

1 Putting Randomized Compiler Testing into 2 Production

3 Alastair F. Donaldson 

4 Google, UK

5 Imperial College London, UK

6 afdx@google.com

7 Hugues Evrard

8 Google, UK

9 hevrard@google.com

10 Paul Thomson

11 Google, UK

12 paulthomson@google.com

13 — Abstract —

14 We describe our experience over the last 18 months on a compiler testing technology transfer project:
15 taking the GraphicsFuzz research project on randomized metamorphic testing of graphics shader
16 compilers, and building the necessary tooling around it to provide a highly automated process
17 for improving the Khronos Vulkan Conformance Test Suite (CTS) with test cases that expose
18 fuzzer-found compiler bugs, or that plug gaps in test coverage. We present this tooling for test
19 automation—`gfauto`—in detail, as well as our use of differential coverage and test case reduction
20 as a method for automatically synthesizing tests that fill coverage gaps. We explain the value
21 that GraphicsFuzz has provided in automatically testing the ecosystem of tools for transforming,
22 optimizing and validating Vulkan shaders, and the challenges faced when testing a tool ecosystem
23 rather than a single tool. We discuss practical issues associated with putting automated metamorphic
24 testing into production, related to test case validity, bug de-duplication and floating-point precision,
25 and provide illustrative examples of bugs found during our work.

26 **2012 ACM Subject Classification** Software and its engineering → Compilers; Software and its
27 engineering → Software testing and debugging

28 **Keywords and phrases** Compilers, metamorphic testing, 3D graphics, experience report

29 **Digital Object Identifier** 10.4230/LIPIcs.ECOOP.2020.22

30 **Category** Experience Report

31 **1** Introduction

32 Graphics processing units (GPUs) provide hardware-accelerated graphics in many scenarios,
33 such as 3D and 2D games, applications, web browsers, and operating system user interfaces.
34 To utilize GPUs, developers must use a graphics programming API, such as OpenGL,
35 Direct3D, Metal or Vulkan, and write *shader programs* that execute on the GPU in an
36 embarrassingly-parallel manner. Shaders are written in a *shading language* such as GLSL,
37 HLSL, MetalSL, or SPIR-V (associated with the OpenGL, Direct3D, Metal and Vulkan
38 APIs, respectively), and are usually portable enough to run on many different GPU models.
39 A graphics driver contains one or more *shader compilers* to translate shaders from portable
40 shading languages to machine code specific to the system's GPU.

41 Functional defects in graphics shader compilers can have serious consequences. Clearly,
42 as with any bug in any application, it is undesirable if a mis-compiled graphics shader causes
43 unintended visual effects. Furthermore, since shaders are compiled at *runtime* (because
44 the GPU and driver that will be present when an application executes is not known at the



© Alastair F. Donaldson, Hugues Evrard and Paul Thomson;
licensed under Creative Commons License CC-BY

34th European Conference on Object-Oriented Programming (ECOOP 2020).

Editors: Robert Hirschfeld and Tobias Pape; Article No. 22; pp. 22:1–22:30

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

45 application’s compile time), a shader compiler crash can lead to an overall application crash.
46 Additionally, developers cannot feasibly test for or workaround these issues, as the driver
47 version that crashes may not have even been written at the time of application development.
48 Worse still, because the graphics driver usually drives the whole system’s display, if a shader
49 compiler defect leads to the state of the driver being corrupted, the entire system may become
50 unstable. This can lead to device freezes and reboots, display corruption and information
51 leakage; see [15] for a discussion of some examples, including information leak bugs in iOS [2]
52 (CVE-2017-2424) and Chrome [5] caused by GPU driver bugs, and an NVIDIA machine
53 freeze [38] (CVE-2017-6259).

54 One way to provide a degree of graphics driver quality—and shader compiler quality in
55 particular—is via standardized test suites. An example is the Khronos Vulkan Conformance
56 Test Suite (Vulkan CTS, or just CTS for short) [25]. This is a large set of functional
57 tests for implementations of the Vulkan graphics API. The Khronos Group, who define
58 various standards and APIs including Vulkan, requires a Vulkan implementation (such as a
59 GPU driver and its associated GPU hardware) to demonstrate that they pass the Vulkan
60 CTS tests in order for the vendor to use the official Vulkan branding. Google’s Android
61 Compatibility Test Suite incorporates the Vulkan CTS, so that Android devices that provide
62 Vulkan capabilities must include drivers that pass the Vulkan CTS. Improving the quality
63 and thoroughness of the Vulkan CTS is thus an indirect method for improving the quality of
64 Vulkan graphics drivers in general, and on Android in particular.

