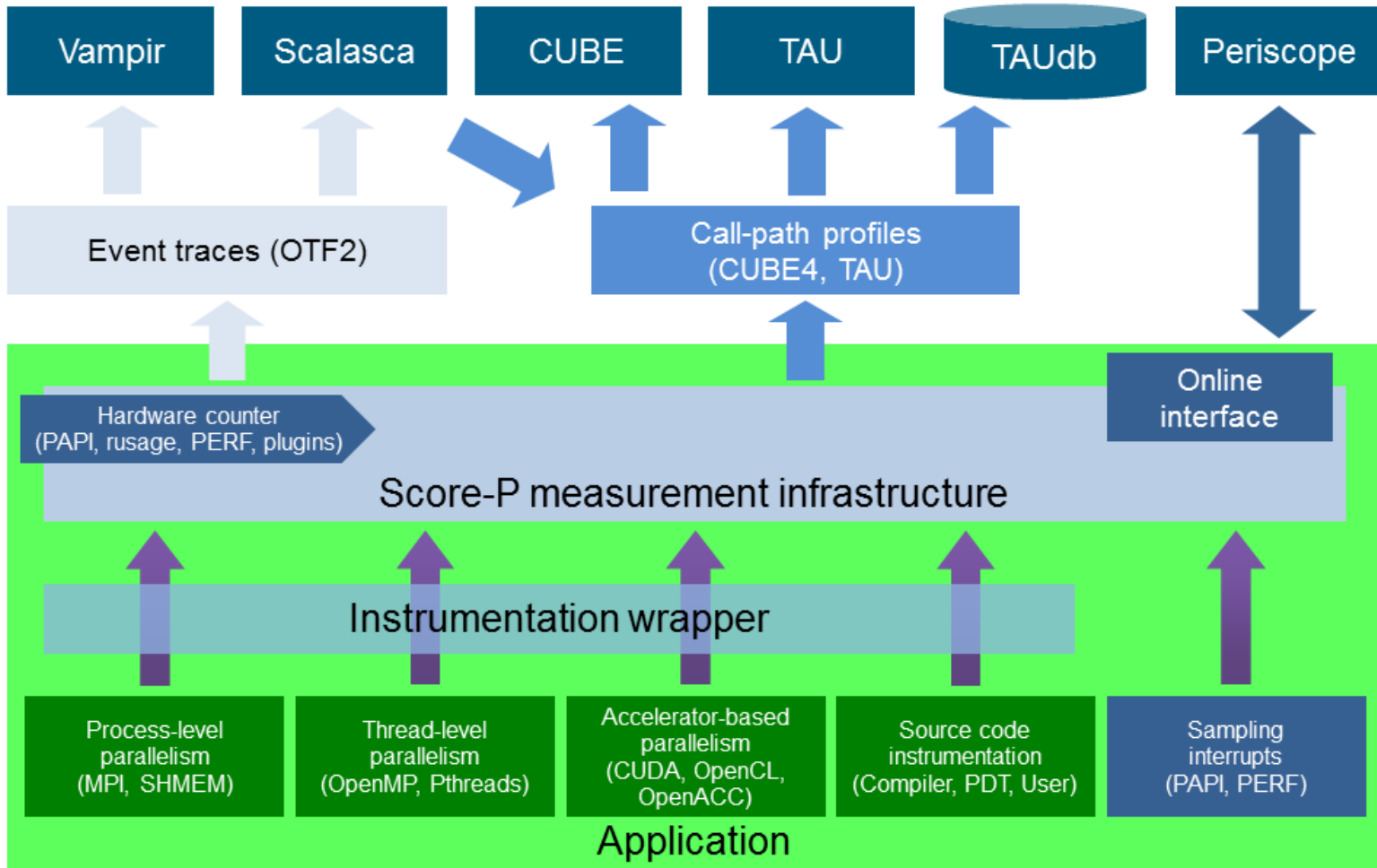


Performance Tools for MPI

PPCES 2017

Hristo Iliev
IT Center / JARA-HPC

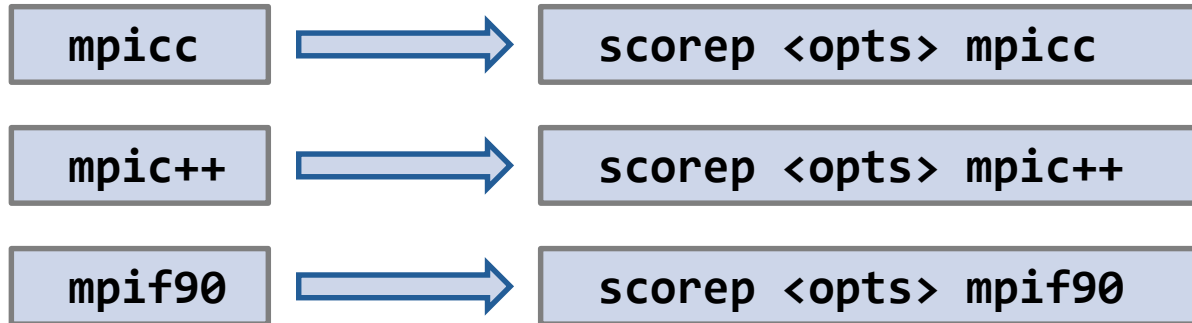
- **MPI programs do not always behave as expected**
 - Communication delays
 - Synchronisation overhead
 - Inefficient communication patterns
- **Tools exist for profiling and tracing the program execution**
 - Profiling
 - How often certain functions are called?
 - How much time is spent in certain parts?
 - How much data is exchanged in certain parts?
 - Tracing
 - A timeline of events happening during the execution



Source: [Score-P website](#) (ZIH)

■ Code first has to be instrumented accordingly:

→ Recompile with instrumentation



→ When run, the *instrumented binary* produces event data in several formats

■ Instrumentation type

- `--compiler` compiler assisted instrumentation (default)
very detailed; huge trace files
- `--user --nocompiler` manual tracing using the Score-P API
- `--mpp=mpi` traces MPI events (auto)
- `--thread=openmp` traces OpenMP parallel constructs (auto)

