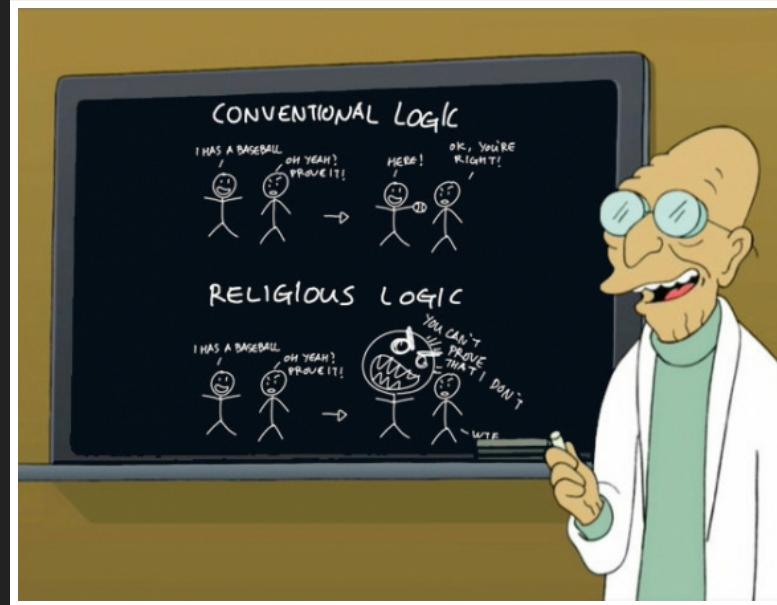


DISTRIBUTED STORAGE FOR CLOUD APPLICATIONS

A (very fast) overview

presented by Dimitri Pertin @ Polytech Nantes - 20/11/15

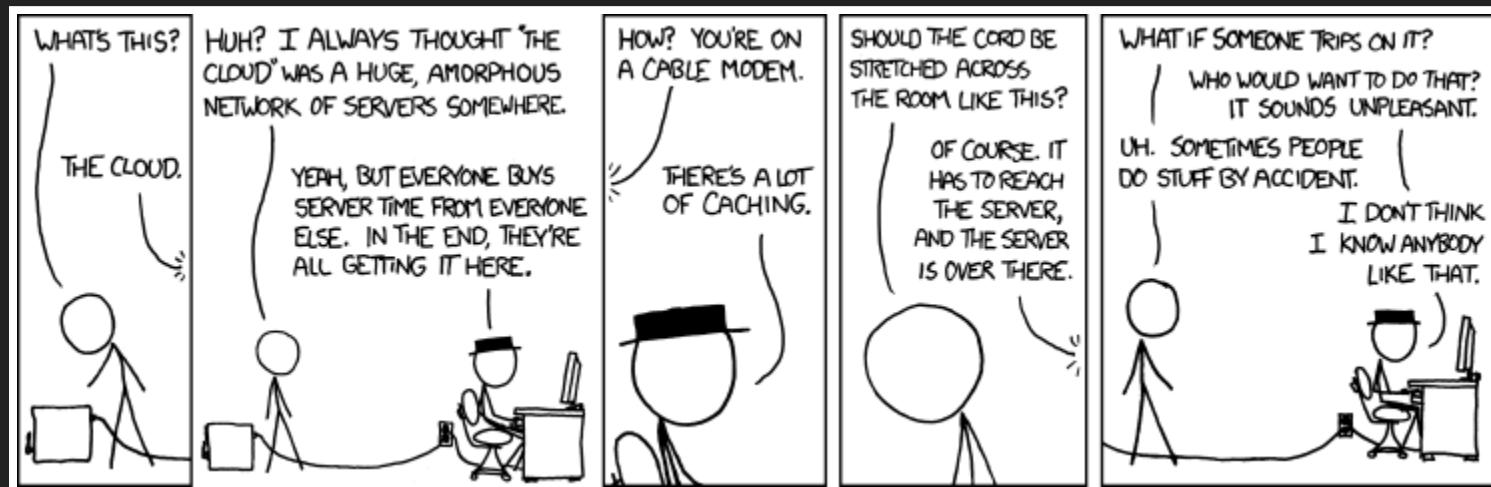
WHOAMI



*"Geometrical approach to design
distributed-storage-oriented erasure
codes"*

CLOUD STORAGE

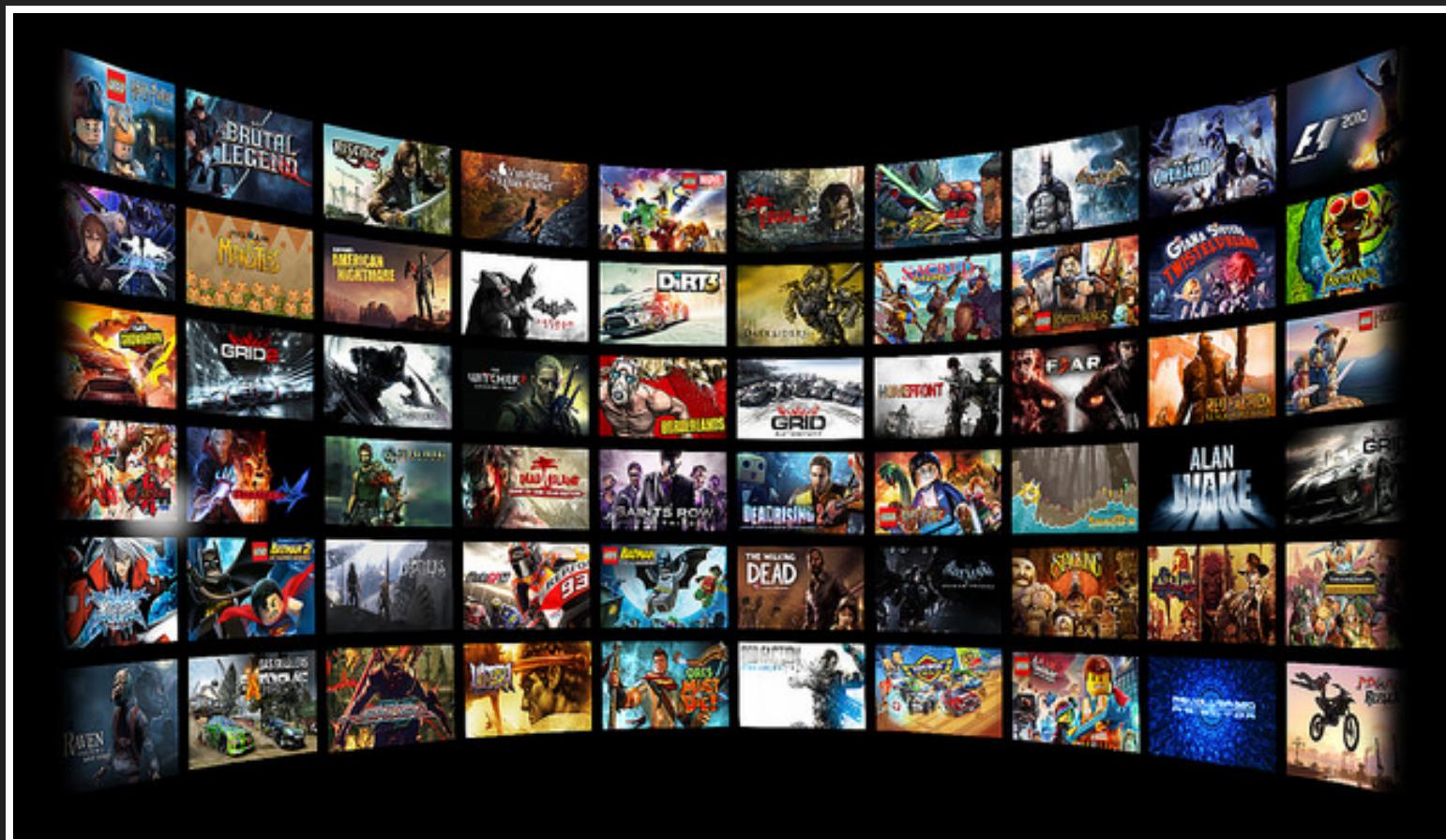
"There is no cloud. It's just someone else's computer"



from [xkcd](#)

CLOUD APPLICATIONS

MULTIMEDIA



BIG DATA



HIGH-PERFORMANCE COMPUTING



LONG-TERM STORAGE



PROBLEMS



RESOURCE ALLOCATION

DIFFERENT TYPES OF HARDWARE



Windows

An error has occurred. To continue:

