

# Map-Reduce with CouchDB

Kore Nordmann  
<kore@php.net>  
@koredn

March 21, 2010

## Our documents

Introduction

Joins

Tree structures

QA

## ► Example: Issue tracker bug document

```
1 { "_id":      "issue-42",  
2   "type":    "issue",  
3   "title":   "Bug in module foo",  
4   "text":    "Your software sucks!",  
5   "creator": "user-bar",  
6   "state":   "new",  
7   "edited":  1935678239,  
8   ...  
9 }
```

## ▶ Example: Issue tracker bug comment

```
1 { "_id": "80572348asdfg789342",
2   "type": "issue-comment",
3   "issue": "issue-42",
4   "edited": 1935678239,
5   "text": "I second that!",
6   ...
7 }
```

Our documents

Introduction

Joins

Tree structures

QA



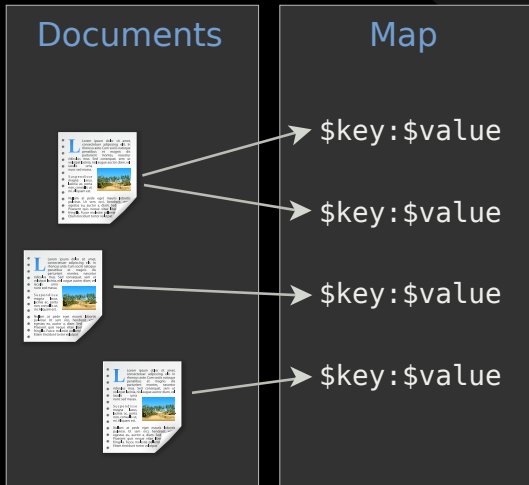
- ▶ “MapReduce is a software framework introduced by Google to support distributed computing on large data sets on clusters of computers.” [Wik09]
- ▶ Used by CouchDB to implement views

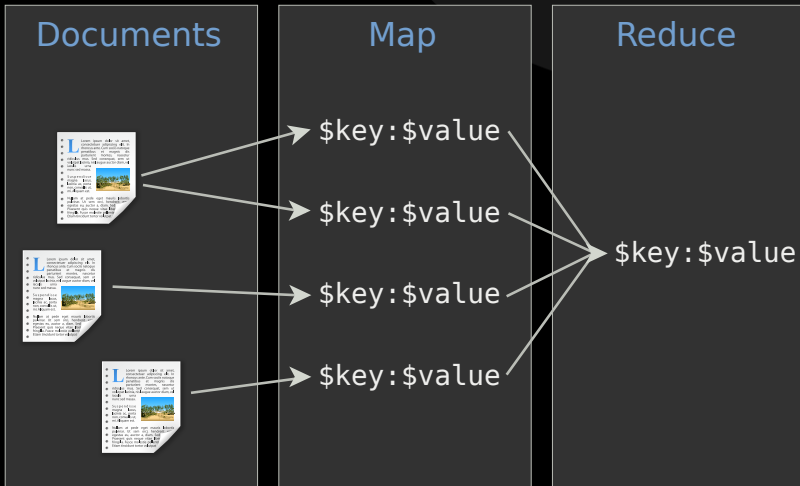
- ▶ “MapReduce is a software framework introduced by Google to support distributed computing on large data sets on clusters of computers.” [Wik09]
- ▶ Used by CouchDB to implement views
- ▶ Just a framework / pattern: You can implement “any” algorithm using map-reduce.

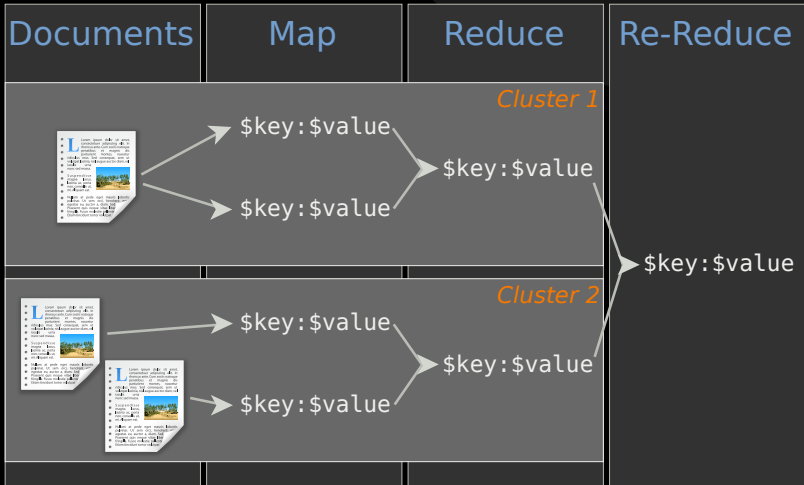


## Documents









- ▶ Map and reduce functions are custom

- ▶ Map and reduce functions are custom
- ▶ Reduce is optional, plain view serves as a document index

