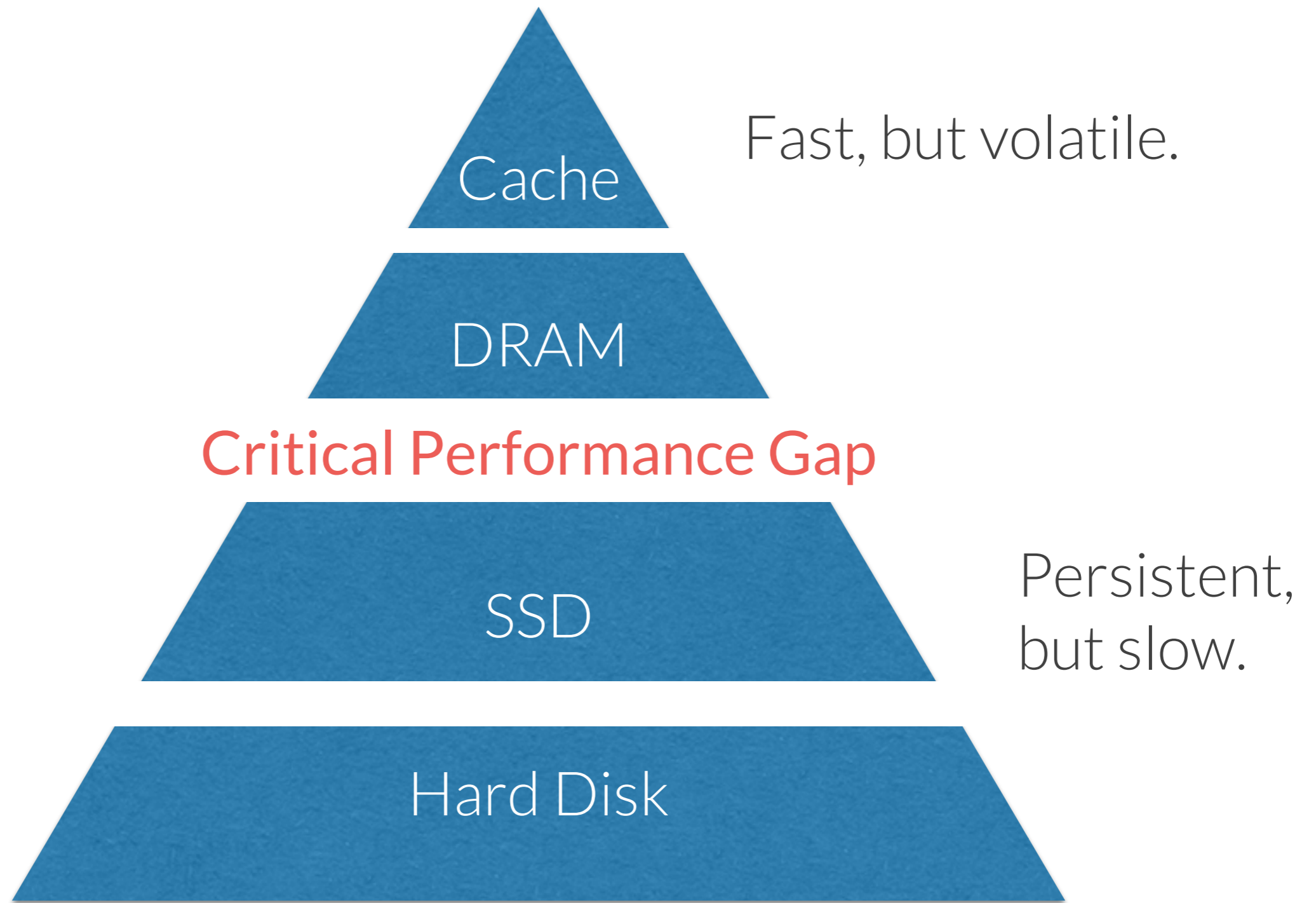


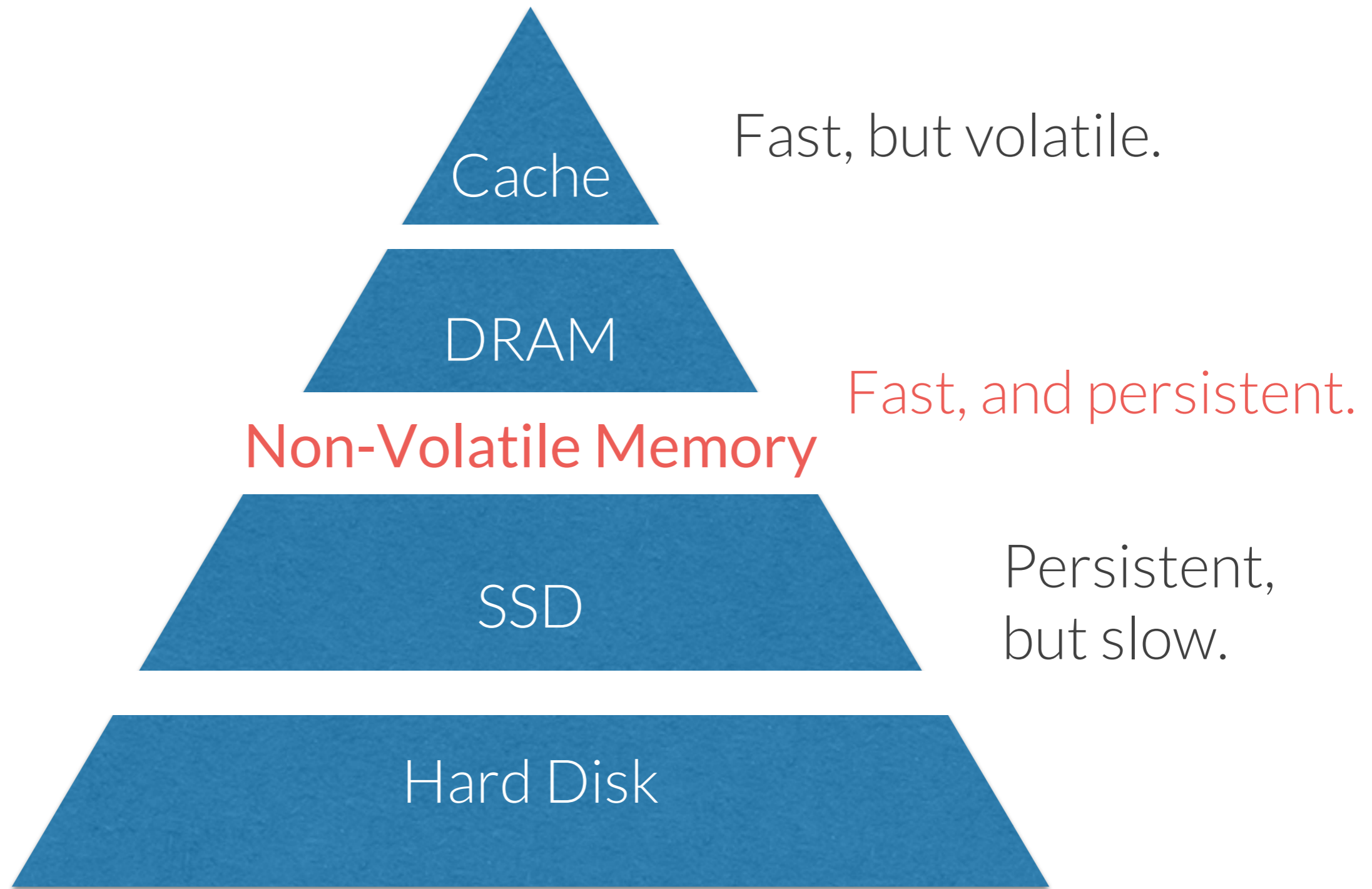
NVMOVE: Helping Programmers Move to Byte-based Persistence

