# Using Lua in the Ceph distributed storage system

Noah Watkins

noahwatkins@gmail.com

@noahdesu

Michael Sevilla

mikesevilla3@gmail.com
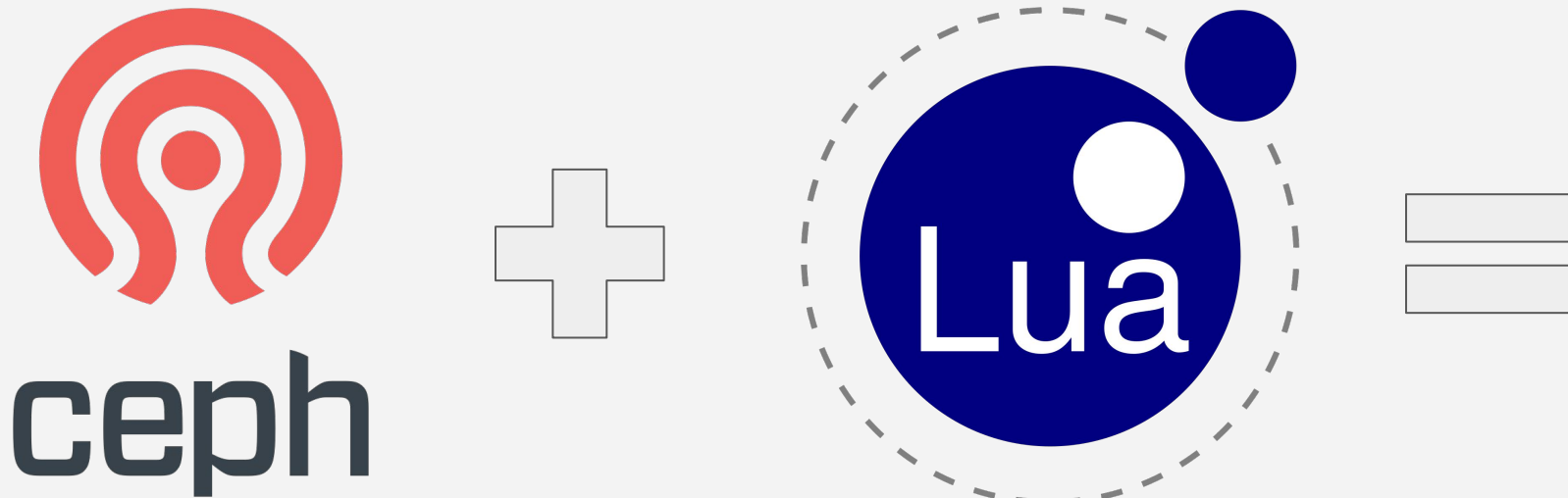
Lua Workshop, October 2017, San Francisco

# yet another Lua embedding



- Ceph is massively **distributed**, highly dynamic **storage system**
- Lua is embedded throughout Ceph
  - Defining policies
  - Domain-specific I/O interfaces
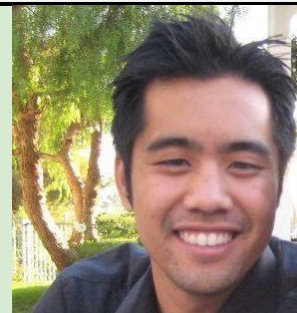  - Remote data processing

# Who we are

## Noah Watkins

- PhD candidate
- UC Santa Cruz
- Experience
  - LANL, HGST
  - IBM, Red Hat
- Contact
  - noahwatkins@gmail.com
  - @noahdesu

- Project Lead
  - Distributed shared-log on SDS
    - https://github.com/noahdesu/zlog
  - **Lua integration in Ceph object store**

## Michael Sevilla

- PhD candidate
- UC Santa Cruz
- Experience
  - HP Enterprise
  - LANL
- Contact
  - mikesevilla3@gmail.com

- Project Lead
  - Mantle: Programmable load balancer in Ceph
  - **Lua integration in Ceph file system**

# Where are we from

UC SANTA CRUZ

SRL
UCSC Systems Research Lab

- Funding from: CR⊘SS | CENTER FOR RESEARCH IN OPEN SOURCE SOFTWARE

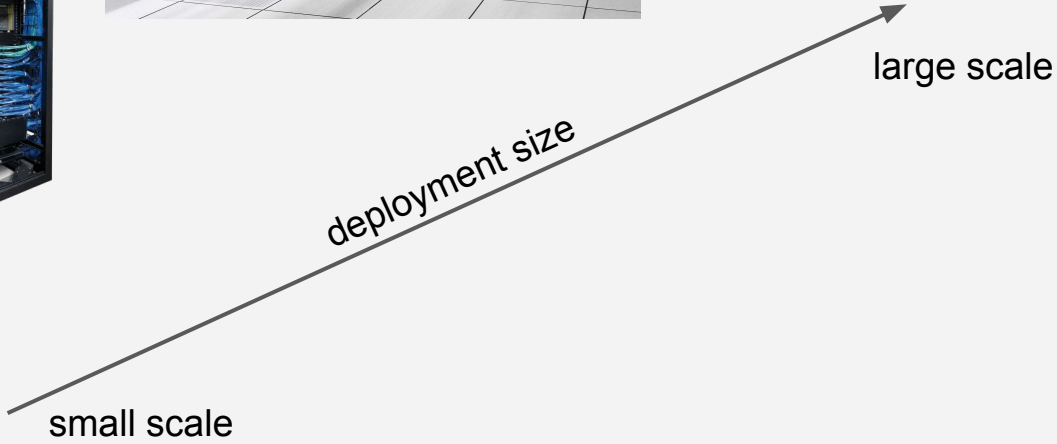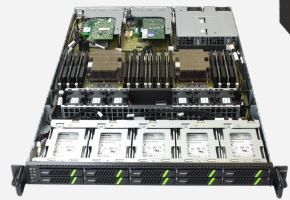| Doug Cutting | Karen Sandler | Sage Weil | Charlie McDowell | Carlos Maltzahn |
|---|---|---|---|---|
| Lucene, Nutch, Avro, Hadoop | GNOME Foundation, SW Freedom Law | WebRing, Dreamhost, Inktank (Ceph) | Associate Dean (Undergrad) @ UCSC | Director, Adjunct Professor |

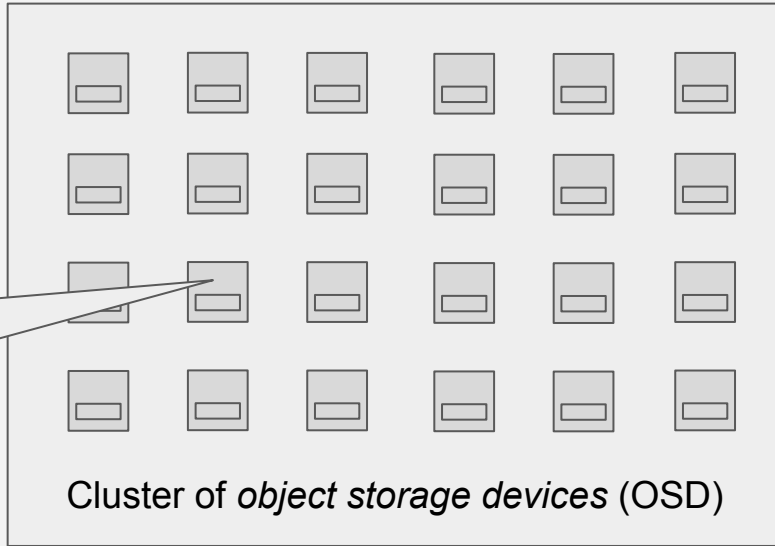- Bridges gap between research & open source
  - Incubator projects, Research projects, College courses
- Proposals considered every 6 months

4

# Ceph is a distributed storage system

- Horizontal scalability
- No single point of failure
- Software controlled
- Hardware agnostic
  - COTS

large scale

deployment size

small scale

# Ceph is a distributed storage system

Cluster of *object storage devices* (OSD)

# Ceph scales horizontally. Just add more nodes!



Cluster of *object storage devices* (OSD)

# Ceph scales horizontally. Just add more nodes!



10x - 1000x

Cluster of *object storage devices* (OSD)

# Ceph works with all sorts of hardware configurations



10x - 1000x

Cluster of *object storage devices* (OSD)

# Ceph is an *object-storage system*



- Cloud
- Image, Video, Music, Etc...

