# Gluster Technical Overview

Kaleb KEITHLEY
Red Hat
11 April, 2015

# Concepts:

- Translator — a shared library that implements storage semantics

  - e.g. distribution, replication, performance, etc.

- Translator Stack — a directed graph of translators

- Brick — a server and a directory

- Volume — a collection of bricks

**Kaleb KEITHLEY**

# Gluster Processes

- Server: glusterfsd —  the heart of it all
- Clients:
  - FUSE, glusterfs
  - NFS server, glusterfs (gluster client, NFS server)
- Management: glusterd
- Healing: glustershd
- And more——

**Kaleb KEITHLEY**

# Let's get real — creating a volume

```
Srv1% gluster peer probe srv1

...

srv1% gluster volume create my_vol
srv1:/bricks/my_vol

srv1% gluster volume start my_vol


...

client1% mount -t glusterfs srv1:my_vol
/mnt
```
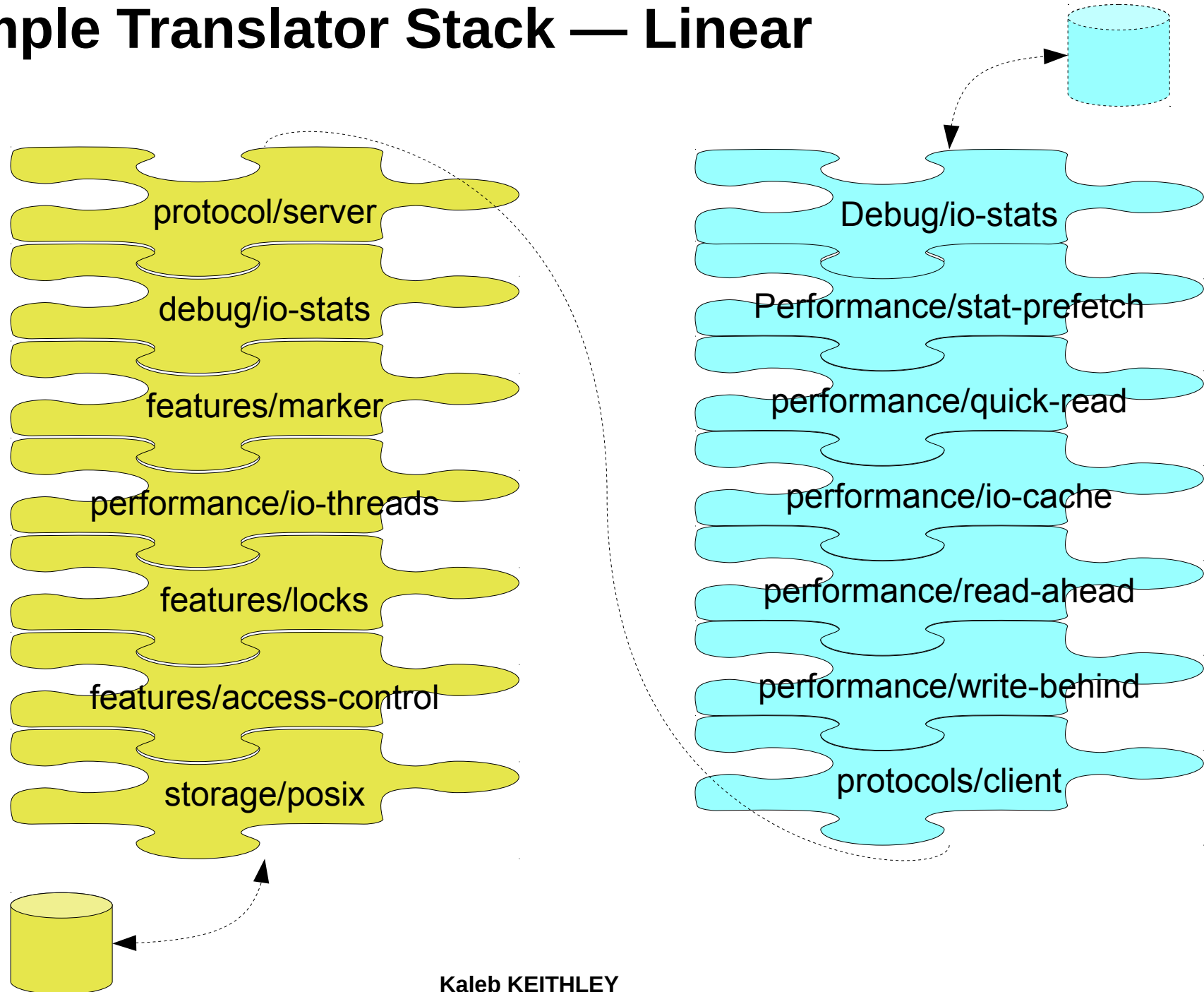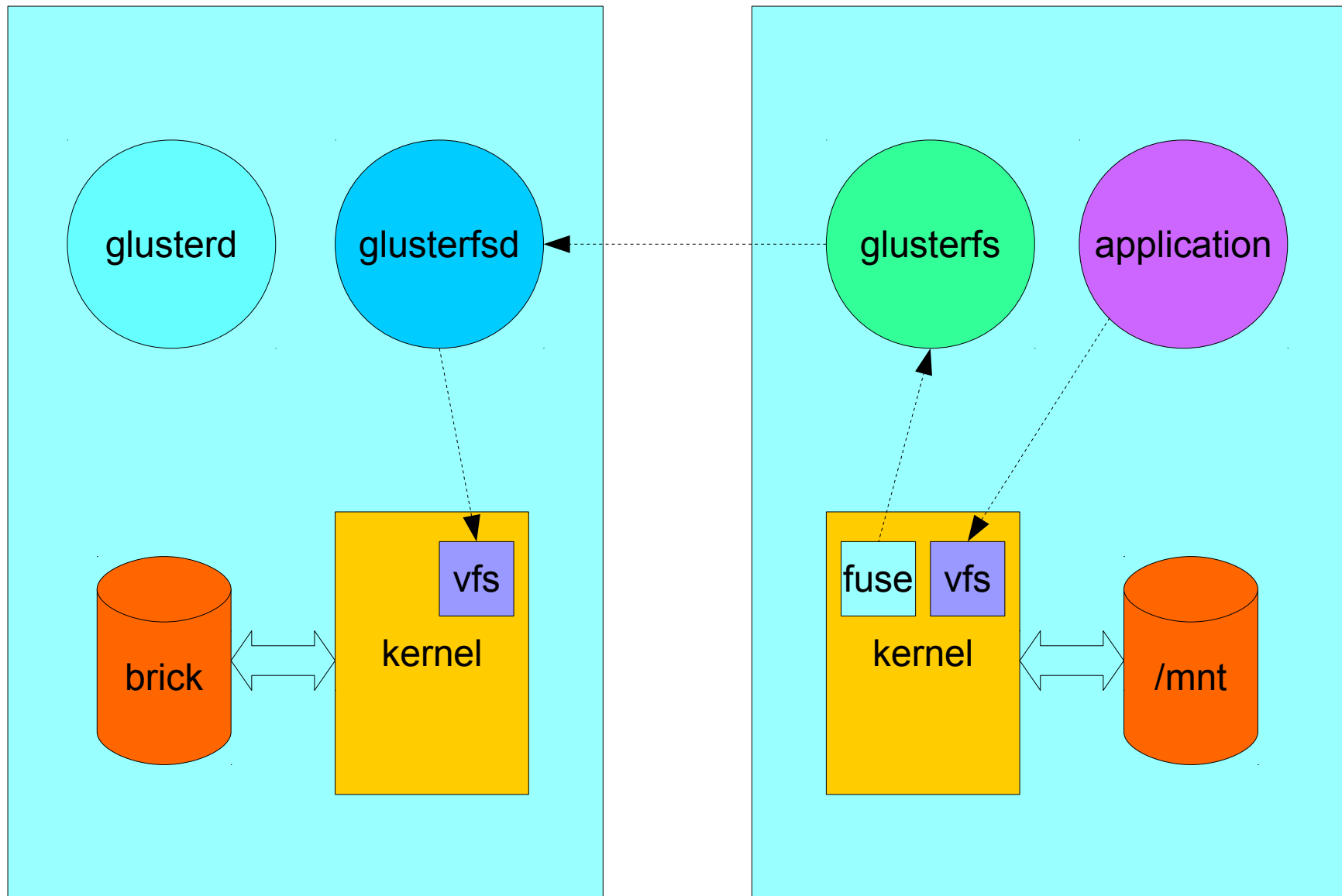
**Kaleb KEITHLEY**

