

Cross the boundary

Hung-Wei Tseng

Recap: From the software perspective

In-Storage Processing

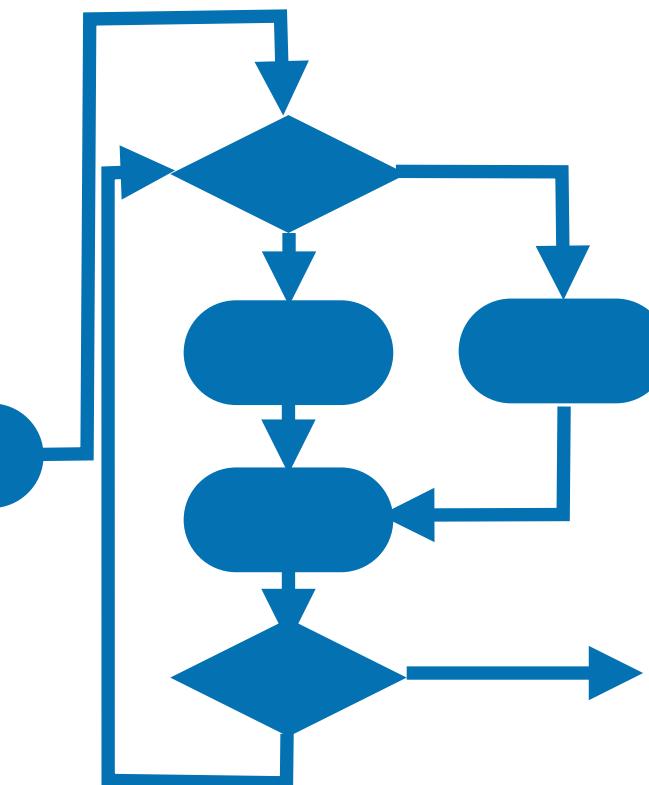
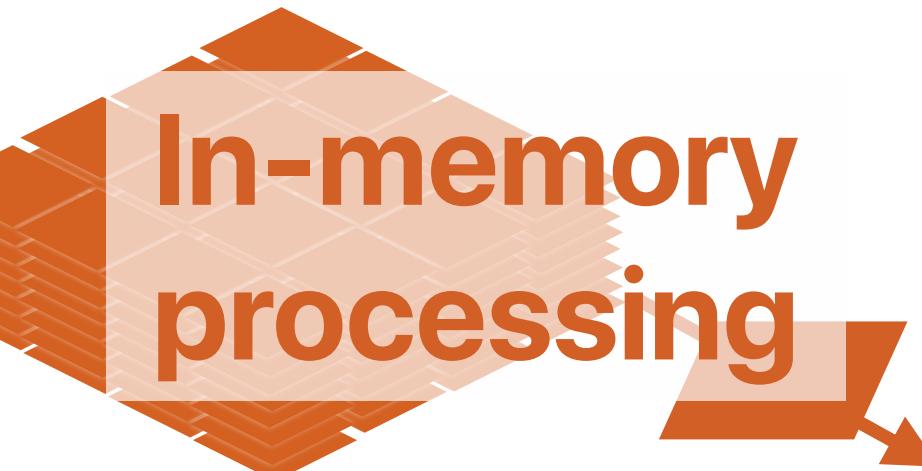
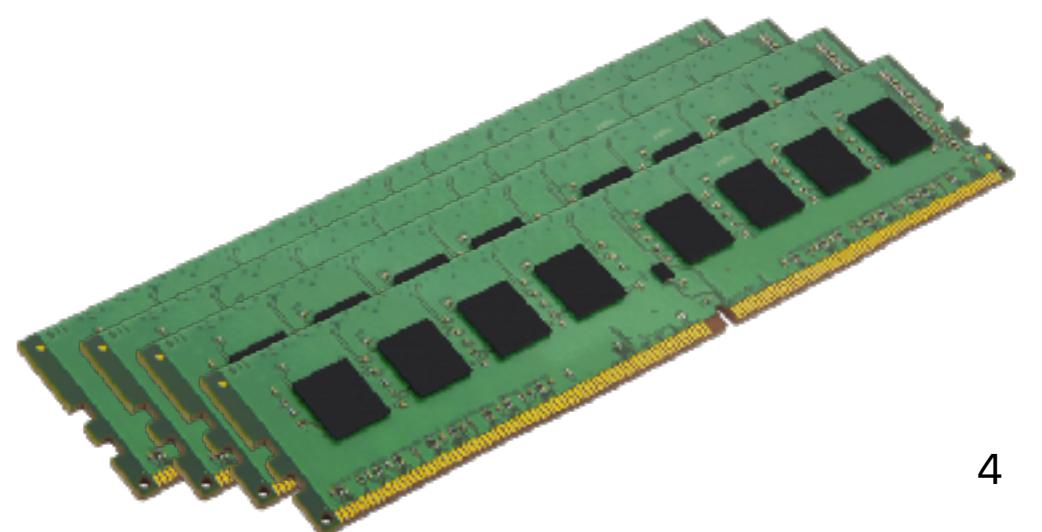


How about this process?

Source "data"

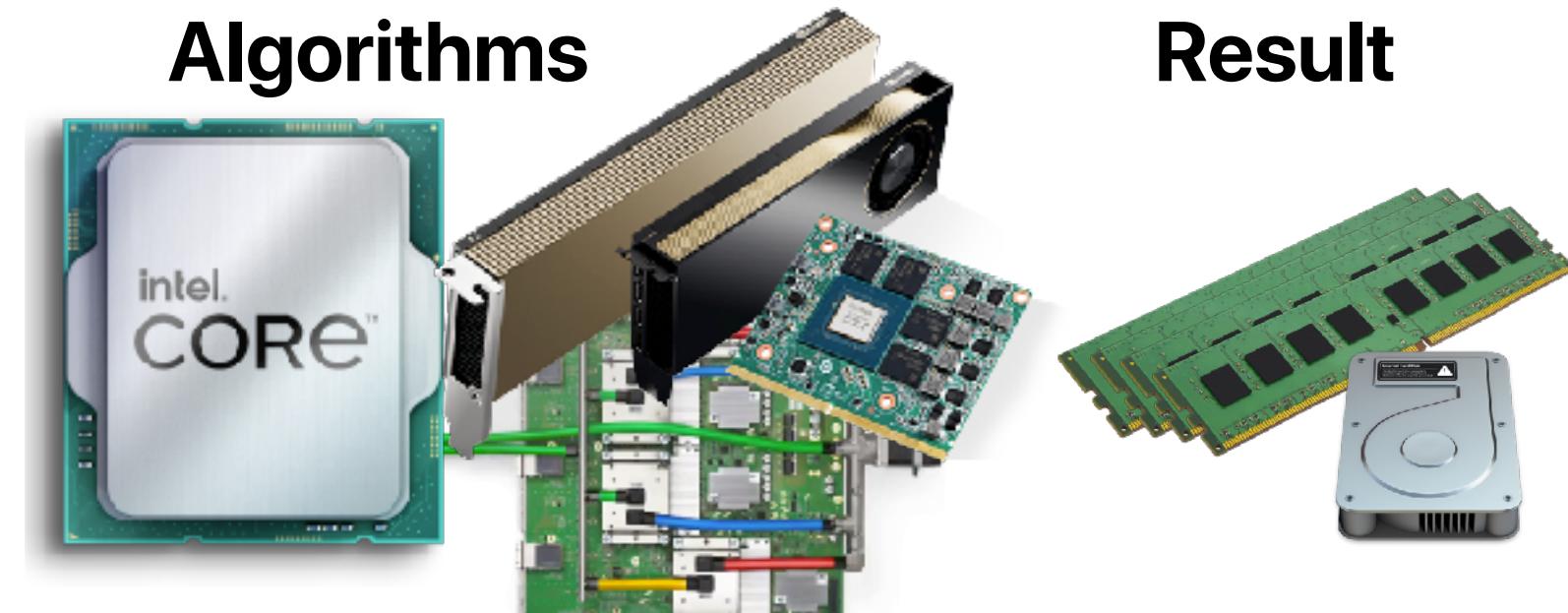


"Data" structures



Hardware Accelerators

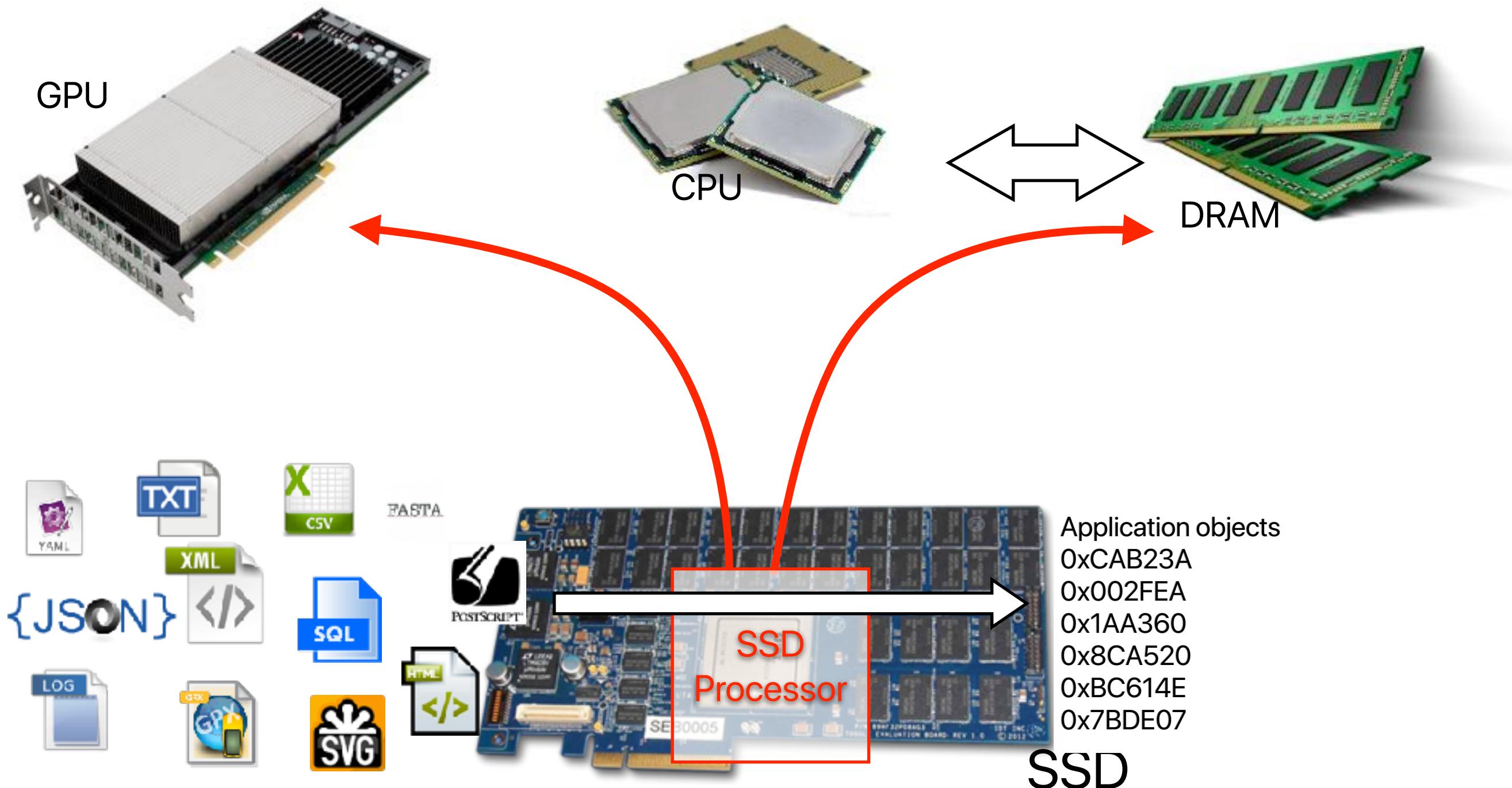
Algorithms



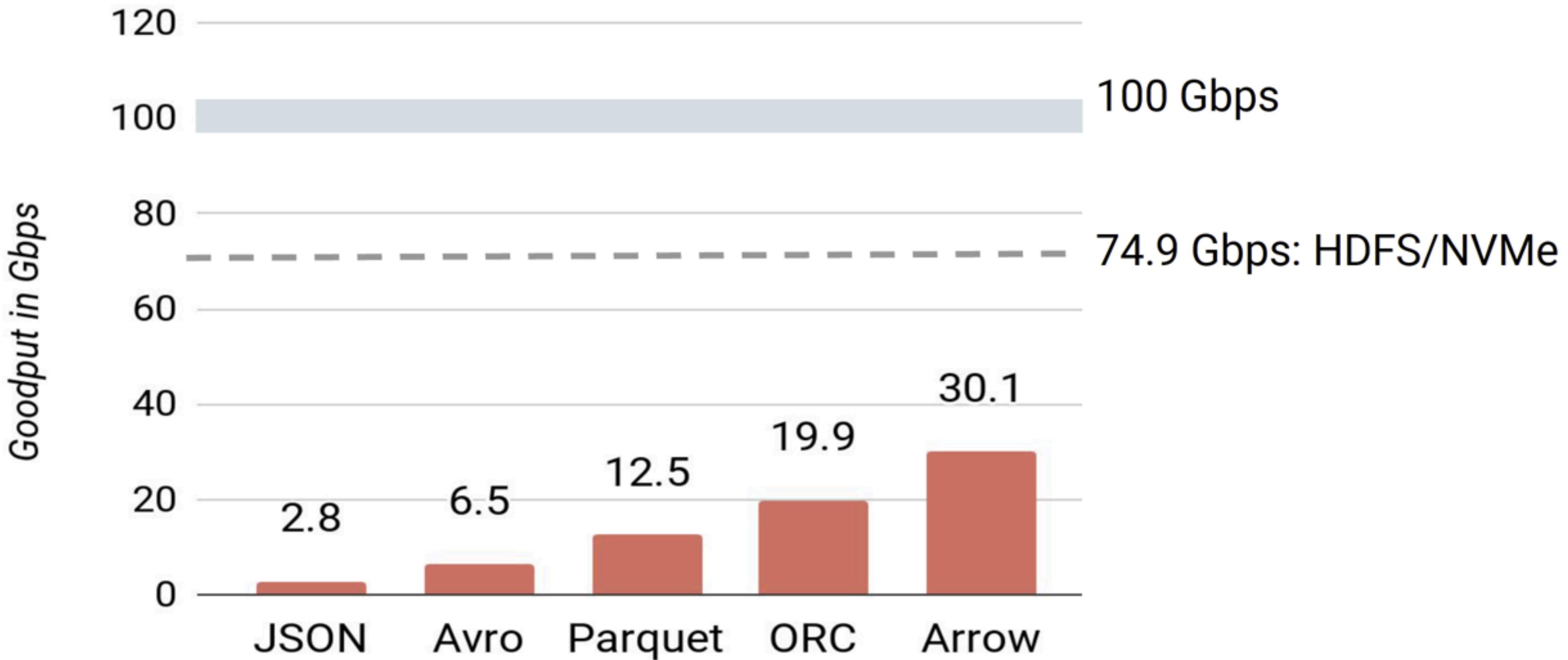
Result



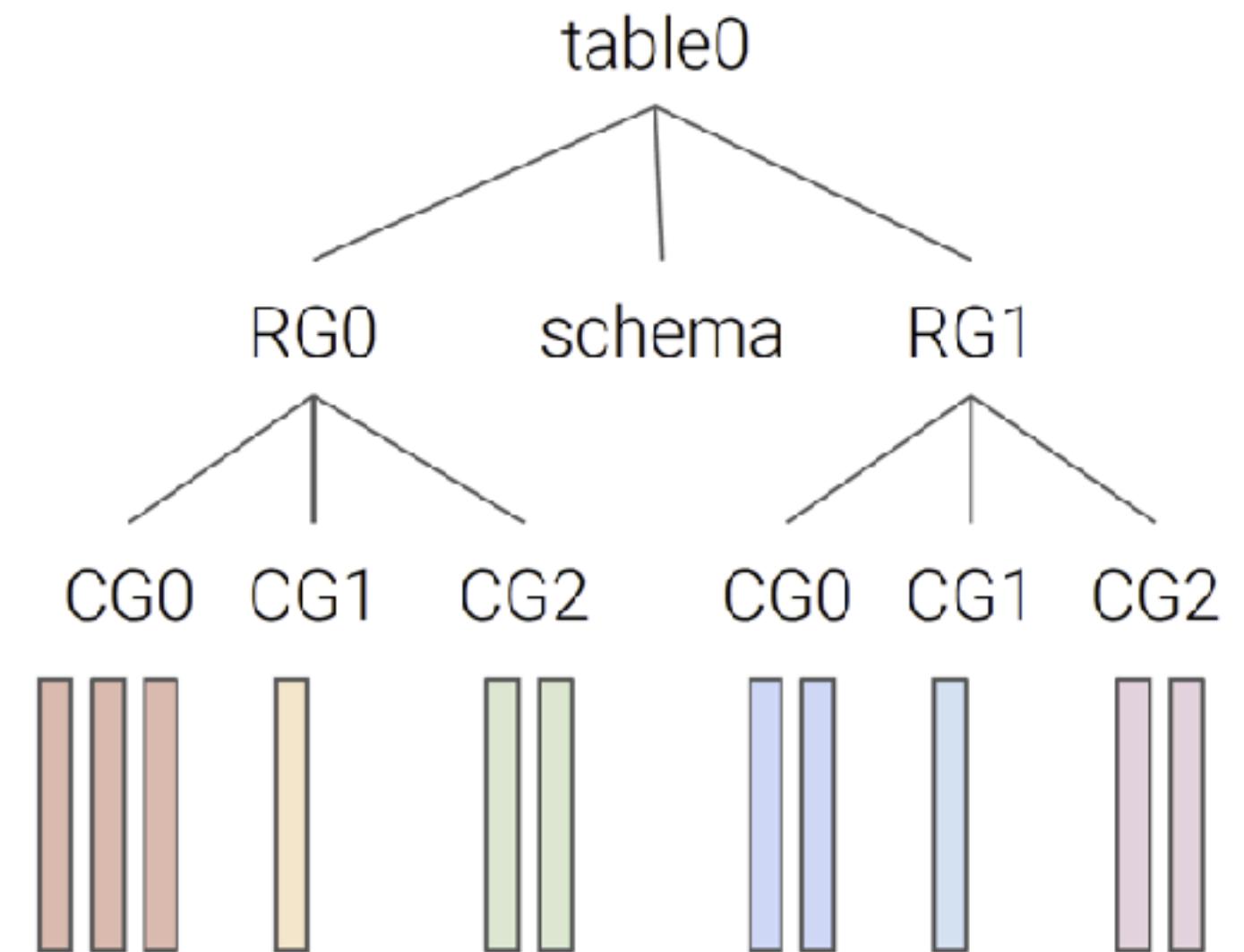
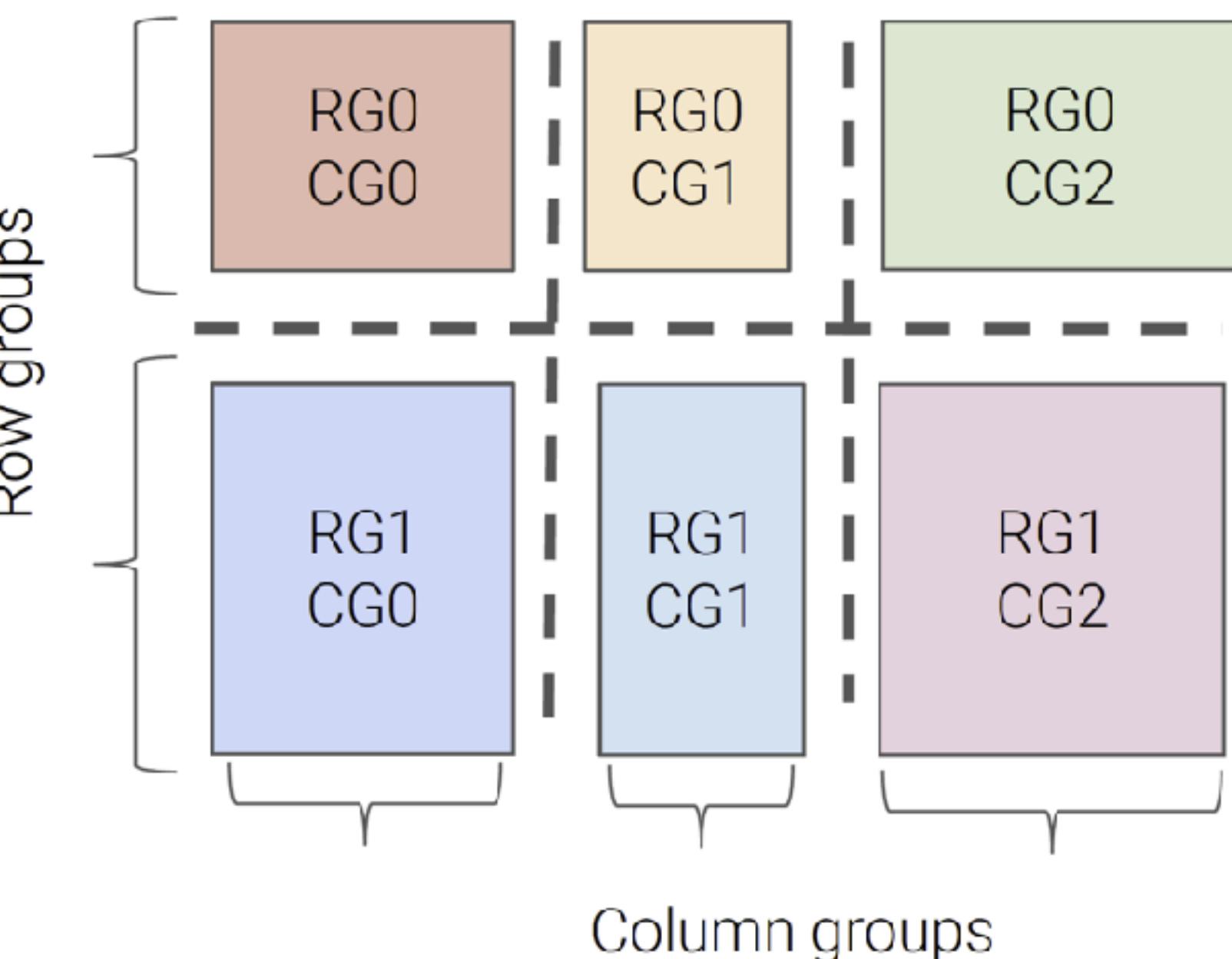
Morpheus