65 GraphicsFuzz (originally called GLFuzz) [16, 15] is a technique and tool chain for auto-
66 matically finding crash and miscompilation bugs in shader compilers using metamorphic
67 testing [9, 42]. Whenever GraphicsFuzz synthesizes a test that exposes a bug in a conformant
68 Vulkan driver, this demonstrates a gap in the Vulkan CTS: the driver has passed the con-
69 formance test suite despite exhibiting this bug. If GraphicsFuzz synthesizes a test that covers
70 a part of a conformant driver’s source code, but the driver does not crash, and the code is
71 not covered by any existing CTS tests, then this also exposes a CTS gap (albeit arguably a
72 less severe one): it demonstrates that part of the driver’s source code *can* be covered but is
73 *not* covered by the CTS; bugs that creep into such code in the future would not be caught.

74 In this experience report we describe our activities at Google over the last 18 months put-
75 ting the GraphicsFuzz tool chain into production, with the aim of improving implementations
76 of the Vulkan API. We have set up a process whereby the randomized metamorphic testing
77 capabilities of GraphicsFuzz are used to find tests that expose driver bugs or CTS coverage
78 gaps, shrink such tests down to small examples that are simple enough for humans to read
79 and debug, and package the resulting tests into a form whereby they are almost completely
80 ready to be added to the Vulkan CTS. So far, this has led to 122 tests that exposed driver
81 and tooling bugs and 113 that exposed driver coverage gaps being added to CTS. The bugs
82 affect a range of mobile and desktop drivers, as well as tools in the SPIR-V ecosystem. Our
83 contribution of CTS tests that expose them means that future conformant Vulkan drivers
84 cannot exhibit them (at least not in a form that causes these tests to fail).

85 We start by presenting relevant background on graphics programming APIs, shader
86 processing tools, the Vulkan CTS, and the GraphicsFuzz testing approach (§2). We then
87 describe how we set up a pathway for incorporating tests that expose bugs found by
88 GraphicsFuzz into the CTS, and various practical issues we had to solve to ensure valid
89 tests (§3). With this pathway in place we were empowered to build a fuzzing framework, *gfauto*,
90 for running GraphicsFuzz against a range of drivers and shader processing tools, automatically
91 shrinking tests that find bugs and getting them into a near-CTS-ready form (§4). To aid in
92 finding coverage gaps, we have built tooling for *differential* coverage analysis; we describe

93 how—by treating coverage gaps as bugs—`gfauto` can be used to synthesize tests that expose
 94 such gaps in a highly automatic fashion (§5). A strength of `GraphicsFuzz` is that it facilitates
 95 testing not only vendor graphics drivers, but also a variety of translation, optimization and
 96 validation tools that are part of the Vulkan ecosystem. We explain how this also presents a
 97 challenge: it can be difficult to determine which component of the ecosystem is responsible
 98 for a bug (§6). Throughout, we provide illustrative examples of noteworthy bugs and tests
 99 found and generated by our approach, including bugs that affect core infrastructure (such as
 100 LLVM), bugs that affect multiple tools simultaneously, and bugs for which the responsible
 101 tool is non-trivial to identify. We conclude by discussing related (§7) and future (§8) work.

102 **Main takeaways** We hope this report is simply interesting for researchers and practitioners
 103 to read as an example of successful technology transfer of research ideas to industrial practice.
 104 In addition, we believe the following aspects could provide inspiration for follow-on research:

- 105 ■ The pros and cons of fuzzing a low level language via a program generator for a higher
 106 level language and a suite of translation and optimization tools, including the problem of
 107 how to determine *where* in a tool chain a fault has occurred (§3.1 and §6);
- 108 ■ The need for image differencing algorithms that are well-suited to tolerating the degree
 109 of variation we expect from graphics shaders due to floating-point precision (§3.4);
- 110 ■ Threats to test validity caused by undefined behavior, long-running loops and floating-
 111 point precision, where more advanced program analyses have the potential to be ap-
 112 plied (§3.5);
- 113 ■ The difficulty of correctly maintaining a test case generator and a corresponding test case
 114 reducer, especially when test case reduction needs to be semantics-preserving (also §3.5)
- 115 ■ The challenge of de-duplicating bugs that do not exhibit distinguished characteristics,
 116 such as wrong image bugs and message-free compile and link errors (§4.2);
- 117 ■ The idea of using differential coverage analysis and test case reduction to fill coverage
 118 gaps (§5), and the challenge of going beyond synthesizing tests that trivially cover new
 119 code to tests that are also equipped with meaningful oracles (§5.4 specifically).

