

# Multigrid solver for elliptic problem on GPU

Using dolfinx (and other things)

Igor, Chris, and Adrian

Cambridge University and Edinburgh University

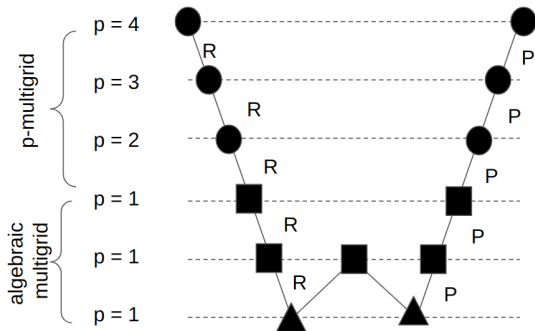
April 21, 2023





# P-multigrid

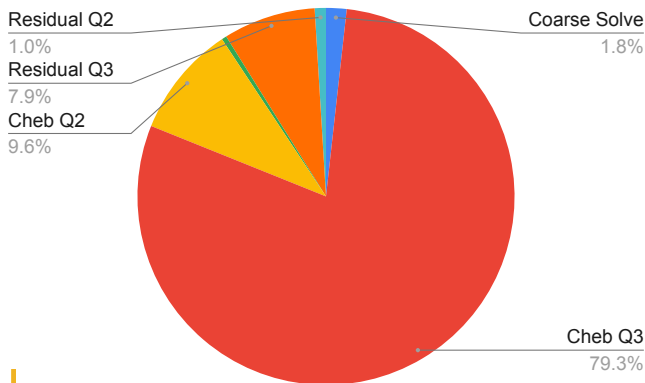
(2)



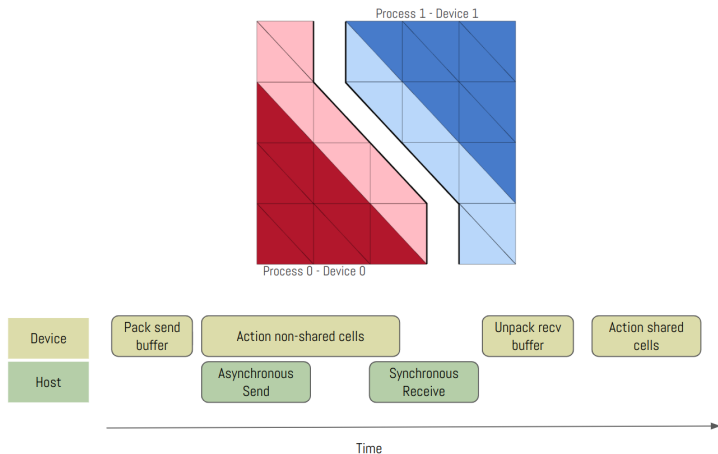
# Timing 1 GPU

## Hierarchy:

- o Level Q1: dofs 191,748 - nnz 4,999,696 – coarse level (AMG?)
- o Level Q2: dofs 1,494,425 - nnz 93,773,201
- o Level Q3: dofs 4,999,696 - nnz 61,623,0976



# Overlapping communication and computation



# Distributed MatVec I

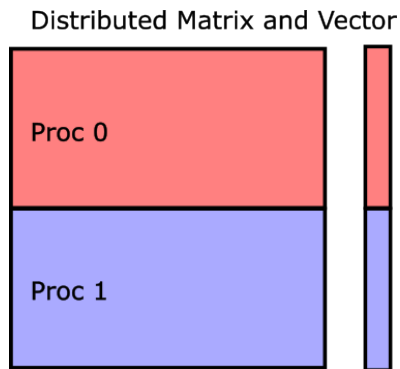
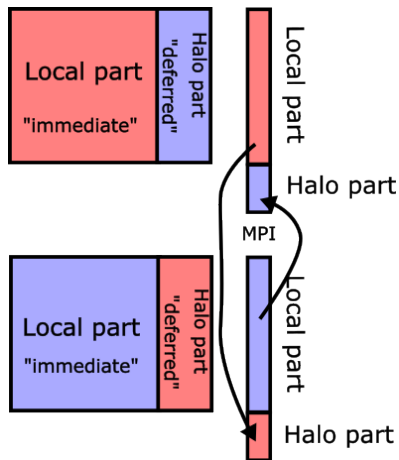
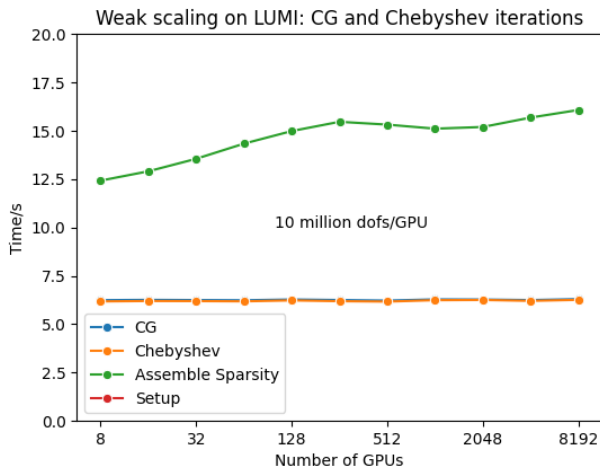


