# Riak Search

A Full-Text Search
and Indexing Engine
based on Riak



Erlang Factory · London · June 2010

Basho Technologies

Rusty Klophaus - @rklophaus

Why did we build it?

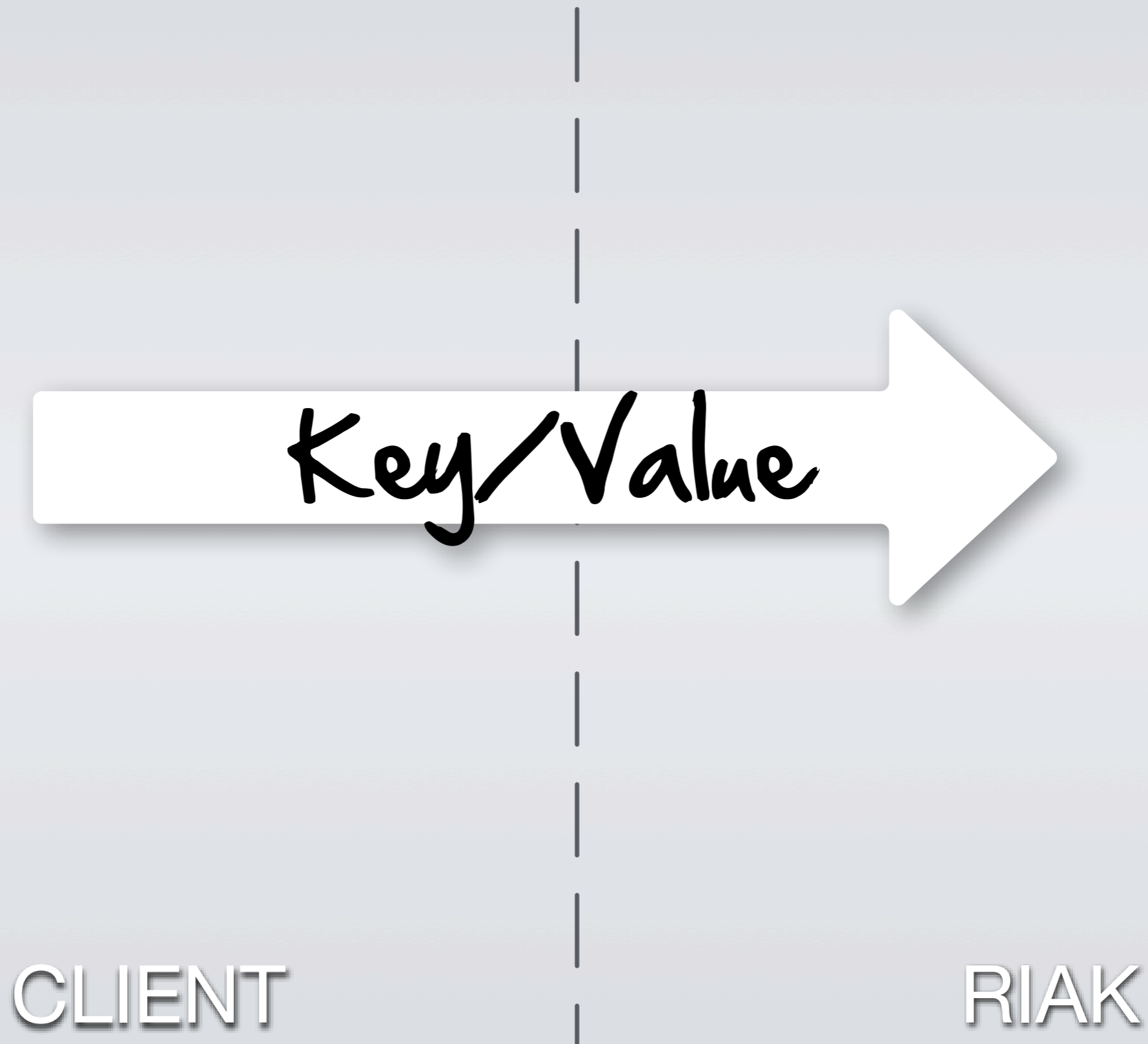What are the major goals?

