

Accelerating Intelligence

John D. Davis, Ph.D.

Hardware Accelerators Break Through the Processing Wall




On track to millions of servers w/FPGAs




Introduces FPGA powered server instances




TPUs help avoid cost of 12-15 data centers



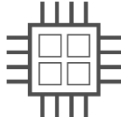

\$16.7B acquisition of Altera



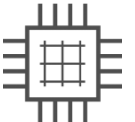

\$1B in datacenter revenues



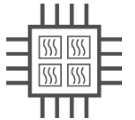

193% growth in datacenter segment



FPGA



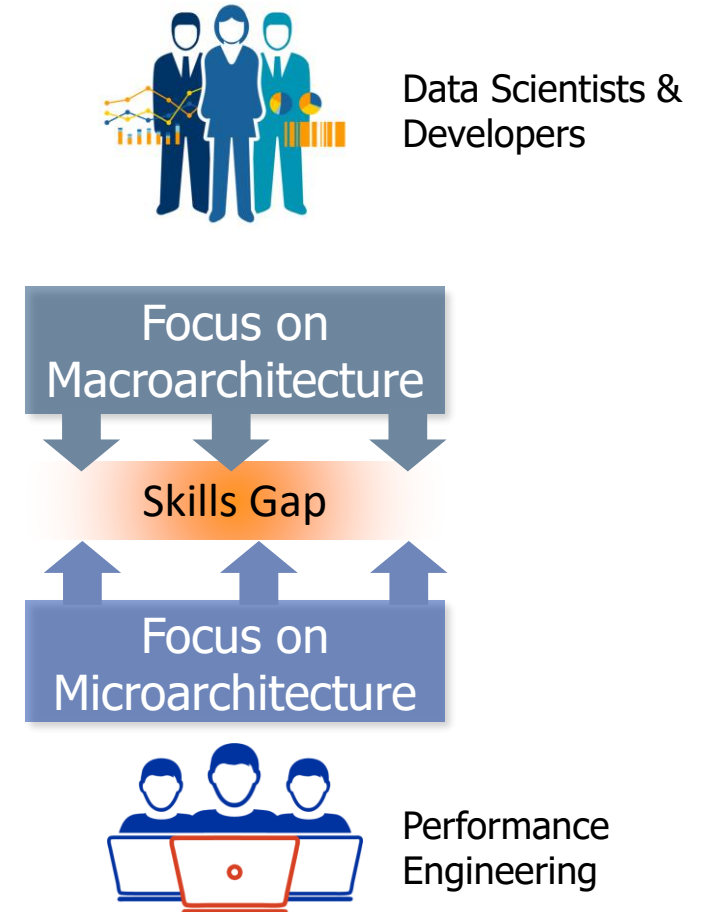
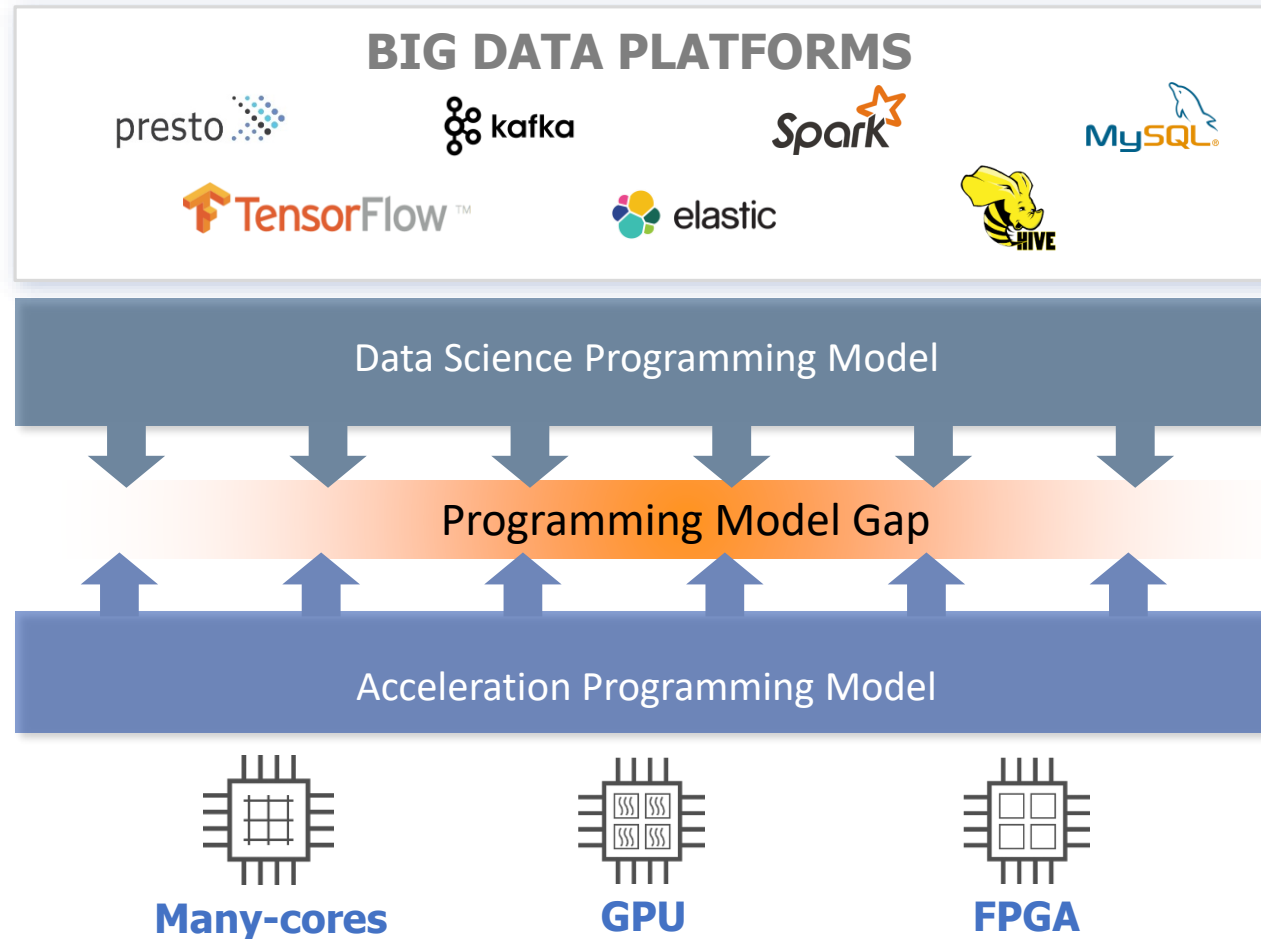
ASIC



GPU

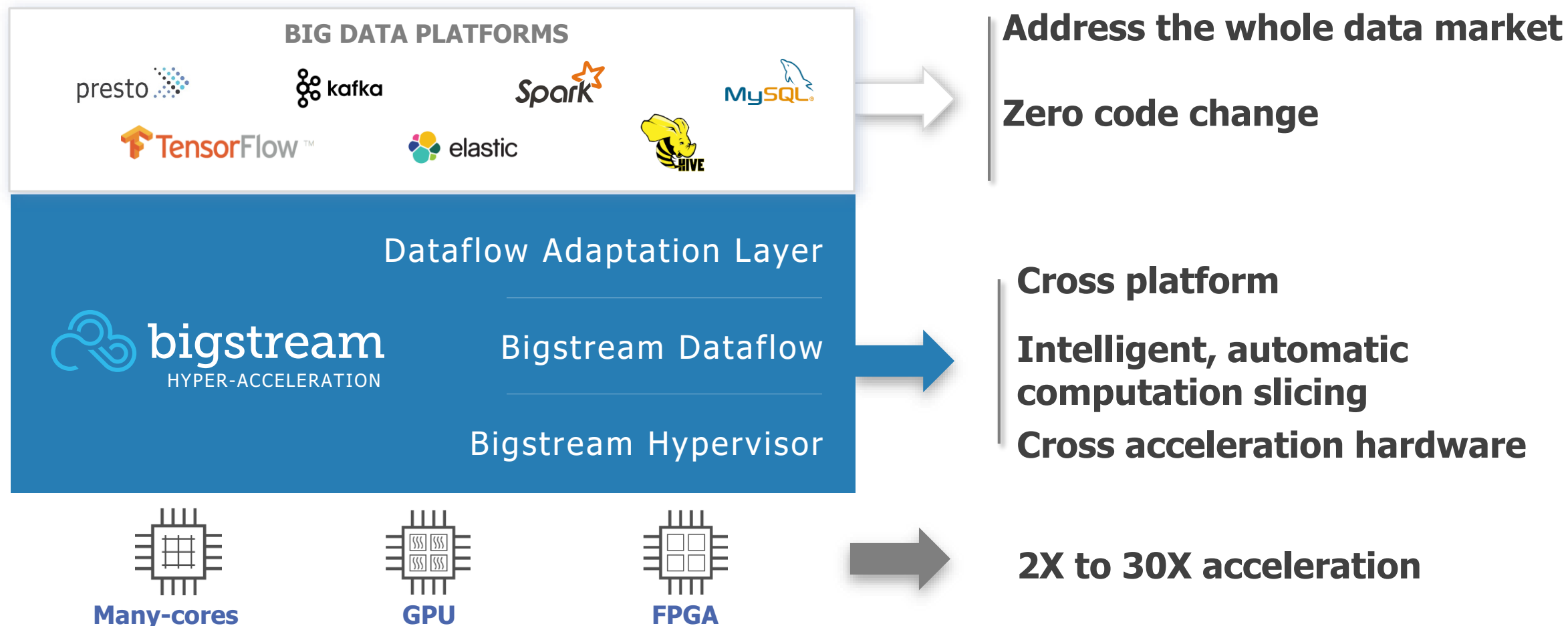
What about "Data Gravity"?
Push compute near the Data