Figure: Matrix-vector Product (two processes)

# Distributed MatVec II



# Weak scaling







# rocpfrof wrapper script

```
#!/bin/bash
set -euo pipefail
name="$1"
if [[ -n ${OMPI_COMM_WORLD_RANK+z} ]]; then
    export MPI_RANK=${OMPI_COMM_WORLD_RANK}
elif [[ -n ${MV2_COMM_WORLD_RANK+z} ]]; then
    export MPI_RANK=${MV2_COMM_WORLD_RANK}
elif [[ -n ${SLURM_PROCID+z} ]]; then
    export MPI_RANK=${SLURM_PROCID}
else
    echo "Unknown MPI layer detected! Must use OpenMPI, MVAPICH, or Slurm"
    exit 1
fi
export ROCR_VISIBLE_DEVICES=${SLURM_LOCALID}
rocpfrof="/opt/rocm/bin/rocpfrof"
pid="$$_"
outdir="rank_${pid}_${MPI_RANK}"
outfile="${name}_${pid}_${MPI_RANK}.csv"
${rocpfrof} -d ${outdir} -o ${outdir}/${outfile} --hsa-trace
--hip-trace --roctx-trace "${@:2}"
```

Wrapper scripted called like this:

```
srun -N ${SLURM_NNODES} -n ${SLURM_NTASKS} ${cpu_bind} ./rocpfrof-wrapper.sh app_name app
```

Note the moving of the `ROCM_VISIBLE_DEVICES` into the wrapper script.



# rocm-smi API calls

```
rsmi_status_t err;
uint64_t total, usage;
float return_value = 0.0;
uint32_t num_devices;
MPI_Comm_rank(MPI_COMM_WORLD, &rank);
num_devices = num_monitored_devices();
if(num_devices == 1 || rank == 0){
    for(uint32_t i = 0; i < num_devices; ++i) {
        std::cout << text << " MPI Rank " << rank << " GPU " << i << " % memory used:";
        for (uint32_t mem_type = RSMI_MEM_TYPE_FIRST; mem_type <= RSMI_MEM_TYPE_LAST; ++mem_type) {
            err = rsmi_dev_memory_total_get(i, static_cast<rsmi_memory_type_t>(mem_type), &total);
            if (err != RSMI_STATUS_SUCCESS) {
                return return_value;
            }
            err = rsmi_dev_memory_usage_get(i, static_cast<rsmi_memory_type_t>(mem_type), &usage);
            if (err != RSMI_STATUS_SUCCESS) {
                return return_value;
            }
            if(mem_type == RSMI_MEM_TYPE_VRAM){
                return_value = (static_cast<float>(usage)*100)/total;
            }
            std::cout << " " << kDevMemoryTypeNameMap.at(static_cast<rsmi_memory_type_t>(mem_type));
            std::cout << " " << std::setprecision(2) << (static_cast<float>(usage)*100)/total;
        }
        std::cout << "\n";
    }
}
```

# Whole application profile



Figure: omnitrace profile of the whole application, 2 nodes, 8 GPU per node, 30m DoF per GPU. Profile from MPI rank 0