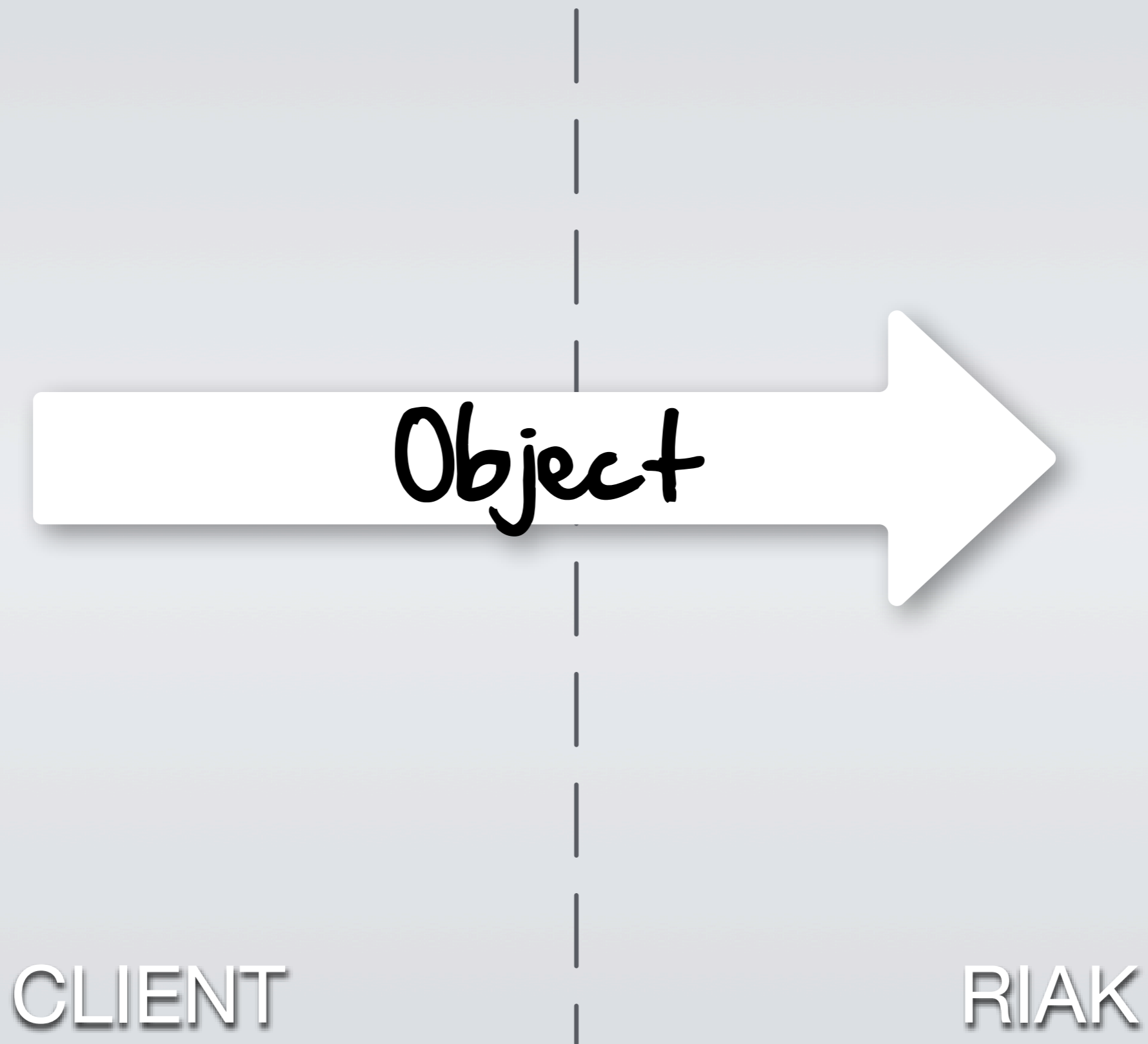How does it work?

# Part One

Why did we build

Riak Search?