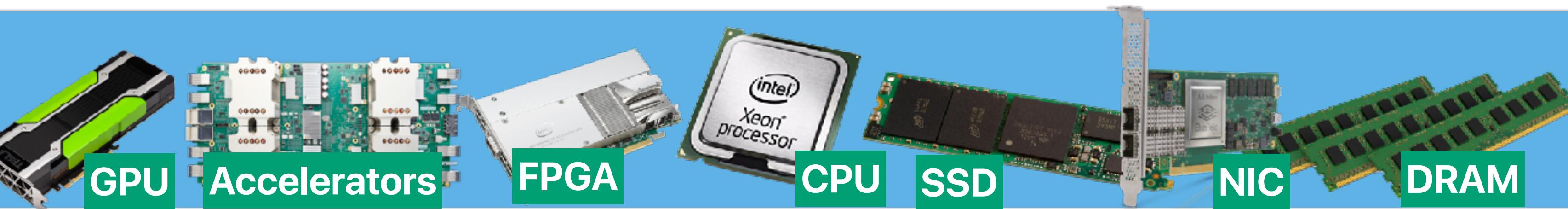
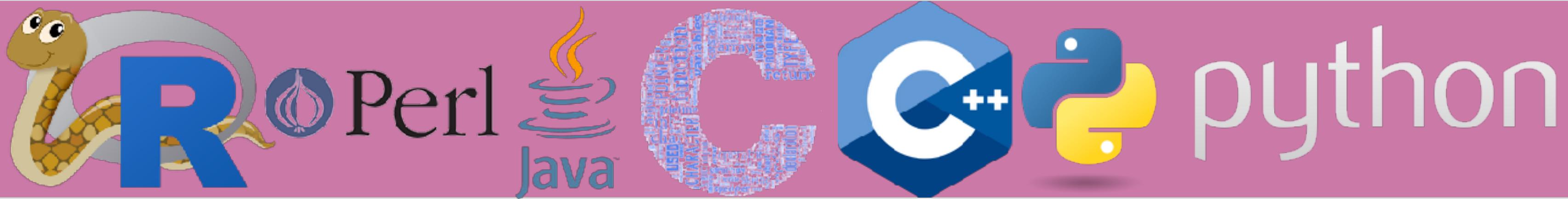
The deserialization performance on popular formats



Albis: Not just row groups, but also column groups



DO NOT CROSS



Is layering design a good approach? Why or why not?

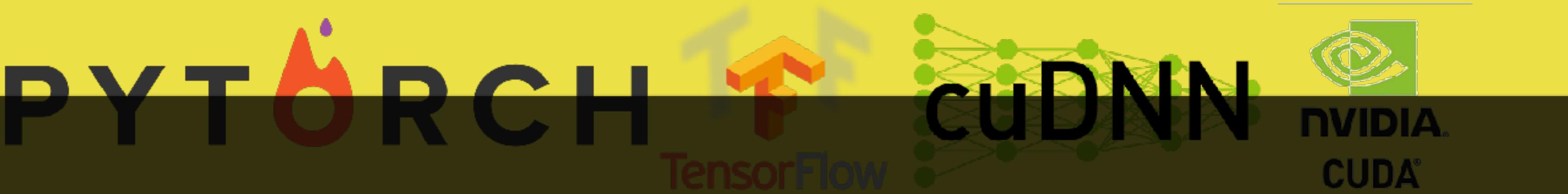
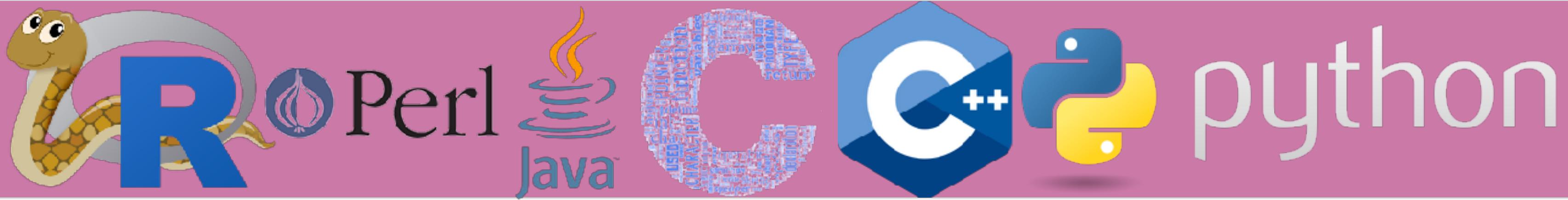
Outline

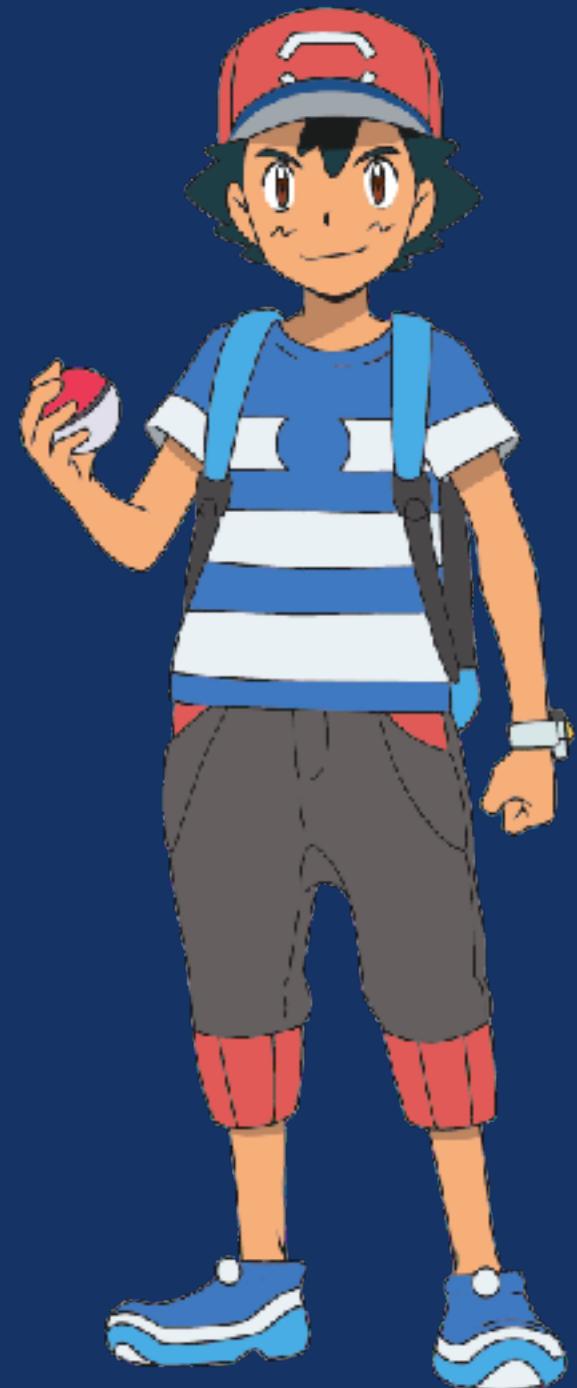
- Pros and cons of layering
- Cross-layer design example
- Pros and cons of approximate computing
- Cross-layer design for approximate computing

Pros/Cons of Layering

Pros/Cons of Layering

- Good
 - Ease of Debugging (T.H.E. by Dijkstra)
- Bad
 - Performance
 - Sub-optimal optimizations

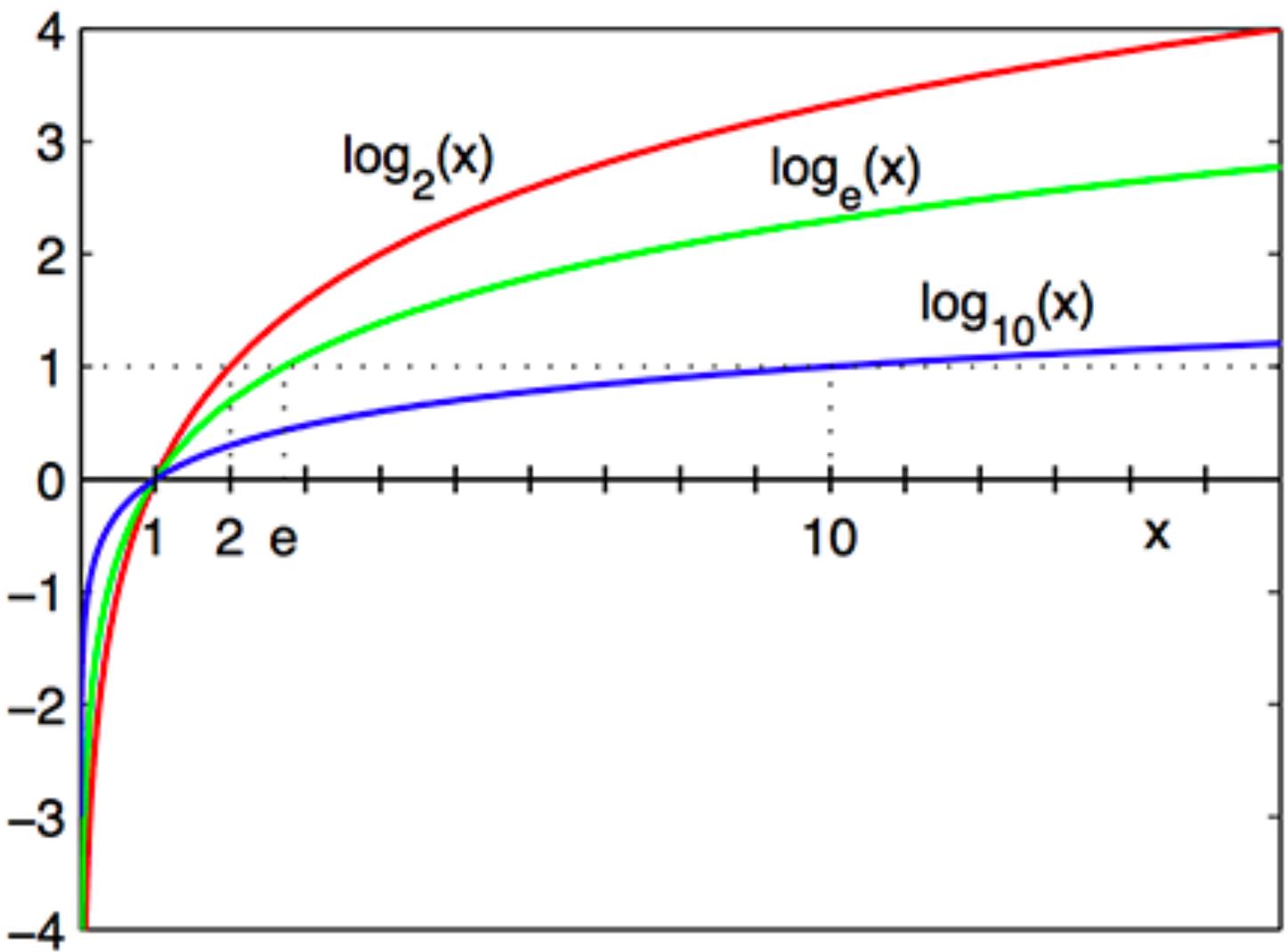




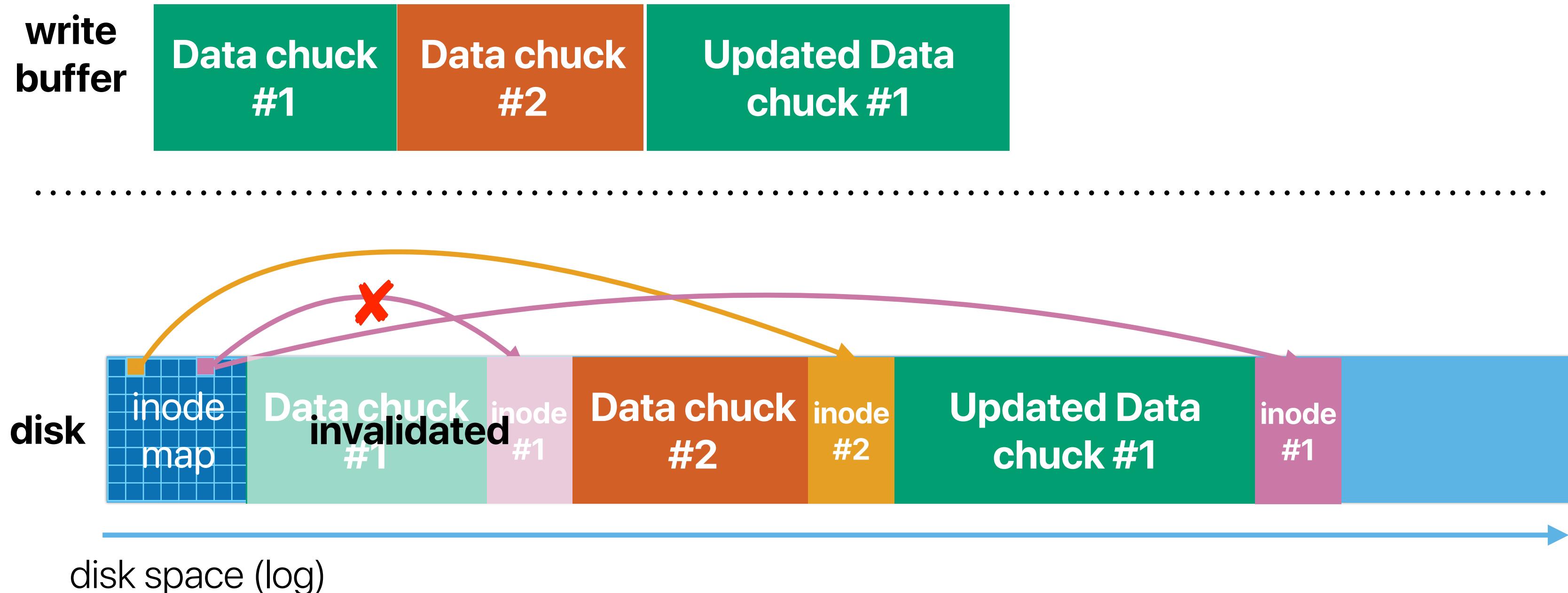
?????

Don't stack your log on my log

**Jingpei Yang, Ned Plasson, Greg Gillis, Nisha Talagala, and
Swaminathan Sundararaman
SanDisk Corporation**



LFS in motion

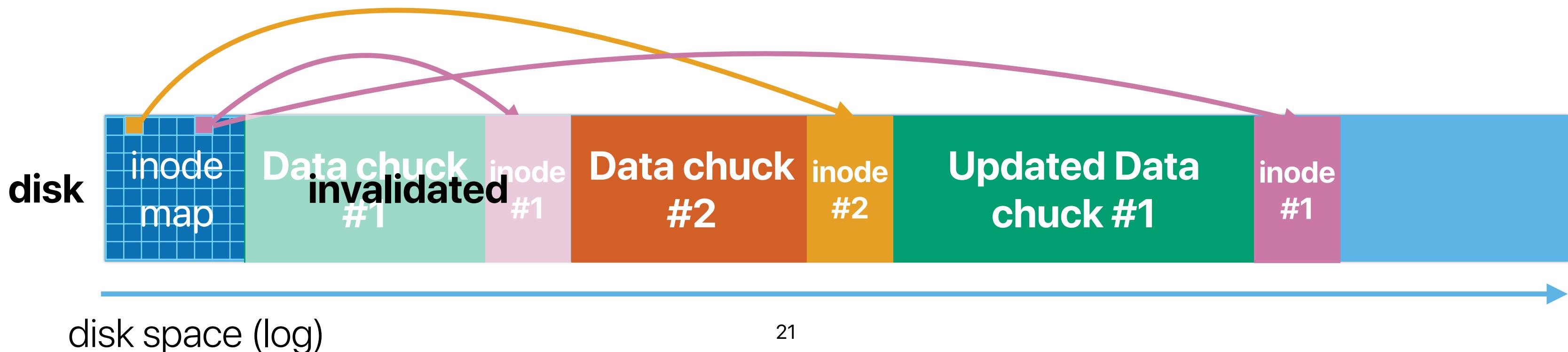


**What's the good things and bad
things about "Logs"**

Pros/Cons of Logs

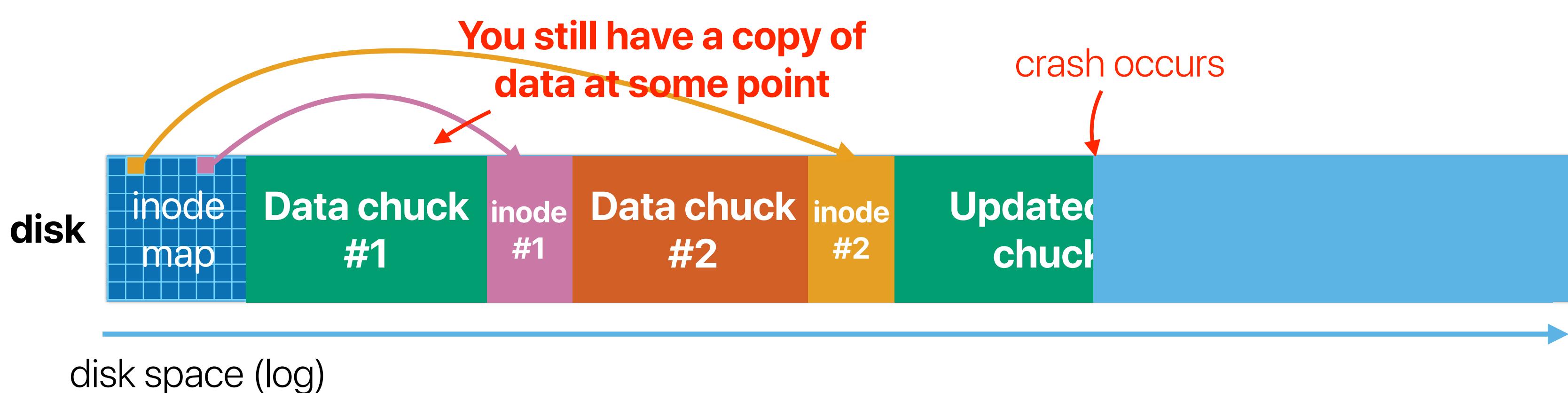
Log

- An append only data structure that records all changes
- Advantages of logs
 - Better performance — always sequential accesses
 - Faster writes — you just need to append without sanitize existing data first
 - Ease of recovery — you can find previous changes within the log
- Disadvantage of logs — you will need to explicit perform garbage collections to reclaim spaces

