

Feb 9, 2023 | 📅 RISC-V Perf Analysis SIG Meeting

Attendees: Beeman Strong tech.meetings@riscv.org Marc Casas

Notes

- **Attendees:** Beeman, Bruce, DmitriyR, StasN, Fei Wu, DmitriyP, RobertC, AllenB, OlegL, Ved
- **Slides/video** [here](#)
- Opens
 - Supervisor Counter Delegation approved for fast track, still waiting to hear on Instret & Cycle Privilege Mode Filtering approval
- Stas and others worked for many years on Intel Vtune
- Developing a performance analysis tool for RISC-V, akin to Vtune
- Need precise events, with little/no skid (ala Intel)
 - Allows precise attribution of sample to an instruction
- PC is most important state to capture precisely, don't care how
- Which events are precise?
 - Tend to be retire time events, or events that get plumbed to be counted at retire
 - Frontend events, e.g., are usually speculative and then not precise. In fact, usually skid "backwards" (sample before the inst that causes, say, an icache miss).
- Need a uarch methodology, ala Intel's top-down
 - Quickly identify uarch bottlenecks via generic, hierarchical metrics
 - Collected over time, but can do it at function level rather than system-level
 - Dick Sites has asked for always available counters for the same top-level events
 - Essentially asking for standardized events for these
 - Do the encodings need to be standardized, or just the events definitions?
 - Want to be able to ask for these events from Linux perf, so if perf supports an alternate approach (e.g., event list) that's fine
- Looking to port Vtune to RISC-V?
 - No, Intel continues developing Vtune. Have ambitions for something Vtune-like.
- Memory access analysis
 - Intel provides DLA (data linear address) for precise memory events
 - Akin to xtval on exceptions. But interrupts don't update xtval. Could have xtval set on LCOFI?
- SoC Monitoring
 - Measuring traffic at DRAM, PCIe, etc
 - Discoverable and unified interface
 - Can be just counters, don't need sampling
 - Beeman: planned to collaborate with DTPM SIG to develop standards here
 - IOMMU spec defines standard HPMs for IOMMU, and set of standard events and MSI on overflow
 - What is purpose of interrupt on overflow here?

- To collect ratios, e.g. misses per X cycles
 - Counter freeze can help for ratios like this
 - Intel's Uncore PMI not used by tools, hasn't worked reliably
 - Can sample on CPU events (cycles, time), don't need samples from external source
 - Hard to get standard interface into external IPs
- Call stacks
 - Issues unwinding call-stack (no frame pointers, etc)
 - LBR call-stack has limited depth, SW not always well-behaved
 - For flame graphs, fragmented if don't have the full stack
 - If you can guarantee to get 32 depth, but LBR doesn't now due to SW exceptions
 - SW jumps and adjusts stack ptr, HW can't account for this
 - Maybe could identify these cases with a calling convention?
 - Could we make a fast way to adjust the LBR/CTR TOS that SW unwind could adjust?
 - Shadow stack (part of CFI) provides an unwind mechanism. Can be used in place of call-stack unwind. Being developed in TG.
 - Security concerns with accessing SS? Just can't write it
- HW Tracing
 - Used for anomaly detection (e.g., the 1% of transactions that are bad)
 - RV has E-trace & N-trace standards, would be nice to integrate HPMs
 - Vtune uses PTWRITE inst to insert counter values into trace, used to identify causes of anomalies
 - N-trace can insert counter values into trace based on trigger match, without instrumentation
 - Would like to get top-down metrics per transaction
- Need to be able to profile within a VM, and need eBPF working

Action items

- ☐ Atish Kumar Patra - Aug 25, 2022 - check on how to read multiple counters in perf when taking an interrupt on one
- ☐ Beeman Strong - Jul 28, 2022 - Reach out about proprietary performance analysis tools
- ☐ Beeman Strong - Jul 28, 2022 - Reach out to VMware about PMU enabling
- ☐ Beeman Strong - Jul 28, 2022 - Talk to security HC about counter delegation