Riak is

a scalable, highly-available, networked,
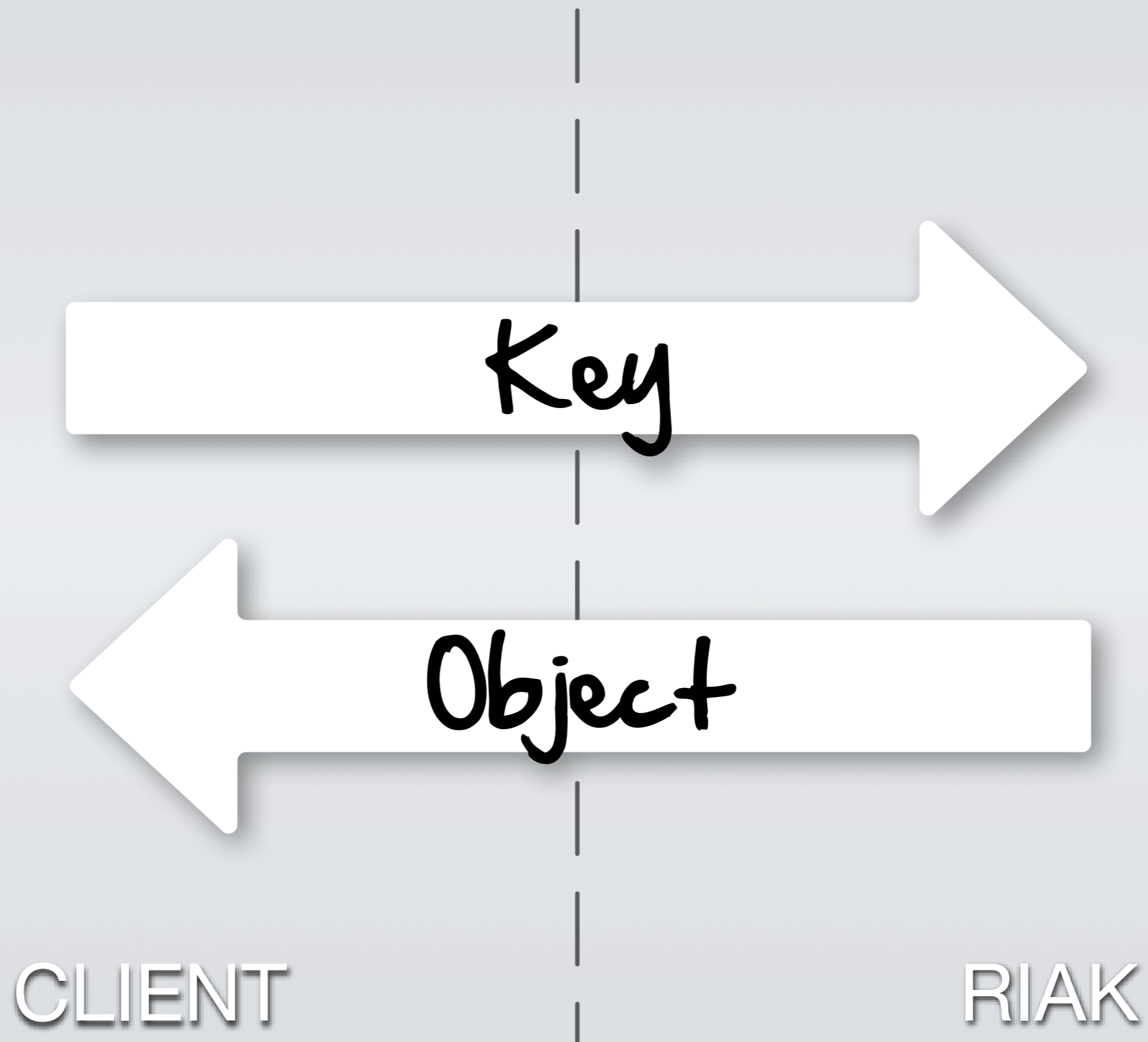
open-source key/value store.
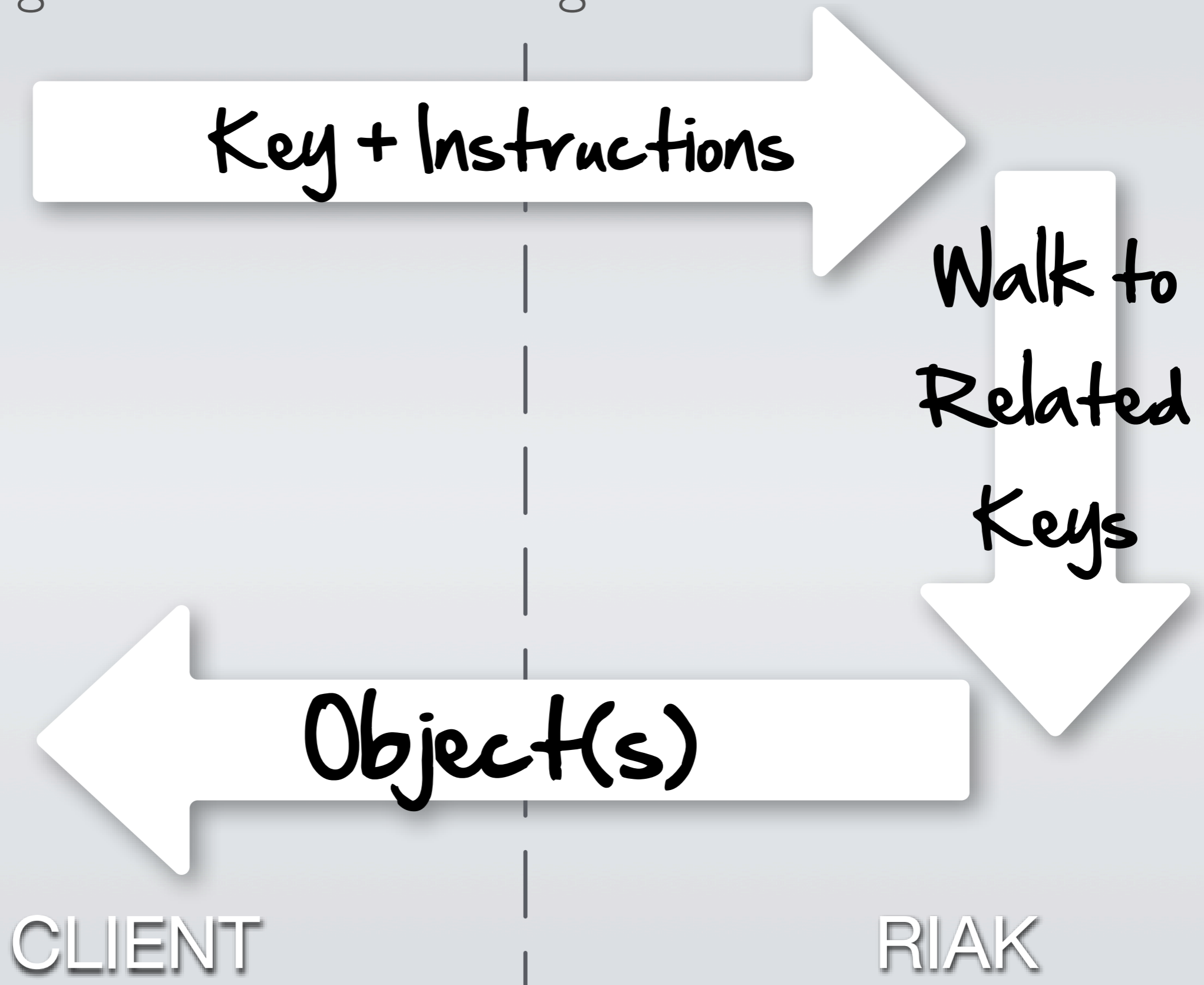
# Writing to a Key/Value Store