Xen, KVM, VMware

- File systems
- Databases

10x - 1000x

Cluster of *object storage devices* (OSD)

# Ceph is an *object-storage system*



- Cloud
- Image, Video, Music, Etc...

Xen, KVM, VMware

- File systems
- Databases

10x - 1000x

- Peer-to-peer
- Self-healing

Cluster of *object storage devices* (OSD)

# Ceph is an *object-storage system*



Cloud
- Image, Video, Music, Etc...

Xen, KVM, VMware

- File systems
- Databases

- Billions of objects
- Object API

10x - 1000x

- Peer-to-peer
- Self-healing

Cluster of *object storage devices* (OSDs)

12

# Ceph is an *object-storage system*

network

Clients access Ceph via **services**, or directly through the **object API**.

| S3/Swift | Block Device | File System |

**Object API**

Recovery and failures are transparent to clients.

Cluster of *object storage devices* (OSDs)

# Ceph is widely deployed and developed

- 100% open-source
  - Hundreds of contributors
  - Industry, academics, government
- https://github.com/ceph/ceph
- Recently tested at CERN
  - 10,800 OSDs
  - 65 petabytes
- Embedded devices
  - 504 WDLabs converged Ethernet drives (4 PB)
  - ARM CPU, memory, network, storage
  - OSD and disk plug directly into network

# questions?

**<u>Recap</u>**: Ceph 🌀 is a distributed storage system, that:

- scales horizontally. Just add more nodes!
- works with all sorts of hardware configurations
- is an object-storage system
- is widely deployed and developed

# Lua in your objects

# Targeting object APIs can be challenging

**API:**
write(**object**, bytes)
blob = read(**object**)

OSD

S3/Swift

Block Device

File System

Object API

Cluster of *object storage devices* (OSDs)

# Targeting object APIs can be challenging

**API:**
write(**object**, bytes)
blob = read(**object**)

OSD

S3/Swift

**Object API**

Partitioning (sharding) large data over many objects

Cluster of *object storage devices* (OSDs)

# Targeting object APIs can be challenging



**API:**
write(**object**, bytes)
blob = read(**object**)

single object strong consistency

Partitioning (sharding) large data over many objects

Maintaining consistency across objects in a data set

S3/Swift

Object API

OSD

Cluster of *object storage devices* (OSDs)

# A rich object API makes things better

**API:**

| Enumeration | I/O hints | Atomic |
| Watch/Notify | Extended attr. | compounds |
| Sync, Async | Key-value DB | Locking |
| Caching | Snapshots | *Extensions...* |

OSD

S3/Swift

Block Device

File System

**Object API**

Cluster of *object storage devices* (OSDs)

# A rich object API makes things better



**API:**

| | | |
|---|---|---|
| Enumeration | I/O hints | Atomic |
| Watch/Notify | Extended attr. | compounds |
| Sync, Async | Key-value DB | Locking |
| Caching | Snapshots | *Extensions...* |

OSD

S3/Swift

Block Device

File System

Object API

Cluster of *object storage devices* (OSDs)

# A rich object API makes things better

**API:**

| | | |
|---|---|---|
| Enumeration | I/O hints | Atomic compounds |
| Watch/Notify | Extended attr. | |
| Sync, Async | Key-value DB | Locking |
| Caching | Snapshots | *Extensions...* |

OSD

S3/Swift

Block Device

File System

**Object API**

Cluster of *object storage devices* (OSDs)

# A rich object API makes things better

**API:**

| | | |
|---|---|---|
| Enumeration | I/O hints | Atomic compounds |
| Watch/Notify | Extended attr. | |
| Sync, Async | Key-value DB | Locking |
| Caching | Snapshots | Extensions... |

OSD

S3/Swift    Block Device    File System

Object API

Cluster of *object storage devices* (OSDs)

# Object classes *extend* the object API



**API:**

| | | |
|---|---|---|
| Enumeration | I/O hints | Atomic |
| Watch/Notify | Extended attr. | compounds |
| Sync, Async | Key-value DB | Locking |
| Caching | Snapshots | *Extensions...* |

```
void
sha1(bytes in, bytes out)
{
  bytes data;

  int ret = read_obj(data);
  if (ret)
    return ret;

  out = apply_checksum(data);
  return 0;
}
```

*C++ Module*

OSD

S3/Swift

Block Device

File System

**Object API**

Cluster of *object storage devices* (OSDs)

# Object classes *extend* the object API

**API:**

| | | |
|---|---|---|
| Enumeration | I/O hints | Atomic |
| Watch/Notify | Extended attr. | compounds |
| Sync, Async | Key-value DB | Locking |
| Caching | Snapshots | *Extensions...* |

```
void
sha1(bytes in, bytes out)
{
  bytes data;

  int ret = read_obj(data);
  if (ret)
    return ret;

  out = apply_checksum(data);
  return 0;
}
```

*C++ Module*

OSD

S3/Swift | Block Device | File System

**Object API**

Cluster of *object storage devices* (OSDs)

Atomic compound operation
- easy consistency model → all or nothing
- make building blocks: *compare-and-swap*

25

# Applications want to build object classes



**26**

# Object classes exist in a walled garden

**API:**

| | | |
|---|---|---|
| Enumeration | I/O hints | Atomic |
| Watch/Notify | Extended attr. | compounds |
| Sync, Async | Key-value DB | Locking |
| Caching | Snapshots | *Extensions...* |

```
void
sha1(bytes in, bytes out)
{
  bytes data;

  int ret = read_obj(data);
  if (ret)
    return ret;

  out = apply_checksum(data);
  return 0;
}
```

*C++ Module*

OSD

S3/Swift

Block Device

File System

**Object API**

Cluster of *object storage devices* (OSDs)

# Object classes exist in a walled garden

**API:**

Enumeration
Watch/Notify
Sync, Async
Caching

I/O hints
Extended attr.
Key-value DB

Atomic
compounds
Locking

**S3/Swift**

**Block Device**

**File System**

Object API

C++ modules
- Compile against every ABI deployed
- Deploy and version every .so on every machine

```
void
sha1(bytes in, bytes out)
{
  bytes data;

  int ret = read_obj(data);
  if (ret)
    return ret;

  out = apply_checksum(data);
  return 0;
}
```

*C++ Module*

OSD

Cluster of *object storage devices* (OSDs)

# Object classes exist in a walled garden

**API:**

Enumeration    I/O hints     Atomic
Watch/Notify    Extended attr.   compounds
Sync, Async
Caching

S3/Swift    Block Device    File System

C++ modules
- Compile against every ABI deployed
- Deploy and version every .so on every machine

```
void
sha1(bytes in, bytes out)
{
  bytes data;

  int ret = read_obj(data);
  if (ret)
    return ret;

  out = apply_checksum(data);
  return 0;
}
```

*C++ Module*

OSD

Restarts completely blow the cache of the workload. No thanks.

Cluster of *object storage devices* (OSDs)

# Object classes exist in a walled garden

**API:**

Enumeration
Watch/Notify
Sync, Async
Caching

I/O hints
Extended attr.
Key-value DB

Atomic
compounds
Locking

S3/Swift

Block Device

File System

```
void
sha1(bytes in, bytes out)
{
  bytes data;

  int ret = read_obj(data);
  if (ret)
    return ret;

  out = apply_checksum(data);
  return 0;
}
```
*C++ Module*

OSD

C++ modules
- Compile against every ABI deployed
- Deploy and version every .so on every machine

Restarts completely blow the cache of the workload. No thanks.

Have a cool idea for an interface? An app that dynamically changes interfaces? Good luck with RHEL/Debian/Alpine/XYZ packaging!

