

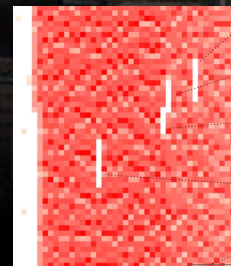
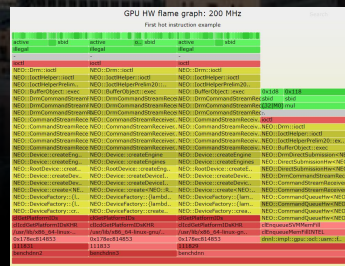
LINUX PLUMBERS CONFERENCE

TOKYO, JAPAN / DECEMBER 11-13, 2025

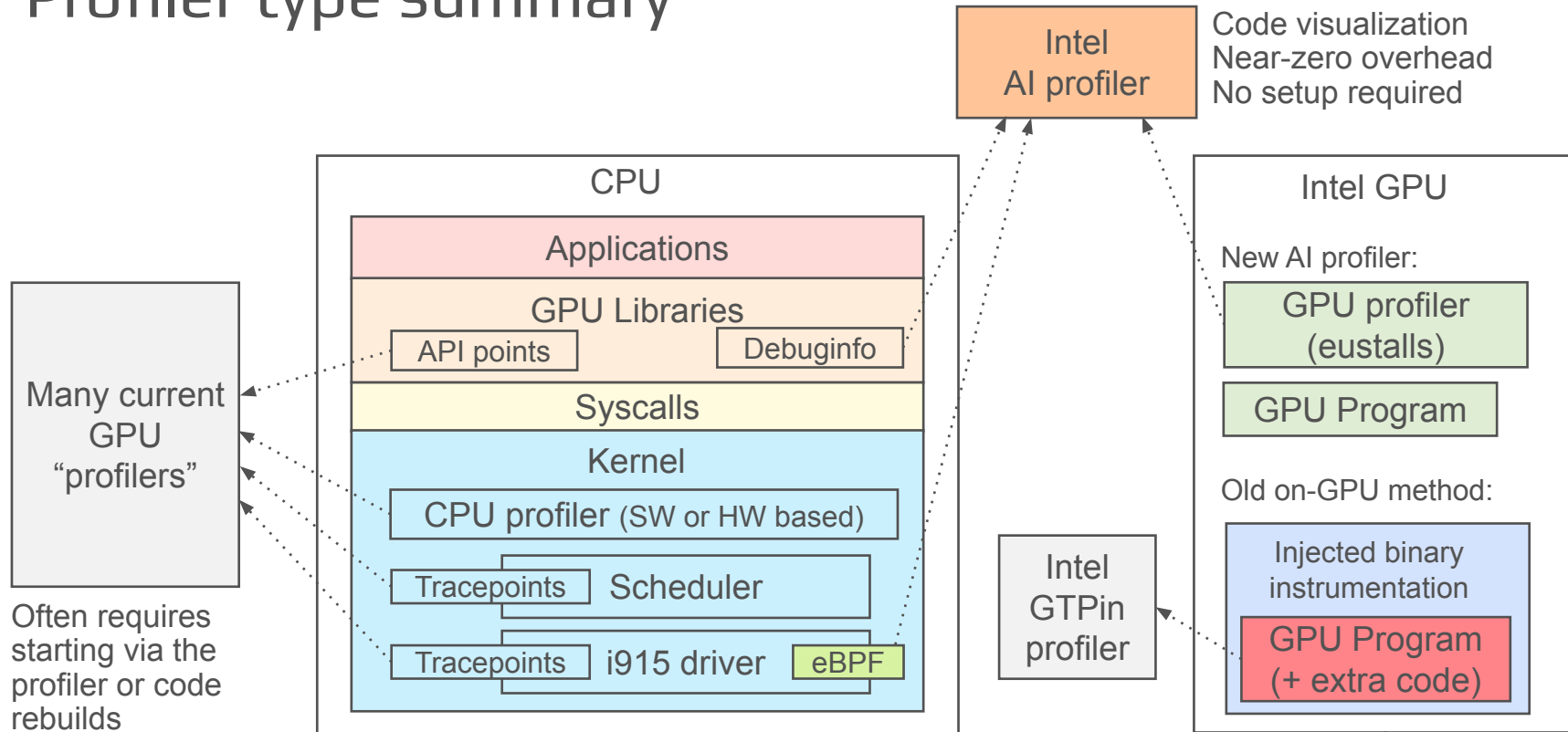
AI Flame Graphs With eBPF

Brendan Gregg, Ben Olson

Dec 2025

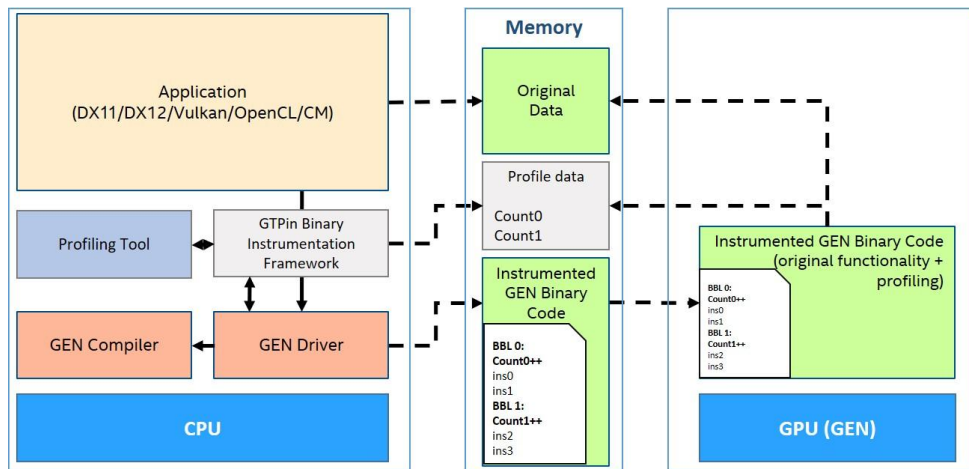


Profiler type summary



Overhead

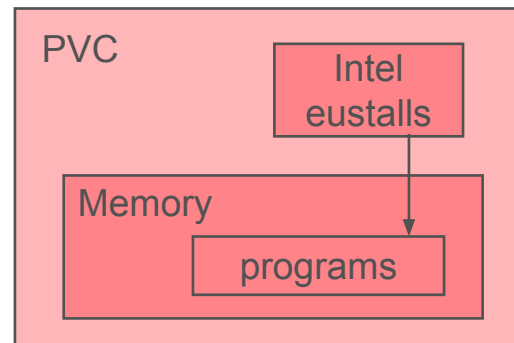
Other profilers (e.g., GTPin)



Source: <https://software.intel.com/sites/landingpage/gtpin/index.html>

Does show GPU software internals,
but costs high overhead.

Intel AI Profiler



Programs run as-is.
HW-based profiling: Intel eustalls
(execution unit stalls).

AI/
GPU

Instruction offset
Stall reason
Instruction
mnemonic
Stack/source
functions
Source file

Legend

GPU HW
GPU source
OS Kernel
CPU C++
CPU C
Python

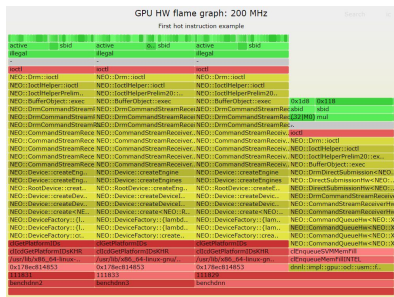
CPU call
stack

PID
Program

AI Flamegraph				Search	ic
Intel					
0x220	0xf8	0x250			
sbid	sbid	sbid			sbid
mov	mul	send			sync
CopyBufferRectBytes2d line 3					typeinfo n..
src/compute-runtime/shared/source/built_ins/kernels/copy_buffer_rect.builtin_kernel					benchmarks..
-					-
ioctl					ioctl
NEO::SysCalls::ioctl(int, unsigned long, void*)					NEO::SysCal..
NEO::Drm::ioctl(NEO::DrmIoctl, void*)					NEO::Drm::..
NEO::IoctlHelper::ioctl(NEO::DrmIoctl, void*)					NEO::Ioctl..
NEO::IoctlHelperPrelim20::execBuffer(NEO::ExecBuffer*, unsigned long, unsigned long)					NEO::Ioctl..
NEO::BufferObject::exec(unsigned int, unsigned long, unsigned int, bool, NEO::OsContext*,..					NEO::Buffe..
