

Project X: Testing GPU Memory Consistency At Large (Experience Paper)

Anonymous Author(s)

ABSTRACT

Memory consistency specifications (MCSs) are a difficult, yet critical, part of a concurrent programming framework. Existing MCS testing tools are not immediately accessible so have only been applied to a limited number of platforms. However, in the post-Dennard scaling landscape there has been an explosion of new architectures and frameworks, especially for GPUs. Studying the shared memory behaviors of different devices (across vendors and architecture generations) is important to ensure conformance and to understand the extent that devices show different behaviors.

In this paper we present Project X, a widescale GPU MCS testing tool. Project X has two interfaces: a web interface and an Android app. Using Project X, we deployed a testing campaign to that checks conformance and characterizes weak behaviors. We advertised our web app on forums and social media, allowing us to collect testing data from 106 devices, spanning seven vendors. In terms of devices tested, this constitutes the largest study on weak memory behaviors by at least 10 \times , and our conformance tests identified two new bugs on embedded Arm and NVIDIA devices. Analyzing our characterization data yields many insights, including quantifying weak behavior occurrence rates (e.g., AMD GPUs show 25.3 \times more weak behaviors on average than Intel) and showing how devices can be clustered according to stress testing sensitivity. We conclude with a discussion on how to further scale our tools and the impact it has on software development for these performance-critical devices.

ACM Reference Format:

Anonymous Author(s). 2023. Project X: Testing GPU Memory Consistency At Large (Experience Paper). In *Proceedings of ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA 2023)*. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

The end of Dennard Scaling has brought about an explosion of multi-core architectures that improve application performance through large-scale parallelism. Graphics Processing Units (GPUs) exemplify this trend and are integral components of many systems, from smartphones to large HPC super computers. While GPUs were previously only used for graphics applications, they now have applications in a variety of areas including machine learning [40] and particle simulations used in drug development [36]. GPUs are even being used for security and safety critical applications such

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
ISSTA 2023, 17-21 July, 2023, Seattle, USA

© 2023 Association for Computing Machinery.
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM... \$15.00
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

Table 1: The GPU vendors and devices included in the study. Overall, we ran almost 400 billion iterations of weak memory tests on 106 devices, of which 58 we have confirmed to be unique models. We observed over 35 million weak behaviors, with the rates per device and vendor characterized in Sec. 4.

Framework	Vendor	Devices (Unique)	Tests	Weak Behaviors
WebGPU	Intel	26 (17)	105.3b	0.2m
	Apple	26 (6)	104.4b	9.7m
	NVIDIA	31 (18)	125.3b	10.8m
	AMD	15 (9)	60.4b	14.7m
Vulkan	Arm	2 (2)	51.6m	18.2k
	Qualcomm	4 (4)	17.6m	27.2k
	PowerVR	1 (1)	6.1m	0
	NVIDIA	1 (1)	49.6m	454
Total:		106 (58)	395.5b	35.4m

as encryption [35] and self-driving cars [12], making correctness an increasing concern on these devices.

Because GPUs are produced by several vendors (NVIDIA, AMD, Intel, etc.) and evolve rapidly, many different devices are currently deployed. Devices vary in their performance and even their functional behavior. To account for this, the community has developed portable GPU programming frameworks, such as Vulkan [21] and WebGPU [49], as a unified abstraction for these diverse devices.

Memory consistency specifications (MCSs), which define the semantics of shared memory operations, are an important part of these abstractions. While MCSs provide many guarantees, such as atomicity and coherence, they often allow an architecture to implement *weak* memory behaviors to improve efficiency [32]. For example, x86's relaxed MCS [42] allows store buffering behaviors, in which a processor may buffer local store operations; in turn, another process may observe the buffered store occurring out-of-order.

Because relaxed MCSs can be complex and nuanced, there is a history of platforms (compilers and architectures) containing MCS conformance bugs [1, 3, 24, 28, 29]. That is, the MCS provides a guarantee that the implementation does not honor. Due to the non-determinism of concurrency, MCS bugs may occur extremely rarely, or only when provoked, e.g., by side-channel stress [44]. Apart from to conformance, a device's weak behavior *profile*, i.e., the frequency at which allowed weak behaviors occur and how system stress influences this frequency, is also a useful metric. For example, this profile aids in developing conformance testing strategies [26] and enables developers to reason about tradeoffs between accuracy and performance in approximate computing that elides occasional synchronization [33, 39, 41].

Unfortunately, previous GPU MCS testing work had limited scale, only testing a few devices [24, 26], with the largest study testing eight [1]. These approaches could not scale widely due to the difficulty of portable GPU application development, e.g., while

frameworks like OpenCL [20] are portable in theory, there are many difficulties in practice [45]. Consequently, little is known about the prevalence of conformance bugs and weak behavior profiles *at large*. This is especially problematic as portable GPU frameworks depend upon many layers and environments (e.g., architectures, compilers, runtimes, operating systems, etc.); it is difficult to extrapolate insights from a small number of platforms tested in controlled environments to the diverse universe of deployed GPUs.

1.1 Project X

In this paper, we present a large-scale study of GPU MCS testing, which, to the best of our knowledge, tests 10× more devices than previous studies. Figure 1 summarizes our study, including the number of GPUs that we tested (106), broken down by two frameworks (WebGPU and Vulkan) and seven vendors (Intel, Apple, NVIDIA, AMD, Arm, Qualcomm, and Imagination). This scale is empowered by Project X¹, a new cross-platform GPU MCS testing tool suite. Project X includes two front-ends, a browser web app (using WebGPU) and an Android app (using Vulkan). We advertised our web app on campus forums and social media to obtain a significant number of WebGPU results. We test far fewer Vulkan devices as our Android app is not yet widely accessible on the Google Play Store, but in Sec. 6 we discuss how we will enable larger mobile studies on both Android and iOS.

Project X uses *litmus tests*, small concurrent programs that check for load/store reordering corresponding to weak memory behaviors. Current GPU MCS testing tools execute litmus tests many times in succession to check conformance and characterize devices [1, 24, 44]. However, these prior approaches have several shortcomings: (1) they are implemented in vendor-specific languages, e.g., CUDA, or (2) they require expert users to build, configure, and execute tests on each device, e.g., as is the case for OpenCL, or (3) litmus tests were embedded in vendor-specific memory stress and thus would not execute efficiently on other devices. This cumbersome litmus testing workflow made it infeasible to perform a large scale study. In contrast, Project X defines litmus tests using a neutral configuration (written in JSON), which it compiles to a portable shading language (WGSL [48] or SPIR-V [19]). The resulting litmus testing application then tunes the testing stress automatically. The net result is a fully automated and easy-to-use tool for GPU MCS testing at large. Table 1 shows how many weak memory litmus test iterations were run and how many weak behaviors were observed in our study.

We perform the following two investigations on our data set: (1) we examine the results of MCS conformance tests and find two new bugs in mobile device GPUs from Arm and NVIDIA, and (2) we characterize weak memory behavior profiles, e.g., the rates at which allowed weak behaviors occur and their sensitivity to system stress. Additionally, we provide several analyses on the weak memory profiles. First, we comment on how per-vendor average profiles compare; for example, AMD shows an average percentage of 1.5% weak behaviors on the tests we characterize, while Intel only shows .06%. We then cluster different GPUs, and find that, surprisingly, cross-vendor devices often have similar profiles, while devices from the *same* vendor sometimes have vastly different profiles. Finally,

¹We do not give the true name of our project to retain anonymity

we discuss how the wide range of different profiles we observed can impact application development and testing strategies.

Contributions. In summary, our contributions are:

- (1) **Tooling:** We introduce Project X, a new cross-platform GPU MCS testing tool with an accessible web and Android interface. Project X enables large-scale litmus tests due to its cross-platform support (Sec. 3).
- (2) **GPU MCS Weak Behavior Characterization:** We conduct the largest-ever GPU weak memory characterization and conformance testing study, collecting data from 106 GPUs (Sec. 4).
- (3) **Conformance Testing and Analysis:** We analyze the conformance testing data and discover two unreported bugs in Arm and NVIDIA devices (Sec. 5.1). We then analyze the rate at which weak behaviors occur and report on statistical similarities across the wide set of GPUs in our study (Sec. 5.2). We discuss how these insights may impact developing applications for these devices (Sec. 5.3).