Key/Value

CLIENT

RIAK

basho

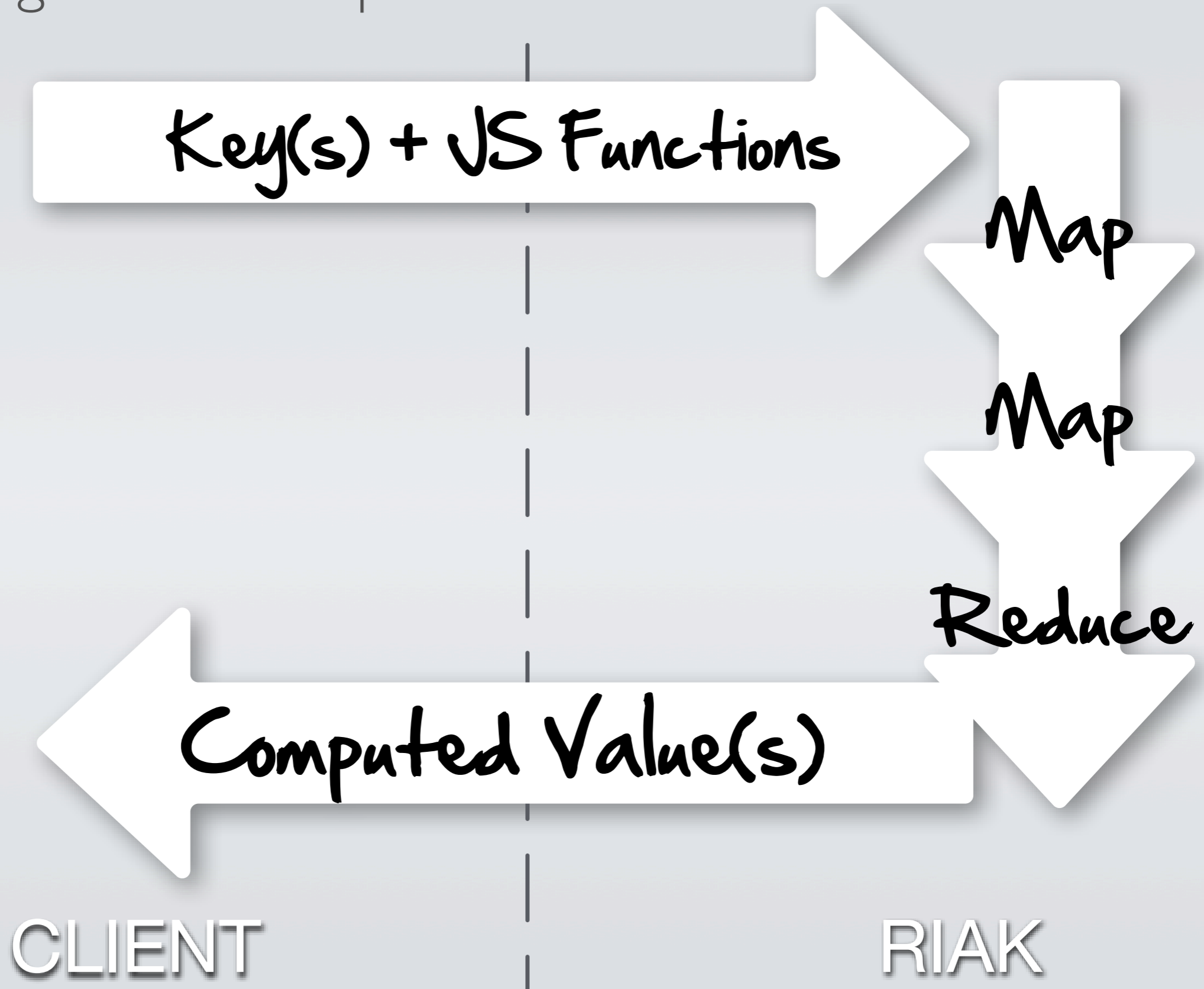# Writing to a Key/Value Store



CLIENT

RIAK

basho

# Querying a Key/Value Store

Key →

← Object

CLIENT | RIAK

Querying Riak via Map/Reduce

Key(s) + JS Functions

Map

Map

Reduce

Computed Value(s)

CLIENT

RIAK

basho

9

# Key/Value Stores
# like
# Key-Based Queries

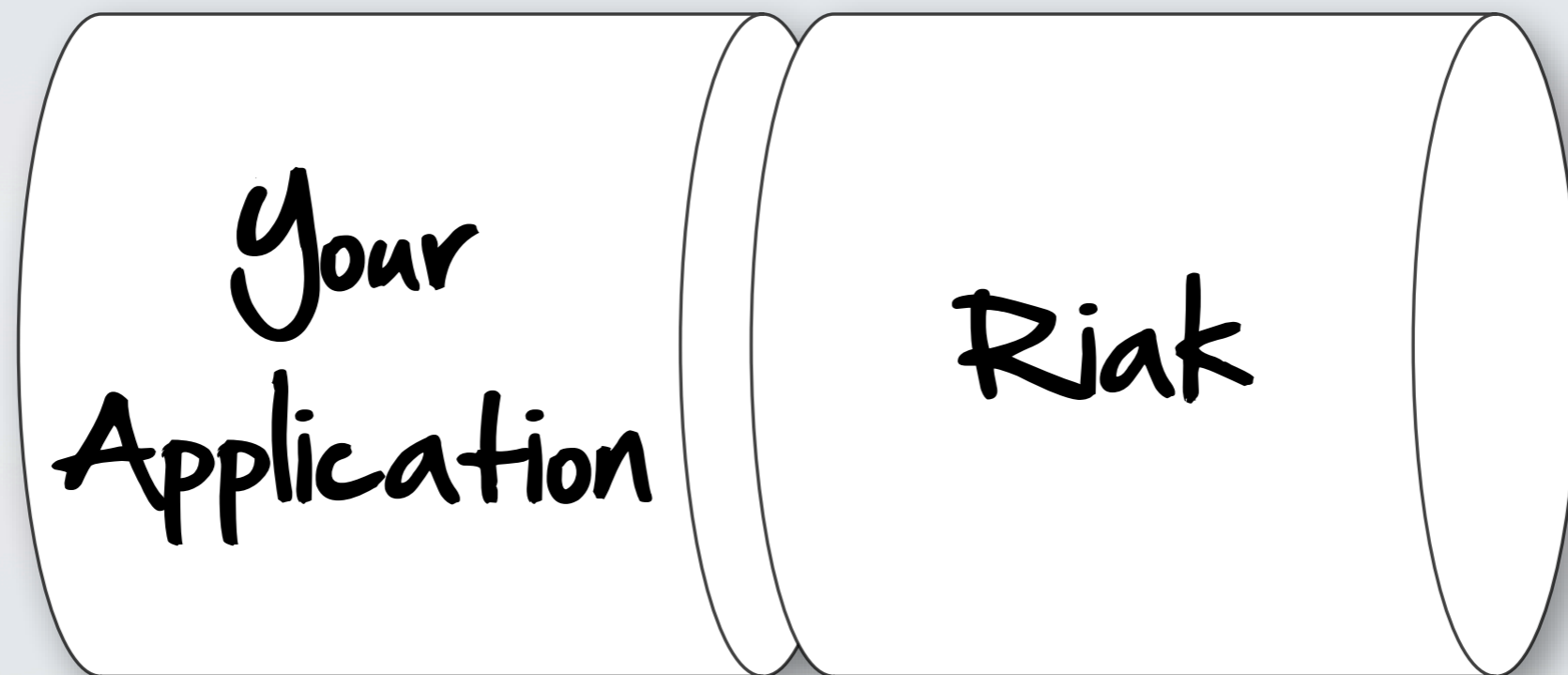"Converse AND Shoes"

