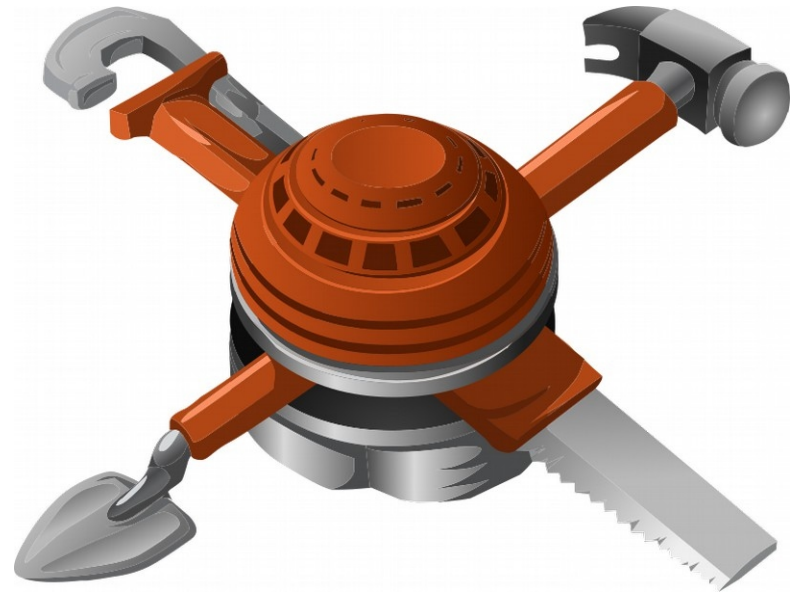


# CS 470 Spring 2025

Mike Lam, Professor



## Performance Tools

# Software Tools

- **Software tool**: computer program used by developers to create, debug, maintain or support other programs

# Traditional Software Tools

- Text editors
- Version control
- Debuggers
- Profilers
- Test automation frameworks
- Deployment frameworks
- Integrated development environments (IDEs)

# Traditional Software Tools

- **Debuggers**
  - Purpose: finding and removing software defects
  - Often done via a process monitoring interface
- **Profilers**
  - Purpose: detecting performance characteristics and identifying **bottlenecks**
  - Often done via instrumentation (added code that tracks the program's execution)
- Both of these are difficult in parallel and distributed systems

# Traditional Debugging

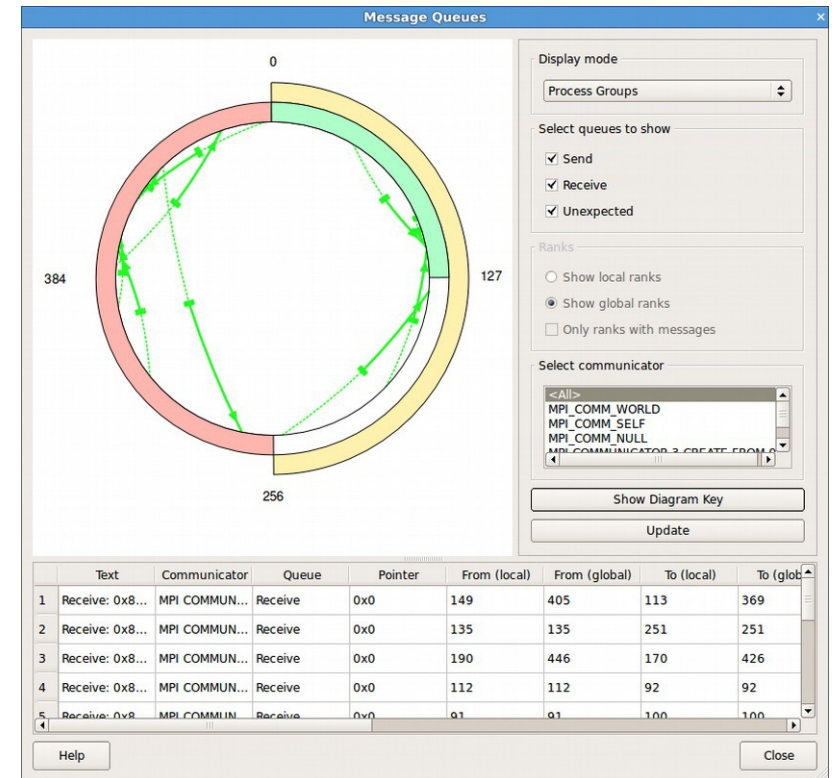
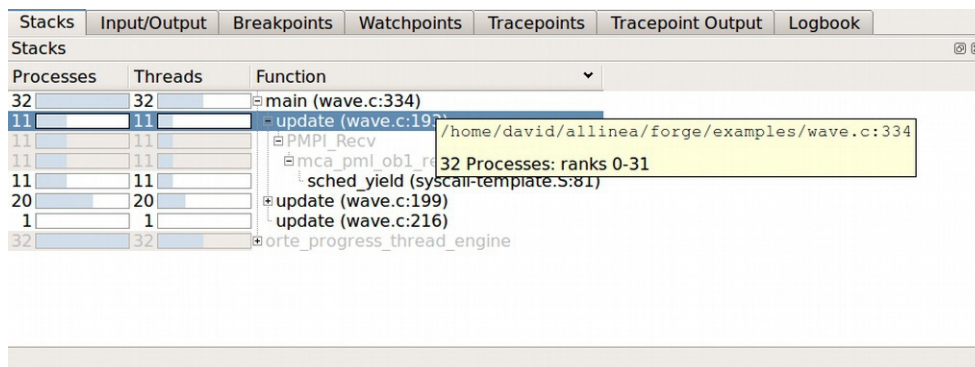
- Mechanisms
  - `ptrace`: system call that allows one process to control another
  - Simulation: slower, but safer
- Common features
  - Breakpoints and watchpoints
  - Single-stepping (by instruction or line of code)
  - Variable examination and modification
  - In newer debuggers: reverse-stepping
- Free debuggers: [gdb](#), [lldb](#), [Eclipse](#), [Valgrind \(Memcheck\)](#)