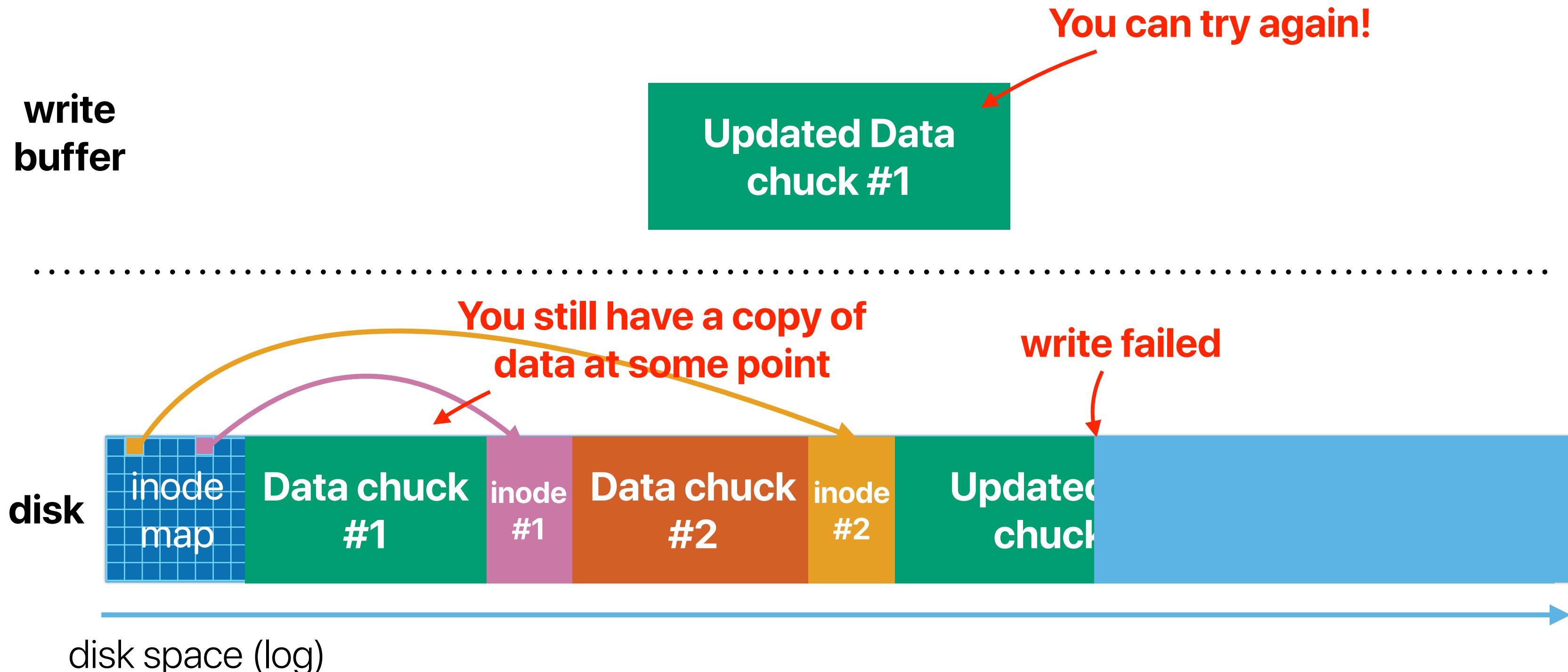


LFS v.s. crash

write
buffer

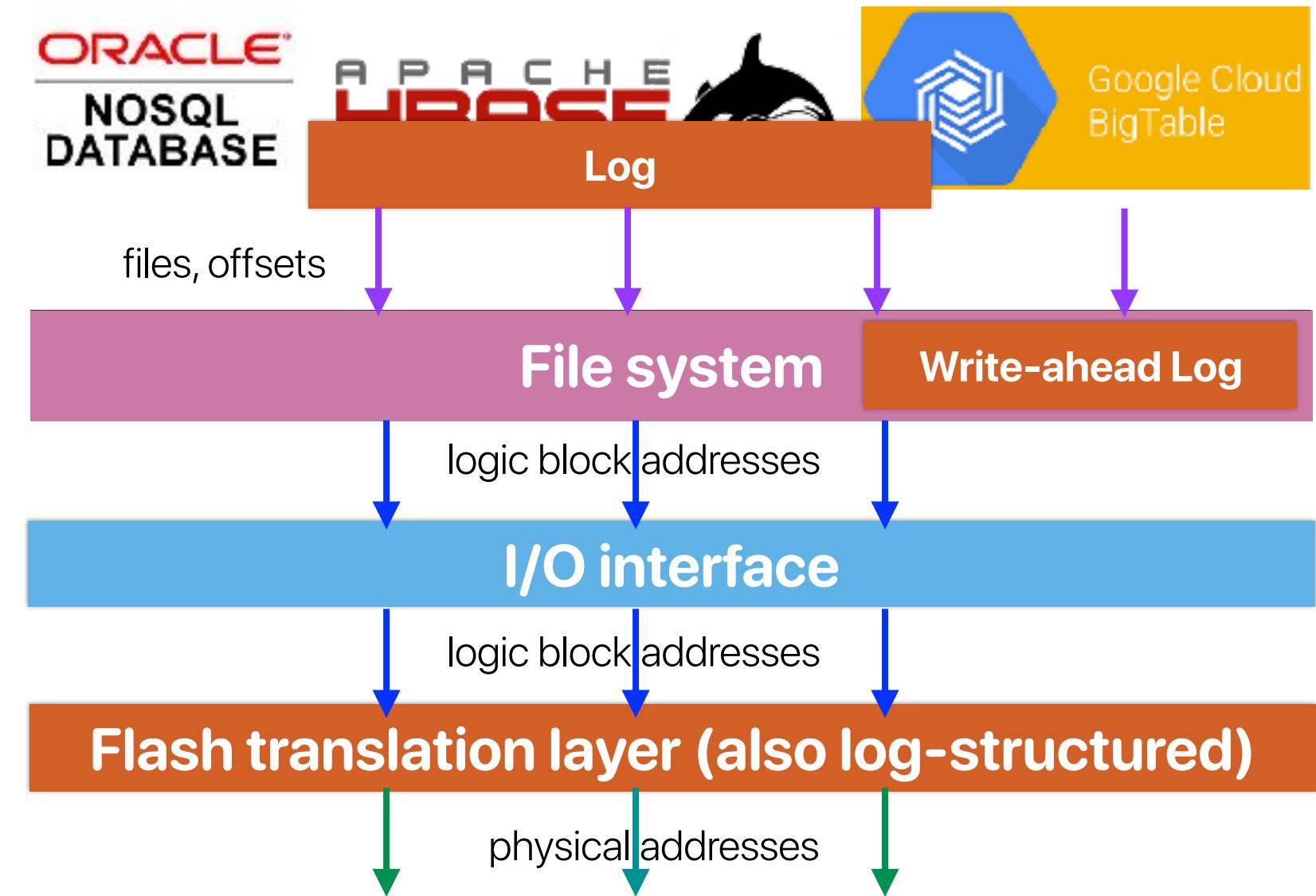


LFS v.s. write failed

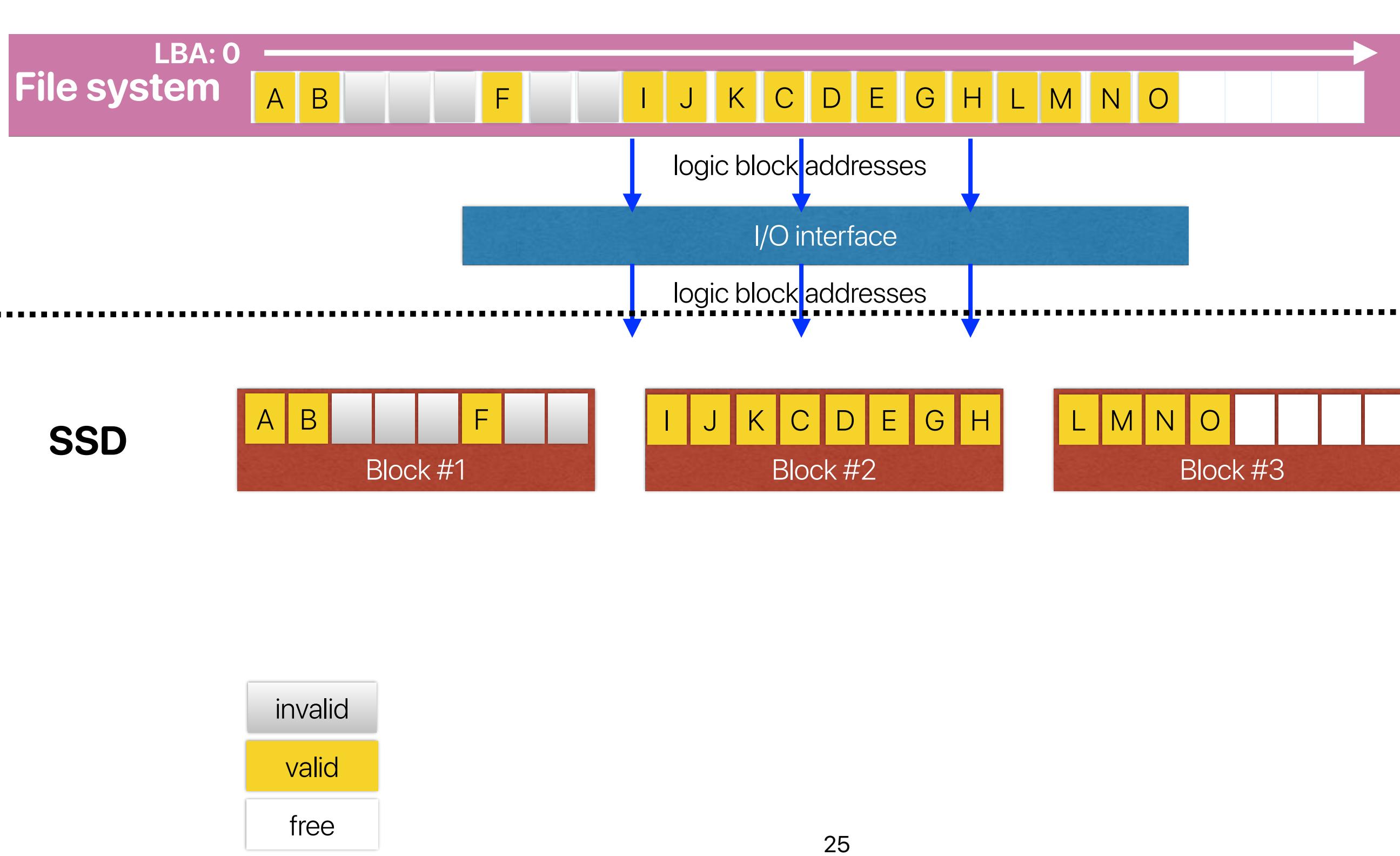


Why should we care about this paper?

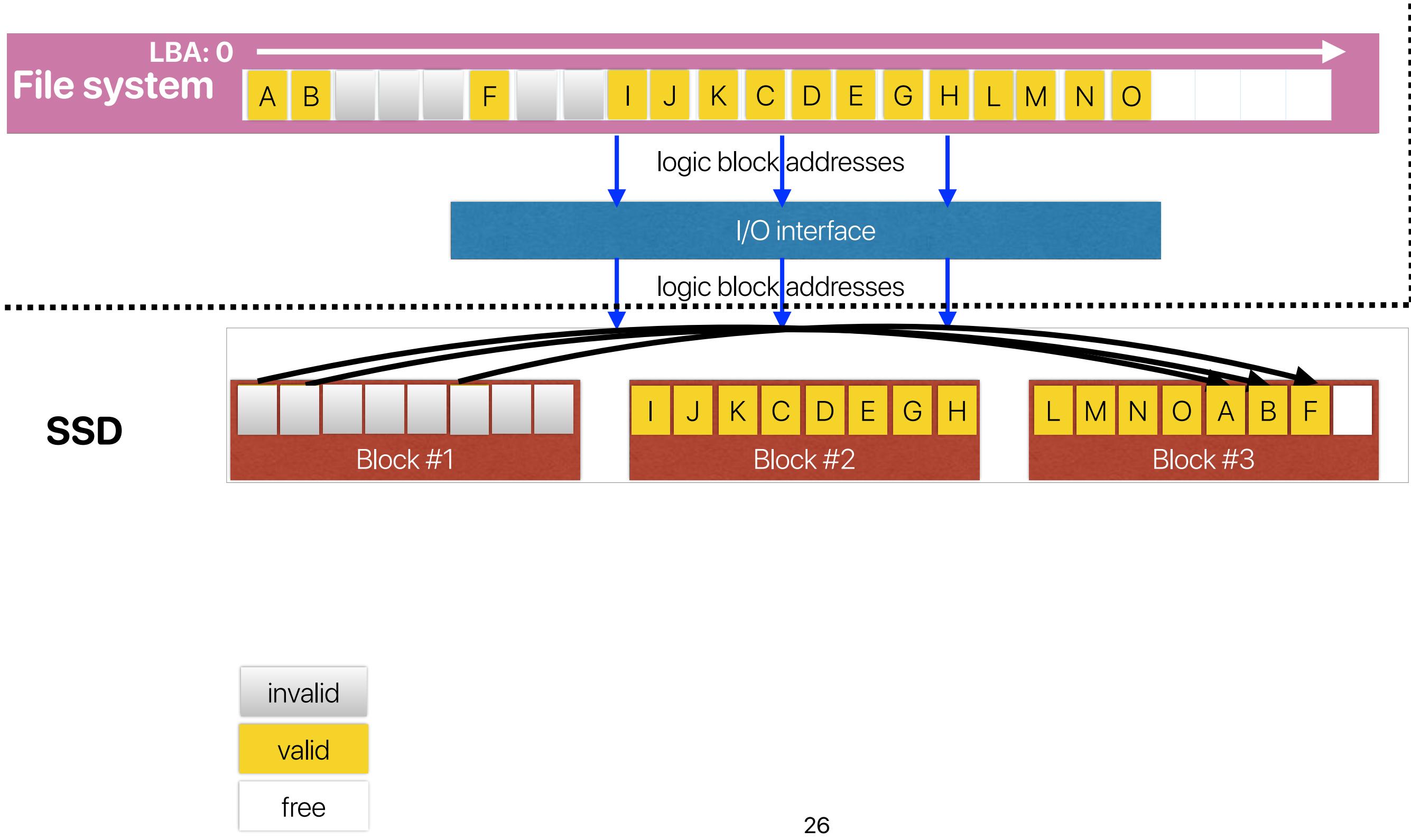
- Log is everywhere
 - Application: database
 - File system
 - Flash-based SSDs
- They can interfere with each other!
- An issue with software engineering nowadays



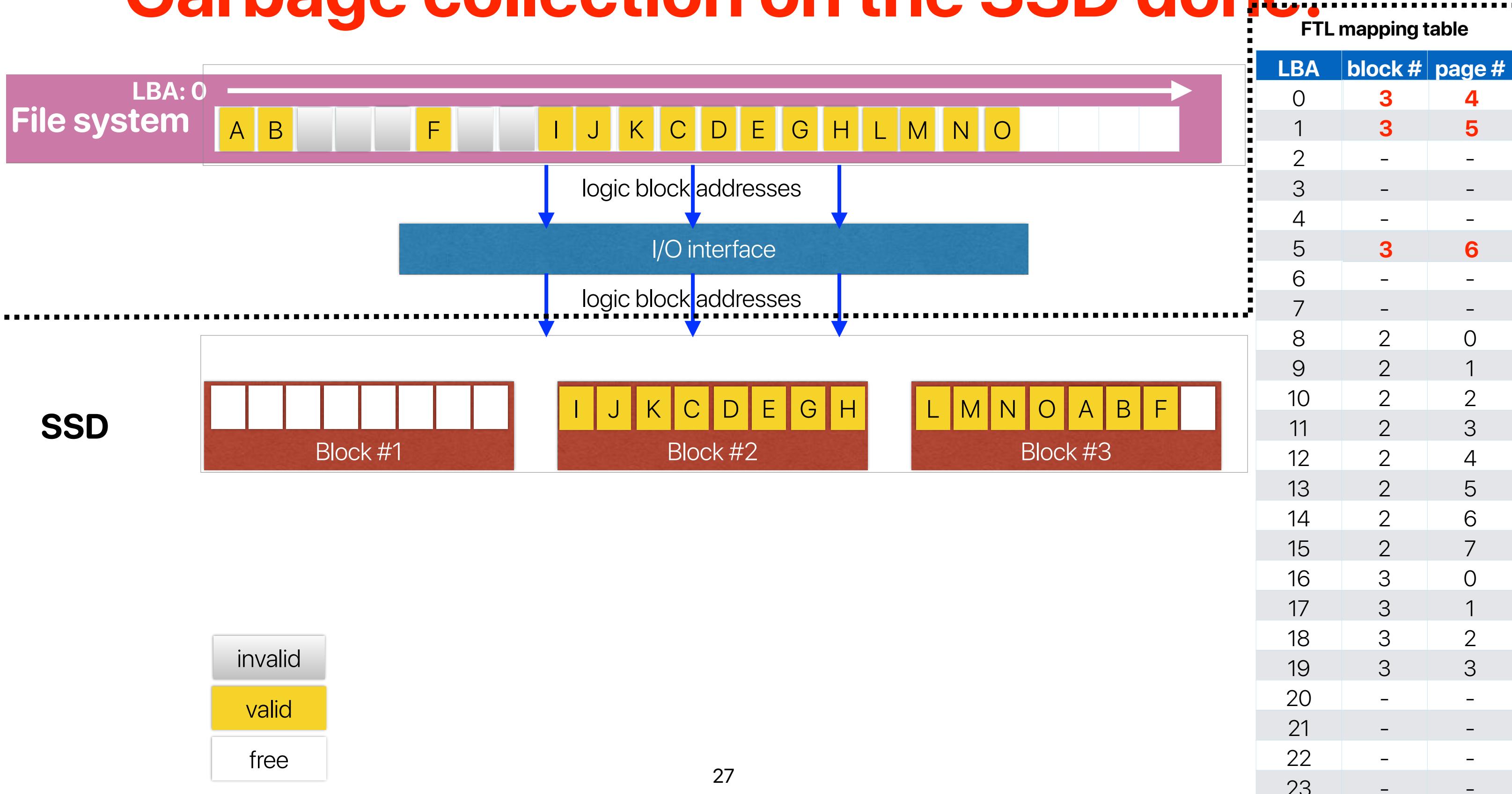
For example, garbage collection



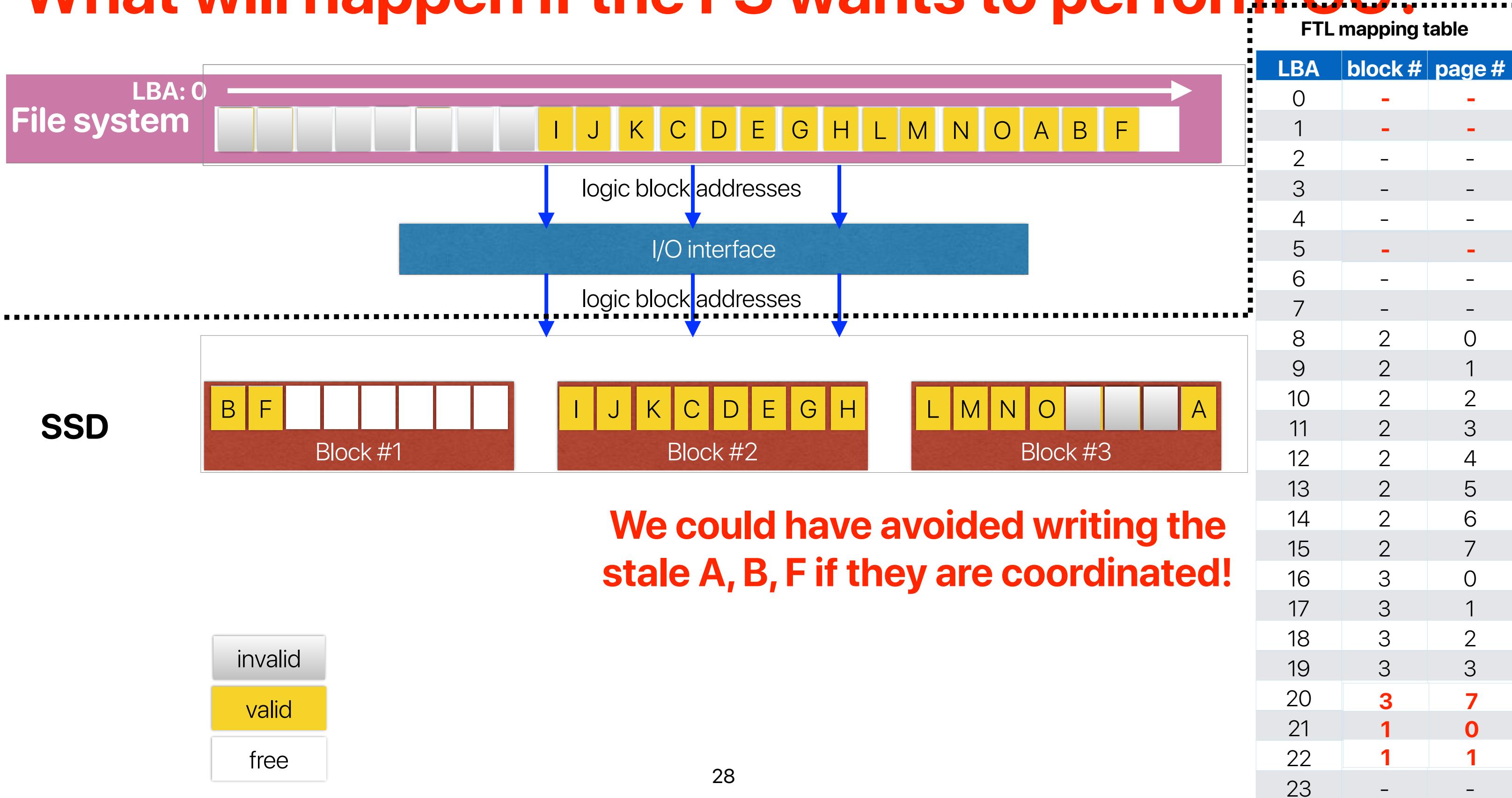
Now, SSD wants to reclaim a block



Garbage collection on the SSD done!



