# Debugging and Analyzing Programs using the Intercept Layer for OpenCL™ Applications

Ben Ashbaugh

IWOCL 2018
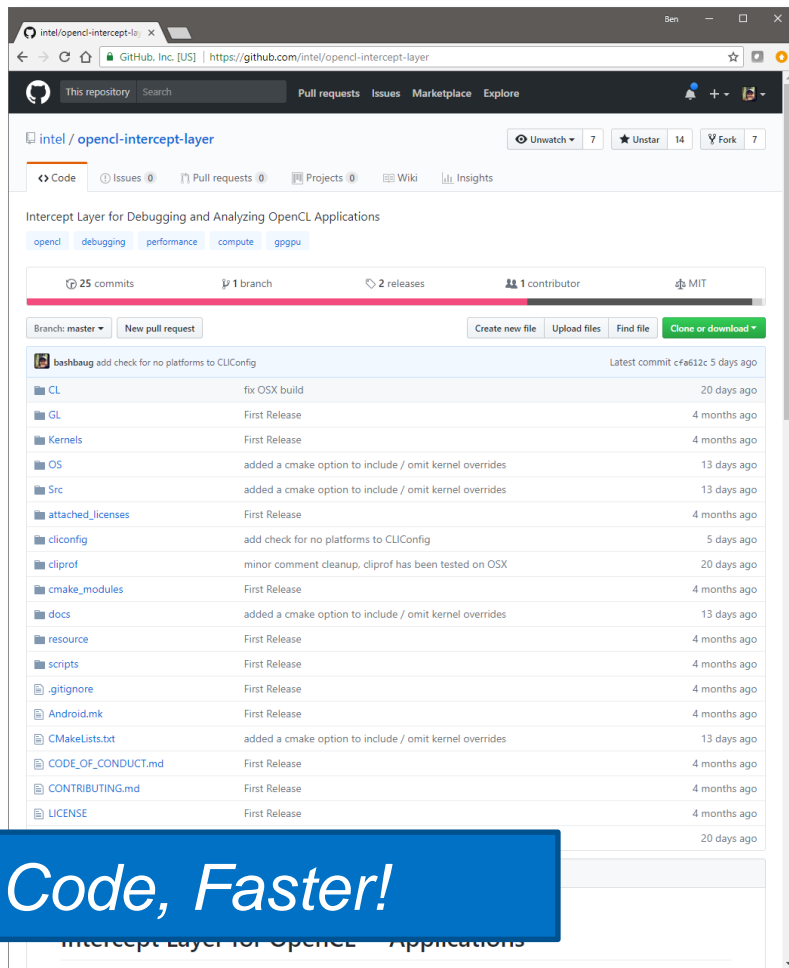
https://github.com/intel/opencl-intercept-layer

# Why am I here?

## Intercept Layer for OpenCL Applications

- Debug and Analyze OpenCL Applications

- Open Source, Permissive License

- Works with Any* OpenCL Implementation

- Requires No Application Modifications

- Thin, Fast, Easy to Install / Uninstall

- Community Contributions are Welcome and Encouraged!



*Develop Fast OpenCL Code, Faster!*

# Agenda

History

How it Works

What it Can Do

Implementation Details

Possible Next Steps

Wrap Up

# History

(2009-Present)

# Initial Requests:

I'm debugging an application.  Can you modify the driver to print the OpenCL APIs that are called?

Yeah, no problem.

One week later…

Can you print the API arguments too?

Sure, I think we can add that.

# More requests:

I'm debugging the GPU compiler.  Can you modify the driver to dump OpenCL kernels to a file?

Yeah, that's not too hard.

Great.

Can you also make it work for the CPU OpenCL implementation?

I'm not sure – I think so?

