

An Autonomous Execution Model for GPUs: When CPUs Take a Back Seat


Assoc. Prof. Didem Unat

dunat@ku.edu.tr

Koç University, Istanbul, Türkiye

- **Data path** has moved to the GPU
- **Control path** still on the CPU
- **CPU involvement** adds latencies

```
// ① Time loop on CPU
while (iter++ < num_iterations) {
  // ② Launch compute kernel
  stencil_kernel<<<..., comp_stream>>>(...)
  // ③ Sync with neighboring GPUs
  wait_neighbors(comm_stream)
  // ④ Compute and communicate boundaries
  compute_boundaries<<<..., comm_stream>>>(...)
  write_neighbors(comm_stream)
  // ⑤ Sync comm and comp streams
  sync(comm_stream, comp_stream)
}
```



What if we move everything to the GPU?

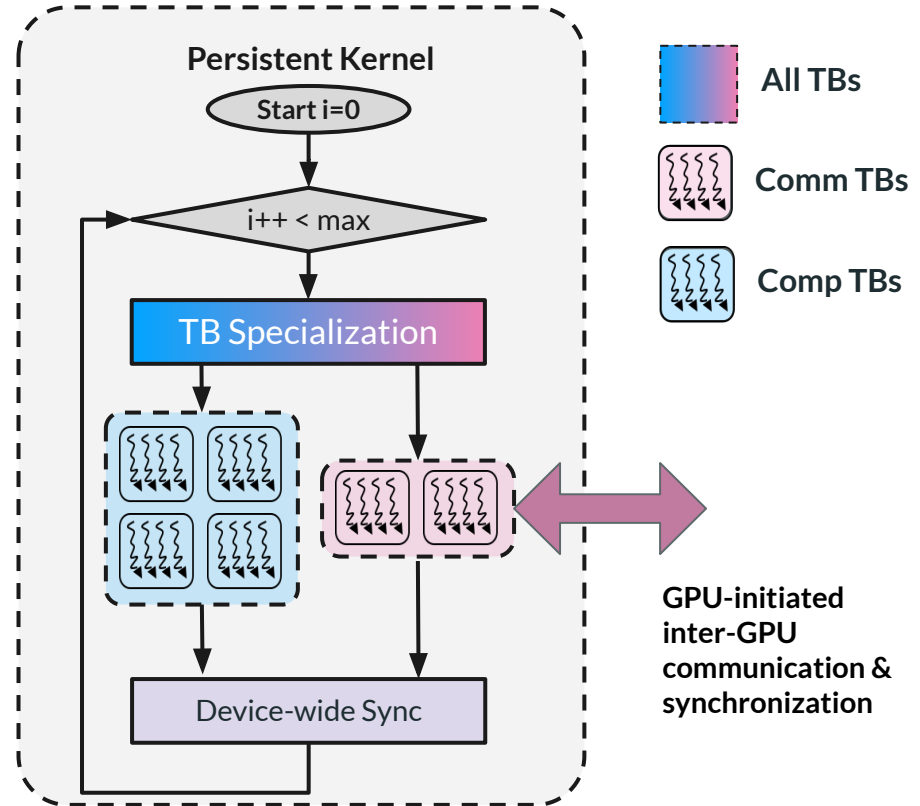
- Move the time loop to the GPU
- Initiate communication from the GPU
- Explicitly overlap on the GPU
- Synchronize within and across devices directly from the GPU

```
__global__ void CPU_Free_Jacobi(...) {  
    // Time loop on GPU  
    while (iter++ < num_iterations) {  
        //a Compute boundary using top neighbor's values  
        if (TB_index == 0) {  
            //① Wait for top neighbor to signal  
            wait_top_neighbor(...)  
            //② Compute top boundary using halos  
            top_boundary = compute(top_halo, south, ...)  
            //③ Write to top neighbor's bottom halo  
            write_top_neighbor(top_boundary)  
            //④ Signal top neighbor that iteration is done  
            signal_top_neighbor(...)  
        }  
        //b Compute boundary with bottom neighbor's values  
        if (TB_index == 1) { ... }  
        //c Remaining TBs compute the inner part  
        if (TB_index == <rest>) { ... }  
        //⑤ Synchronize all  
        grid.sync()  
    }  
}
```