# Parallel Debugging

- Multithreaded debugging can be difficult
  - Must attach to the correct thread
  - Must control other threads as well
  - Nondeterminism means unpredictability
  - GDB does include support for multithreading:
    - <http://sourceware.org/gdb/current/onlinedocs/gdb/Threads.html>
  - Valgrind also provides the [Helgrind](#) error detector
- Distributed debugging is even harder
  - Hundreds or thousands of nodes; millions of processes
  - Enormous launch overhead
  - Control and visualization issues
  - [tmpi](https://github.com/Azrael3000/tmpi): OSS tmux hack (<https://github.com/Azrael3000/tmpi>)

# Commercial debuggers

- Microsoft Visual Studio
- Intel Debugger
- Rational Purify
- RogueWave TotalView
- Allinea DDT



# Performance Profiling

- Goal: gain insights concerning a program's performance characteristics
- Common **metrics**
  - Wall or CPU time
  - Memory use, page faults, and cache misses
  - Network traffic and saturation
  - Energy use
- Common **scopes**
  - Function
  - Basic block
  - Instruction
  - Source code line



# Measurement

- **Instrumentation**: inserting analysis code
  - Binary vs. source
  - Static vs. dynamic
  - Best for event-based monitoring (e.g., function calls)
- **Sampling**: polling an analysis source
  - Hardware counters
    - Performance Application Programming Interface (**PAPI**)
  - Randomized vs. periodic
  - Averaging vs. min/max
  - Best for continuous monitoring (e.g., memory usage)

# Measurement

- **Context**
  - Flat vs. call graph
  - Partial vs. full context
- **Profiling vs. tracing** (latter builds time-series)
- **Issues**
  - **Overhead**: added run time due to profiling software
  - **Perturbation**: skewing of behavior due to profiling software
  - **Skid**: execution may not stop immediately on sample
- **Tradeoff: better information vs. lower overhead**
  - Instrumentation: more instrumentation points
  - Sampling: higher frequency or less aggregation

# GNU Profiler (gprof)

- Compile with “-pg” flag
- Run as usual; generates “gmon.out” file
- View results with “gprof” utility
  - “gprof <executable>”
- See <https://sourceware.org/binutils/docs/gprof/> for more documentation
- Google also has a multi-threaded profiler:
  - <https://github.com/gperftools/gperftools>