NEO::DrmDirectSubmission<NEO::XeHpcCoreFamily, NEO::RenderDispatcher<NEO::XeHpcC..					NEO::DrmD..
NEO::DirectSubmissionHw<NEO::XeHpcCoreFamily, NEO::RenderDispatcher<NEO::XeHpcCor..					NEO::Direct..
NEO::DirectSubmissionHw<NEO::XeHpcCoreFamily, NEO::RenderDispatcher<NEO::XeHpcCo..					NEO::Direct..
NEO::DrmCommandStreamReceiver<NEO::XeHpcCoreFamily>::flush(NEO::BatchBuffer&, std..					NEO::Dr..
NEO::CommandStreamReceiverHw<NEO::XeHpcCoreFamily>::flushHandler(NEO::BatchBuff..					NEO::Co..
NEO::CommandStreamReceiverHw<NEO::XeHpcCoreFamily>::handleImmediateFlushSendB..					NEO::Com..
NEO::CommandStreamReceiverHw<NEO::XeHpcCoreFamily>::flushImmediateTask(NEO::Line..					NEO::Co..
LO::CommandListCoreFamilyImmediate<(GFXCORE_FAMILY)3080>::flushImmediateRegular..					LO::Comma..
LO::CommandListCoreFamilyImmediate<(GFXCORE_FAMILY)3080>::executeCommandList..					LO::Comma..
LO::CommandListCoreFamilyImmediate<(GFXCORE_FAMILY)3080>::executeCommandListI..					LO::Comma..
LO::CommandListCoreFamilyImmediate<(GFXCORE_FAMILY)3080>::flushImmediate(_ze_res..					LO::Comma..
LO::CommandListCoreFamilyImmediate<(GFXCORE_FAMILY)3080>::appendMemoryCopyReg..					LO::Comma..
LO::zeCommandListAppendMemoryCopyRegion(_ze_command_list_handle_t*, void*, _ze_c..					LO::zeCom..
zeCommandListAppendMemoryCopyRegion					zeComman..
enqueueMemCopyRectHelper(ur_command_t, ur_queue_handle_t*, void const*, void*, ur_r..					urEnqueue..
1409108					
matrix.dpcpp					
all					