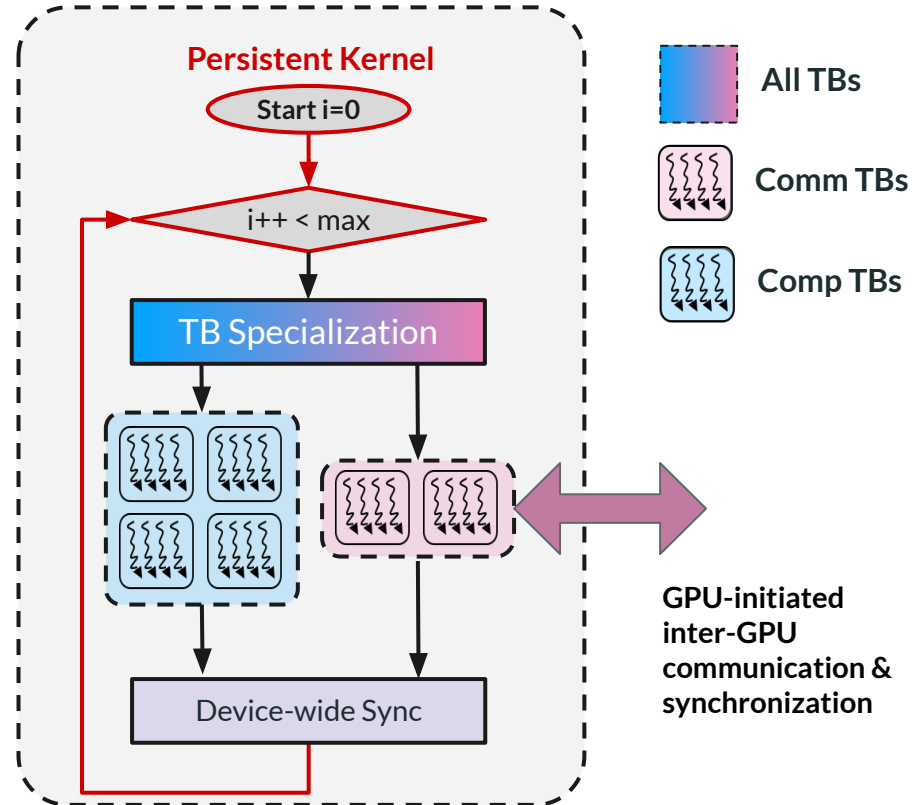
CPU-Free

Allow the GPU to take the reins of the data and *control* paths!

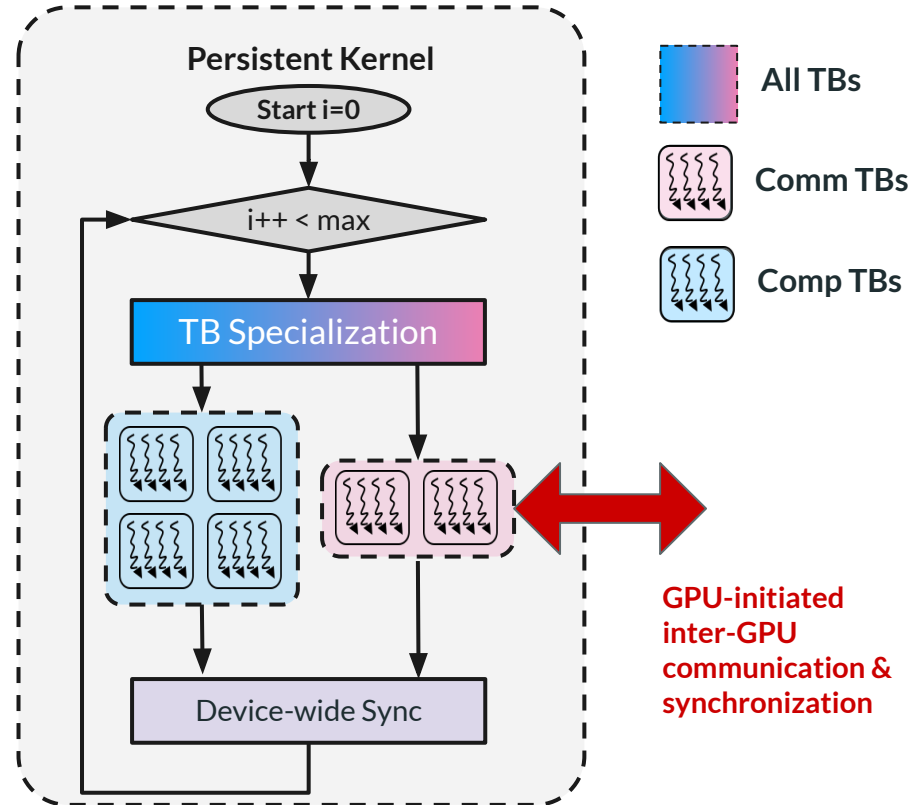
- CPU-Free model grants complete autonomy to the GPU
- Both data and control paths move to device-side
- **Components**
 - Persistent kernels
 - GPU-initiated data movement
 - TB specialization
 - Device-side synchronization