- ▶ Map and reduce functions are custom
- ▶ Reduce is optional, plain view serves as a document index
- ▶ Reduce may be applied to subsets of the documents

- ▶ Map and reduce functions are custom
- ▶ Reduce is optional, plain view serves as a document index
- ▶ Reduce may be applied to subsets of the documents
- ▶ Reduce may be grouped



## ► Example: Issue tracker bug document

```
1 { "_id": "issue-42",  
2   "type": "issue",  
3   "title": "Bug in module foo",  
4   "text": "Your software sucks!",  
5   "creator": "user-bar",  
6   "state": "new",  
7   "edited": 1935678239,  
8   ...  
9 }
```

## ▶ Count issue states

```
1 function( doc )
2 {
3     if ( doc.type == "issue" )
4     {
5         emit( doc.state , 1 );
6     }
7 }
```

- ▶ The simplest reduce function is just `count()`
  - ▶ Often used for statistics

```
1 function( keys , values , combine )
2 {
3     return sum( values );
4 }
```

## ► The mapping result

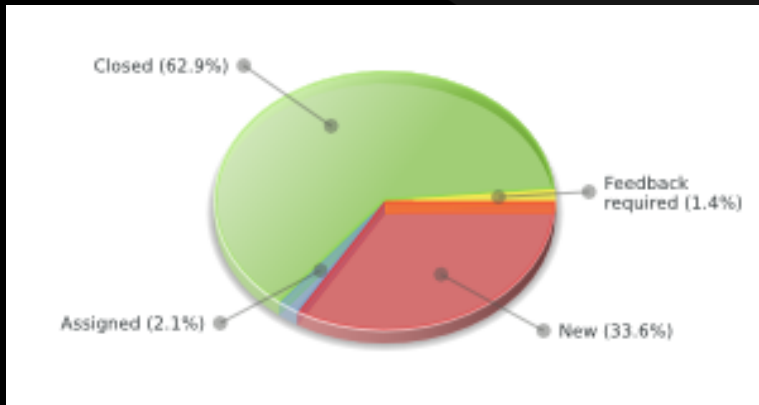
```
1 "new"      => 1,  
2 "assigned" => 1,  
3 "new"      => 1,  
4 "closed"   => 1,  
5 "new"      => 1,  
6 "new"      => 1,  
7 "closed"   => 1,  
8 "assigned" => 1,  
9 "assigned" => 1,
```

▶ The reduce result

1 **null** => 12

► The grouped reduce result

```
1 "new"      => 5 ,  
2 "assigned" => 23 ,  
3 "closed"   => 42 ,
```



## ► The map function

```
1 function( doc )
2 {
3     if ( doc.type == "wiki" )
4     {
5         date = new Date();
6         date.setTime( doc.edited * 1000 );
7         emit( [
8             date.getUTCFullYear(),
9             date.getUTCMonth() + 1,
10            date.getUTCDate(),
11            date.getUTCHours(),
12            date.getUTCMinutes(),
13            date.getUTCSeconds(),
14            ], 1 );
15        // You could also emit the whole doc as value
16    }
17 }
```



## ► The mapping result

```
1 [2008, 10, 11, 9, 11, 12] => 1
2 [2008, 10, 11, 9, 11, 12] => 1
3 [2008, 10, 11, 9, 11, 12] => 1
4 [2008, 10, 11, 9, 13, 8] => 1
5 [2008, 10, 11, 9, 13, 44] => 1
6 [2008, 10, 11, 9, 14, 2] => 1
7 [2008, 10, 12, 17, 46, 15] => 1
8 [2008, 10, 12, 17, 57, 52] => 1
9 [2008, 10, 12, 18, 0, 45] => 1
10 [2008, 10, 14, 8, 36, 29] => 1
11 [2008, 10, 14, 19, 33, 21] => 1
12 [2008, 10, 14, 19, 33, 35] => 1
```

