>
<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>INDEX A</b>

Monday, November 26, 12

Indexing

Monday, November 26, 12

/\_changes

# **DATABASE A**

Monday, November 26, 12

/\_changes

**1**

**DATABASE A**

Monday, November 26, 12

/\_changes

<b>2</b>
<b>1</b>
<b>DATABASE A</b>

Monday, November 26, 12

/\_changes

<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

Monday, November 26, 12

/\_changes

<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

Monday, November 26, 12

/\_changes



<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

Monday, November 26, 12

/\_changes

<b>6</b>
<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

Monday, November 26, 12

/\_changes

<b>7</b>
<b>6</b>
<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

Monday, November 26, 12

/\_changes

<b>8</b>
<b>7</b>
<b>6</b>
<b>5</b>
<b>4</b>
<b>3</b>
<b>2</b>
<b>1</b>
<b>DATABASE A</b>

Monday, November 26, 12

/\_changes

Monday, November 26, 12

## Compaction

# **DATABASE A**

Monday, November 26, 12

Compaction

**1. DOC\_A**

**DATABASE A**

Monday, November 26, 12

Compaction

**2. DOC\_B**

**1. DOC\_A**

**DATABASE A**

Monday, November 26, 12

Compaction



<b>3. DOC_C</b>
<b>2. DOC_B</b>
<b>1. DOC_A</b>
<b>DATABASE A</b>

Monday, November 26, 12

Compaction

<b>4. DOC_A</b>
<b>3. DOC_C</b>
<b>2. DOC_B</b>
<b>1. DOC_A</b>
<b>DATABASE A</b>

Monday, November 26, 12

Compaction

<b>5. DOC_D</b>
<b>4. DOC_A</b>
<b>3. DOC_C</b>
<b>2. DOC_B</b>
<b>1. DOC_A</b>
<b>DATABASE A</b>

Monday, November 26, 12

Compaction

<b>6. DOC_B</b>
<b>5. DOC_D</b>
<b>4. DOC_A</b>
<b>3. DOC_C</b>
<b>2. DOC_B</b>
<b>1. DOC_A</b>
<b>DATABASE A</b>

Monday, November 26, 12

Compaction

<b>7. DOC_F</b>
<b>6. DOC_B</b>
<b>5. DOC_D</b>
<b>4. DOC_A</b>
<b>3. DOC_C</b>
<b>2. DOC_B</b>
<b>1. DOC_A</b>
<b>DATABASE A</b>

Monday, November 26, 12

Compaction

<b>8. DOC_G</b>
<b>7. DOC_F</b>
<b>6. DOC_B</b>
<b>5. DOC_D</b>
<b>4. DOC_A</b>
<b>3. DOC_C</b>
<b>2. DOC_B</b>
<b>1. DOC_A</b>
<b>DATABASE A</b>

Monday, November 26, 12

Compaction

<b>8. DOC_G</b>
<b>7. DOC_F</b>
<b>6. DOC_B</b>
<b>5. DOC_D</b>
<b>4. DOC_A</b>
<b>3. DOC_C</b>
<b>2. DOC_B</b>
<b>1. DOC_A</b>
<b>DATABASE A</b>

**COMPACT A**

<b>8. DOC_G</b>
<b>7. DOC_F</b>
<b>6. DOC_B</b>
<b>5. DOC_D</b>
<b>4. DOC_A</b>
<b>3. DOC_C</b>
<b>2. DOC_B</b>
<b>1. DOC_A</b>
<b>DATABASE A</b>

<b>3. DOC_C</b>
<b>COMPACT A</b>

Monday, November 26, 12

Compaction



<b>8. DOC_G</b>
<b>7. DOC_F</b>
<b>6. DOC_B</b>
<b>5. DOC_D</b>
<b>4. DOC_A</b>
<b>3. DOC_C</b>
<b>2. DOC_B</b>
<b>1. DOC_A</b>
<b>DATABASE A</b>

<b>4. DOC_A</b>
<b>3. DOC_C</b>
<b>COMPACT A</b>

Monday, November 26, 12

Compaction

<b>8. DOC_G</b>
<b>7. DOC_F</b>
<b>6. DOC_B</b>
<b>5. DOC_D</b>
<b>4. DOC_A</b>
<b>3. DOC_C</b>
<b>2. DOC_B</b>
<b>1. DOC_A</b>
<b>DATABASE A</b>

<b>5. DOC_D</b>
<b>4. DOC_A</b>
<b>3. DOC_C</b>
<b>COMPACT A</b>

Monday, November 26, 12

Compaction

<b>8. DOC_G</b>
<b>7. DOC_F</b>
<b>6. DOC_B</b>
<b>5. DOC_D</b>
<b>4. DOC_A</b>
<b>3. DOC_C</b>
<b>2. DOC_B</b>
<b>1. DOC_A</b>
<b>DATABASE A</b>

<b>6. DOC_B</b>
<b>5. DOC_D</b>
<b>4. DOC_A</b>
<b>3. DOC_C</b>
<b>COMPACT A</b>

Monday, November 26, 12

Compaction

<b>8. DOC_G</b>
<b>7. DOC_F</b>
<b>6. DOC_B</b>
<b>5. DOC_D</b>
<b>4. DOC_A</b>
<b>3. DOC_C</b>
<b>2. DOC_B</b>
<b>1. DOC_A</b>
<b>DATABASE A</b>