Inhibitor: Programming Model Gap for Hardware Accelerators

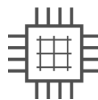
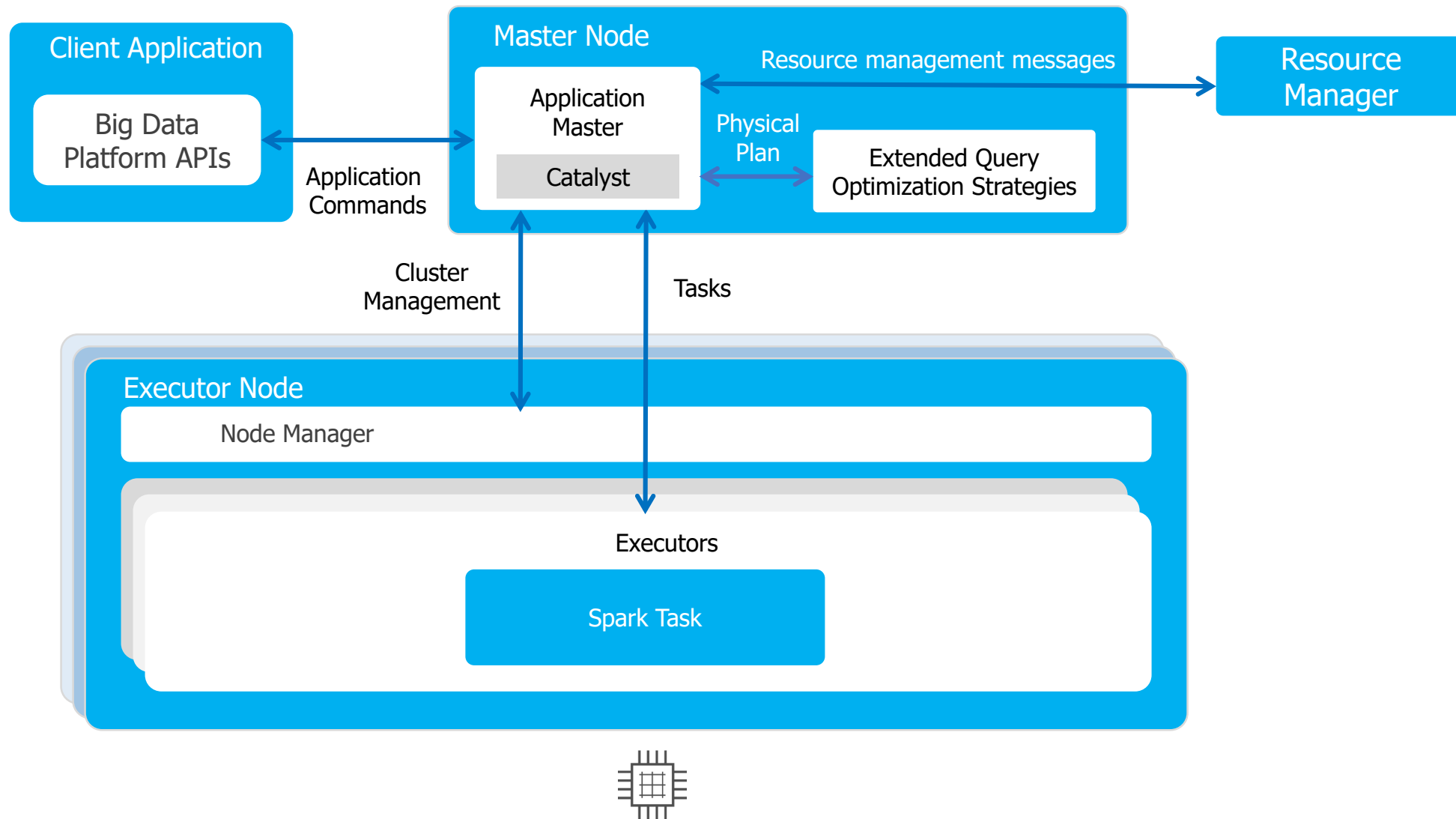


Introducing: Bigstream Hyper-acceleration Layer

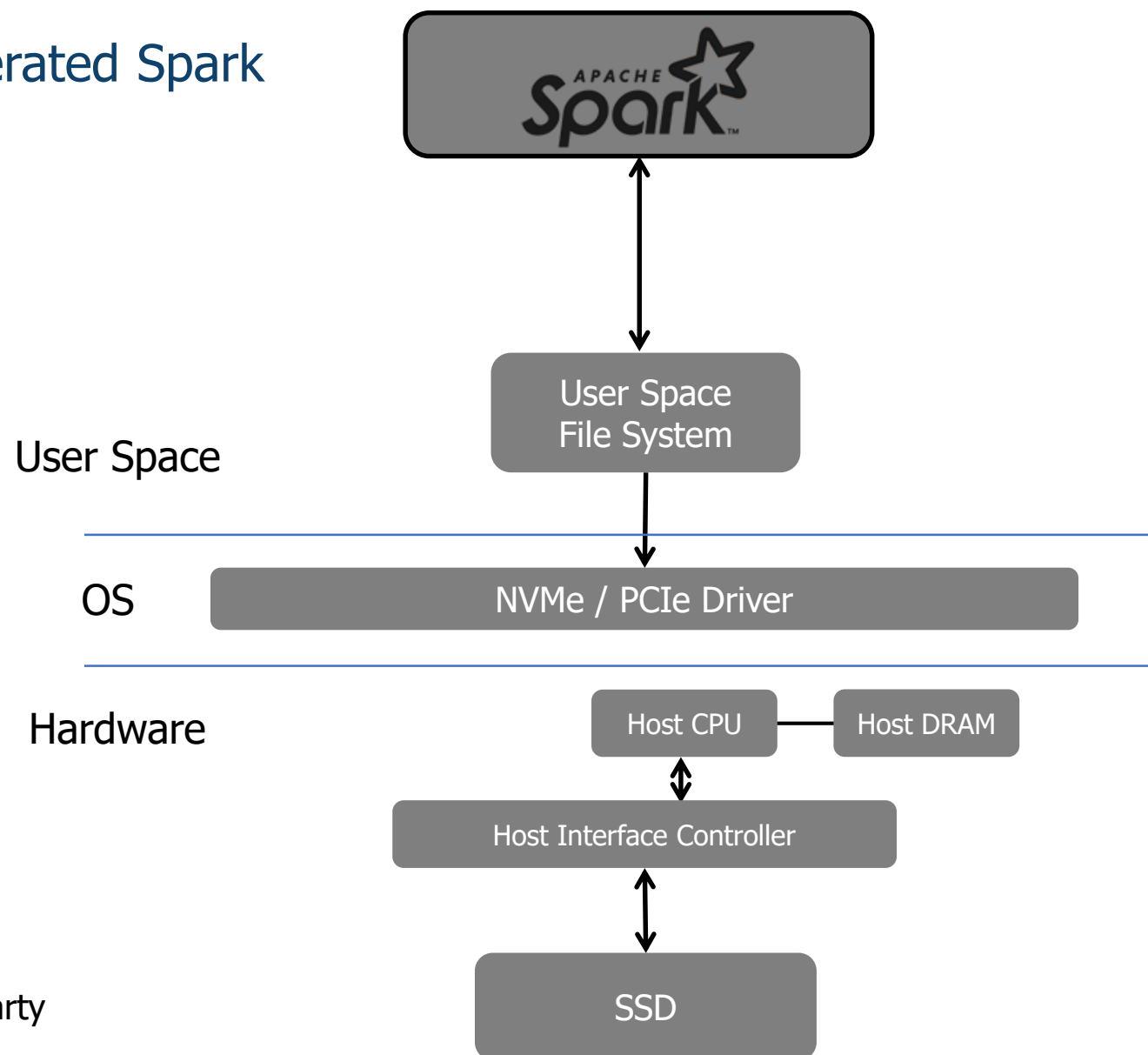
Goal: Provider of Open-ecosystem for Hyper-acceleration