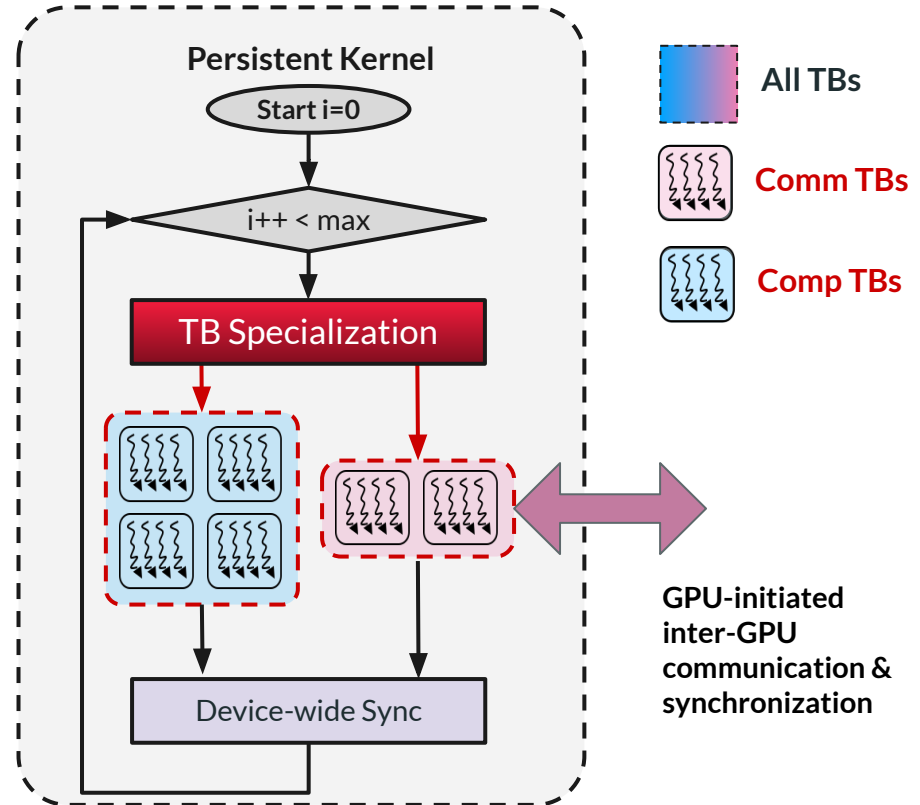
- Move time loop to GPU using a persistent kernel
- *Discrete* kernel - Launched repeatedly
- *Persistent* kernel - Launched once
- Long running persistent kernel grants GPU more autonomy



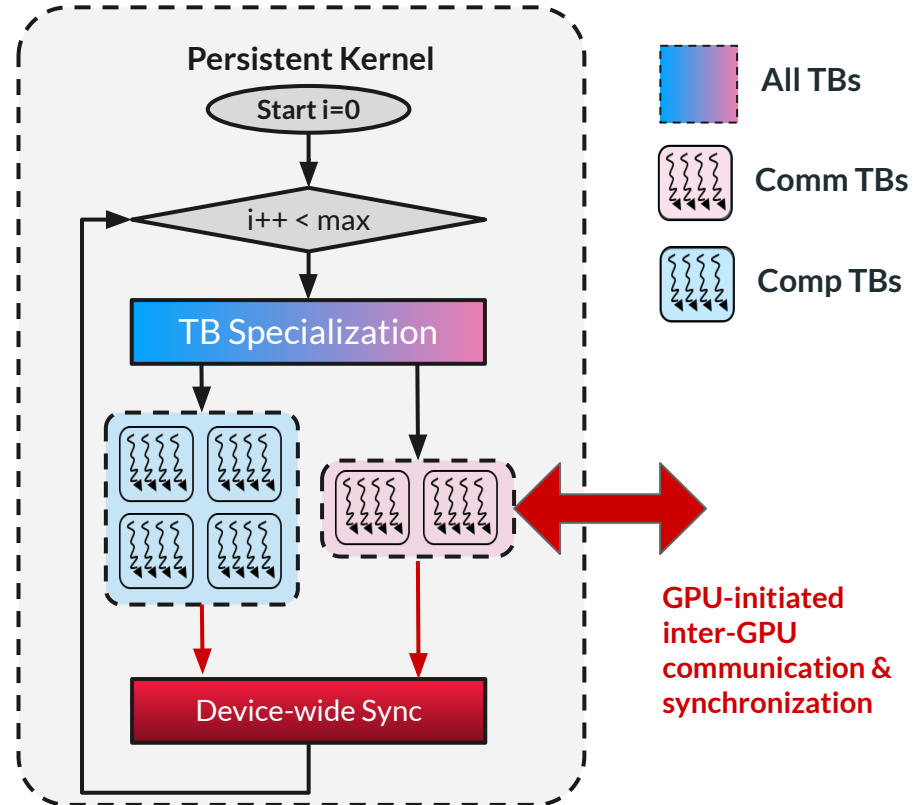
- Initiate communication directly from GPU
- Communication:
 - Data transmission
 - Multi-GPU synchronization
- NVSHMEM used as communication mechanism
- NVSHMEM provides efficient GPU-side communication primitives