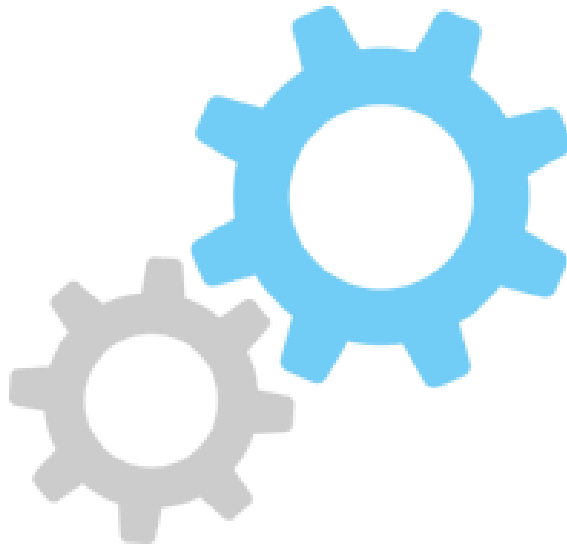Fantastic.  Can you make it work for [third party competitor]?

# Meanwhile:

Our Driver Team was also adding instrumentation:

- Flush or Finish After Enqueue
- Assert on OpenCL Errors
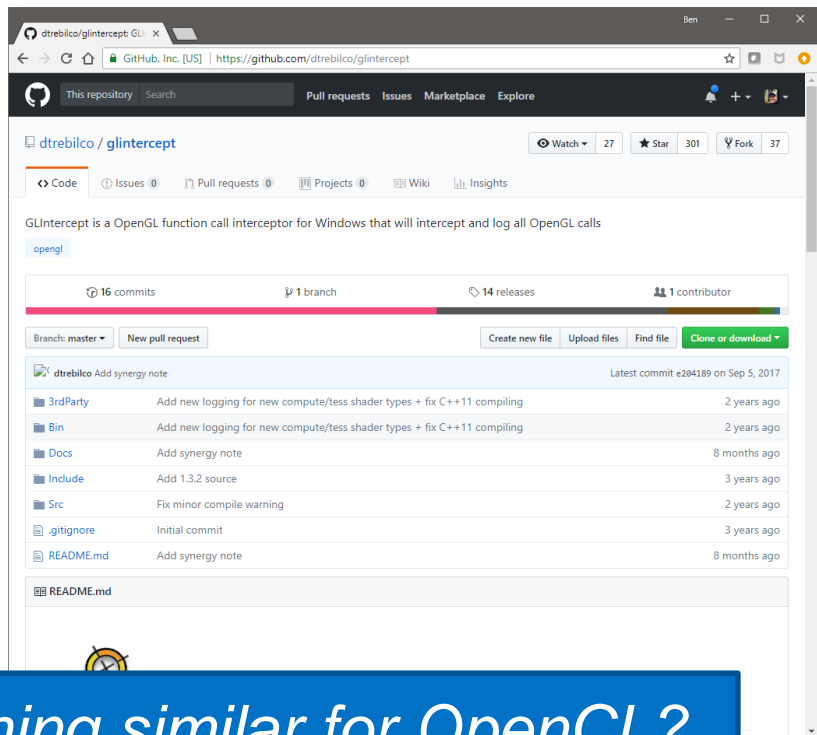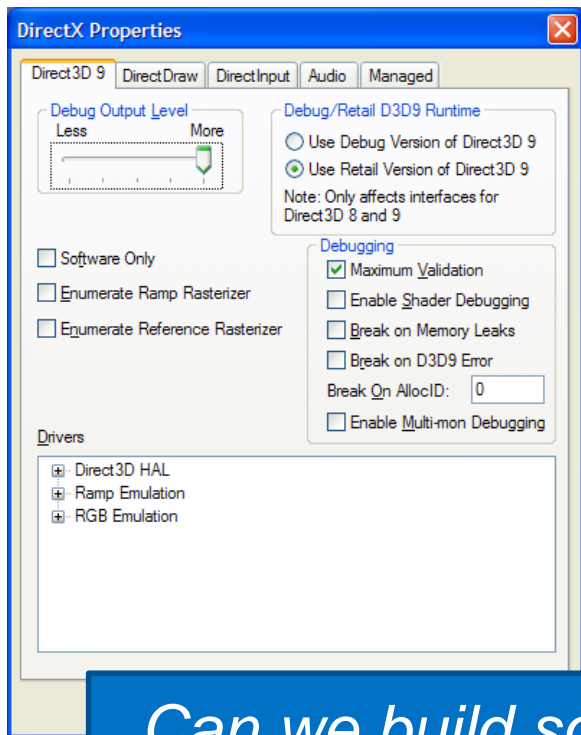- Timing API Calls
- More …

But:

- Required driver modifications!