Once our paper is published we will open source our datasets, along with the code and tools we used to collect them, so that the wider community can perform their own analysis and draw further insights.

2 BACKGROUND

We first provide an overview of memory consistency specifications. Next, we overview litmus tests and how they allow reasoning about relaxed memory models. Finally, we introduce GPU programming concepts from the WebGPU and Vulkan GPU frameworks, including descriptions of their MCSs.

2.1 Memory Consistency Specifications

Today, the MCSs for both architectures like x86 [42] and languages like C++ [6] are formalized using mathematical logic, representing shared memory program executions as a set of memory operations, e.g., reads, writes, and read-modify-writes, and relations between these events, e.g., happens-before (*hb*). *Hb* is an acyclic relation constraining the allowed behaviors of a program. The strongest MCS is sequential consistency (SC) [25], which states that concurrent program executions must correspond to a total *hb* order such that the order respects the per-thread program order, allowing events from multiple threads to be interleaved. In relaxed MCSs, the *hb* relation is a partial order, allowing various weak behaviors (i.e. executions that are *not* SC) if shared memory operations on multiple threads are not synchronized.

There is a large body of work focused on formalizing relaxed MCSs, including a model of Vulkan’s MCS [18]. WebGPU currently defines a non-normative model, with prior work [26] formalizing portions of its MCS necessary for describing the semantics of simple litmus tests. However, for this work it is not necessary to understand the full formalization of the WebGPU and Vulkan MCSs, so we describe the necessary subset of the specification briefly and informally. In addition, we follow prior work on MCS testing [24, 26] and consider only *trivially data-race-free* programs where all operations are atomic, as our intention is not to test the behavior of programs with undefined semantics.

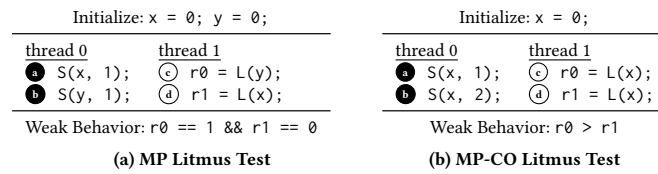


Figure 1: The weak behavior in the MP litmus test is allowed by many relaxed MCSs. On the other hand, the weak behavior in the closely-related MP-CO litmus test violates coherency and is disallowed by every major MCS that we know of. Prior work [26] has shown that tuning system stress for the allowed MP test can be used to design effective tests for finding bugs related to the MP-CO test.

Our Target MCS. Because Vulkan is one of several backends to WebGPU, the MCS for WebGPU is a subset of the MCS for Vulkan. In order to provide a unified study across both frameworks, we target only the WebGPU MCS, which we then map to its Vulkan counterpart. The WebGPU MCS provides very little inter-workgroup synchronization due to the diversity of backends it targets, with the weakest backend being Apple’s Metal [5], which provides only relaxed atomic operations.

The one property provided by WebGPU atomics is *coherence*, which states that memory accesses to a single location must respect sequential consistency; this property is also known as SC-per-loc in some works [3]. However, memory accesses to disjoint addresses are allowed to be reordered. Mapping these WebGPU atomics to Vulkan is straightforward; all WebGPU atomic accesses are simply mapped to SPIR-V atomic accesses with a relaxed memory order. While our testing campaign considers only relaxed memory accesses, Vulkan allows additional memory orders; specifically, acquire and release. While the precise semantics of these memory orders is complex, especially when combined with other relaxed atomics, we note that they are required to implement the synchronization required by a mutex². The lock() method needs to execute an acquire atomic operation when the mutex is obtained and the unlock() method requires executing a release atomic operation. If a mutex is implemented without these memory orders, it is possible to violate mutual exclusion, as we show in Sec. 5.3.

2.2 Litmus Tests

Litmus tests are small concurrent programs which are used both to illustrate and compare MCSs [27, 43, 46] and to empirically test MCS implementations [1, 2, 24]. These tests contain a final condition on local variables and memory values that check for a weak behavior. For example, the program in Fig. 1a is known as the message passing (MP) litmus test, in which one thread writes to a memory location x (Ⓐ) followed by a write to y (Ⓑ), while a second thread reads from y (Ⓒ) and then x (Ⓓ).

As mentioned earlier, in this work, we assume that all of the memory operations in a litmus tests are *atomic*, which in languages

²Because WebGPU does not provide inter-workgroup acquire and release memory orders, it is not currently possible to implement a well-specified mutex in that framework

that follow the C11 style MCS [6] ensures that the semantics of shared memory operations are well-defined. Additionally, unless explicitly noted otherwise, we consider these atomic operations to have a relaxed memory order, which allows compilers and hardware to aggressively optimize their execution.

The condition underneath the test shows an outcome that only occurs in relaxed executions. In this case, the behavior corresponds to an execution in where the read of y returns 1 but the read of x returns 0. While some relaxed MCSs do not allow this behavior, e.g., the x86 MCS [42], many other relaxed MCSs, especially ones for high level languages (e.g., C++ [6]), do allow the behavior. As mentioned earlier, our target WebGPU MCS does not provide any guarantees outside of coherence, and thus the two memory accesses per thread (which target disjoint addresses) can be reordered. In cases where the weak behavior is allowed (both by the MCS and the implementation), the rate at which this behavior is observed on real systems is highly dependent on *system stress*. Early GPU application development work did not observe any weak behaviors, despite specifications allowing them [11]. However, later work added specialized system stress around the the test execution, and revealed many cases of surprising weak behaviors [1, 44].

Executing litmus tests on deployed systems can be used for two purposes, which we will illustrate using a litmus test L that can exhibit a weak behavior e , and an MCS S .

- (1) **Conformance testing:** if e is disallowed on S then we can check implementations of S . That is, if a platform p claims to implement S , then we can execute L many times on p , checking for e . If e is ever observed, then there is a platform MCS bug.
- (2) **Profiling weak behaviors:** if e is *allowed* on S , and a platform p claims to implement S , then we can execute L many times on p to understand the extent to which that platform allows e . In some cases, p might not show e empirically, or maybe e appears more frequently under a certain configuration of system stress. A collection of this type of data creates a weak memory profile for p .

Prior work [26] has utilized weak memory profiles in highly tuned conformance testing. In that work, it was shown that allowed MP executions could be used to tune system stress for disallowed behaviors in associated conformance tests. For example, the MP-CO litmus test, shown in Fig. 1b, is similar to MP, except that every memory access targets the same memory address and different values are stored, so that a weak behavior can be identified. Given that there is only one address used in MP-CO, the weak behavior in this test disallowed under coherence, and thus in the WebGPU MCS. Prior work showed that if certain system stress revealed weak behaviors in the allowed MP litmus test, then, in the case where a platform contained a bug, it was likely to reveal the buggy behavior in the MP-CO conformance test. In Sec. 5.1 we show how we used the same methodology to find two bugs in mobile GPUs, while we discuss all the litmus tests used in our large experimental campaign in Sec. 3.1.

2.3 GPU Programming

This study targets two cross-platform GPU frameworks, Vulkan and WebGPU. Vulkan is a modern graphics and compute API that

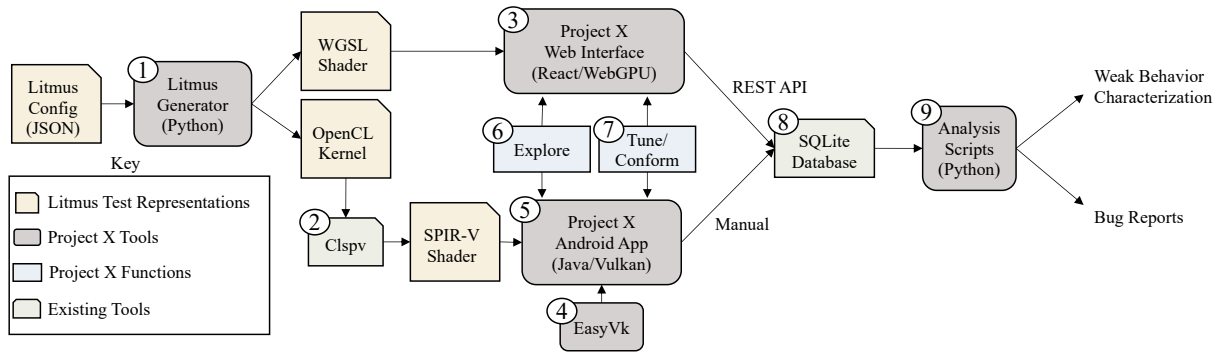


Figure 2: An overview of our tooling for testing the WebGPU and Vulkan MCS.

can be run on many Linux, Android, and Windows devices, and can target Apple devices through the MoltenVK [22] portability layer. WebGPU is designed to run in browser environments, and is compiled to different backends depending on the operating system of the device (Direct3D [31] on Windows, Vulkan on Linux/Android, and Metal [5] on Apple devices).