- Specialize few thread blocks to handle communication
- Remaining thread blocks compute
- Traditional overlap uses concurrent GPU streams
- Streams are synced through GPU events



- *Within GPU* - Device-side barriers
- *Across GPUs* - Device-initiated signal / flag mechanisms
- CUDA Cooperative Groups API
- NVSHMEM signal operations / barriers for multi-GPU sync



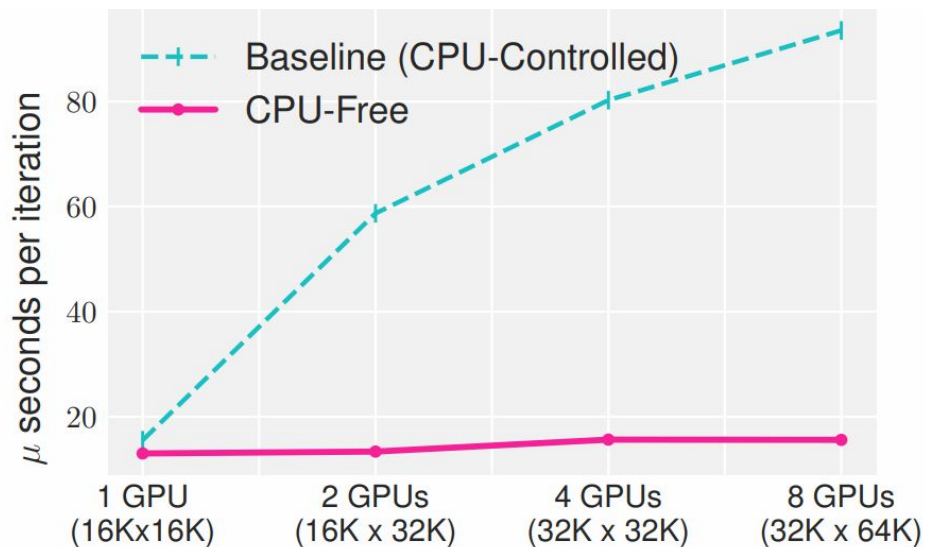
CPU-controlled execution needs multiple API calls

- Kernel launches
- Communication calls
- Multiple streams
- GPU events
- Global barriers (OpenMP, MPI barriers)

One fused kernel eliminates these overheads

Reduces Communication Overheads and Overlaps Better

- Can initiate communication as soon as the data is ready
- Can inline communication with computation
- When device not saturated, API call latencies dominate
- CPU-Free can overlap in small domains
- Especially relevant in *strong scaling*
- More asynchrony



(b) Communication overhead with no computation

- Shared memory has lifetime of kernel
- CPU-Free execution can reuse shared memory across iterations
- We integrate PERKS kernel into our model