*Is there a better way to add these capabilities?*

# Prior Work from Graphics APIs:



*Can we build something similar for OpenCL?*
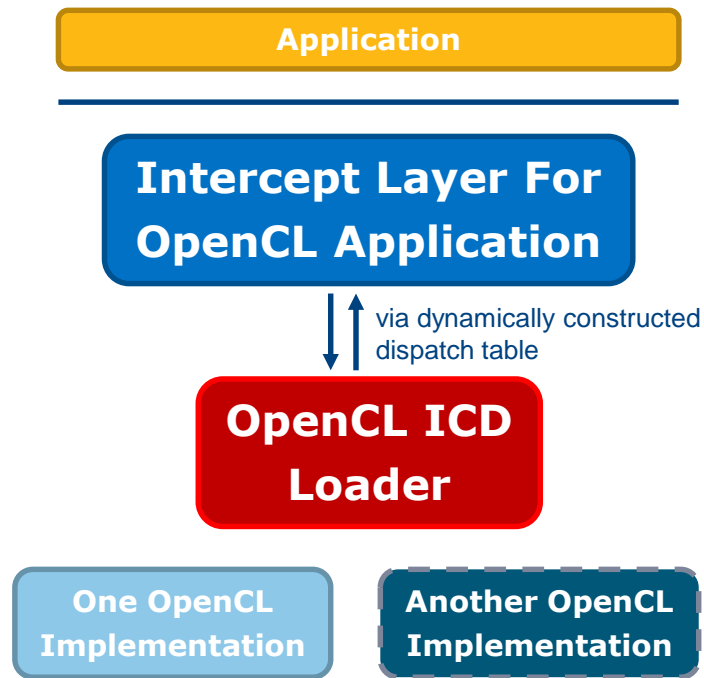
# How It Works

# Intercept Layer for OpenCL Applications

## Architecture: How it Works*:

- Inserts between Application and OpenCL ICD Loader

- Constructs Dispatch Table During Initialization

- Passes Through API Calls... or not!

## Philosophies:

- Focus on Features that Solve Problems
  - For OpenCL Implementers
  - For OpenCL Developers