# Callgrind/Cachegrind

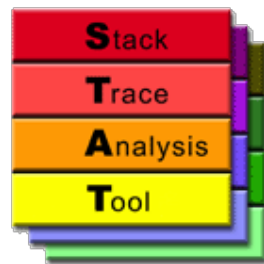
- Run with Valgrind
  - Callgrind: `valgrind --tool=callgrind <executable>`
  - Cachegrind: `valgrind --tool=cachegrind <executable>`
    - `--cache-sim=yes` or `--branch-sim=yes` to enable cache/CPU simulations
  - This will produce a `*.out.xxxx` file with raw results (could be large!)
  - Remember to call `mpirun` first if it's an MPI program
    - (And use `cg_merge` to merge multiple Cachegrind output files)
- Post-process results
  - Callgrind: `callgrind_annotate <output-files>`
    - GUI alternative: `kcachegrind` (or `qcachegrind` on Mac OS X)
  - Cachegrind: `cg_annotate <output-file>` (`--auto=yes` for code)
    - Dx = data cache (level X)      Ix = instruction cache (level X)
    - 1 = L1 cache                      L/LL = lowest level (on the cluster, this is L3)
    - r = read                      w = write                      m = miss                      Ir = Instructions read
- See <http://valgrind.org/docs/manual> for more documentation

# Perf\_events

- Sample-based performance profiler
  - Kernel module reads performance counters
    - More lightweight than Valgrind-based analysis
    - Can sample many different events
  - User space utility `perf` to interface with kernel
    - `perf record -F 49 <command>`
      - Generates `perf.data` file
    - `perf report -n [--stdio]`
    - `perf annotate [--stdio]`
  - Cheat sheet link on resources page

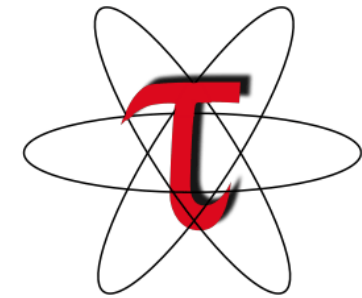
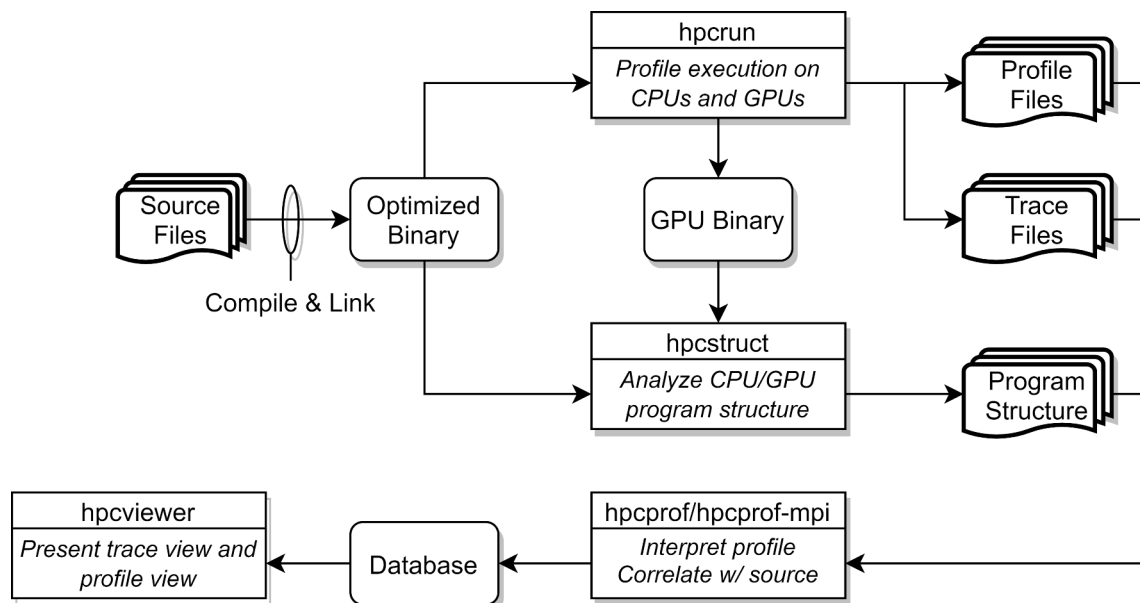
# Distributed Analysis

- Lots of data!
  - Collect at each rank but only store compressed or aggregated data
  - Aggregate using a tree-based reduction structure to reduce communication overhead
    - Research projects: [STAT](#) and [MRNet](#)



# Other HPC analysis tools

- [HPCToolkit](#) – Rice University
- [Tuning and Analysis Utilities \(TAU\)](#) – University of Oregon
- [Open|SpeedShop](#) - Krell Institute
- [Scalasca](#)
- [Paraver](#)



