# HyperFlow: An Efficient Dataflow Architecture for Multi CPU-GPU Systems
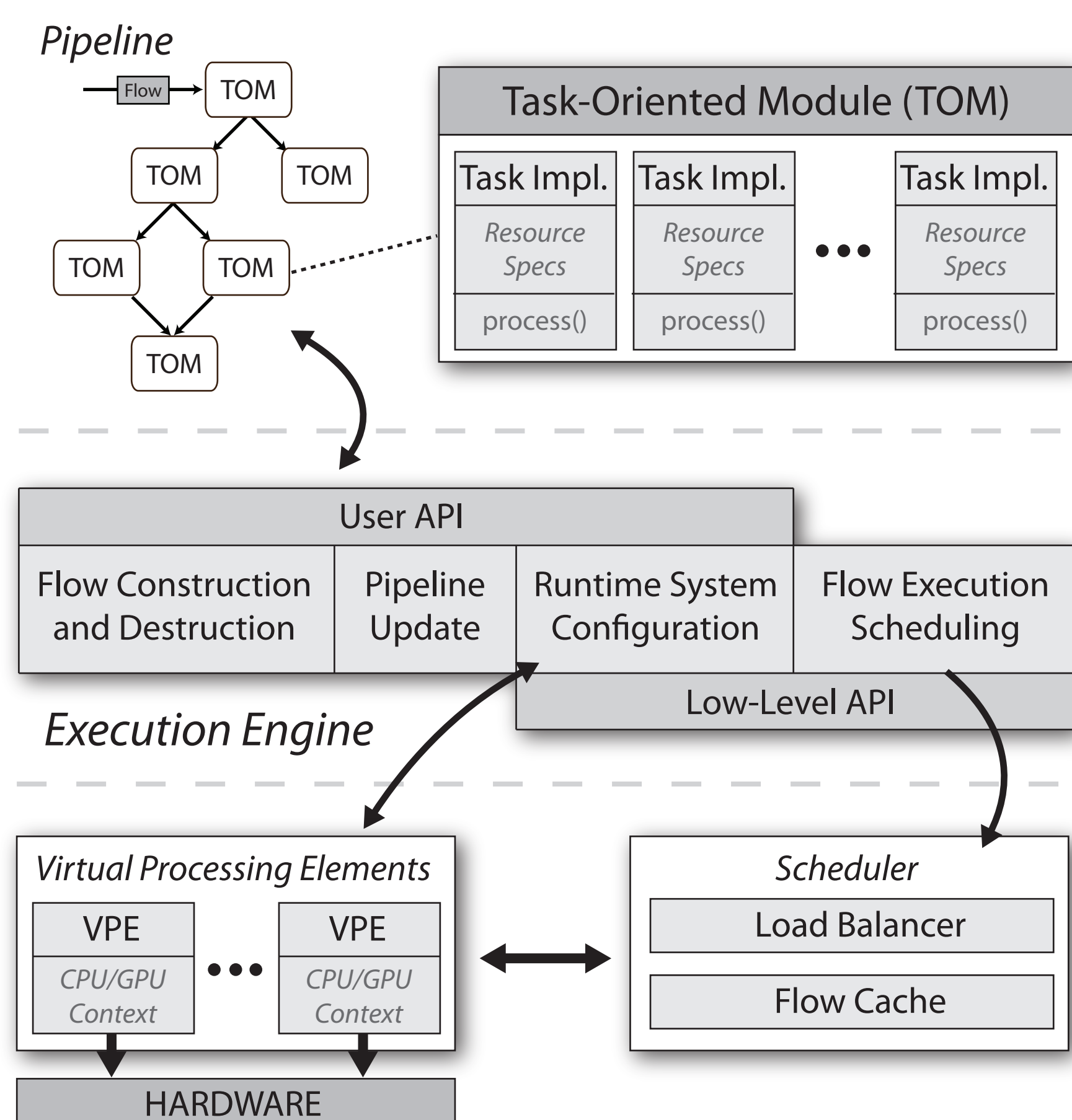
Huy Vo[1], Daniel Osmari[2], Luiz Scheidegger[2], João Comba[2], Jason Shepherd[3] and Cláudio Silva[1]