What will happen if the FS wants to perform GC?

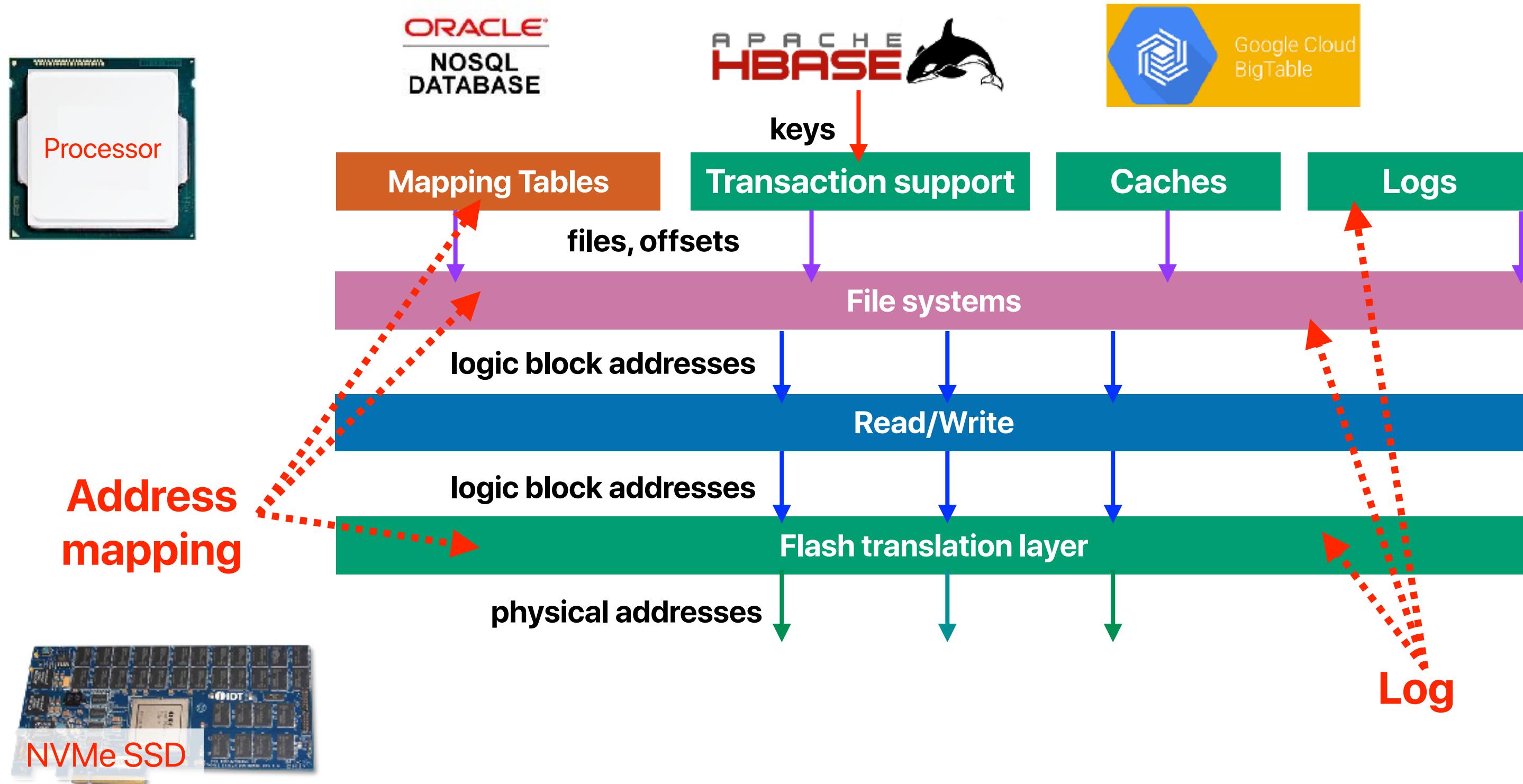


All problems in computer science can be solved by another level of
indirection

–David Wheeler

...except for the problem of too many layers of indirection.

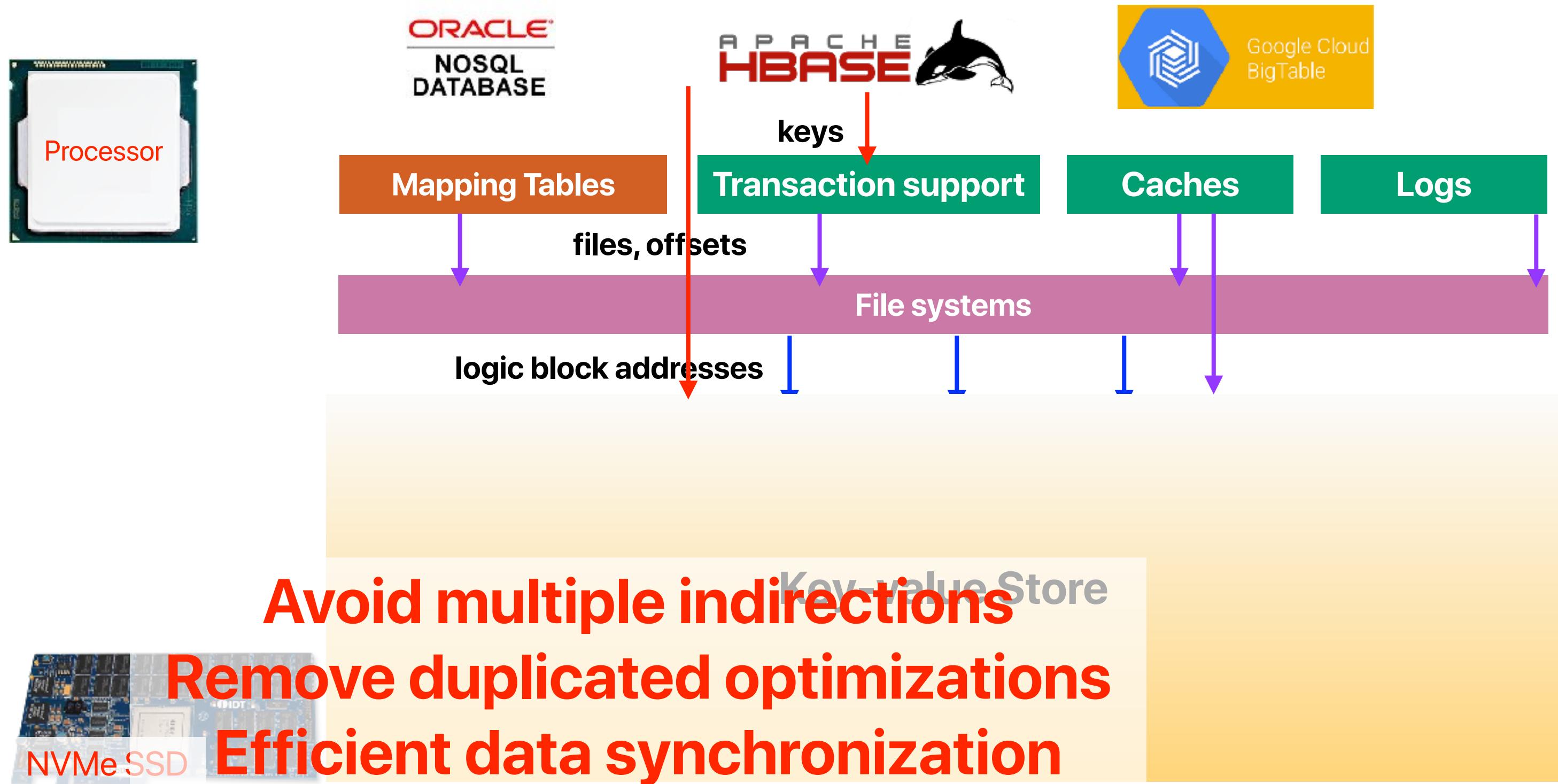
When the conventional interface meets modern applications



When the conventional interface meets modern applications

- Log on log problems
 - Database logs, FS and FTL are all log-structured
 - Redundant functions can cause suboptimal performance
- Data translations/lookup in each layer
 - Waste resources to maintain the mapping
 - Addition latency for each data access

When the conventional interface meets modern applications



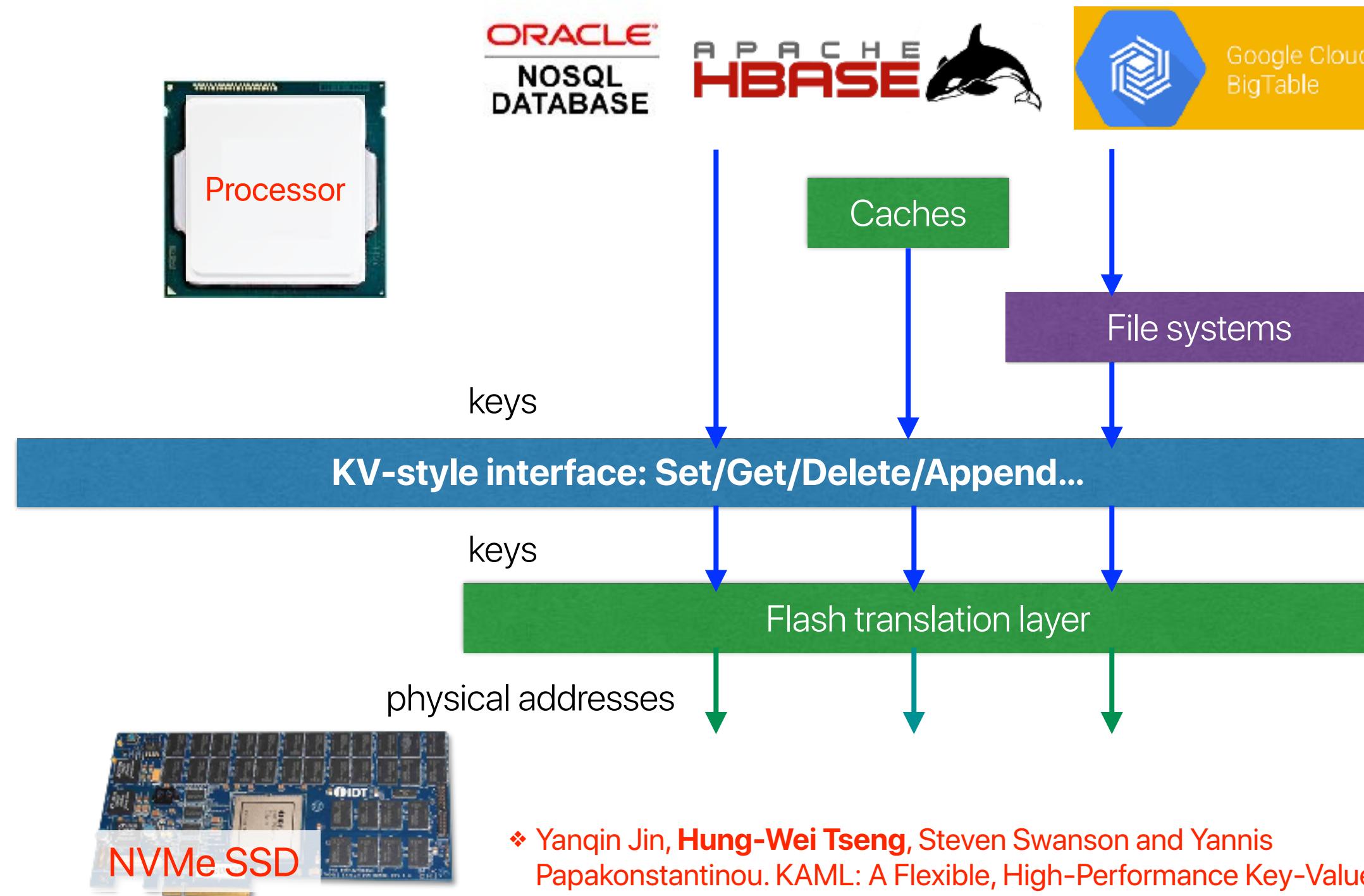
What if the SSD understands key-value?

- Remove the indirections
 - No need to maintain a mapping table
 - No need to translate between logical block addresses and physical addresses
- Avoid duplicated functions
 - Flash translation layer already does address remapping
 - Flash translation layer already maintains the physical memory as logs
- More efficient transaction support in hardware
 - No software synchronization overhead
 - The storage device itself exploits internal parallelism better

KAML

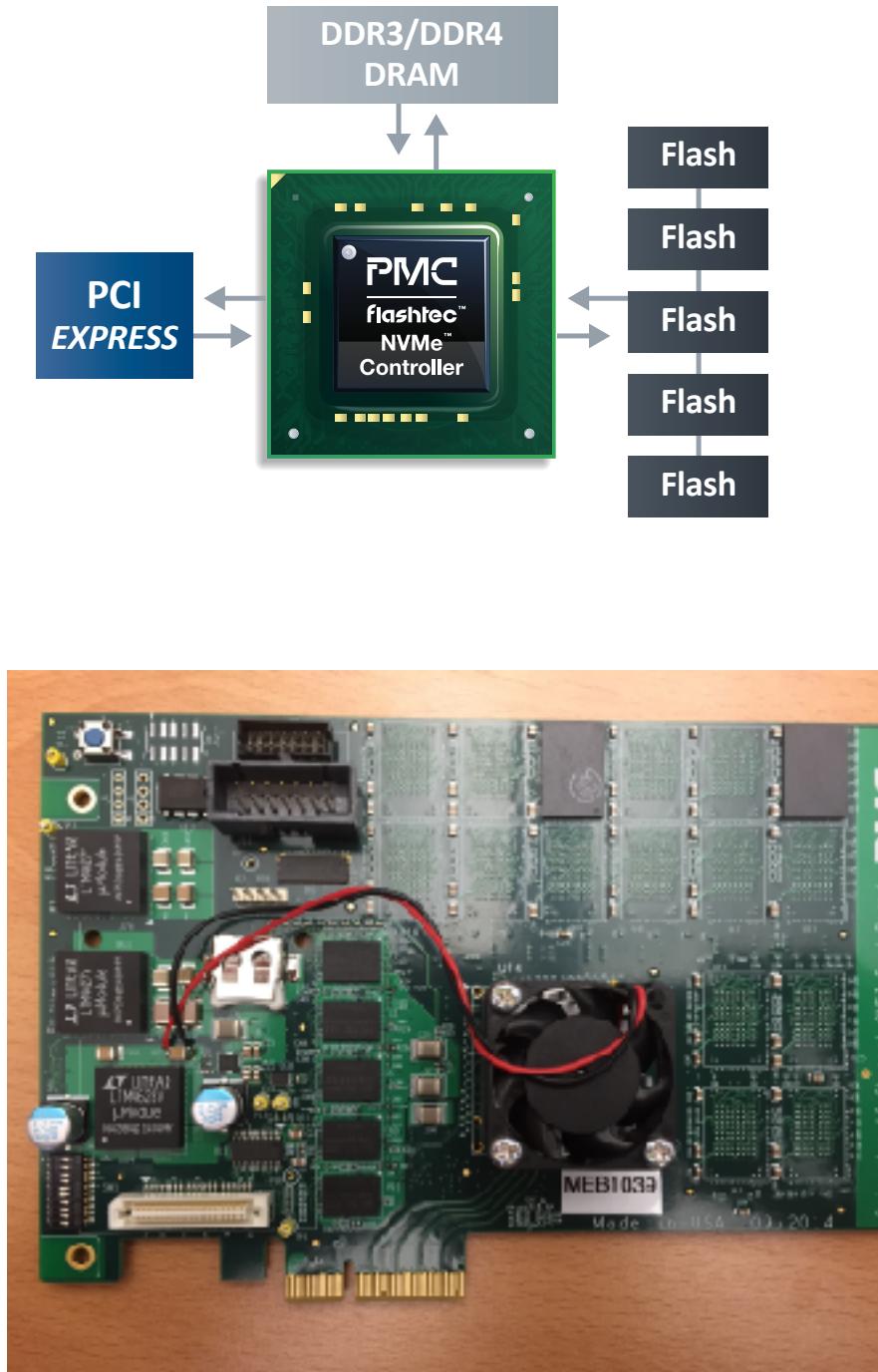
❖ Yanqin Jin, **Hung-Wei Tseng**, Steven Swanson and Yannis Papakonstantinou. KAML: A Flexible, High-Performance Key-Value SSD. In HPCA 2017.