This is getting old.

CLIENT

RIAK

basho

These kinds of queries

need an Index.


*Market Opportunity!*
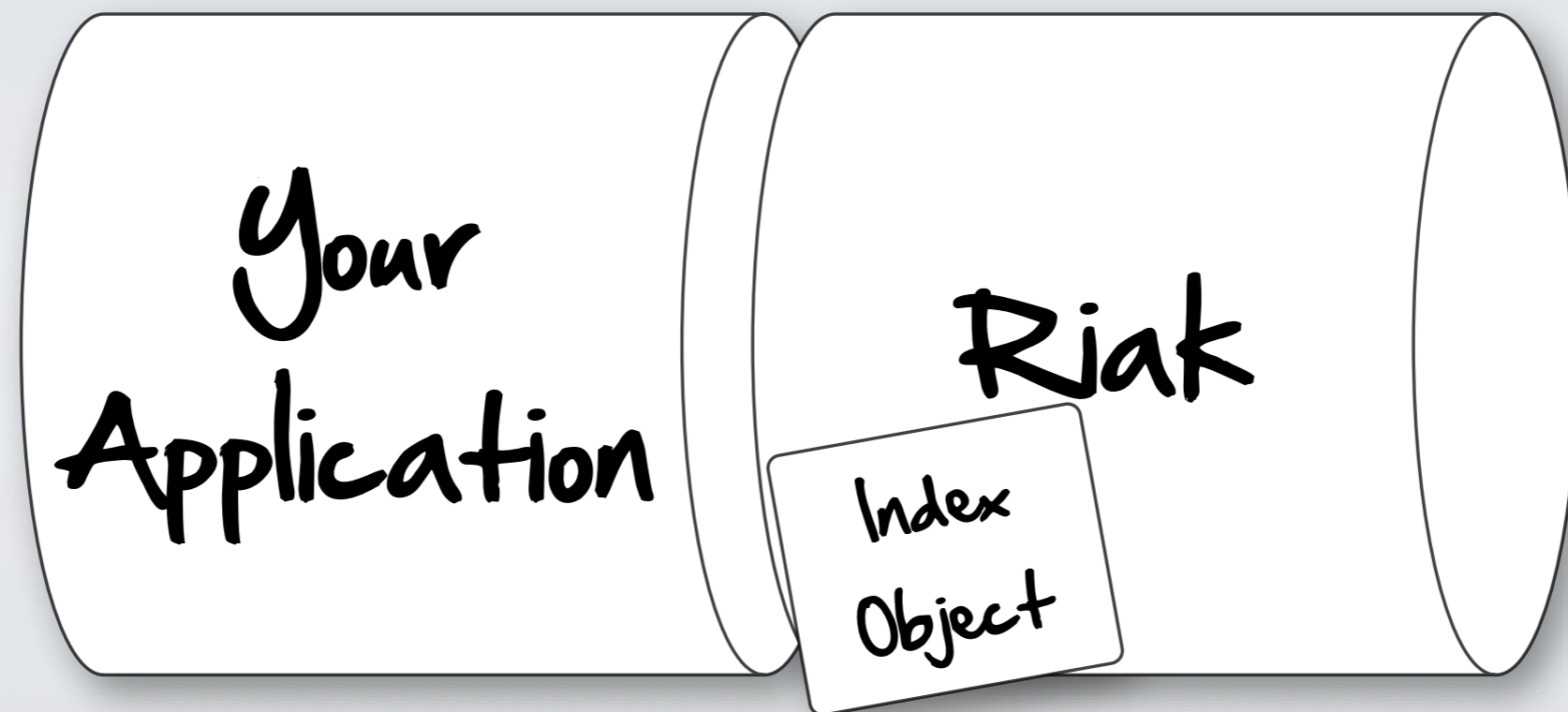
# Part Two

What are the major

goals of Riak Search?

