# Topics

- **Events for CPU performance measurement**

- **Kernel sampling**

- **Context recycling for dynamic threads**

# Sample Sources - I

## Linux thread-centric timers

- **CPUTIME (DEFAULT if no sample source is specified)**
  - — **CPU time used by the thread in microseconds**
  - — **does not include time blocked in the kernel**
    - – **disadvantage: completely overlooks time a thread is blocked**
    - – **advantage: a blocked thread is never unblocked by sampling**

- **REALTIME**
  - — **real time used by the thread in microseconds**
  - — **includes time blocked in the kernel**
    - – **advantage: shows where a thread spends its time, even when blocked**
    - – **disadvantages**
      - **activates a blocked thread to take a sample**
      - **a blocked thread appears active even when blocked**

Best for analysis of profile data

Produces more intuitive traces

**Note: Only use one Linux timer to measure an execution**

# Sample Sources - II

**Linux perf_event monitoring subsystem**

- **Kernel subsystem for performance monitoring**

- **Access and manipulate**
    - **hardware counters: cycles, instructions, …**
    - **software counters: context switches, page faults, …**

- **Available in modern Linux kernels**

A useful explanation about events available through perf
https://sites.google.com/site/lbathen/research/perf

# perf_event Hardware Event Counters

- **PERF_COUNT_HW_CPU_CYCLES**

- **PERF_COUNT_HW_INSTRUCTIONS**

- **PERF_COUNT_HW_CACHE_REFERENCES**

- **PERF_COUNT_HW_CACHE_MISSES**

- **PERF_COUNT_HW_BRANCH_INSTRUCTIONS**

- **PERF_COUNT_HW_BRANCH_MISSES**

- **PERF_COUNT_HW_BUS_CYCLES**

- **PERF_COUNT_HW_STALLED_CYCLES_FRONTEND**

- **PERF_COUNT_HW_STALLED_CYCLES_BACKEND**

- **PERF_COUNT_HW_REF_CPU_CYCLES**

# perf_event Hardware Cache Events

- **Hardware cache**
  - **PERF_COUNT_HW_CACHE_L1D**
  - **PERF_COUNT_HW_CACHE_L1I**
  - **PERF_COUNT_HW_CACHE_LL**
  - **PERF_COUNT_HW_CACHE_DTLB**
  - **PERF_COUNT_HW_CACHE_ITLB**
  - **PERF_COUNT_HW_CACHE_BPU**
  - **PERF_COUNT_HW_CACHE_NODE**

- **Operations**
  - **PERF_COUNT_HW_CACHE_OP_READ**
  - **PERF_COUNT_HW_CACHE_OP_WRITE**
  - **PERF_COUNT_HW_CACHE_OP_PREFETCH**

- **Results**
  - **PERF_COUNT_HW_CACHE_RESULT_ACCESS**
  - **PERF_COUNT_HW_CACHE_RESULT_MISS**

# perf_event Software Events

- **PERF_COUNT_SW_CPU_CLOCK**

- **PERF_COUNT_SW_TASK_CLOCK**

- **PERF_COUNT_SW_PAGE_FAULTS**

- **PERF_COUNT_SW_CONTEXT_SWITCHES**

- **PERF_COUNT_SW_CPU_MIGRATIONS**

- **PERF_COUNT_SW_PAGE_FAULTS_MIN**

- **PERF_COUNT_SW_PAGE_FAULTS_MAJ**

- **PERF_COUNT_SW_ALIGNMENT_FAULTS**

- **PERF_COUNT_SW_EMULATION_FAULTS**

# Measuring Other Hardware Events

- **See the full list of available events with**
  - hpcrun -L

- **Perf events are grouped by categories indicated by a prefix**
  - **ix86arch::<event>**                // Intel architecture
  - **perf::<event>**                         // perf_event builtin
  - **bdw_ep::<event>**                    // Broadwell EP specific
  - **…**

- **For convenience**
  - **you may omit the category prefix, e.g. "perf::"**
  - **you may specify counter names using lower case**

# Multiplexing Events

- **In a single execution, you can measure more hardware events than the number of hardware counters available per thread**

- **If you specify more events than counters available**
  - — **perf_events will automatically multiplex them**

- **How multiplexing works with Linux perf_event subsystem**
  - — **at any time, the number of events being collected will not exceed the number of hardware counters available per thread**
  - — **the kernel will partition events into sets that can be monitored simultaneously using hardware counter resources**
  - — **the kernel will monitor one set of events for a while then switch to another**
  - — **monitoring of event sets is scheduled in round-robin fashion**
  - — **while multiplexing is convenient, there is some loss of accuracy**
    - – **my advice: multiplexing is fine for casual execution analysis**