# Main solving phase

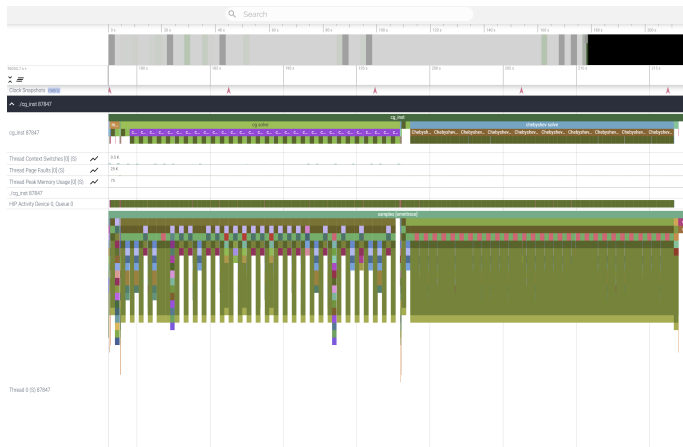


Figure: omnitrace profile of the main solving phase of the application, 2 nodes, 8 GPU per node, 30m DoF per GPU. Profile from MPI rank 0



# Chebyshev phase

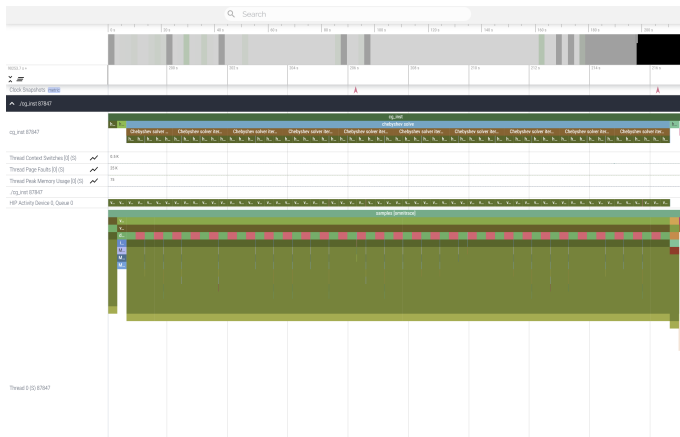


Figure: omnitrace profile of the Chebyshev solving phase





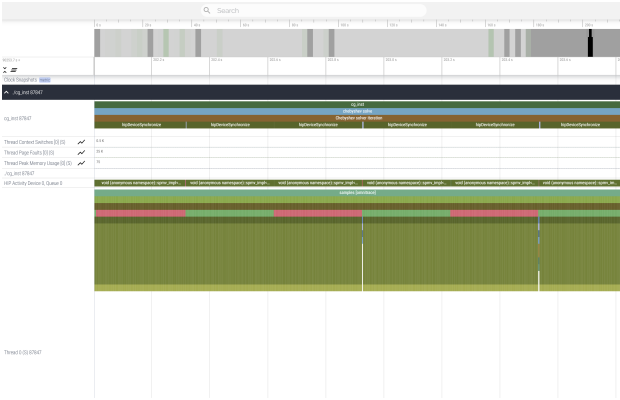








# Chebyshev iteration



# Chebyshev iteration

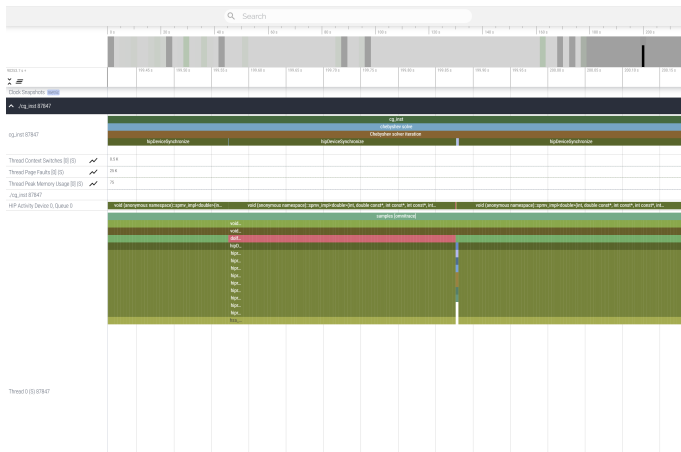


Figure: omnitrace profile of the Chebyshev iteration in more detail