Apache Spark

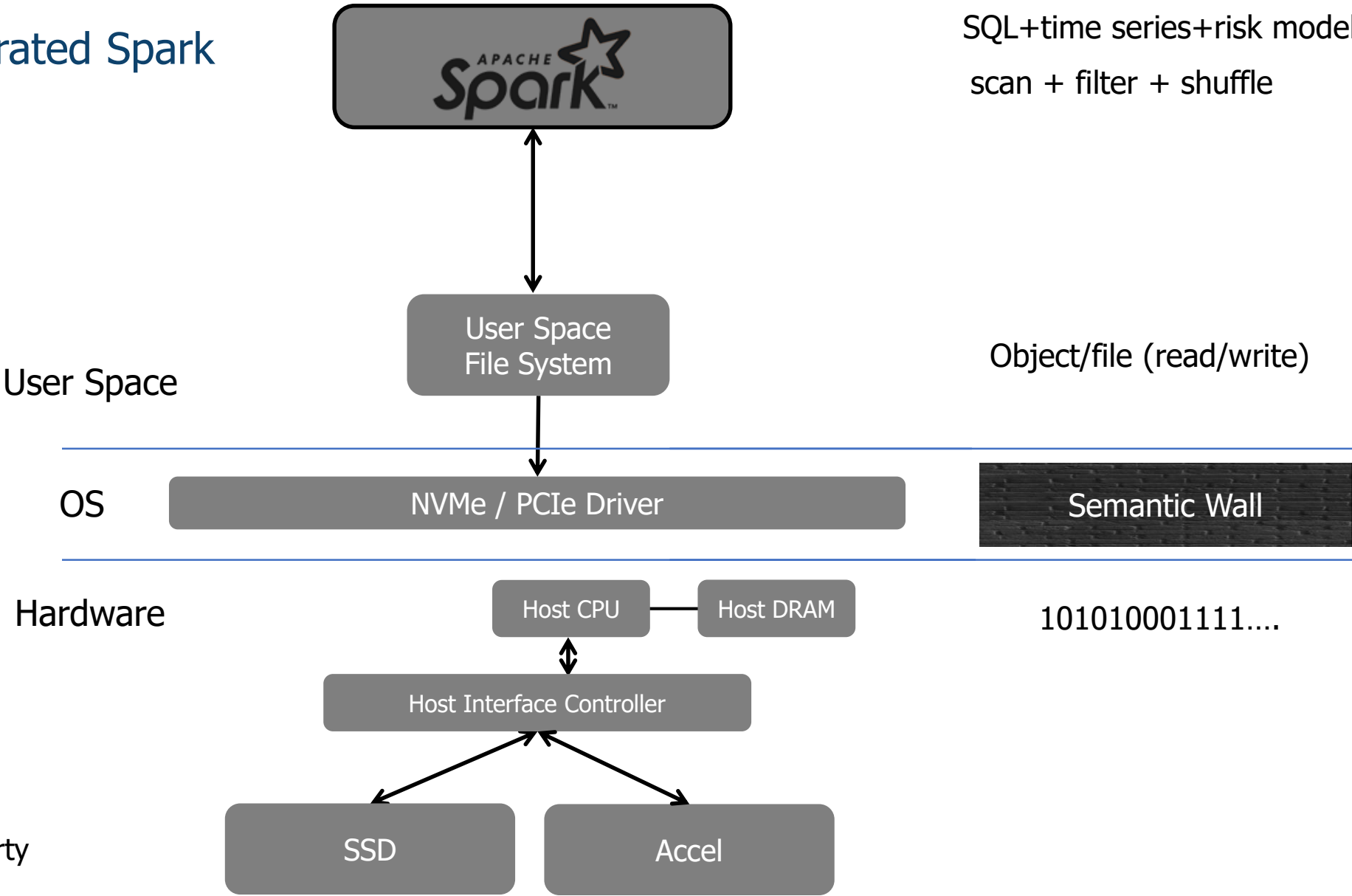


Non-accelerated Spark



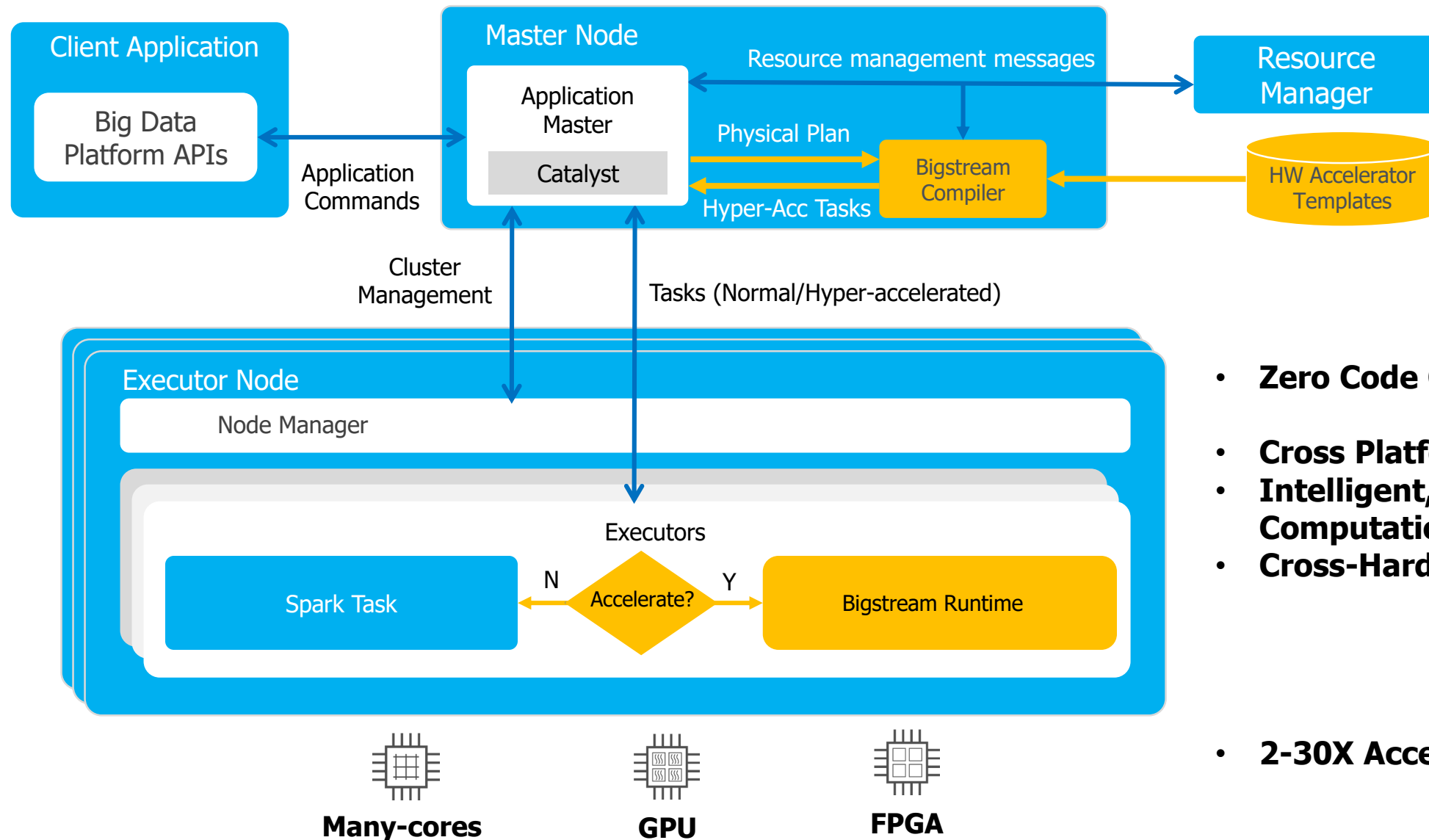
Non-accelerated Spark

SQL+time series+risk models + ...
scan + filter + shuffle



101010001111....

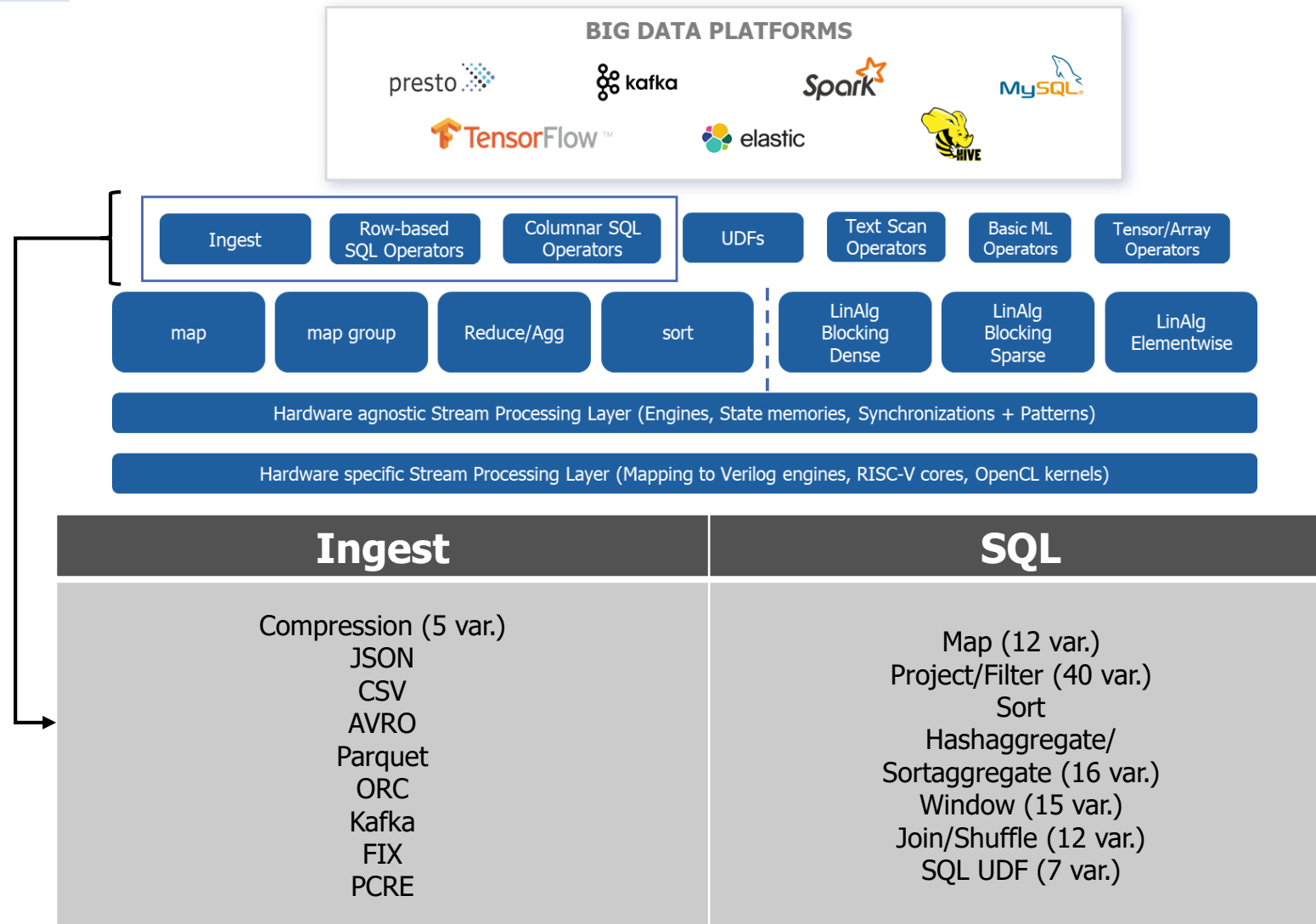
Bigstream Seamless Acceleration of Apache Spark











- **Zero Code Change**
- **Cross Platform**
- **Intelligent, Automatic Computation Slicing**
- **Cross-Hardware Acceleration**