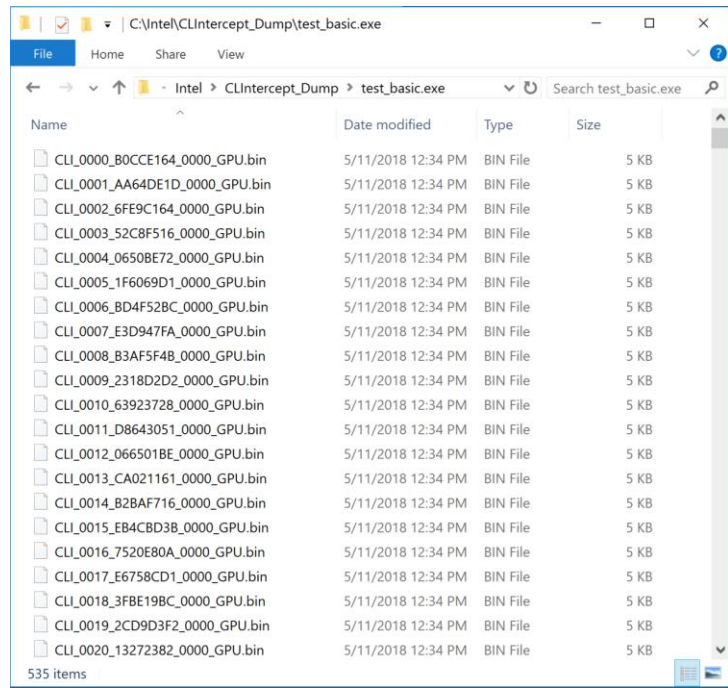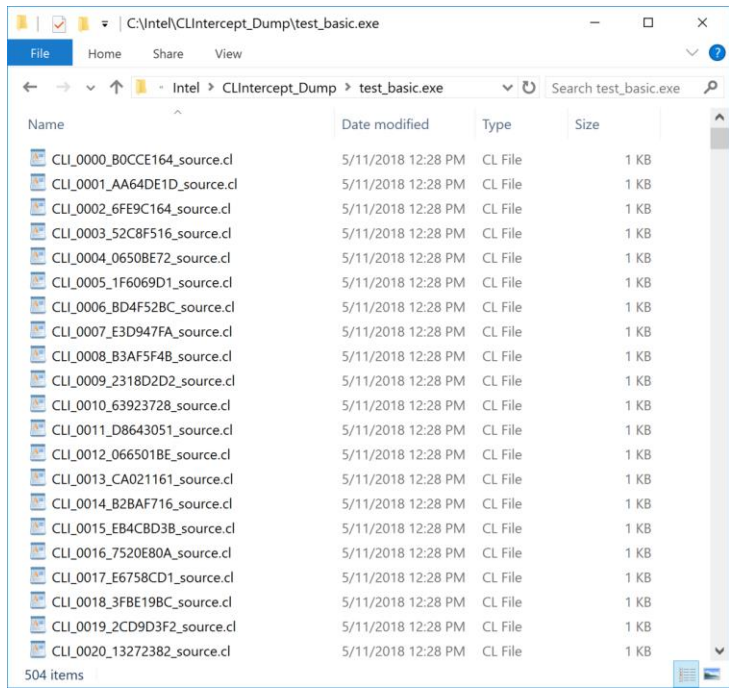- Support Any OpenCL Device on Any Platform

- Be Invisible By Default

**Application**

**Intercept Layer For OpenCL Application**

via dynamically constructed dispatch table

**OpenCL ICD Loader**

**One OpenCL Implementation**

**Another OpenCL Implementation**

* Typical usage on Windows and Linux, OSX is a little different.

# What It Can Do - Examples

# Call and Error Logging

```
>>>> clGetPlatformIDs
<<<< clGetPlatformIDs
>>>> clGetPlatformIDs
<<<< clGetPlatformIDs
>>>> clGetDeviceIDs: platform = [ NVIDIA CUDA ], device_type = CL_DEVICE_TYPE_ALL (FFFFFFFF)
<<<< clGetDeviceIDs
>>>> clGetDeviceIDs: platform = [ Intel(R) OpenCL ], device_type = CL_DEVICE_TYPE_ALL (FFFFFFFF)
<<<< clGetDeviceIDs
>>>> clCreateContextFromType: properties = [ CL_CONTEXT_PLATFORM = Intel(R) OpenCL ], ...
ERROR! clCreateContextFromType returned CL_DEVICE_NOT_FOUND (-1)
<<<< clCreateContextFromType: returned 00000000
>>>> clCreateContextFromType: properties = [ CL_CONTEXT_PLATFORM = Intel(R) OpenCL ], ...
<<<< clCreateContextFromType: returned 00E97068
>>>> clGetContextInfo: param_name = CL_CONTEXT_DEVICES (00001081)
<<<< clGetContextInfo
>>>> clGetContextInfo: param_name = CL_CONTEXT_DEVICES (00001081)
<<<< clGetContextInfo
>>>> clCreateCommandQueue: device = [ Intel(R) Core(TM) i7-2600K CPU @ 3.40GHz (CL_DEVICE_TYPE_CPU) ]
<<<< clCreateCommandQueue: returned 05B038F8
>>>> clGetContextInfo: param_name = CL_CONTEXT_DEVICES (00001081)
<<<< clGetContextInfo
>>>> clGetContextInfo: param_name = CL_CONTEXT_DEVICES (00001081)
<<<< clGetContextInfo
>>>> clCreateProgramWithSource: context = 00E97068, count = 1
<<<< clCreateProgramWithSource: returned 04572EA8, program number = 0000
```