KAML: A Flexible, High-Performance Key-Value SSD

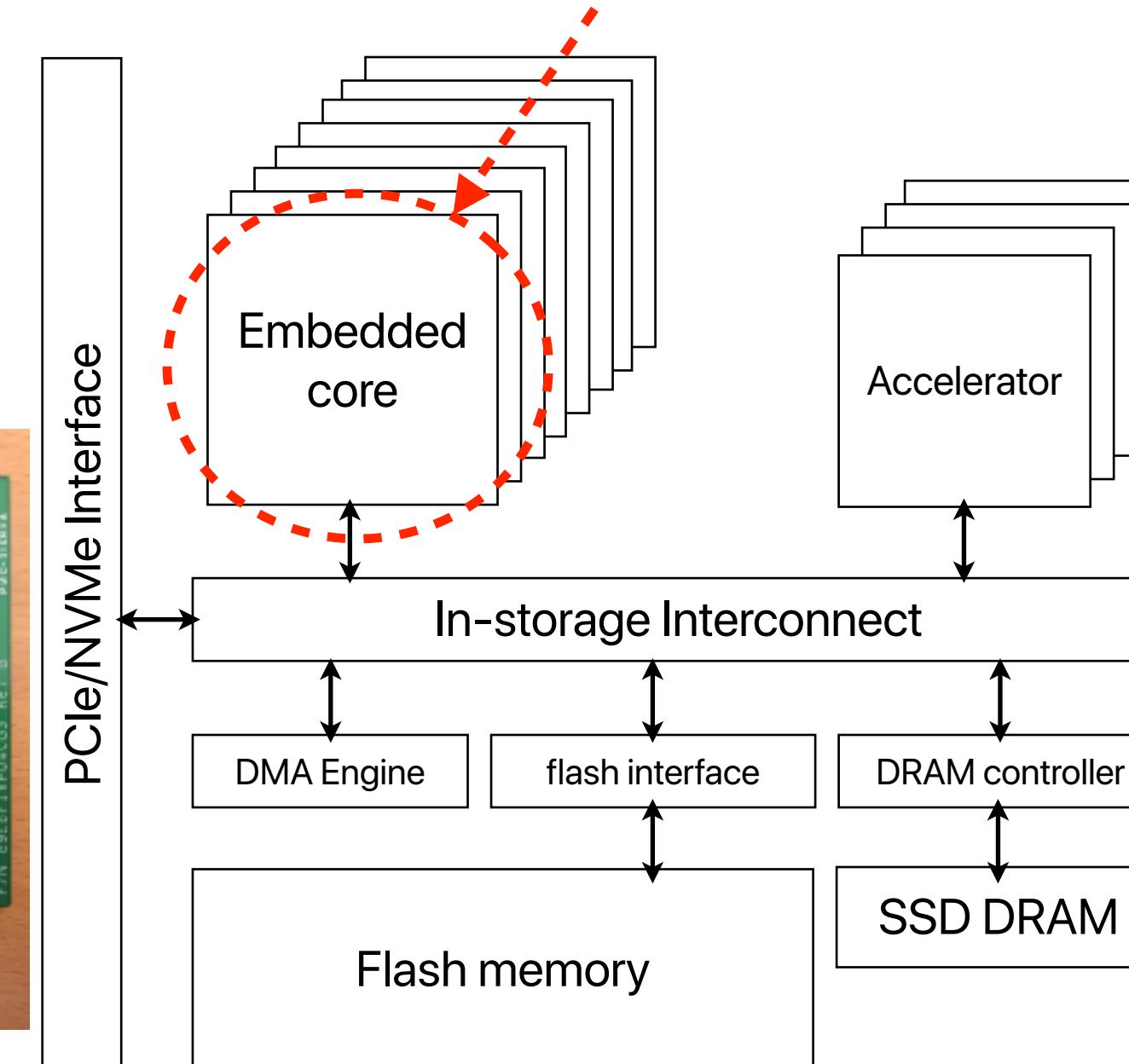


❖ Yanqin Jin, Hung-Wei Tseng, Steven Swanson and Yannis Papakonstantinou. KAML: A Flexible, High-Performance Key-Value SSD. In HPCA 2017.

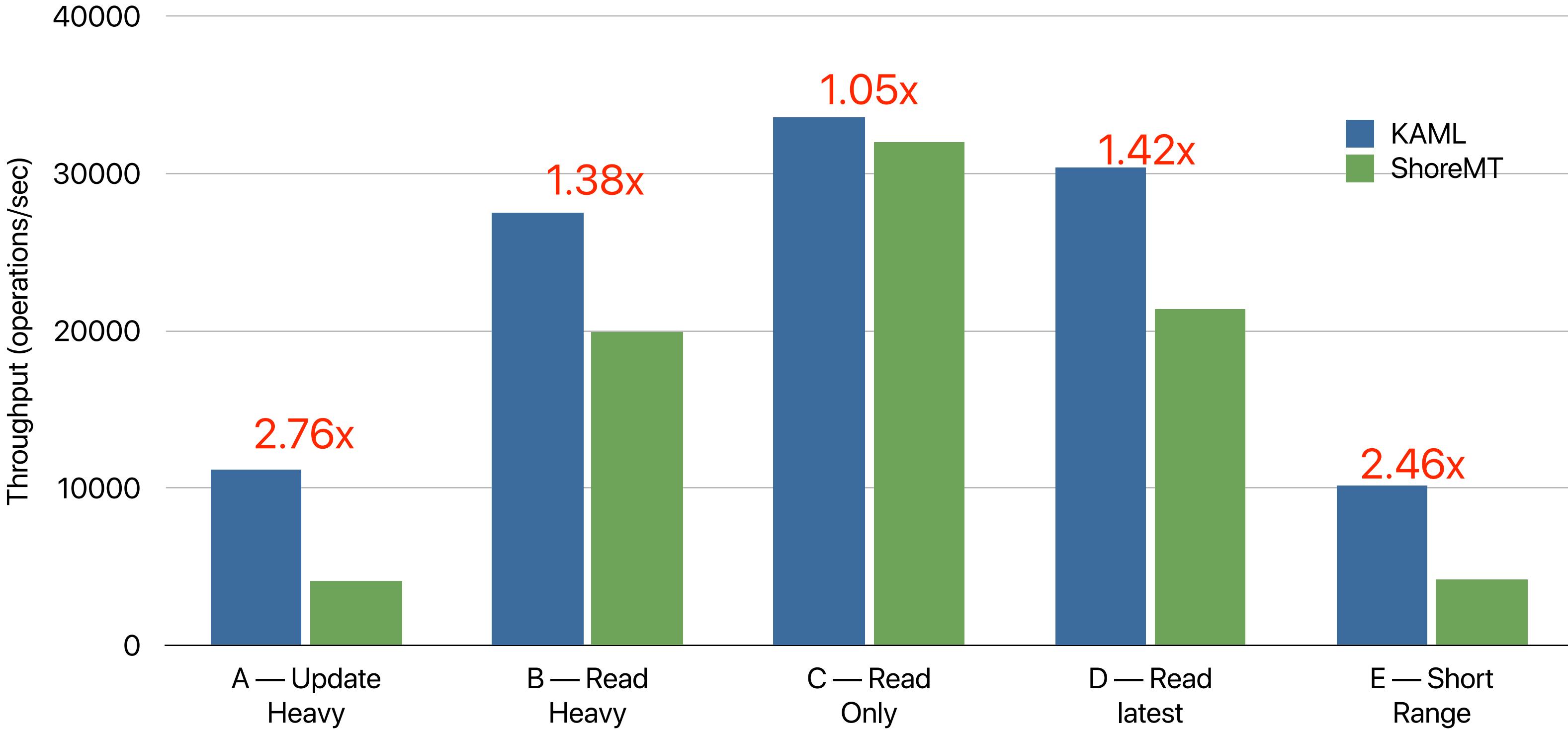
Implementing KAML



Managing key-value-like commands



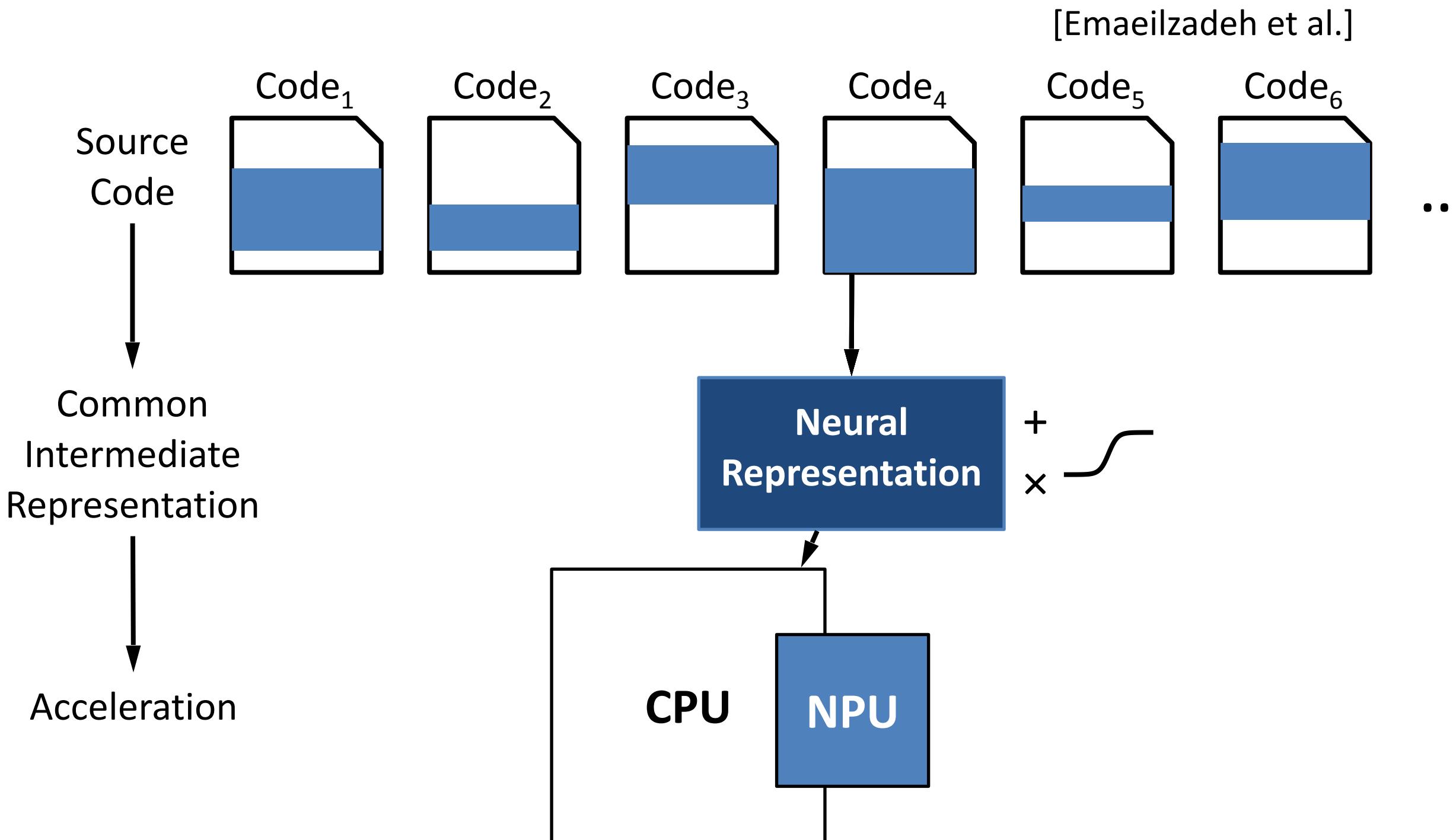
YCSB



Approximate Computing



Neural Algorithmic Transformation

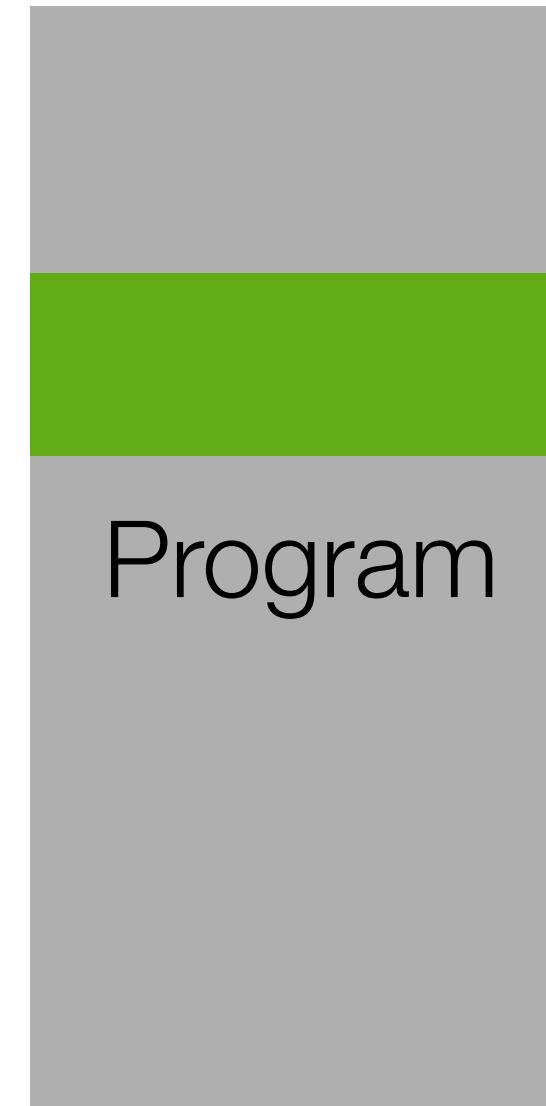


Program execution as a learning problem [Esmaeilzadeh NPU]



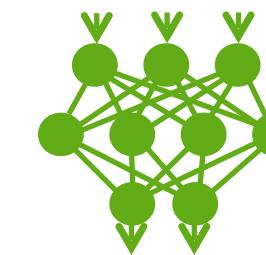
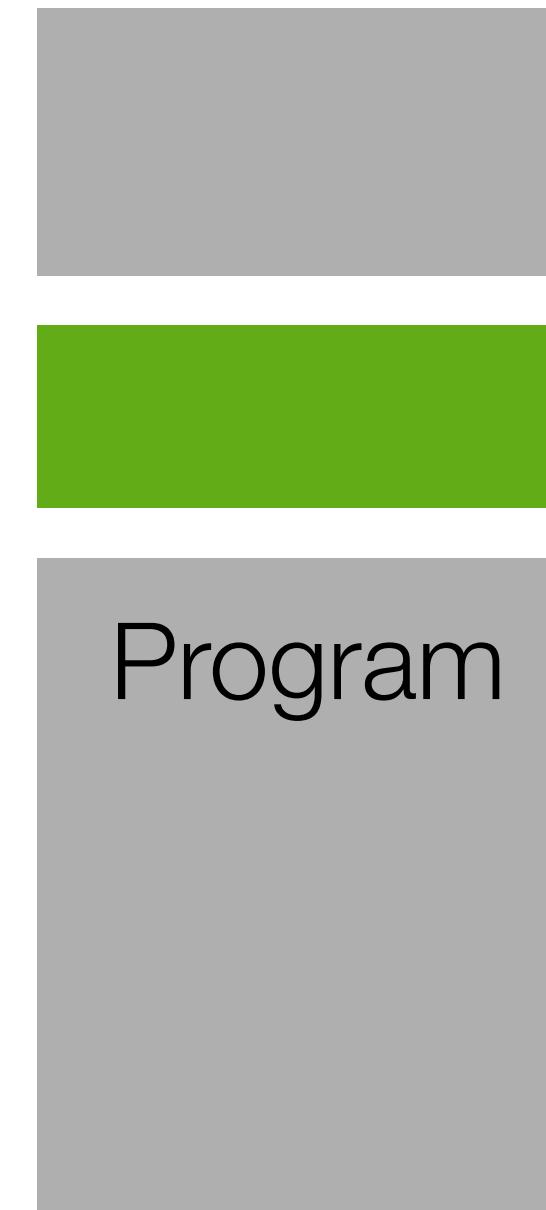
Program

Program execution as a learning problem [Esmaeilzadeh NPU]



Find an approximate
program component

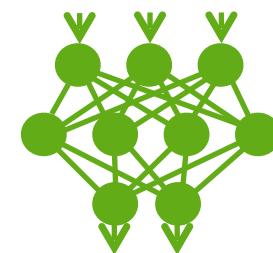
Program execution as a learning problem [Esmaeilzadeh NPU]



Find an approximate
program component

Compile the program
and train a neural network

Program execution as a learning problem [Esmaeilzadeh NPU]



Program

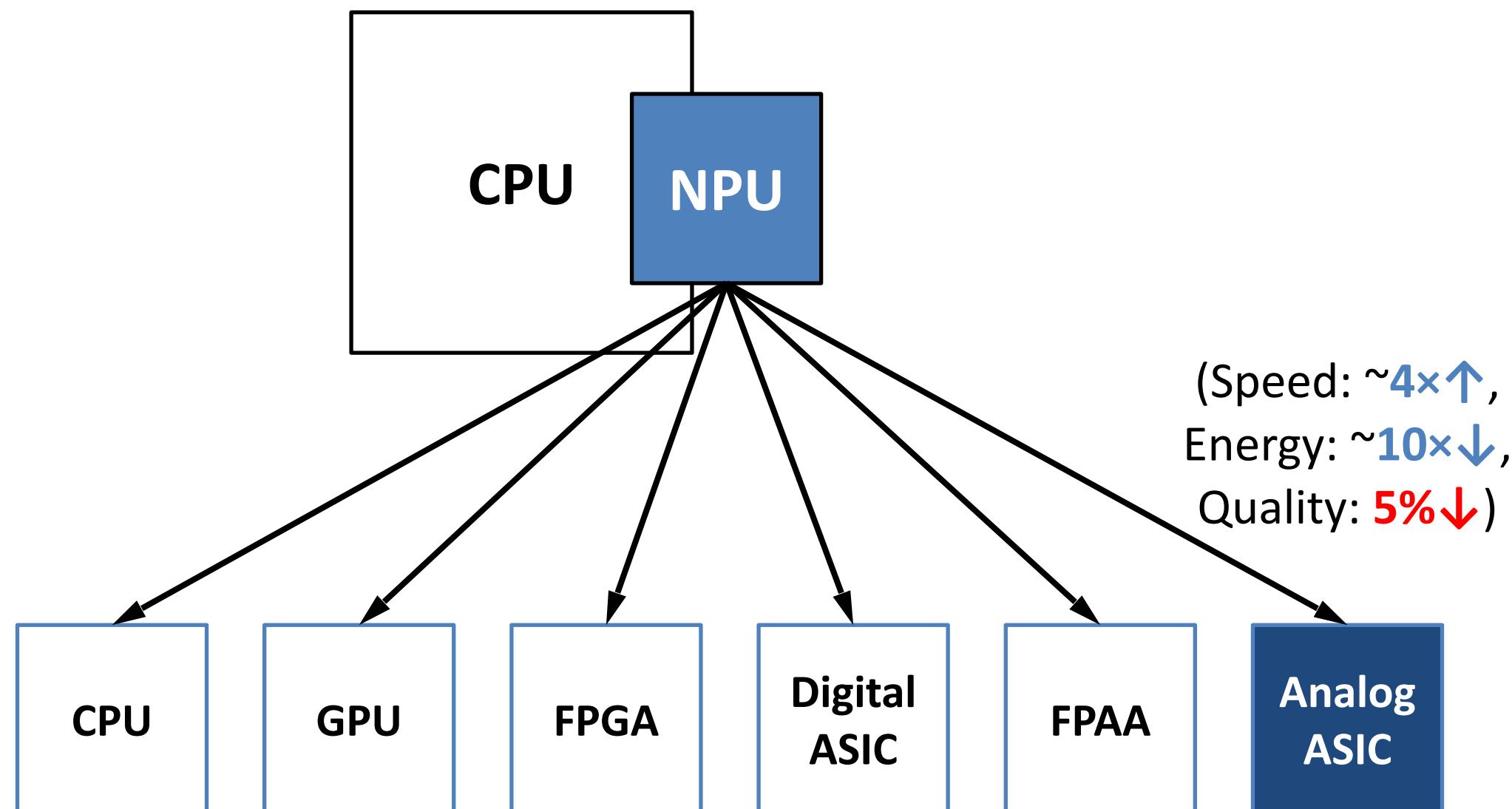
Find an approximate
program component

Compile the program
and train a neural network

Execute on a fast Neural
Processing Unit (NPU)

Neural Acceleration

[Emaeilzadeh et al.]



Remind yourselves Approximate
Computing — why/when is it a good
idea? why/when is it a bad idea?

Good things about Approximate Computing?

Good things about Approximate Computing

- Reduce circuit size
- Avoid complex computation
- Smaller data
- Can use “unreliable” hardware

Bad things about Approximate Computing?