Both Vulkan and WebGPU define their own programming languages, called SPIR-V and WGSL respectively. Programs written in these languages are called *shaders* and run on the GPU, while the APIs used to allocate memory on the GPU and dispatch shaders are written in the language of the host device, commonly C++ for Vulkan and JavaScript for WebGPU. In this work, we discuss the complexities of writing tools that must be implemented in different languages and how future development (Sec. 6) could ease the difficulty of cross platform GPU MCS testing.

GPU Execution Model. GPUs run thousands of concurrent threads (*invocations* in Vulkan and WebGPU) organized hierarchically and executed in a single-instruction, multiple-thread (SIMT) format. To support this execution model, in WGSL and SPIR-V threads are partitioned into discrete *workgroups*, with built in identifiers used to query a thread’s workgroup id. Workgroups are limited in size (e.g. 1024 in CUDA, with limits varying depending on the device in WGSL/SPIR-V) and have access to an efficient shared memory region. A group of threads organized into workgroups and running on the device is called a *grid*, with the number of threads per workgroup and the number of workgroups specified at dispatch time. All threads in the same dispatch have access to a global memory region.

While our target MCS was discussed in the previous section, we note that GPU atomic operations can be annotated with a *memory scope*, that indicates how wide their synchronization occurs. Two common scopes in Vulkan and WebGPU are *workgroup*, which specifies the synchronization occurs only between threads in the same workgroup, and *device*, which specifies that the synchronization occurs across all threads executing on the device.

Threads within workgroups generally have access to efficient primitive barrier operations, e.g., `workgroupBarrier` in WebGPU. However, highly optimized implementations of important parallel routines (e.g. inter-workgroup prefix scans [30]) rely on fine-grained inter-workgroup communication. Thus, like prior work [24, 26], we see a more imminent need for testing MCS properties at

the inter-workgroup level; which we keep as our sole scope for this work. Similarly, GPU programs have several different memory types, e.g. whether it is shared-workgroup memory or device-wide memory. Given that we consider only inter-workgroup interactions, we can only consider device-wide memory.

3 SYSTEM OVERVIEW

Building on approaches in prior work [26], we now discuss: (1) our testing campaign (Sec. 3.1), and (2) the development of our MCS testing tools that are easily accessible on a wide range of devices, which is summarized in Figure 2. We overview each stage of the tooling, starting with litmus test generation (Sec. 3.2), moving on to the design of Project X’s web interface and Android app (Sec. 3.3). We end the section by describing our data collection process (Sec. 3.4).

3.1 Litmus Test Selection

The tests we utilize in our study build off of the MCS mutation testing strategy used in [26]. We use 32 *mutants*, out of which 24 are litmus tests with weak behaviors allowed by the WebGPU MCS. The mutants are used to find effective system stress to then run the conformance tests. Our results analysis focuses on characterizing the rates of weak behaviors of six of the mutants, one of which is **MP** (Fig. 1a), with the other five shown in Fig. 3. These tests enumerate all the combinations of four instructions on two threads that can lead to weak behaviors. Thus, they capture testing for all pair-wise memory reorderings. For example, the **SB** test checks for store-load reorderings, while the **LB** test checks for load-store reorderings. Additionally, these tests capture synchronization patterns used in common concurrency algorithms like a compare-and-swap lock. Because of this, prior work has also focused on these tests and has shown their utility in finding bugs in both applications and MCS implementations [24, 44].

Once the mutants are run, we use the weak behavior profile of a device to determine an effective system stress configuration to run conformance tests under. We utilize the 20 conformance tests from [26]. As a concrete illustration using one mutant and conformance test, we would run the **MP** test under many different system stress configurations to build a weak behavior profile. We then use the most effective configuration at revealing **MP** weak behaviors to run a closely related conformance test, e.g., **MP-CO** (Fig. 1b). This approach was shown to be effective in prior work [26]

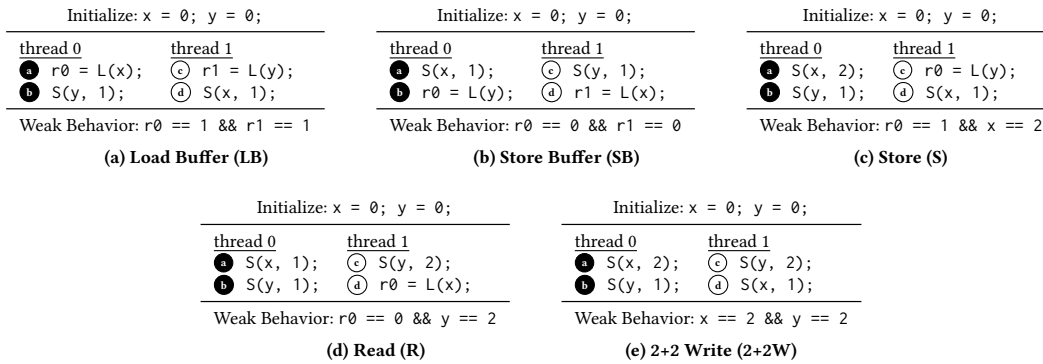


Figure 3: These litmus tests, along with MP from Fig. 1a, represent six classic weak behaviors allowed by relaxed MCSs. S and L signify a relaxed atomic store and load, respectively.

and this work continues to show its effectiveness by discovering new bugs: a violation of MP-CO on Arm devices and a violation of MP-CO on an NVIDIA device (see Sec. 5.1).

3.2 Litmus Test Generation

We now discuss our tooling efforts to generate and run our testing and characterization campaign. Litmus test behaviors are non-deterministic and highly sensitive to system stress. Due to this, the shaders that run the litmus tests contain not only the actual litmus test instructions, like those in Fig. 3, but take in a number of parameters and define functions that are used to construct the system stress around the test.

To provide a standardized interface for defining litmus tests in different GPU languages, we built a tool, LITGEN (1 in Fig. 2), which is similar to previous litmus testing tools [2] but is specifically targeted to create GPU programs with system stress, as was shown is necessary for testing GPU MCSs [1, 24, 26]. Litmus tests can then be written in an abstract format, currently structured using JSON, that specify the actions of the test (e.g. loads and stores) and the possible behaviors of the test, including weak behaviors. LITGEN then outputs a *test shader*, which runs the test along side system stress developed in prior work [24, 26], and a *result shader*, which aggregates the behaviors of the test.

The result shader is generated separately from the test shader for several reasons:

- (1) Some tests, like 2+2W (Fig. 3e), need to examine the memory locations themselves for weak behaviors. To avoid this examination interfering with the test process itself, along with relying on some of the same memory synchronization features that we are trying to test, we pass the buffer used in the test into the result shader where the values can be examined.
- (2) We implemented the parallel testing strategy described in [26] to run hundreds thousands of *instances* of each litmus test concurrently. Thus it is only natural to leverage the inherent parallelism of the GPU to aggregate the results, which otherwise may be time consuming to do on the host CPU, especially since it may require explicitly copying memory from the GPU to the host.

Currently, two backends exist for LITGEN. The tool outputs WGSL shaders directly, as WGSL is a text based language. SPIR-V, on the other hand, is a low-level representation similar to LLVM, increasing its flexibility but making code generation more complex. Therefore, for Vulkan backends LITGEN first outputs OpenCL, another compute focused GPU language which is similar in syntax to C++. Then, it utilizes Clspv [14] (2), a prototype compiler from OpenCL to SPIR-V, to generate the shader used in the Android app.