## ► The reduce function

```
1 function( keys , values , combine )  
2 {  
3     return sum( values );  
4 }
```

▶ The reduce result

1 `null`  $\Rightarrow$  12

## ► The grouped reduce result

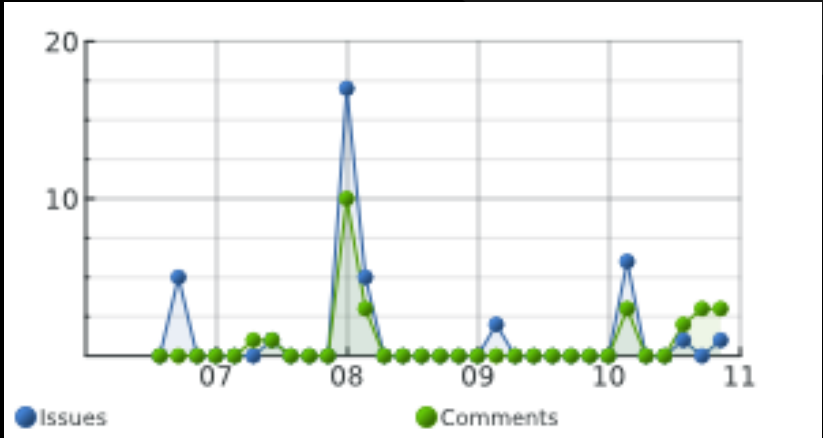
```
1 [2008, 10, 11, 9, 11, 12] => 3
2 [2008, 10, 11, 9, 13, 8] => 1
3 [2008, 10, 11, 9, 13, 44] => 1
4 [2008, 10, 11, 9, 14, 2] => 1
5 [2008, 10, 12, 17, 46, 15] => 1
6 [2008, 10, 12, 17, 57, 52] => 1
7 [2008, 10, 12, 18, 0, 45] => 1
8 [2008, 10, 14, 8, 36, 29] => 1
9 [2008, 10, 14, 19, 33, 21] => 1
10 [2008, 10, 14, 19, 33, 35] => 1
```

- ▶ The filtered grouped reduce result
- ▶ `startkey=[2008,10,11]` and `endkey=[2008,10,12]`

```
1 [2008, 10, 11, 9, 11, 12] => 3
2 [2008, 10, 11, 9, 13, 8] => 1
3 [2008, 10, 11, 9, 13, 44] => 1
4 [2008, 10, 11, 9, 14, 2] => 1
```

- ▶ The grouped reduce result, with group level
- ▶ `group-level=3`

1	[2008, 10, 11]	=>	6
2	[2008, 10, 12]	=>	3
3	[2008, 10, 14]	=>	3



## ► Index all documents by all their words

```
1 function( doc ) {
2   if ( doc.type == "issue" ) {
3     // Simple word indexing, does not respect overall
4     // occurrences of words,
5     // stopwords, different word separation characters,
6     // or word variations.
7     var text = doc.title.replace( /\s:.,!?-]+/g, "_" )
8     +
9     doc.text.replace( /\s:.,!?-]+/g, "_" )
10    ;
11    var words = text.split( "_" );
12    for ( var i = 0; i < words.length; ++i ) {
13      value = {};
14      value[doc._id] = 1;
15      emit( words[i].toLowerCase(), value );
16    }
17  }
18 }
```



## ► Index all documents by all their words

```
1  ...
2  "a"      => {issue -8: 1}
3  "a"      => {issue -8: 1}
4  "a"      => {issue -8: 1}
5  "a"      => {issue -8: 1}
6  "a"      => {issue -81: 1}
7  "a"      => {issue -83: 1}
8  "a"      => {issue -83: 1}
9  "able"   => {issue -39: 1}
10 "able"   => {issue -56: 1}
11 "able"   => {issue -73: 1}
12 "able"   => {issue -80: 1}
13 "about"  => {issue -24: 1}
14 "about"  => {issue -43: 1}
15 "about"  => {issue -85: 1}
16  ...
```

## ► Reduce by word count

```
1 function( keys , values ) {
2   var count = {};
3   for ( var i in values ) {
4     for ( var id in values[i] ) {
5       if ( count[id] ) {
6         count[id] = values[i][id] + count[id];
7       } else {
8         count[id] = values[i][id];
9       }
10    }
11  }
12  return count;
13 }
```

## ► Index all documents by all their words

```
1  ...
2  "a"    => {
3      issue -68: 6,
4      issue -66: 6,
5      issue -22: 4,
6      issue -63: 3,
7      issue -60: 2,
8      issue -35: 2,
9      issue -34: 1,
10     issue -31: 1,
11     ...
12     }
13  "able" => {issue -86: 1, issue -80: 1, issue -73: 1,
14           issue -56: 1, issue -39: 1}
15  "about" => {issue -85: 1, issue -43: 1, issue -24: 1}
16  ...
```