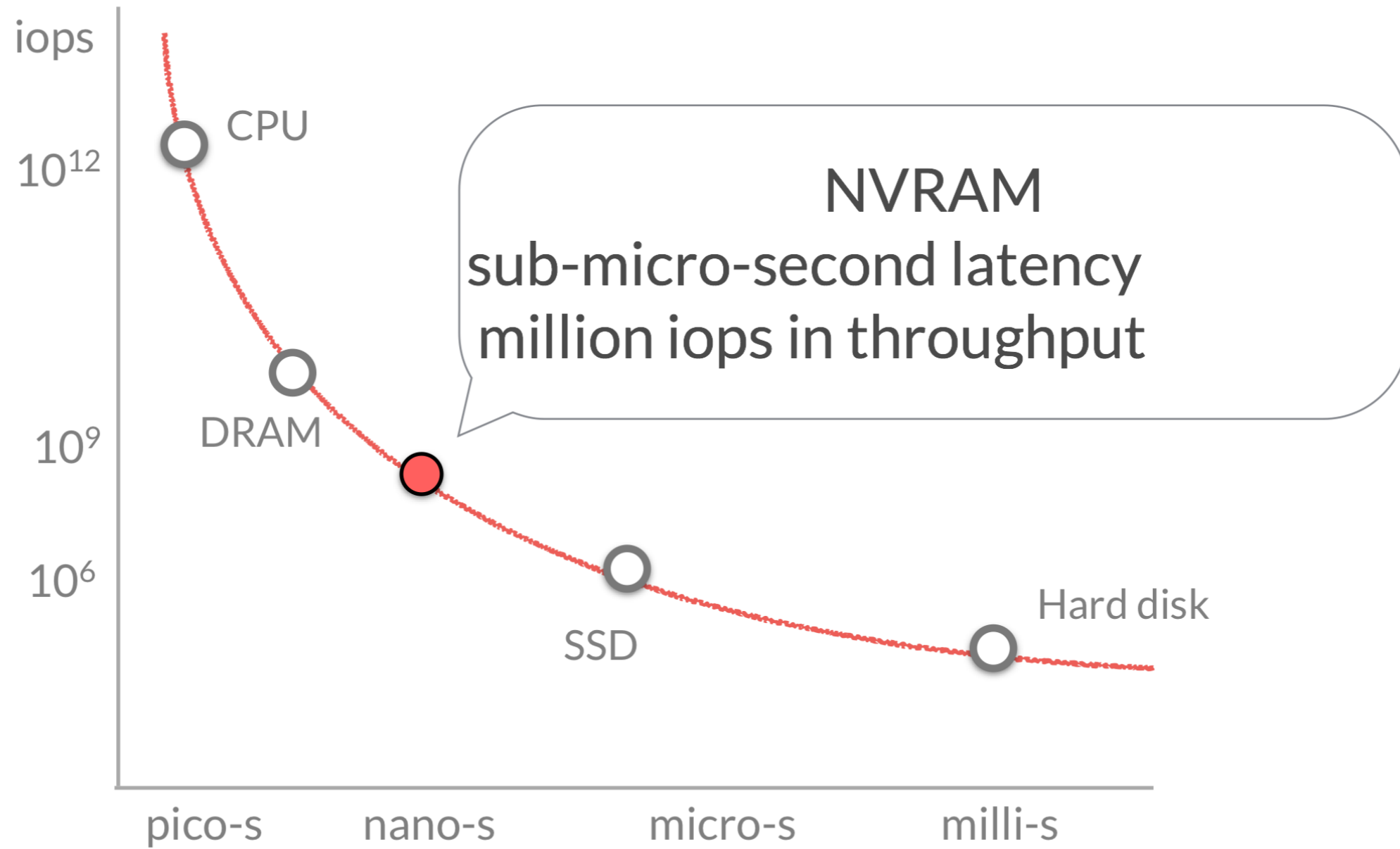
Himanshu Chauhan

with

Irina Calciu, Vijay Chidambaram,
Eric Schkufza, Onur Mutlu, Pratap Subrahmanyam





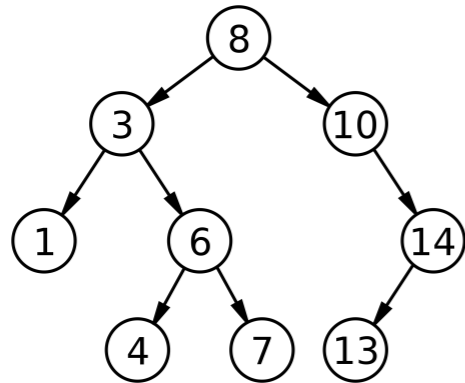


Persistent Programs

```
typedef struct {  
  
} node
```

1. allocate from memory
2. data read/write + program logic
3. save to storage

Persistence Today



In-memory binary search tree

Serialization

Flat Buffer

```
printf(buf, "%d:%s", node->id, node->value)
```

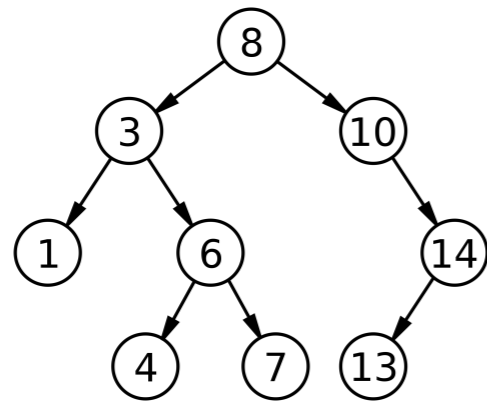
File

```
write(fd, buf, sizeof(buf))  
fsync(fd)
```

Block-sized Writes

Block-based Storage

Persistence with NVM



In-memory binary search tree

```
node->id = 10  
memcpy(node->value, myvalue)  
pmemobj_persist(node)
```

Byte-sized Writes

Byte-based NVM

Changing Persistence Code

Present

```
/* allocate from volatile memory*/  
node n* = malloc(sizeof(...))  
...
```

```
node->value = val //volatile  
update  
...
```

```
/* persist to block-storage*  
char *buf= malloc(sizeof(...))  
int fd = open("data.db",O_WRONLY)  
sprintf(buf,"...", node->id,  
         node->value);  
write(fd, buf, sizeof(buf));
```



NVM

```
/* allocate from non-volatile memory*/  
node n* = pmalloc(sizeof(...))  
...
```

```
node->value = val //persistent  
update  
...
```

```
flush cache and commit*/  
cache_flush + __commit
```


Porting to NVM: Tedious

- Identify data structures that should be on NVM
- Update them in a consistent manner

Redis: simple key-value store (~50K LOC)

- Industrial effort to port Redis is on-going after two years
- Open-source effort to port Redis has minimum functionality

Changing Persistence Code

Present

```
/* allocate from volatile memory*/  
node n* = malloc(sizeof(...))
```

```
node->value = val //volatile  
update  
...
```

```
/* persist to block-storage*/  
char *buf= malloc(sizeof(...))  
int fd = open("data.db",O_W  
sprintf(buf,"...", node->id,  
node->value);  
write(fd, buf, sizeof(buf));
```



NVM

```
/* allocate from non-volatile memory*/  
node n* = pmalloc(sizeof(...))
```

```
node->value = val //persistent  
update  
...
```

```
flush cache and commit*/  
cache_flush + __commit
```

Goal:

Port existing applications to
NVM with minimal programmer
involvement.



DRAG CHUTE

DATE: _____
EMERGENCY FLOW AT SEA LEVEL
MIL. 171M. RPM AT 17°C

RADIO CALL
53-504

INSTRUMENTATION CONTROL

GUNS
SIGHT
CAMERA
& RADAR
OFF
SIGHT
CAMERA
& RADAR

FIRE WARNING
PUSH TO TEST

HYDRAULIC
PRESSURE



EMERGENCY UP
OFF LANDING & TAXI LIGHT
EMERGENCY JETTISON
LANDING GEAR POSITION
L R
LANDING GEAR
HORN CUTOUT

OIL COOLER PUMP
AIR START
ENG MASTER
FUEL PURGE
RAM TURB ON SYS 2

EXT. EXT.
FIN HTP
POWER OFF AC

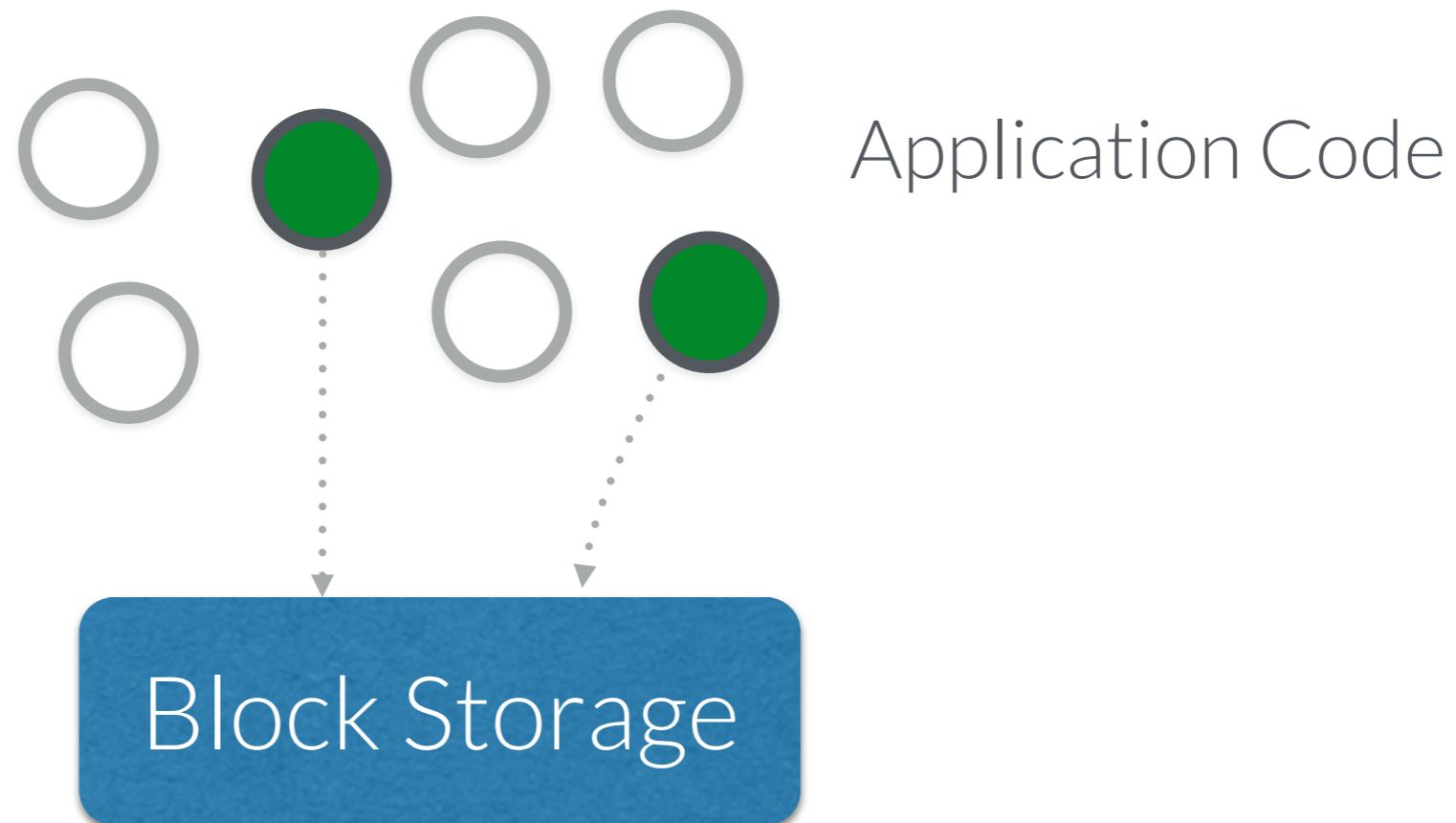
OXYGEN REG. PRESSURE
FLOW
EMERGENCY TEST



By Kiko Alario Salom [CC BY 2.0 (<http://creativecommons.org/licenses/by/2.0>)], via Wikimedia Commons

Persistent Types in Source

User defined source types (structs in C) that are persisted to block-storage.



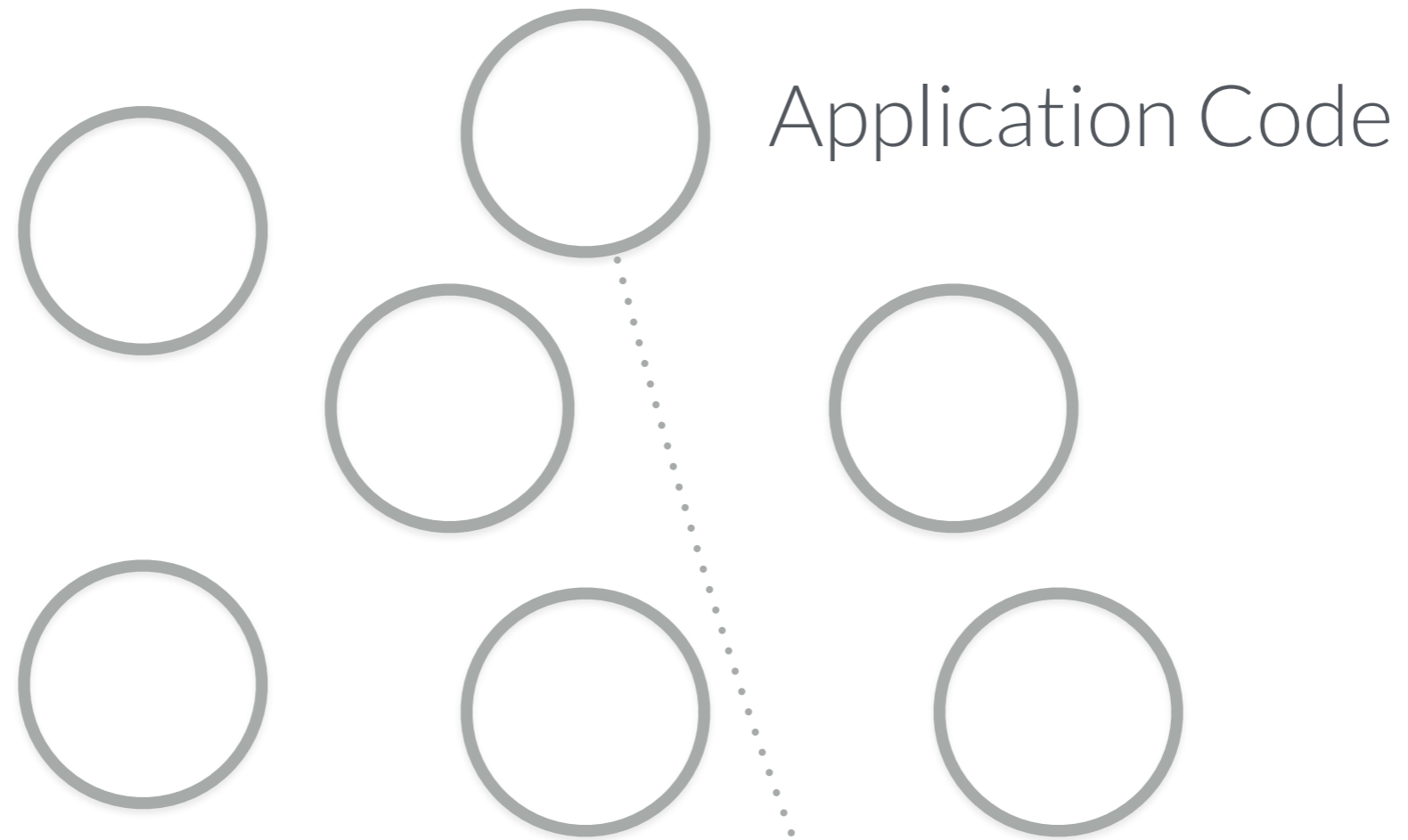
First Step:

Identify persistent types in
application source.