# Controlling perf_event Sampling Frequency

- **Automatic** <span style="color:orange">Recommended</span>
  - **HPCToolkit samples perf_event counters min(300x/second, maximum Linux allows)**
    - **may be higher than necessary for long executions**
      reducing the frequency will reduce measurement overhead

- **Specify frequency**
  - **use the @f<freq> suffix for an event to specify frequency**
    - **hpcrun -e CYCLES@f100 -e INSTRUCTIONS@f200 …**
  - **Specify a different default frequency using the -c option**
    - **example: sample both CYCLES and INSTRUCTION 200x per second**
      hpcrun -c f200 -e CYCLES -e INSTRUCTIONS

- **Specify period**
  - **Use the @<period> suffix for an event to specify a period**
    - **hpcrun -e CYCLES@1000000 -e INSTRUCTIONS@5000000 …**

# Topics

- **Events for CPU performance measurement**
- **Kernel sampling**
- **Context recycling for dynamic threads**

# Topics

- **Events for CPU performance measurement**
- **Kernel sampling**
- **Context recycling for dynamic threads**

# Kernel Sampling in HPCToolkit

- **When sampling using the Linux perf_event subsystem**
  - **sample user space activity**
  - **sample kernel space activity**

- **When a thread is active in the kernel, the user calling context is frozen**

- **Attribute kernel activity to the point where it occurred in the user calling context**
  - **form a calling context that has**
    - **user calling context as the prefix**
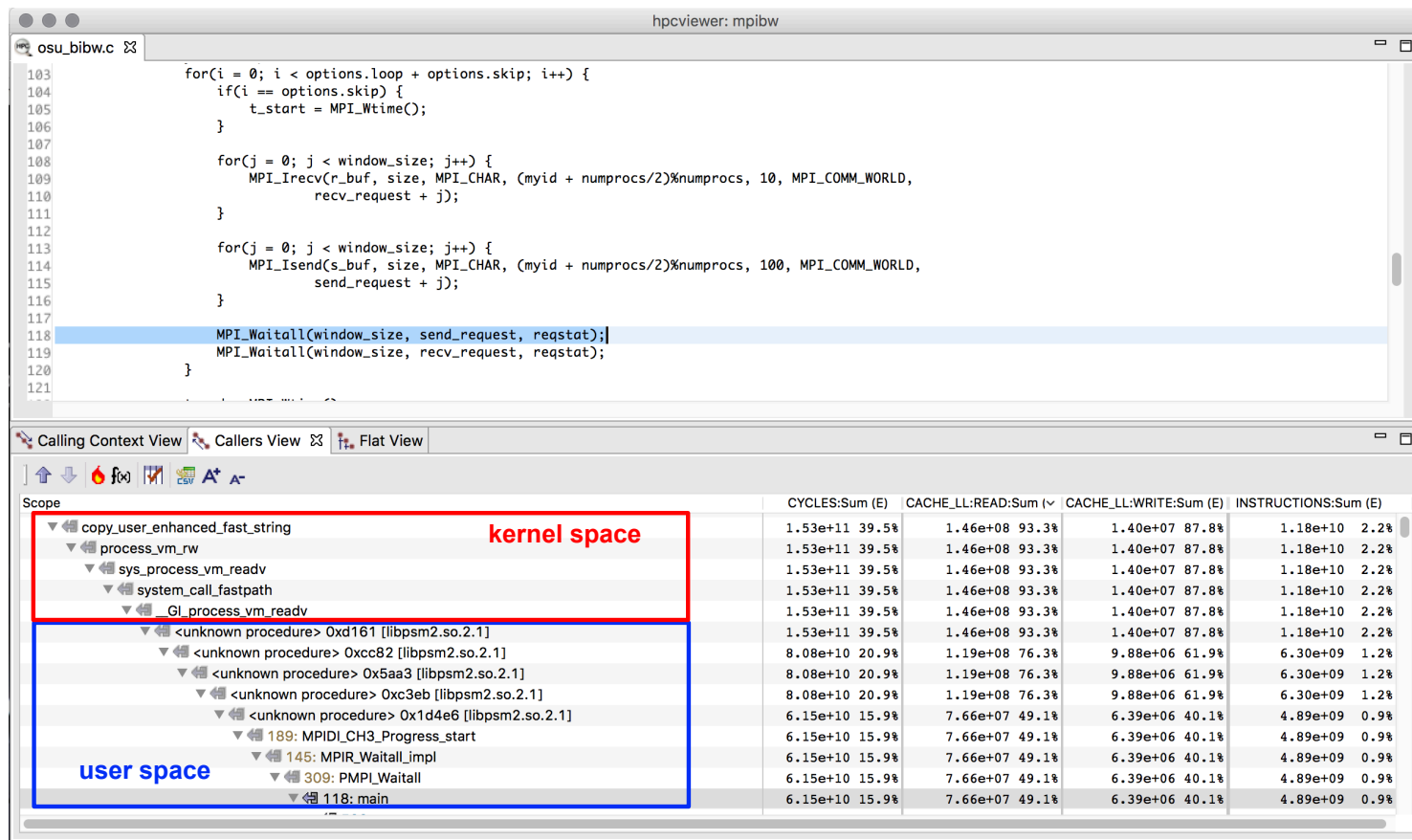    - **kernel calling context as the suffix**