- **2-30X Acceleration**

Hyper-Acceleration Layered Compilation Approach



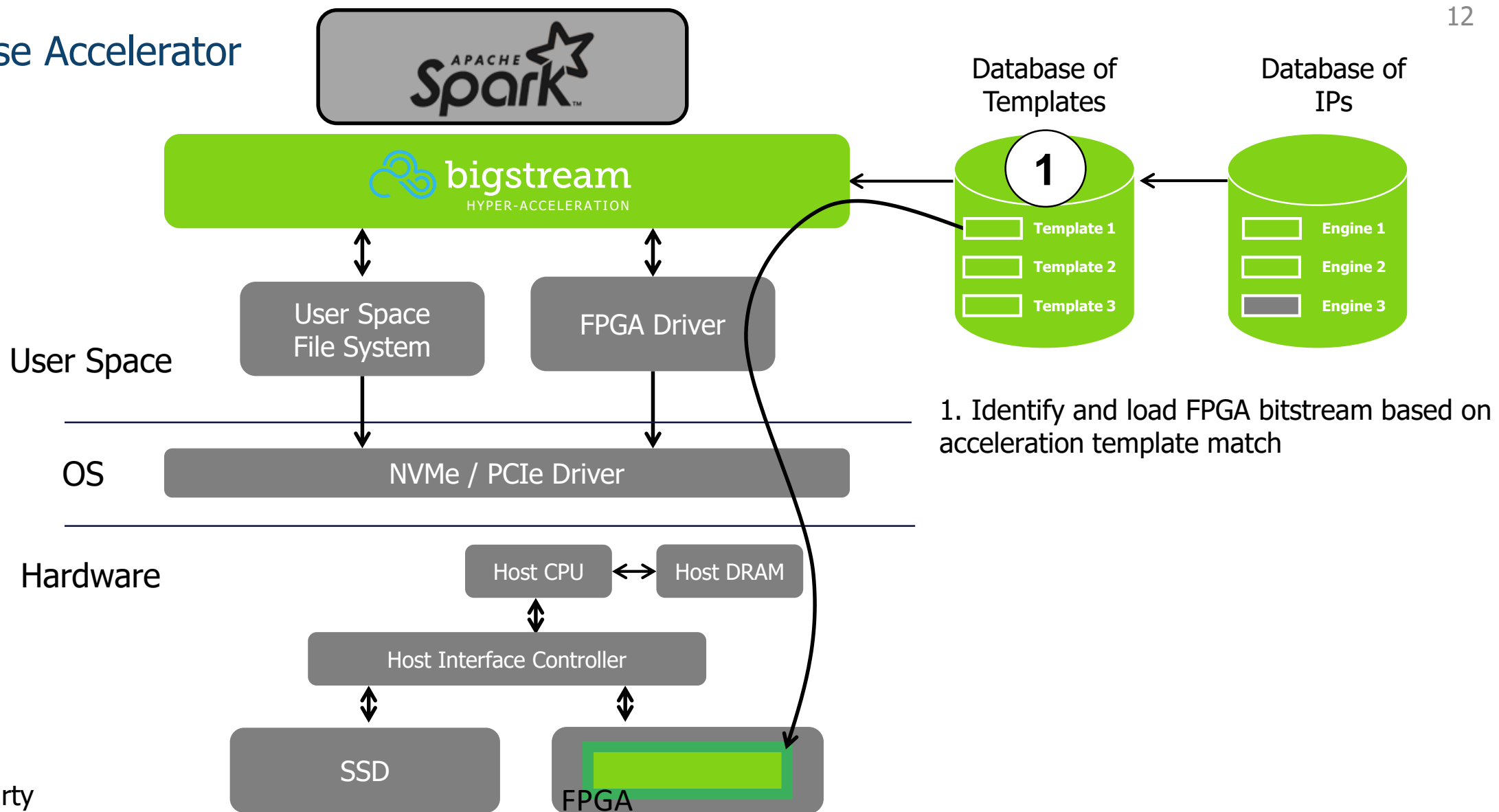
HW Accelerator Engines

 Deserialization	JSON	(.*) Search/Regex	PCRE (3 rd Party)
	CSV	 CPU Cores	RISC-V (3 rd Party)
	Parquet (under development)		
	FIX (under development)		
 Decompression	Snappy	 Machine Learning	Linear/Logistic Regression
	GZIP (3 rd Party)		K-means
 Encryption/Decryption	AES	 Deep Learning	CNN (3 rd Party)
 SQL	Project		RNN (3 rd Party)
	Filter	 Networking	IP/UDP
	Sort		IP/TCP (3 rd Party)
	Hash Aggregate		

General Application and HW Characteristics

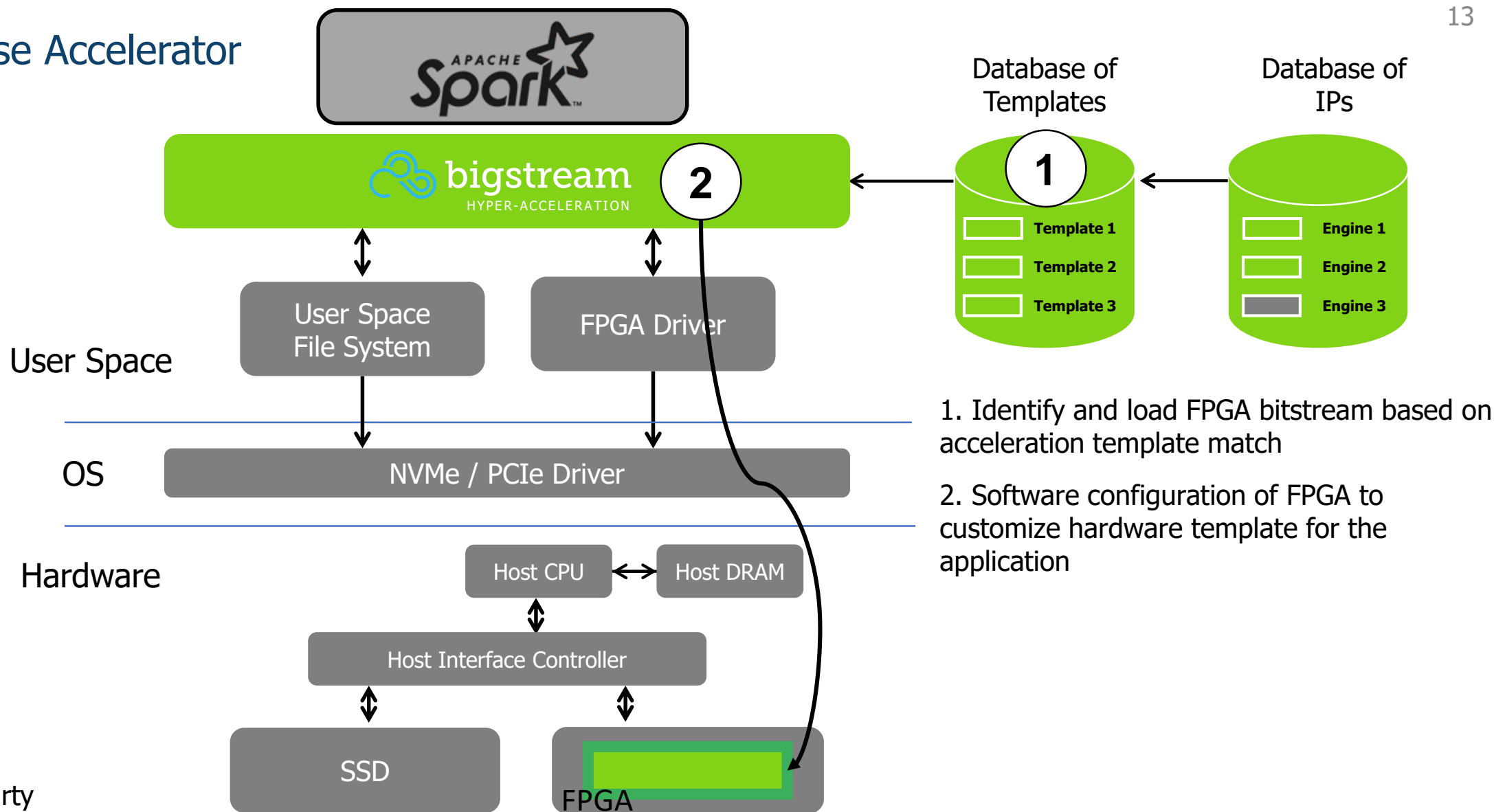
- Application characteristics:
 - Batch and streaming analytics applications
 - Collection of computation/transformation stages
 - Terabytes to Petabytes or more of data
 - Large clusters: 10's – 10,000's servers
 - I/O (network and/or storage) and compute bound
- How we use the FPGA:
 - Build out pipelines and import as kernels
 - Computational overlays (Time Division Multiplexing)
 - Offload and Inline (Storage and Network) acceleration
 - Moving between Streaming, Memory Mapped, and P2P DMA with Software
 - Bringing Computation to the Data.