Solution: Static Analysis

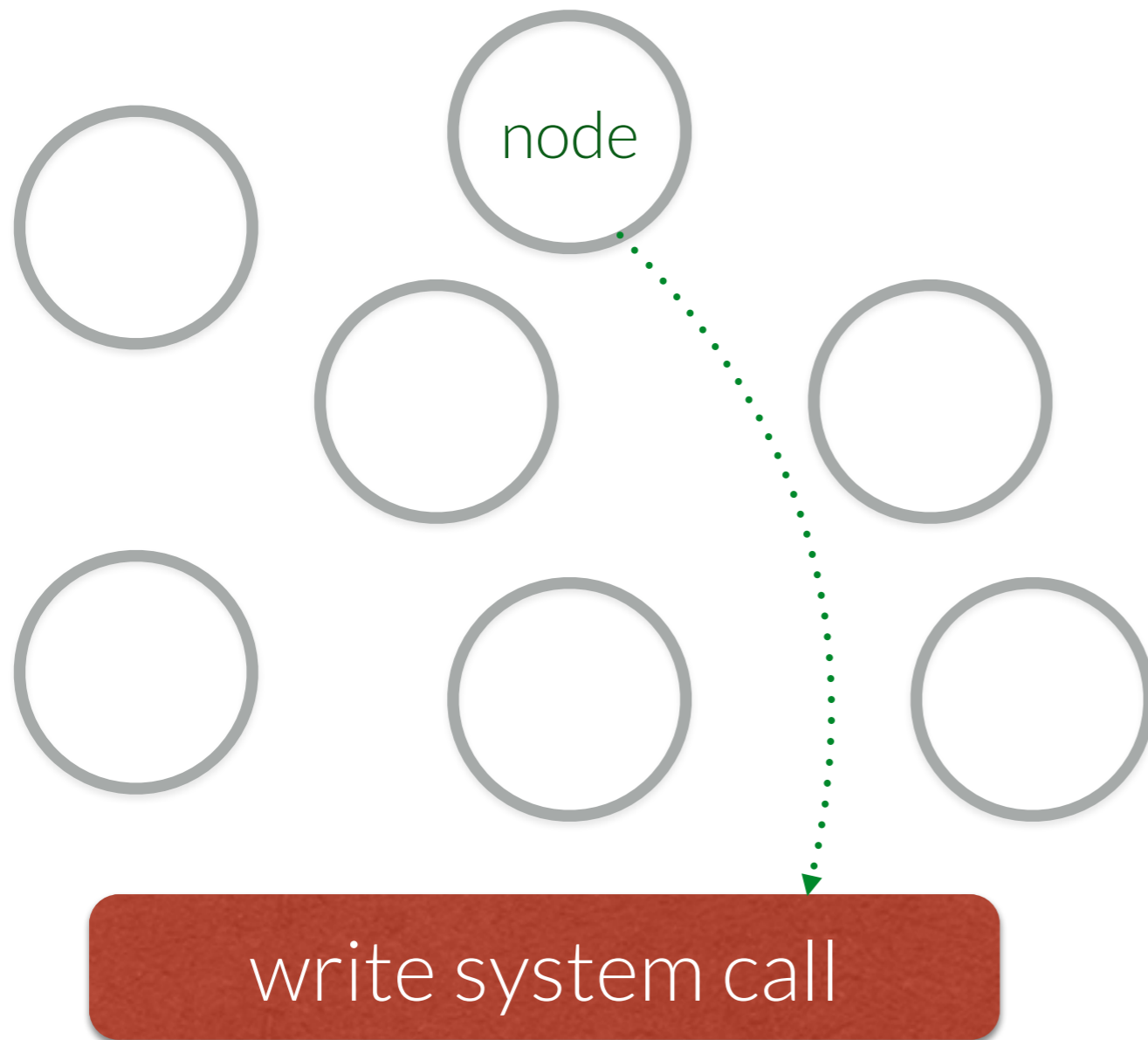
Current Focus: C

types = structs



write system call

Block Storage



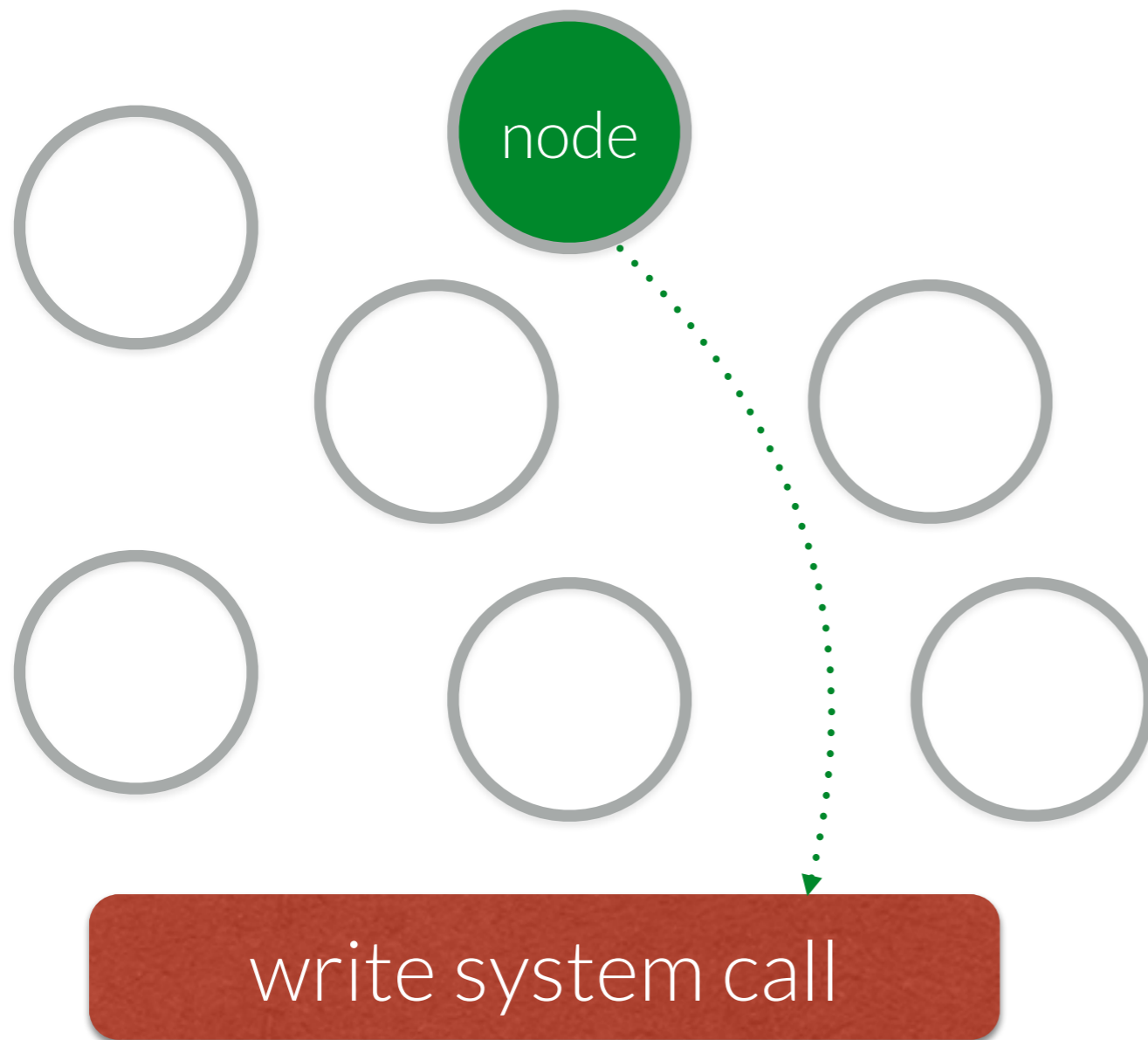
```
node *n = malloc(sizeof(node))
```

```
iter *it = malloc(sizeof(iter))
```

```
/* persist to block-storage*/  
char *buf = malloc(...)  
int fd = open(...)
```

```
sprintf(buf, "...", node->value)
```

```
write(fd, buf, ...)
```