Reading a full (inverted) GPU Flame Graph

“why”



“how”

Program name
PID
CPU call stack (user)
CPU call stack (kernel)
-
GPU source directories
GPU source file
GPU function stack
GPU instruction mnemonic
GPU stall reason
GPU instruction offset

“This application

...uses the GPU for *this* reason

...via this driver

...and runs *this* GPU application

...from *this* source file

...and runs *this* function

...which runs *this* instruction

...and is slow for *this* reason

...for these *exact* instructions.”

AI Profiler: Requirements

Near-zero overhead

FAANGs target <0.1% overhead for profilers. A mere 5% overhead can trigger “bad instance” detection and auto termination. Methods involving baked-in instruction instrumentation are non-starters. We keep overhead low using eBPF and eustalls.

Easy to deploy

The profiler should not require special deployments such as requiring the developer SSH to their servers and modify their application start scripts. Most developers don't SSH anymore and some sites have completely removed SSH, plus many sites use continuous delivery UIs leaving the developer with no idea how to manually modify the application start process to include some Intel library/layer. We used advanced eBPF to automatically instrument the kernel, no restarts of anything required. This method is becoming known nowadays as “**zero instrumentation**.”

Easy to understand

The tool must speak the language of the developer: show them their own code function/method names.

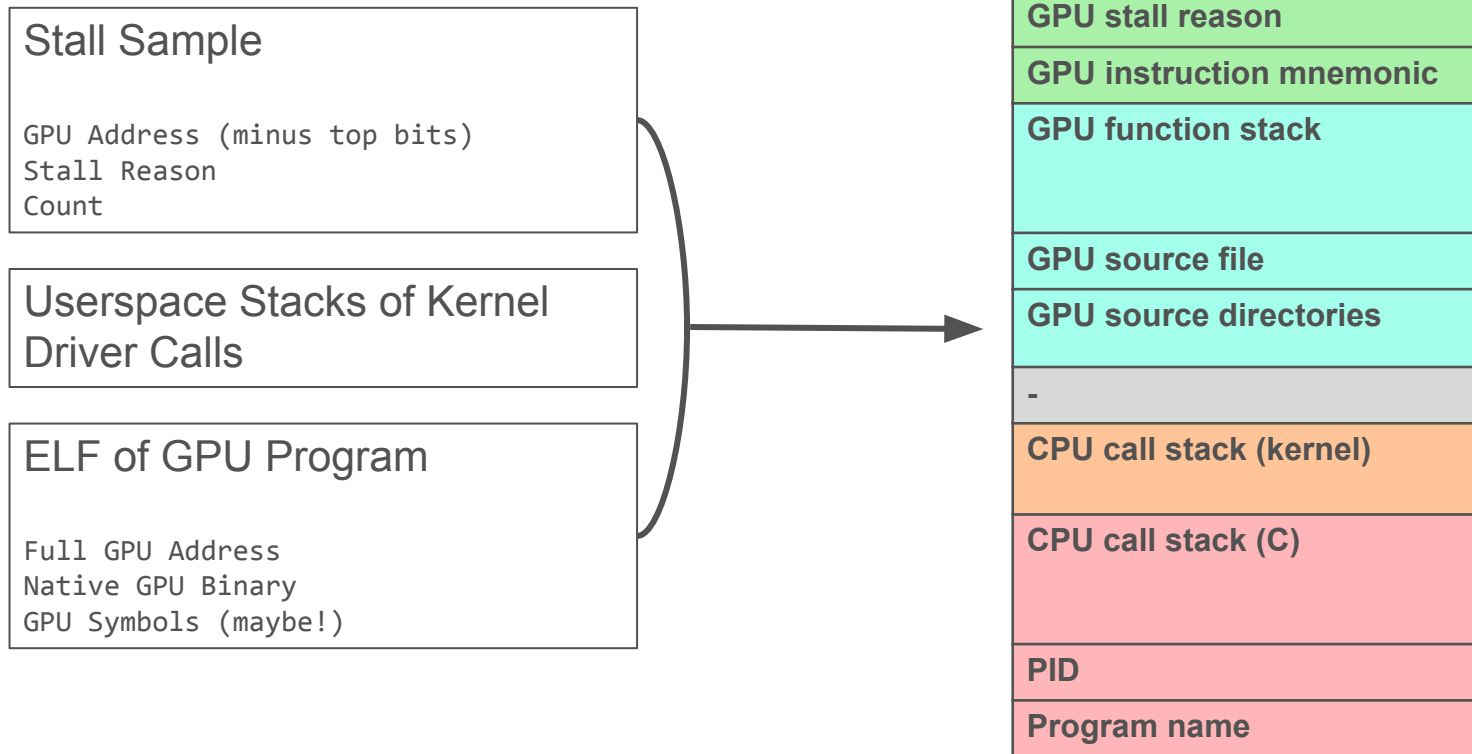
Actual GPU HW profiling