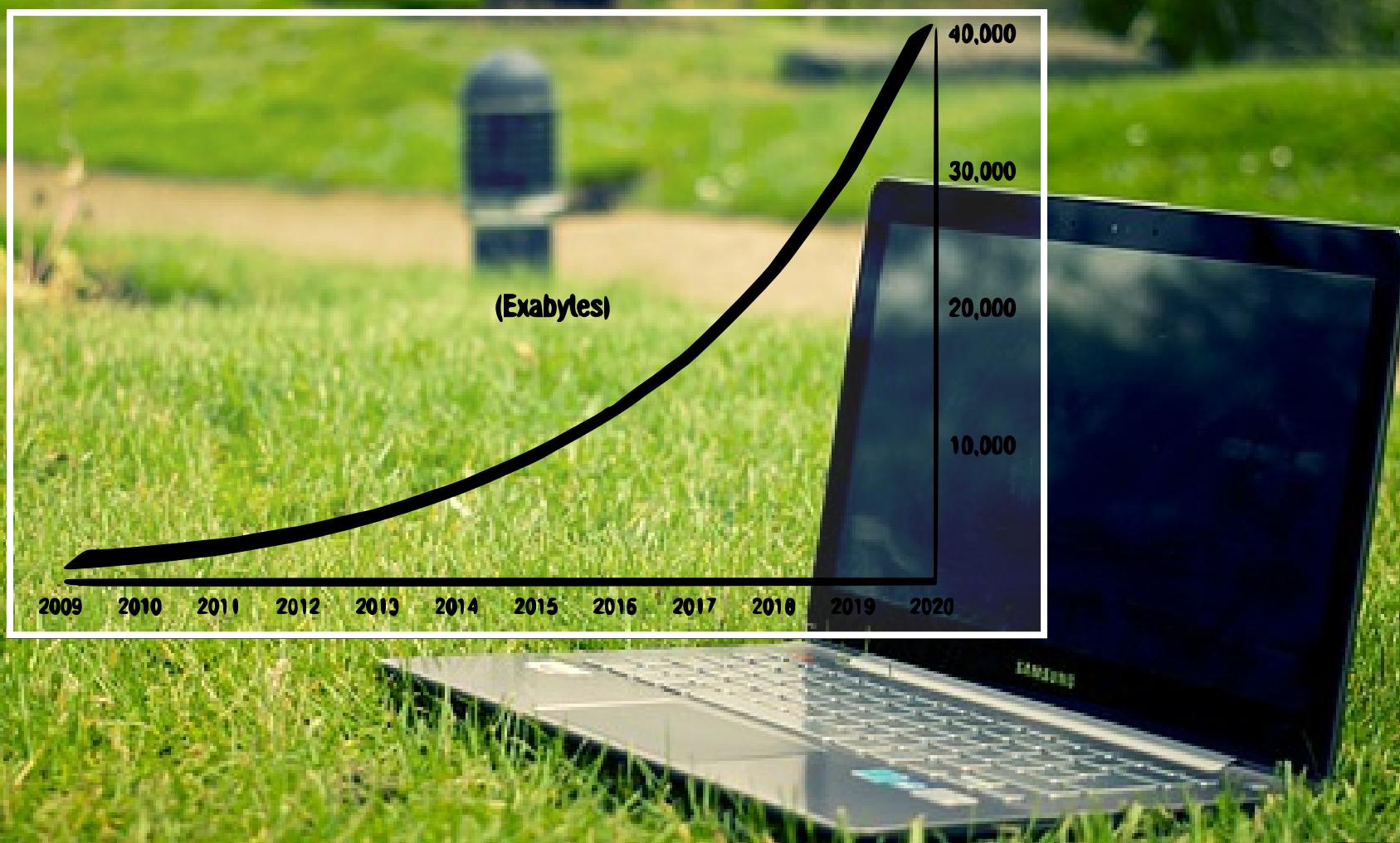
Press Enter to return to Windows

Press CTRL+ALT+DEL to restart your computer. If you do this,
you will lose any unsaved information in all open applications.

Error: 0E : 016F : BFF9B3D4

Press any key to continue _

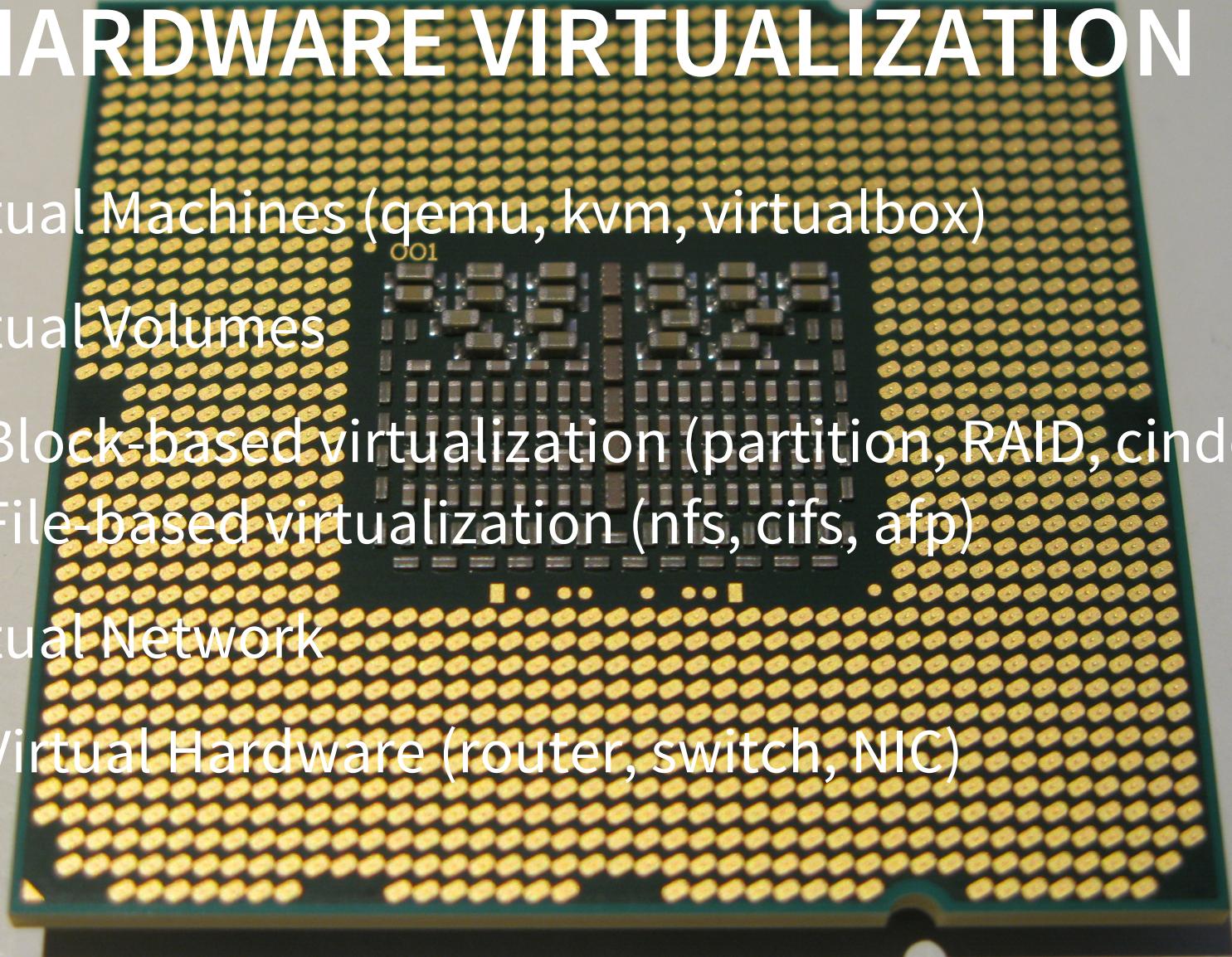
COST / ENERGY EFFECTIVE (\$\$\$)



HOLY GRAIL: VIRTUALIZATION

HARDWARE VIRTUALIZATION

- Virtual Machines (qemu, kvm, virtualbox)
- Virtual Volumes
 - Block-based virtualization (partition, RAID, cinder)
 - File-based virtualization (nfs, cifs, afp)
- Virtual Network
 - Virtual Hardware (router, switch, NIC)



010 Fizians SAS. <<http://www.fizians.com>>
 software; you can redistribute it and/or modify
 terms of the GNU General Public License as published
 by the Free Software Foundation, version 2.
 Distributed in the hope that it will be useful, but
 WITHOUT ANY WARRANTY; without even the implied warranty of
 MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GN
 License for more details.
 received a copy of the GNU General Public License
 along with this program. If not, see
<http://www.gnu.org/licenses/>.