# Simple Translator Stack — Linear

**Left stack (yellow):**

- protocol/server
- debug/io-stats
- features/marker
- performance/io-threads
- features/locks
- features/access-control
- storage/posix

**Right stack (cyan):**

- Debug/io-stats
- Performance/stat-prefetch
- performance/quick-read
- performance/io-cache
- performance/read-ahead
- performance/write-behind
- protocols/client

**Kaleb KEITHLEY**

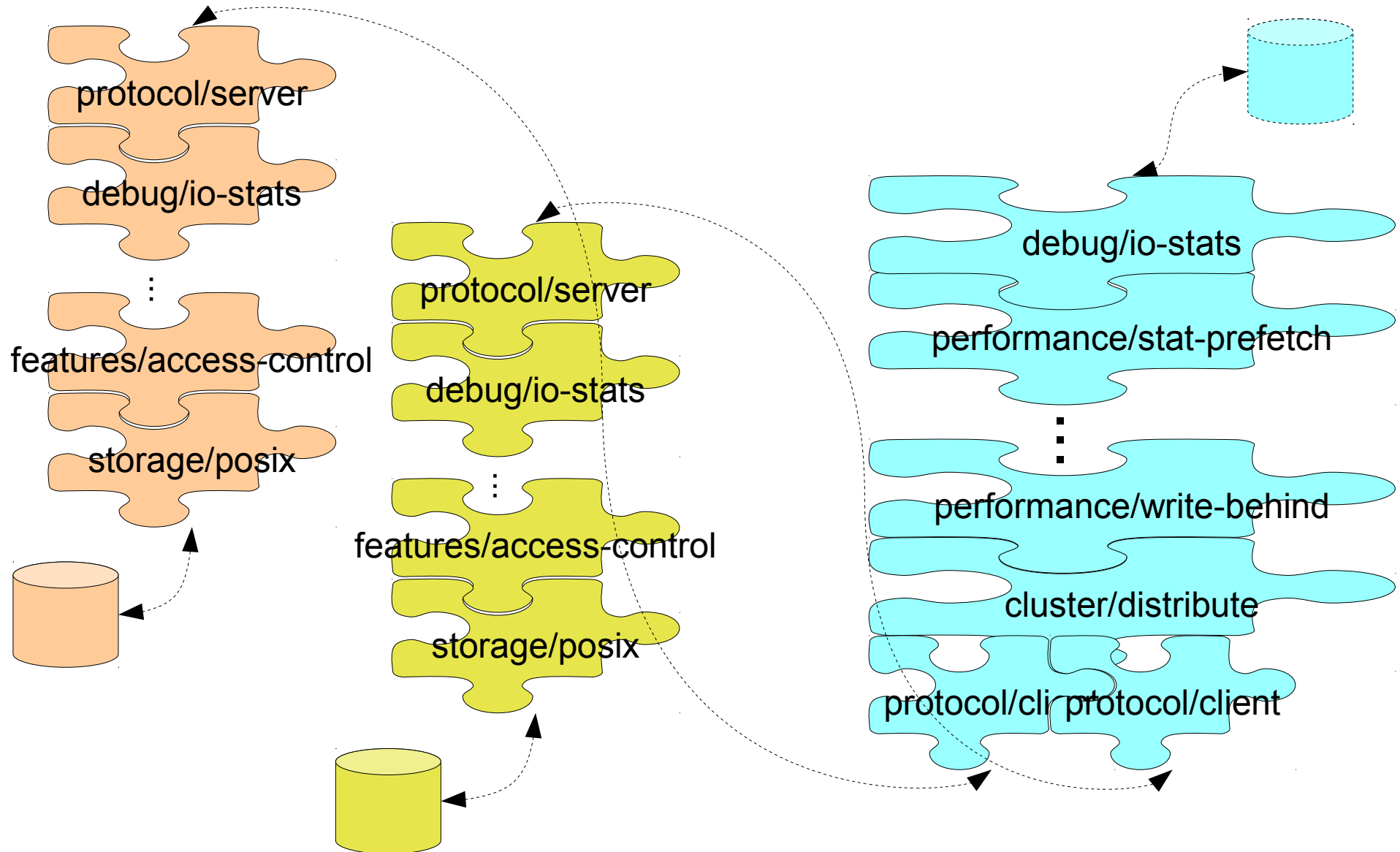**Kaleb KEITHLEY**

# Let's get real — creating a distribute volume

```
srv1% gluster volume create my_dist
srv1:/bricks/my_vol srv2:/bricks/my_vol

srv1% gluster volume start my_dist


...

client1% mount -t glusterfs srv1:my_dist /mnt
```
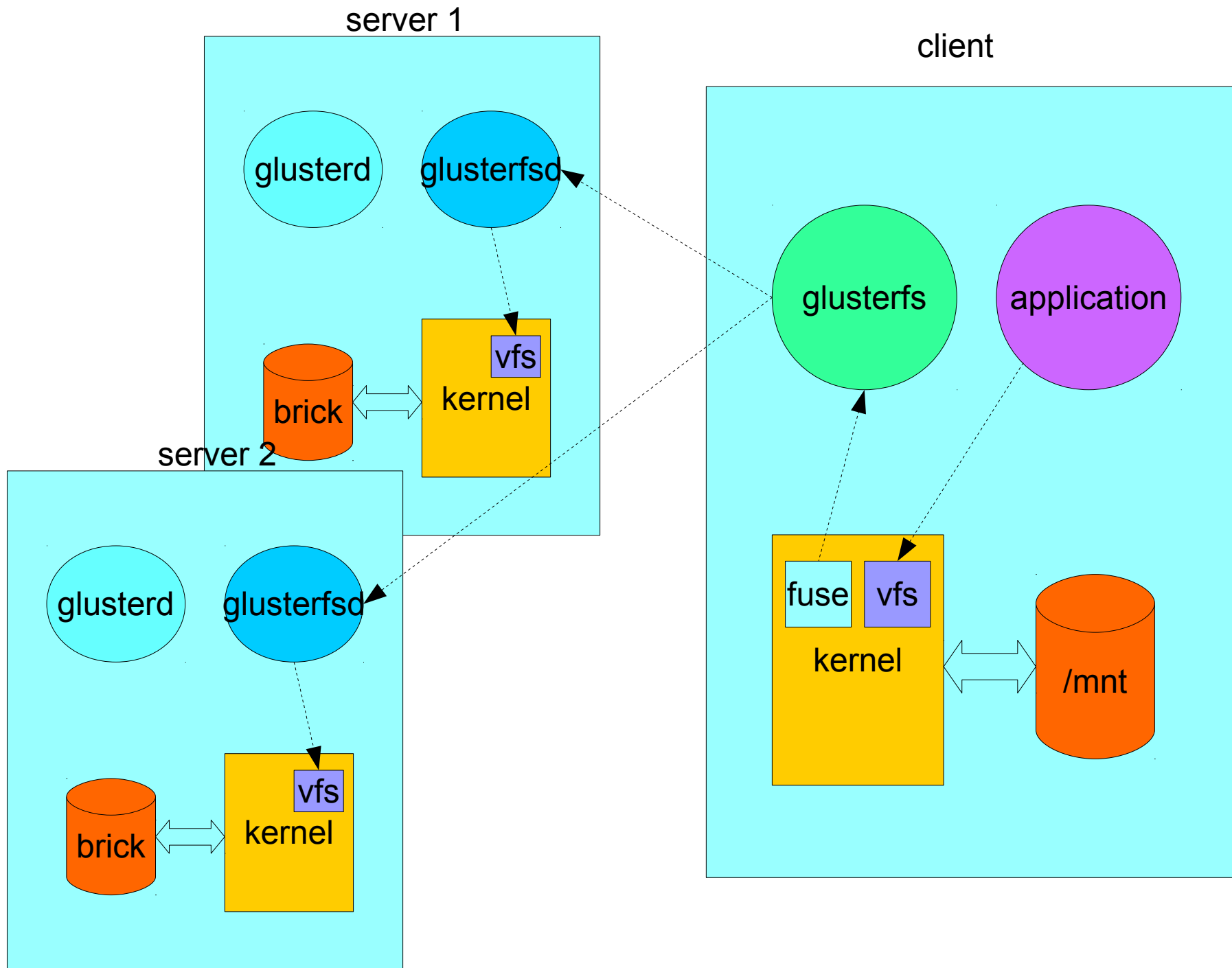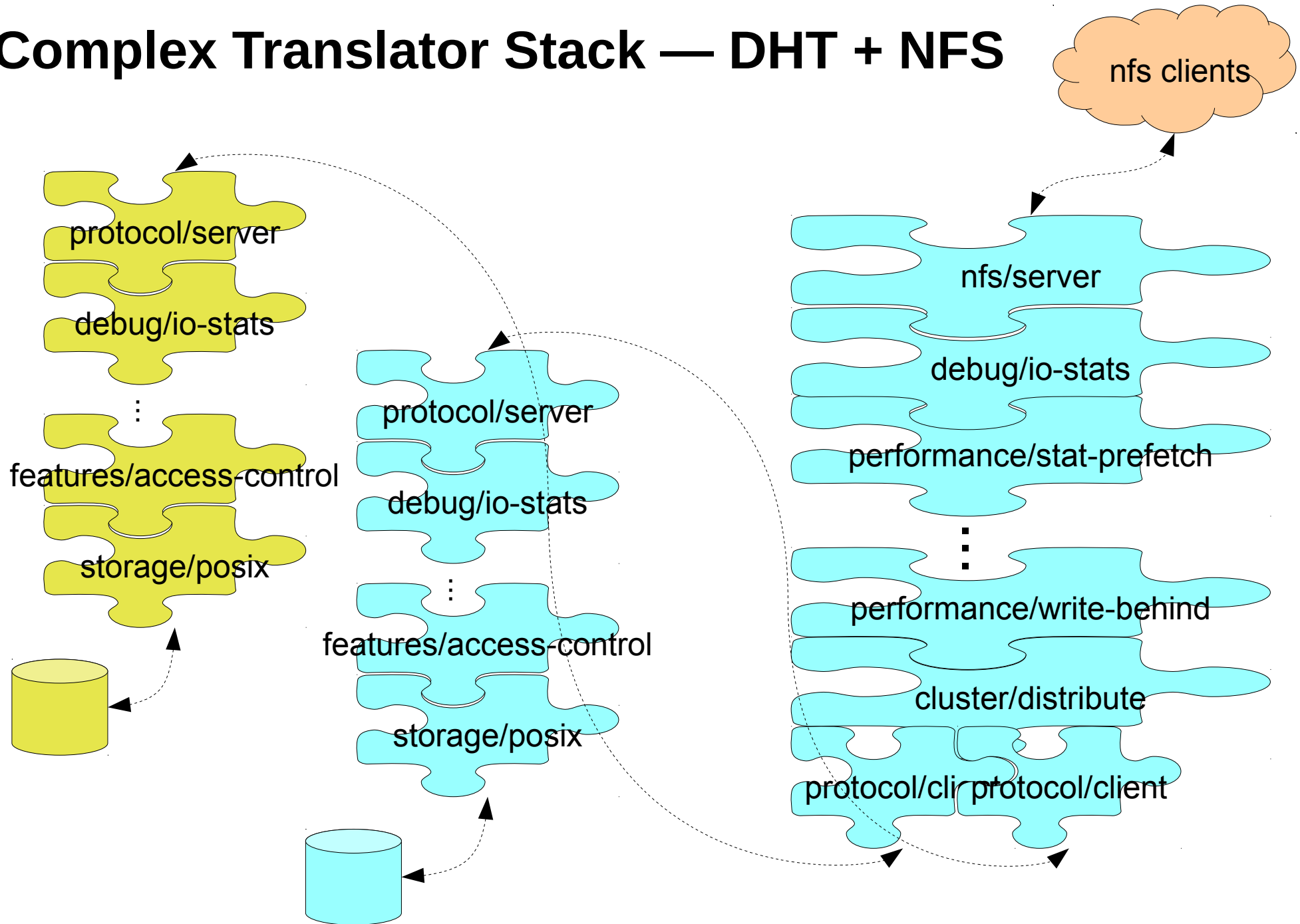
**Kaleb KEITHLEY**