# Lua object class: tunnel requests through LuaVMs



**API:**

| | | |
|---|---|---|
| Enumeration | I/O hints | Atomic |
| Watch/Notify | Extended attr. | compounds |
| Sync, Async | Key-value DB | Locking |
| Caching | Snapshots | *Extensions...* |

script

Lua

OSD

S3/Swift  Block Device  File System

Object API

Cluster of *object storage devices* (OSDs)

# Example: remotely compute md5 hash

### C++ API

```
int compute_md5(cls_method_context_t hctx, bufferlist *in,
  bufferlist *out)
{
  size_t size;
  int ret = cls_cxx_stat(hctx, &size, NULL);
  if (ret < 0)
    return ret;                Read object

  bufferlist data;
  ret = cls_cxx_read(hctx, 0, size, data);
  if (ret < 0)
    return ret;

  byte digest[AES::BLOCKSIZE];
  MD5().CalculateDigest(digest, (byte*)data.c_str(),
    data.length());

  out->append(digest, sizeof(digest));
  return 0;                     Compute
}                               result
```

### Equivalent Lua

```
???
```

# Example: remotely compute md5 hash

### C++ API

```cpp
int compute_md5(cls_method_context_t hctx, bufferlist *in,
  bufferlist *out)
{
  size_t size;
  int ret = cls_cxx_stat(hctx, &size, NULL);
  if (ret < 0)
    return ret;

  bufferlist data;
  ret = cls_cxx_read(hctx, 0, size, data);
  if (ret < 0)
    return ret;

  byte digest[AES::BLOCKSIZE];
  MD5().CalculateDigest(digest, (byte*)data.c_str(),
    data.length());

  out->append(digest, sizeof(digest));
  return 0;
}
```

*Read object*

*Compute result*

### Equivalent Lua

```lua
local md5 = require 'md5'

function compute_md5(input, output)
  local data = objclass.read()
  output = md5.sumhexa(data)
end
```

# Error handling is easy and robust

## C++ API

```cpp
int compute_md5(cls_method_context_t hctx, bufferlist *in,
  bufferlist *out)
{
  size_t size;
  int ret = cls_cxx_stat(hctx, &size, NULL);
  if (ret < 0)
    return ret;

  bufferlist data;
  ret = cls_cxx_read(hctx, 0, size, data);
  if (ret < 0)
    return ret;

  byte digest[AES::BLOCKSIZE];
  MD5().CalculateDigest(digest, (byte*)data.c_str(),
    data.length());

  out->append(digest, sizeof(digest));
  return 0;
}
```

## Equivalent Lua

```lua
local md5 = require 'md5'

function compute_md5(input, output)
  local data = objclass.read()
  output = md5.sumhexa(data)
end
```

- Low-level errors are usually not caught
- Common error handling patterns
  - transparent in Lua integration
- Sometimes you want the error
  - wrap in *ok, ret, args = pcall(...)*

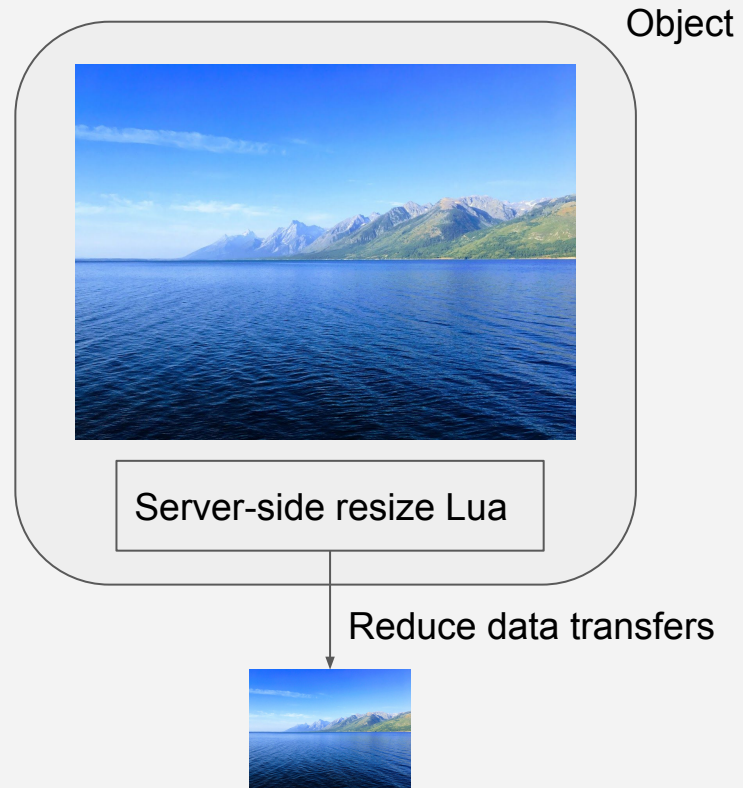# One more example: image thumbnail generation

```
local magick = require "magick"

function thumb(input, output)
  local blob = objclass.read()
  local img = assert(
    magick.load_image_from_blob(blob:str()))

  local spec_string = input:str()
  img = magick.thumb(img, spec_string)

  output:append(img)
end

objclass.register(thumb)
```

Object

Server-side resize Lua

Reduce data transfers

35

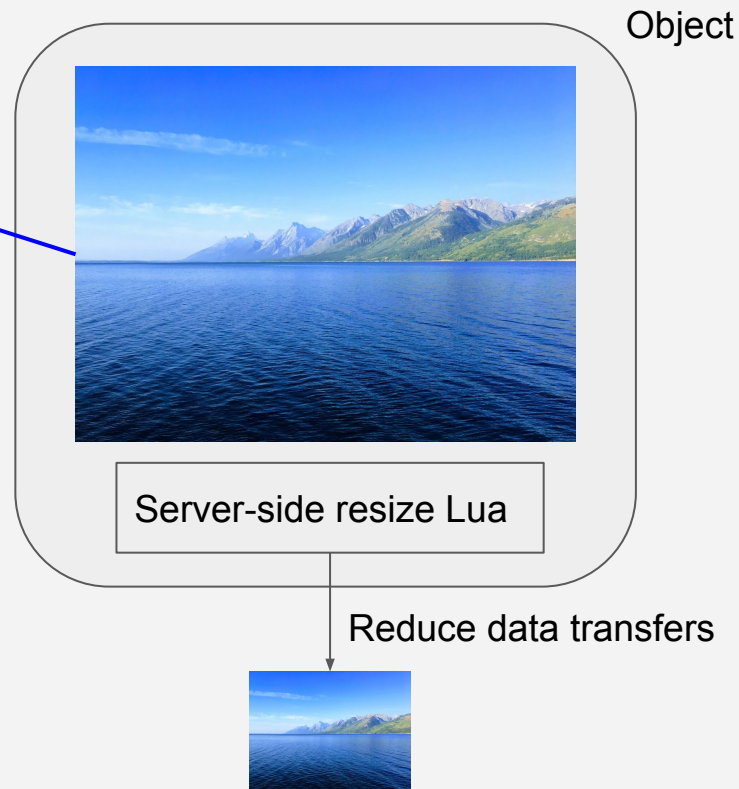# One more example: image thumbnail generation

```lua
local magick = require "magick"

function thumb(input, output)
  local blob = objclass.read()
  local img = assert(
    magick.load_image_from_blob(blob:str()))

  local spec_string = input:str()
  img = magick.thumb(img, spec_string)

  output:append(img)
end

objclass.register(thumb)
```

Object

Server-side resize Lua

Reduce data transfers

# One more example: image thumbnail generation

```
local magick = require "magick"

function thumb(input, output)
  local blob = objclass.read()
  local img = assert(
    magick.load_image_from_blob(blob:str()))

  local spec_string = input:str()
  img = magick.thumb(img, spec_string)

  output:append(img)
end

objclass.register(thumb)
```

Object

Server-side resize Lua

Reduce data transfers

37

# One more example: image thumbnail generation

```
local magick = require "magick"

function thumb(input, output)
    local blob = objclass.read()
    local img = assert(
        magick.load_image_from_blob(blob:str()))

    local spec_string = input:str()
    img = magick.thumb(img, spec_string)

    output:append(img)
end

objclass.register(thumb)
```
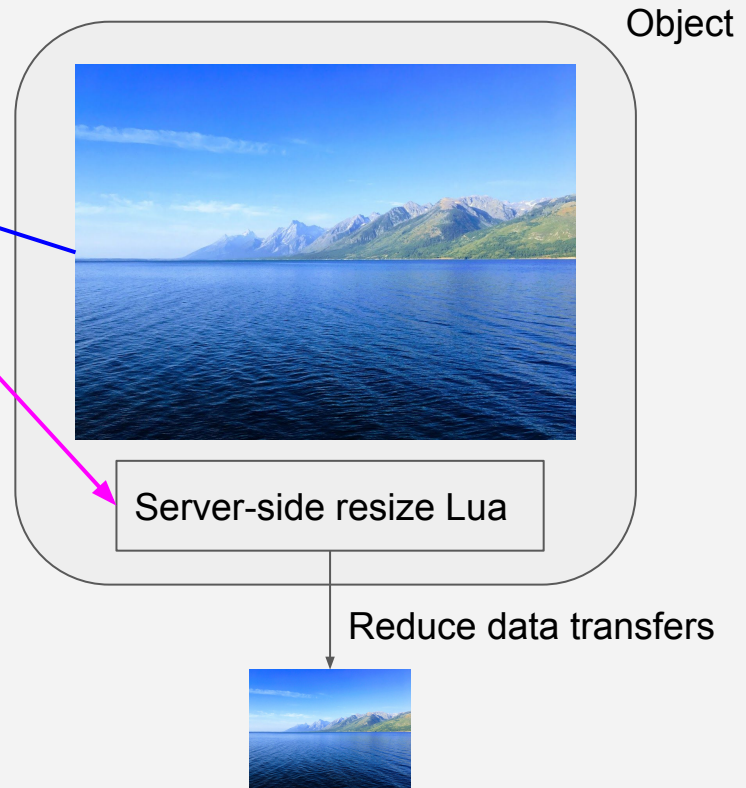
Object