SOFTWARE-DEFINED VIRTUALIZATION

- Software-Defined Storage (SDS)
 - File Systems (rozofs, cephfs, hdfs)
- Software-Defined Network (SDN)
 - Virtual elements (firewall, load balancer)
 - OpenFlow (protocol to access forwarding planes)
 - Virtual LAN

```

189 }
190 /*-----*/
191 */
192 /**
193 * Perform a mojette transform forward
194
195 @param thread_ctx_p: pointer to the thread context
196 @param msg : address of the message received
197
198 */
199 static inline void storcli_mojette_forward(rozofs_mojette_thread_ctx_t
200 *thread_ctx_p,rozofs_stcmoj_thread_msg_t * msg)
201 {
202     struct timeval timeDay;
203     unsigned long long timeBefore, timeAfter;
204     unsigned long long cycleBefore, cycleAfter;
205     rozofs_storcli_ctx_t * working_ctx_p;
206     rozofs_storcli_ingress_write_buf_t *wr_proj_buf_p;
207     storcli_write_arg_no_data_t *storcli_write_rq_p;
208     uint8_t layout;
209     int i;
210     int block_count = 0;
211
212     gettimeofday(&timeDay, (struct timezone *)0);
213     timeBefore = MILLEND(timeDay);
214     cycleBefore = rdtsc();
215
216     /*
217      ** update statistics
218      */
219     thread_ctx_p->stat.MojetteForward_count++;
220     working_ctx_p = msg->working_ctx;
221     wr_proj_buf_p = working_ctx_p->wr_proj_buf;
222     storcli_write_rq_p = &working_ctx_p->storcli_write_arg;
223     Layout = storcli_write_rq_p->layout;
224
225     /*
226      * This function is triggered by OpenFlow when a transform must
227      */
228     for (i = 0; i < ROZOFS_WR_MAX; i++)
229     {
230         if (wr_proj_buf_p[i].state == ROZOFS_WR_ST_TRANSFORM_REQ)
231         {
232             block_count += wr_proj_buf_p[i].number_of_blocks;
233             rozofs_storcli_transform_forward(working_ctx_p->prj_ctx,

```

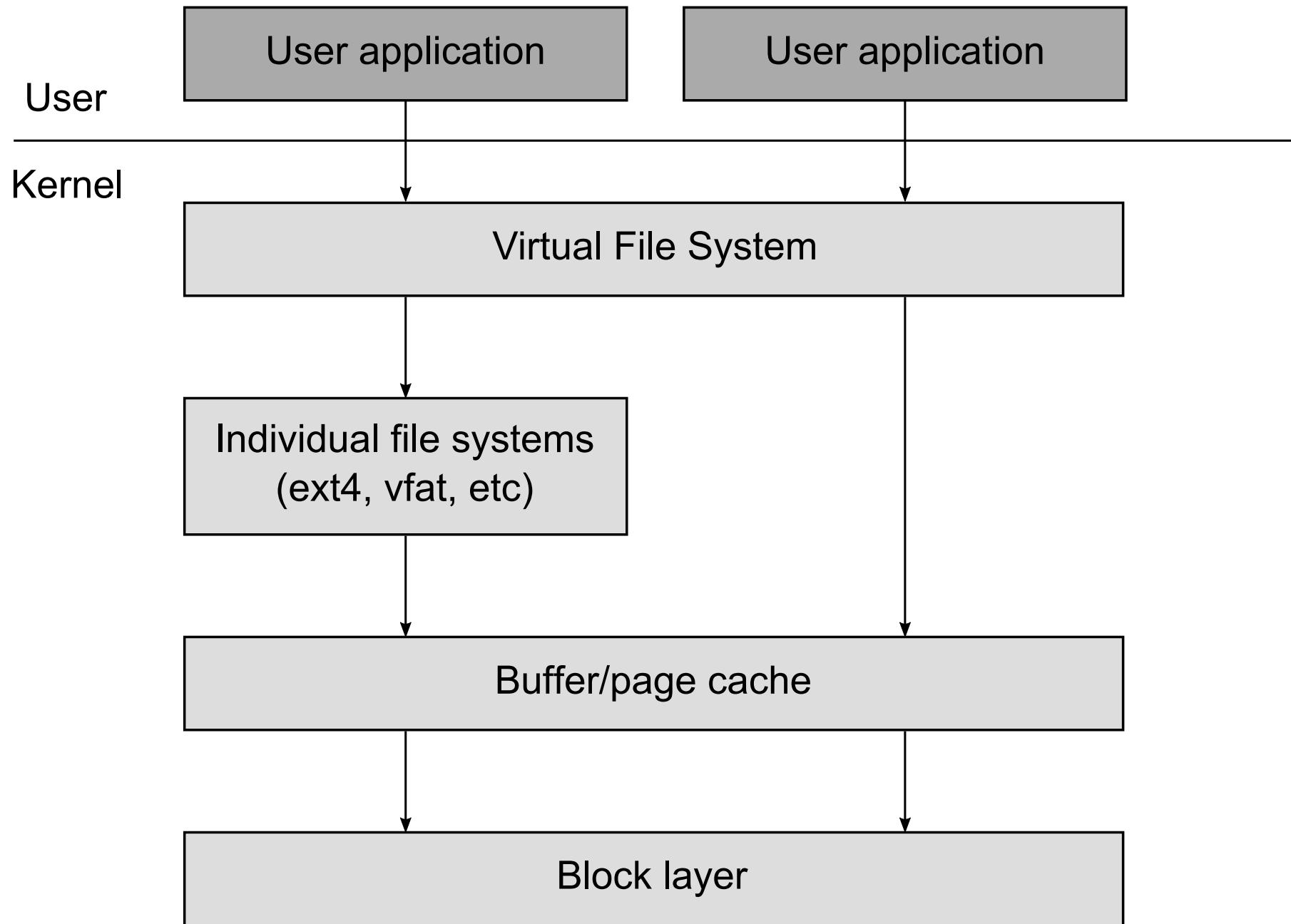
OUTLINE

1. File System Basics
2. Scale-out Network Attached Storage (NAS)
3. Example with RozoFS: a Distributed File System
4. My storage in the Cloud (devops)