# An application built on Riak.

Your
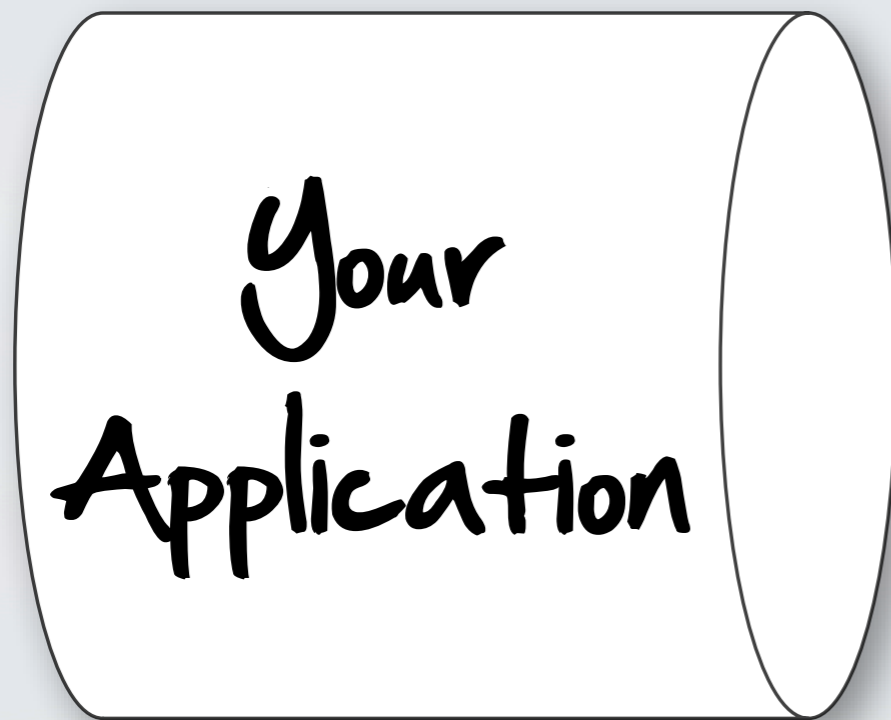Application

Riak

# Hrm... I need an index.

# Hrm... I need an index with more features.

Your Application

???

Riak

# Lucene should do the trick...

# ...replicate to add more throughput.

Your Application

Lucene Lucene Lucene

Lucene Lucene Lucene

Lucene Lucene Lucene

Riak

Operations nightmare!

# What do we really want?



Your Application | Riak-ified Lucene | Riak

# What do we really want?

# Functionality? Be like Lucene (and more).

- Lucene Syntax

- Leverages Java Lucene Analyzers

- Solr Endpoints

- Integration via Riak Post-Commit Hook (Index)

- Integration via Riak Map/Reduce (Query)

- Near-Realtime

- Schema-less

# Operations? Be like Riak.

- No special nodes

- Add nodes, get more compute and storage

- Automatically load balance

- Replicas for durability and performance

- Index and query in parallel

- Swappable storage backends

# Part Three

How do we do it?

# A Gentle Introduction to Document Indexing

# Document

# Inverted Index

**#1** Every dog has his day.

```
day, 1
dog, 1
every, 1
has, 1
his, 1
```

basho

## Documents

**#1** Every dog has his day.

**#2** The dog's bark is worse than his bite.

**#3** Let the cat out of the bag.

**#4** It's raining cats and dogs.

## Combined Inverted Index

```
and, 4
bag, 3
bark, 2
bite, 2
cat, 3
cat, 4
day, 1
dog, 1
dog, 2
dog, 4
every, 1
has, 1
...
```

"dog AND cat"