Server-side resize Lua

Reduce data transfers



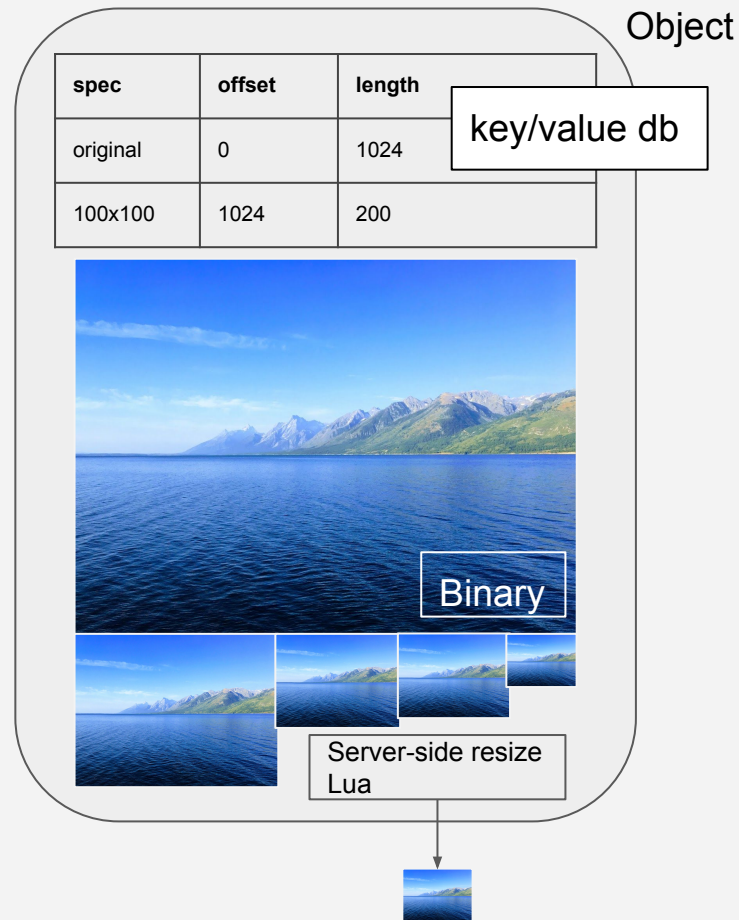Avoid *recomputation* by caching thumbnails!

# Example: image thumbnail generation *with caching*

```lua
local magick = require "magick"

function thumb(input, output)
  local build_thumb = false
  local spec = input:str()
  ok, ret, loc = pcall(objclass.get_map_val, spec)
  if ret == -objclass.ENOENT then
    loc = objclass.get_map_val("original")
    build_thumb = true
  end

  local size, off = string.match(loc:str(), "(%d+)@(%d+)")
  local blob = objclass.read(off, size)
  if not build_thumb then
    output:append(blob)
  else
    img = magick.load_image_from_blob(blob:str()).
    img = magick.thumb(img, spec)

    local obj_size, mtime = objclass.stat()
    loc = #img .. "@" .. obj_size
    objclass.write(off, #img, img)
    objclass.map_set_val(spec, loc)
    output:append(img)
  end
end
objclass.register(thumb)
```
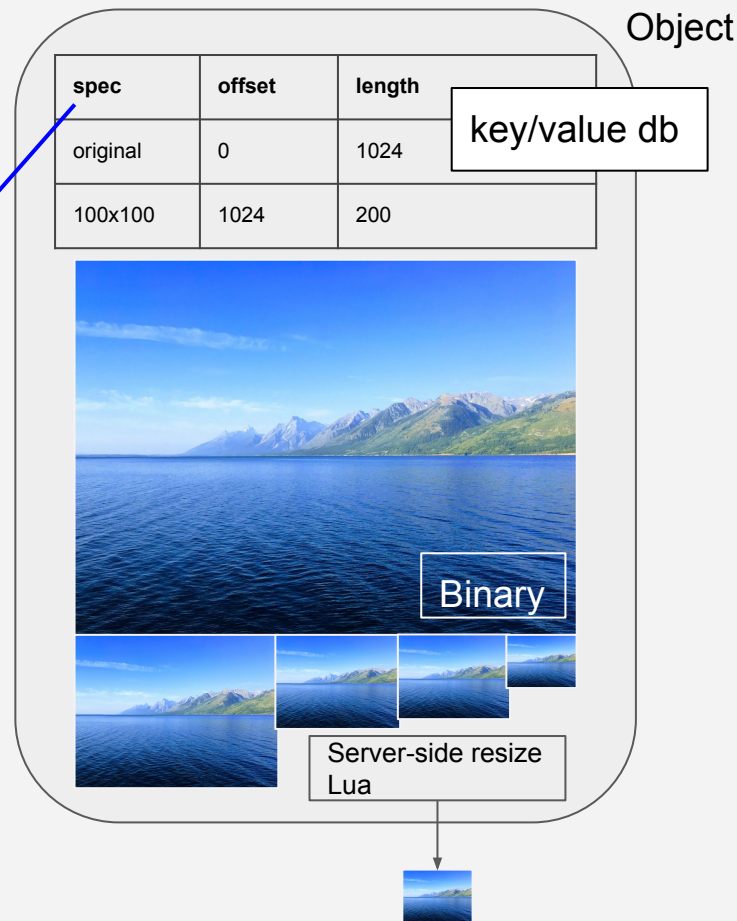
Object

key/value db

| spec | offset | length |
|------|--------|--------|
| original | 0 | 1024 |
| 100x100 | 1024 | 200 |



Binary

Server-side resize
Lua

**39**

# Example: image thumbnail generation *with caching*

```lua
local magick = require "magick"

function thumb(input, output)
  local build_thumb = false
  local spec = input:str()
  ok, ret, loc = pcall(objclass.get_map_val, spec)
  if ret == -objclass.ENOENT then
    loc = objclass.get_map_val("original")
    build_thumb = true
  end

  local size, off = string.match(loc:str(), "(%d+)@(%d+)")
  local blob = objclass.read(off, size)
  if not build_thumb then
    output:append(blob)
  else
    img = magick.load_image_from_blob(blob:str()).
    img = magick.thumb(img, spec)

    local obj_size, mtime = objclass.stat()
    loc = #img .. "@" .. obj_size
    objclass.write(off, #img, img)
    objclass.map_set_val(spec, loc)
    output:append(img)
  end
end
objclass.register(thumb)
```