# Dumping Program Source (and Binaries!)



## Can also Modify and/or Inject Modified Program Source or Binaries!

# Host API Performance Timing

```
Host Performance Timing Results:

                               Function Name, Calls, Average (ns),       Min (ns),      Max (ns)
                              clBuildProgram,     3,    711065926,       22172160, 1634864192
                               clCreateBuffer,    23,      2234125,           2113,     36218573
                         clCreateCommandQueue,     1,        25054,          25054,        25054
                              clCreateContext,     1,    123618277,      123618277,    123618277
                              clCreateImage2D,     2,      8600269,        4682137,     12518402
                               clCreateKernel,     6,         7898,           2113,        14489
                     clCreateProgramWithSource,    3,        24551,           4829,        51617
    clEnqueueNDRangeKernel( AdvancePaths ), 18036,        36967,          22941,     61064301
            clEnqueueNDRangeKernel( Init ),     1,      7529273,        7529273,      7529273
clEnqueueNDRangeKernel( InitFrameBuffer ),     1,      1095145,        1095145,      1095145
     clEnqueueNDRangeKernel( Intersect ), 18036,        25952,          15998,     24253177
       clEnqueueNDRangeKernel( Sampler ), 18036,        29856,          15696,       218847
                        clEnqueueReadBuffer,  2288,      3758695,         123158,     10236648
                                     clFinish,     2,      4723341,         717519,      8729163
                                      clFlush, 18036,        31018,          21432,       374003
                               clGetDeviceIDs,     4,         1811,            301,         5735
...
```
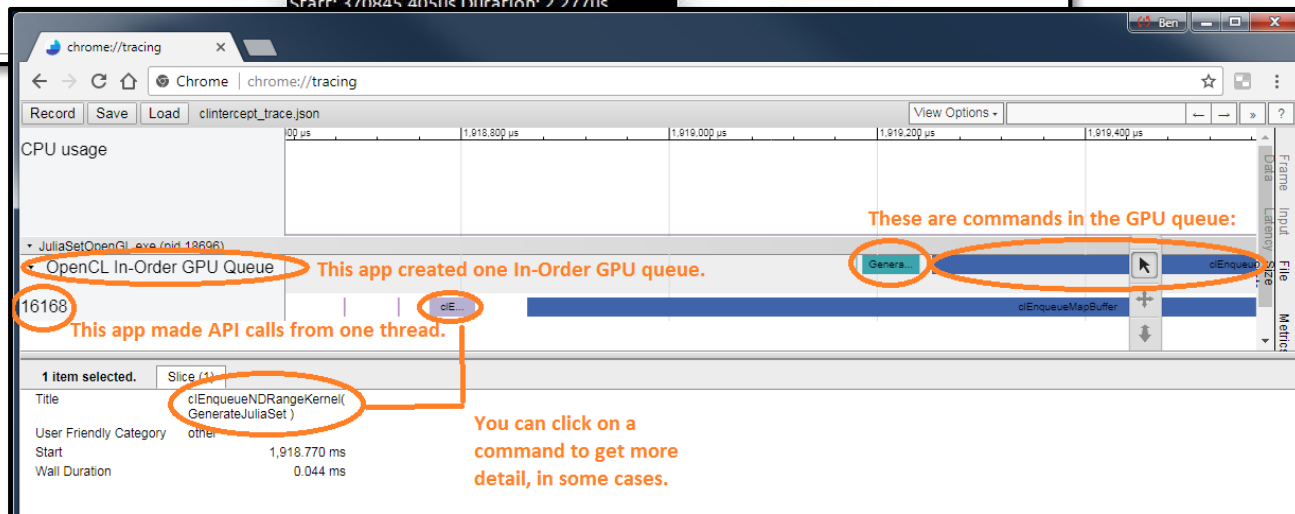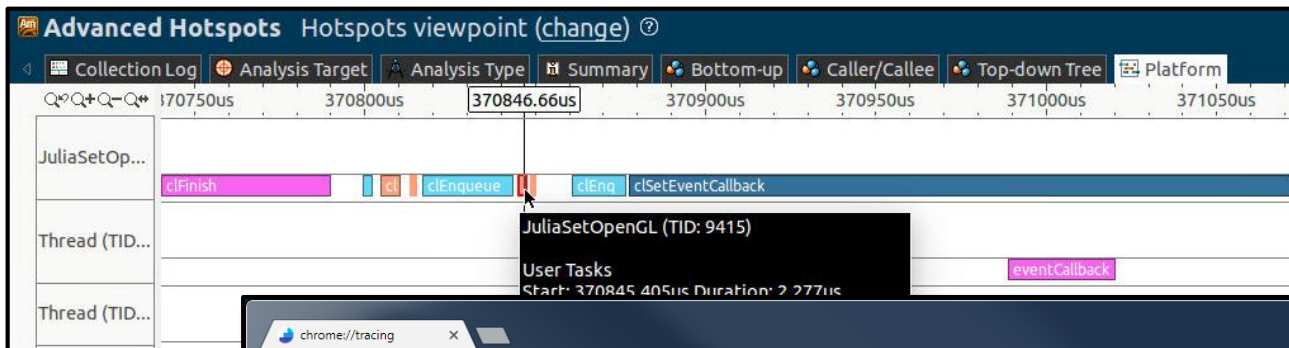
# Device Command Performance Timing

```
Device Performance Timing Results:

Total Time (ns): 123904875200

         Function Name,  Calls,      Time (ns),  Time (%),  Average (ns),  Min (ns),  Max (ns)
          AdvancePaths,  18036, 28203368032,     22.76%,       1563726,   1388096,   1761472
                  Init,      1,     8600000,      0.01%,       8600000,   8600000,   8600000
        InitFrameBuffer,      1,      155712,      0.00%,        155712,    155712,    155712
             Intersect,  18036, 79765237056,     64.38%,       4422556,   3248832,   5297600
               Sampler,  18036, 14307721664,     11.55%,        793286,     75712,   1182400
     clEnqueueReadBuffer,   2288,  1619792736,      1.31%,        707951,     39904,   4220992
```