As WebGPU is primarily browser based while Vulkan runs on native devices, we currently maintain litmus testing driver programs in two languages. WebGPU exposes a relatively simple JavaScript API which we build our web interface (3) on. In this interface, we include features that increase the system stress such as passing buffers to the shader that map built-in thread and workgroup identifiers to a new, randomized identifier. Vulkan’s native C/C++ API is more complex, so to simplify this process, we have built a Vulkan compute wrapper which we call EASYVK (3). This library exposes a simple interface where buffers are defined and shaders dispatched to the GPU in a few lines of code. While EASYVK is only used for our MCS testing purposes currently, in the we believe it may have utility as a library for quick deployment of Vulkan compute shaders and we will make it publicly available along with the rest of our tooling upon publication.

3.3 Project X Design

Previous GPU MCS studies have been limited in reach due to the difficulty in deploying cross platform GPU applications [45]. One of the reasons for this is the fractured landscape of GPU development. NVIDIA’s popular CUDA framework [34] is used for many data science applications but is only fully supported on NVIDIA GPUs. OpenCL was introduced in 2009 by Apple and quickly adopted as a potential cross platform standard, with prior work [24] using it to evaluate MCS testing strategies. However, today Apple no longer supports OpenCL, and instead requires developers to use their proprietary Metal API.

Another issue with previous GPU MCS testing approaches has been a reliance on expert users to run testing campaigns. For example, using OpenCL would require users to install the required drivers, build the application under a specific environment, etc. To

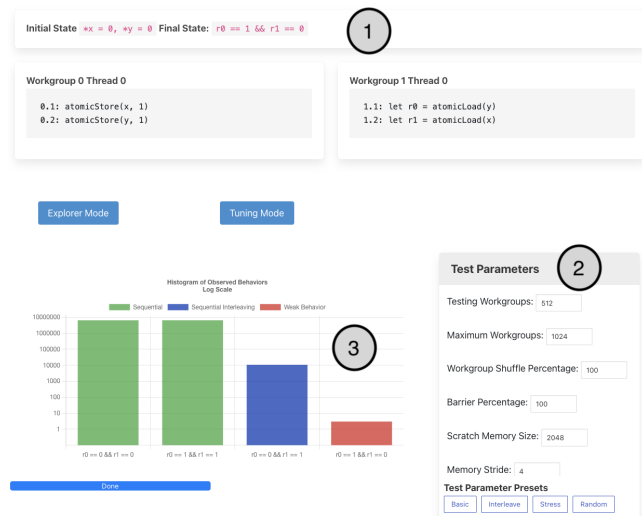


Figure 4: A screenshot of Project X’s web interface showing the “Explore” page for the MP litmus test.

collect data from the diversity of devices necessary to gain confidence in cross platform GPU frameworks, tools must minimize the friction in setting up and running tests. The tools we introduce here are easily distributed applications with a user friendly interface on top of new GPGPU frameworks so that even non-technical users can collect data on their GPU MCS and submit results for analysis. Specifically, we introduce Project X for MCS GPU testing tools, which has two frontends (③ and ⑤) in Fig. 2):

- **Project X Web Interface (③):** A website that runs MCS GPU tests using WebGPU. The web interface includes pages for exploring tests, useful as a pedagogical example and for building intuition on different MCS behaviors, and a page for tuning system stress to increase the rate of weak behaviors and search for bugs in conformance tests.
- **Project X Android App (⑤):** An app that runs MCS GPU tests using Vulkan, which is intended to be supported by all Android devices. While this study includes the largest collection of data on mobile GPU MCS behaviors, in Sec. 6 we discuss future work that could increase the reach of mobile GPU MCS testing even further.

Both Project X’s web interface and Android app have a common design with two functions: exploring (⑥) and tuning/conforming (⑦). Explore pages run specific litmus tests, displays histograms of results, and can adjust various parameters that control system stress. When tuning and conforming, a set of tests are chosen to run with multiple random system stress configurations, searching for configurations that maximize the rate of weak behaviors and discover bugs in MCS implementations.

Exploring. Figure 4 shows a screenshot of Project X’s web interface explore page for the MP litmus test after the test has been run with relatively high systems stress on a MacBook Pro with an integrated Intel Iris GPU. The top of the page (①) includes a description of the test and pseudocode showing the test instructions. The right hand side (②) includes an editable list of the parameters that define

system stress, along with several presets. When the test is running, the histogram (③) updates in real time with the number of times each behavior is observed. The progress bar gives an estimate on how much longer is left to run, based on the speed of previous iterations.

The green bars correspond to sequential behaviors, where one thread runs entirely before the other. The blue bar corresponds to interleaved behaviors, where actions from each thread are interleaved (e.g. leading to the behavior $r0 == 0 \ \&\& \ r1 == 1$ in the MP litmus test). The red bar corresponds to weak behaviors; in this run, three MP weak behaviors were observed out of over 13 million test instances, so the histogram shows behaviors using a log scale.

Tuning and Conforming. Both the web interface and the Android app can be used to tune system stress, as in [26]. When tuning, a set of tests can be selected, with presets available for weak memory tests (e.g. those in Fig. 3) and conformance tests, e.g., to test coherence. Options like the number of configurations, the maximum number of workgroups, and other parameter overrides can be modified to run different experiments and check specific tests without redeploying any code.

To collect data from a diverse set of devices, we want to minimize the options users have to configure, reducing the chances of errors and giving us a standardized dataset to analyze. For this study, this was especially important on Project X’s web interface, as we collected data from many devices and needed to ensure we did not have to manually clean large datasets. The web interface’s tuning page therefore includes a tab that exposes no configuration options, but instead only a few buttons, one that starts a combined tuning/conformance run with default parameters and another that pulls up a form to fill out some (optional) contact information and submit the results of the run. Our results are all anonymized; contact details were only collected if users wanted to be informed about the outcome of the study. Before submitting, users agreed that their anonymized results could be aggregated, reported on, and released as part of this study.

3.4 Data Collection

To submit test data, the web interface communicates with a backend service that exposes an API for submitting results and inserting them into a SQLite database (⑧ in Fig. 2). The data is then analyzed using Python scripts (⑨). The Android app is not yet available on the app store nor is it integrated with the SQLite backend, so results are manually copied off of the device for analysis. In Sec. 6 we discuss how we can increase the scale of future studies, but nevertheless our study of eight devices is the largest testing campaign of mobile GPU MCS behaviors that we are aware of.

While system stress configurations are generated randomly, we would like to ensure that the configurations run on different devices are the same for data analysis purposes. That is, if different GPUs are tested with the same stress configurations, we can compare how the different devices behaved under the same stress. We ensure this by integrating a seedable Park-Miller random number generator [37] into both the web interface and the Android app and using the same seed when running all of our tuning experiments.

By default, browsers only expose limited information about the user’s GPU without turning on vendor-specific development flags

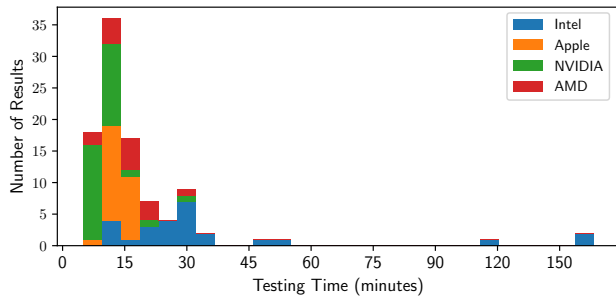


Figure 5: Histogram showing the time spent testing WebGPU’s MCS on each device using Project X’s web interface. Each bin is broken down by vendor.

due to privacy and security concerns around fingerprinting [50]. To increase the accuracy of our results, we included instructions asking users to temporarily enable flags so we could collect detailed GPU information; of the 98 results we collected, 67 included the exact GPU tested. The other 31 results which did not specify the exact GPU included only the vendor and a string describing the GPU architecture, such as “intel gen-9” or “nvidia ampere”. All Apple devices reported an architecture of “common-3”, making it impossible to immediately distinguish M1’s vs M2’s. However, we show in Sec. 5.2 that data collected in the same system stress configurations can be used to infer device information, affecting the ability of browsers to hide the specifics of a user’s GPU.