Object

| spec | offset | length |
|------|--------|--------|
| original | 0 | 1024 |
| 100x100 | 1024 | 200 |

key/value db

Binary

Server-side resize
Lua

# Example: image thumbnail generation *with caching*

```lua
local magick = require "magick"

function thumb(input, output)
  local build_thumb = false
  local spec = input:str()
  ok, ret, loc = pcall(objclass.get_map_val, spec)
  if ret == -objclass.ENOENT then
    loc = objclass.get_map_val("original")
    build_thumb = true
  end

  local size, off = string.match(loc:str(), "(%d+)@(%d+)")
  local blob = objclass.read(off, size)
  if not build_thumb then
    output:append(blob)
  else
    img = magick.load_image_from_blob(blob:str()).
    img = magick.thumb(img, spec)

    local obj_size, mtime = objclass.stat()
    loc = #img .. "@" .. obj_size
    objclass.write(off, #img, img)
    objclass.map_set_val(spec, loc)
    output:append(img)
  end
end
objclass.register(thumb)
```
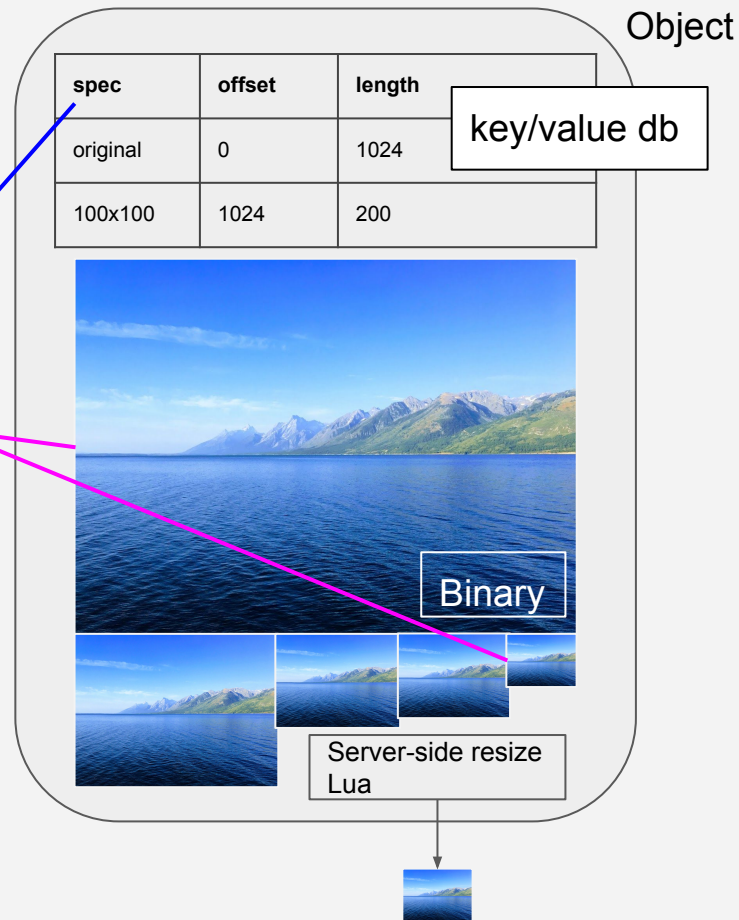
Object

| spec | offset | length |
|------|--------|--------|
| original | 0 | 1024 |
| 100x100 | 1024 | 200 |

key/value db

Binary

Server-side resize
Lua

41

# Example: image thumbnail generation *with caching*

```lua
local magick = require "magick"

function thumb(input, output)
  local build_thumb = false
  local spec = input:str()
  ok, ret, loc = pcall(objclass.get_map_val, spec)
  if ret == -objclass.ENOENT then
    loc = objclass.get_map_val("original")
    build_thumb = true
  end

  local size, off = string.match(loc:str(), "(%d+)@(%d+)")
  local blob = objclass.read(off, size)
  if not build_thumb then
    output:append(blob)
  else
    img = magick.load_image_from_blob(blob:str()).
    img = magick.thumb(img, spec)

    local obj_size, mtime = objclass.stat()
    loc = #img .. "@" .. obj_size
    objclass.write(off, #img, img)
    objclass.map_set_val(spec, loc)
    output:append(img)
  end
end
objclass.register(thumb)
```
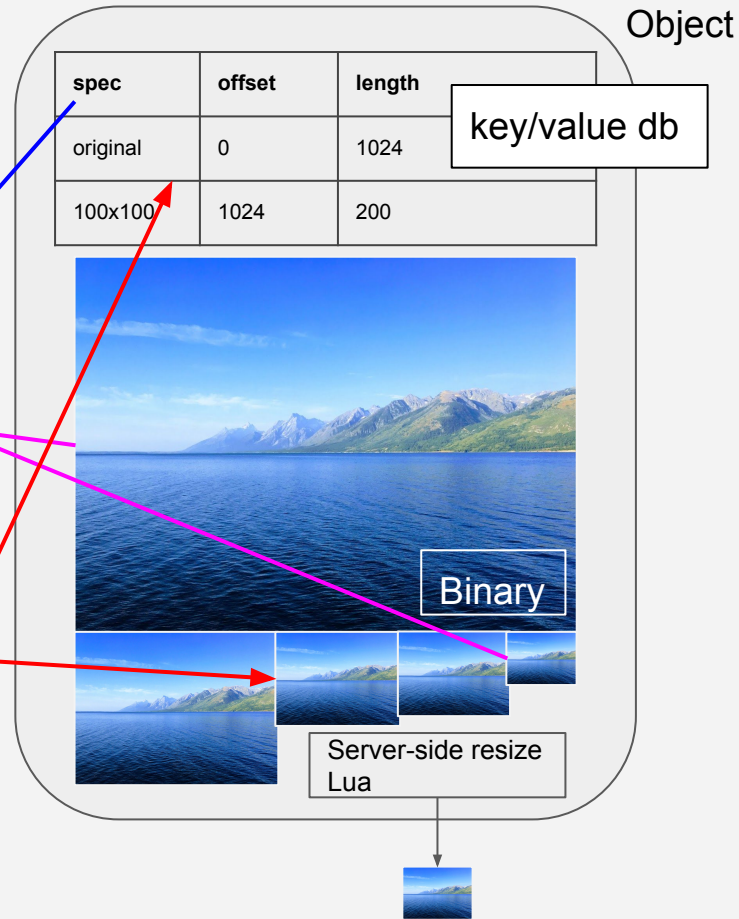
Object

key/value db

| spec | offset | length |
|------|--------|--------|
| original | 0 | 1024 |
| 100x100 | 1024 | 200 |