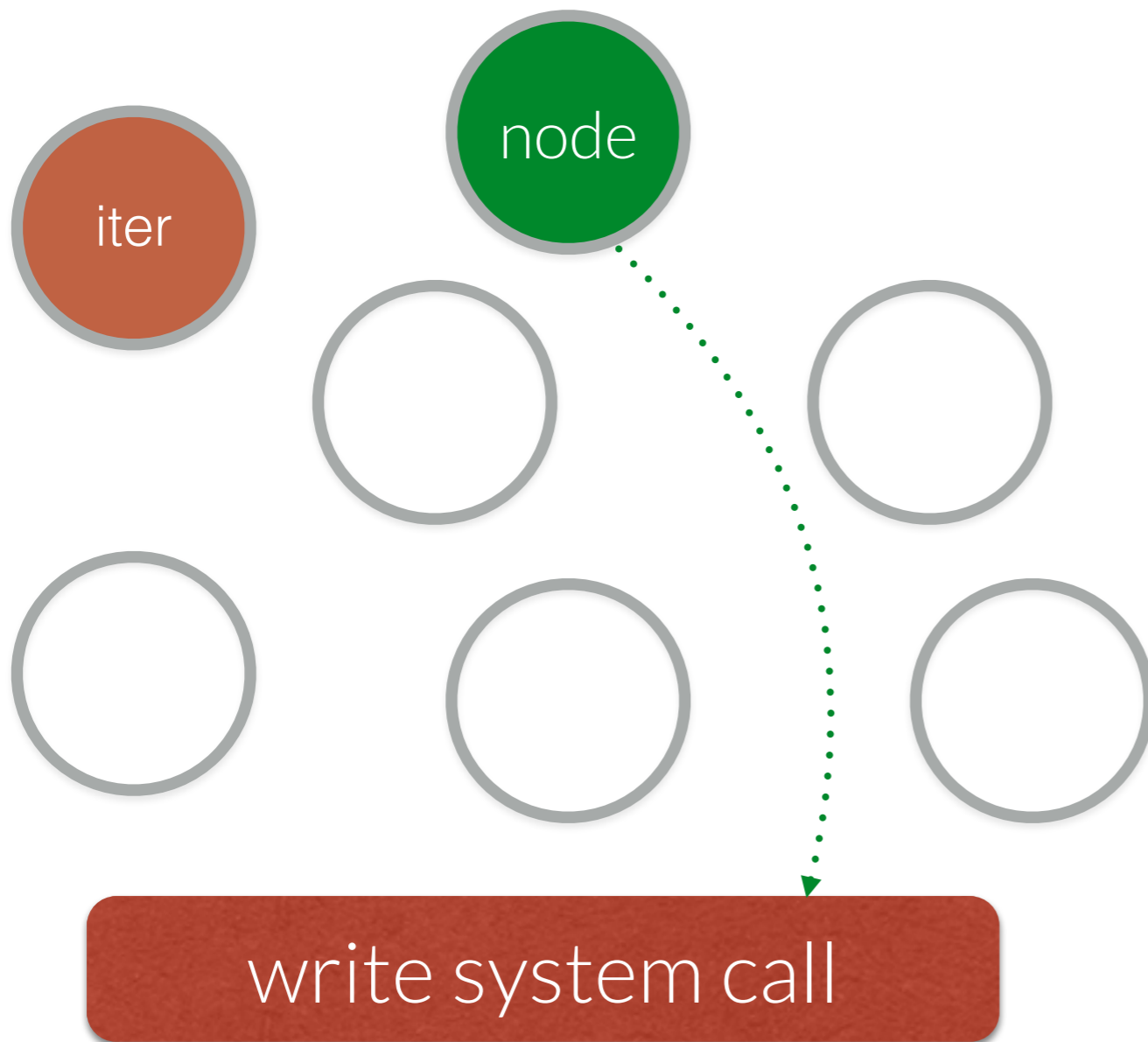
```
node *n = malloc(sizeof(node))
```

```
iter *it = malloc(sizeof(iter))
```

```
/* persist to block-storage */  
char *buf = malloc(...)  
int fd = open(...)
```

```
sprintf(buf, "...", node->value)
```

```
write(fd, buf, ...)
```



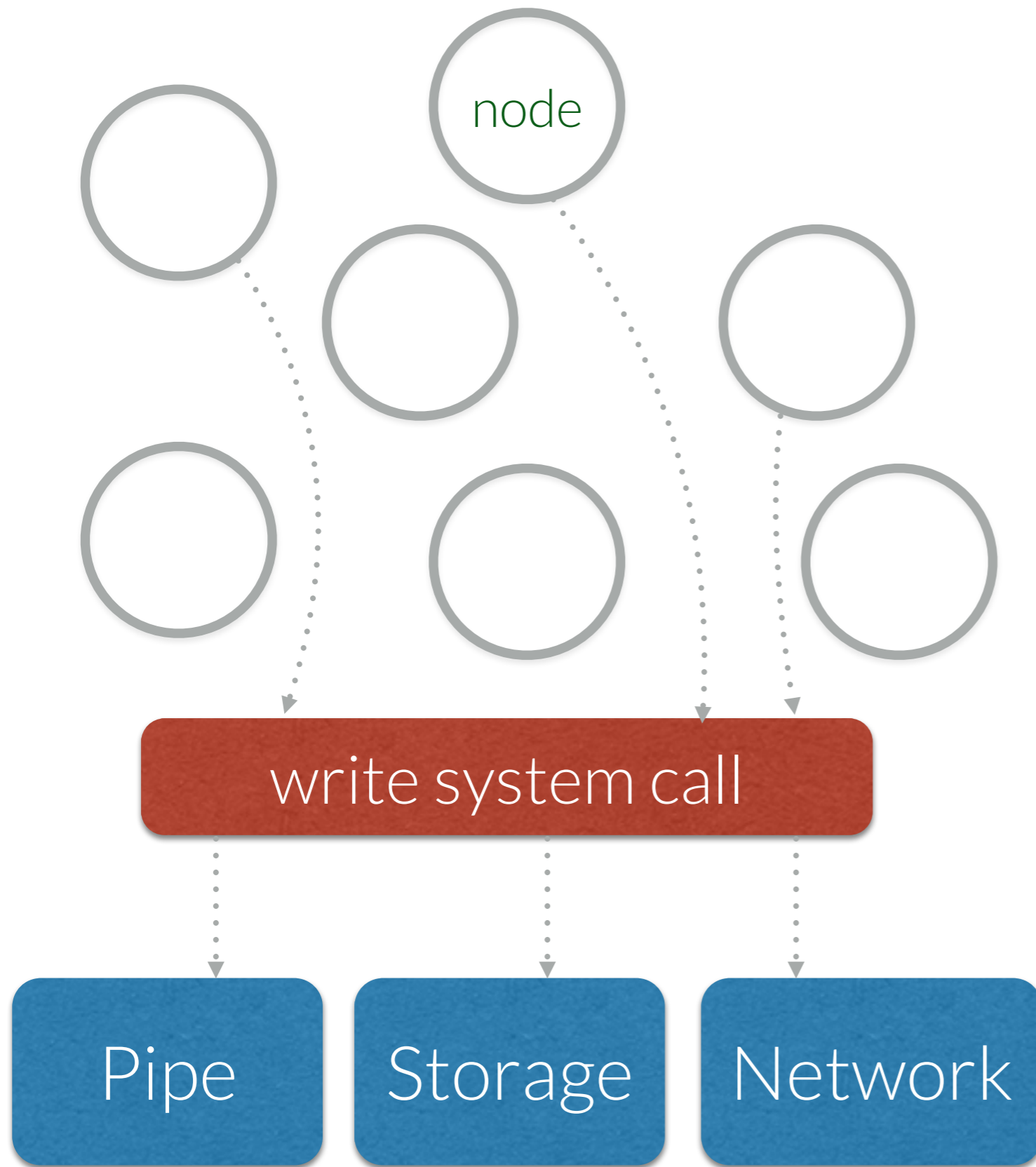
```
node *n = malloc(sizeof(node))
```

```
iter *it = malloc(sizeof(iter))
```

```
/* persist to block-storage */  
char *buf = malloc(...)  
int fd = open(...)
```

```
sprintf(buf, "...", node->value)
```

```
write(fd, buf, ...)
```



```
/* write to network socket*/
```

```
...
```

```
write(socket, "404", ...)
```

```
/* write to error stream*/
```