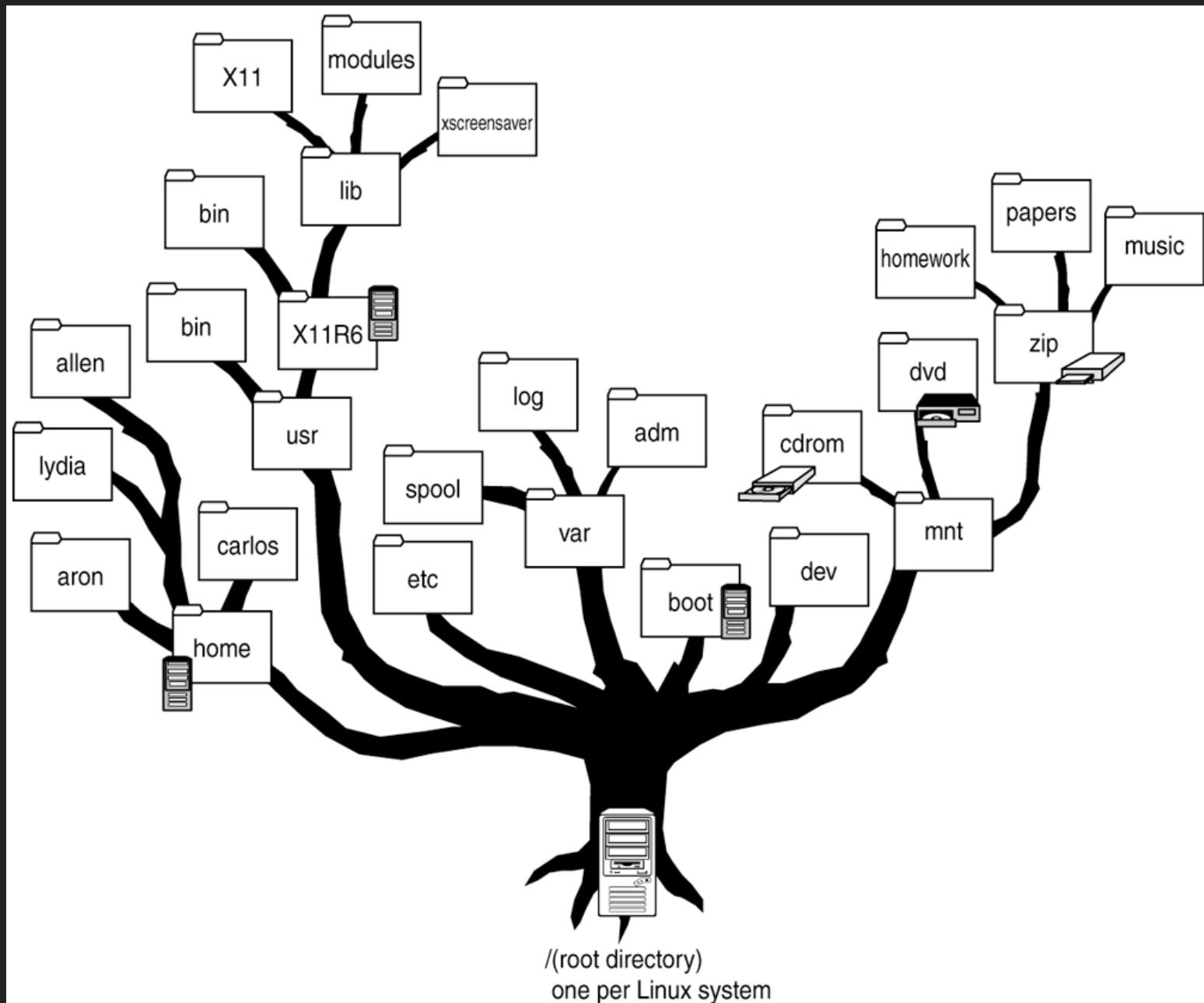
FILE SYSTEM BASICS

DATA REPRESENTATIONS AND ACCESSES

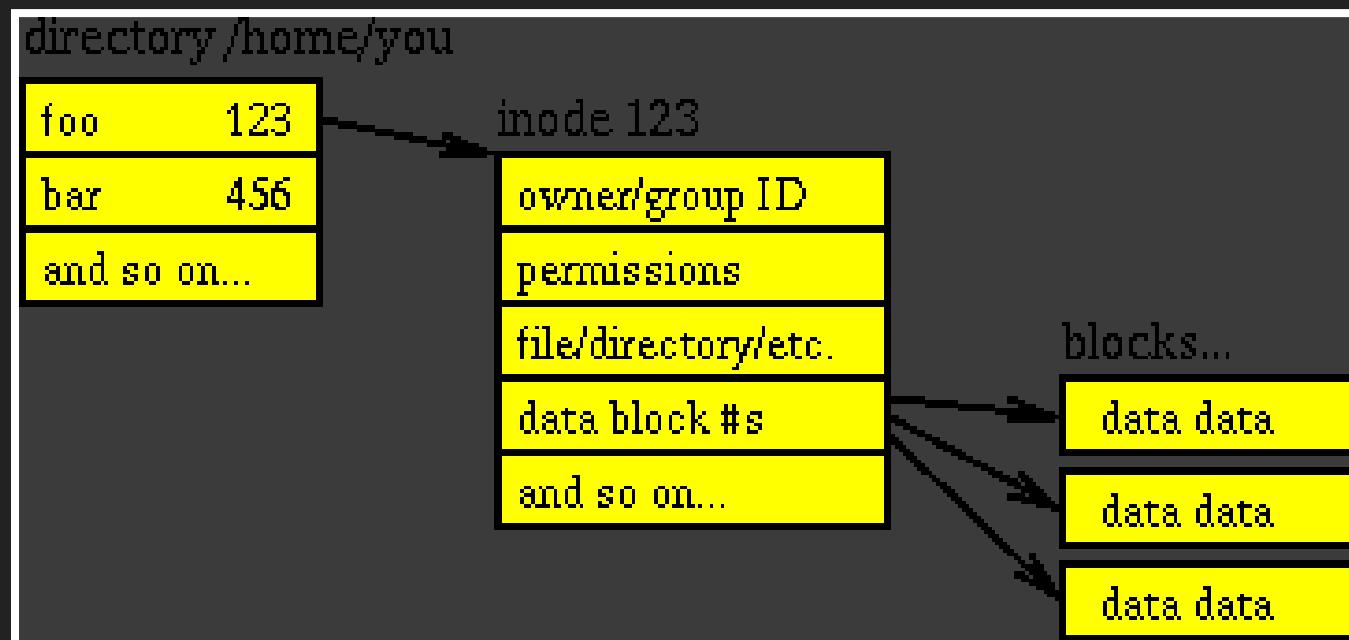
- Block storage
- File systems



TREE-LIKE STRUCTURE



INODE



BLOCK STORAGE

Close to the storage hardware. Consider raw volumes of storage

- (-) requires specific implementations
- (+) more efficient, specific features

used for instance in databases and virtual machines

FILE SYSTEMS

Most commonly used in storage systems

Contained in block devices (*mkfs()*)

- (-) level of abstraction, worse perf
- (+) simple to use, no need to re-implement softwares

e.g. ext4, xfs, btrfs

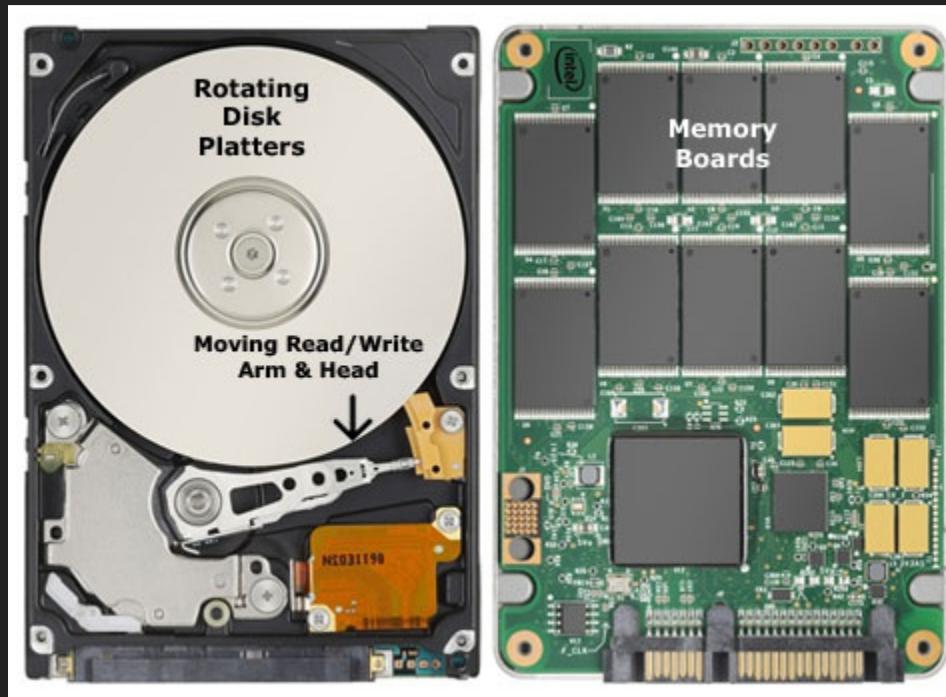
APPLICATION WORKLOADS

ACCESS MODES

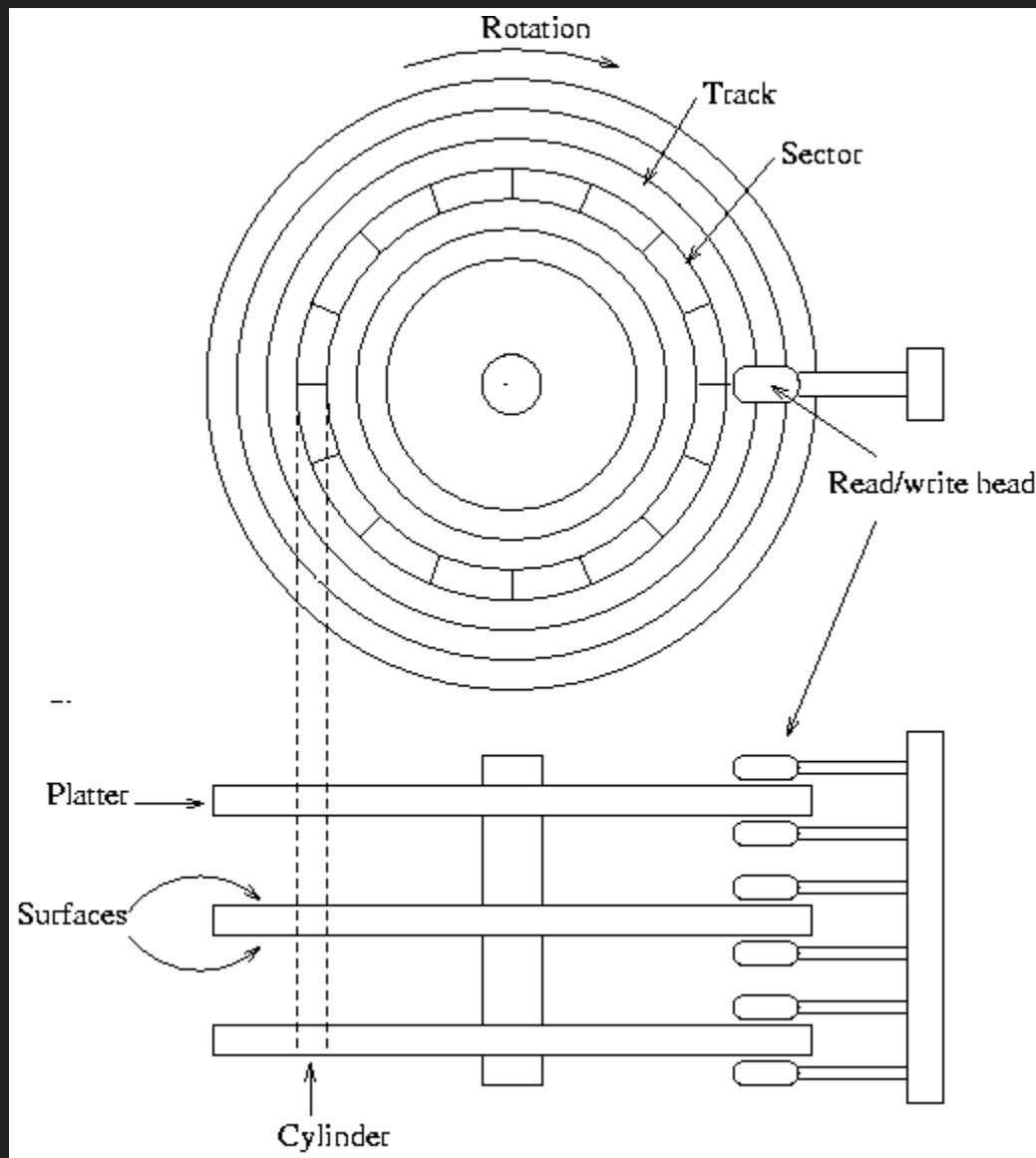
1. Type of operations
 - write
 - read
2. Access pattern
 - sequential
 - random
3. The workload depends on the application
 - video streaming: sequential over large blocks (64 KB)
 - databases: random over small blocks (4 KB)