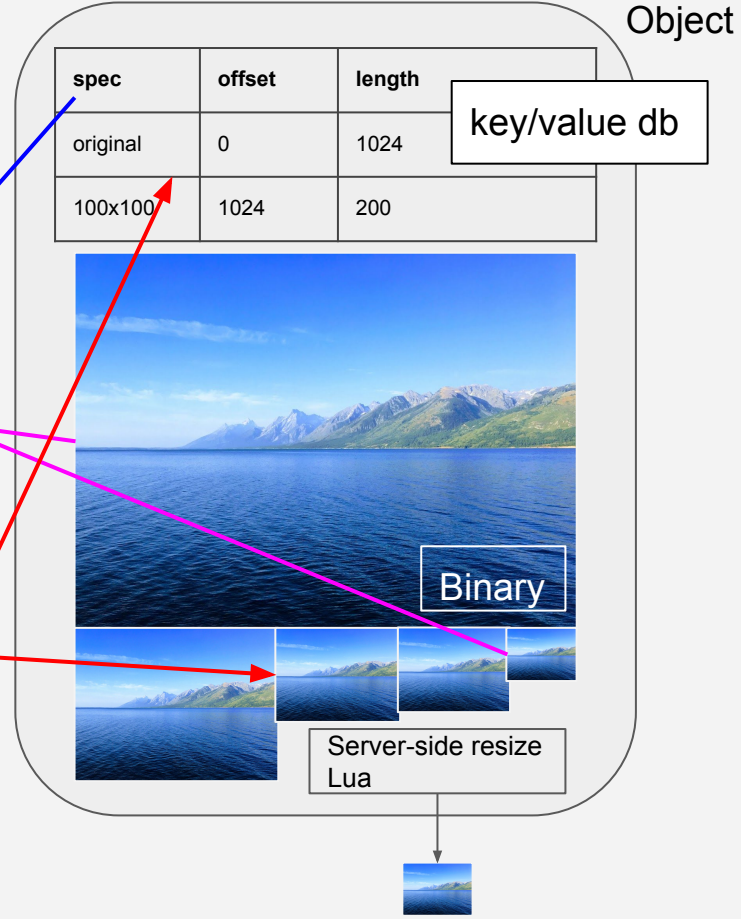Binary

Server-side resize
Lua

**42**

# Example: image thumbnail generation *with caching*
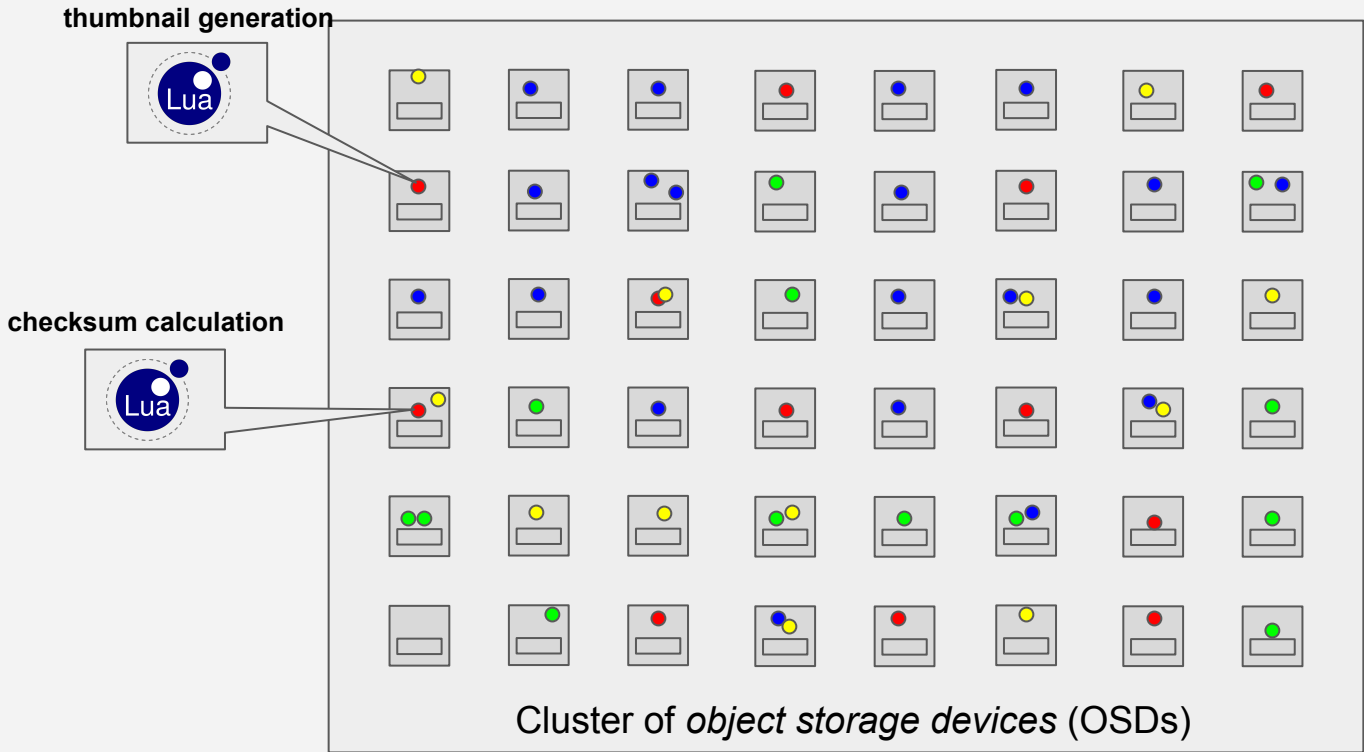
```lua
local magick = require "magick"

function thumb(input, output)
  local build_thumb = false
  local spec = input:str()
  ok, ret, loc = pcall(objclass.get_map_val, spec)
  if ret == -objclass.ENOENT then
    loc = objclass.get_map_val("original")
    build_thumb = true
  end

  local size, off = string.match(loc:str(), "(%d+)@(%d+)")
  local blob = objclass.read(off, size)
  if not build_thumb then
    output:append(blob)
  else
    img = magick.load_image_from_blob(blob:str()).
    img = magick.thumb(img, spec)

    local obj_size, mtime = objclass.stat()
    loc = #img .. "@" .. obj_size
    objclass.write(off, #img, img)
    objclass.map_set_val(spec, loc)
    output:append(img)
  end
end
objclass.register(thumb)
```

Object

key/value db

| spec | offset | length |
|------|--------|--------|
| original | 0 | 1024 |
| 100x100 | 1024 | 200 |

Binary

Server-side resize
Lua

43

# Control billions of objects with Lua



thumbnail generation

checksum calculation

Cluster of *object storage devices* (OSDs)

# Control billions of objects with Lua



data reorganization

thumbnail generation

checksum calculation

app-specific metadata

regex filtering / search

reductions (sum, avg)

Cluster of *object storage devices* (OSDs)

45

# Control billions of objects with Lua

nginx, Luvit, REST APIs, lambdas, etc...

thumbnail generation

data reorganization

regex filtering / search

checksum calculation

reductions (sum, avg)

app-specific metadata

Cluster of *object storage devices* (OSDs)

# Shipping in upstream Ceph as of March '17 (Kraken)

- Available in official packages
  - \>= Kraken release
- Expand API coverage as needed
- Use your favorite Ceph clients
  - C/C++, Python, Lua
  - Java, Rust, Go
- Input methods (Lua + Input Data)
  - C/C++ has efficient data handling
    - Native encoding
  - JSON for compatability