Offload Base Accelerator Ecosystem

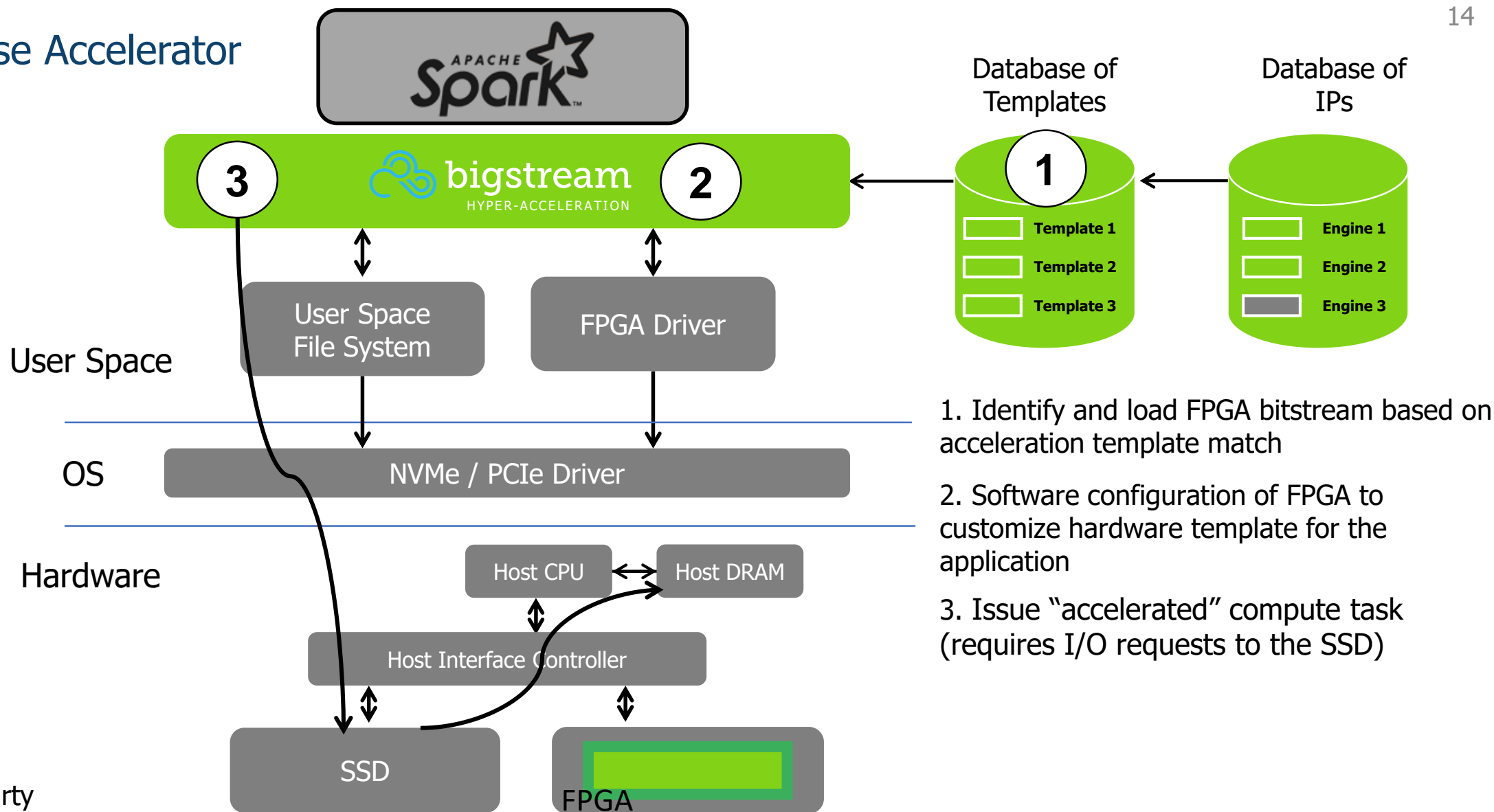


1. Identify and load FPGA bitstream based on acceleration template match

Offload Base Accelerator Ecosystem

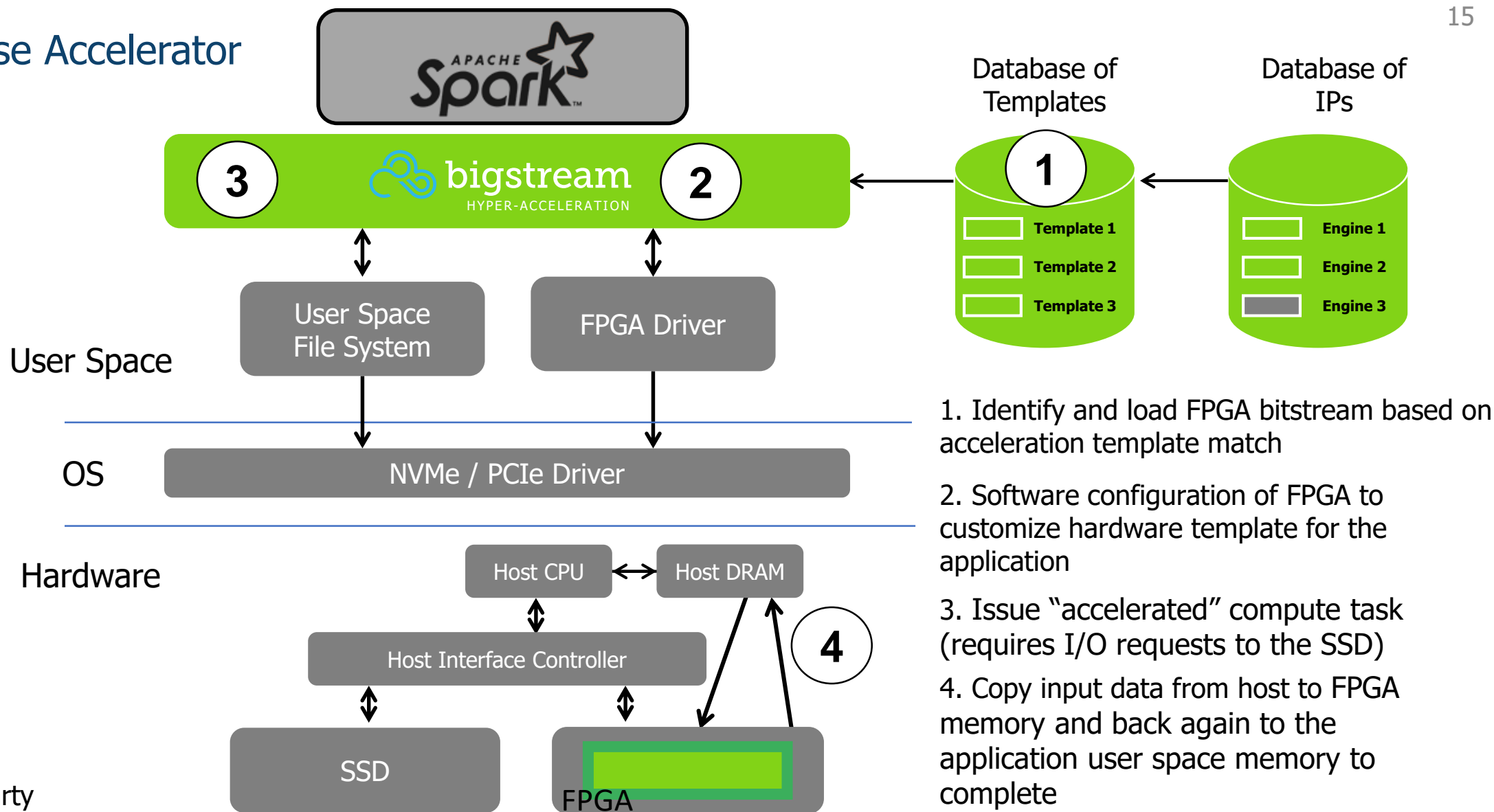


Offload Base Accelerator Ecosystem



1. Identify and load FPGA bitstream based on acceleration template match
2. Software configuration of FPGA to customize hardware template for the application
3. Issue "accelerated" compute task (requires I/O requests to the SSD)

Offload Base Accelerator Ecosystem



SQL+time series+risk models + ...

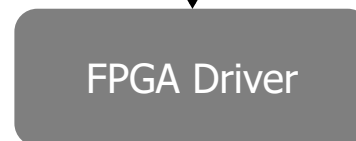
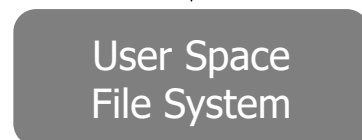
scan + filter + shuffle

Object/file (read/write)

User Space

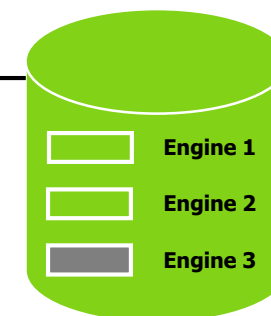
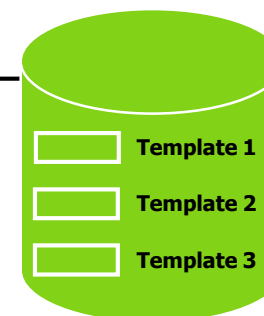
OS

Hardware



Database of
Templates

Database of
IPs



1. Identify and load FPGA bitstream based on acceleration template match

2. Software configuration of FPGA to customize hardware template for the application

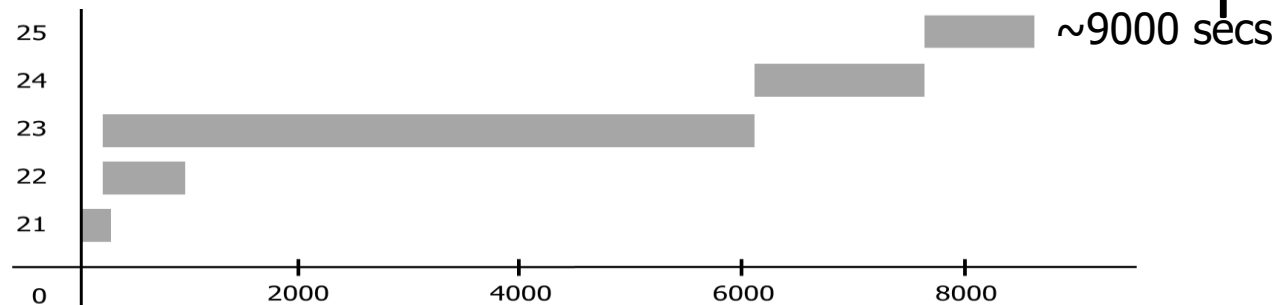
HW Acceleration with
Application Intelligence



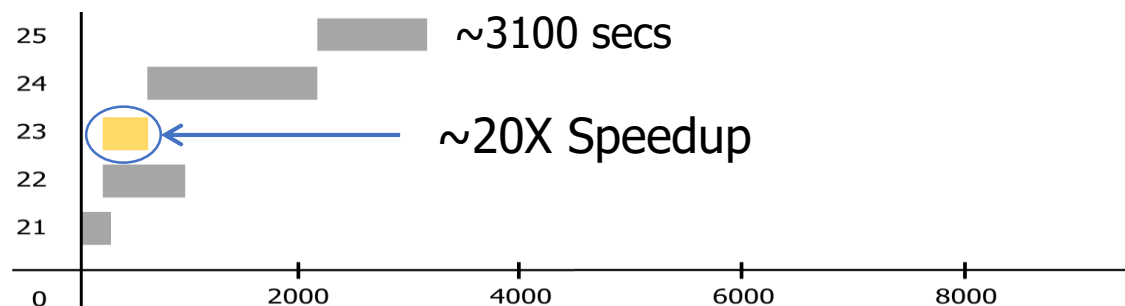
OSS/3rd Party

bigstream (Existing Stack)