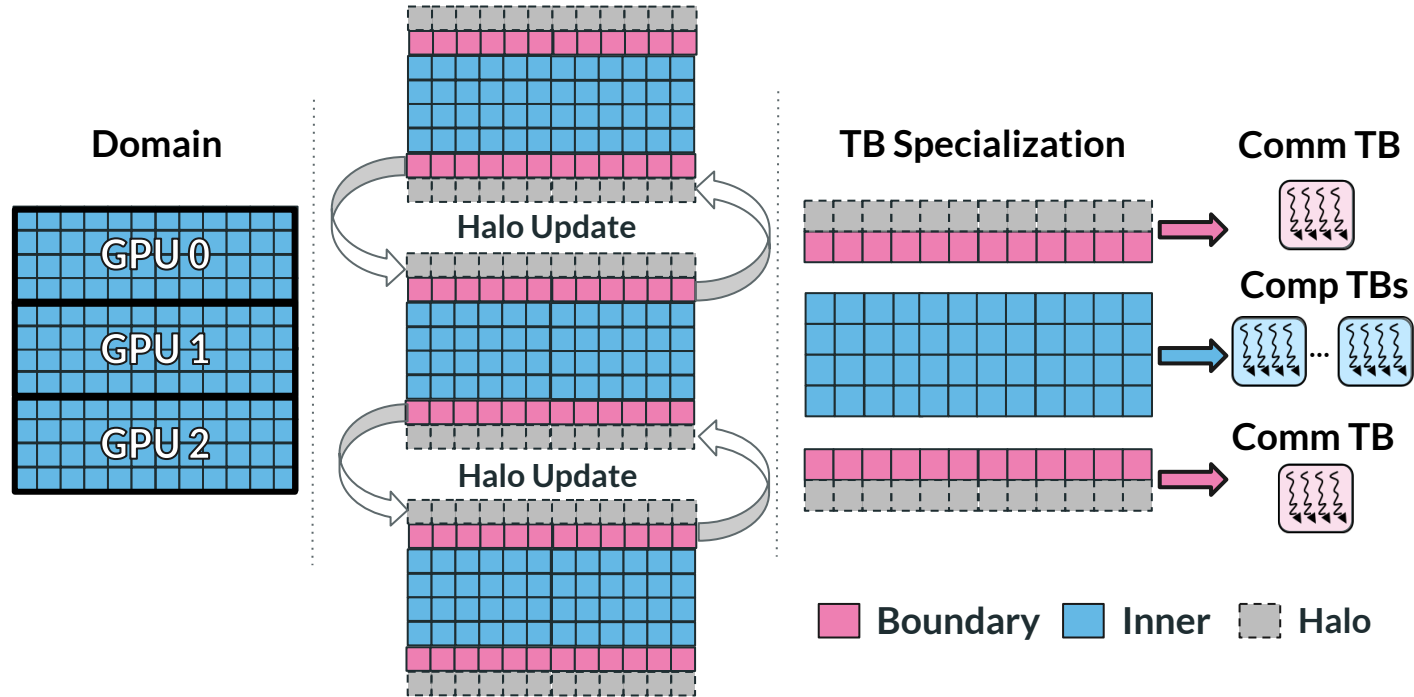
Use-Cases

- Jacobi 2D / 3D
- Conjugate Gradient



*Published at Multi-GPU Communication Schemes for Iterative Solvers: When CPUs are Not in Charge.
In Proceedings of the 37th International Conference on Supercomputing (ICS '23).*

Use-Case 1 - Jacobi Stencil




procedure STANDARD(**A**, **b**, **x**₀)


$$\mathbf{r}_0 = \mathbf{b} - \mathbf{A}\mathbf{x}_0$$


$$\mathbf{p}_0 = \mathbf{r}_0$$


$$\gamma_0 = \mathbf{r}_0 \odot \mathbf{r}_0$$


for $k = 0, 1, \dots$ **do**


 $\mathbf{s} = \mathbf{A}\mathbf{p}_k$

 $\delta_k = \mathbf{p}_k \odot \mathbf{s}$ **and** $\alpha_k = \gamma_k / \delta_k$

 $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k$

 $\mathbf{r}_{k+1} = \mathbf{r}_k - \alpha_k \mathbf{s}$

 $\gamma_{k+1} = \mathbf{r}_{k+1} \odot \mathbf{r}_{k+1}$ **and** $\beta_k = \gamma_{k+1} / \gamma_k$

 $\mathbf{p}_{k+1} = \mathbf{r}_{k+1} + \beta_k \mathbf{p}_k$

end for

end procedure

Saxpy - Local computations; no communication

SpMV - Need entries on other GPUs; need communication

Dot - Global sync point because of global reduction

Each step depends on previous

No possible overlap

```
procedure PIPELINED(A, b, x0)
```

$$\mathbf{r}_0 = \mathbf{b} - \mathbf{A}\mathbf{x}_0$$

$$\mathbf{w}_0 = \mathbf{A}\mathbf{r}_0$$

```
for  $k = 0, 1, \dots$  do
```

$$\gamma_k = \mathbf{r}_k \odot \mathbf{r}_k$$

$$\delta_k = \mathbf{w}_k \odot \mathbf{r}_k$$

$$\mathbf{q}_k = \mathbf{A}\mathbf{w}_k$$

```
if  $k > 0$  then
```

$$\beta_k = \gamma_k / \gamma_{k-1}$$

$$\alpha_k = \gamma_k / (\delta_k - (\beta_k \gamma_k) / \alpha_{k-1})$$

```
else
```

$$\beta_k = 0$$

$$\alpha_k = \frac{\gamma_k}{\delta_k}$$

```
end if
```

$$\mathbf{z}_k = \mathbf{q}_k + \beta_k \mathbf{z}_{k-1}$$

$$\mathbf{s}_k = \mathbf{w}_k + \beta_k \mathbf{s}_{k-1}$$

$$\mathbf{p}_k = \mathbf{r}_k + \beta_k \mathbf{p}_{k-1}$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k$$

$$\mathbf{r}_{k+1} = \mathbf{r}_k - \alpha_k \mathbf{s}_k$$

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \mathbf{z}_k$$

```
end for
```

```
end procedure
```