AND

dog

cat

Result: 4

**AND**
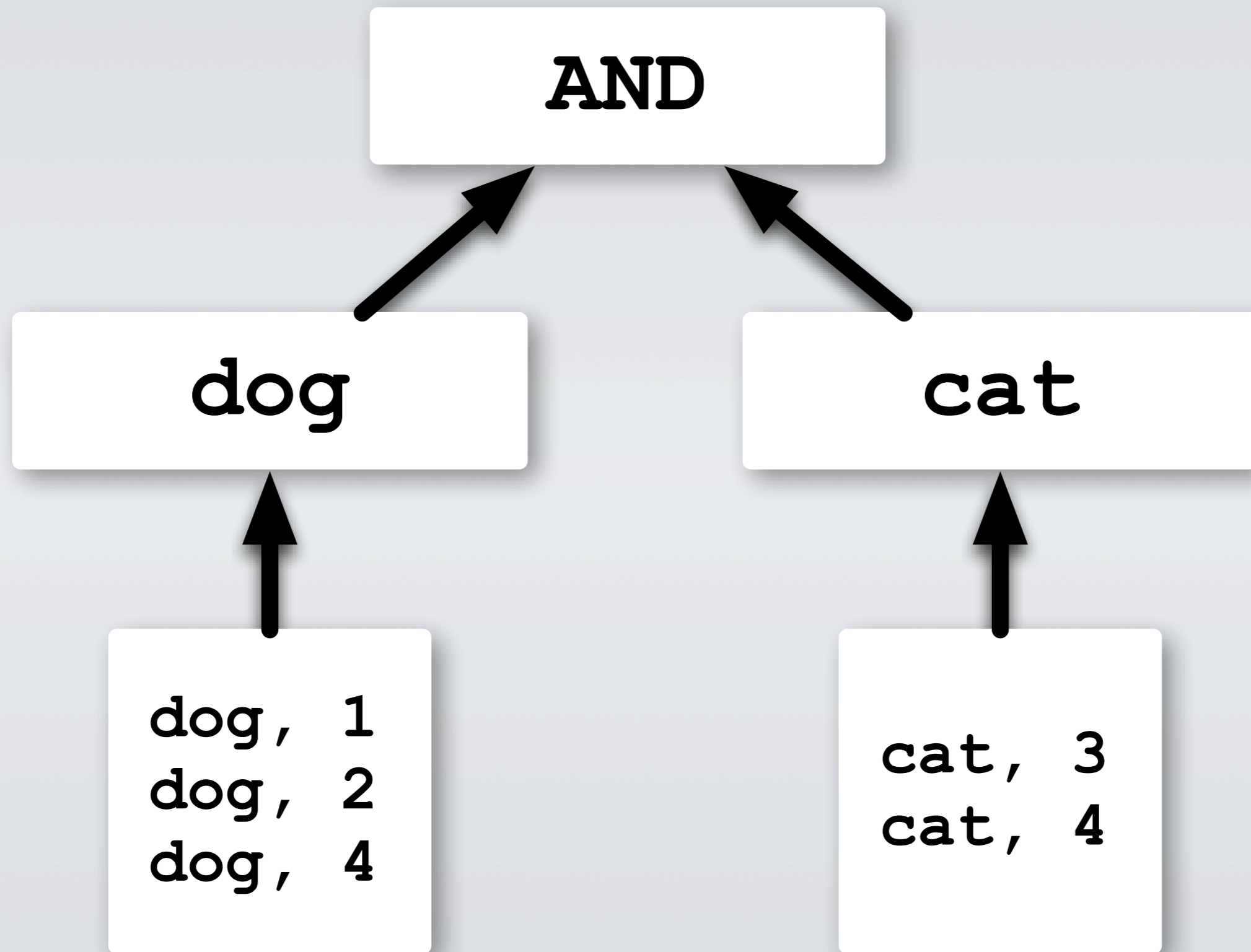**(Merge Intersection)**

1
2
4

3
4

basho

## At Query Time...

Result: 1, 2, 3, 4

OR

(Merge Union)

1
2
4

3
4

# Complex Behavior from Simple Structures

# Storage Approaches...

Riak Search uses
Consistent Hashing
to store data on
Partitions

```
Partitions = 10

Number of Nodes = 5

Partitions per Node = 2

Replicas (NVal) = 2
```

Object

# Document Partitioning

## vs.

# Term Partitioning

# ...and the
# Resulting Tradeoffs

#1

Every dog has his day.

"dog OR cat"

# Term Partitioning @ Index Time

#1 Every dog has
his day.

➤

```
day, 1
dog, 1
every, 1
has, 1
his, 1
```

basho

# Term Partitioning @ Index Time



**day, 1**

**dog, 1**

**has, 1**

**his, 1**

**every, 1**

"dog OR cat"

# Tradeoffs...

## Document Partitioning

+ Lower Latency Queries

- Lower Throughput

- Lots of Disk Seeks

## Term Partitioning

- Higher Latency Queries

+ Higher Throughput

- Hotspots in Ring
  (the "Obama" problem)

# Riak Search: Term Partitioning

**Term-partitioning** is the most viable approach for our beta clients' needs: high throughput on Really Big Datasets.

**Optimizations**:

- Term splitting to reduce hot spots

- Bloom filters & caching to save query-time bandwidth

- Batching to save query-time & index-time bandwidth

Support for either approach eventually.

# Diving Deeper:
# The Lifecycle of a Query

# Parse the Query

# The Query

**meeting AND (face OR phone)**

```
[{land, [
    {term,"meeting",[]},
    {lor,[
        {term,"face",[]},
        {term,"phone",[]}
    ]}
]}]
```

# The Query as a Graph

```
                    #land
                   /      \
                  /        \
              #term        #lor
                |         /      \
                |        /        \
          "meeting"   #term      #term
                        |           \
                        |            \
                    "phone"        "face"
```

52

# Plan the Query

**Term** → **TermID** → **Term Weight & File Offset**

**#land**

**#term**

**#lor**

**"meeting"**

**#term**

**#term**

**23 @ node B**

**"phone"**

**"face"**

**17 @ node A**    **13 @ node C**

# Use Term Weights to Plan the Query

**#node@B**

#land

#term          **#node@A**

"meeting"          #lor

23 @ node B          #term          #term

"phone"          "face"

17 @ node A          13 @ node C

# The Node-Assigned Query as an Erlang Term

```erlang
[{node,
    {land, [
        {node,
            {lor, [
                {term,{"email","body","face"}, [
                    {node_weight,'node_c@127.0.0.1', 13}
                ]},
                {term,{"email","body", "phone"}, [
                    {node_weight,'node_a@127.0.0.1', 17}
                ]}
            ]},
            'node_a@127.0.0.1'
        },
        {term, {"email","body","meeting"}, [
            {node_weight,'node_b@127.0.0.1', 23}
        ]}
    ]},
    'node_b@127.0.0.1'
}]
```

# Execute the Query

# Spawn the Query Processes

```
                    #node@B
                       |
                       v
                    #land
                   /      \
                  v        v
              #term       #node@A
                |            |
                v            v
           "meeting"       #lor
                          /    \
                         v      v
                     #term      #term
                       |          |
                       v          v
                   "phone"      "face"
```

# Spawn the Query Processes

# Spawn the Query Processes

```
          #node@B
             │
             ▼
          #land
          ╱      ╲
     #term        #node@A
       │              │
       ▼              ▼
  "meeting"         #lor
                   ╱     ╲
              #term       #term
                │            │
                ▼            ▼
            "phone"       "face"
```