Our documents

Introduction

**Joins**

Tree structures

QA

- ▶ There is no ensured inter document consistency in CouchDB

ents:

- ▶ There is no ensured inter document consistency in CouchDB
- ▶ Different possibilities of relating documents:

- ▶ There is no ensured inter document consistency in CouchDB
- ▶ Different possibilities of relating documents:
  - ▶ List IDs of related documents in document (n:m)

- ▶ There is no ensured inter document consistency in CouchDB
- ▶ Different possibilities of relating documents:
  - ▶ List IDs of related documents in document (n:m)
  - ▶ ... both directions are feasible



- ▶ There is no ensured inter document consistency in CouchDB
- ▶ Different possibilities of relating documents:
  - ▶ List IDs of related documents in document (n:m)
  - ▶ ... both directions are feasible
  - ▶ Embed the whole related document (1:n)

- ▶ There is no ensured inter document consistency in CouchDB
- ▶ Different possibilities of relating documents:
  - ▶ List IDs of related documents in document (n:m)
  - ▶ ... both directions are feasible
  - ▶ Embed the whole related document (1:n)
- ▶ Solution depends on update-ratio

```
1 { "type": "issue",  
2   "title": "Hello world",  
3   "text": "...",  
4   "comments": [  
5     { "comment": "..."} ],  
6   ],  
7   "creator": "user-foo",  
8 }
```

## ▶ JOIN query

```
1 function( doc )
2 {
3     if ( doc.type == "issue" )
4     {
5         emit( [doc._id , 0], doc._id );
6     }
7
8     if ( doc.type == "issue_comment" )
9     {
10        emit( [doc.issue , doc._id], doc._id );
11    }
12 }
```

## ▶ JOIN query result

```
1 ["tracker_issue -1", 0] => "tracker_issue -1"
2 ["tracker_issue -1", "1eaaa503adbec07b0013c060b5d7b53c"]
   => "1eaaa503adbec07b0013c060b5d7b53c"
3 ["tracker_issue -1", "6e77ebd5ea0383cf1d8ed505f2517ca7"]
   => "6e77ebd5ea0383cf1d8ed505f2517ca7"
4 ["tracker_issue -1", "b0502afb84a80a48fecc7442153c0aa2"]
   => "b0502afb84a80a48fecc7442153c0aa2"
5 ["tracker_issue -10", 0] => "tracker_issue -10"
6 ["tracker_issue -10", "63872cd8a4d5301238dbdb084a4d7a3f"]
   ] => "63872cd8a4d5301238dbdb084a4d7a3f"
7 ["tracker_issue -10", "9af0dd3c4d184c07c2ff98982f98d39f"]
   ] => "9af0dd3c4d184c07c2ff98982f98d39f"
```

▶ Can again be filtered...

▶ Using “?include\_docs=true” also gives you all documents  
(since 0.11)

Our documents

Introduction

Joins

Tree structures

QA

```
1 { "type": "forum-post",  
2   "title": "Hello world",  
3   "text": "...",  
4   "parents": [ "forum-post-23", "forum-post-42" ]  
5   "creator": "user-foo",  
6 }
```

## ▶ Tree view

```
1 function( doc )
2 {
3     if ( doc.type == "forum-post" )
4     {
5         location = doc.parents;
6         location.append( doc._id );
7         emit( location , doc._id );
8     }
9 }
```

## ▶ Tree query result

```
1 ["forum-post-1"]
2   => "forum-post-1"
3 ["forum-post-1", "forum-post-2"]
4   => "forum-post-2"
5 ["forum-post-1", "forum-post-2", "forum-post-4"]
6   => "forum-post-4"
7 ["forum-post-1", "forum-post-3"]
8   => "forum-post-3"
```



- ▶ Tree query subtree result
- ▶ `startkey=["forum-post-1", "forum-post-2"]`

```
1 ["forum-post-1", "forum-post-2"]  
2   => "forum-post-2"  
3 ["forum-post-1", "forum-post-2", "forum-post-4"]  
4   => "forum-post-4"
```

Our documents

Introduction

Joins

Tree structures

QA

- ▶ Open questions?
- ▶ Further remarks?
- ▶ Contact
  - ▶ Mail: <kore@php.net>
  - ▶ Web: <http://kore-nordmann.de/> (Slides will be available here soonish)
  - ▶ Twitter: <http://twitter.com/koredn>

[Wik09] Wikipedia, *Mapreduce* — *wikipedia, the free encyclopedia*, 2009, [Online; accessed 27-August-2009].