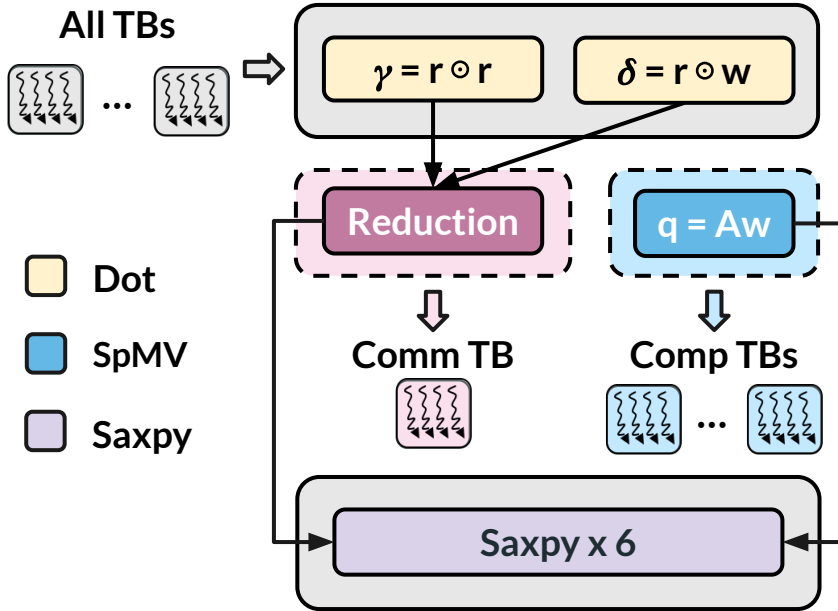
Introduce auxiliary vectors to allow overlap

Can implement dot product global reductions with one reduction

Can overlap dot product global reductions with SpMV

- Pipelined CG can showcase overlap
- We implement CPU-Free for both Standard and Pipelined CG

Use-Case 2 - Conjugate Gradient



```
__global__ void CPU_Free_PipelinedCG(...) {  
    // ① Time loop on GPU  
    while (iter++ < num_iterations) {  
        local_dot(r, r, gamma)  
        local_dot(r, w, delta)  
  
        ...  
        // ② Specialize one TB for communication  
        if (TB_index == 0) {  
            // ③ Multi-GPU reduction  
            sum_reduce(gamma, delta)  
        } else {  
            SpMV(A, w, q)  
        }  
  
        // ④ Sync within device  
        grid.sync()  
  
        ...  
        saxpy(...) //x6  
        // ⑤ Sync across devices  
        multi_gpu.sync()  
    }  
}
```

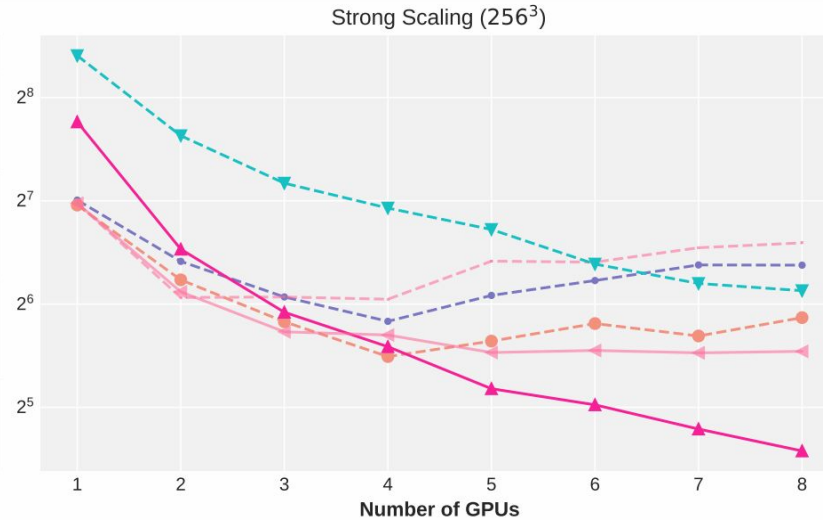
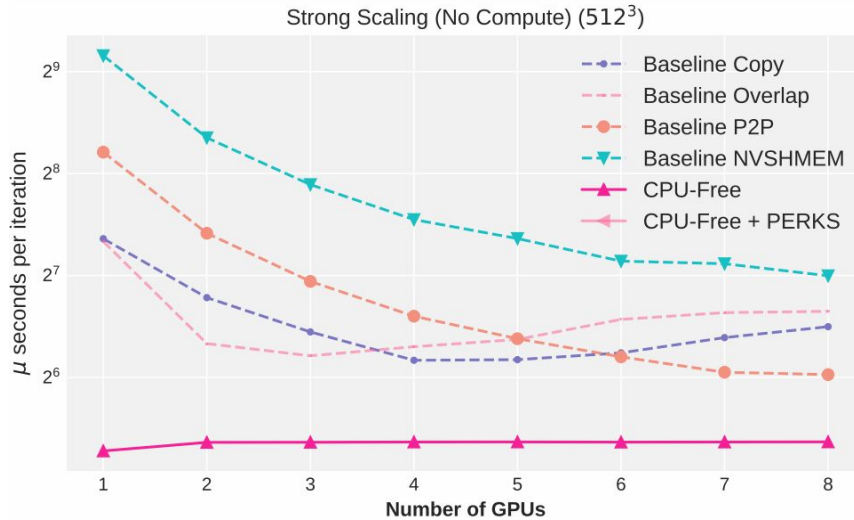
Evaluation