Many current GPU profilers do not see inside the HW, and merely time requests. They are not “profilers” in the common use of the term. Our solution actually sees inside the HW thanks to Intel eustalls.

Complete: CPU & GPU

Full stack visibility, to answer why (GPU) and how (GPU).

How



A Tale of Two Methods

uprobes

Probe runtime writing commands



Extract program addresses and names

kprobes

- Probe where all buffer objects get allocated by the kernel driver
- Probe where all buffer objects get bound to a VM
- Probe where commands get sent to the GPU
- Implement a simulator of a GPU Command Streamer, including:
 - Registers
 - Feature flags
 - Base addresses
 - Per-context instruction pointers
 - 3D commands
 - Program execution
 - Semaphores
 - Jumps
- Ensure the simulator doesn't run forever
- Ensure the simulator sees every single buffer object allocation, binding, and execution
- Ensure overhead of simulator doesn't grind workload to a halt
- Maintain BPF maps of GPU address ranges containing programs
- Send GPU address ranges of programs to userspace
- Parse ELF's from a special driver feature to associate GPU addresses with program names

Caveats

Need to collect ELF's of all GPU programs, requires expensive cooperation with the kernel driver (Xe).

GPU instruction offset	
GPU stall reason	
GPU instruction mnemonic	
GPU function stack	
GPU source file	
GPU source directories	
-	
CPU call stack (kernel)	
CPU call stack (C)	(C++)
PID	
Program name	

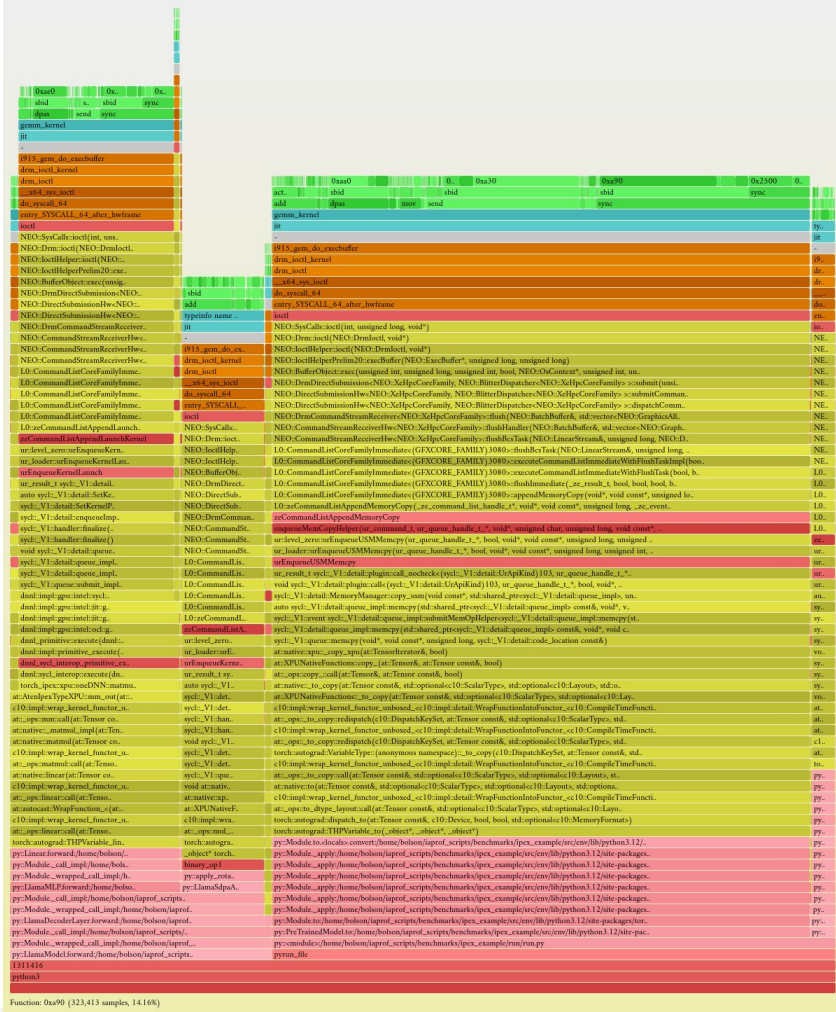
CPU-side callstacks currently require frame pointers built into *all* layers of the stack.