4 INITIAL RESULTS: WEAK BEHAVIOR CHARACTERIZATION

To collect data from as many sources as possible, we disseminated the link to Project X’s web interface to the general public, utilizing campus forums and social media, and ran the Android app on eight devices which we could physically access. As shown in Tab. 1, we collected data from millions of tests; each test used a randomly generated system stress configuration (we used 50 configurations on the web interface and 150 on the Android app). In each configuration tests were run millions of times based on a randomly generated number of workgroups and threads per workgroup.

To ensure data integrity, we implemented a checksum algorithm that verified we saw the expected number of overall behaviors based on the system stress configuration. The testing duration was also recorded on a per-test and system stress configuration basis, but we ran into one issue here; some computers went to sleep in the middle of the tests, suspending the browser’s process and leading to extremely long recorded test times. To profile testing time, we did not include any test/configuration durations over one minute; as each individual test runs quickly (e.g. in less than 5 seconds), these durations were most likely when the computer went to sleep. Then, to approximate the length of the test that was suspended, we used a neighboring test’s time.

One consideration for collecting data from the wider public is that we cannot afford to run tests for hours at a time. Previous work targeted only a few devices, running tests on one device for a minimum of 36 hours [24] or 2 hours [26], but asking users

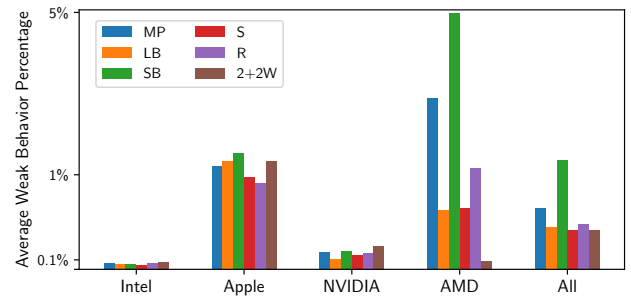


Figure 6: Data showing the average rate of weak behaviors across vendors when running litmus tests using WebGPU.

to leave their browser and computer open for that long would almost certainly decrease the number of submissions. Therefore, we heuristically chose the number of test environments and iterations per environment, aiming for the tests to finish in 10-20 minutes.

Figure 5 shows the distribution of testing time on our web interface, broken down by vendor. The results show that NVIDIA devices were the fastest on average, mostly running all tests in under 15 minutes. On the other hand, Intel devices ran slower, with two older Intel GPUs taking over an hour and a half to complete.

We now analyze our WebGPU and Vulkan data to characterize the rates at which weak behaviors occur on devices from different vendors.

4.1 Weak Behaviors in WebGPU

Figure 6 shows the average rates of observed weak behaviors for the six litmus tests of Fig. 3 (plus MP) in the test environment that maximizes the rate on each device broken down by test and vendor. As described in Fig. 1, we have data from at least 15 devices from each vendor. The overall testing time across all 98 devices was 31.1 hours, an average of 19 minutes per device.

Devices from all vendors showed weak behaviors on each litmus test. In all but two cases, observing weak behaviors was all or nothing; if a device revealed weak behaviors on one litmus test, it revealed weak behaviors on all of them. In contrast, on a device implementing x86’s TSO MCS, we would expect to only see store buffering behaviors. However, unlike x86, GPU devices do not provide low-level details, such as the hardware-level MCS, thus it was not clear what types of weak behaviors we would observe. These results show that many GPUs implement very relaxed memory models, in contrast to stronger architectures like x86 TSO.

Intel devices tended to have the lowest rate of weak behaviors, with just over half of them (15/26) revealing weak behaviors on each test. The median rate of weak behaviors on Intel devices was even lower than their average, around .02% for each test. No Intel device showed a rate of weak behaviors above 1% on any test.

NVIDIA devices revealed weak behaviors at a relatively low rate. Our results include results from NVIDIA’s Kepler (2012), Maxwell (2014), Pascal (2016), Turing (2018), and Ampere (2020) architectures, with a majority being the more recent Ampere. Older devices generally showed fewer weak behaviors, with the minimum on each of the six tests being Kepler and Maxwell devices. However,

Table 2: Data showing the average rate of weak behaviors across vendors when running litmus tests using Vulkan.

Vendor	Device	Litmus Test					
		MP	LB	SB	S	R	2+2W
Qualcomm	Adreno 610	0%	0%	0%	0%	0%	0%
	Adreno 640	0%	2.04%	1.45%	1.65%	0%	2.21%
	Adreno 642L	0.04%	5.75%	5.81%	3.81%	0%	6.38%
	Adreno 660	0.12%	8.5%	14.37%	5.69%	0%	11.5%
Arm	Mali G71	0.04%	0%	0%	0%	0%	0%
	Mali G78	1.56%	0%	0%	0%	0%	0%
Imagination	PowerVR GE8320	0%	0%	0%	0%	0%	0%
NVIDIA	Tegra X1	0.01%	0.05%	0.02%	0.01%	0.01%	0.05%

one outlier is that the maximum rate of **SB** behaviors (.73%) was seen on a Kepler device. Interestingly, that device was also the only device not to observe any weak behaviors on **S**, **LB**, and **2+2W**. The only other device not to reveal weak behaviors on a test was a Quadro K620 with a Maxwell architecture, on **MP**, **R**, and **SB**.

Apple devices were consistently weak, revealing weak behaviors on every device and test, generally at a higher rate on all tests than NVIDIA devices but with less variation than AMD devices. Apple GPUs have only been recently built into non-mobile devices, so these results represent the first comprehensive evaluation of the weak behaviors on Apple GPUs. We don't have the specific name of every Apple device, but we were able to collect enough information to show we had results from Apple M1 (basic, Pro, Max) and Apple M2 (basic, Pro) devices.

AMD devices were also very weak, with 100% of devices showing weak behaviors on every test. The clear highest average rate occurs on the **SB** litmus test on AMD GPUs. Most of the AMD devices show a high rate of weak behaviors on **SB**, approaching 10% and higher, but devices with AMD's Graphics Core Next 5 micro-architecture all showed rates under 1%. This means that even from a single vendor, the behaviors of different architectures can vary widely and past results from one vendor cannot be counted on to predict future behaviors.

4.2 Weak Behaviors in Vulkan

The data in Tab. 2 shows the percentage of weak behaviors in the test environment that maximizes the rate at which they occur for our android devices. In contrast to our web GPUs, in the mobile setting weak behaviors were observed in every test on only one device, the NVIDIA Tegra X1, but the rates on this device were very low, beneath 0.1%. The most difficult test to observe in general was **R**, which checks whether a store is re-ordered with a following load on one thread. We did not observe any weak behaviors on the Imagination GPU; because testing is fundamentally incomplete, this could mean that the device implements a strong MCS, or that our testing approach was not effective. Interestingly, ARM only showed weak behaviors in the **MP** test.

We observe that, in general, the rates of weak behaviors increase as devices become more powerful. This is especially apparent from the four Qualcomm devices we test, as the rate of weak behaviors increases from 0% on the Adreno 610 (which has 96 *shading units*, analogous to NVIDIA's CUDA cores) up to a maximum of 14.37%

in **SB** on the Adreno 660 (with 512 shading units). One intuitive explanation for this might be that smaller GPUs lack the ability to schedule as many threads at once, naturally reducing the rates of weak behaviors despite architectures which might allow them. We see a similar trend on the Arm Mali GPUs, where the smaller G71 (32 shading units) showed a lower rate of weak behaviors than the larger G78 (384 shading units).

5 INSIGHTS AND IMPACTS

We now use our data and characterization of weak behavior rates to answer the following research questions:

- (1) Do bugs exist in the wild, especially in GPUs which are relatively untested (Sec. 5.1)?
- (2) Can we use our data to identify unknown devices, create browser fingerprints, and choose application testing strategies (Sec. 5.2)?
- (3) What are the impacts of weak behaviors on applications that require synchronization, and can our characterization of weak behaviors be used as a programming guide (Sec. 5.3)?

5.1 Bugs on Mobile Devices Running Vulkan

When analyzing results from running conformance tests, we discovered coherency violations in three mobile GPUs from two vendors: an Arm Mali G71, Arm Mali G78, and NVIDIA Tegra X1. While our tested NVIDIA device was on a Shield device (as it runs Android), we note that this SoC is the same one on the popular Nintendo Switch gaming console (however, the switch does not easily run Android). We found the bugs by running the **MP-CO** conformance litmus test under tuned system stress. We choose the system stress using the methodology described in [26], i.e., by selecting a system stress configuration effective at revealing weak behaviors in the **MP** litmus test. While discovering the bugs confirms that tuning system stress configurations using weak memory litmus tests is a good strategy, we perform the same correlation analysis as [26] to empirically validate the methodology.