# Tool frameworks

- Many analysis tools need similar functionality
  - E.g., source/binary parsing or instrumentation
- **Tool framework**: a library that provides common functionality upon which custom tools can be written
  - **Rose** (source-based compiler framework)
  - **LLVM** (binary-based compiler framework)
  - **Intel Pin** (insert just-in-time binary instrumentation)
  - **Dyninst** (insert binary instrumentation)
  - **Valgrind** (track memory accesses)
  - **CRAFT** (instrument floating-point arithmetic)

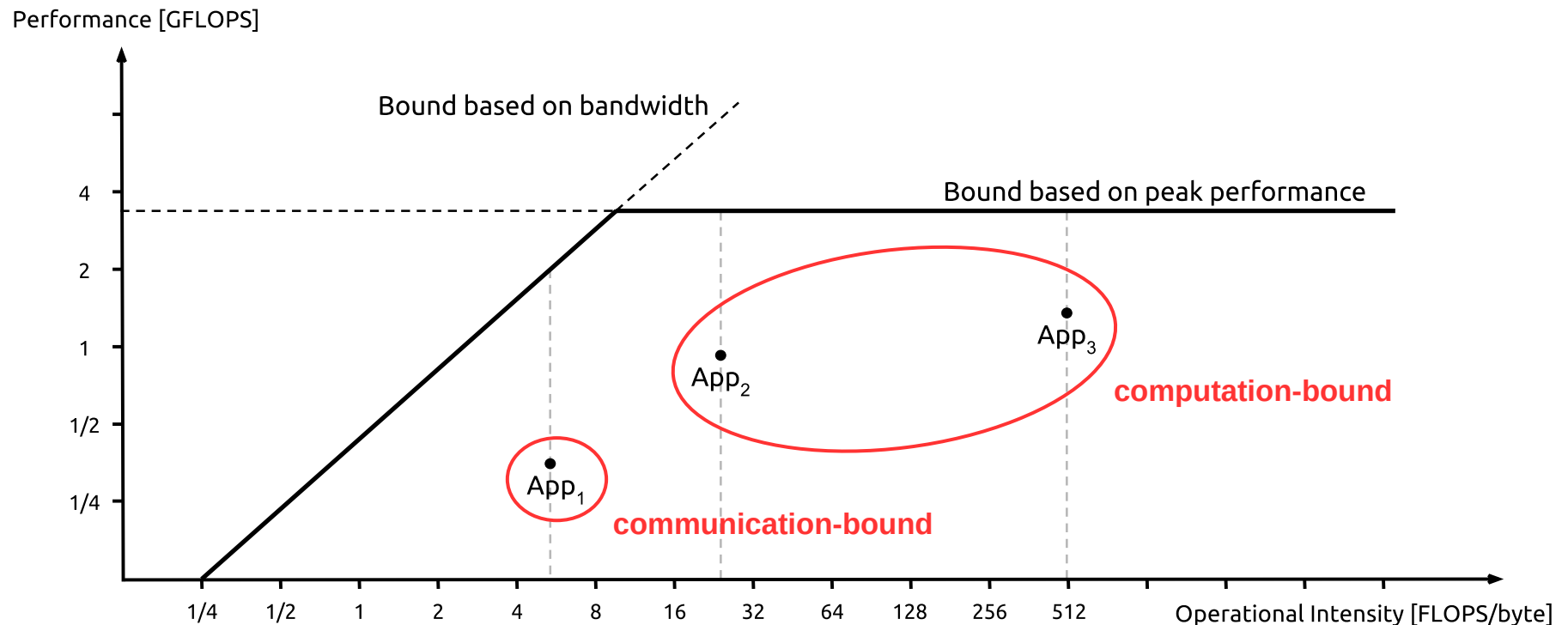


# Modeling and autotuning

- Observation: modern systems have a lot of knobs
  - Message size, block size, # of threads, # of processes
  - Many of these factors influence each other
  - Different runs could require different “optimal” settings
- Idea #1: **autotune** the system at runtime
  - Managed by middleware (the autotuner)
  - Overhead could be expensive
  - Optimizing across multiple dimensions can be difficult
- Idea #2: build a **model** of these interactions
  - Needs training data

# Performance models

- One popular model: **roofline** model
  - Shows theoretical limits on performance
  - Based on computation and communication bounds



# Performance models