STORAGE MEDIA

HARD DRIVE DISK (HDD)



HARD DRIVE DISK (HDD)



SOLID-STATE DRIVE (SSD)

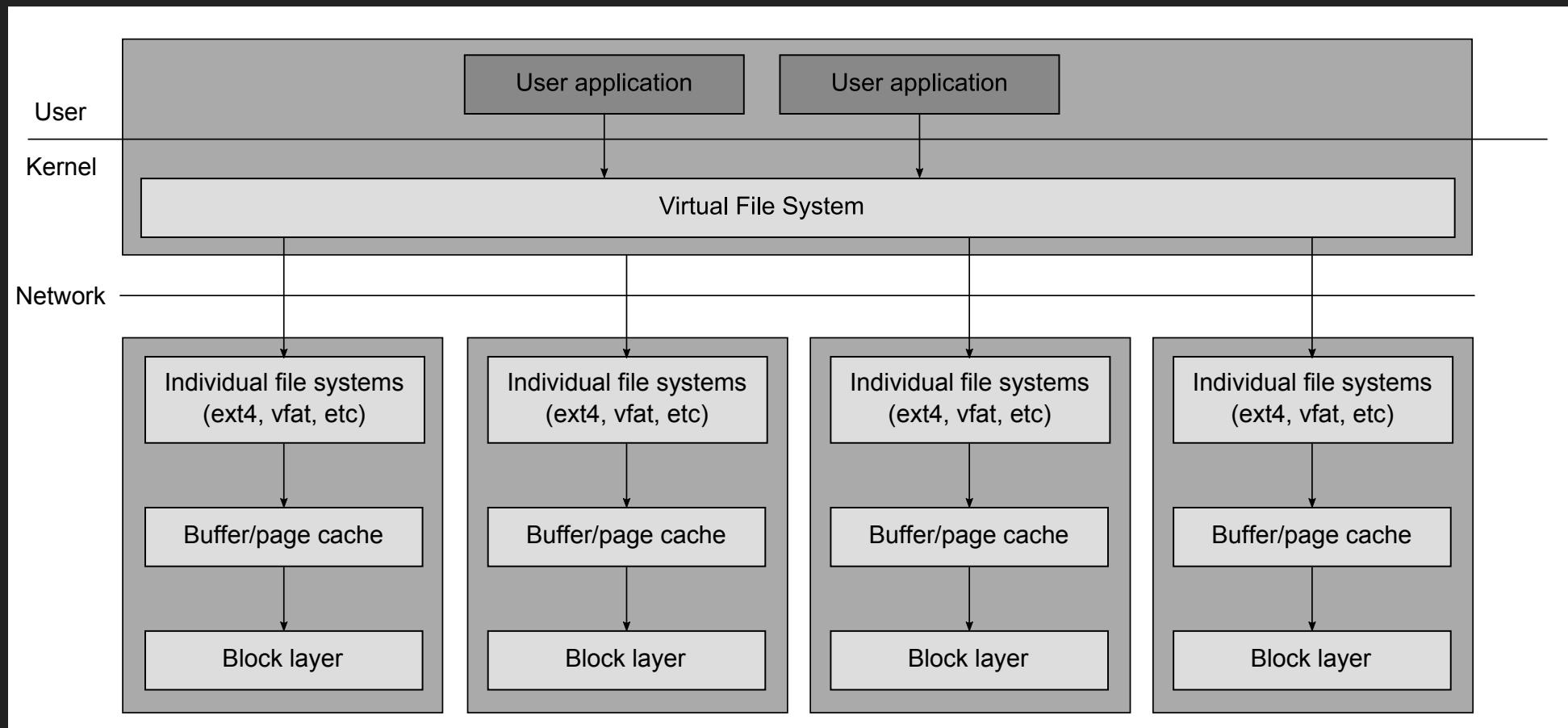
DATA SHARING OVER THE NETWORK

- Storage Area Network (SAN)
 - access block devices through network
 - *Fibre Channel* (expensive, dedicated hw)
 - iSCSI (SCSI over TCP/IP)
- Network Attached Storage (NAS)
 - file system access through network
 - NFS, CIFS, AFP

SCALE-OUT NETWORK ATTACHED STORAGE

ABSTRACTION ON TOP OF NATIVE FILE SYSTEMS

DISTRIBUTED ARCHITECTURE

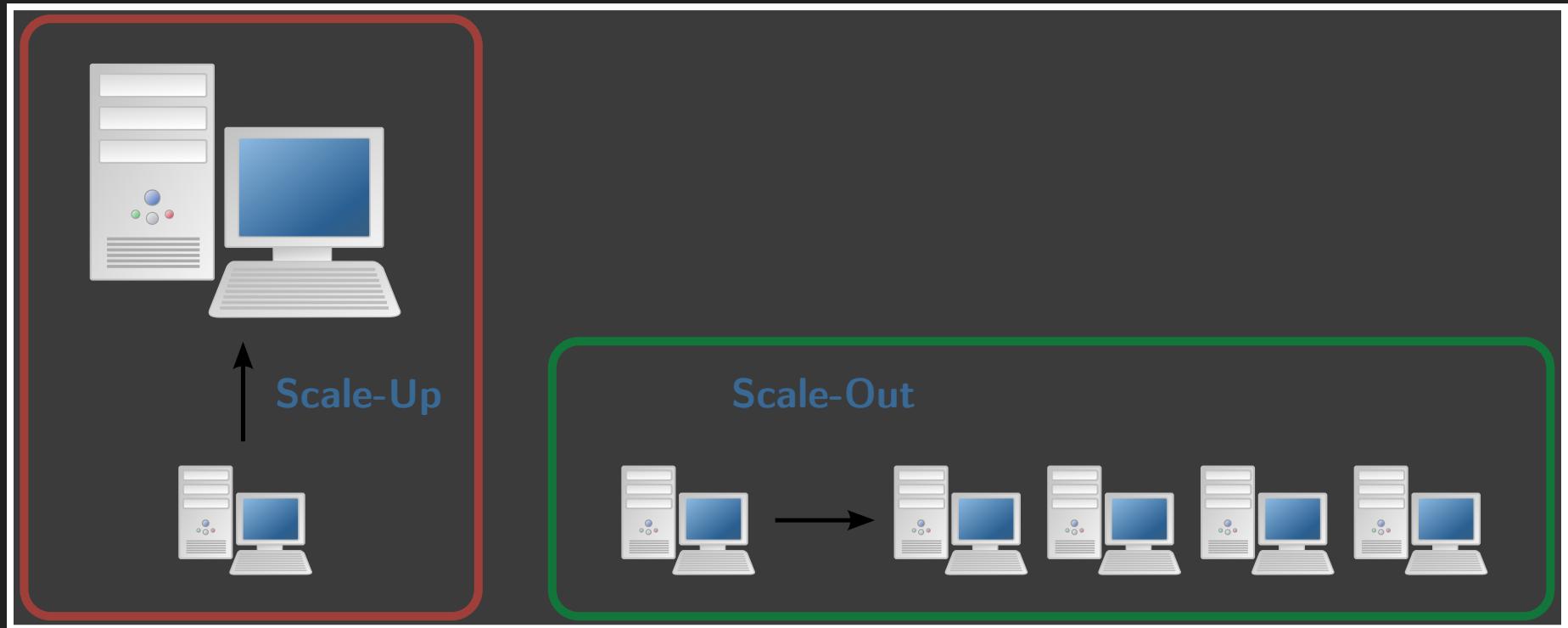


ADDRESS SPACE REMAPPING

- location independence by abstracting physical location
- virtualization system presents to user a logical space
- provides a unique namespace
- mapping for data distribution
- metadata managed within a metadata server

SCALABILITY

SCALE-UP VS SCALE-OUT



METADATA MANAGEMENT

Different implementations

- mapping table
- compute locations using an algorithm (CRUSH)

DATA MANAGEMENT

- The virtualization software redirects I/O requests
 1. receive an I/O to a location on logical disk (LD)
 2. translate location on LD into location on physical disk
 3. translate the request
- Data distribution