Bad things about Approximate Computing

- Result
- Re-computation due to low quality
- Limited applications

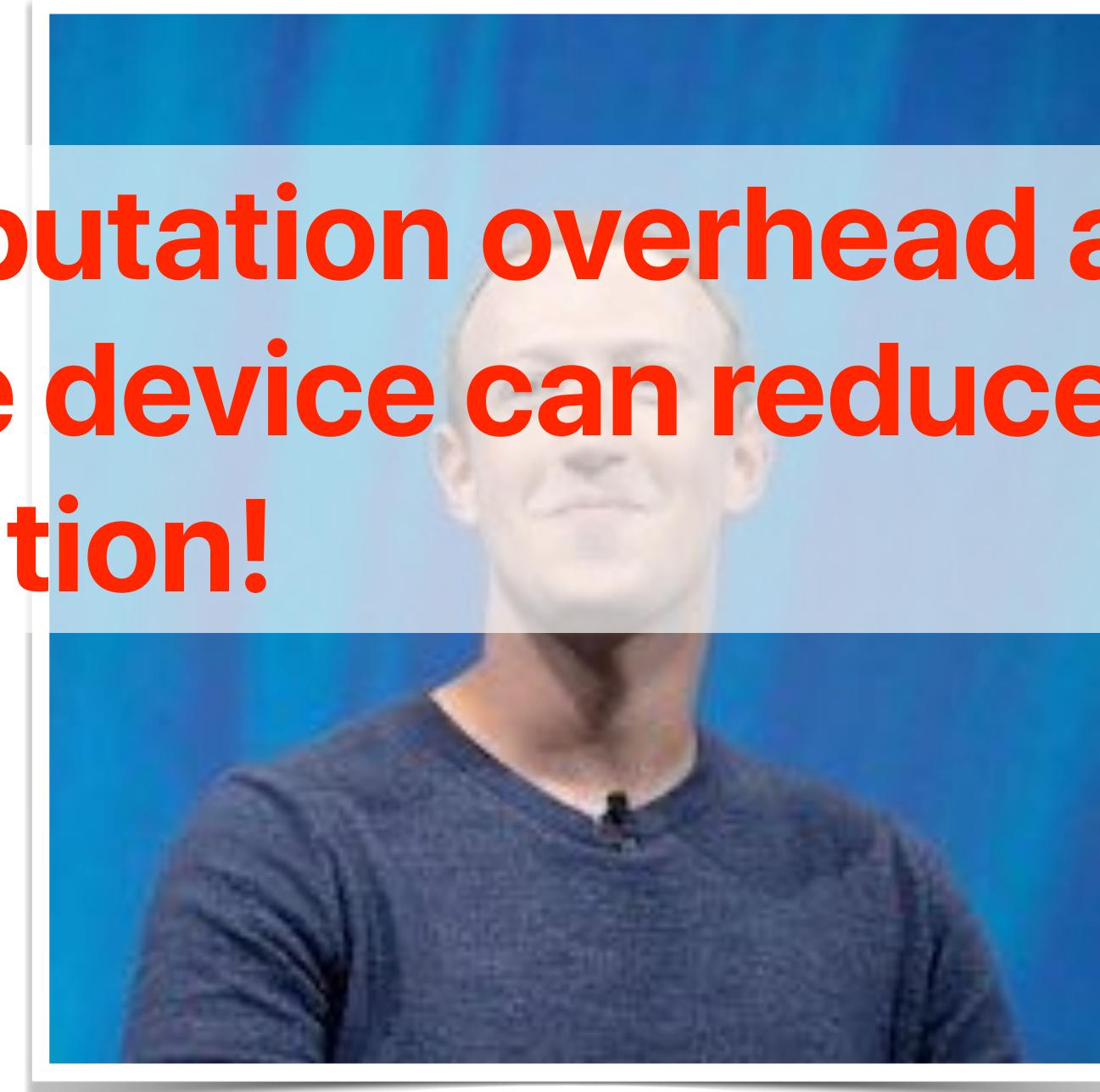
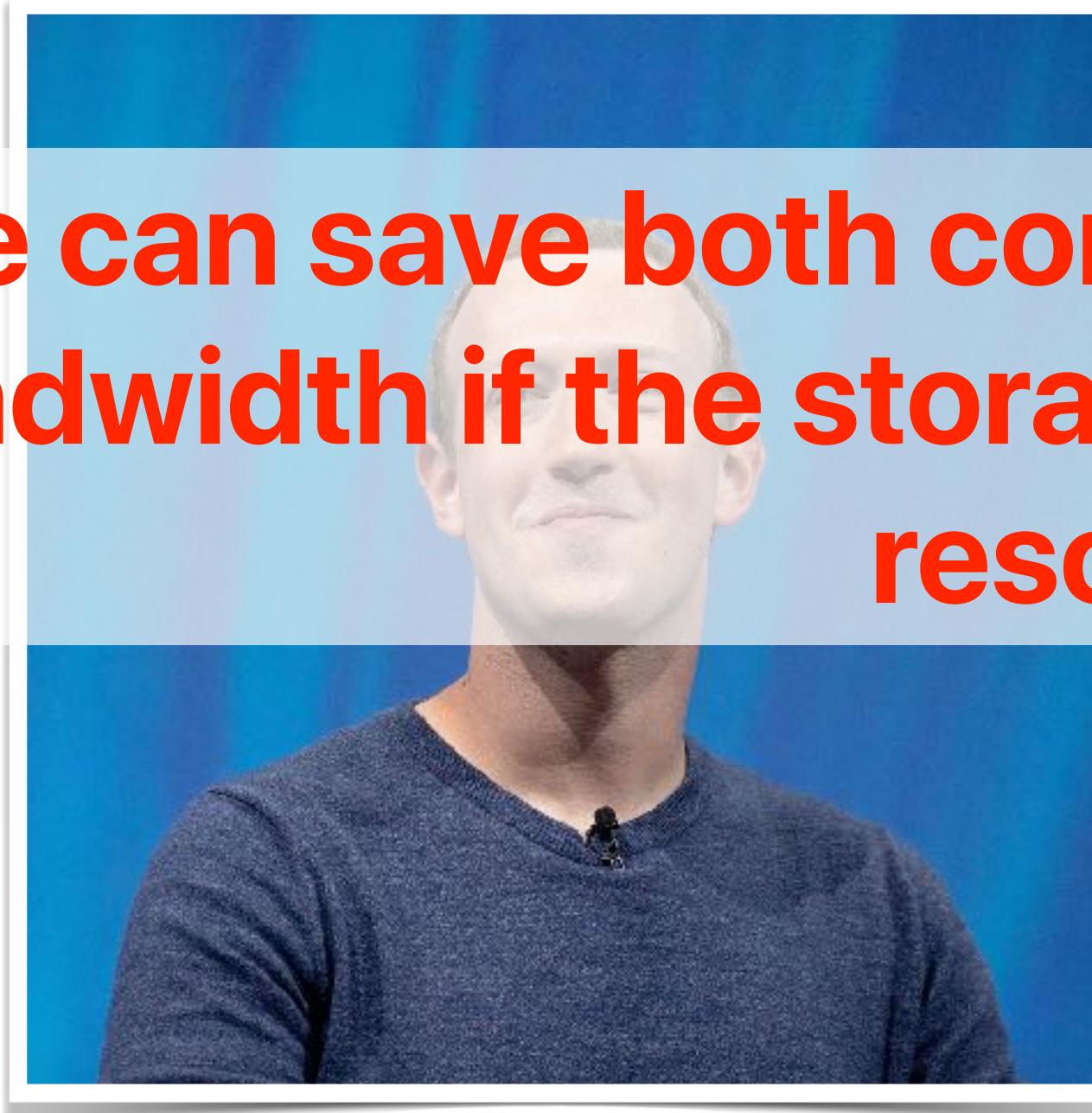
What is “approximate computing”?

- Building *acceptable* systems out of unreliable/inaccurate hardware and software components
 - Compute, storage and communication
 - Inaccuracy can be deterministic or non-deterministic
- Trade fidelity of results to achieve better efficiency and performance
 - Enable new applications or meet constraints (cost, power, etc.)
 - In summary, embrace widespread output quality trade-off and “closer to physics”, non-deterministic hardware

We don't need full resolution in many cases

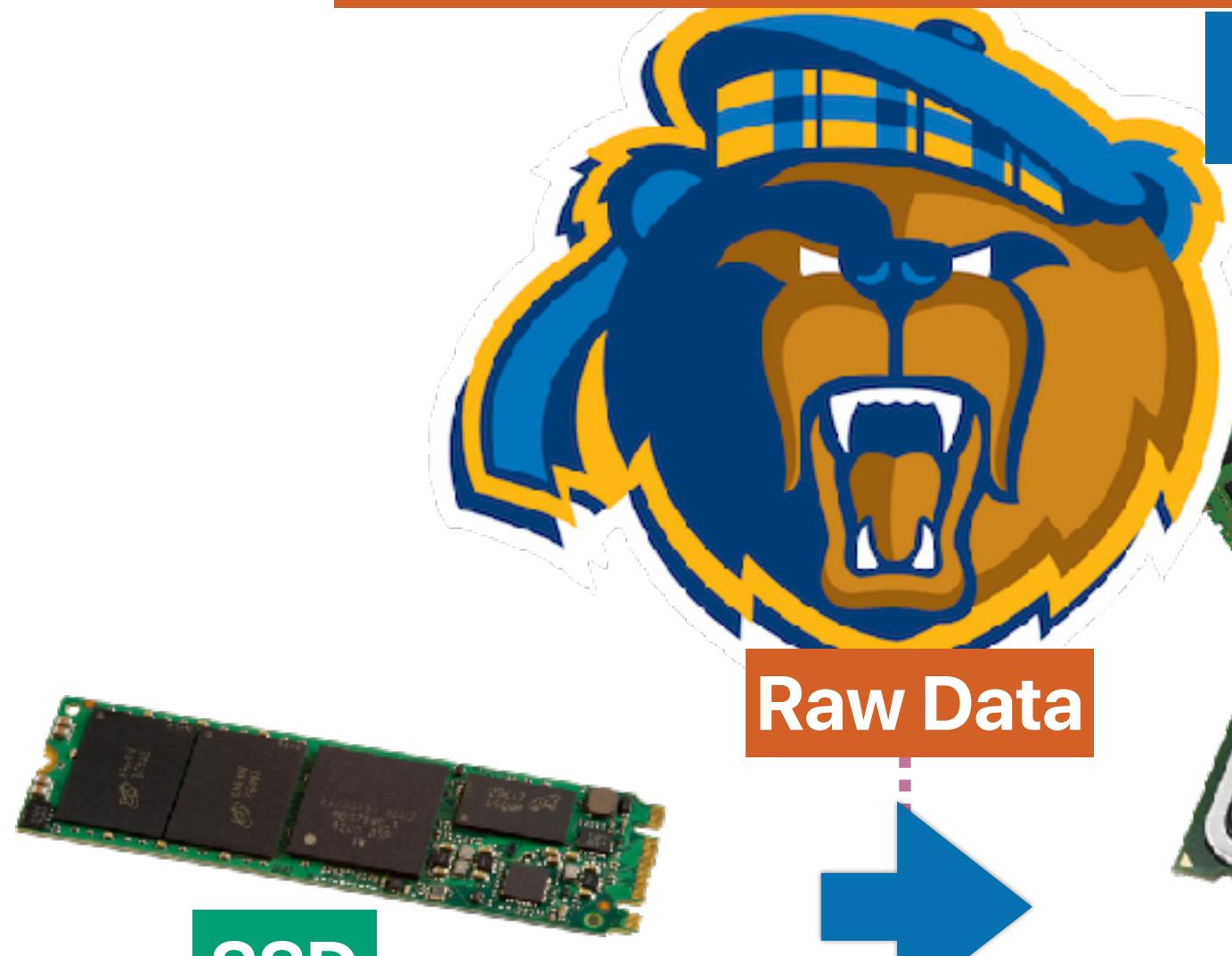
Only need 25% bytes to transfer

We can save both computation overhead and bandwidth if the storage device can reduce the resolution!