# Complex Translator Stack — Distribute (DHT)



protocol/server

debug/io-stats

⋮

features/access-control

storage/posix

protocol/server

debug/io-stats

⋮

features/access-control

storage/posix

debug/io-stats

performance/stat-prefetch

⋮

performance/write-behind

cluster/distribute

protocol/cli protocol/client

**Kaleb KEITHLEY**

**Kaleb KEITHLEY**

# Complex Translator Stack — DHT + NFS

nfs clients

protocol/server

debug/io-stats

⋮

features/access-control

storage/posix

protocol/server

debug/io-stats

⋮

features/access-control

storage/posix

nfs/server

debug/io-stats

performance/stat-prefetch

⋮

performance/write-behind

cluster/distribute

protocol/client protocol/client

**Kaleb KEITHLEY**

**Kaleb KEITHLEY**

# Let's get real — creating a replica volume

```
srv1% gluster volume create replica 2 my_repl
srv1:/bricks/my_vol srv2:/bricks/my_vol

srv1% gluster volume start my_repl


...

client1% mount -t glusterfs srv1:my_repl /mnt
```
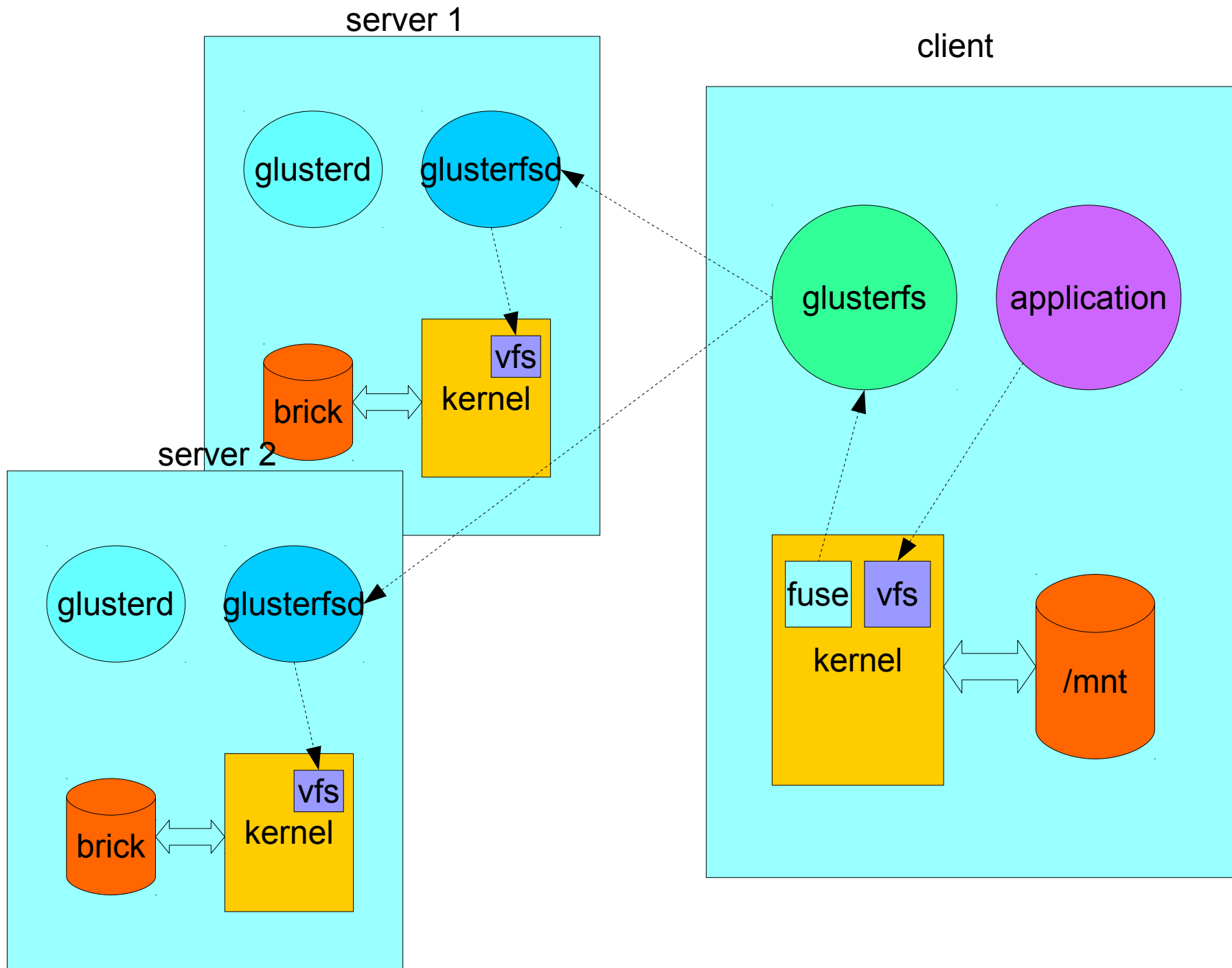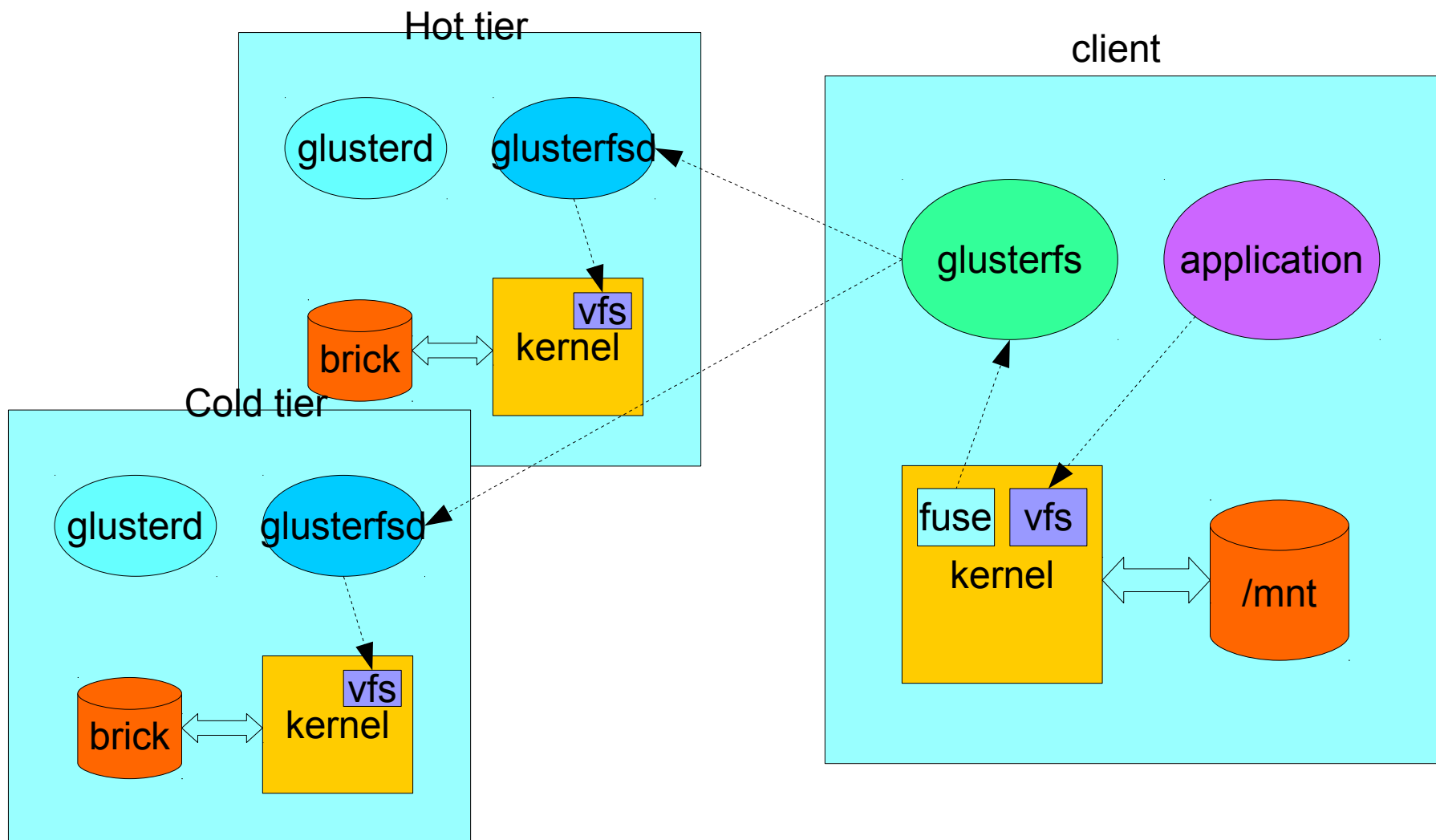
**Kaleb KEITHLEY**
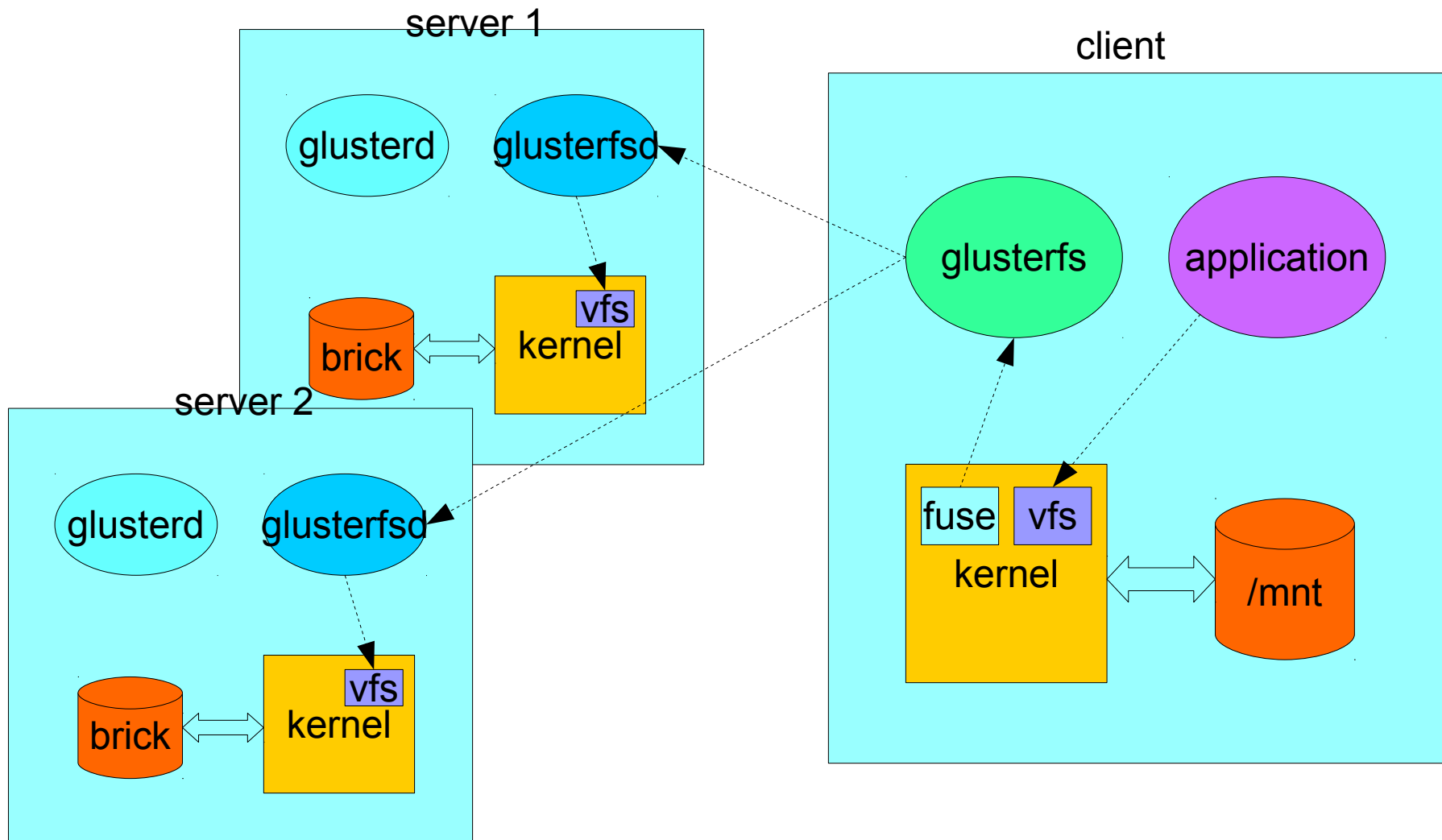
# Complex Translator Stack — Replica (AFR)