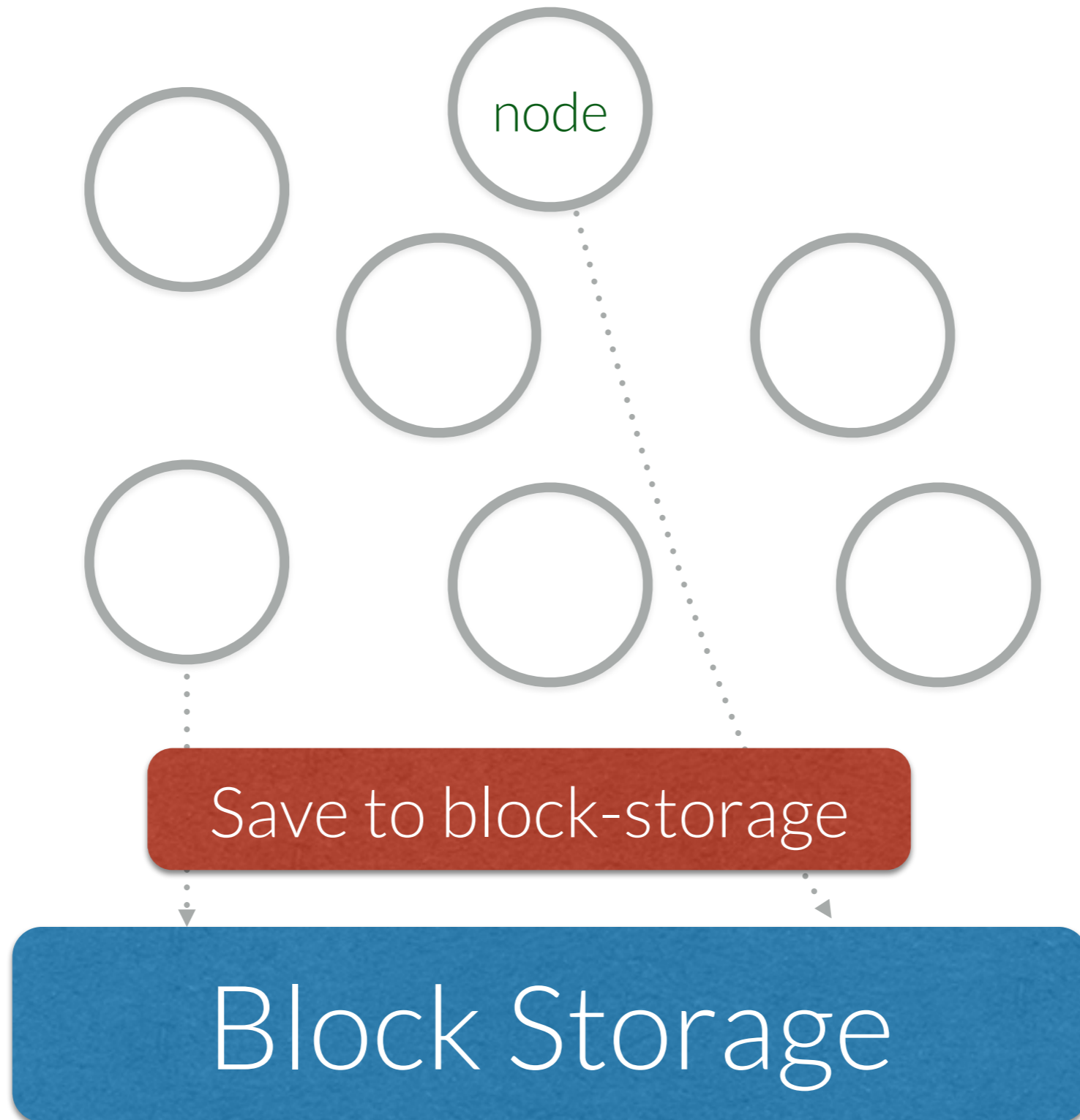
```
...
```

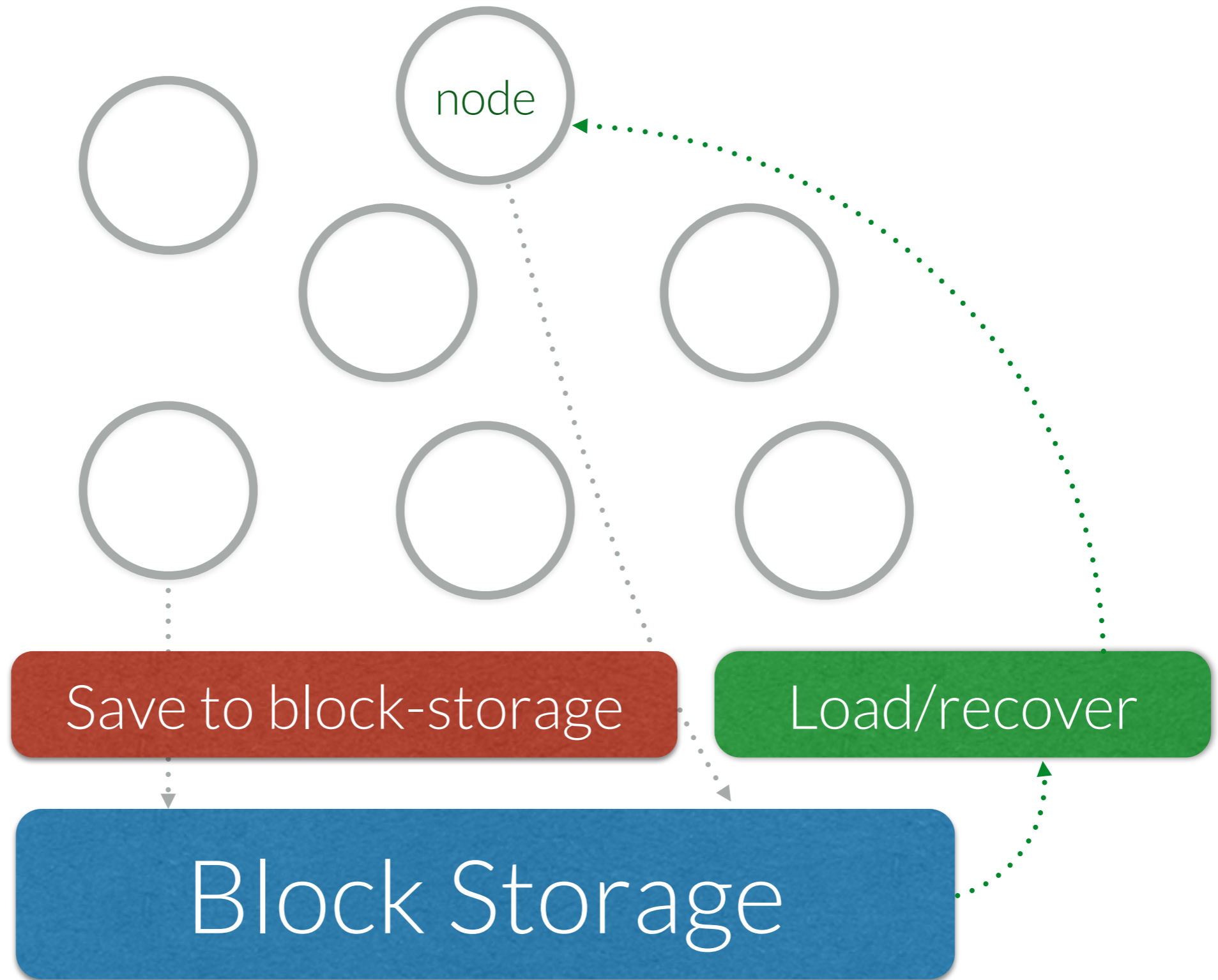
```
write(stderr, "All is lost.", ...)
```

```
/* persist to block-storage*/
```

```
...
```

```
write(fd, buf, ...)
```





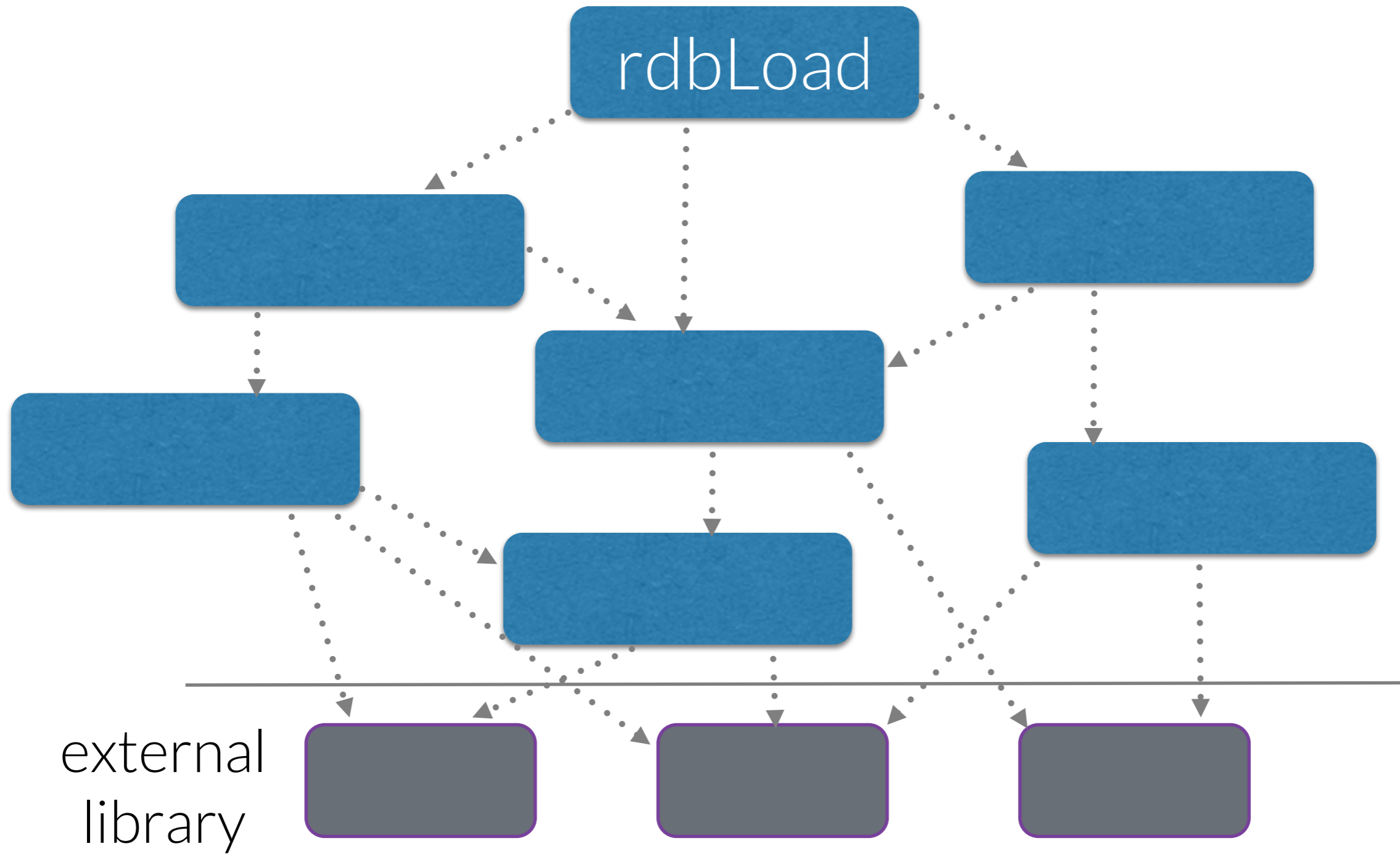
“rdbLoad” is the load/recovery function.

```
1257 int rdbLoad(char *filename) {
1258     uint32_t dbid;
1259     int type, rdbver;
1260     redisDb *db = server.db+0;
1261     char buf[1024];
1262     long long expiretime, now = mstime();
1263     FILE *fp;
1264     rio rdb;
1265
1266     if ((fp = fopen(filename,"r")) == NULL) return C_ERR;
1267
1268     rioInitWithFile(&rdb,fp);
1269     rdb.update_cksum = rdbLoadProgressCallback;
1270     rdb.max_processing_chunk = server.loading_process_events_interval_bytes;
1271     if (rioRead(&rdb,buf,9) == 0) goto eoferr;
1272     buf[9] = '\0';
1273     if (memcmp(buf,"REDIS",5) != 0) {
1274         fclose(fp);
1275         serverLog(LL_WARNING,"Wrong signature trying to load DB from file");
1276         errno = EINVAL;
1277         return C_ERR;
1278     }
1279     rdbver = atoi(buf+5);
```

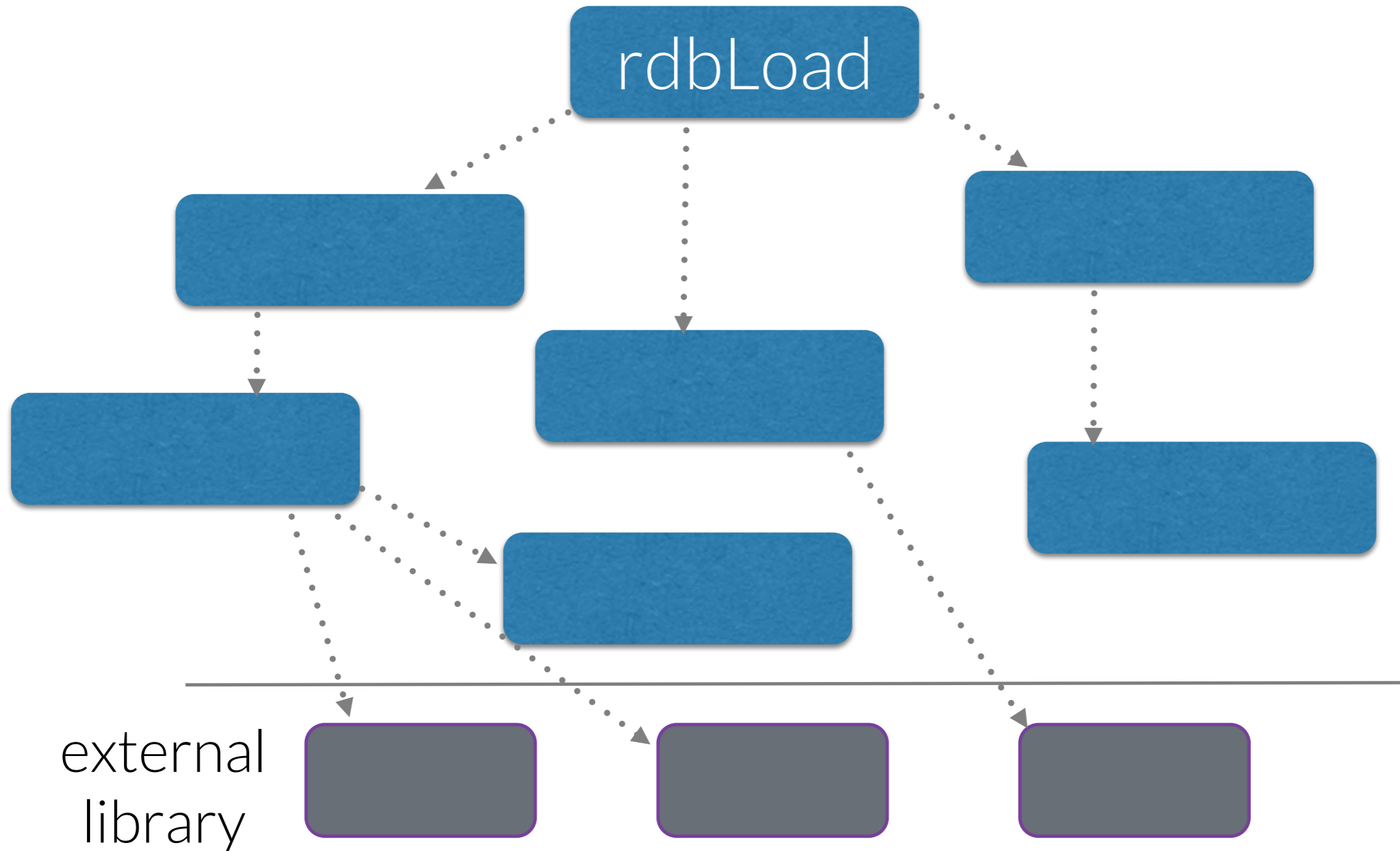
Mark every type that can be created during the recovery.