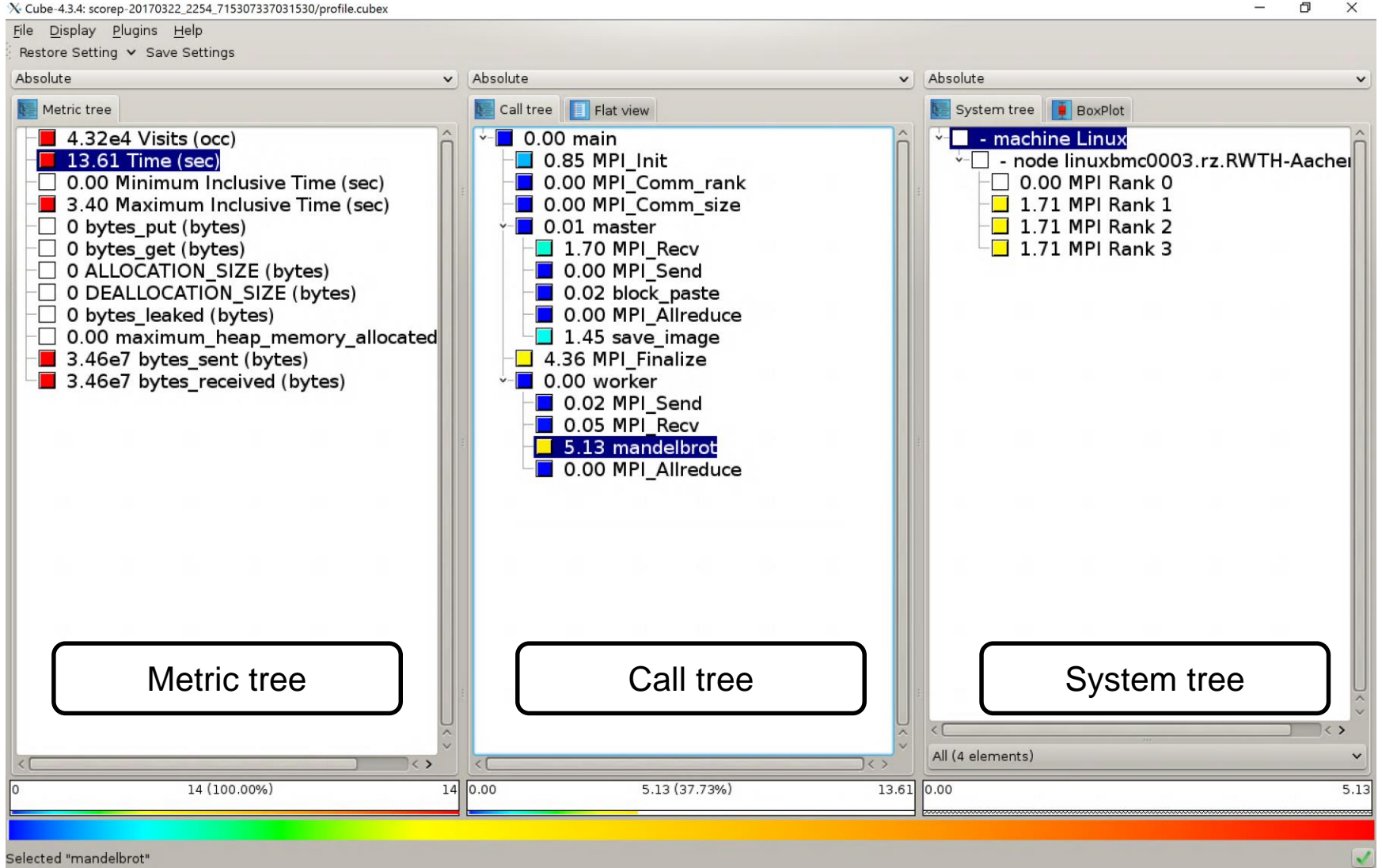
■ Score-P is controlled by environment variables