*API:*

| | | |
|---|---|---|
| Enumeration | I/O hints | Atomic compounds |
| Watch/Notify | Extended attr. | |
| Sync, Async | Key-value DB | Locking |
| Caching | Snapshots | Read |
| Write | Stat | Create |
| Remove | | |

```
{
  "script": "function echo(input, output) output:append(input:str()); end
              objclass.register(json_echo)",
  "input": "omg it works",
  "handler": "echo"
}
```
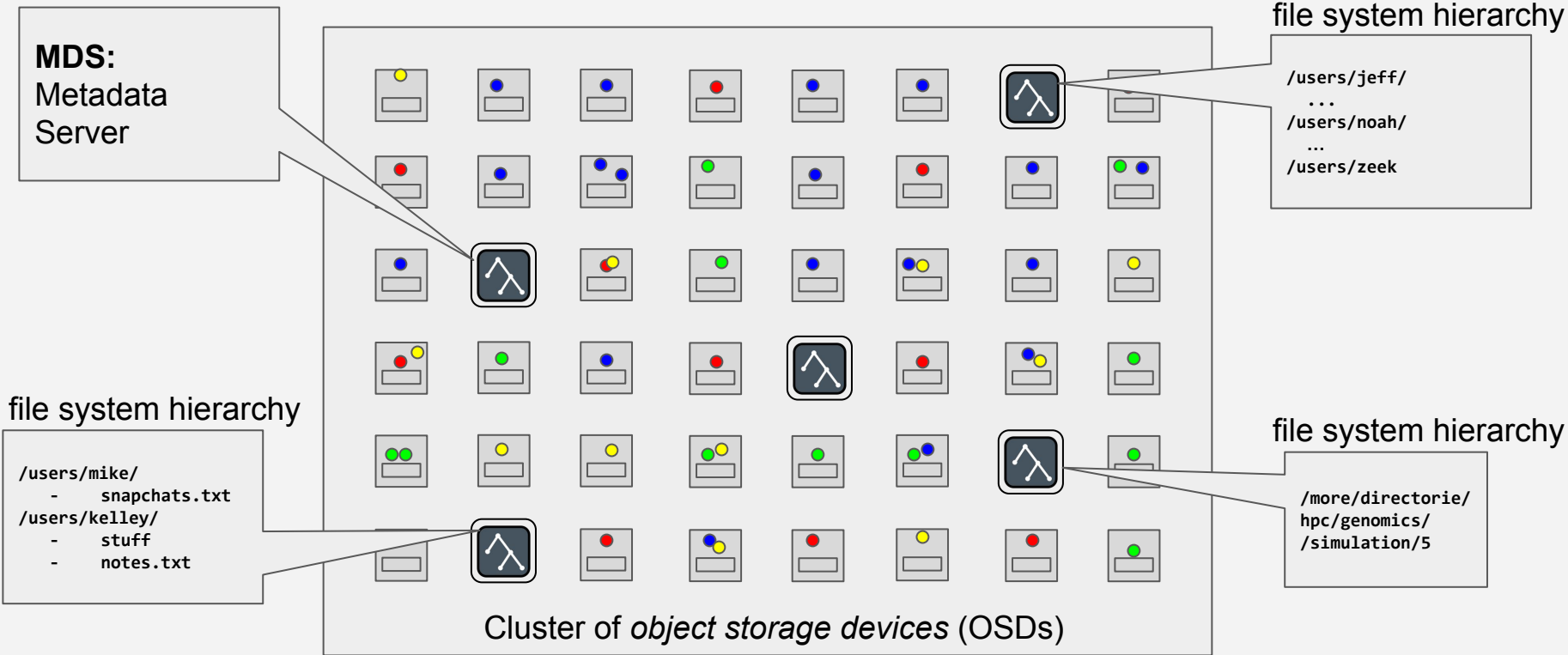
# Lua in your file systems
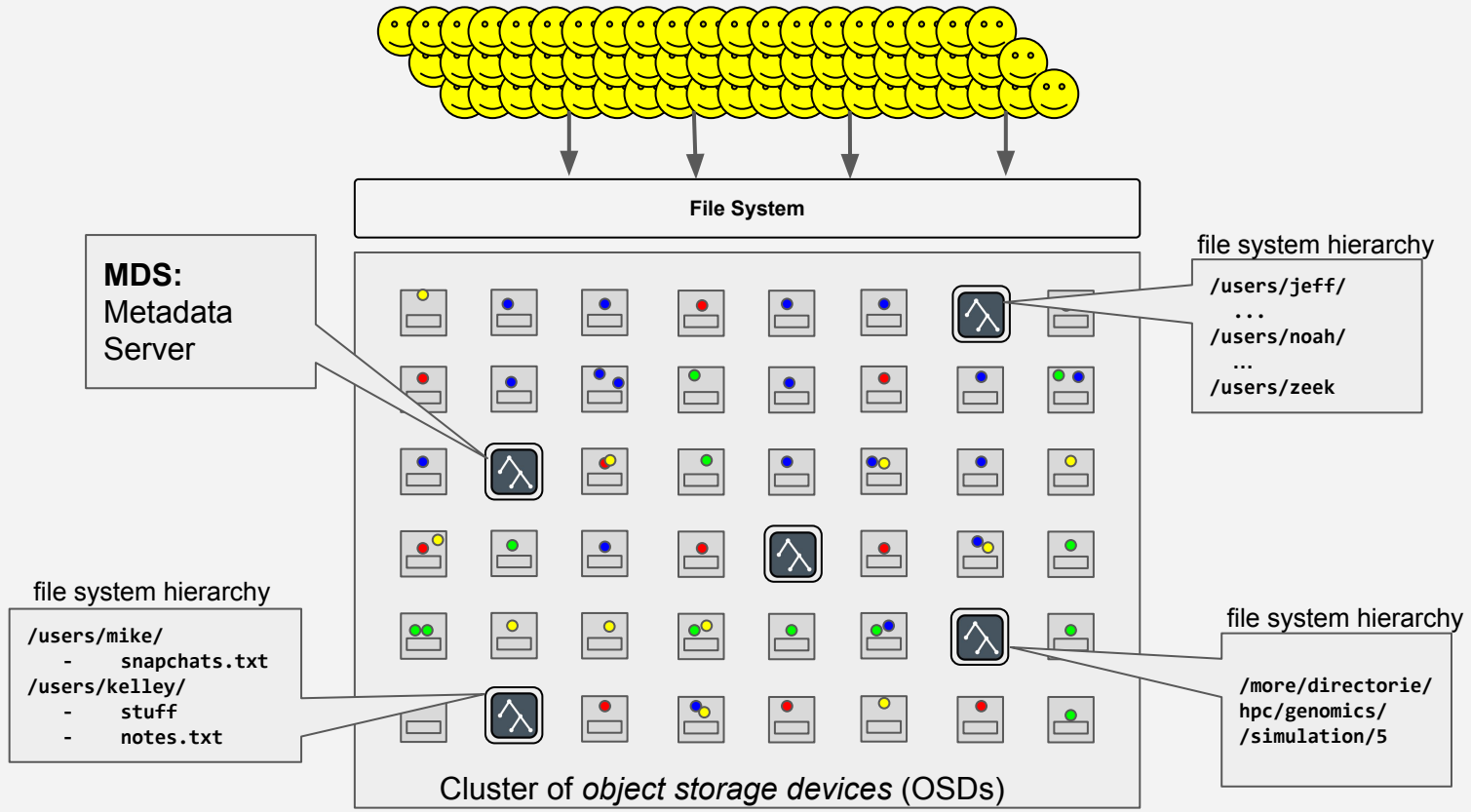
Mantle: programmable metadata load balancer
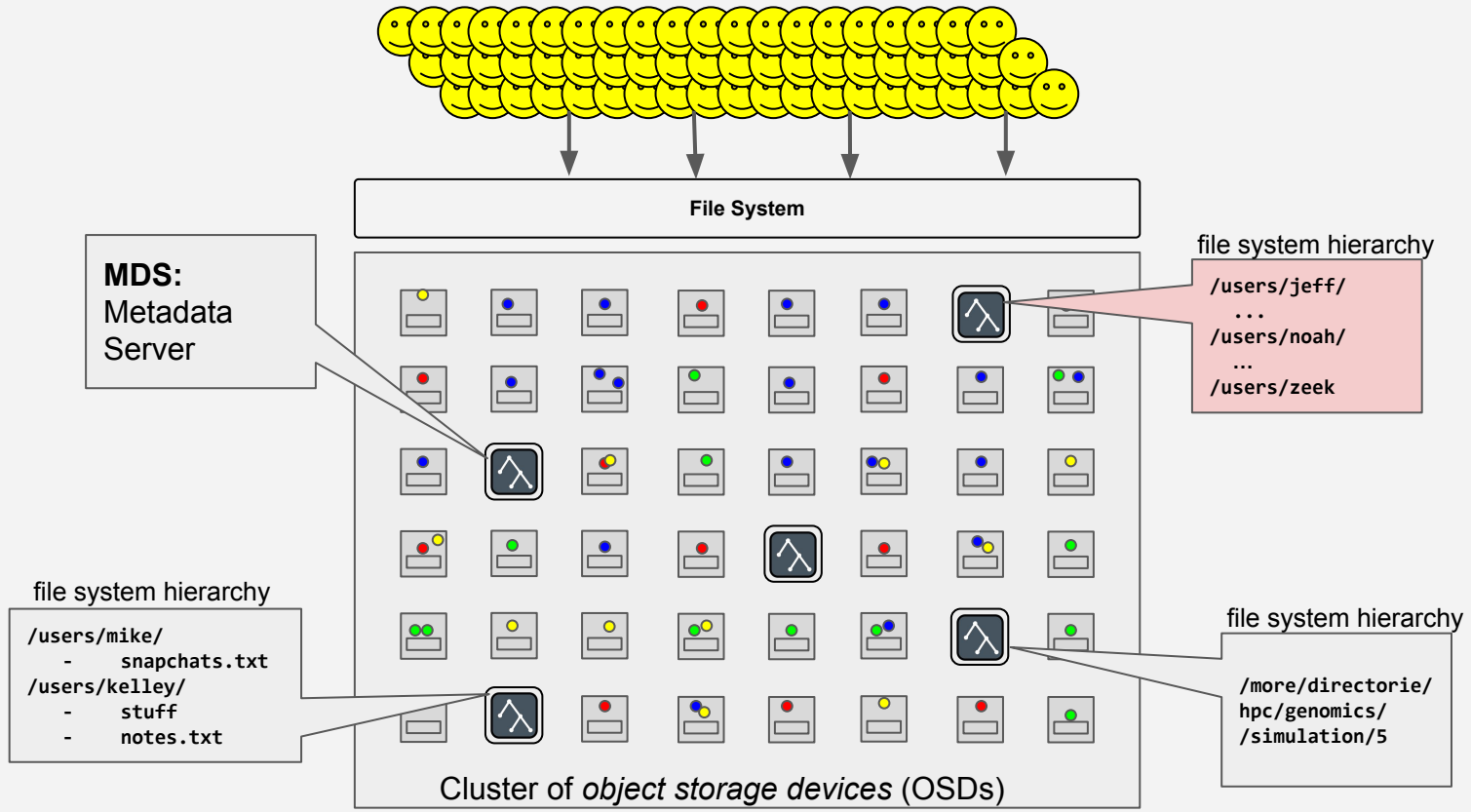
## Michael Sevilla
– http://programmability.us