- Jacobi 2D/3D
- Conjugate Gradient

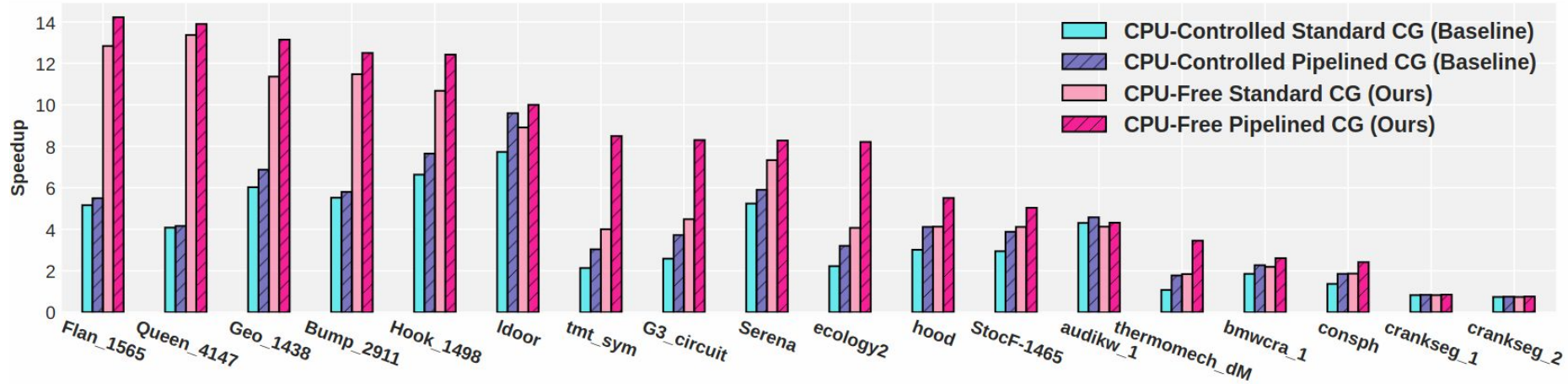


8 NVLink all-to-all connected NVIDIA A100 GPUs

*Published at Multi-GPU Communication Schemes for Iterative Solvers: When CPUs are Not in Charge.
In Proceedings of the 37th International Conference on Supercomputing (ICS '23).*

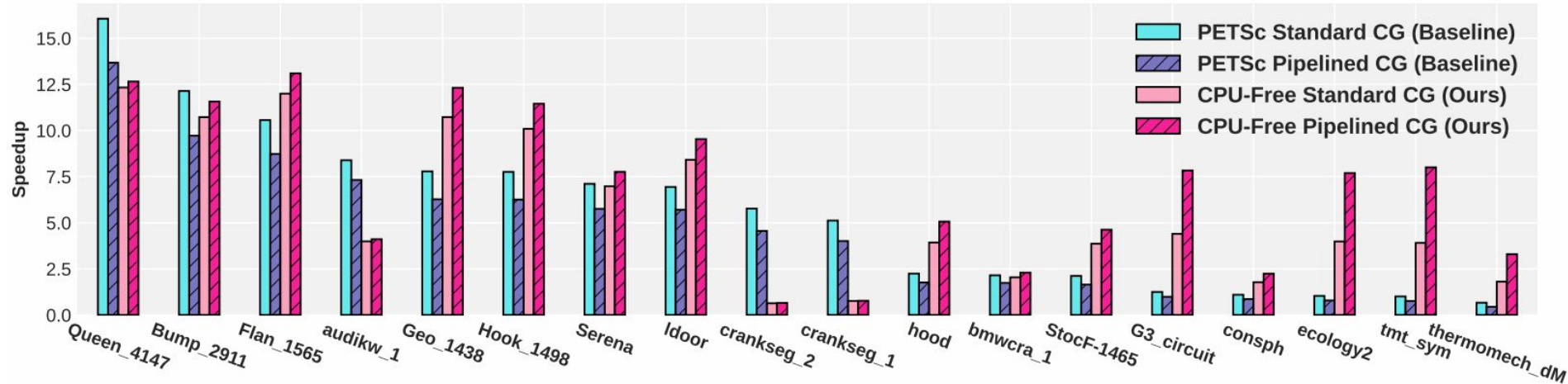


- Consistently lower communication overheads
- Excels in strong scaling scenarios
 - Underperforms at small numbers of GPUs (compute-bound)
 - As GPU count increases, overheads start to dominate
 - CPU-Free pulls ahead at larger GPU counts