Underlying hardware sampling on Intel hardware, eustalls, does not provide a context ID to differentiate between simultaneous workloads. Workarounds in SW, but need HW changes.

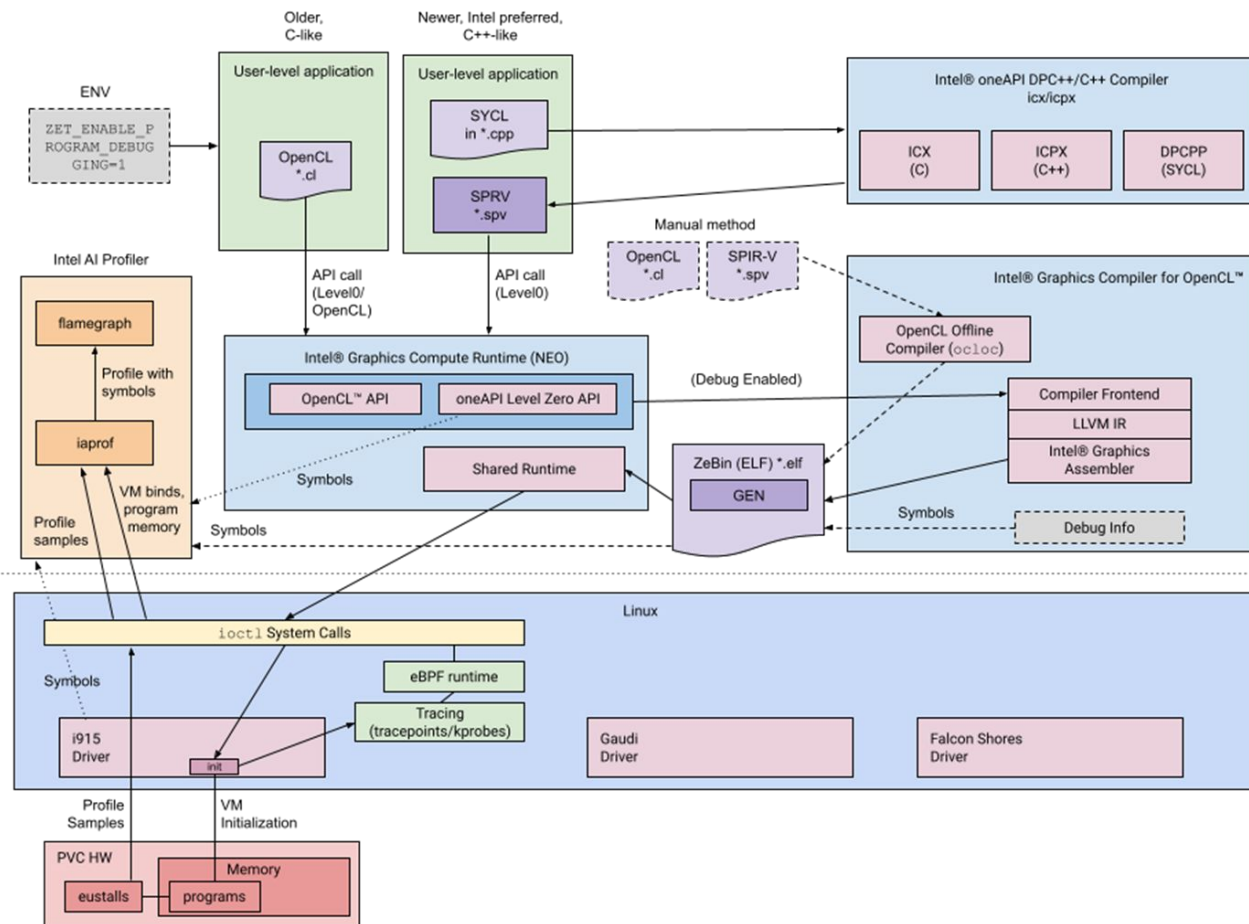
Currently using GPU kernel name, with very little support for actual callstacks in GPU code.