# Ceph file system metadata load balancing



**MDS:**
Metadata
Server

file system hierarchy

```
/users/jeff/
    ...
/users/noah/
    ...
/users/zeek
```

file system hierarchy

```
/users/mike/
    -    snapchats.txt
/users/kelley/
    -    stuff
    -    notes.txt
```

file system hierarchy

```
/more/directorie/
hpc/genomics/
/simulation/5
```

Cluster of *object storage devices* (OSDs)

**49**

# Ceph file system metadata load balancing



**MDS:**
Metadata
Server

**File System**

file system hierarchy

```
/users/jeff/
    ...
/users/noah/
    ...
/users/zeek
```

file system hierarchy

```
/users/mike/
    -    snapchats.txt
/users/kelley/
    -    stuff
    -    notes.txt
```

file system hierarchy

```
/more/directorie/
hpc/genomics/
/simulation/5
```

Cluster of *object storage devices* (OSDs)

# Ceph file system metadata load balancing



**MDS:** Metadata Server

**File System**

file system hierarchy

```
/users/jeff/
    ...
/users/noah/
    ...
/users/zeek
```

file system hierarchy

```
/users/mike/
    -    snapchats.txt
/users/kelley/
    -    stuff
    -    notes.txt
```

file system hierarchy

```
/more/directorie/
hpc/genomics/
/simulation/5
```

Cluster of *object storage devices* (OSDs)

# Ceph file system metadata load balancing



**MDS:** Metadata Server

file system hierarchy

```
/users/jeff/
    ...
/users/noah/
    ...
/users/zeek
```

file system hierarchy

```
/users/mike/
    -    snapchats.txt
/users/kelley/
    -    stuff
    -    notes.txt
```

file system hierarchy

```
/more/directorie/
hpc/genomics/
/simulation/5
```

Cluster of *object storage devices* (OSDs)

# Ceph file system metadata load balancing



**MDS:** Metadata Server

file system hierarchy

```
/users/jeff/
    ...
/users/noah/
    ...
/users/zeek
```

**Rebalance:** When and what to migrate?

file system hierarchy

```
/users/mike/
    -    snapchats.txt
/users/kelley/
    -    stuff
    -    notes.txt
```

file system hierarchy

```
/more/directorie/
hpc/genomics/
/simulation/5
```

Cluster of *object storage devices* (OSDs)

# Ceph file system metadata load balancing

# Ceph file system metadata load balancing



**MDS:**
Metadata
Server

**Rebalance:**

Lua

file system hierarchy

```
/users/jeff/
    ...
/users/noah/
    ...
/users/zeek
```

file system hierarchy

```
/users/mike/
    -    snapchats.txt
/users/kelley/
    -    stuff
    -    notes.txt
```

file system hierarchy

```
/more/directorie/
hpc/genomics/
/simulation/5
```

**File System**

Cluster of *object storage devices* (OSDs)

# Mantle load balancing policy example

**Policy model**



**Example**

```lua
local function when()
    if      server[whoami]['load']>0.1 and
            server[whoami + 1]['load']<0.1 then
            return true
    end
    return false
end
```

EXAMPLE!

# Lua in your databases

SkyhookDB: leveraging programmable storage toward DB elasticity

Jeff LeFevre and Noah Watkins
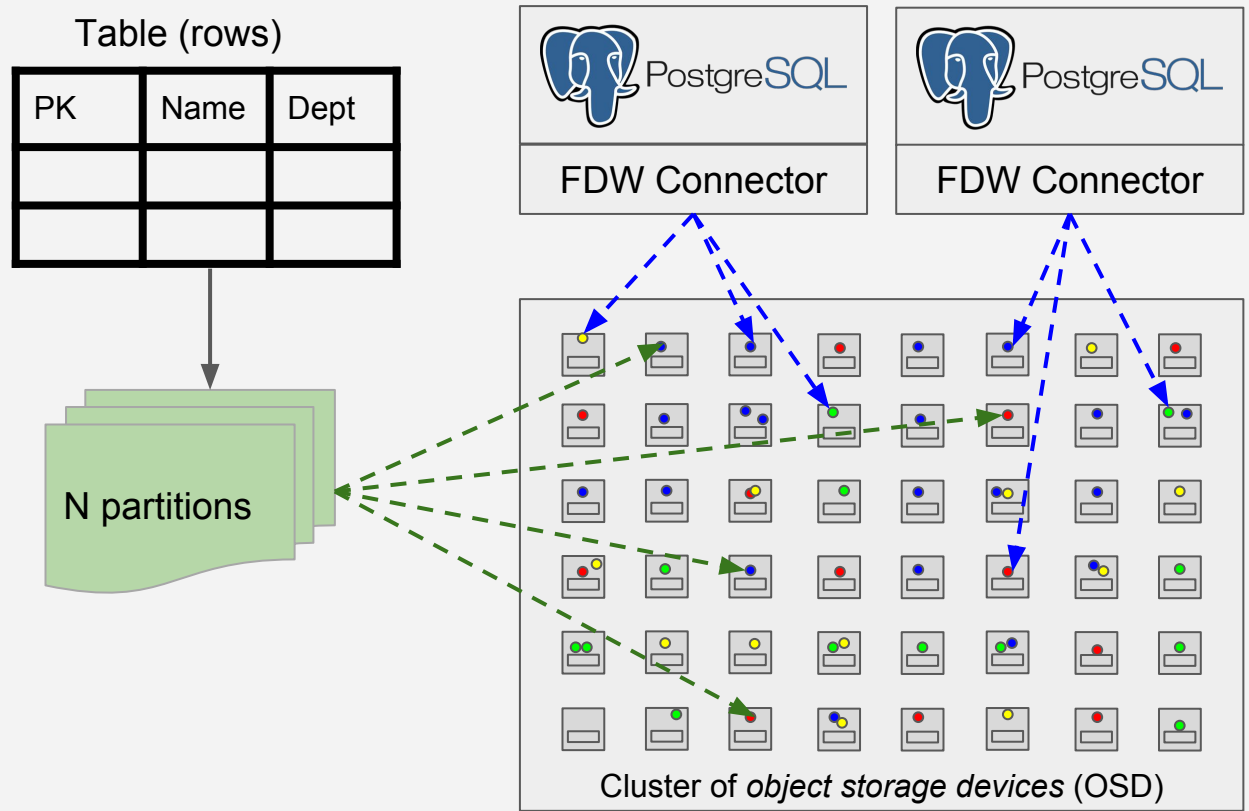- – https://sites.google.com/site/skyhookdb/
- – https://cross.ucsc.edu
- – jlefevre@ucsc.edu

# Skyhook Storage Engine (Database System)



Table (rows)

| PK | Name | Dept |
|----|------|------|
|    |      |      |
|    |      |      |

N partitions

Cluster of *object storage devices* (OSD)

# Skyhook Storage Engine (Database System)



Table (rows)

| PK | Name | Dept |
|---|---|---|
|  |  |  |
|  |  |  |

N partitions

FDW Connector

FDW Connector

Cluster of *object storage devices* (OSD)

# Skyhook Storage Engine (Database System)



Table (rows)

| PK | Name | Dept |
|----|------|------|
|    |      |      |
|    |      |      |

N partitions

FDW Connector

FDW Connector

Cluster of *object storage devices* (OSD)

- Scale-out
  - CPU
  - I/O
- Push-downs on steroids
- Predicate evaluation
- Regex search
- Projection

# Skyhook Storage Engine (Database System)



Table (rows)

| PK | Name | Dept |
|----|------|------|
|    |      |      |
|    |      |      |

N partitions

Cluster of *object storage devices* (OSD)

- Scale-out
  - CPU
  - I/O
- Push-downs on steroids
- Predicate evaluation
- Regex search
- Projection

# Recapping

- Ceph distributed storage system
  - Massively scalable
  - Object storage system
- Embedding Lua
  - Flexible policy injection
  - Application I/O interfaces
  - Remote computation
- Resources
  - https://github.com/ceph/ceph
  - https://ceph.com/rados/dynamic-object-interfaces-with-lua/
    - Basic tutorial

# What we'd like to do next

- We aren't Lua or programming language experts
  - Our APIs are simple and mimic the C APIs
  - What are better interfaces for our use cases?
- Script and dependency management (big issue)
  - Currently
    - Scripts sent with request
    - Dependencies at host level
  - Ideas
    - Install scripts "into" the cluster
    - Caching VMs to reduce latency

# Thanks. And questions?

Noah Watkins

noahwatkins@gmail.com

@noahdesu

Michael Sevilla

mikesevilla3@gmail.com