# Spawn the Query Processes



65

# Spawn the Query Processes

```
#node@B
```

```
#land
```

```
#term
```

```
#node@A
```

```
"meeting"
```

```
#lor
```

```
#term
```

```
#term
```

```
"phone"
```

```
"face"
```

#node@B

#land

#term

'disconnect'

#node@A

#lor

#term

'disconnect'

#term

'disconnect'

# Message Format

# The Message Format

```
Message ::
    {results, [Result]} |
    {results, disconnect}

Result ::
    {DocID, Properties}

DocID ::
    term()

Properties ::
    proplist()
```

# The Message Format

```
{results, [
    {375, []},
    {961, [{color, "red"}]},
    {155, [{pos, [1,2,5]}]}
]}
```

# Yay for Erlang!

- Clean lines between load balancing and logic, single- and multi-node look the same

- Easy to create new operators, rapid development of experimental features

- Linked processes make cleanup a breeze

- Significant code reduction over early Java prototypes

# Part Four

Review

# Riak Search turns this...



"Converse AND Shoes"

WTF!? I'm a KV store!
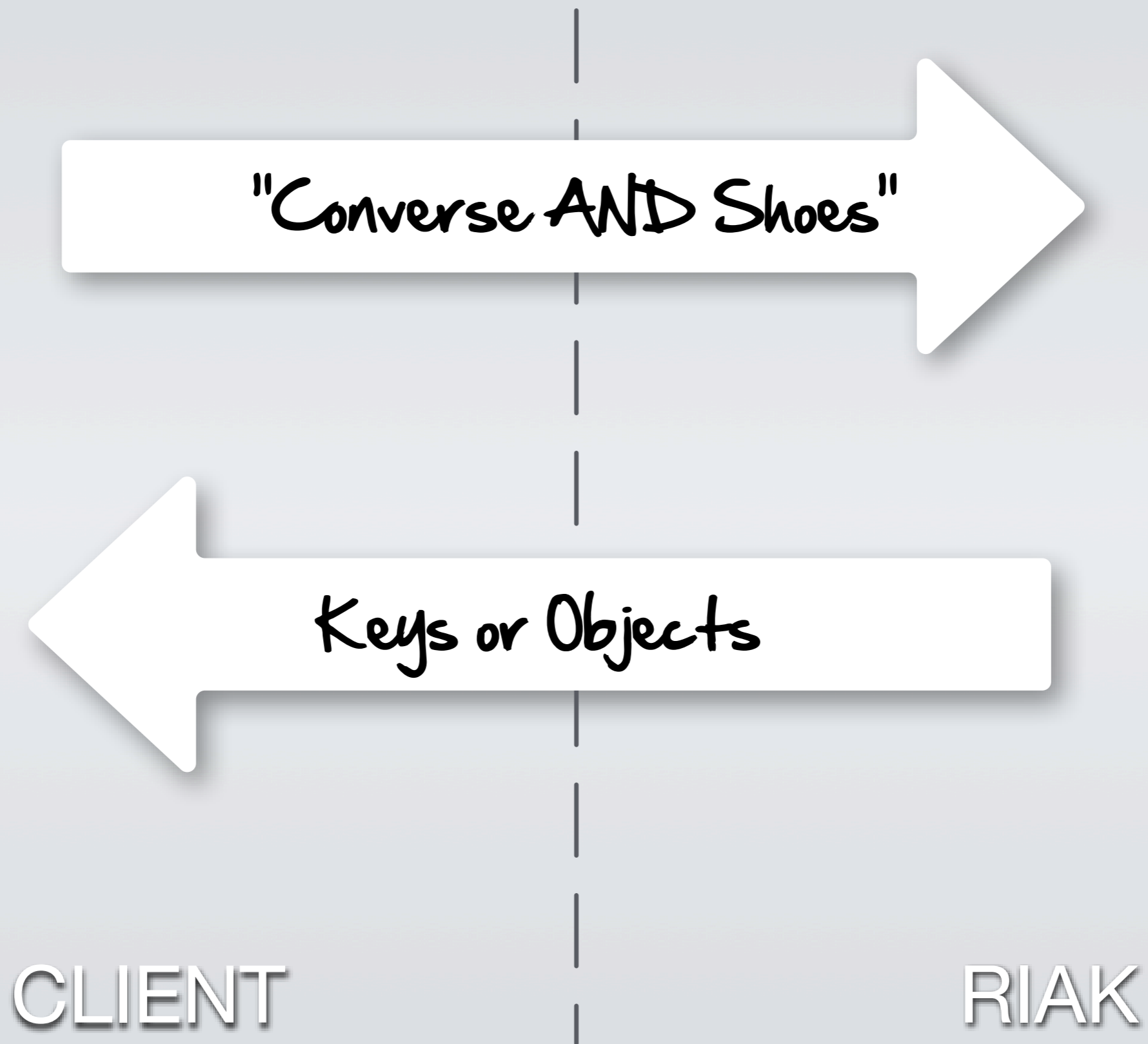
CLIENT                          RIAK

...into this...

"Converse AND Shoes"

Keys or Objects

CLIENT

RIAK

# ...while keeping operations easy.

# Thanks! Questions?

**Search Team:**

John Muellerleile - @jrecursive

Rusty Klophaus - @rklophaus

Kevin Smith - @kevsmith

Currently working with a small set of Beta users.

Open-source release planned for Q3.

www.basho.com