*if defined in application source.

Call Graph from Load

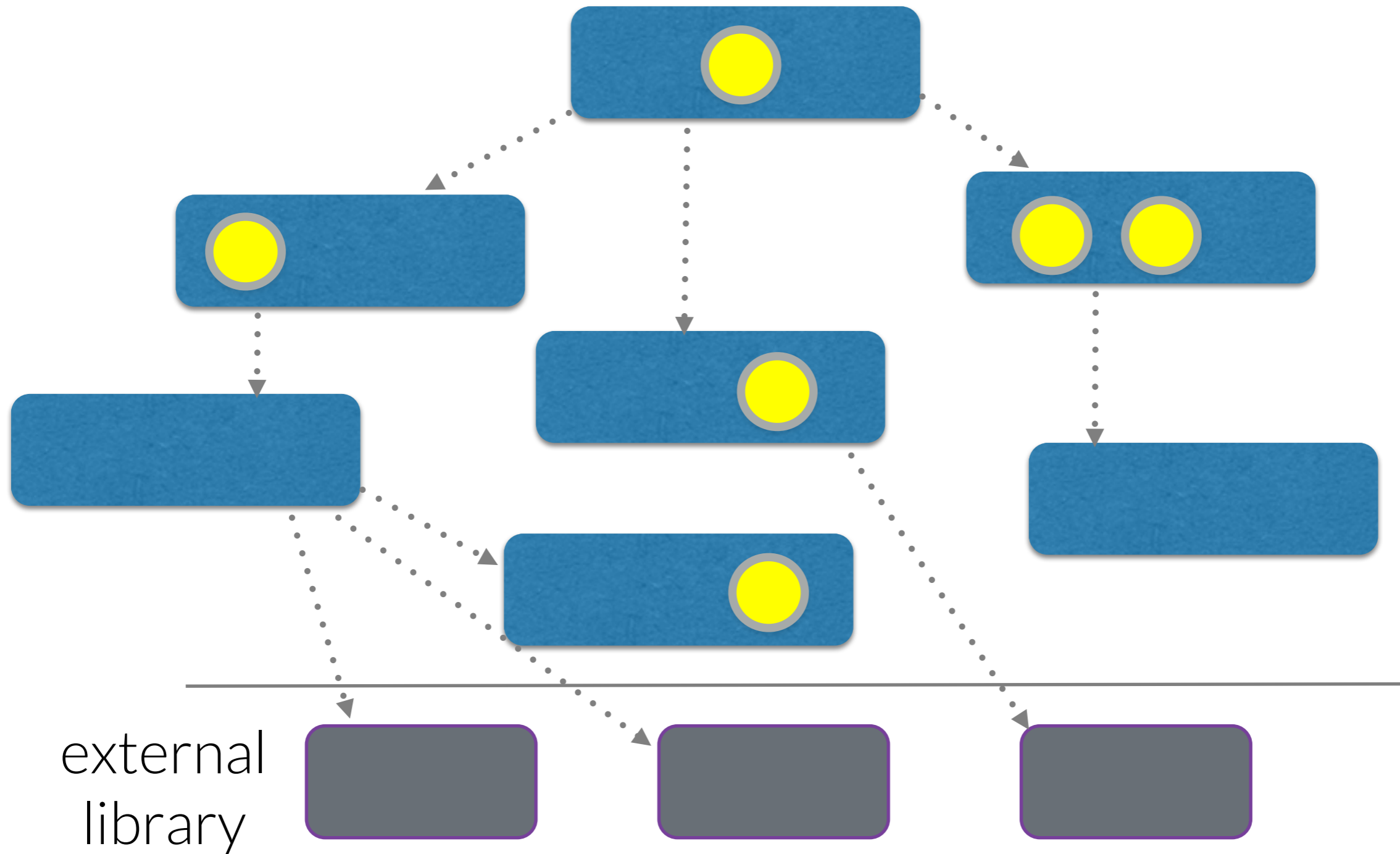


BFS on Call Graph from Load



BFS on Call Graph from Load

 Application type created/modified



NVMove Implementation

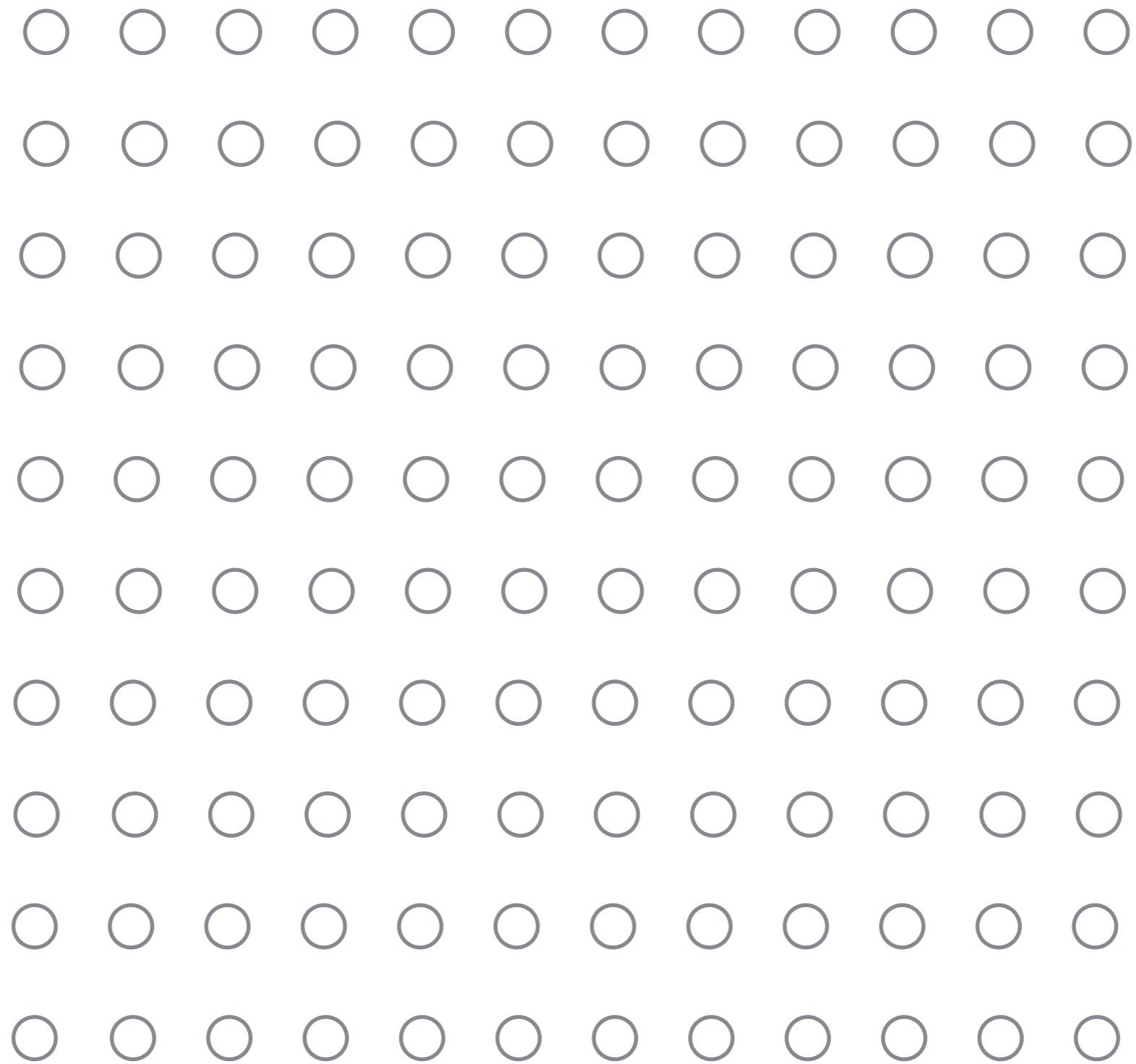
- Clang
 - Frontend Parsing
- Parse AST and Generate Call Graph
 - Find all statements that create/modify application types in graph
- Currently supports C applications

Evaluation



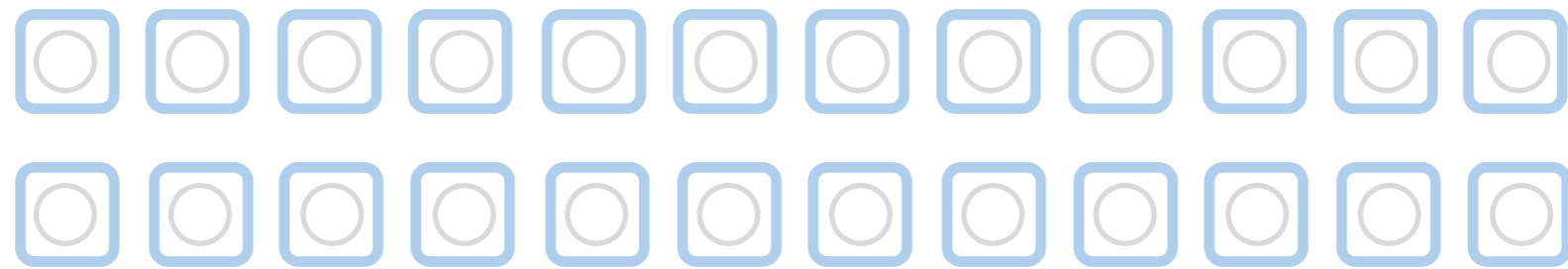
- In-memory data structure store
 - strings, hashes, lists, sets, indexes
- On-disk persistence
 - data-snapshots(RDB),
 - command-logging (AOF)
- ~50K lines-of-code

Identification Accuracy



122 types (structs) in Redis Source

Identification Accuracy



Total types	122
NVM _{OVE} identified persistent types	25
True positives (manually identified)	14
False positives	11
False negatives	0



Performance Impact

Redis Persistence

Snapshot (RDB)

- Data snapshot per second
- Not fully durable

Logging (AOF)