**Kaleb KEITHLEY**

server 1

client

glusterd  glusterfsd

glusterfs  application

vfs
kernel

brick

server 2

glusterd  glusterfsd

fuse  vfs

kernel

/mnt

vfs
kernel

brick

14

**Kaleb KEITHLEY**

# Tiering: a variation on distribution

- Hot and cold data
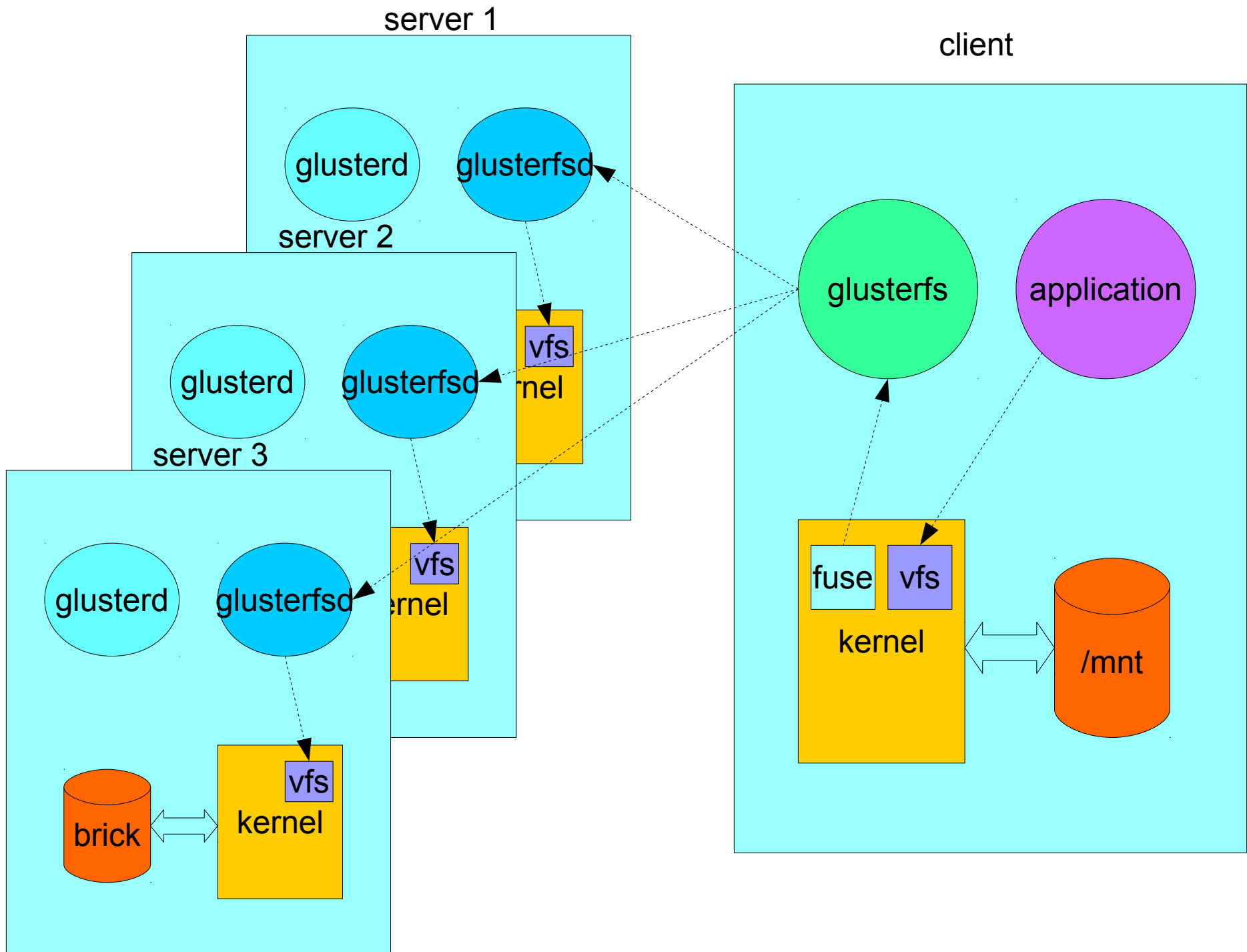
**Kaleb KEITHLEY**

# Adding a brick — distributed and rebalancing

**Kaleb KEITHLEY**

# Adding a brick — replica and healing

**Kaleb KEITHLEY**

**Kaleb KEITHLEY**

# Isn't FUSE slow?

client

glusterfsd ← glusterfs    application

fuse   vfs

kernel   ↔   /mnt

- Count the copies
- And context switches

**Kaleb KEITHLEY**

# gfapi to the rescue

client

glusterfsd

application

- POSIX-like API
  - glfs_open()
  - glfs_close()
  - glfs_read(), glfs_write()
  - glfs_seek()
  - glfs_stat()
  - etc.

fuse   vfs

kernel

/mnt

# libgfapi — hellogluster.c

```c
#include <gluster/api/glfs.h>

int main (int argc, char** argv)

{

    fs = glfs_new ("fsync");

    ret = glfs_set_volfile_server (fs, "tcp", "localhost",
    24007);

    /* or ret = glfs_set_volfile (fs, "/tmp/foo.vol"); */

    ret = glfs_init (fs);

    fd = glfs_creat (fs, filename, O_RDWR, 0644);

    ret = glfs_write (fd, "hello gluster", 14, 0);

    glfs_close (fd);

    return 0;

}
```

# Writing a translator in Python with GluPy

**Kaleb KEITHLEY**

# SwiftOnFile — RESTful API, examples with curl

- Create a container: `curl -v -X PUT -H 'X-Auth-Token: $authtoken' https://$myhostname:443/v1/AUTH_$myvolname/$mycontainername -k`

- List containers: `curl -v -X GET -H 'X-Auth-Token: $authtoken' https://$myhostname:443/v1/AUTH_$myvolname -k`

- Copy a file into a container (upload): `curl -v -X PUT -T $filename -H 'X-Auth-Token: $authtoken' -H 'Content-Length: $filelen' https://$myhostname:443/v1/AUTH_$myvolname/$mycontainername/ $filename -k`

- Copy a file from a container (download): `curl -v -X GET -H 'X-Auth-Token: $authtoken' https://$myhostname:443/v1/AUTH_$myvolname/ $mycontainername/$filename -k > $filename`

**Kaleb KEITHLEY**

# Translator basics

- Translators are shared objects (shlibs)

  - Methods

    - int32_t init(xlator_t *this);
    - void fini(xlator_t *this);