[1] University of Utah, USA    [2] UFRGS, Brazil    [3] Sandia National Laboratories, USA
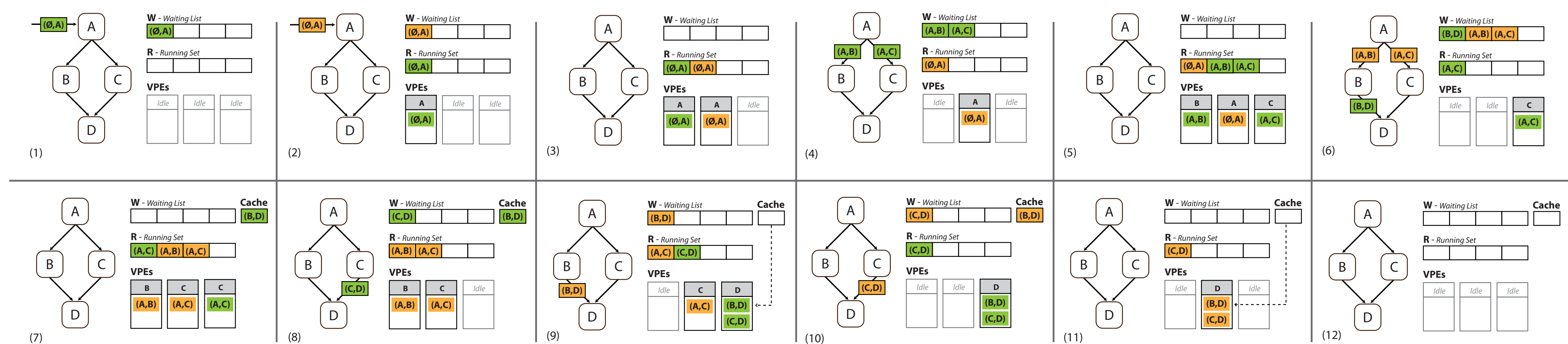
## WHY HYPERFLOW?

• Dataflow is widely used in scientific computing
• Current dataflow systems lack parallelism support:
  - Mostly designed for single-core machines
  - Some only support data-parallelism with MPI
  - Ad-hoc for multi-core and (multi-)GPU platforms
• There exists parallel frameworks such as TBB, StarPU Stream SDK, etc. but difficult to apply to dataflow due to being tightly coupled with task-graph execution.
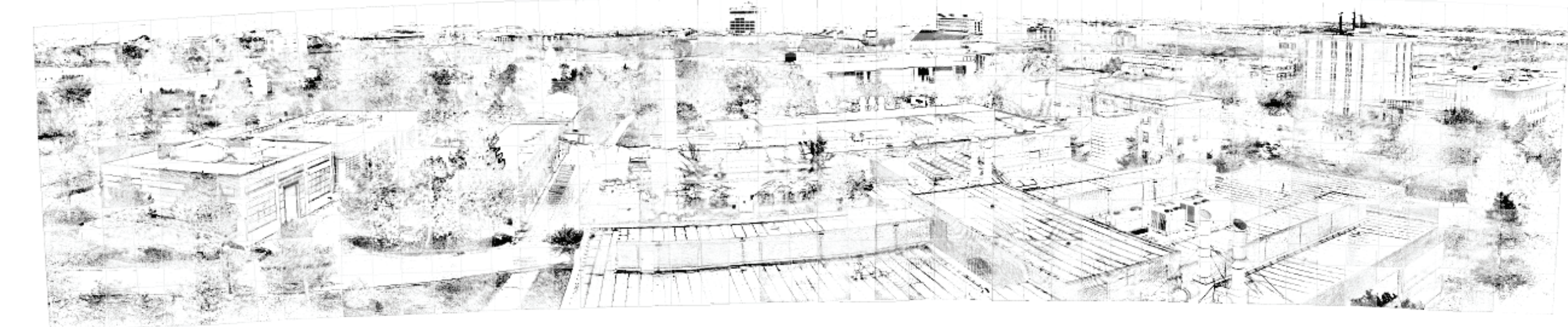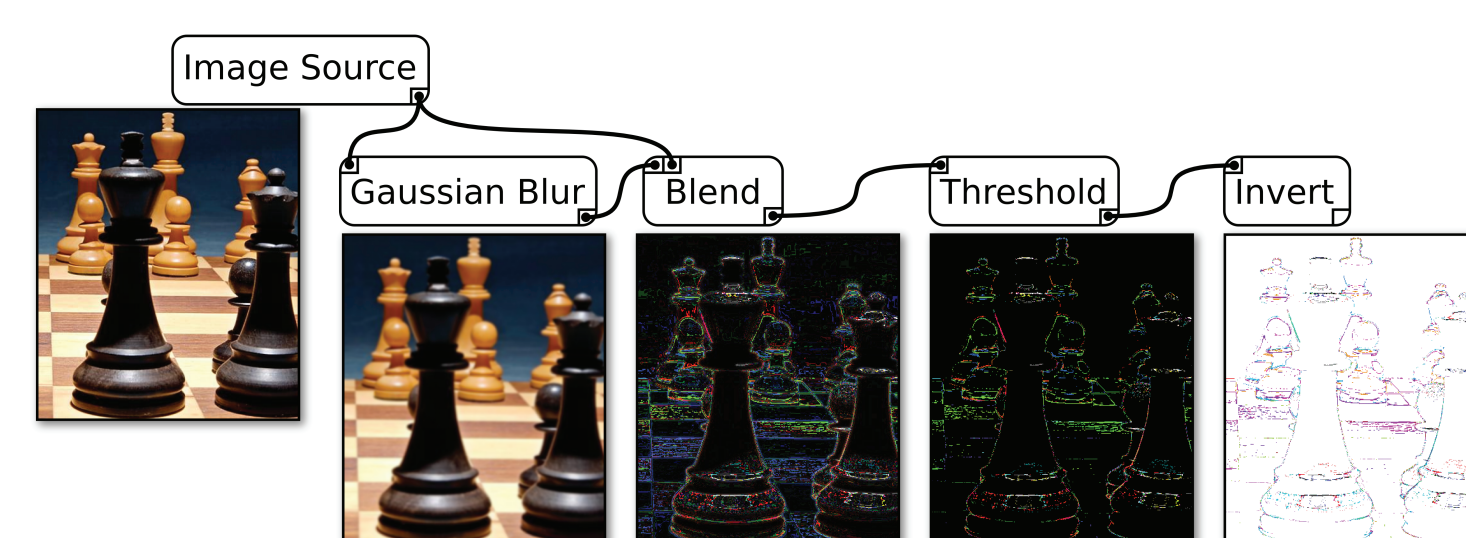
## SYSTEM ARCHITECTURE