- Append each update command to a file
- Slow

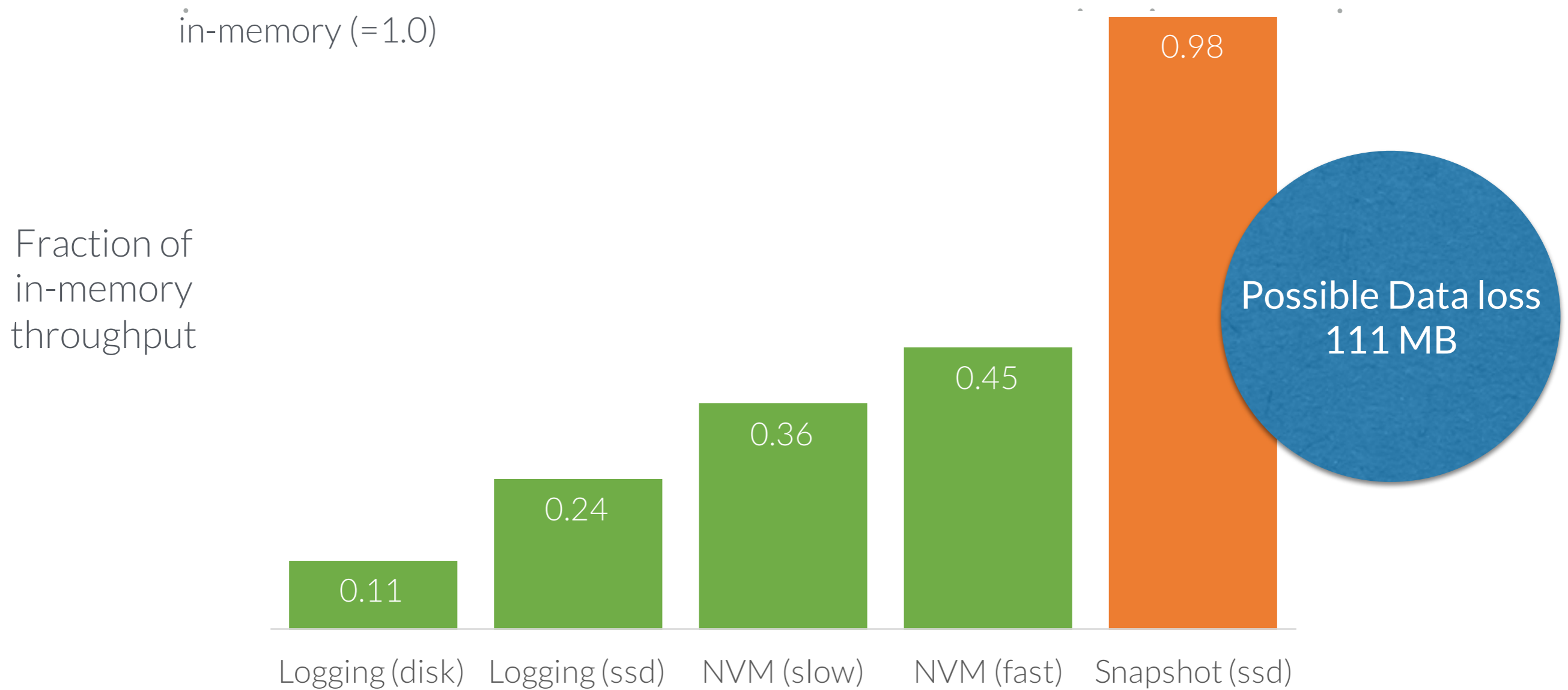
Both performed by forked background process.

NVM Emulation

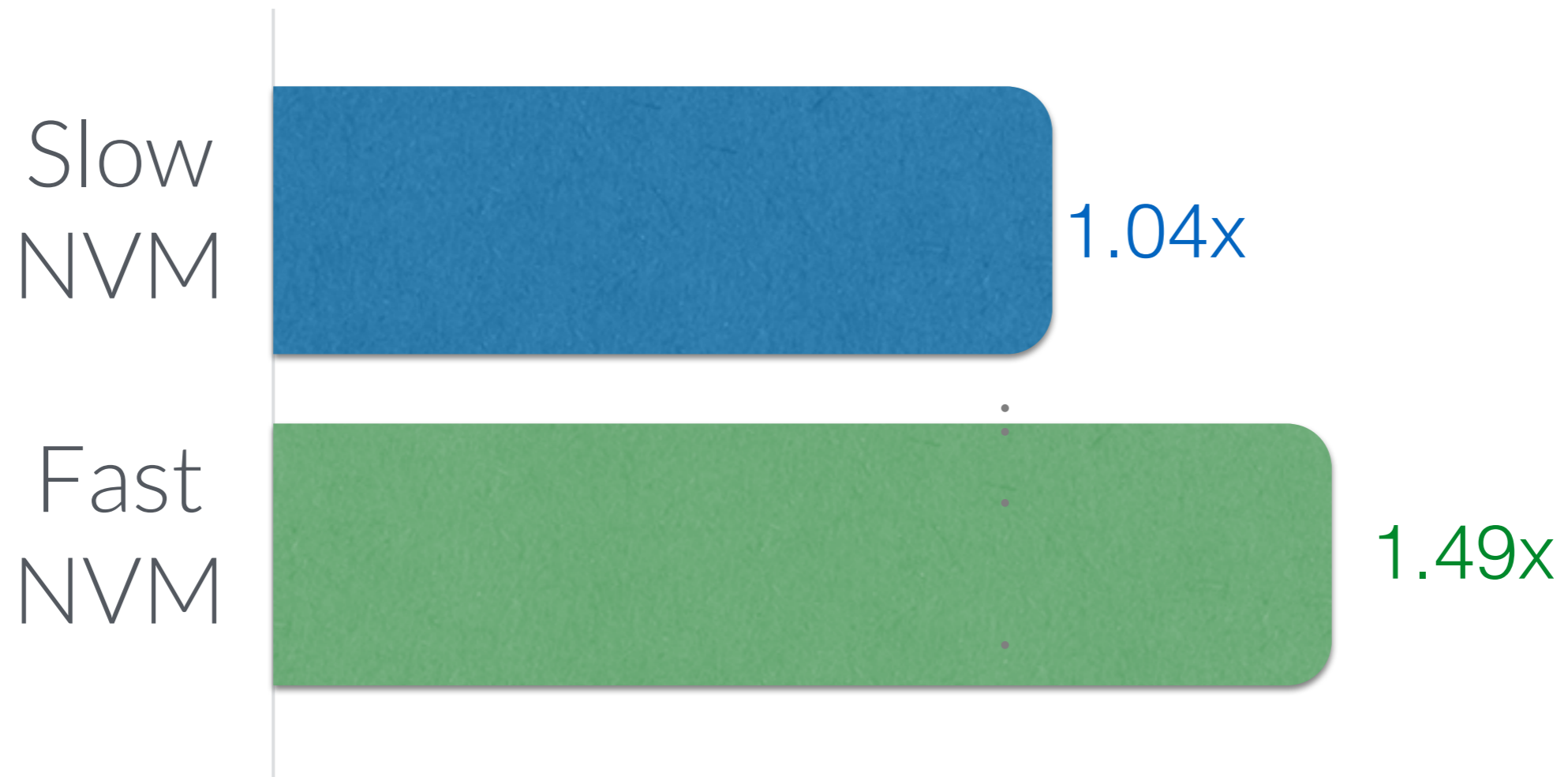
- Emulate allocation of NVMove identified types on NVM heap
 - Slow and Fast NVM
 - Inject delays for load/store of all NVM allocated types.
 - Worst-case performance estimate.
- Compare emulated NVM throughput against logging, and snapshot based persistence.

YCSB Benchmark Results

write-heavy (90% updated, 10% read ops)



Performance without False-Positives



1.0

Speedup in throughput

First Step:

Identify persistent types in
application source.

Next steps:

- Improve identification accuracy.
- Evaluate on other applications.

Backup

Throughputs (ops/sec)

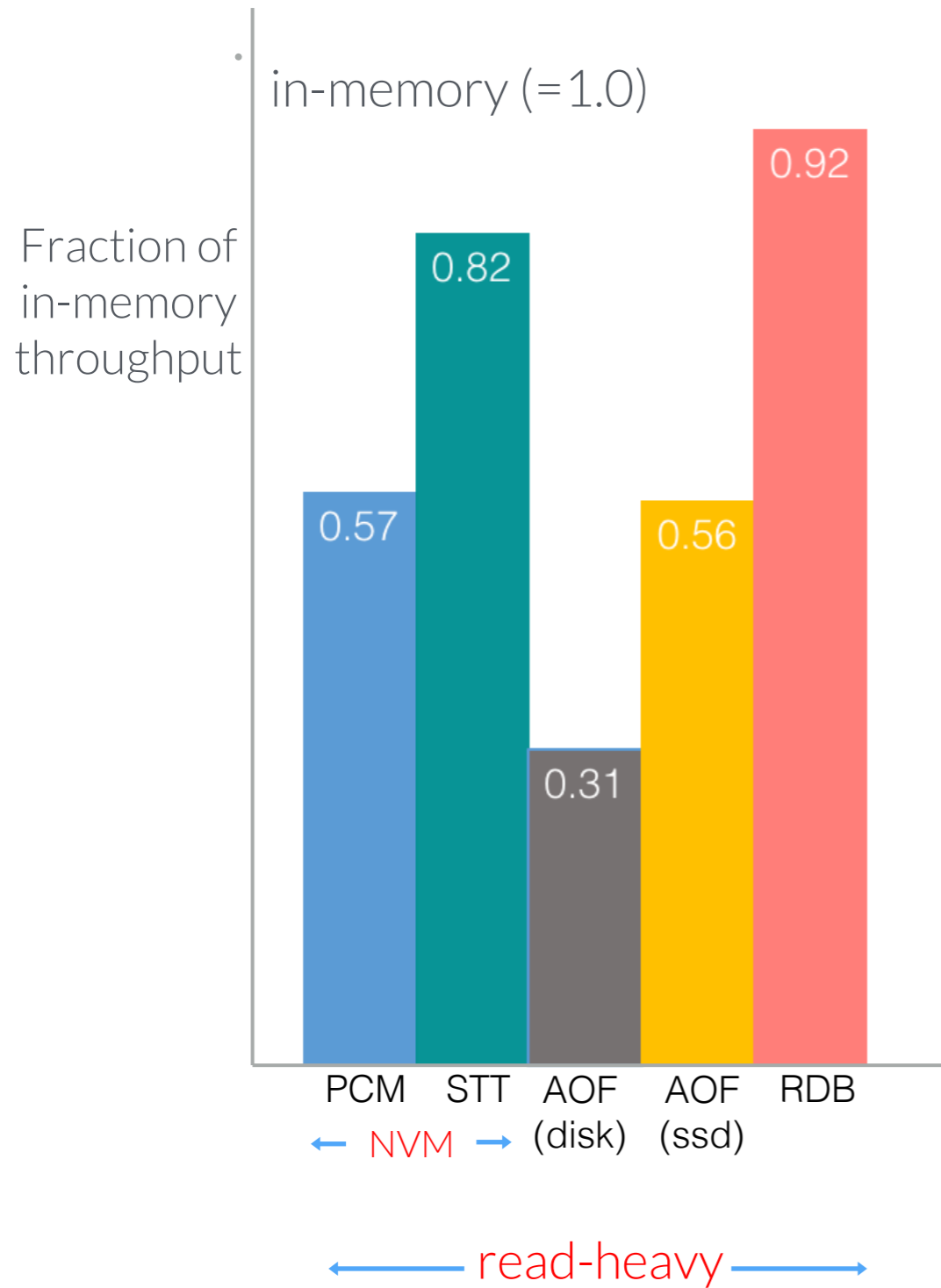
	readheavy	balance	writeheavy
PCM	28399	25,302	9759
STTRam	41213	38,048	12155
AoF (disk)	15634	6,457	2868
AoF (SSD)	27946	17,612	6605
RDB	46355	47,609	26605
Memory	50163	48,360	27156