# Kernel Sampling Yields Insight

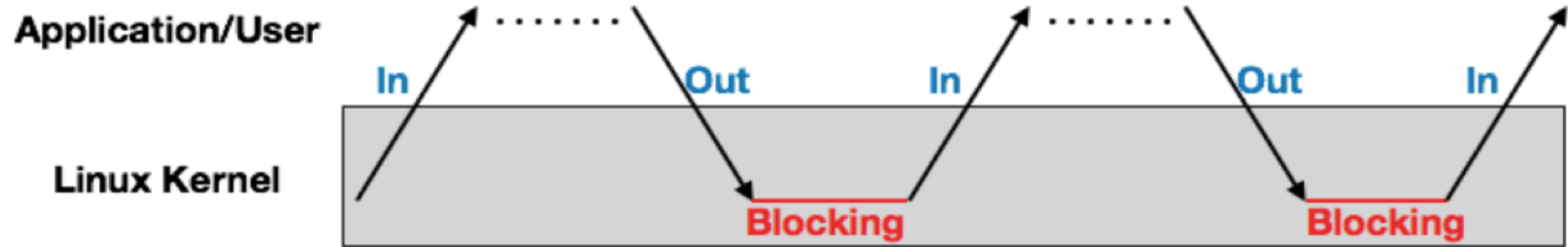## Investigating MPI Performance with Kernel Sampling

**Platform**
- **Intel Broadwell**
- **Infiniband network**

- **Q: Why is MPI communication bandwidth so low on node (6-9 GB/s)?**

- **A: Bounded by single thread memory bandwidth**

- **Memcpy 12 GB/s**
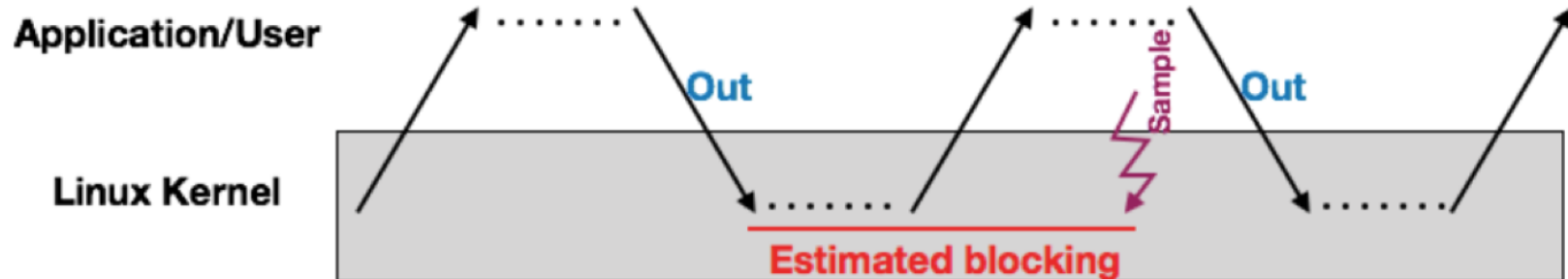
- **Stream (1T) 8-9 GB/s**