# Performance Timing on VTune and Chrome*

# Platform and Device Queries



*Explore how applications respond to different query responses!*

# Implementation Details:
# OpenCL API Learnings and Insights

# OpenCL APIs from a Layering Perspective

Most things went really well!

Features that made life easy:

- Built-in Reference Counting and Object Queries
- Standard Event Profiling, Standard Program Binaries
- Online Compilation

Features that made things complicated:

- NULL Local Work Size: Need "what happened" queries!
- Out-of-Order Queues – especially with Device Timing
- Device-side Only Controls (kernel attributes)
    - Easier to permute Host-side Controls (build options)

# Intel-Specific Enhancements

# Intel Specific Enhancement: Driver Diagnostics

cl_intel_driver_diagnostics: Intel Extension

- Extends Context Callback to Include **GOOD** / **BAD** / **INFORMATIONAL** Diagnostics

- Use the Intercept Layer for OpenCL Applications to Enable and Log Diagnostics

```
>>>> clCreateBuffer: flags = CL_MEM_WRITE_ONLY | CL_MEM_ALLOC_HOST_PTR (12), ...
=======> Context Callback (private_info = 00AFF728, cb = 4):
Performance hint: clCreateBuffer needs to allocate memory for buffer. For subsequent
operations the buffer will share the same physical memory with CPU.
<======= End of Context Callback
<<<< clCreateBuffer: returned 0573E000
...
>>>> clEnqueueMapBuffer: [ map count = 0 ] queue = 03254850, buffer = 0573E000, ...
=======> Context Callback (private_info = 00AFF214, cb = 4):
Performance hint: clEnqueueMapBuffer call on a buffer 0573E000 will not require any data
copy as buffer shares the same physical memory with CPU.
<======= End of Context Callback
<<<< clEnqueueMapBuffer: [ map count = 1 ] returned 04702000
```

# Future Work

# Future Work – Short Term

Continue to stay use-case driven: How to find bugs and fix them faster?