NVM Emulation

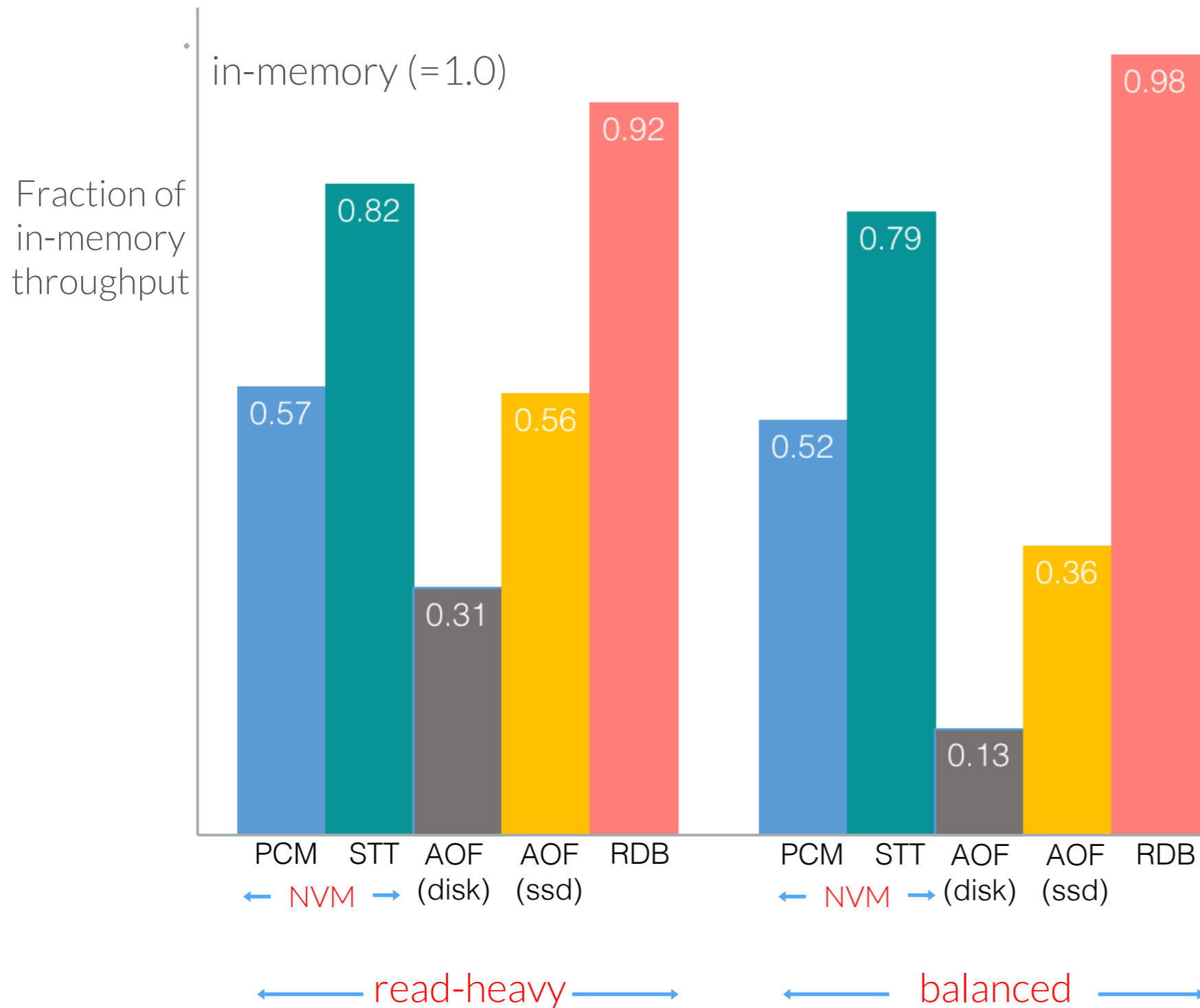
	Read Latency	Cache-line Flush Latency	PCOMMIT Latency
STT-RAM (Fast NVM)	100 ns	40 ns	200 ns
PCM (Slow NVM)	300 ns	40 ns	500 ns

*Xu & Swanson, NOVA: A Log-structured File System for Hybrid Volatile/Non-volatile Main Memories, FAST16.

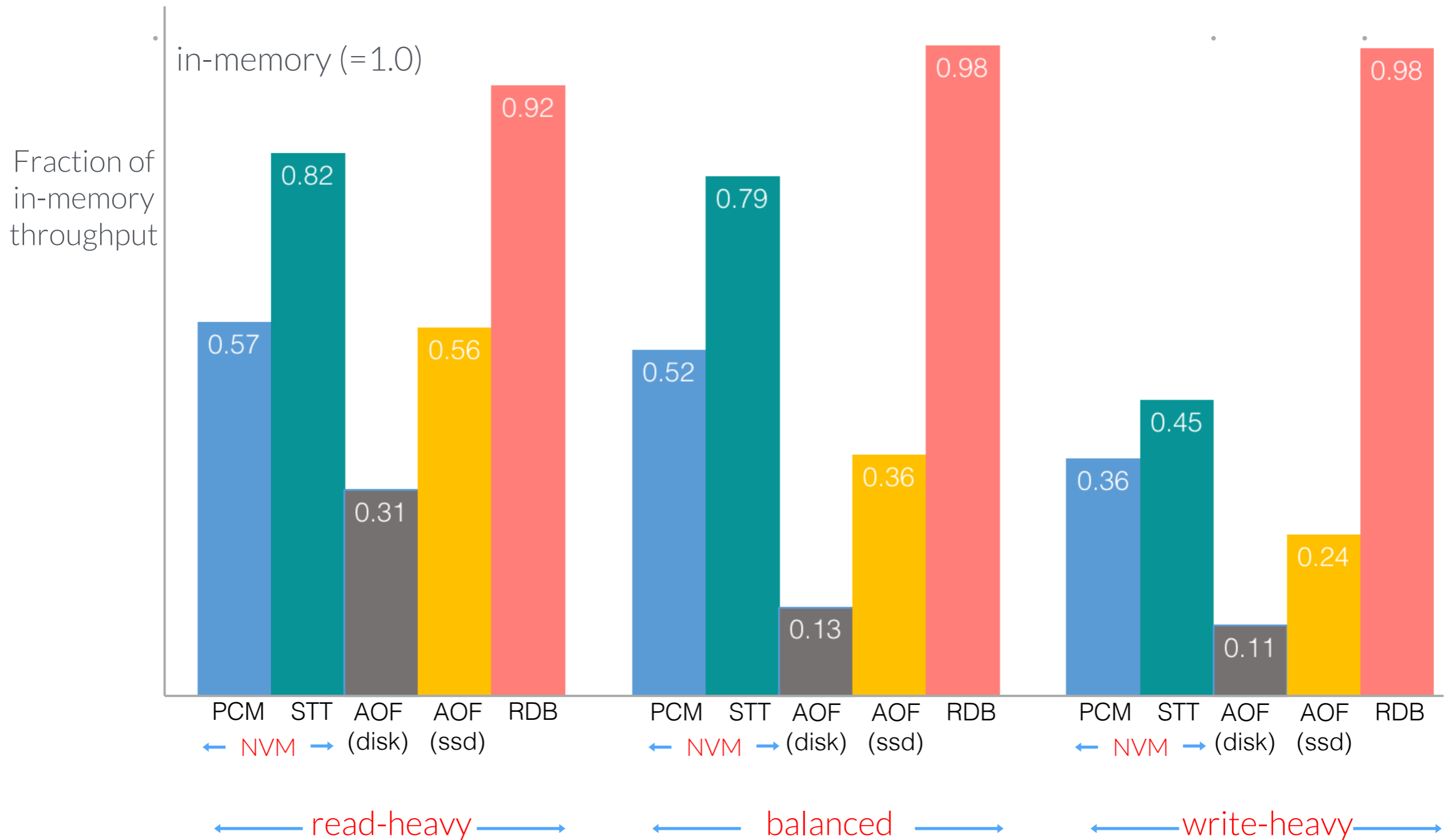
YCSB Benchmark Results



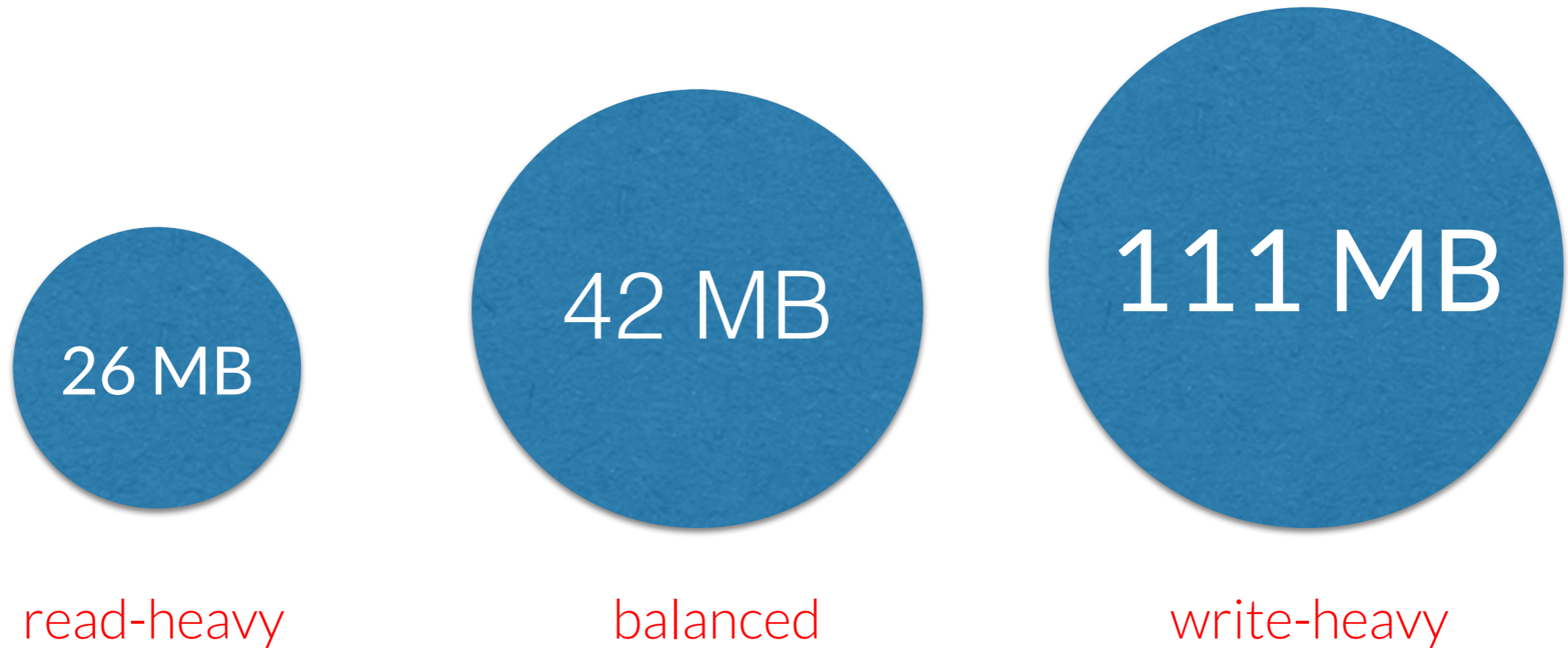
YCSB Benchmark Results



YCSB Benchmark Results



RDB Data Loss



Performance without False-Positives