- **Stream (OMP) 60 GB/s**

hpcviewer: mpibw

osu_bibw.c

```
103   for(i = 0; i < options.loop + options.skip; i++) {
104       if(i == options.skip) {
105           t_start = MPI_Wtime();
106       }
107
108       for(j = 0; j < window_size; j++) {
109           MPI_Irecv(r_buf, size, MPI_CHAR, (myid + numprocs/2)%numprocs, 10, MPI_COMM_WORLD,
110                     recv_request + j);
111       }
112
113       for(j = 0; j < window_size; j++) {
114           MPI_Isend(s_buf, size, MPI_CHAR, (myid + numprocs/2)%numprocs, 100, MPI_COMM_WORLD,
115                     send_request + j);
116       }
117
118       MPI_Waitall(window_size, send_request, reqstat);
119       MPI_Waitall(window_size, recv_request, reqstat);
120   }
121
```

Calling Context View    Callers View    Flat View

| Scope | CYCLES:Sum (E) | | CACHE_LL:READ:Sum (∨) | | CACHE_LL:WRITE:Sum (E) | | INSTRUCTIONS:Sum (E) | |
|---|---|---|---|---|---|---|---|---|
| ▼ copy_user_enhanced_fast_string     **kernel space** | 1.53e+11 | 39.5% | 1.46e+08 | 93.3% | 1.40e+07 | 87.8% | 1.18e+10 | 2.2% |
| ▼ process_vm_rw | 1.53e+11 | 39.5% | 1.46e+08 | 93.3% | 1.40e+07 | 87.8% | 1.18e+10 | 2.2% |
| ▼ sys_process_vm_readv | 1.53e+11 | 39.5% | 1.46e+08 | 93.3% | 1.40e+07 | 87.8% | 1.18e+10 | 2.2% |
| ▼ system_call_fastpath | 1.53e+11 | 39.5% | 1.46e+08 | 93.3% | 1.40e+07 | 87.8% | 1.18e+10 | 2.2% |
| ▼ __GI_process_vm_readv | 1.53e+11 | 39.5% | 1.46e+08 | 93.3% | 1.40e+07 | 87.8% | 1.18e+10 | 2.2% |
| ▼ <unknown procedure> 0xd161 [libpsm2.so.2.1] | 1.53e+11 | 39.5% | 1.46e+08 | 93.3% | 1.40e+07 | 87.8% | 1.18e+10 | 2.2% |
| ▼ <unknown procedure> 0xcc82 [libpsm2.so.2.1] | 8.08e+10 | 20.9% | 1.19e+08 | 76.3% | 9.88e+06 | 61.9% | 6.30e+09 | 1.2% |
| ▼ <unknown procedure> 0x5aa3 [libpsm2.so.2.1] | 8.08e+10 | 20.9% | 1.19e+08 | 76.3% | 9.88e+06 | 61.9% | 6.30e+09 | 1.2% |
| ▼ <unknown procedure> 0xc3eb [libpsm2.so.2.1] | 8.08e+10 | 20.9% | 1.19e+08 | 76.3% | 9.88e+06 | 61.9% | 6.30e+09 | 1.2% |
| ▼ <unknown procedure> 0x1d4e6 [libpsm2.so.2.1] | 6.15e+10 | 15.9% | 7.66e+07 | 49.1% | 6.39e+06 | 40.1% | 4.89e+09 | 0.9% |
| ▼ 189: MPIDI_CH3_Progress_start | 6.15e+10 | 15.9% | 7.66e+07 | 49.1% | 6.39e+06 | 40.1% | 4.89e+09 | 0.9% |
| ▼ 145: MPIR_Waitall_impl | 6.15e+10 | 15.9% | 7.66e+07 | 49.1% | 6.39e+06 | 40.1% | 4.89e+09 | 0.9% |
| ▼ 309: PMPI_Waitall   **user space** | 6.15e+10 | 15.9% | 7.66e+07 | 49.1% | 6.39e+06 | 40.1% | 4.89e+09 | 0.9% |
| ▼ 118: main | 6.15e+10 | 15.9% | 7.66e+07 | 49.1% | 6.39e+06 | 40.1% | 4.89e+09 | 0.9% |