NEW DATA REPRESENTATION: OBJECT

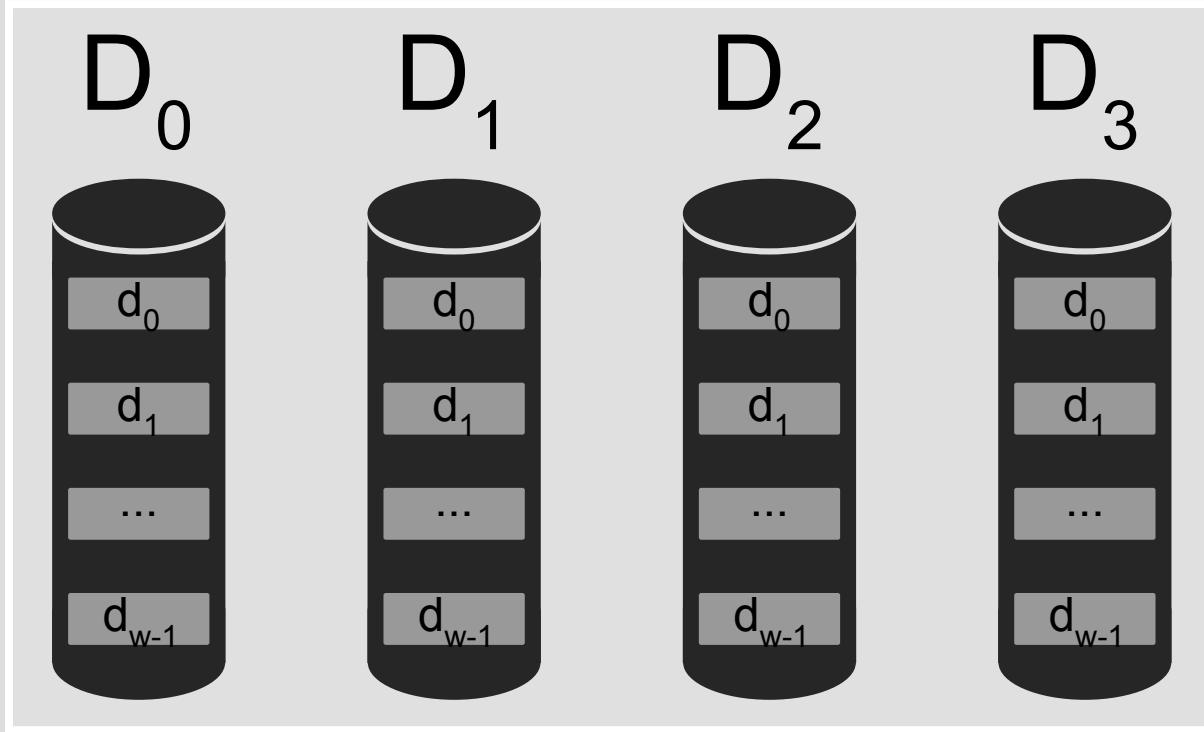
Key/Value store

FAULT-TOLERANCE

A vast expanse of yellow rubber ducks floating on blue water. The ducks are arranged in several large, dense clusters, with many more scattered individually across the surface. Some ducks have black sunglasses and pink beaks.

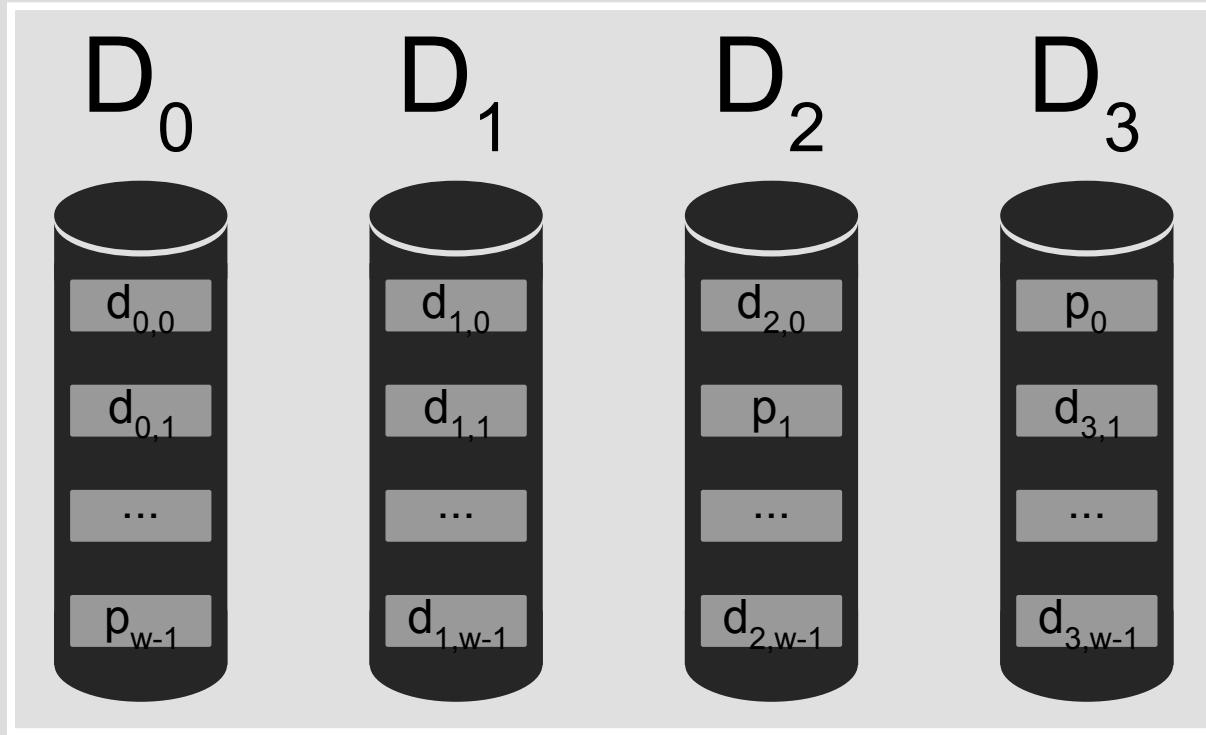
HIGH-AVAILABILITY BY REDUNDANCY

RAID-1: REPLICATION



Four-way replication here implies significant storage overhead of 300%

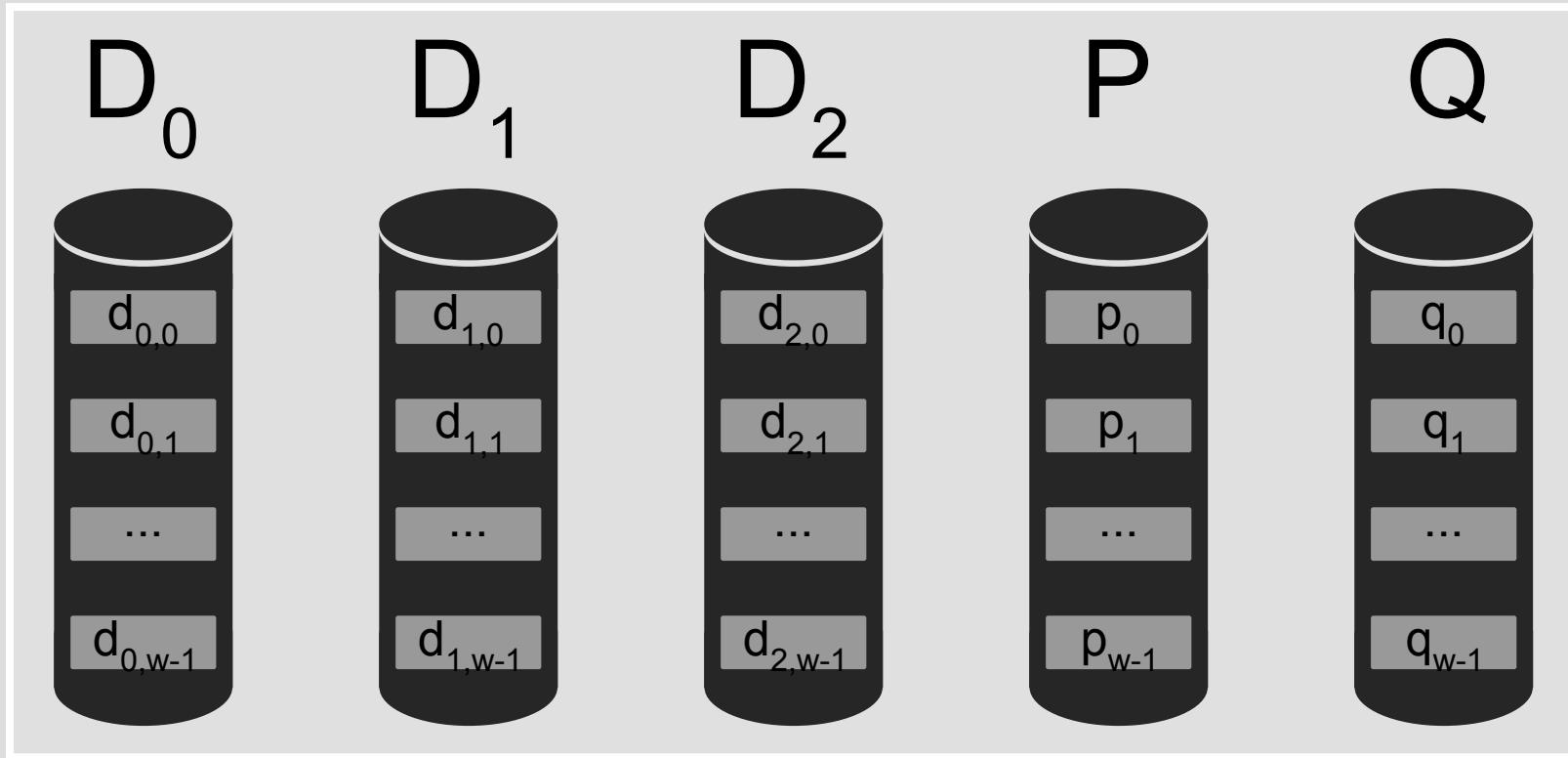
RAID-4: SINGLE PARITY



The parity disk protect the system against a single failure.