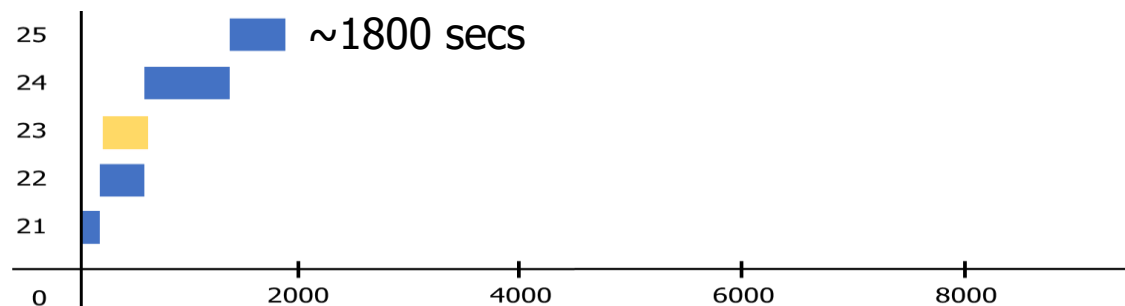
Overall Acceleration Comparison



Baseline – No Acceleration

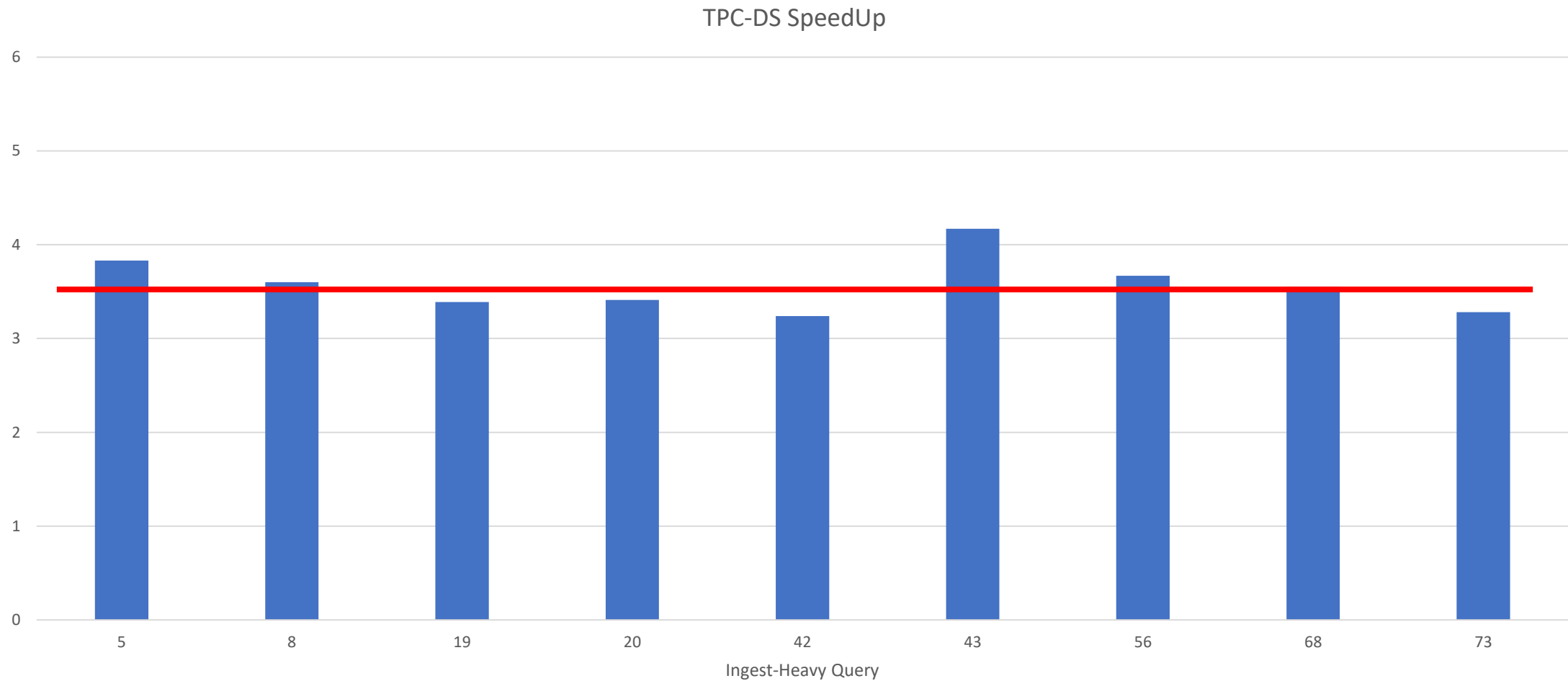


Hardware Acceleration Only: **~3x**



Software + Hardware Acceleration: **~5x**

Offload Performance: 5 Node Cluster

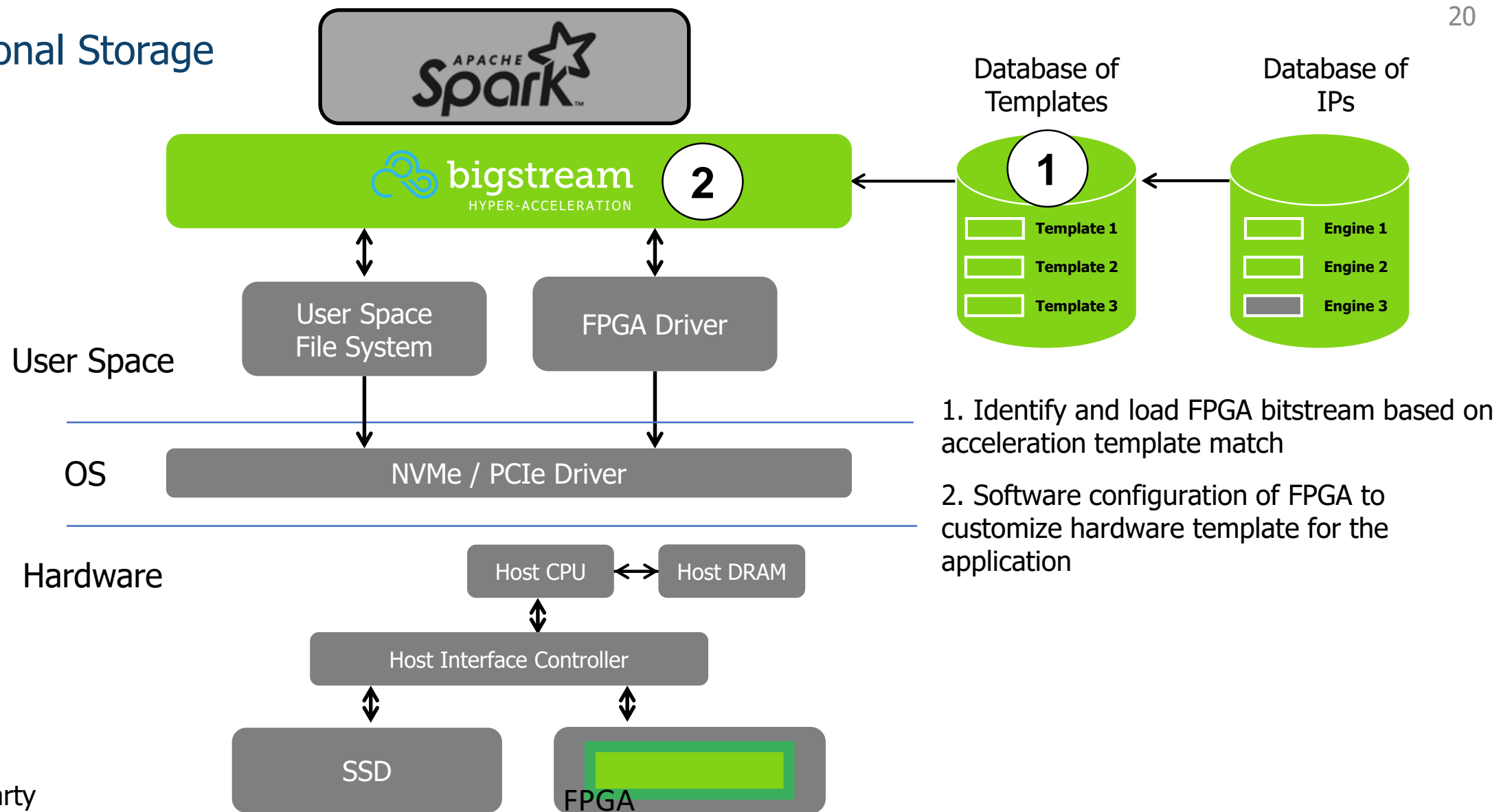


1 Master Node + 4 Worker Nodes
Spark baseline vs FPGA acceleration with Zero Code Change

- Rewind for Computational Storage



Computational Storage



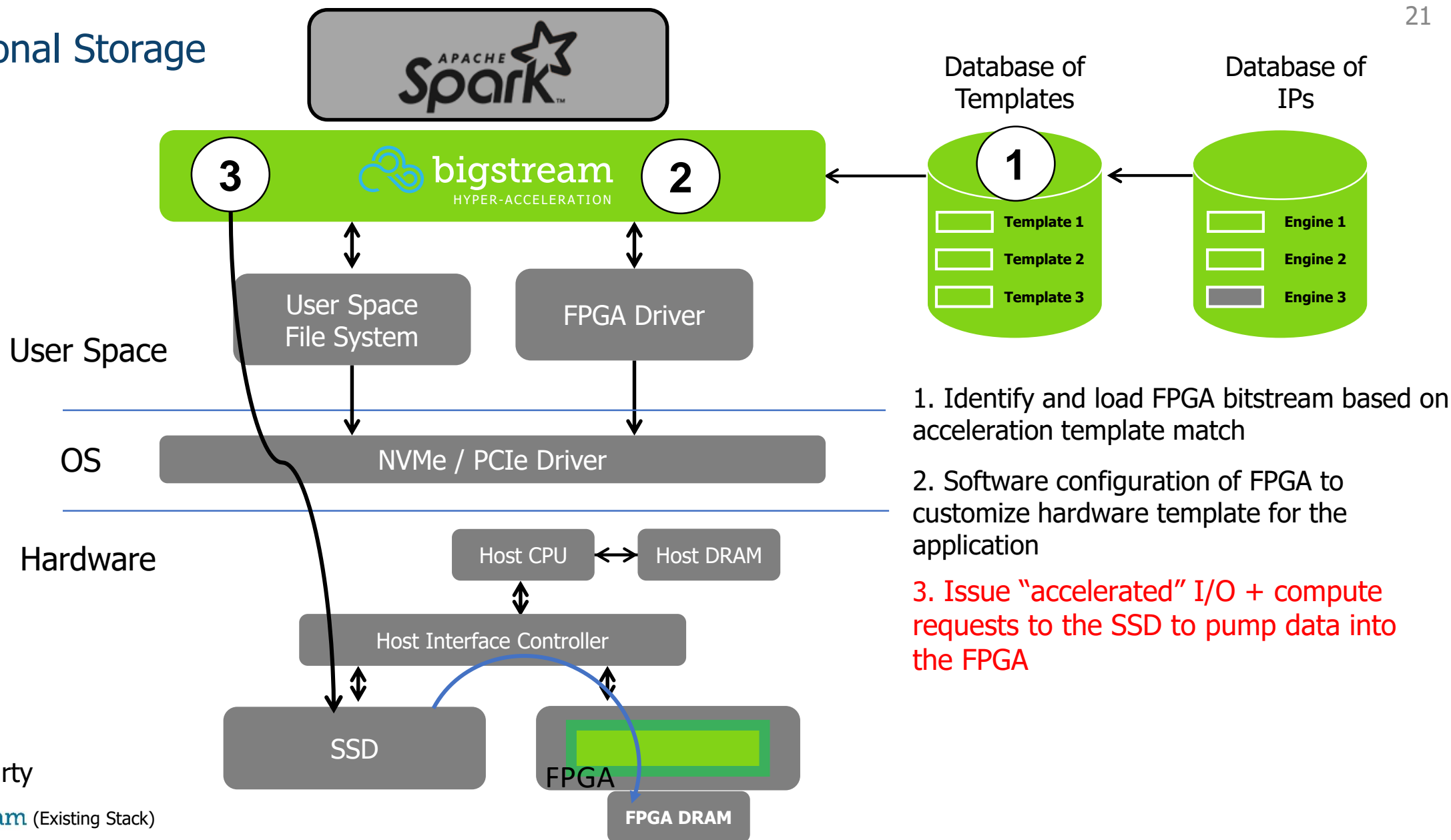
1. Identify and load FPGA bitstream based on acceleration template match