# Measure Thread Blocking using perf_events



Original idea: Kernel **blocking** time

Application/User

In    Out    In    Out    In

Linux Kernel

Blocking    Blocking

Our approach: Estimated kernel blocking time
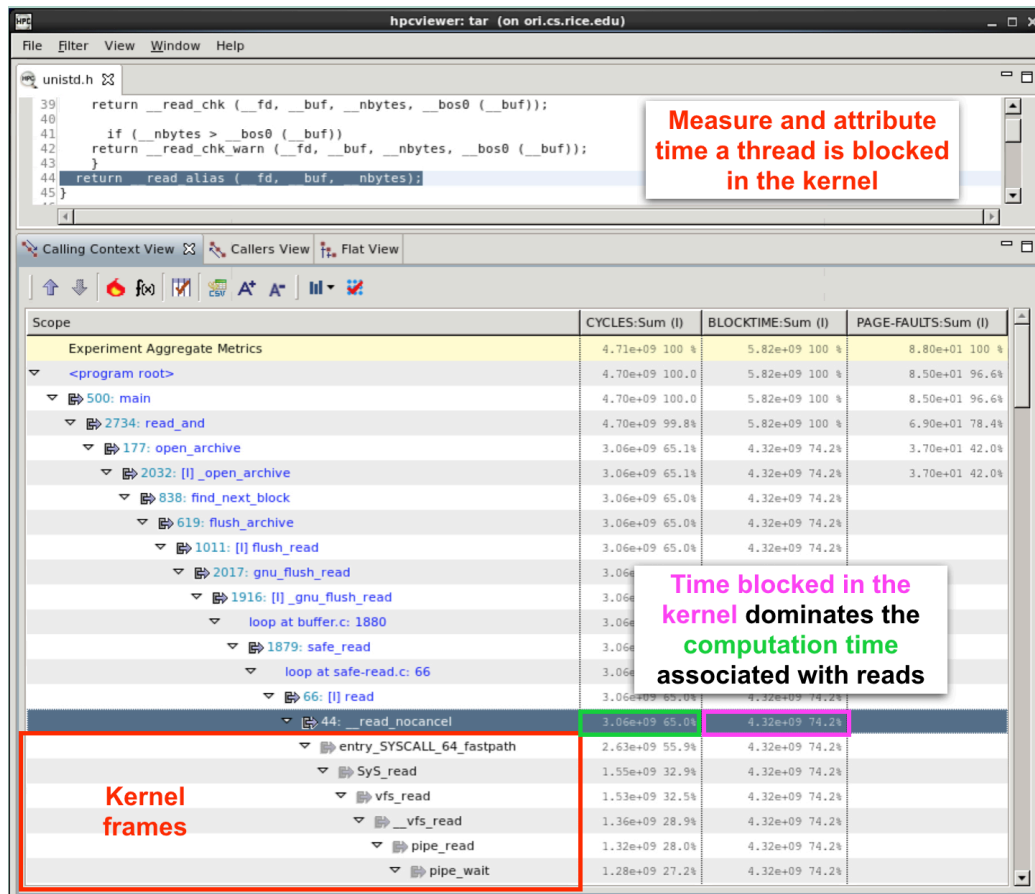
Application/User

Out    Sample    Out

Linux Kernel

**Estimated blocking**

# Example: Thread Blocking in "tar"

hpcrun -e CYCLES -e BLOCKTIME -e PAGE-FAULTS tar xzf \
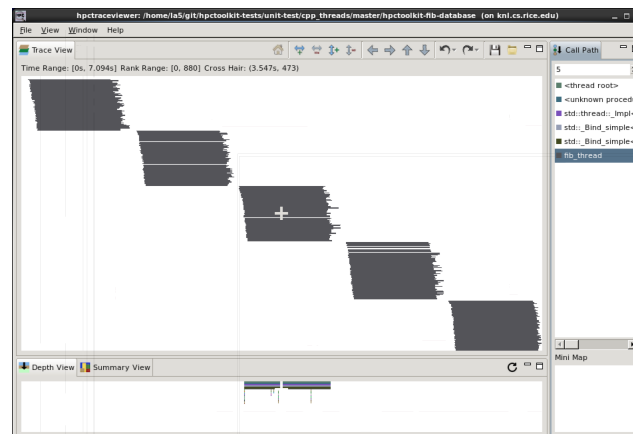~/Downloads/eclipse-rcp-indigo-linux-gtk-x86_64.tar.gz

# A Few More Things

- **Events for CPU performance measurement**

- **Differential performance analysis (useful for CPU and GPU)**

- **Kernel sampling**

- **Context recycling for dynamic threads**

# Context Recycling for Short-lived Threads

- **Problem**
  - **some codes create many short-lived threads**
    - **DCA+ 160 ranks generated 1.2M thread profiles and traces**
  - **time-centric views of such codes are problematic**

- **Solution**
  - **when a thread completes, put its (CCT, trace) in a free list**
  - **when a new thread starts, look for an available (CCT, trace) pair to augment**
  - **create a new one only if needed**

Credit: Laksono Adhianto

18

# DCA+ using Context Recycling

## DCA+ 10 ranks, 12 threads each with context recycling