## TASK-ORIENTED MODULE (TOM)

• Hold meta-data for its task and parameters
• Manage multiple implementations for a specific task
• Each implementation specifies its own resources
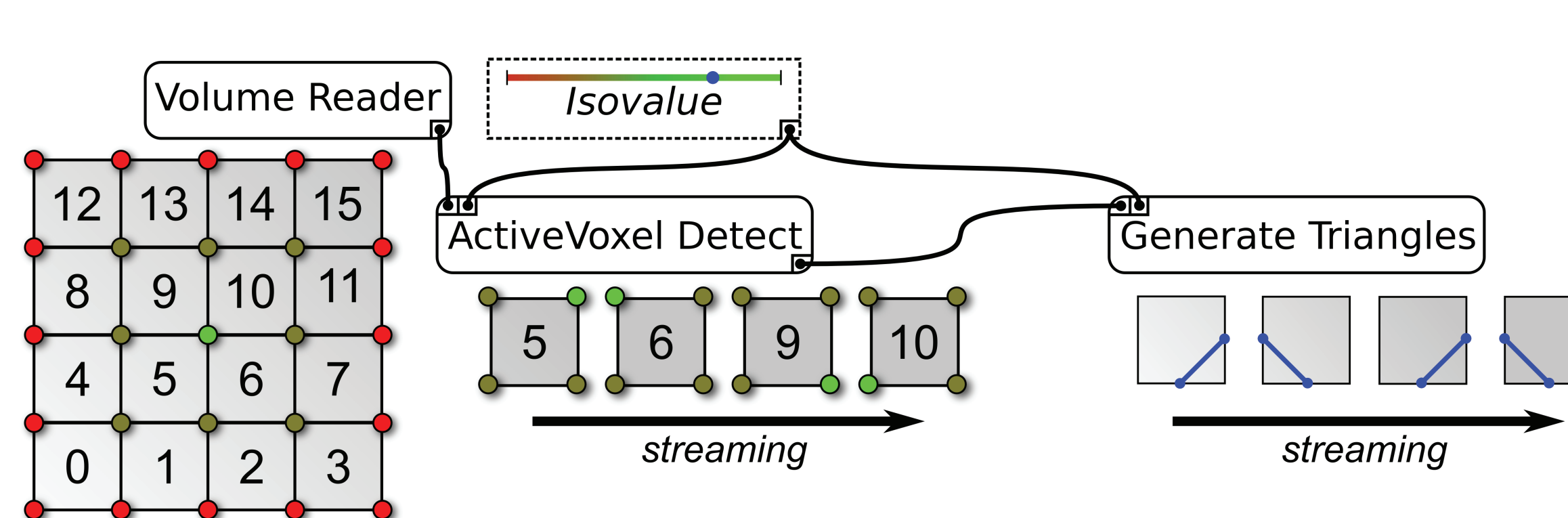
```
1 class TaskOrientedModule {
2 public:
3   TaskOrientedModule(int numberOfInputPorts,
4                      int numberOfOutputPorts);
5   void addImplementation(TaskImplementation* impl);
6   static void connect(
7       TaskOrientedModule* srcModule, int srcPort,
8       TaskOrientedModule* dstModule, int dstPort
9   );