2. Software configuration of FPGA to customize hardware template for the application

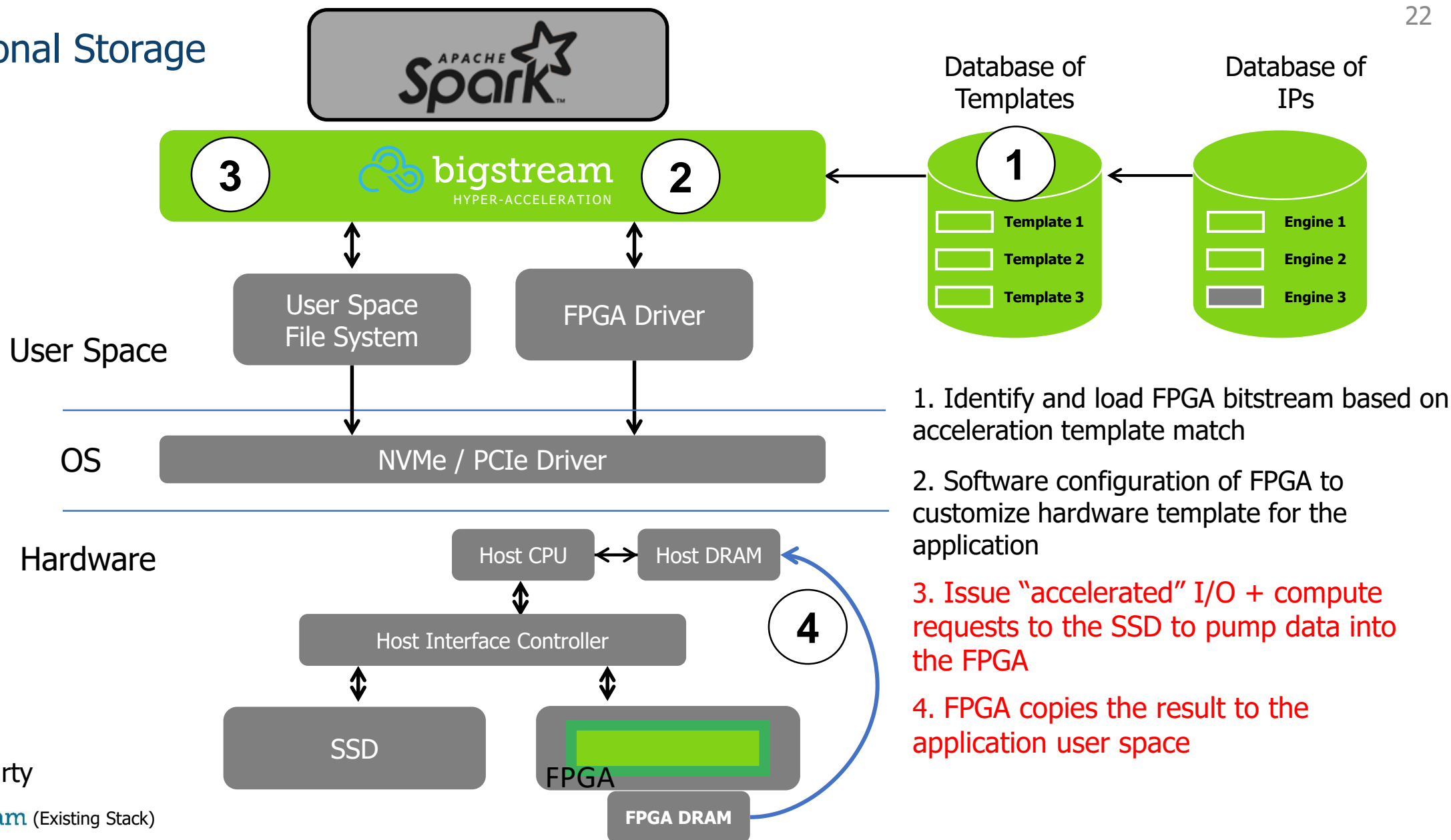
OSS/3rd Party

bigstream (Existing Stack)

Computational Storage



Computational Storage



OpenCL Code

- Memory-mapped DMA

> P2P DMA

// Read from SSD part

```
buffer = (char*) malloc (size);
```

```
bytes_read = fread (buffer, 1, size, _fd);
```

//WRITE PART- from host to fpga

```
err = clEnqueueWriteBuffer (commands, buffer, CL_TRUE, 0,
(data_actual_size + 96), json_write_host_mem, 0, NULL, & write_event);
```

```
clWaitForEvents(1, & write_event);
```

//Actual compute task

```
err = clEnqueueTask (commands, kernel, 0, NULL, & kernel_event);
```

```
clWaitForEvents(1, & kernel_event);
```

//READ PART - from fpga to host

```
err |= clEnqueueReadBuffer (commands, json_data_ptr,
CL_TRUE, s_read_offset, bytes_to_read + 32, son_read_host_mem, 0, NULL, & readevent2);
```

// Read from SSD & write to fpga

```
clCreateBuffer (context[1], CL_MEM_READ_ONLY | CL_MEM_EXT_PTR_XILINX,
(aligned_size), & xmem_flags, & err);
```

```
_mapped_virtual_addr = clEnqueueMapBuffer (commands, d_axi00_ptr0_wr,
CL_TRUE, CL_MAP_READ | CL_MAP_WRITE, 0, (aligned_size), 0, NULL, & map_event, & err);
```

```
clWaitForEvents(1, & map_event);
```

```
bytes_read = read (_fd, ((char *) _mapped_virtual_addr), size);
```

//Actual compute task

```
err = clEnqueueTask (commands, kernel, 0, NULL, & kernel_event);
```

```
clWaitForEvents(1, & kernel_event);
```

//READ PART -from fpga to host

```
err |= clEnqueueReadBuffer (commands, json_data_ptr, CL_TRUE,
s_read_offset, bytes_to_read + 32, son_read_host_mem, 0, NULL, & readevent2);
```

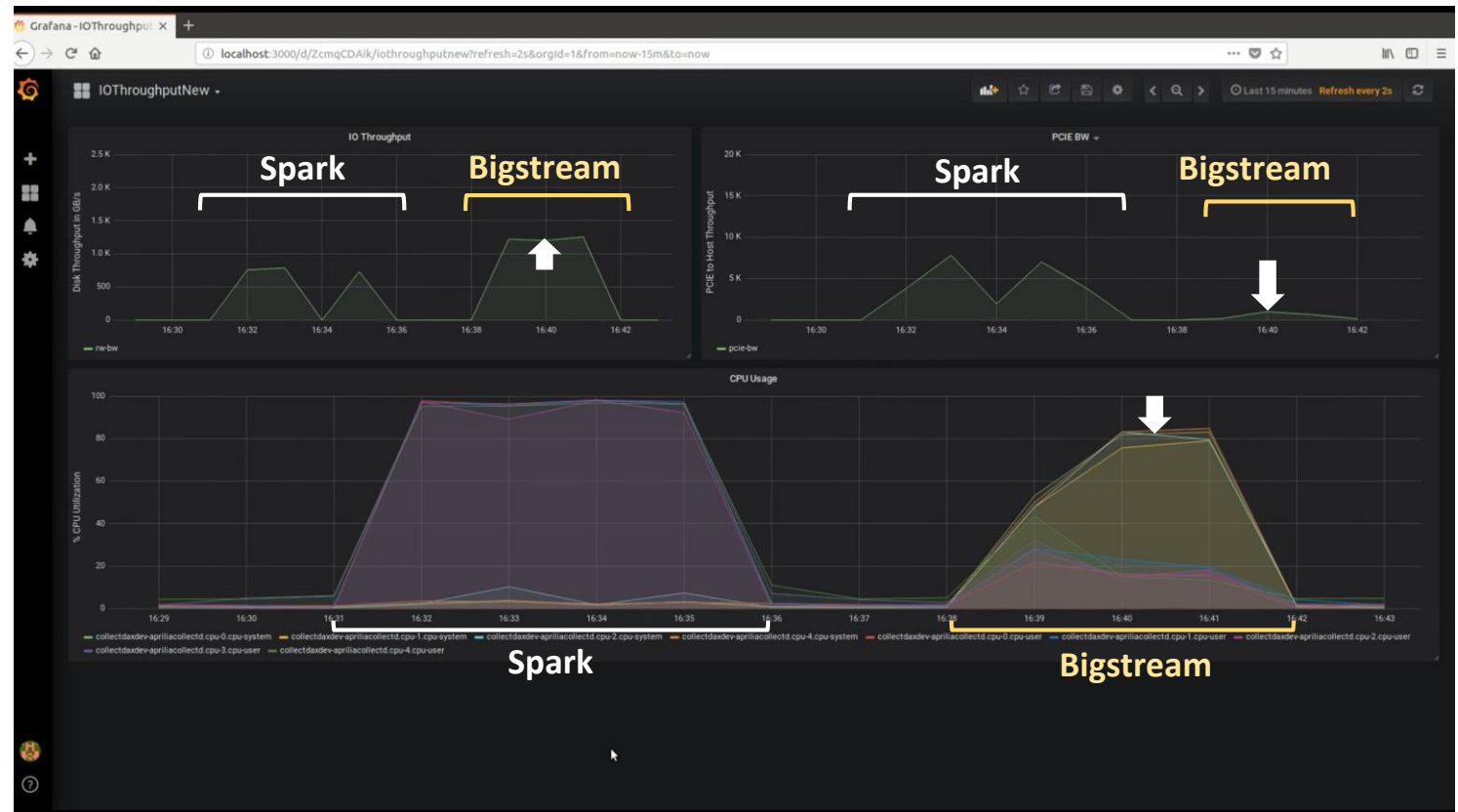
Same code can run in VU9P+ and ZU19, just change the target device