A few choices on which CPU callstack to use: buffer object allocation, execbuffer, or maybe something else?

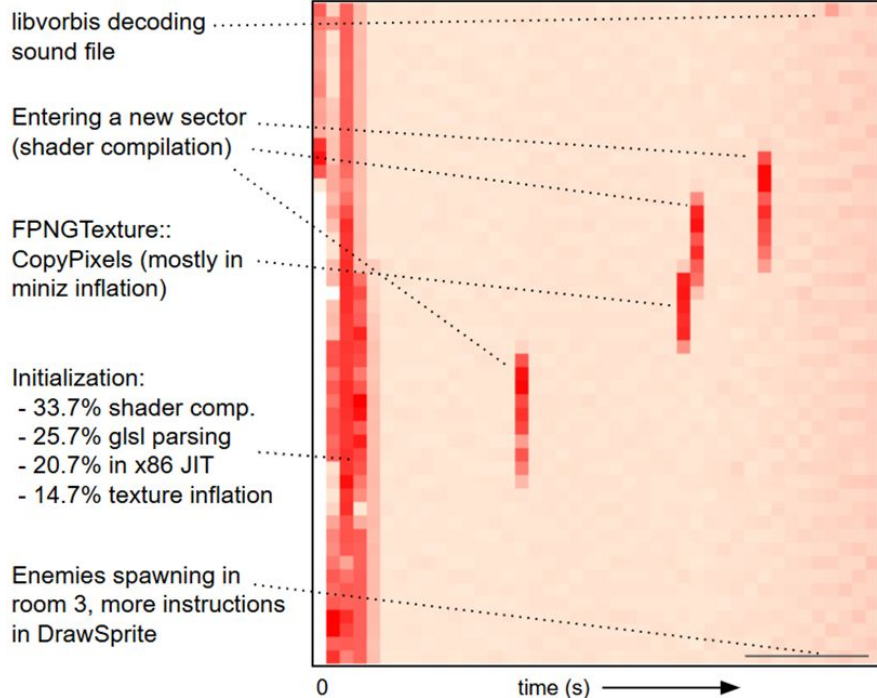
PyTorch