10  // Additional parameters if subclassed
11  ...
12 };
13
14 // Construction of the pipeline
15 TaskOrientedModule Source(0, 1, "Image Reader")
16 Source.addImplementation(new SourceImpl());
17
18 TaskOrientedModule Filter(1, 1, "Gaussian Blur")
19 Filter.addImplementation(new FilterImpl());
20
21 TaskOrientedModule Sink(1, 0, "Viewer")
22 Sink.addImplementation(new SinkImpl());
23
24 TaskOrientedModule::connect(&Source,0,&Filter,0);
25 TaskOrientedModule::connect(&Filter,0,&Sink,0);
```
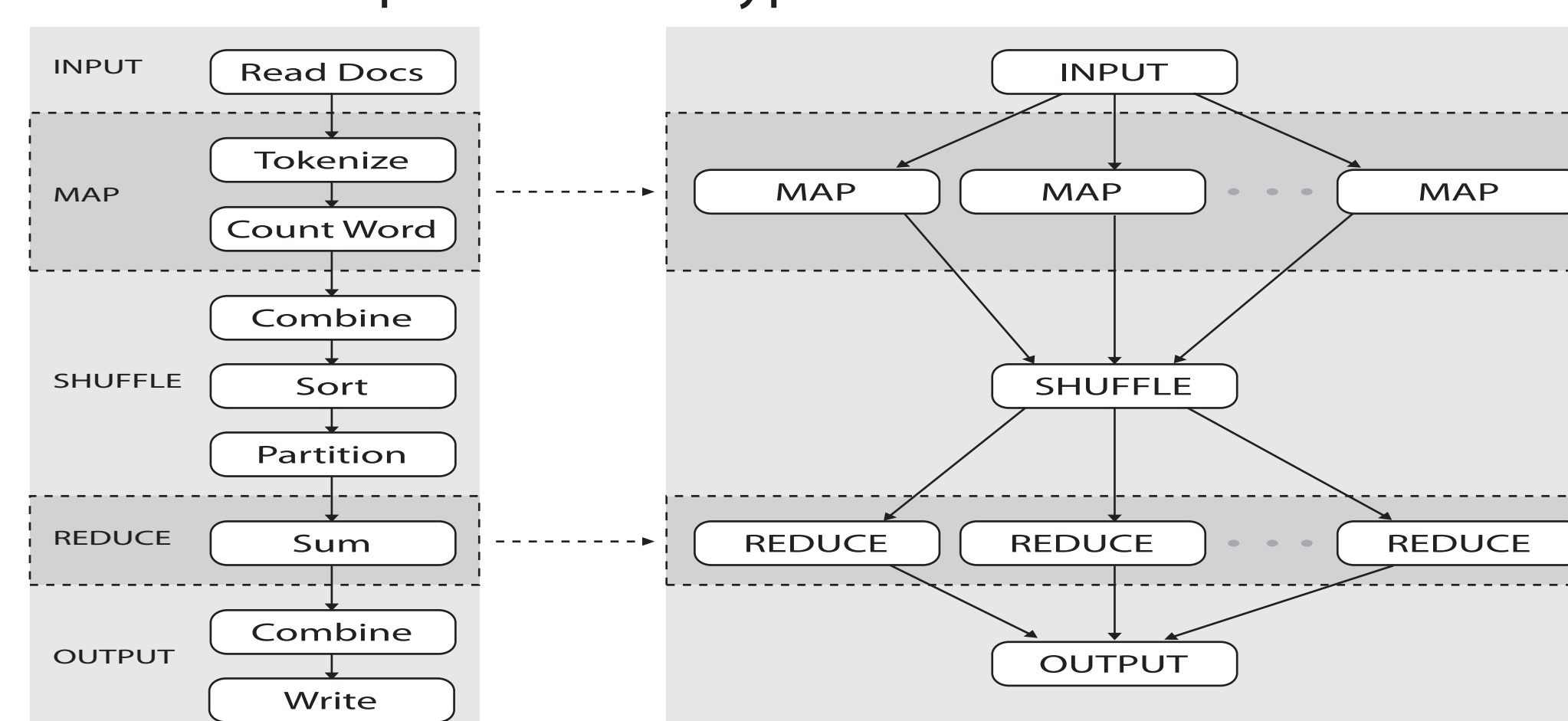
## MULTI-MEDIUM DATA INTERFACE