- Example: OpenCL Object Leak Detection

Improve Usability and Accessibility
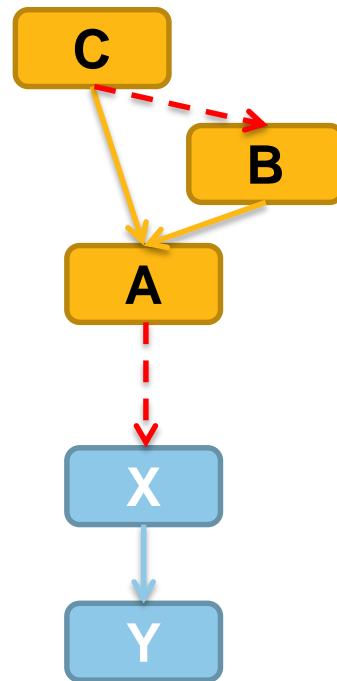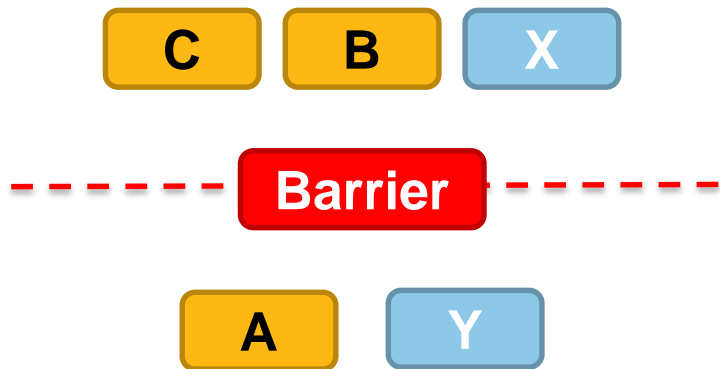
# Future Work – Longer Term

Android Support ☺

Intercept Layer as an ICD: Work with applications that statically link to the ICD loader.

# Future Work – Longer Term

## Log and Analyze Graphs of OpenCL Commands

- Especially Important with Out-of-Order Queues
- Can we plot graphs of commands?
- Can we time device execution of subgraphs?

# Future Work – Longer Term

## Automatic Reproducer Generation

- Very Limited Capture-Playback
- One Kernel + Inputs + Params
- See Fossilize for Vulkan

## Speaking of Vulkan…

- Lots of Layer Prior Art
- Steal with Pride?

# Wrap Up

# Summary and Call to Action

Try the Intercept Layer for OpenCL Applications!

- Debug and Analyze OpenCL programs faster!

- Send Issues and Pull Requests!

Grow the OpenCL Ecosystem with Layers

- Layers are an important part of the OpenCL ecosystem

To the Khronos OpenCL Working Group: Design the API with layers in mind!

To OpenCL Users: Use layers, evangelize layers, build layers!

Thank you!

- ben.ashbaugh@intel.com

# Acknowledgements

Thanks to Michal Mrozek, Michael Carroll, Mike Kinsner, and Adam Herr for reviewing these slides.

Thanks to everyone from Intel who has used or contributed to the Intercept Layer for OpenCL Applications!

# Useful Links:

Intercept Layer for OpenCL Applications:

https://github.com/intel/opencl-intercept-layer

Vulkan Loader and Layers:

https://github.com/KhronosGroup/Vulkan-LoaderAndValidationLayers

# Legal Notice and Disclaimers

Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software or service activation. Learn more at intel.com, or from the OEM or retailer.

No computer system can be absolutely secure.

Tests document performance of components on a particular test, in specific systems. Differences in hardware, software, or configuration will affect actual performance. Consult other sources of information to evaluate performance as you consider your purchase.  For more complete information about performance and benchmark results, visit **http://www.intel.com/benchmarks**.

Intel, the Intel logo and others are trademarks of Intel Corporation in the U.S. and/or other countries. *Other names and brands may be claimed as the property of others.

© 2018 Intel Corporation.

*OpenCL and the OpenCL logo are trademarks of Apple Inc. used by permission by Khronos.*

*\* Other names and brands may be claimed as the property of others.*

# Legal Disclaimer and Optimization Notice