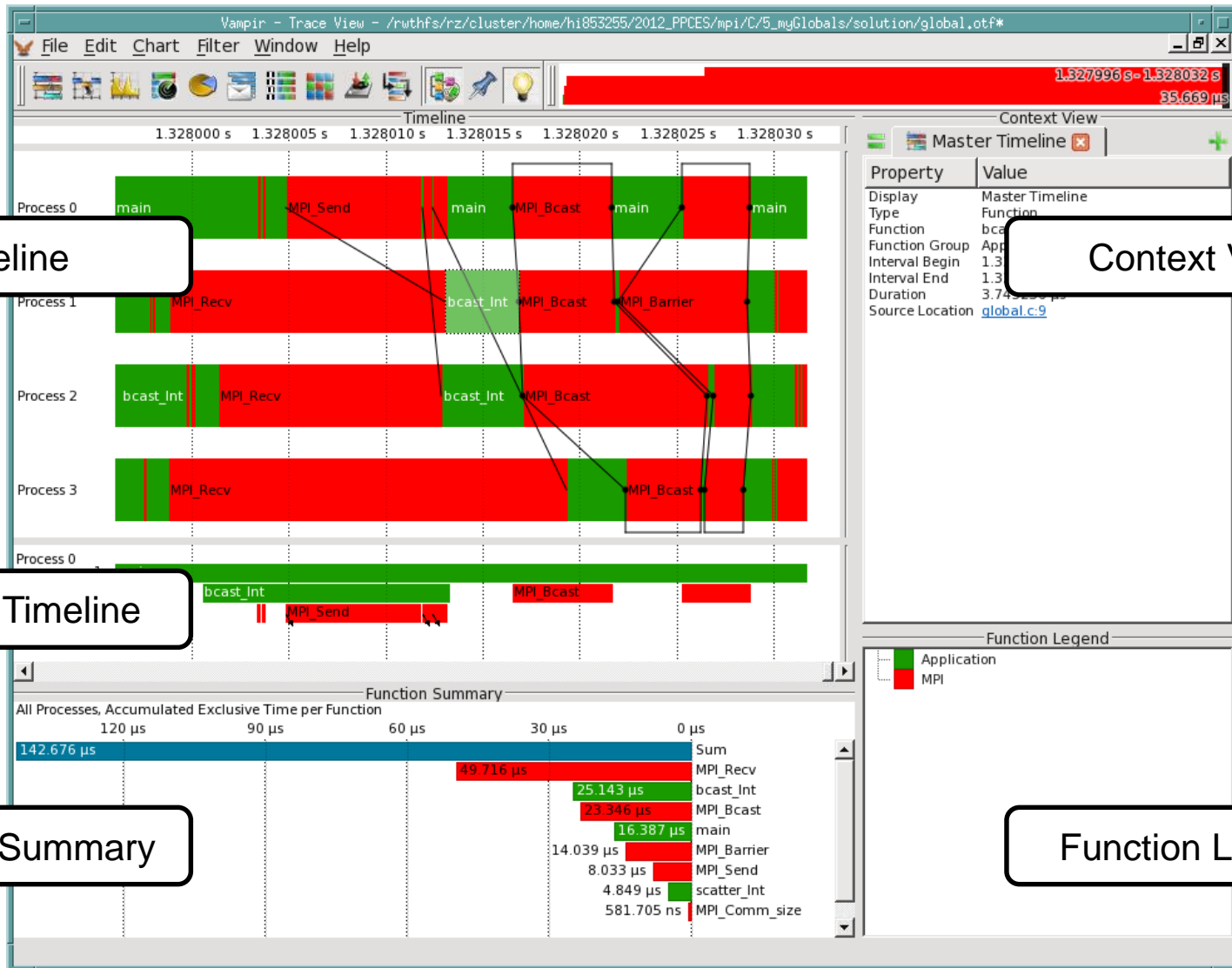
- `SCOREP_ENABLE_TRACING` – enables writing of trace files when set to **1 / true**
- `SCOREP_ENABLE_PROFILING` – enables writing of profiling information
- `SCOREP_TOTAL_MEMORY` – size of the event buffers, e.g. **200M**
- `SCOREP_EXPERIMENT_DIRECTORY` – location of the output files
- `SCOREP_FILTERING_FILE` – location of a file containing filtering rules