• Data explicitly denotes which platform to reside in
• Data transferring is triggered automatically but conversions are done by user-define functions

```
1 class DImageData : public Data {
2 public:
3   DImageData(BaseImage *img, DMTYPE medium) {
4     // Store image with the medium type
5     // ...
6   }
7   virtual Data* createInstance(DMTYPE targetMedium) {
8     // Convert from CPU Image to CUDA Image
9     if (this->medium==DM_CPU_MEMORY &&
10        targetMedium==DM_GPU_MEMORY) {
11      CUDAImage *newImage = new CUDAImage();
12      cudaMemcpy(newImage, this->image,
13               cudaMemcpyHostToDevice);
13      return new DImageData(newImage, DM_GPU_MEMORY);
14    }
15    // Convert from CUDA Image to CPU Image
16    else if (this->medium==DM_GPU_MEMORY &&
17            targetMedium==DM_CPU_MEMORY) {
18      CPUImage *newImage = new CPUImage();
19      cudaMemcpy(newImage, this->image,
               cudaMemcpyDeviceToHost);
20      return new DImageData(newImage, DM_CPU_MEMORY);
21    }
22    // Use default transfer methods for CPU-CPU, etc.
23    else
24      return Data::createInstance(medium);
25  }
```
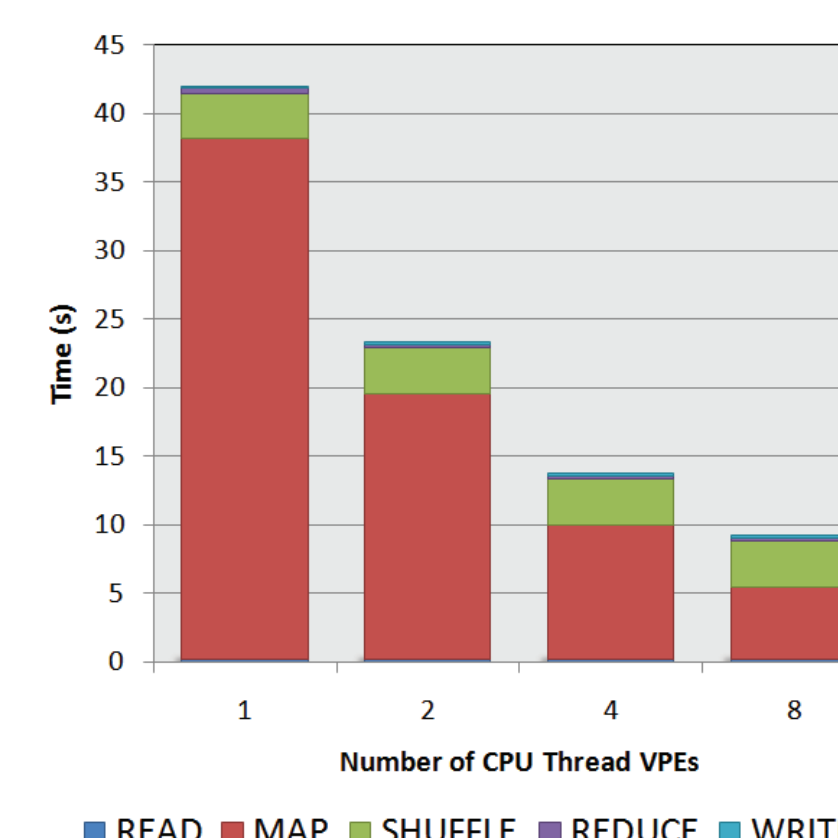
## PIPELINE EXECUTION