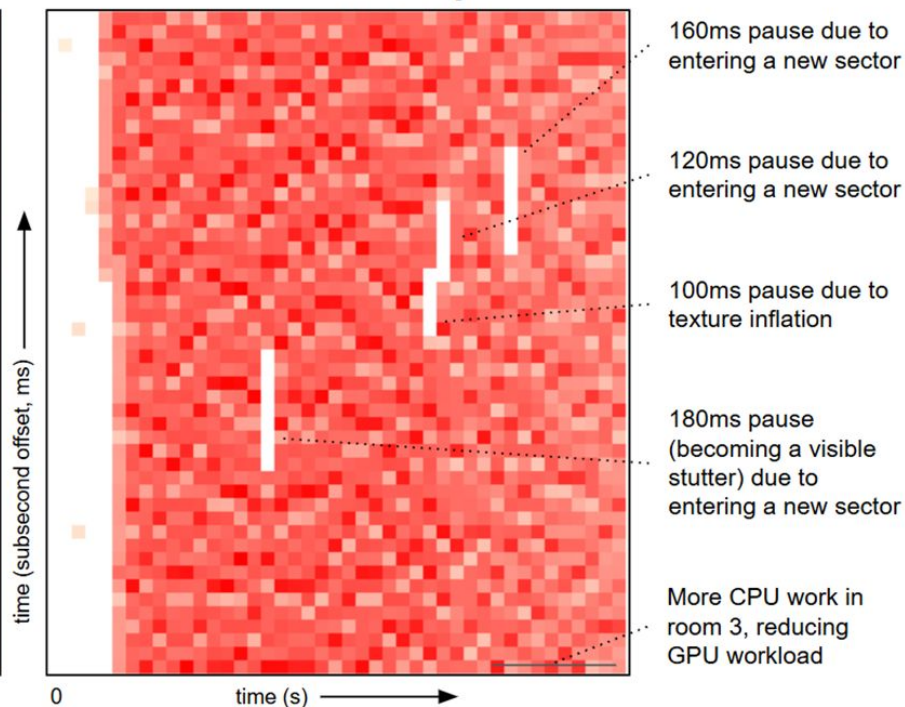
Intel AI Profiler Internals



CPU FlameScope

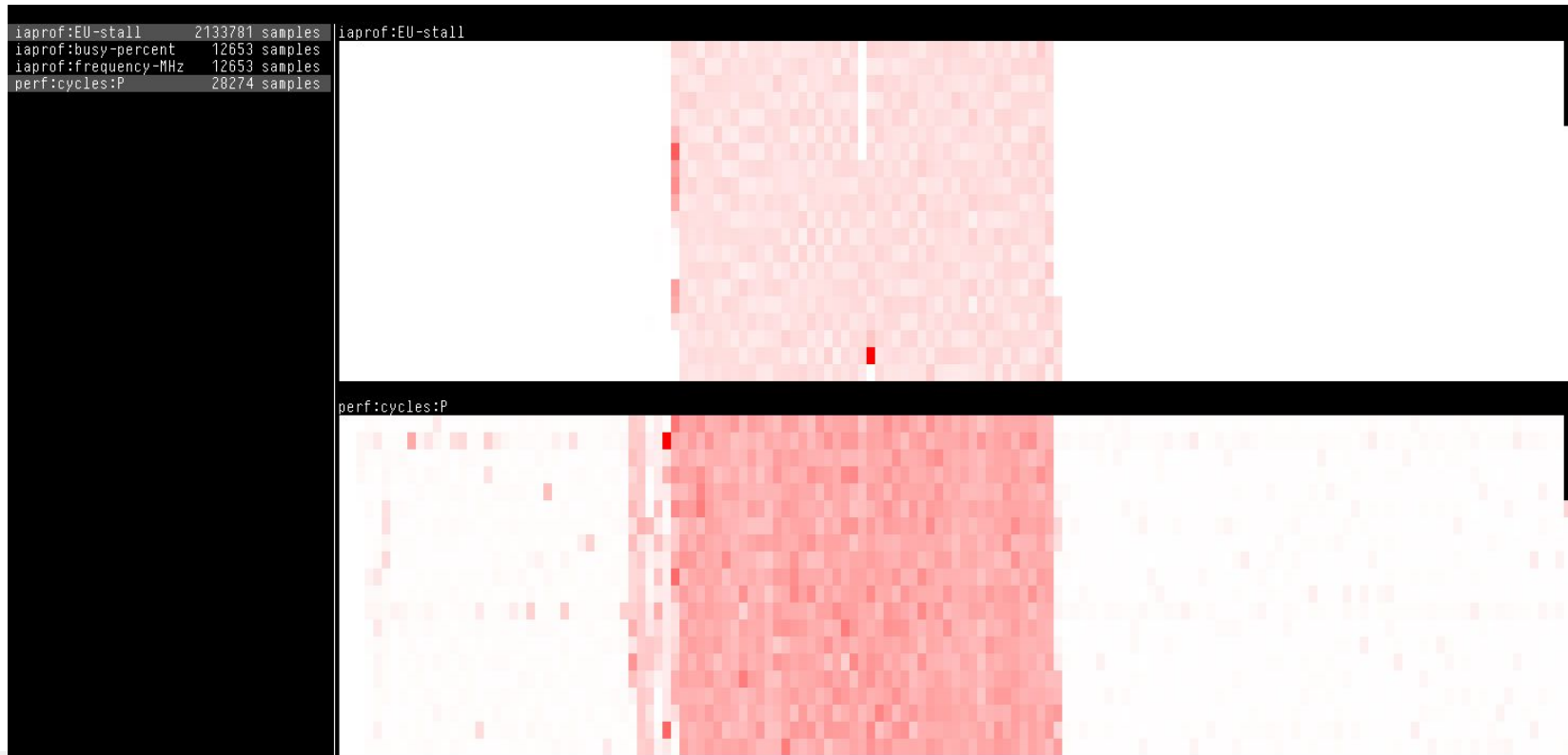


GPU FlameScope



visual correlation

Proviz (demo?)