We run both the **MP** and **MP-CO** litmus tests in 150 randomly generated system stress configurations on each device, recording the rate of weak behaviors in **MP** and the rate of buggy (i.e. non-coherent) behaviors in **MP-CO**. Our results show that the Pearson Correlation Coefficient (PCC) between the rate of weak and buggy behaviors is 0.732 on the Arm Mali G71, 0.759 on the Arm Mali G78, and 0.832 on the NVIDIA Tegra X1. Since these behaviors are recorded from 150 samples (i.e. system stress configurations), we have 148 degrees of freedom, and running a Student's t-test leads to a p-value less than $10^{-5}\%$ on each device. This shows that the PCC between weak behaviors and bugs is certainly not due to random chance, further validating that configurations tuned using weak behaviors are effective at revealing bugs in conformance tests.

5.2 GPU Similarity

All of the data was collected by running the tests with pseudo-randomly generated system stress configurations, but as mentioned in Sec. 3.4 the generator is seeded with a known value. The similarity of different testing runs can then be calculated by comparing the behaviors of different devices. Each testing run is represented as a vector of the non-sequential (i.e. where one thread runs entirely

Table 3: Each row shows the cosine similarity statistics between all pairs of devices from that vendor. The last row shows the similarity statistics across all pairs of devices.

Vendor	Avg	Median	Min	Max
Intel	0.87	0.891	0.683	0.985
Apple	0.903	0.913	0.699	0.993
NVIDIA	0.903	0.931	0.67	0.996
AMD	0.904	0.927	0.661	0.989
All	0.84	0.862	0.477	0.996

Table 4: Device clustering shows that choosing one device from each vendor is not an optimal way to test applications that require shared memory operations.

Device	Cluster					
	A	B	C	D	E	F
Intel	3	0	0	5	18	0
Apple	12	4	0	10	0	0
NVIDIA	1	0	19	8	2	1
AMD	13	1	0	1	0	0

before another one) behaviors of every test in each configuration. Ignoring the sequential behaviors gives the data a degree of freedom, which is necessary for calculating a valid similarity measure.

For our similarity metric, we choose *cosine similarity*, which measures the cosine of the angle between two vectors and ranges from -1 to 1. We chose cosine similarity because it is a relative metric, not an absolute like Euclidean distance, meaning that devices which show different absolute rates of behaviors, but at a similar relativity, are classified as more closely related.

Device Identification. Table 3 shows a summary of the similarity between devices in our study. All similarities are positive, with the minimum being 0.477 between an Intel and AMD device. This is not surprising, since effective system stress is likely to reveal weak behaviors on many devices. However, the average and median similarities between devices from each vendor is higher than the overall average and medians, showing that in general devices from the same vendor tend to have more similar MCS behaviors.

For Apple and NVIDIA, we confirmed that the maximum similarity occurs between identical GPUs: two Apple M1 Max’s and two NVIDIA GeForce RTX 3080’s. For AMD, we observe a maximum similarity of 0.989 between two devices, one of which is a Radeon Pro 5500M (A) while the other device (B) did not report a model and instead only indicated that it was from the same architectural generation as A. However, we observed a high similarity (0.985) between A and another Radeon Pro 5500M (C), as well as a similarity of 0.984 between B and C, so it seems likely that A, B, and C are all the same device. We do a similar analysis with Intel to determine that an unknown device is most likely an Intel Xeon Graphics. Given these results, we conclude that GPU MCS behavior data is can be used to fingerprinting GPU devices, despite the specification hiding this information for security reasons.

Clustering Based Testing Strategies. K-means clustering attempts to minimize the *distortion*, or the sum of squared distances between vectors and a centroid. Applying k-means clustering to GPU MCS

Table 5: Analysis of three locking algorithms, Test-and-Set (TAS), Test and Test-and-Set (TTAS), and Compare-and-Swap (CAS), showing in how many test runs (out of 1000) we observed failures of unfenced (UF) and fenced (F) lock implementations to protect a critical section. The total time to run all tests on each device is also recorded.

Device	Time (min)	TAS		TTAS		CAS	
		UF	F	UF	F	UF	F
Qualcomm Adreno 610	3.5	0	0	0	0	0	0
Arm Mali G78	67.9	18	0	11	0	7	0

behavior has implications for testing strategies; when developing cross-platform GPU applications that rely on shared memory operations, testing these applications on a number of devices can increase the confidence in the correctness of the implementation. A naive strategy might be to choose one device from each major vendor, but our results show that this is not necessarily an optimal strategy.

Table 4 shows the result of running k-means with six clusters on the similarity data from Tab. 3. The “elbow-method” heuristic showed that the rate of decrease in distortion leveled off at 6 clusters on our data. The clustering data shows that devices from the same vendor are generally placed into the same cluster, but there are outliers in each case. The only NVIDIA Kepler device in our study was dissimilar enough from other devices that it was placed in its own cluster. Kepler is also the oldest NVIDIA architecture in our study, showing that special testing attention might be needed when supporting older devices in cross platform GPU frameworks.

When selecting which devices to test a shared memory GPU application on, the strategy should be to first choose a number of clusters based on the rate of decrease in distortion, and then select at least one device from each cluster. Using our data, this might mean selecting an AMD device from A, Apple devices from clusters B and D, NVIDIA devices from clusters C and F (the only device in that cluster), and an Intel device from cluster E. Despite choosing more Apple and NVIDIA devices than AMD and Intel, the similarity data ensures the tests maximally cover devices with different behavior profiles.

5.3 Writing Correct Locking Algorithms

We now discuss a use case of how the diversity of weak memory profiles across these different GPUs can impact software development. Locking algorithms are implemented using atomic operations to synchronize access to critical sections. Implementation of locks depends on careful placement of memory fences to avoid compilers and hardware from re-ordering memory accesses, which can cause critical section failures. In this section, we implemented three common spin-locks: test-and-set (TAS), test-and-test-and-set (TTAS) and compare-and-swap (CAS). Each of these locks specifically need to disallow MP behaviors using acquire/release memory fences. However, our results in Sec. 4 show that on some mobile devices MP weak behaviors never occur, meaning that if the locks are tested on these devices they may run correctly despite being incorrectly implemented.

To investigate this, we tested our three locks on two Android devices, an Arm Mali G78 and a Qualcomm Adreno 610. The locks were implemented with and without appropriate acquire/release memory fences. In these tests, threads from different workgroups acquire the lock 10k times and increment a non-atomic memory location in the critical section. We ran this test for 1k iterations and recorded the number of critical section violations we observed for each device and each lock.

On the Arm Mali G78, a larger GPU which exhibits a relatively high rate of **MP** behaviors, we observed critical section failures on unfenced versions all three locks; in every failure case except one the value was 189,999 instead of 190,000, meaning that just one of the increments was not reflected. In the remaining failure case, the value was 189,998. On the Qualcomm Adreno 610, which exhibited no **MP** behaviors in our study, we saw no failures. Both devices exhibited no failures when locks were run with the correct fences.

To confirm that locks on these devices are actually necessary, we also ran a test where the same threads incremented a non-atomic memory location 10k times for 1k iterations without acquiring any locks. As expected, this racy program led to very incorrect results, with 90% of the increments not reflected on the Mali G78, and 68% not reflected on the Adreno 610. Therefore, when writing applications that require synchronization, careful attention must be paid to make sure tests are run on devices where incorrect implementations will lead to failures, highlighting the importance of collecting and characterizing MCS behavior data.

6 FUTURE WORK

This work has spent significant engineering effort enabling the testing of many different GPUs. However, given the difficulty of cross GPU programming, we were still unable to test mobile Apple GPUs, which appear in some of the most widely used mobile devices. Additionally, our web interface and Android app contain disjoint user interfaces and GPU setup code, causing duplicate efforts and maintenance. In this section, we outline a path forward, with Flutter as a fitting match for these goals.