CPU-Free achieves **1.63x** and **1.54x** geo mean speedup over CPU-Controlled for Pipelined and Standard CG, respectively

Comparison with Petsc

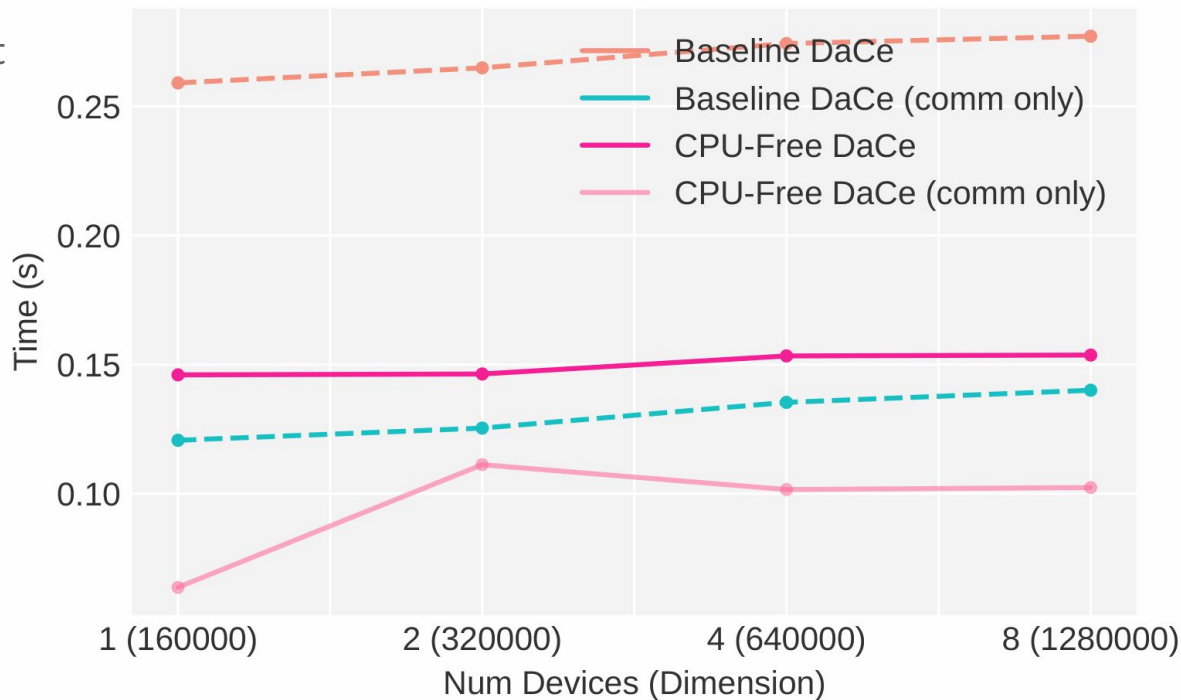


Outperforms PETSc for 13 out 18 sparse matrices

- Compiling autonomous CPU-free code from Python
- Developing profiling and monitoring tools for inter-GPU communication
- Building a runtime system that uses CudaGraphs for CPU-free execution - particularly useful for irregular applications
- Extending the CPU-free model to AMD GPUs

Extending the DACE compiler

- 44.5% performance improvement at 8 GPUs
- 26.8% improvement in communication latency
- Little overall communication, improvements attributed to synchronization overheads



- **Hard-to-scale** management overheads can be moved to GPU
- GPUs can be more autonomous
 - Managing their own communication, synchronization etc
- Less reliance on CPU
- Need to automate this process and support it with tools
- Do we need fat expensive CPUs to be attached to GPUs?
 - Modular supercomputers
 - Use light CPUs for GPUs nodes
 - Reduce cost and energy consumption

This project has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant agreement No 949587).



- Outperforms CPU-controlled baselines for both Jacobi 2D/3D and CG
- Especially suited for **communication latency-bounded and strong scaling scenarios**
 - Any application suggestions?
- Code available at <https://github.com/ParCoreLab/CPU-Free-model>