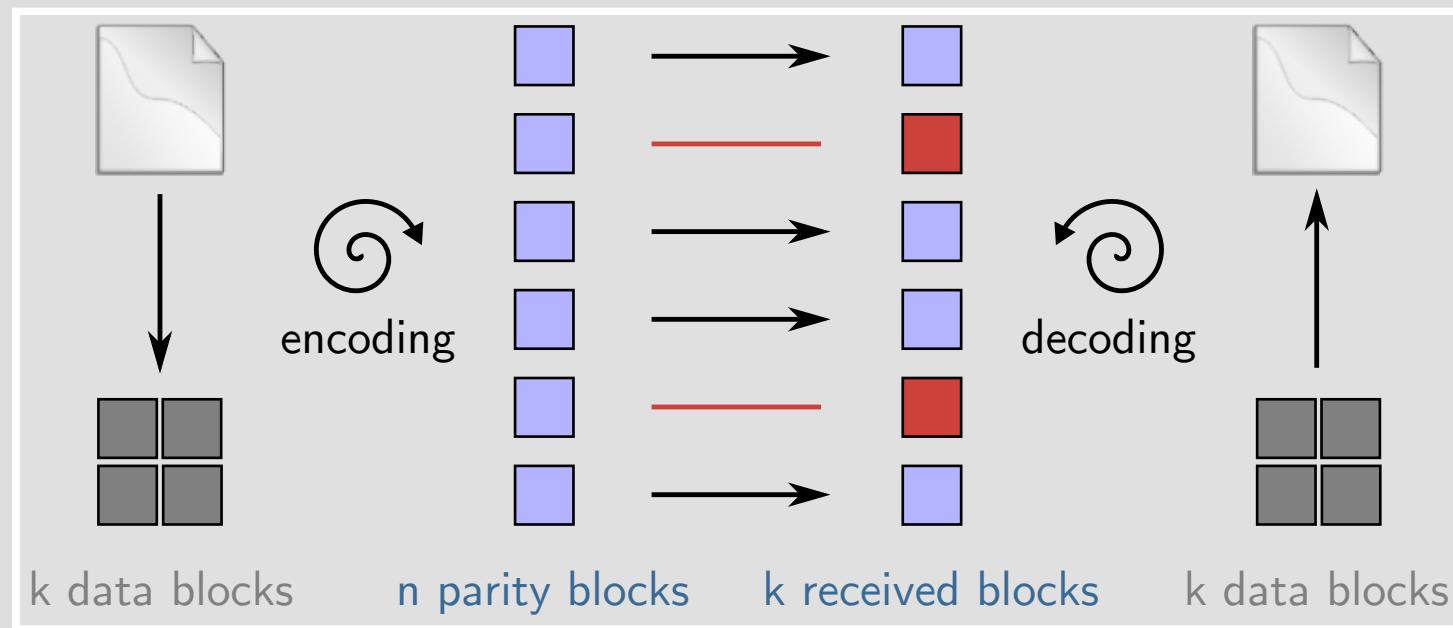
$$p_j = \bigoplus_{i=0}^{k-1} d_{i,j}$$

RAID-6: DOUBLE PARITY



The two parity disks protect the system against two failure.
Multiple implementations exist to compute Q.

BEYOND: (N,K) ERASURE CODE

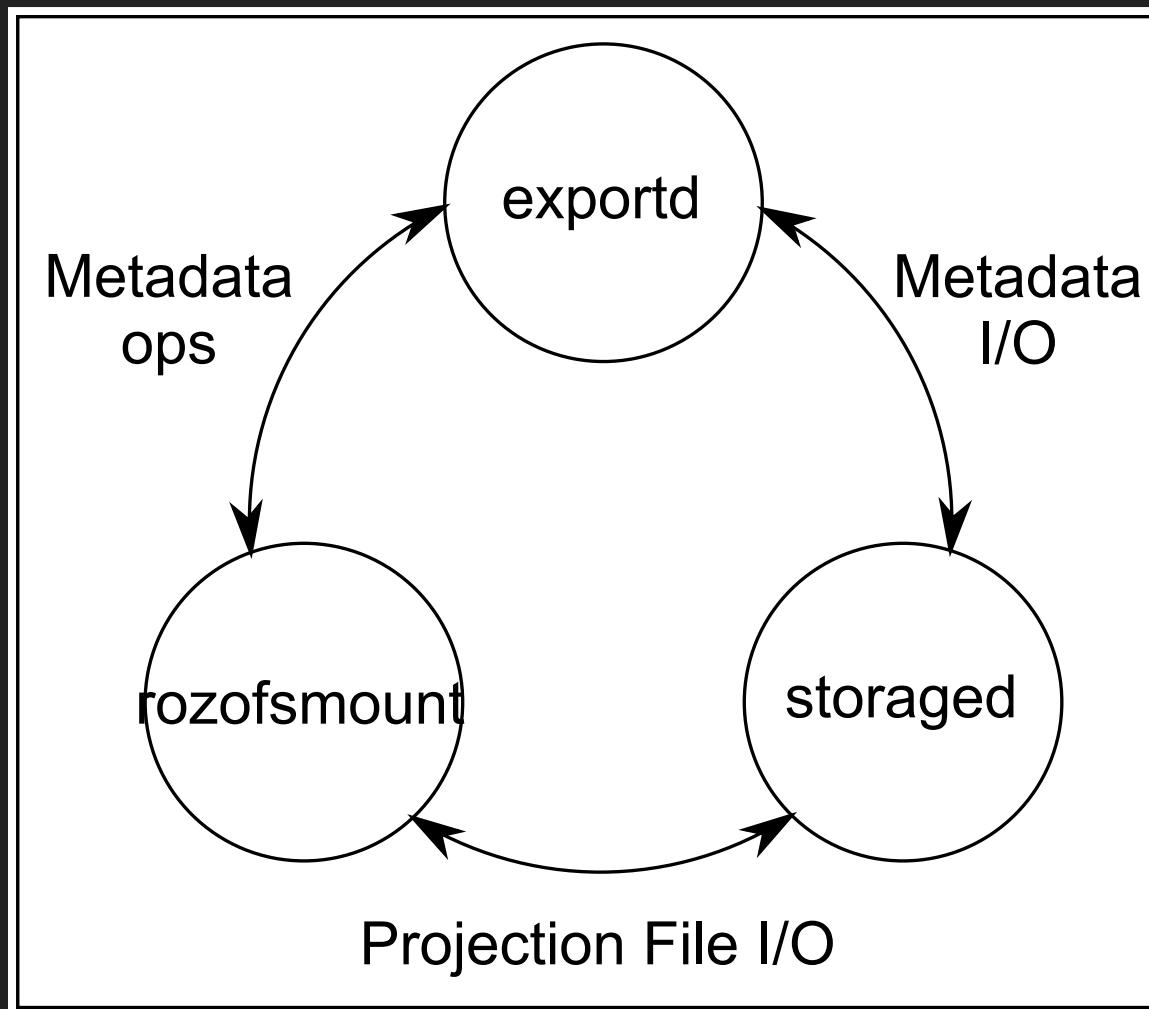


An $(n = 12, k = 8)$ erasure code can protect against up to four failures while providing only 50% of storage overhead