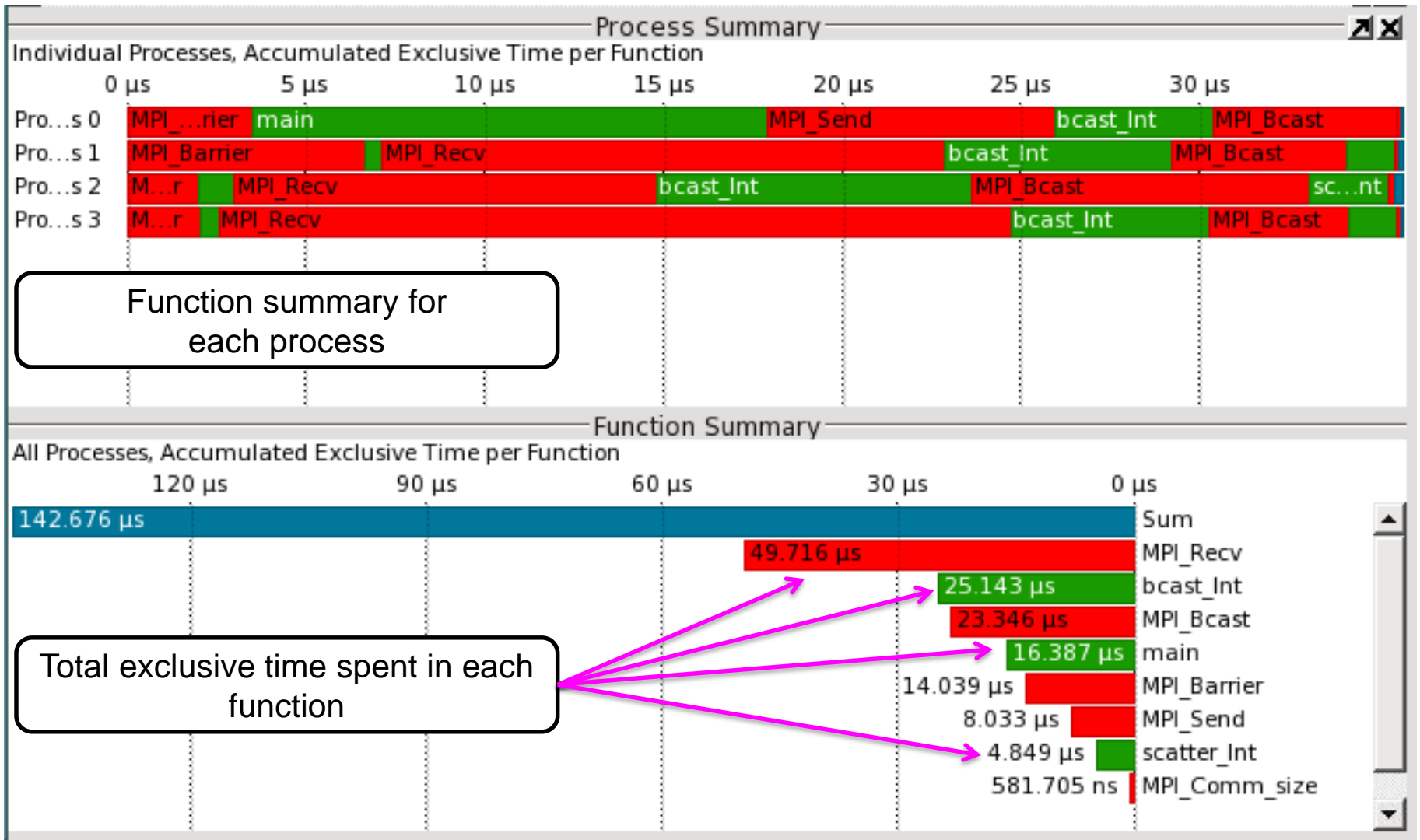
■ Event buffers are flushed when full, which takes time

- Significant skew in program's performance profile possible
- Set buffer size accordingly to prevent flushes or limit what is being recorded

```
> module load UNITE scorep cube vampir
> scorep $MPICC -o a.out program.c
--- an instrumented executable produced ---
> mpiexec -n 4 -x SCOREP_ENABLE_TRACING=0 -x SCOREP_ENABLE_PROFILING=1 \
    a.out
--- program output ---
--- profiling data written in scorep-<timestamp>/profile.cubex ---
> cube scorep-<timestamp>/profile.cubex
--- Cube GUI opens ---
> scorep-score scorep-<timestamp>/profile.cubex
Estimated aggregate size of event trace:                26kB
Estimated requirements for largest trace buffer (max_buf): 11kB
Estimated memory requirements (SCOREP_TOTAL_MEMORY):    4097kB
                                                         ^^^^^^^
> mpiexec -n 4 -x SCOREP_ENABLE_TRACING=1 -x SCOREP_ENABLE_PROFILING=0 \
    -x SCOREP_TOTAL_MEMORY=5M a.out
--- traces written in scorep-<timestamp>/traces.otf2 ---
> vampir scorep-<timestamp>/traces.otf2
--- Vampir GUI opens ---
```