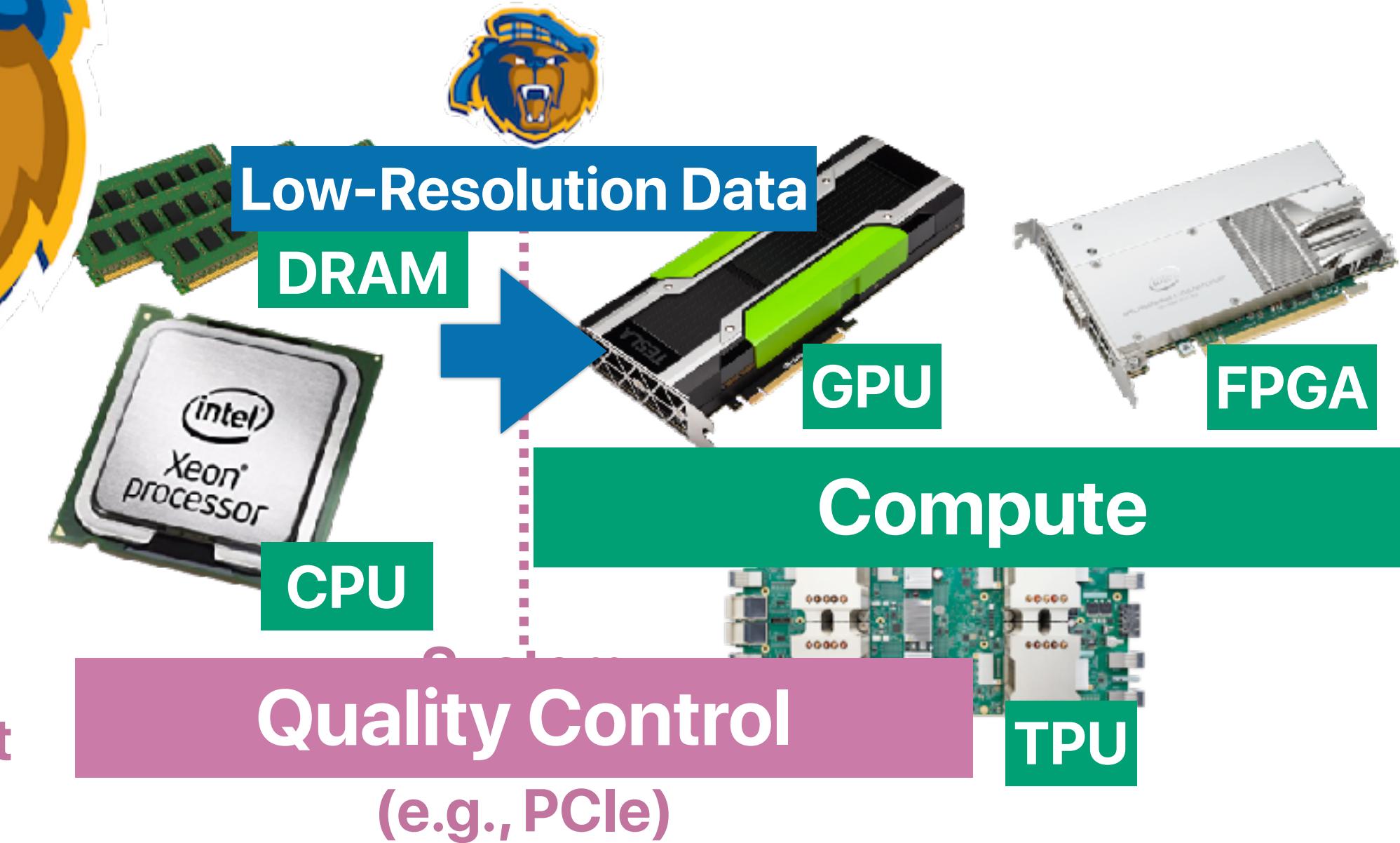


Accelerators in General-Purpose Computers

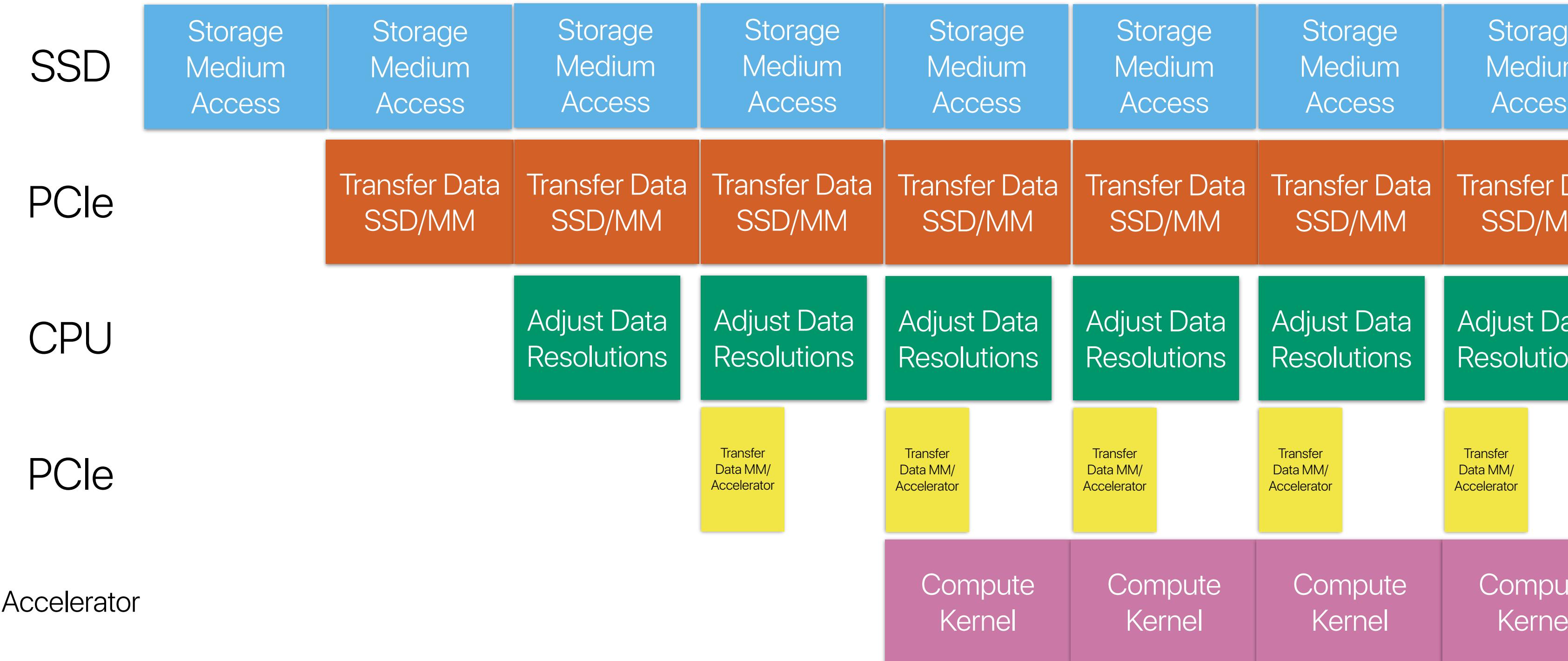
Retrieve Raw Data



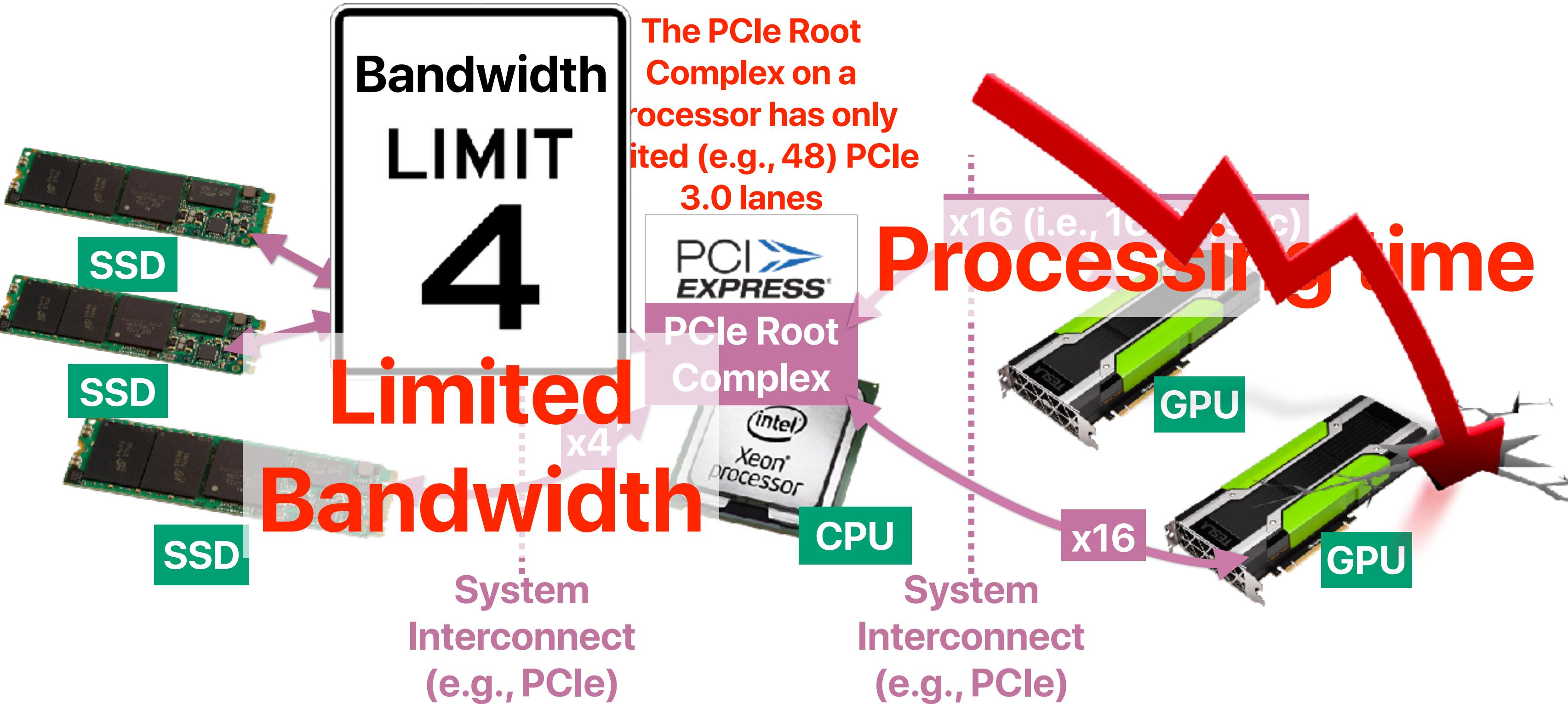
Adjust Data Resolutions



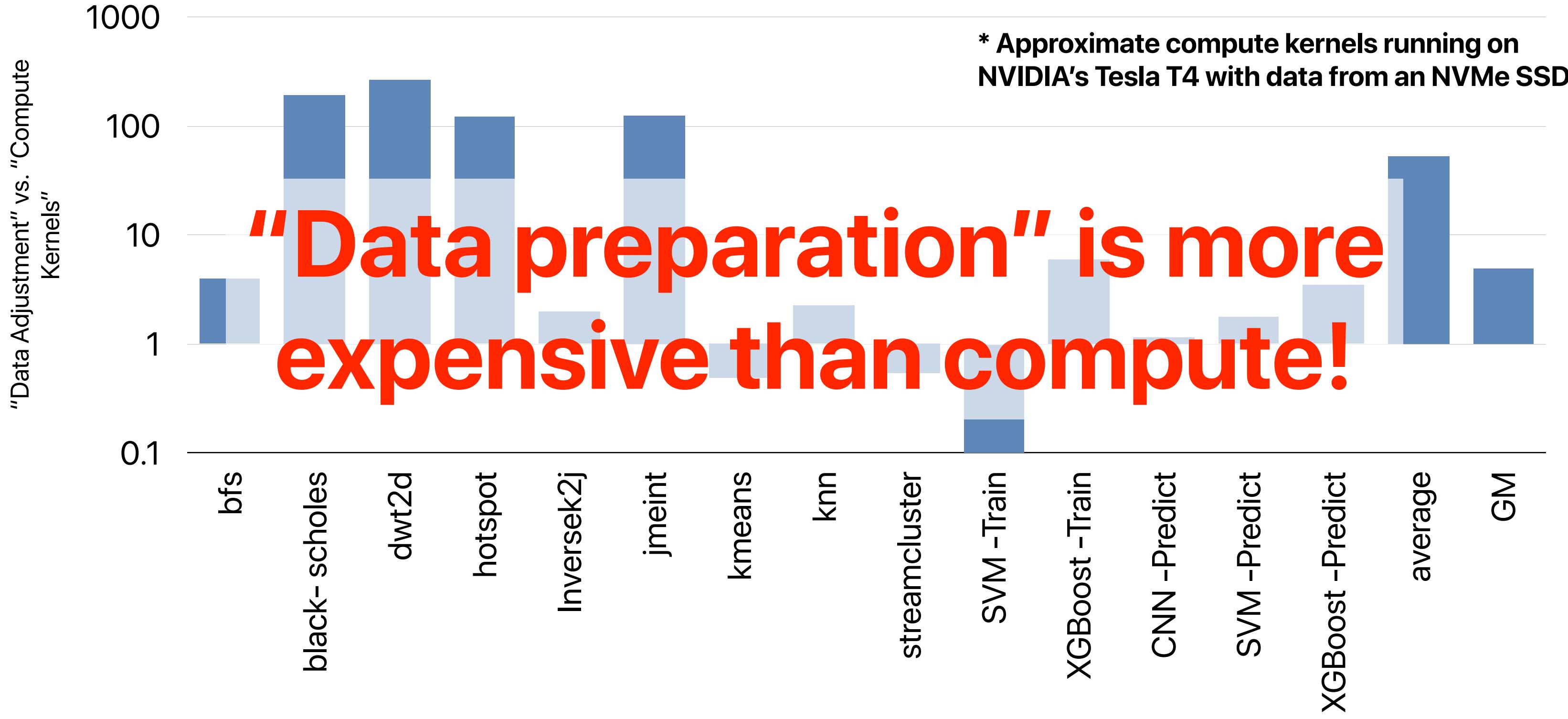
Data processing pipeline in modern computers



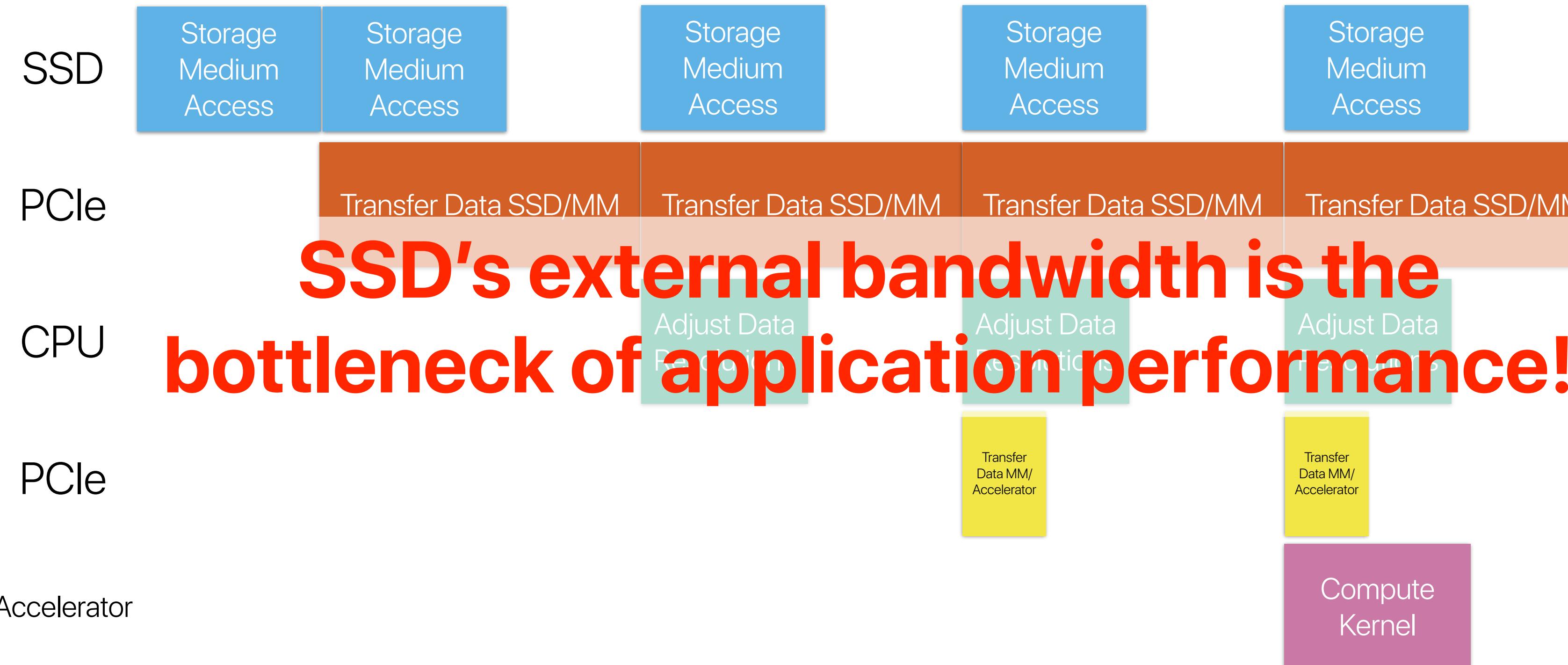
Limited Interconnect Bandwidth



Data Preparation Can be 100x More Than Compute Kernels



The “actual” pipeline ...



Example — Varifocal Storage*

- * Yu-Ching Hu, Murtuza Lokhandwala, Te I and Hung-Wei Tseng. Dynamic Multi-Resolution Data Storage. In the 52nd Annual IEEE/ACM International Symposium on Microarchitecture, MICRO 2019 (Best paper nominee)

Varifocal Storage: A Holistic System Architecture

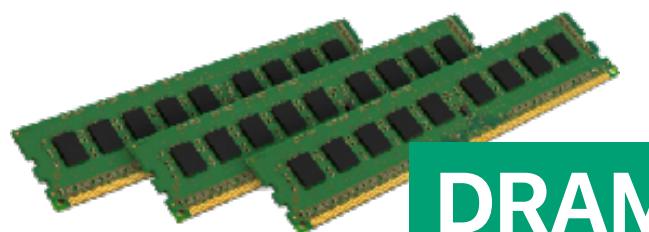
Retrieve Raw Data

Adjust Data Resolutions



Varifocal Storage (VS)

System
Interconnect
(e.g., PCIe)



DRAM



CPU

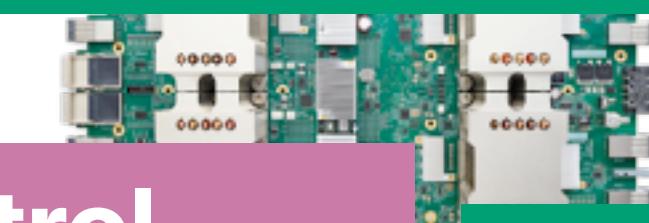


GPU



FPGA

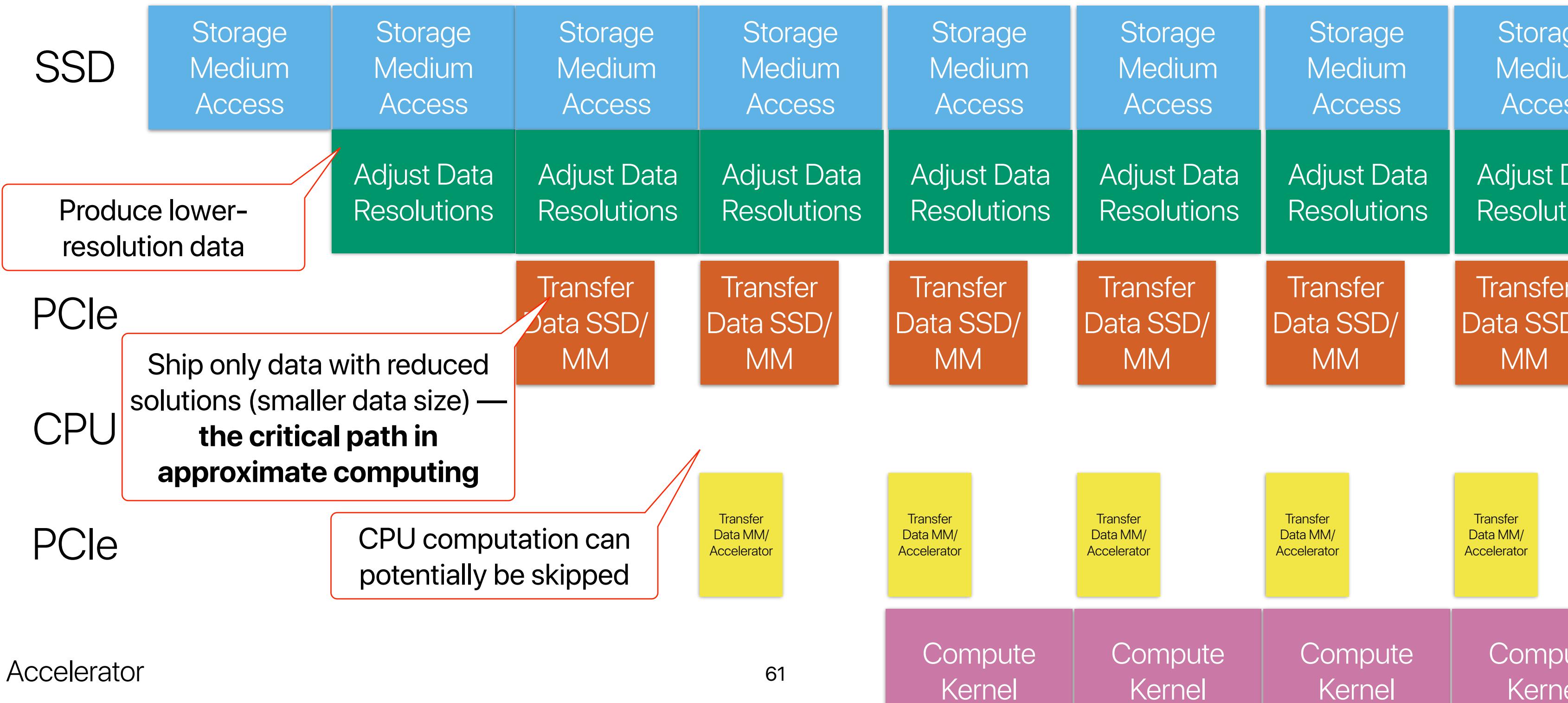
Compute



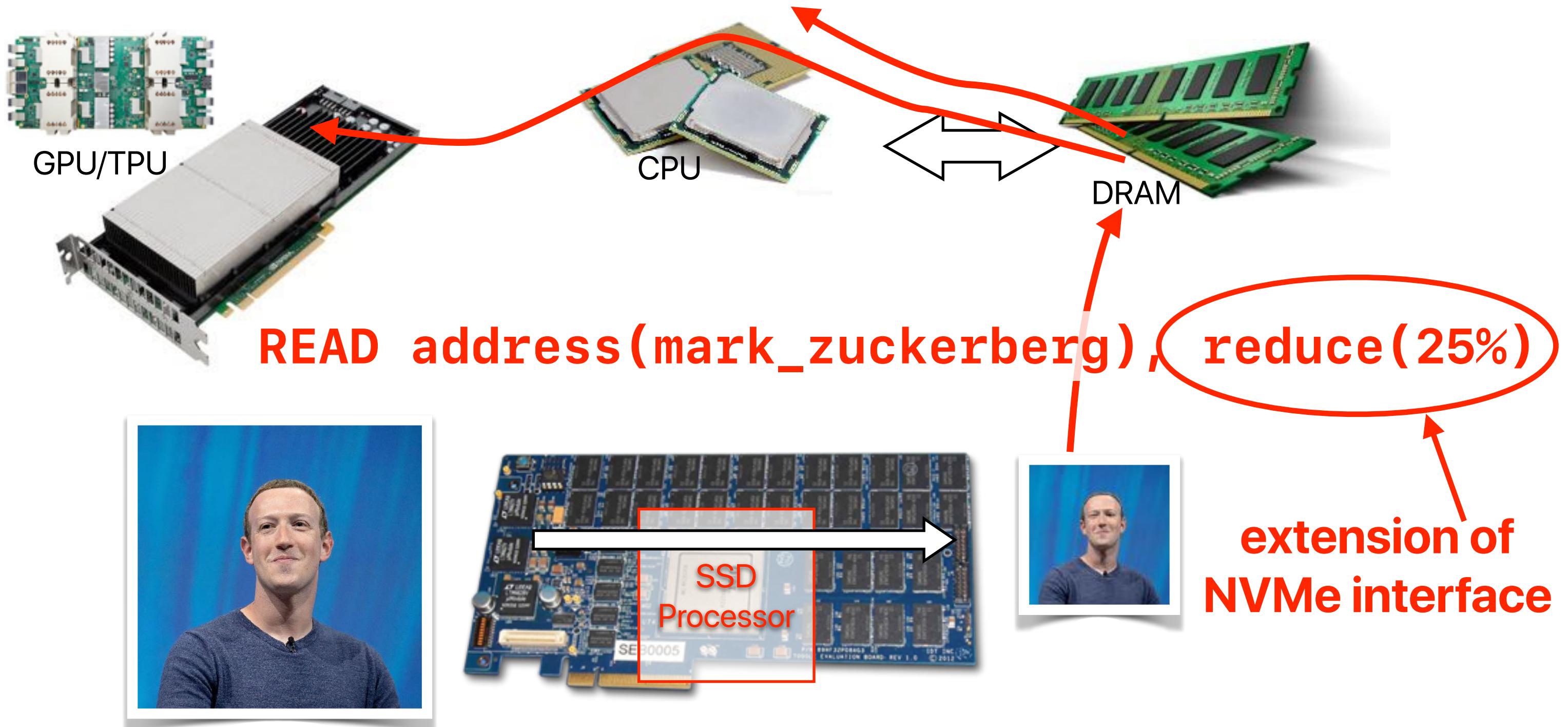
TPU

Quality Control
(e.g., PCIe)

Data processing with ISP



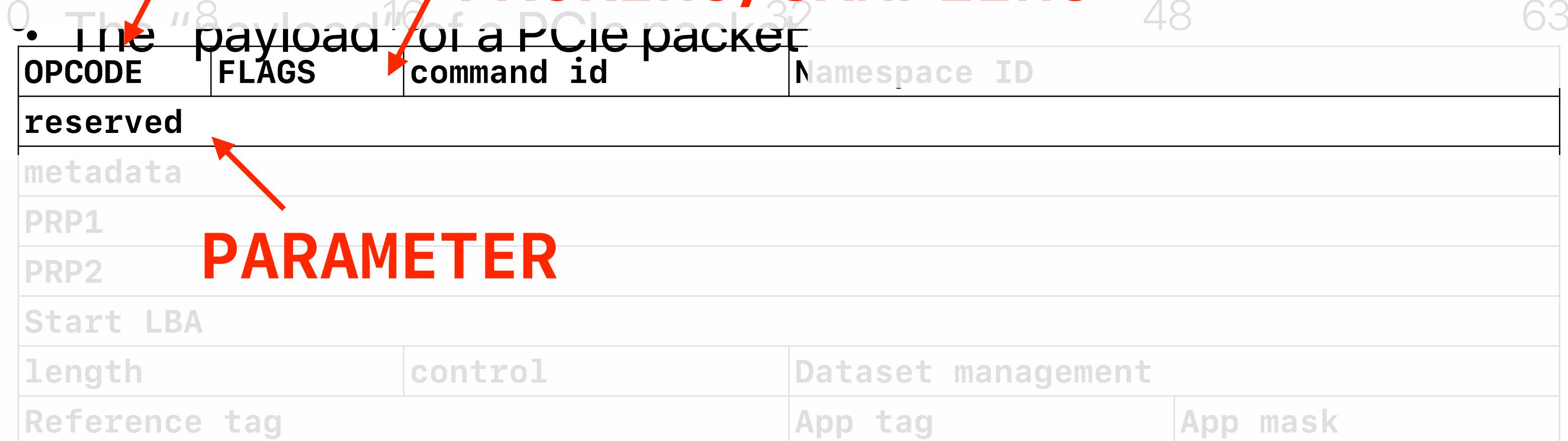
The “concept” of Varifocal Storage



NVMe

- The standard of PCIe SSD devices now
 - Provides multiple command queues to better support multithreading