<b>7. DOC_F</b>
<b>6. DOC_B</b>
<b>5. DOC_D</b>
<b>4. DOC_A</b>
<b>3. DOC_C</b>
<b>COMPACT A</b>

Monday, November 26, 12

Compaction

<b>8. DOC_G</b>
<b>7. DOC_F</b>
<b>6. DOC_B</b>
<b>5. DOC_D</b>
<b>4. DOC_A</b>
<b>3. DOC_C</b>
<b>2. DOC_B</b>
<b>1. DOC_A</b>
<b>DATABASE A</b>

<b>8. DOC_G</b>
<b>7. DOC_F</b>
<b>6. DOC_B</b>
<b>5. DOC_D</b>
<b>4. DOC_A</b>
<b>3. DOC_C</b>
<b>COMPACT A</b>

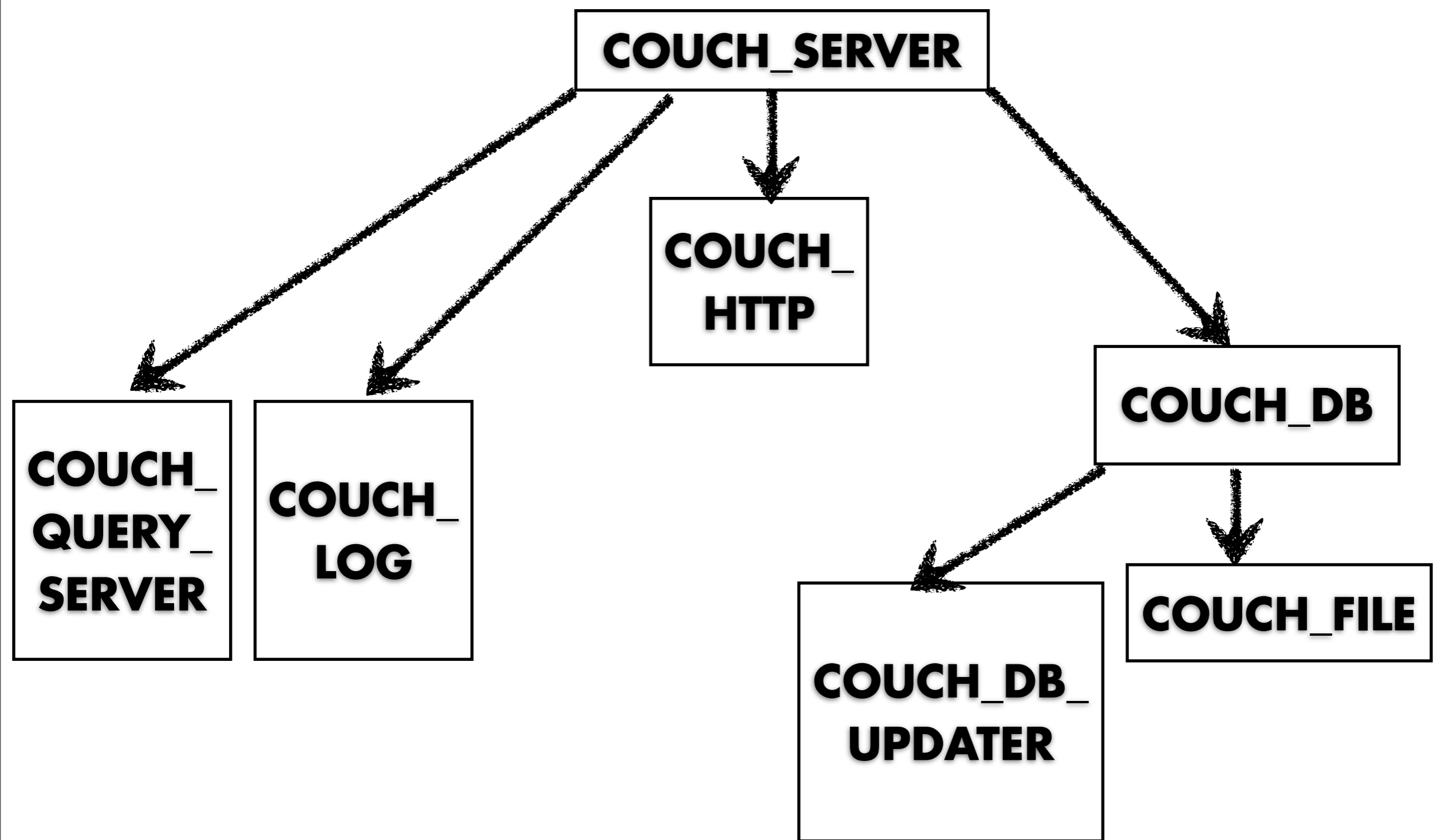
Monday, November 26, 12

Compaction

# Erlang

Monday, November 26, 12

- Small codebase
- Efficient in small teams
  
- Supervision tree
- Isolated processes
- Concurrency
  
- Portable runtime
  
- Hard to recruit for
- Steep ramp-on



Monday, November 26, 12

- Small codebase
- Efficient in small teams

- Supervision tree
- Isolated processes
- Concurrency

- Portable runtime

- Hard to recruit for
- Steep ramp-on

# Potential Improvements

Monday, November 26, 12

- Smarter compactor
- Smarter file-storage
- Less custom HTTP handling
- More indexers



**The End**

**Thanks!**

Monday, November 26, 12