Flutter [16] is an open-source software development kit developed by Google that provides deployment options to desktop platforms (such as Windows, macOS, and Linux), mobile platforms (Android, iOS), and even web deployment from a single frontend codebase. With a unified codebase for the MCS testing frontend, development work can be focused on designing backend implementations specific to those platforms. Underlying Flutter is Dart [15], a language also developed by Google for cross-platform app development. For each supported platform, Flutter provides an interface to backend code native to the specific platform. On the Android end, GPU access is provided through Dart’s foreign function interface (FFI) library to load a dynamically linked C library, compiled against the version of Vulkan provided by Android’s Native Development Kit (NDK) [13]. The Dart FFI library can be used similarly on all supported platforms except for web, for which GPU access will involve calls to JavaScript code utilizing WebGPU.

Vulkan, while well-supported on Windows, Linux, and Android devices, is not officially supported by macOS and iOS clients. For these platforms, there are two possible options. For a more native-friendly option, Vulkan backend code could be instead rewritten to

depend on Apple’s Metal [5] API, with SPIR-V shaders transpiled to the Metal Shading Language (MSL) using SPIRV-Cross [23], a tool developed by Khronos Group. However, to reduce development time and duplicate code across multiple platforms, Vulkan backend code can be passed through MoltenVK. MoltenVK [22] is a Khronos Group implementation of a large subset of Vulkan 1.2 on top of Metal, and provides a portability layer with which to run Vulkan applications on iOS and macOS platforms.

We also plan on integrating our new tools with the current web interface backend, allowing us to collect data from devices we do not have physical access using a simple API interface. With a single source for interface design, GPU setup, and data collection, it is expected that future work will be able to deploy MCS testing at a wider scale, and collect results from GPU hardware previously inaccessible in related work.

7 RELATED WORK

Testing MCSs. Work on testing MCS dates back to tools like ARCHTEST [47] and TSOTool [17], which each generated test programs containing sequences of loads and stores and then looked for violations of sequential consistency. With the introduction of formal MCSs, researchers developed tools like LITMUS [2], which runs litmus tests generated from formal models directly on ISAs (namely x86, Power, and Arm) and includes stress parameters that make weak behaviors more likely.

Techniques for CPU MCS testing have been extended to GPUs [1, 24]. Weak behaviors on GPUs are notoriously difficult to reveal, leading to work that statistically analyzed tuning techniques and reproducibility of results when running litmus tests on GPUs [24]. To better evaluate the efficacy of test environments and provide confidence in MCS implementations, [26] introduced a methodology based on black-box mutation testing [7], finding bugs in several WebGPU MCS implementations.

Previous studies have been limited in the number of devices they were able to test. In contrast, this study introduces tooling that allows us to conduct the largest ever GPU MCS testing campaign, running tests across 2 frameworks, 7 vendors, and 106 devices.

Testing at Scale. Other studies have tested large numbers of devices, searching for bugs in compilers and hardware. In [10], 17 GPU and driver combinations were tested for compiler bugs. Our approach, distributing the GPU MCS testing experiment using a web interface, is a form of *volunteer computing*, where the general public volunteers their computing resources for research studies. Volunteer computing has been used for many compute-intensive tasks, including searching for extraterrestrial life [4], training neural networks [9], sequencing genomes [38], and climate modeling [8].

8 CONCLUSION

We introduce Project X, a tool suite with a web interface and Android app for cross platform GPU MCS testing. We utilize Project X to perform a large-scale study on weak behaviors in 106 GPUs from seven vendors and find two bugs in GPUs running on mobile devices. Our results show the importance of scaling previous MCS testing strategies in order to characterize the behavior of different devices, perform conformance testing, and design application testing strategies.