## SHOW CASE 1 - HIGH-THROUGHPUT IMAGE PROCESSING

Edge Detection Pipeline



Machine Specs:
(1) 2 x Xeon w5580, 8-core
    24GB RAM
    Quadro 5800
(2) Intel 970i, 4-core
    6GB RAM
    GTX 295, Tesla C1060

Timing Results
Input: 312 images ~ 6GB

| #CPU Cores | #GPU Cores | Time (s) |
|---|---|---|
| 1 | 0 | 3261 |
| 2 | 0 | 1747 |
| 4 | 0 | 1049 |
| 8 | 0 | 668 |
| 8 | 1 | 440 |
| 8 | 2 | 272 |
| 4 | 3 | 151 |

## SHOW CASE 2 - HYBRID ISOSURFACE EXTRACTION

The pipeline starts with a dataset reader that loads blocks of the volume and stream them to a TOM classifier. The classifier module traverses each block to generate a list of active voxels, i.e., only voxels that may contain the desired isosurface (voxels that do not contain v can be safely discarded). The list of active voxels is progressively streamed to a triangle generation TOM, which is designed to inspect each voxel and, based on a standard Marching Cubes case table, construct a set of triangles that approximates the desired isosurface inside that voxel. This set of triangles can also be streamed to other modules for further post-processing.

Using HyperFlow, it becomes very easy to optimize the usage of CPU and GPU computational resources. We implemented the dataset reader to divide the input volume into blocks. The execution scheduler automatically instantiates separate, concurrent TOMs to perform voxel classification on each block. Since there is no data dependency between different instances, this parallelism greatly improves algorithm performance. The TOMs designed to construct the actual isosurface can also be trivially parallelized by HyperFlow.

Isosurface Extraction Pipeline



Timing Results
Input: 1 billion voxels

| #CPU Cores | #GPU Cores | Time (s) |
|---|---|---|
| 1 | 1 | 132 |
| 1 | 2 | 101 |
| 2 | 1 | 125 |
| 2 | 2 | 72 |
| 4 | 1 | 127 |
| 4 | 2 | 71 |
| 4 | 3 | 49 |
| 8 | 1 | 113 |
| 8 | 2 | 69 |

## SHOW CASE 3 - MAP/REDUCE WORD COUNT

Map/Reduce in HyperFlow



Timing Results
Input: 83 million words

Two execution scenarios need to be addressed: (1) an upstream module, which executes sequentially, but generates output for concurrent execution downstream and (2) a downstream module which can only be executed sequentially and has to collect flows from upstream modules executed in parallel. Scenario (1) occurs when data are moved from the reader to mappers, as well as from the shuffling phase to reducers. This is necessary to enable data parallelism in the Map and Reduce phases. Scenario (2) happens when the execution comes back from the Map and Reduce to the Shuffling and Writer. HyperFlow supports both scenarios: (1) an upstream can generate multiple flows with unique identifications to trigger data-parallelism in the pipeline; (2) TOM can limit task implementations to run sequentially by using a global lock. These are illustrated in the source code on the right.

```
1 class IReader: public TaskImplementation {
2   virtual int process(ExecutionContext *ctx,
3                       TaskData *input, TaskData *output) {
4     Data *data[N];
5     // Read input file and place them into data[]
6     ...
7     // Then construct flows downstream
8     for (int i=0; i<N; i++) {
9       // We are not passing any reference flow so that
10      // newly created flows will have unique id
11      Flow *flow = ctx->createFlowForOutputPort(0);
12      flow->data = data[i];
13      output->addFlow(flow);
14    }
15    return 0;
16  }
17 };
18 class MCombine: public TaskOrientedModule {
19 public:
20   int        numberOfMappers;
21   SafeCounter doneCount;
22   Mutex       mutex;
23   ListOfData  data;
24 };
25
26 class ICombine: public TaskImplementation {
27   virtual int process(ExecutionContext *ctx,
28                       TaskData *input, TaskData *output) {
29     MCombiner *M = (MCombiner*)ctx->module;
30     // Locking resource to accumulate
31     M->mutex.lock();
32     M->data.push_back(input->flows[0]->data);
33     M->mutex.unlock();
34
35     if (M->doneCounter->increase() < M->numberOfMappers)
36       return -1;
37
38     // Process data in module->data
39     ...
40     return 0;
41  }
42 };
```