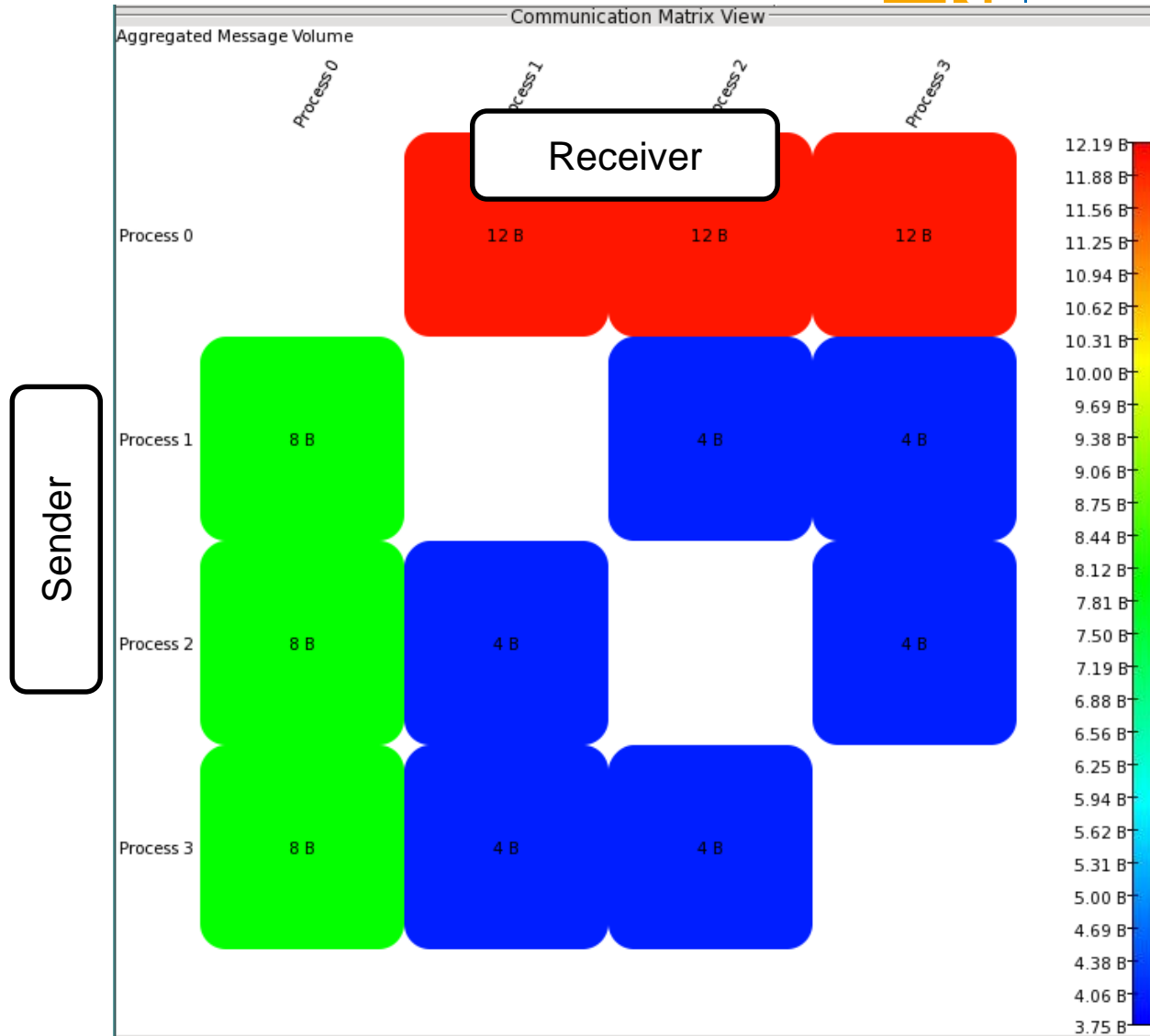




Function summary for each process

Total exclusive time spent in each function

Vampir: Communication Matrix





Live Demonstration

Thank you for your attention!