VS _READ **REDUCE/QUTIZATION/**
PACKING/SAMPLING



Programming model of Varifocal Storage

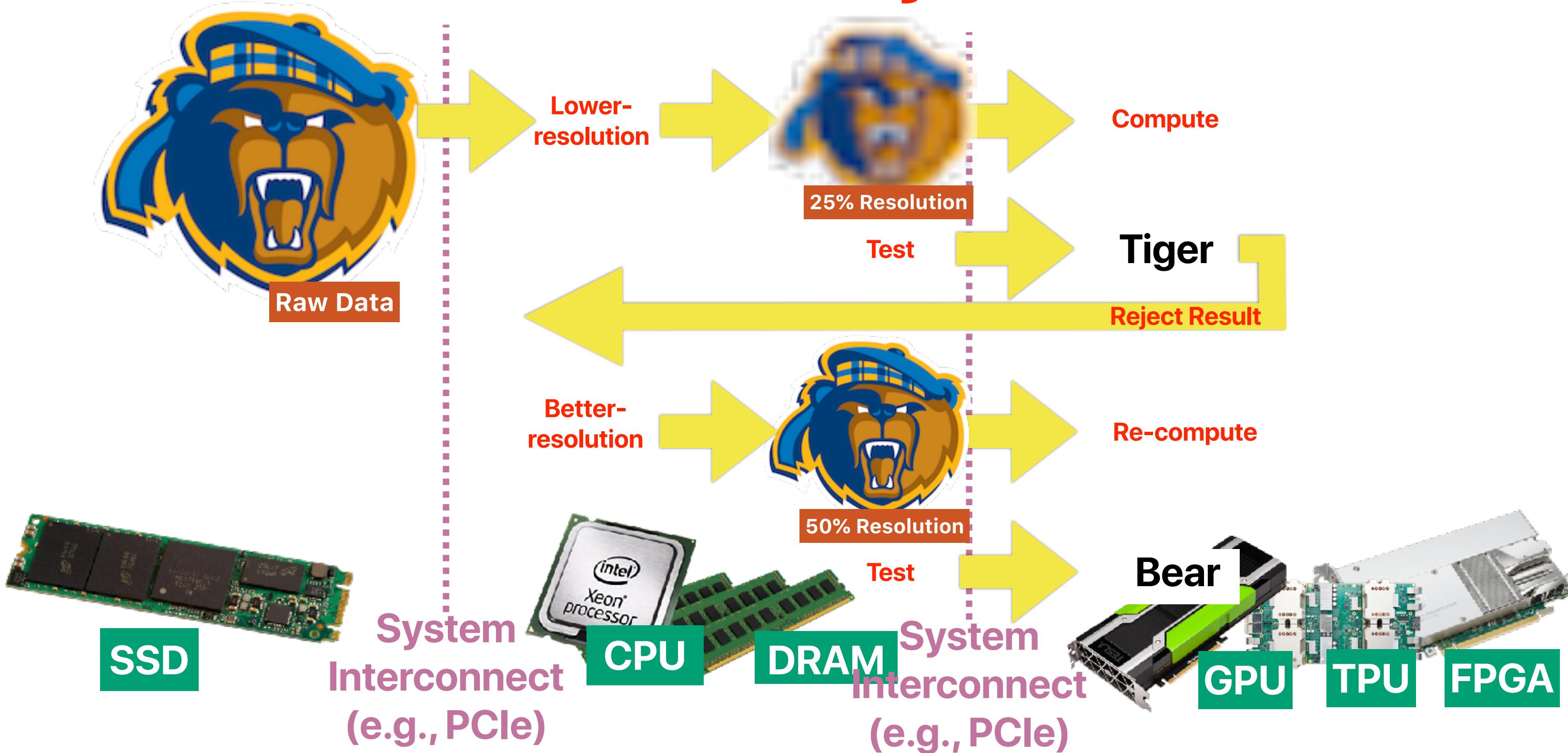
```
int setup(int argc, char **argv) {  
    // Skip – the rest of code ...  
    int infile;  
  
    // VS: Declare VS variables  
    struct vs_operator op[1];  
    struct vs_feedback fb[1];  
  
    // Open a file descriptor  
    infile = open(filename, O_RDONLY, "0600");  
  
    // Read precise data from the file descriptor  
    read(infile, &npoints, sizeof(int));  
    read(infile, &nfeatures, sizeof(int));  
  
    // Skip – some other initialization code ...  
  
    // VS: set parameters for desired operator  
    // PACKING(default)/PACKING_AF(autofocus)/VS_IF(iFilter)  
    op[0].op = PACKING;  
    op[0].resolution = HALF;  
  
    // Skip – some other initialization code ...  
    // VS: apply the desired VS operator for the file  
    vs_setup(infile, &op, "%f");  
    // VS: read data processed by the VS operator  
    vs_read(infile, buf, npoints*nfeatures*sizeof(float), &fb);  
    // VS: disable the usage of VS operator for the file  
    vs_release(infile);  
    // Skip – the rest of code ...  
    // VS: use approximate kernel if the operator succeed  
    if(fb[0].resolution == op[0].resolution)  
        cluster_approximate(...);  
    else  
        cluster(...);  
    // Skip – the rest of code ...  
}
```

```
void test_distributed_page_rank(char* graphfilename,  
                                int num_ofVertex, int num_ofEdges, int iterations)  
{  
    FILE *fin;  
    vs_stream input_stream;  
    void **arg_list;  
    fin = fopen(graphfilename, "r");  
    vs_input_stream = vs_stream_create(fin);  
    Edge *edge_array = (Edge *)malloc(  
                           sizeof(Edge)*num_ofEdges);  
    inputApplet(input_stream, edge_array);  
    // The rest of code ...  
}
```

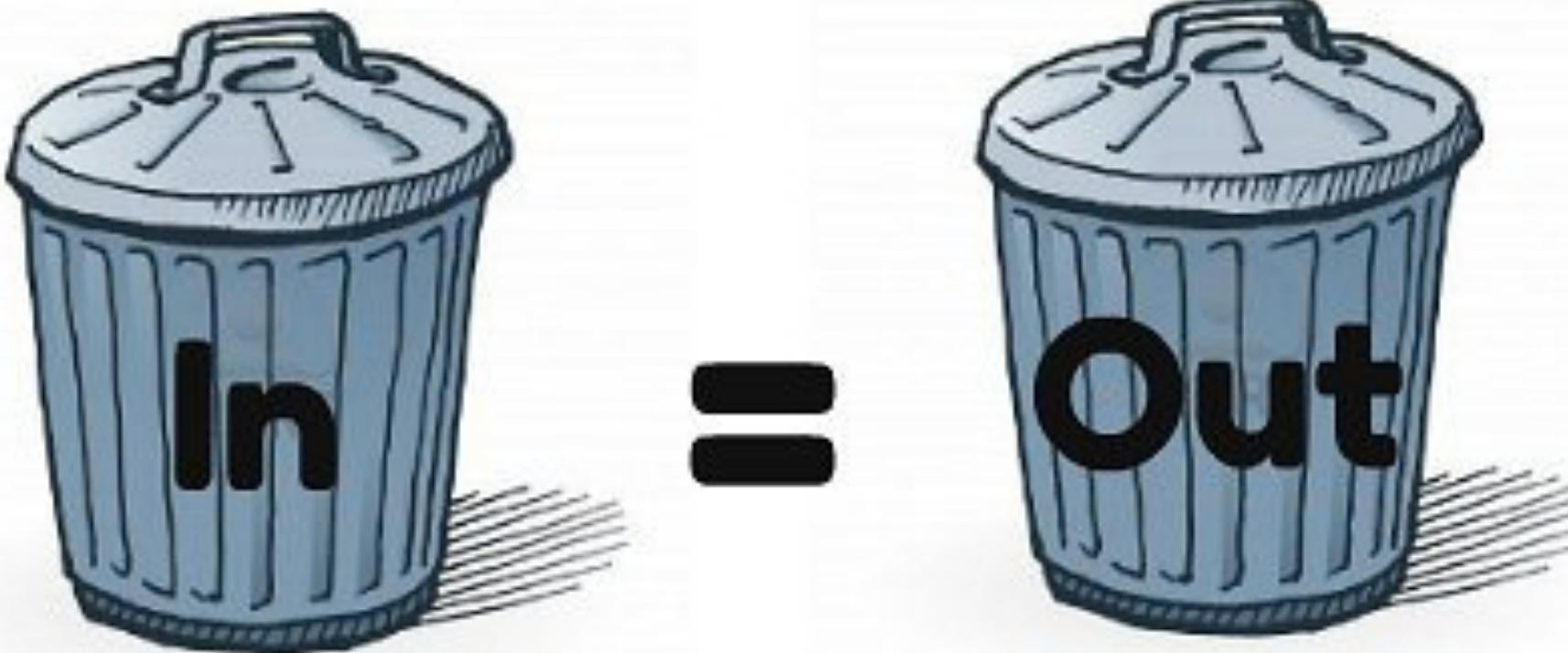
```
StorageApp int inputApplet  
(ms_stream ssd_input_stream, void *edge_array)  
{  
    Edge ssd_edge_array[4096];  
    int i = 0;  
    while(vs_scanf(ssd_input_stream, "%d %d",  
                  &ssd_edge_array[i%4096].first,  
                  &ssd_edge_array[i%4096].second) == 2)  
    {  
        if(i % 4096 == 0)  
        {  
            vs_memcpy(edge_array, ssd_edge_array,  
                      sizeof(Edge)*4096);  
            edge_array += sizeof(Edge)*4096;  
        }  
        vs_memcpy(edge_array, ssd_edge_array,  
                  sizeof(Edge)*(i%4096));  
    }  
    return i;  
}
```

Programmer has
to implement

Traditional Quality Control



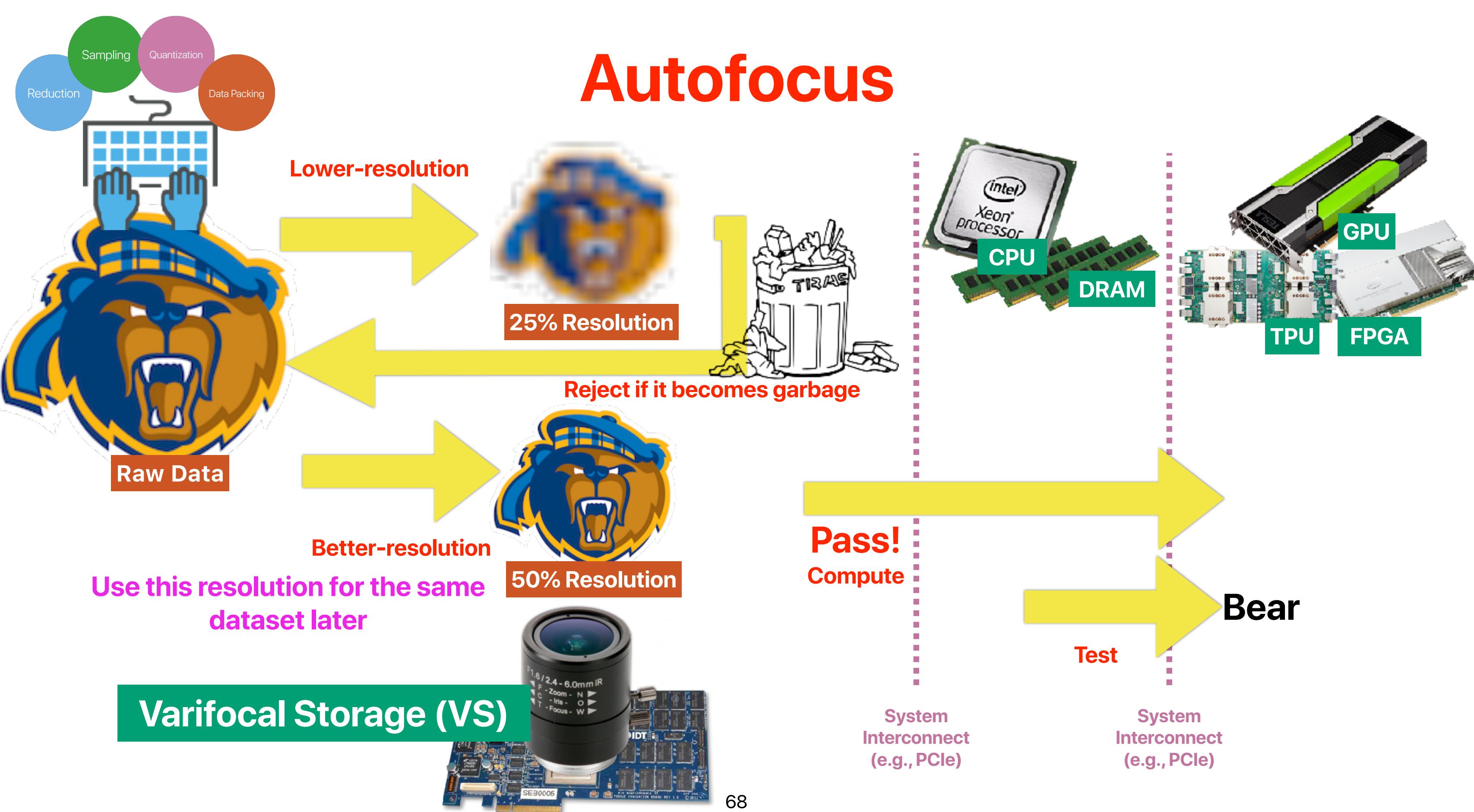
Garbage in, garbage out



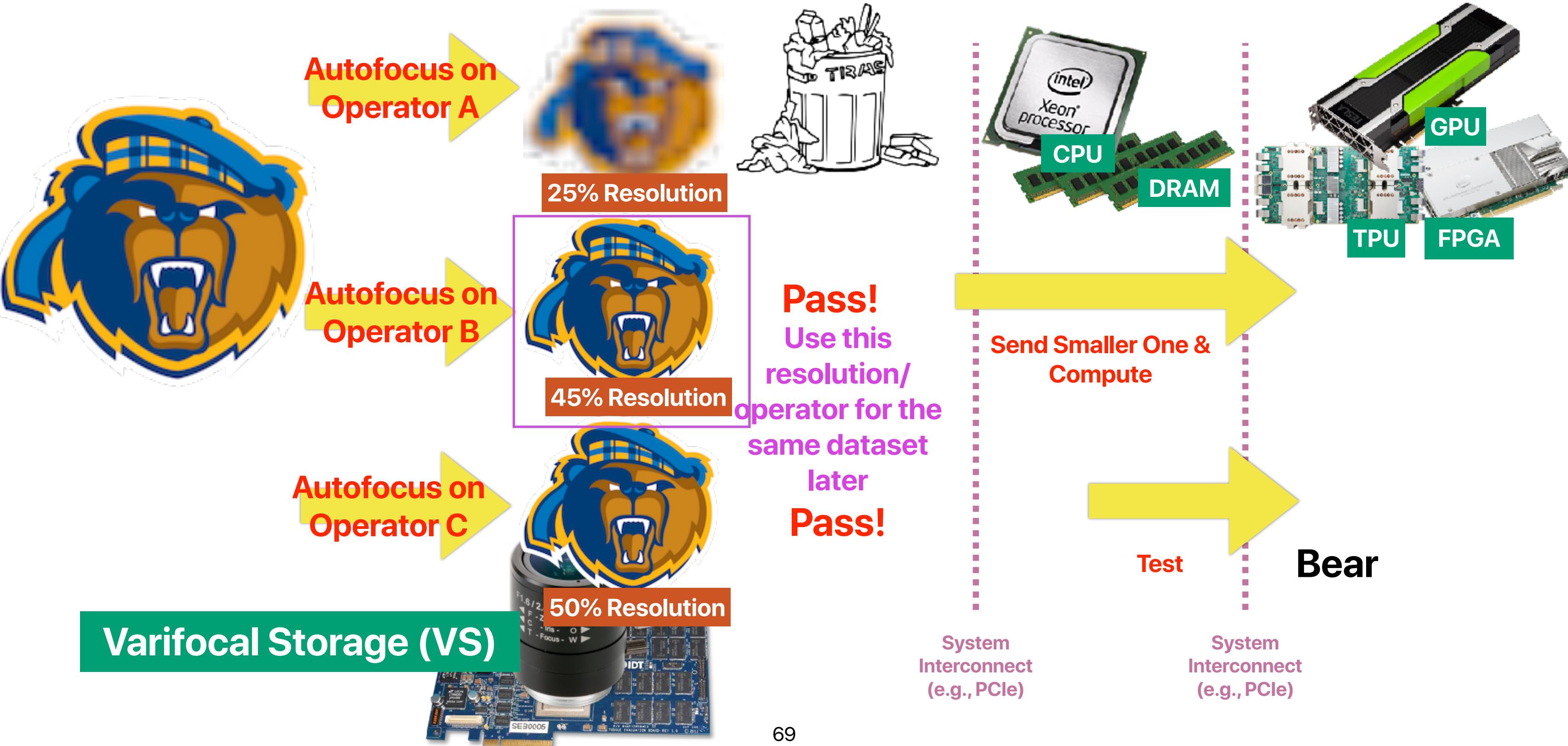
Quality Control



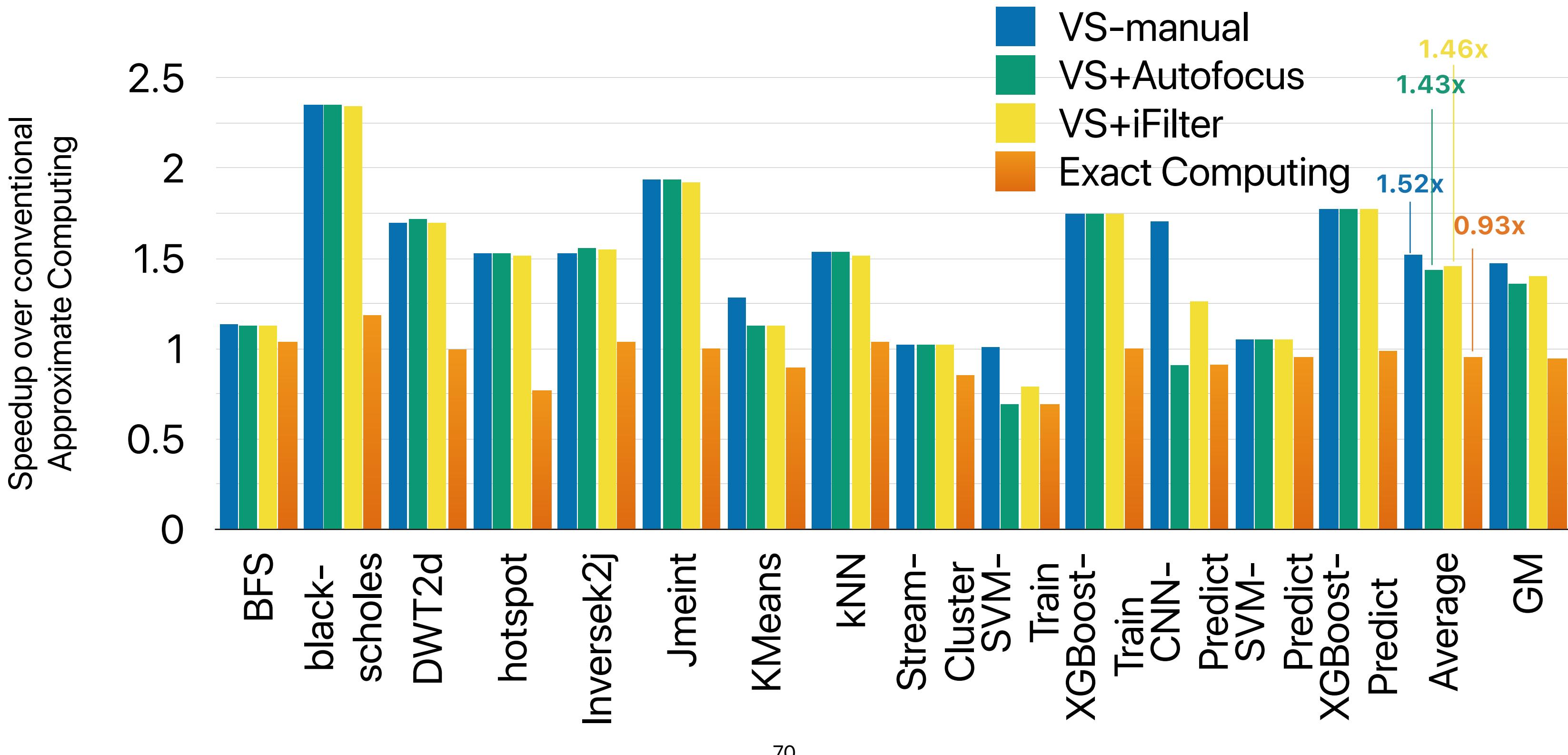
Autofocus



iFilter



Performance of Varifocal Storage (VS)



Announcement

- Final project presentation — 3/21 11:30a-2:30p (SSC 125)
 - Please book your order on the calendar <https://calendar.google.com/calendar/u/0/selfsched?sstoken=UURHeVJHWDhIUUhffGRIZmF1bHR8ZGQ5ODhhOGRhYjU1NGRjOWY1ODM3MmZiZjlmNmE0ODU>
 - Pizza provided
 - No lecture next Tuesday, but we'll be back on Thursday

Electrical Computer Science Engineering

277

つづく