REFERENCES

- [1] Jade Alglave, Mark Batty, Alastair F. Donaldson, Ganesh Gopalakrishnan, Jeroen Ketema, Daniel Poetzl, Tyler Sorensen, and John Wickerson. 2015. GPU concurrency: Weak behaviours and programming assumptions. In *International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS) (ASPLOS '15)*. Association for Computing Machinery, 577–591. <https://doi.org/10.1145/2694344.2694391>
- [2] Jade Alglave, Luc Maranget, Susmit Sarkar, and Peter Sewell. 2011. Litmus: Running tests against hardware. In *International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, Vol. 6605. 41–44. https://doi.org/10.1007/978-3-642-19835-9_5
- [3] Jade Alglave, Luc Maranget, and Michael Tautschnig. 2014. Herding cats: Modelling, simulation, testing, and data mining for weak memory. *Trans. Program. Lang. Syst. (TOPLAS)* 36, 2, Article 7 (July 2014), 74 pages. <https://doi.org/10.1145/2627752>
- [4] David P. Anderson, Jeff Cobb, Eric Korpela, Matt Lebofsky, and Dan Werthimer. 2002. SETI@home: An experiment in public-resource computing. *Commun. ACM* 45, 11 (nov 2002), 56–61. <https://doi.org/10.1145/581571.581573>
- [5] Apple. 2023. Metal. <https://developer.apple.com/documentation/metal/>. Retrieved February 2023.
- [6] Mark Batty, Scott Owens, Susmit Sarkar, Peter Sewell, and Tjark Weber. 2011. Mathematizing C++ concurrency. In *Symposium on Principles of Programming Languages (POPL) (POPL '11)*. Association for Computing Machinery, 55–66. <https://doi.org/10.1145/1926385.1926394>
- [7] Timothy A. Budd and Ajei S. Gopal. 1985. Program testing by specification mutation. *Computer Languages* 10, 1 (1985), 63–73. [https://doi.org/10.1016/0096-0551\(85\)90011-6](https://doi.org/10.1016/0096-0551(85)90011-6)
- [8] C. Christensen, T. Aina, and D. Stainforth. 2005. The challenge of volunteer computing with lengthy climate model simulations. In *First International Conference on e-Science and Grid Computing (e-Science'05)*. 8 pp.–15. <https://doi.org/10.1109/E-SCIENCE.2005.76>
- [9] Travis Desell. 2017. Large scale evolution of convolutional neural networks using volunteer computing. In *Proceedings of the Genetic and Evolutionary Computation Conference Companion (GECCO '17)*. Association for Computing Machinery, 127–128. <https://doi.org/10.1145/3067695.3076002>
- [10] Alastair F. Donaldson, Hugues Evrard, Andrei Lascu, and Paul Thomson. 2017. Automated testing of graphics shader compilers. *Proc. ACM Program. Lang.* 1, OOPSLA, Article 93 (oct 2017), 29 pages. <https://doi.org/10.1145/3133917>
- [11] Wu-chun Feng and Shucai Xiao. 2010. To GPU synchronize or not GPU synchronize?. In *2010 IEEE International Symposium on Circuits and Systems (ISCAS)*. 3801–3804. <https://doi.org/10.1109/ISCAS.2010.5537722>
- [12] Esther Francis. 2014. Autonomous cars: no longer just science fiction. (2014).
- [13] Google. 2023. Android NDK. <https://developer.android.com/ndk>.
- [14] Google. 2023. Clspv. <https://github.com/google/clspv>.
- [15] Google. 2023. Dart. <https://dart.dev/>.
- [16] Google. 2023. Flutter. <https://flutter.dev/>.
- [17] S. Hangal, D. Vahia, C. Manovit, J.-Y.J. Lu, and S. Narayanan. 2004. TSOtool: A program for verifying memory systems using the memory consistency model. In *International Symposium on Computer Architecture (ISCA)*, 2004. 114–123. <https://doi.org/10.1109/ISCA.2004.1310768>
- [18] Jeff Bolz. 2022. Vulkan memory model. <https://www.khronos.org/registry/vulkan/specs/1.1-extensions/html/vkspec.html#memory-model>.
- [19] Khronos Group. 2021. SPIR-V specification version 1.6, revision 1. <https://www.khronos.org/registry/SPIR-V/specs/unified1/SPIRV.html>.
- [20] Khronos Group. 2022. The OpenCL C Specification. https://registry.khronos.org/OpenCL/specs/3.0-unified/html/OpenCL_C.html.
- [21] Khronos Group. 2022. Vulkan 1.3 Core API.
- [22] Khronos Group. 2023. MoltenVK. <https://github.com/KhronosGroup/MoltenVK>.
- [23] Khronos Group. 2023. SPIRV-Cross. <https://github.com/KhronosGroup/SPIRV-Cross>.
- [24] Jake Kirkham, Tyler Sorensen, Esin Tureci, and Margaret Martonosi. 2020. Foundations of empirical memory consistency testing. *Proc. ACM Program. Lang.* 4, OOPSLA, Article 226 (Nov. 2020), 29 pages. <https://doi.org/10.1145/3428294>
- [25] Leslie Lamport. 1978. Time, clocks, and the ordering of events in a distributed system. *Commun. ACM* 21, 7 (July 1978), 558–565. <https://doi.org/10.1145/359545.359563>
- [26] Reese Levine, Tianhao Guo, Mingun Cho, Alan Baker, Raph Levien, David Neto, Andrew Quinn, and Tyler Sorensen. 2023. MC mutants: Evaluating and improving testing for memory consistency specifications. In *Proceedings of the 28th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2 (ASPLOS 2023)*. Association for Computing Machinery, 473–488. <https://doi.org/10.1145/3575693.3575750>
- [27] Sela Mador-Haim, Rajeev Alur, and Milo M K. Martin. 2010. Generating litmus tests for contrasting memory consistency models. In *Proceedings of the 22nd International Conference on Computer Aided Verification (CAV'10)*. Springer-Verlag, 273–287. https://doi.org/10.1007/978-3-642-14295-6_26
- [28] Yatin A. Manerkar, Daniel Lustig, Margaret Martonosi, and Michael Pellauer. 2017. RTLcheck: Verifying the memory consistency of RTL designs. In *Proceedings of the 50th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO-50 '17)*. Association for Computing Machinery, 463–476. <https://doi.org/10.1145/3123939.3124536>
- [29] Yatin A. Manerkar, Caroline Trippel, Daniel Lustig, Michael Pellauer, and Margaret Martonosi. 2016. Counterexamples and proof loophole for the C/C++ to POWER and ARMv7 trailing-sync compiler mappings. arXiv:1611.01507 2016.
- [30] Duane Merrill and Michael Garland. 2016. Single-pass parallel prefix scan with decoupled lookback. https://research.nvidia.com/publication/2016-03_single-pass-parallel-prefix-scan-decoupled-look-back
- [31] Microsoft. 2020. Programming guide for Direct3D 11. <https://docs.microsoft.com/en-us/windows/win32/direct3d11/dx-graphics-overviews>.
- [32] Vijay Nagarajan, Daniel J. Sorin, Mark D. Hill, David A. Wood, and Natalie Enright Jerger. 2020. *A primer on memory consistency and cache coherence* (2nd ed.). Morgan & Claypool Publishers.
- [33] Feng Niu, Benjamin Recht, Christopher Re, and Stephen J. Wright. 2011. HOG-WILD! A lock-free approach to parallelizing stochastic gradient descent. In *Proceedings of the 24th International Conference on Neural Information Processing Systems (NIPS'11)*. Curran Associates Inc., 693–701.
- [34] NVIDIA. 2023. CUDA C++ programming guide. <https://docs.nvidia.com/cuda/cuda-c-programming-guide/>.
- [35] Özgün Özerk, Can Elgezen, Ahmet Can Mert, Erdinç Öztürk, and Erkan Savaş. 2022. Efficient number theoretic transform implementation on GPU for homomorphic encryption. *J. Supercomput.* 78, 2 (feb 2022), 2840–2872. <https://doi.org/10.1007/s11227-021-03980-5>
- [36] Mohit Pandey, Michael Fernandez, Francesco Gentile, Olexandr Isayev, Alexander Tropsha, Abraham C Stern, and Artem Cherkasov. 2022. The transformational role of GPU computing and deep learning in drug discovery. *Nature Machine Intelligence* 4, 3 (2022), 211–221.
- [37] S. K. Park and K. W. Miller. 1988. Random number generators: Good ones are hard to find. *Commun. ACM* 31, 10 (oct 1988), 1192–1201. <https://doi.org/10.1145/63039.63042>
- [38] S. Pellicer, N. Ahmed, Yi Pan, and Yao Zheng. 2005. Gene sequence alignment on a public computing platform. In *2005 International Conference on Parallel Processing Workshops (ICPPW'05)*. 95–102. <https://doi.org/10.1109/ICPPW.2005.35>
- [39] Lakshminarayanan Renganarayanan, Vijayalakshmi Srinivasan, Ravi Nair, and Daniel Prener. 2012. Programming with relaxed synchronization. In *Proceedings of the 2012 ACM Workshop on Relaxing Synchronization for Multicore and Manycore Scalability (RACES '12)*. Association for Computing Machinery, 41–50. <https://doi.org/10.1145/2414729.2414737>
- [40] Albert Reuther, Peter Michaleas, Michael Jones, Vijay Gadepally, Siddharth Samsi, and Jeremy Kepner. 2019. Survey and benchmarking of machine learning accelerators. In *2019 IEEE High Performance Extreme Computing Conference (HPEC)*. 1–9. <https://doi.org/10.1109/HPEC.2019.8916327>
- [41] Mehrzad Samadi, Janghaeng Lee, D. Anoushe Jamshidi, Amir Hormati, and Scott Mahlke. 2013. SAGE: Self-tuning approximation for graphics engines. In *Proceedings of the 46th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO-46)*. Association for Computing Machinery, 13–24. <https://doi.org/10.1145/2540708.2540711>
- [42] Susmit Sarkar, Peter Sewell, Francesco Zappa Nardelli, Scott Owens, Tom Ridge, Thomas Braibant, Magnus O. Myreen, and Jade Alglave. 2009. The semantics of x86-CC multiprocessor machine code. In *Proceedings of the 36th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '09)*. Association for Computing Machinery, 379–391. <https://doi.org/10.1145/1480881.1480929>
- [43] Dennis Shasha and Marc Snir. 1988. Efficient and correct execution of parallel programs that share memory. *ACM Trans. Program. Lang. Syst.* 10, 2 (April 1988), 282–312. <https://doi.org/10.1145/42190.42277>
- [44] Tyler Sorensen and Alastair F. Donaldson. 2016. Exposing errors related to weak memory in GPU applications. In *Proceedings of the 37th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '16)*. Association for Computing Machinery, 100–113. <https://doi.org/10.1145/2908080.2908114>
- [45] Tyler Sorensen and Alastair F. Donaldson. 2016. The hitchhiker's guide to cross-platform OpenCL application development. In *Proceedings of the 4th International Workshop on OpenCL (IWOCCL '16)*. Association for Computing Machinery, Article 2, 12 pages. <https://doi.org/10.1145/2909437.2909440>
- [46] John Wickerson, Mark Batty, Tyler Sorensen, and George A. Constantinides. 2017. Automatically comparing memory consistency models. In *Proceedings of the 44th ACM SIGPLAN Symposium on Principles of Programming Languages (POPL '17)*. Association for Computing Machinery, 190–204. <https://doi.org/10.1145/3009837.3009838>
- [47] William W. Collier. 1994. ARCHTEST. <http://www.mpdia.com/archtest.html>.
- [48] World Wide Web Consortium (W3C). 2022. WebGPU shading language: Editor's draft. <https://gpuweb.github.io/gpuweb/wgsl/>.
- [49] World Wide Web Consortium (W3C). 2023. WebGPU: W3C working draft. <https://www.w3.org/TR/webgpu/>.

1277	[50] World Wide Web Consortium (W3C). 2023. WebGPU: W3C working draft: Privacy considerations. https://www.w3.org/TR/webgpu/#privacy-considerations .	1335
1278		1336
1279		1337
1280		1338
1281		1339
1282		1340
1283		1341
1284		1342
1285		1343
1286		1344
1287		1345
1288		1346
1289		1347
1290		1348
1291		1349
1292		1350
1293		1351
1294		1352
1295		1353
1296		1354
1297		1355
1298		1356
1299		1357
1300		1358
1301		1359
1302		1360
1303		1361
1304		1362
1305		1363
1306		1364
1307		1365
1308		1366
1309		1367
1310		1368
1311		1369
1312		1370
1313		1371
1314		1372
1315		1373
1316		1374
1317		1375
1318		1376
1319		1377
1320		1378
1321		1379
1322		1380
1323		1381
1324		1382
1325		1383
1326		1384
1327		1385
1328		1386
1329		1387
1330		1388
1331		1389
1332		1390
1333		1391
1334		1392