  - Data

    - struct xlator_fops fops { ... };
    - struct xlator_cbks cbks { };
    - struct volume_options options [] = { ... };

- Client, Server, Client/Server

- Threads: write MT-SAFE

- Portability: GlusterFS != Linux only

- License: GPLv2 or LGPLv3+

**Kaleb KEITHLEY**

# Every method has a different signature

- Open fop method and callback

```
typedef int32_t (*fop_open_t) (call_frame_t *, xlator_t *,
loc_t *, int32_t, fd_t *, dict_t *);

typedef int32_t (*fop_open_cbk_t) (call_frame_t *, void *,
xlator_t *, int32_t, int32_t, fd_t *, dict_t *);
```

- Rename fop method and callback

```
typedef int32_t (*fop_rename_t) (call_frame_t *, xlator_t *,
loc_t *, loc_t *, dict_t *);

typedef int32_t (*fop_rename_cbk_t) (call_frame_t , void *,
xlator_t *, int32_t, int32_t, struct iatt *, struct iatt *,
struct iatt *, struct iatt *, struct iatt *);
```

**Kaleb KEITHLEY**

# Data Types in Translators

- call_frame_t —

- xlator_t — translator context

- inode_t — represents a file on disk; ref-counted

- fd_t — represents an open file; ref-counted

- iatt_t — ~= struct stat

- dict_t — ~= Python dict (or C++ std::map)

- client_t — represents the connect client

**Kaleb KEITHLEY**

# fop methods and fop callbacks

```
uidmap_writev (...)

{

    ...

    STACK_WIND (frame, uidmap_writev_cbk,

            FIRST_CHILD (this),
            FIRST_CHILD (this)->fops->writev,
            fd, vector, count, offset, iobref);
    /* DANGER ZONE */

    return 0;

}
```

- Effectively lose control after STACK_WIND

  - Callback might have already happened

  - Or might be running right now

  - Or maybe it's not going to run 'til later

**Kaleb KEITHLEY**

# fop methods and fop callback methods, cont.

```
uidmap_writev_cbk (call_frame_t *frame, void *cookie, ...)

{

    ...

    STACK_UNWIND_STRICT (writev, frame

        op_ret, op_errno, prebuf, postbuf);

    return 0;

}
```

- The I/O is complete when the callback is called

**Kaleb KEITHLEY**

# STACK_WIND versus STACK_WIND_COOKIE

- Pass extra data to the cbk with STACK_WIND_COOKIE

```
quota_statfs (call_frame_t *frame,

    xlator_t *this, loc_t *loc)

{

    inode_t *root_inode = loc->inode->table->root;

    STACK_WIND_COOKIE (frame, quota_statfs_cbk,

        root_inode, FIRST_CHILD (this),

        FIRST_CHILD (this)->fops->statfs, loc, xdata);

    return 0;

}
```

- There is also frame->local

  - shared by all STACK_WIND callbacks

**Kaleb KEITHLEY**

# STACK_WIND, STACK_WIND_COOKIE, cont.

- Pass extra data to the cbk with STACK_WIND_COOKIE

```
quota_statfs_cbk (call_frame_t *frame, void *cookie, ...)

{

    inode_t *root_inode = cookie;

    ...

}
```

**Kaleb KEITHLEY**

# STACK_UNWIND versus STACK_UNWIND_STRICT

- ## STACK_UNWIND_STRICT uses the correct type

  ```
  /* return from function in a type-safe way */

  #define STACK_UNWIND (frame, params ...)

      do {

          ret_fn_t fn = frame->ret;

      ...
  ```

  versus

  ```
  #define STACK_UNWIND_STRICT (op, frame, params ...)

      do {

          fop_##op##_cbk_t fn = (fop_##op##_cbk_t)frame->ret;

      ...
  ```

- And why wouldn't you want strong typing?

**Kaleb KEITHLEY**

# Calling multiple children (fan out)

```
afr_writev_wind (...)

{

    ...

    for (i = 0; i < priv->child_count; i++) {

        if (local->transaction.pre_op[i]) {

            STACK_WIND_COOKIE (frame, afr_writev_wind_cbk,
                             (void *) (long) i,
                             priv->children[i],
                             priv->children[i]->fops->writev,
                             local->fd, ...);

        }

    }

    return 0;

}
```

**Kaleb KEITHLEY**

# Calling multiple children, cont. (fan in)

```
afr_writev_wind_cbk (...)

{

    LOCK (&frame->lock);

    callcnt = --local->call_count;

    UNLOCK (&frame->lock);

    if (callcnt == 0) /* we're done */

    ...

}
```

- failure by any one child means the whole transaction failed?

    - And needs to be handled accordingly

**Kaleb KEITHLEY**

# Dealing With Errors: I/O errors

```
uidmap_writev_cbk (call_frame_t *frame, void *cookie,

    xlator_t *this, int32_t op_ret, int32_t op_errno, ...)

{

    ...

    STACK_UNWIND_STRICT (writev, frame, -1, EIO, ...);

    return 0;

}
```

- op_ret: 0 or -1, success or failure
- op_errno: from <errno.h>
  - Use an op_errno that's valid and/or relevant for the fop

**Kaleb KEITHLEY**

# Dealing With Errors: FOP method errors

```
uidmap_writev (call_frame_t *frame, xlator_t *this, ...)

{

    ...

    if (horrible_logic_error_must_abort) {

        goto error; /* glusterfs idiom */

    }

    STACK_WIND(frame, uid_writev_cbk, ...);

    return 0;

error:

    STACK_UNWIND_STRICT (writev, frame, -1, EIO, NULL, NULL);

    return 0;

}
```

**Kaleb KEITHLEY**

# Call to action

- Go forth and write applications for GlusterFS!

**Kaleb KEITHLEY**

# £©€é¥¢ßŒ

**Kaleb KEITHLEY**