ROZOFs: SCALE-OUT NAS OPEN-SOURCE DISTRIBUTED FILE SYSTEM

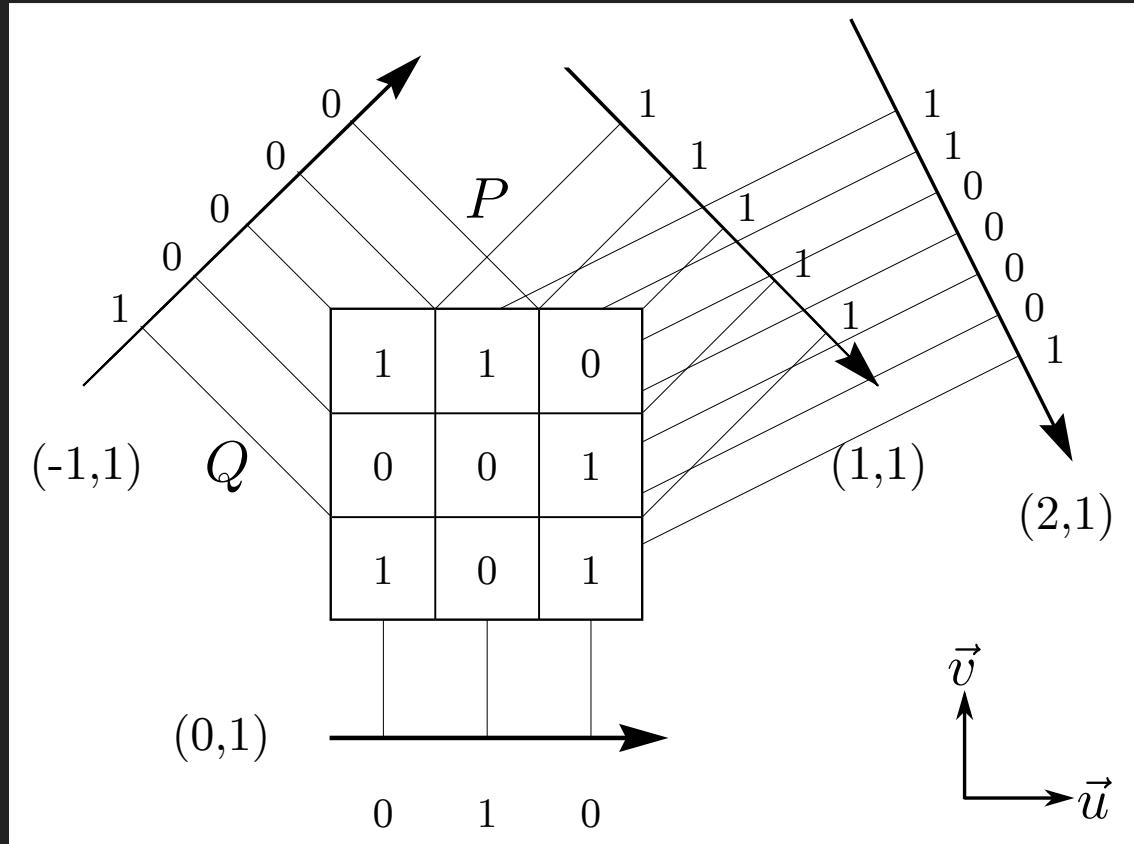
ROZOFS ARCHITECTURE

THREE COMPONENTS



FAULT-TOLERANCE BY THE MOJETTE ERASURE CODE

MOJETTE TRANSFORM



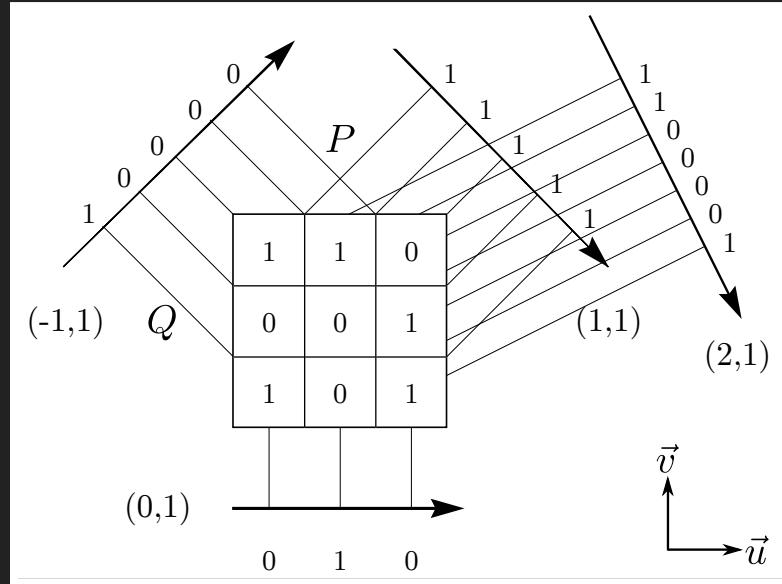
$$[\mathcal{M}f](b, p_i, q_i) = \sum_{k=0}^{Q-1} \sum_{l=0}^{P-1} f(k, l) \Delta(b + kq_i - lp_i)$$

KATZ CRITERION

$$\sum_{i=1}^I |p_i| \geq P \text{ or } \sum_{i=1}^I |q_i| \geq Q$$

Following simplification: $|q_i| = 1$

(n, k) MOJETTE ERASURE CODE

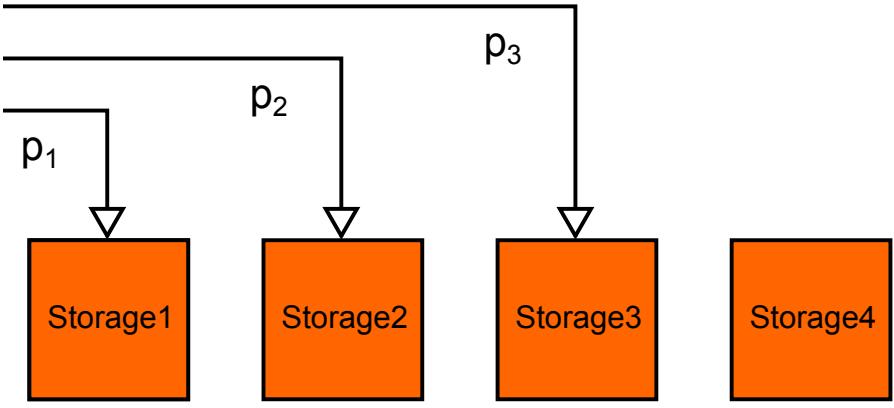


1. consider n projections, computed from $(k = Q)$ -height grid
2. reconstruct from k projections
3. protection against $(n - k)$ failures

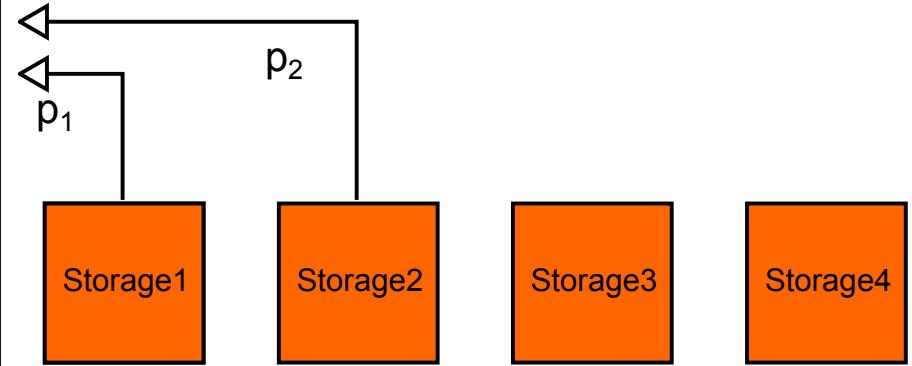
DATA I/O

READ AND WRITE SCENARIOS

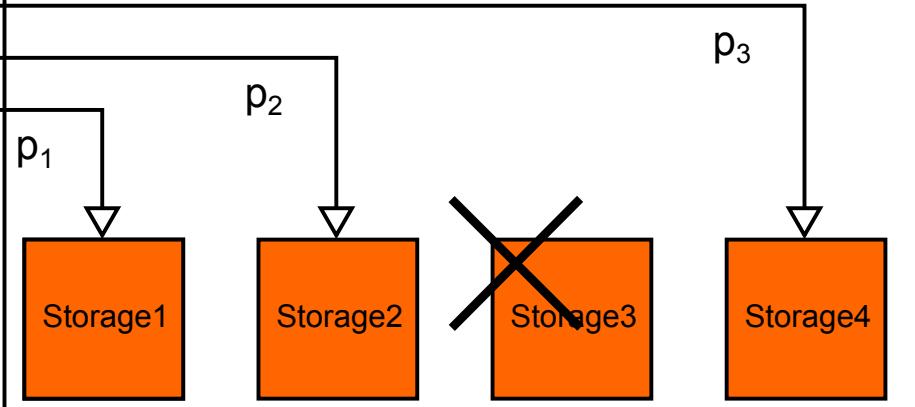
Fault-Free Write



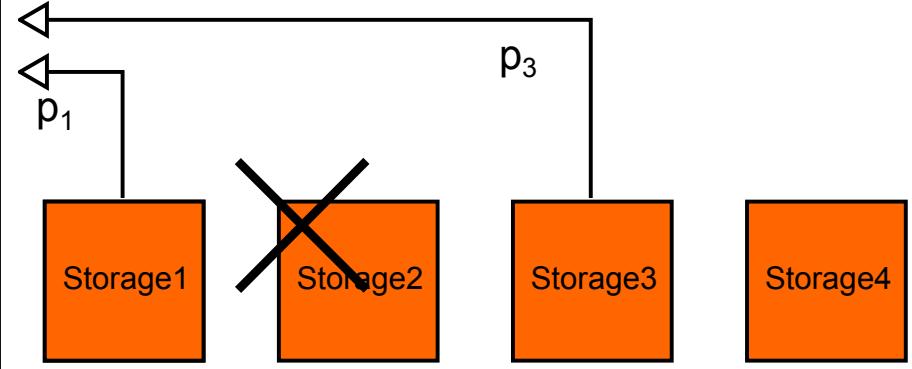
Fault-free Read



Degraded Write



Degraded Read



ROZOFS INTERACTIONS