Proviz

```
0
active                                dist.. pip..    active          sbid          di... .. active    .. d    a.. a    .. a s s s m
[failed decode]                      mov
typeinfo name for dequantize_mul_mat_vec_q4_0_sycl(void const*, float const*, float*, int, int, sycl::V1::queue*)::'lambda'(sycl::V1::nd_item<3>)
/data/projects/iaprof/llama/llama.cpp/ggml/src/ggml-sycl/dmmv.cpp
L0::CommandListCoreFamily<(GFXCORE_FAMILY)3081>::appendLaunchKernelWithParams(L0::Kernel*, _ze_group_count_t const&, L0::Event*, L0::CmdListKernelLaunchParams&)
L0::CommandListCoreFamily<(GFXCORE_FAMILY)3081>::appendLaunchKernel(_ze_kernel_handle_t*, _ze_group_count_t const&, _ze_event_handle_t*, unsigned int, _ze_eve... L0::CommandListFor...
L0::CommandListCoreFamilyImmediate<(GFXCORE_FAMILY)3081>::appendLaunchKernel(_ze_kernel_handle_t*, _ze_group_count_t const&, _ze_event_handle_t*, unsigned int... L0::CommandListFor...
L0::zeCommandListAppendLaunchKernel(_ze_command_list_handle_t*, _ze_kernel_handle_t*, _ze_group_count_t const&, _ze_event_handle_t*, unsigned int, _ze_event.h... L0::CommandListFor...
zeCommandListAppendLaunchKernel
ur_command_list_manager::appendKernelLaunchUnlocked(ur_kernel_handle_t*, unsigned int, unsigned long const*, unsigned long const*, unsigned long const*, unsi... zeC... zeCommandListAppen...
ur_command_list_manager::appendKernelLaunch(ur_kernel_handle_t*, unsigned int, unsigned long const*, unsigned long const*, unsigned long const*, unsigned int... ur... ur_command_list_ma...
v2::ur_queue_immediate_in_order_t::enqueueKernelLaunch(ur_kernel_handle_t*, unsigned int, unsigned long const*, unsigned long const*, unsigned long const*, u... ur... v2::ur_queue_immed...
ur::level_zero::ur_enqueue_kernel_launch(ur_queue_handle_t*, ur_kernel_handle_t*, unsigned int, unsigned long const*, unsigned long const*, unsigned long const... ur::level_zero::ur...
sycl::V1::detail::enqueue_kernel_launch(sycl::V1::detail::queue_impl&, sycl::V1::detail::NDRDescT&, std::vector<sycl::V1::detail::ArgDesc, std::allocator<sycl... urE... ur_enqueue_kernel lau...
sycl::V1::handler::finalize():'lambda'():operator() const::'lambda'():operator()() const
sycl::V1::handler::finalize():'lambda'():operator()() const
sycl::V1::handler::finalize() (.localalias)
std::shared_ptr<sycl::V1::detail::event_impl> sycl::V1::detail::queue_impl::finalizeHandlerInOrderNoEventsUnlocked<sycl::V1::handler>(sycl::V1::handler&)
sycl::V1::detail::queue_impl::submit_impl(sycl::V1::detail::type_eraser_cgo_ty const&, bool, sycl::V1::detail::code_location const&, bool, sycl::V1::deta... std... std::shared_ptr<sy...
sycl::V1::queue::submit_with_event_impl(sycl::V1::detail::type_eraser_cgo_ty const&, sycl::V1::detail::v1::SubmissionInfo const&, sycl::V1::detail::code... syc... sycl::V1::detail::...
sycl::V1::event sycl::V1::queue::submit_with_event<sycl::V1::ext::oneapi::experimental::properties<sycl::V1::ext::oneapi::experimental::detail::properties... syc... sycl::V1::queue::...
ggml_sycl_op_dequantize_mul_mat_vec(ggml_backend_sycl_context&, ggml_tensor const*, ggml_tensor const*, ggml_tensor*, char const*, float const*, char const*, .. syc... sycl::V1::event s...
void ggml_sycl_op_mul_mat<no_quantize_q8_1>(ggml_backend_sycl_context&, ggml_tensor const*, ggml_tensor const*, ggml_tensor*, void (*)(ggml_backend_sycl_conte... ggml... ggml_sycl_cpy(ggml...
ggml_sycl_mul_mat(ggml_backend_sycl_context&, ggml_tensor const*, ggml_tensor const*, ggml_tensor const*)
ggml_backend_sycl_graph_compute_impl(ggml_backend_sycl_context*, ggml_cgraph*)
ggml_backend_sycl_graph_compute(ggml_backend*, ggml_cgraph*)
ggml_backend_sched_graph_compute_async
llama_context::graph_compute(ggml_cgraph*, bool)
llama_context::process_ubatch(llama_ubatch const&, llama_graph_type, llama_memory_context_i*, ggml_status&)
llama_context::decode(llama_batch const&)
llama_decode
main
__libc_init_first
__libc_start_main
tstart
1338469
llama-cli
bll
```

Current Release

GPU profiler is called `iaprof` (Intel Accelerator Profiler)

- Needs frame pointers for full CPU stacktraces
- Overhead is high by default, needs to parse all batchbuffers
- <https://github.com/Intel/iaprof>

Proviz

- <https://github.com/kammerdienerb/proviz>
- Also supports `perf` script output

```
# sudo iaprof record -i 10 1> profile.txt 2> output.txt &  
# [Run workload]  
# sudo kill $(pidof iaprof)  
# proviz report profile.txt
```


The Intel logo is centered on a solid blue background. It features the word "intel" in a white, lowercase, sans-serif font. A small, light blue square is positioned above the letter "i". To the right of the word "intel" is a small white registered trademark symbol (®).

intel®