P2P Comparison

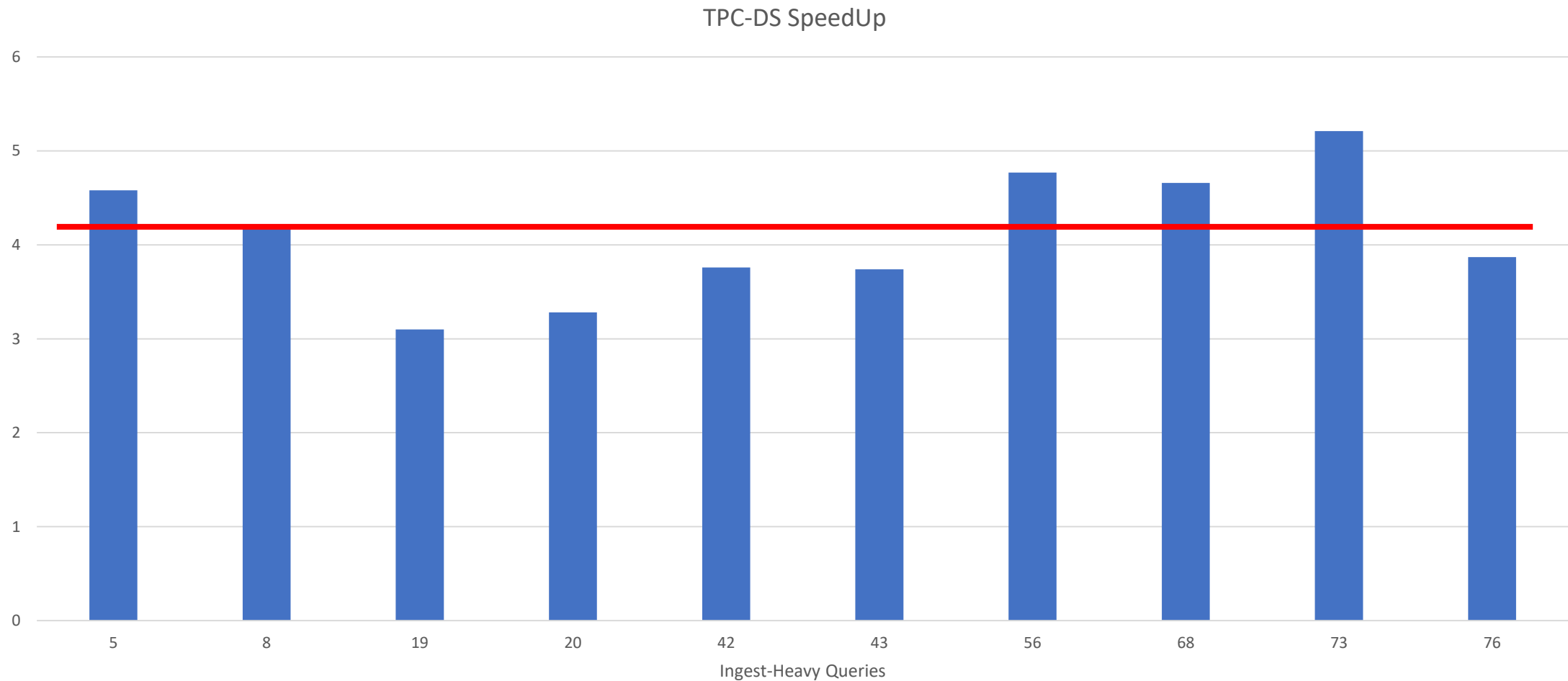
- SmartSSD (SSD + FPGA)
 - 45% Utilization (BRAM, LUTs, DSPs)
 - System level results 6c Xeon Workstation
- Spark vs P2P DMA
- Single FPGA, FPGA Pooling, & Clusters of FPGAs
 - TPC-DS with 340 GB to 2.7 TB
 - Up to 2 SmartSSDs/server
 - 5 node cluster (1 master node and 4 worker nodes)

System-Level Results

- Row-based engine with TPC-DS queries
 - Lower CPU utilization and fewer cores
 - Higher device I/O bandwidth
 - Lower host PCIe bandwidth
 - Reduced CPU DRAM Utilization
- Opportunities
 - Not all DMAs are the same.
 - Flexibility to mix and match DMA
 - When to use P2P?



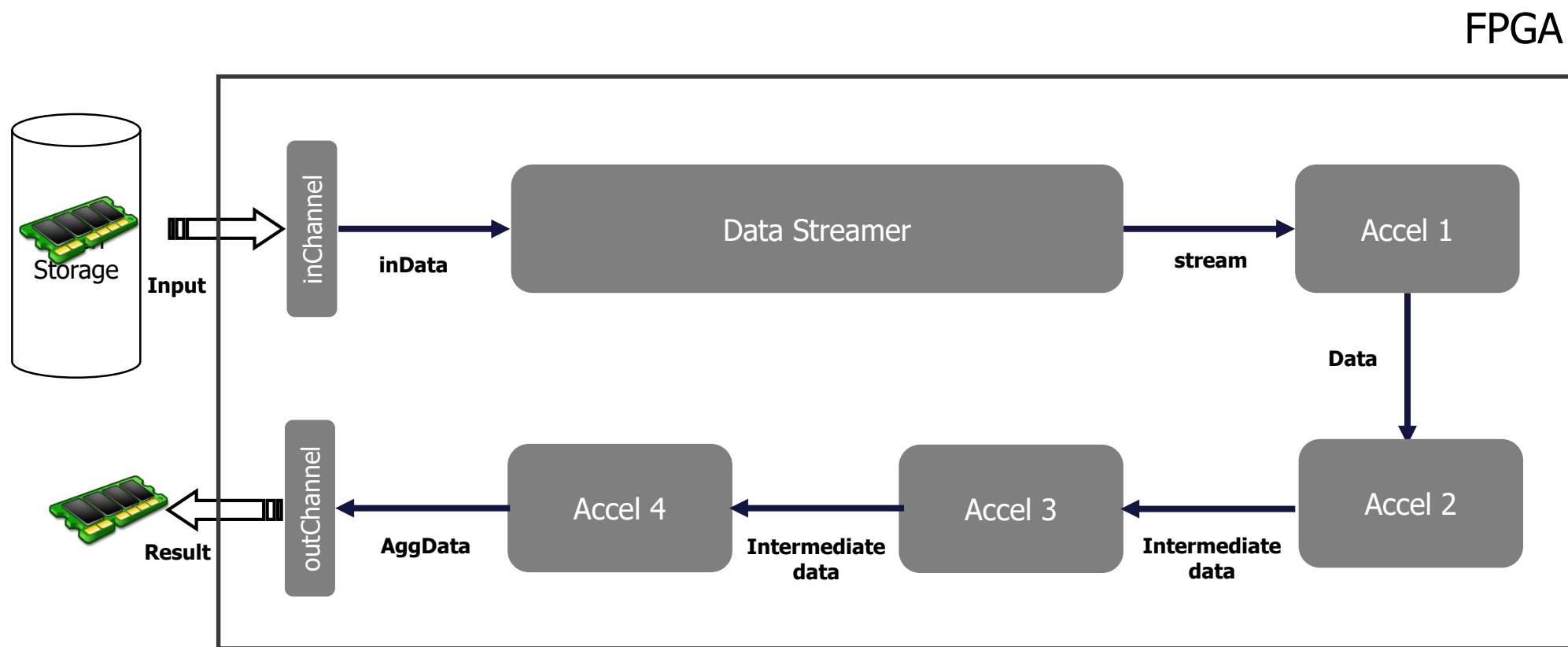
P2P Performance: 5 Node Cluster



1 Master Node + 4 Worker Nodes

Spark baseline vs FPGA acceleration with Zero Code Change

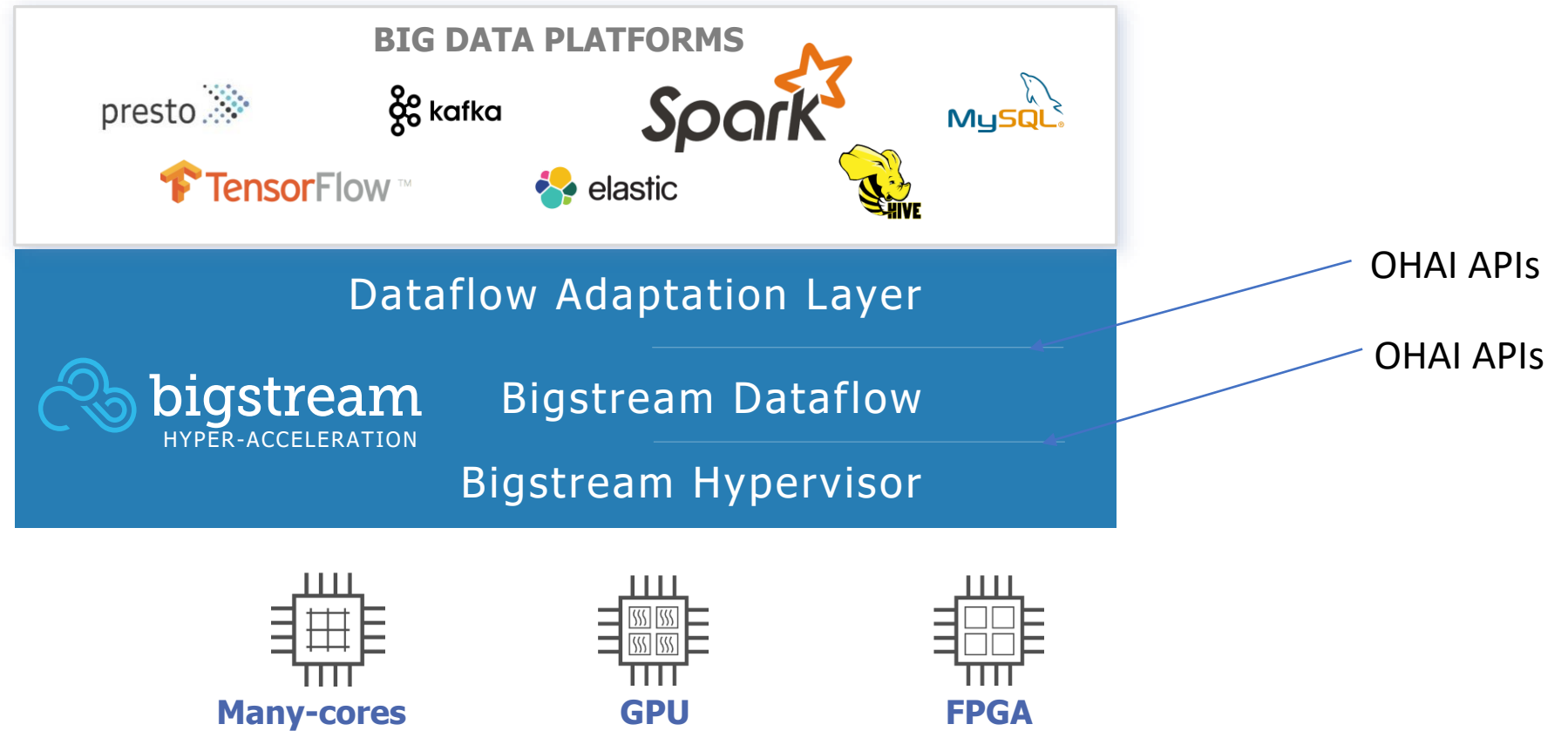
Accelerators Everywhere!



Open Hyper-Acceleration Initiative

- Goal: Extend the development community around hyper-acceleration
- Enable different developers from completely differently skill set
 - HW Template and engine developers
 - HW platform developers (add new FPGAs and FPGA cards)
 - Add new Platforms (i.e., Presto, Hive, TF, etc.)
- Invitation only contributions at this time
 - If interested, please come talk to me

OHAI Open APIs



Conclusions

- Flexibility
 - Accelerating more than just compute: moving compute to the data vs. data to compute
- Programmability
 - Time division multiplexing/hardware overlays
 - Software and Hardware configuration
- Generality
 - Combining HW and SW acceleration transparently embedded in Big Data Platforms
- Automatically
 - Bigstream Hyper-Acceleration with Zero Code Change



Thank You

John D. Davis, Ph.D.
john@bigstream.co