120 **Open sourcing** Our extensions to the `GraphicsFuzz` tool, the new `gfauto` tool, and our
 121 infrastructure for differential code coverage, are open source.¹ The tests we have contributed
 122 to Vulkan CTS are also open source.²

123 2 Background

124 2.1 The GLSL and SPIR-V Shading Languages

125 **GLSL** The OpenGL Shading Language (GLSL) [22] is the main shading language in the
 126 OpenGL graphics API [41] (analogous to HLSL and the Direct3D API). It is used for
 127 rendering hardware-accelerated 2D and 3D graphics. OpenGL ES [32], and its associated
 128 shading language GLSL ES [43], is a subset of the OpenGL API supported by mobile devices,
 129 including Android devices.³ We focus on the GLSL ES shading language version 3.10 and
 130 later, and henceforth drop the ES suffix and version number for brevity. Figure 1 shows an

¹ <https://github.com/google/graphicsfuzz>

² <https://github.com/KhronosGroup/VK-GL-CTS/tree/master/external/vulkancts/data/vulkan/amber/graphicsfuzz>

³ Strictly, OpenGL ES is not quite a subset of OpenGL: over time it has evolved with some features that have been deemed specifically important for mobile platforms.

22:4 Putting Randomized Compiler Testing into Production

```
1 precision highp float;
2
3 layout(location = 0) out vec4 _GLF_color;
4
5 void main()
6 {
7     vec2 a = vec2(1.0);
8     vec4 b = vec4(1.0);
9     pow(vec4(a, vec2(1.0)), b);
10    // Added manually to ensure that the shader writes red
11    _GLF_color = vec4(1.0, 0.0, 0.0, 1.0);
12 }
```

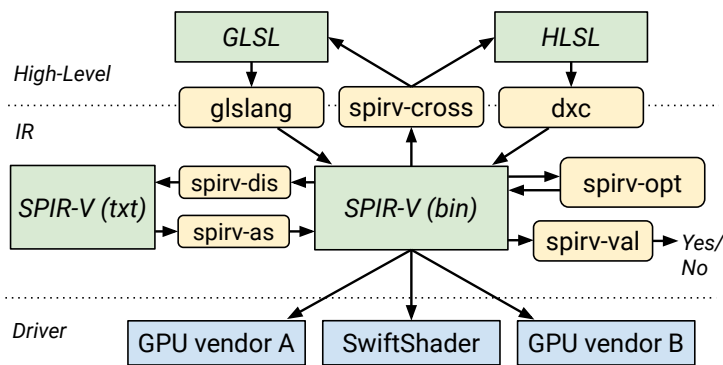
■ **Figure 1** A reduced GLSL ES fragment shader that, after translation to SPIR-V, triggered a bug leading to a crash in the GPU shader compiler for a popular Android device.

131 example of a GLSL *fragment* shader (also known as a pixel shader in Direct3D). The code
132 is C-like, but with some additional features useful for graphics programming. The code in
133 `main` is conceptually executed n times on the GPU for each of the n pixels rendered into a
134 *framebuffer* (which stores the image) using the shader. The `precision highp float;` line
135 causes all subsequent floating-point values to be represented with 32 bits of precision by
136 default (lower precision can be specified via the `mediump` and `lowp` qualifiers on a per-variable
137 basis). Note that `main` has no parameters and a `void` return type; in GLSL, inputs and
138 outputs to the shader are instead expressed using special global variables. Global variable
139 `_GLF_color`⁴ is an output variable into which the fragment shader writes the RGBA colour
140 value that will be rendered at the pixel coordinate for which the shader is running. The `vec2`
141 and `vec4` types are built-in float vector types (with 2 and 4 float components respectively).
142 The vector constructor form that takes one floating-point value (e.g. `vec2(1.0)`) creates a
143 vector with all components set to that value. A vector constructor can also take a combination
144 of vectors and/or scalars (e.g. `vec4(a, vec2(1.0))`) to construct a vector made up of each
145 component in order, as long as the total number of components matches the vector type.
146 The `pow(x,y)` function yields an approximation of x^y , and is an example of one of the many
147 built-in math functions provided in GLSL. The example of Figure 1 was minimized with
148 the aim of reproducing a shader compiler crash bug (discussed further as Example 1 below),
149 and is not representative of a practically useful graphics shader: the `main` function performs
150 some redundant computation and then writes the colour red (`vec4(1.0, 0.0, 0.0, 1.0)`)
151 to the output colour variable. Thus, every pixel rendered by this shader will be red.

152 Shaders can also define *uniform* global variables (not shown in the example), using the
153 `uniform` keyword. These are shader inputs that yield the same value for every pixel
154 being shaded during a single shader invocation. For example, a uniform declaration
155 `uniform float time;` could be used to pass a representation of the current time into
156 a shader, allowing it to produce a time-varying visual effect.

157 **SPIR-V** In comparison to OpenGL, Vulkan [20] is a newer, lower-level graphics API. It is
158 widely supported by modern desktop GPUs, as well as being available on newer Android
159 devices. Standard, Portable Intermediate Representation - V (SPIR-V; the “V” does not
160 stand for anything) is the Vulkan shading language [23]. Unlike GLSL, SPIR-V was designed
161 as an intermediate representation to be stored in a binary form, and thus is not usually written

⁴ The `_GLF` prefix comes from the fact that the tool was originally called GLFuzz. This prefix is used as a default for any special variable, function or macro names used by GraphicsFuzz.



■ **Figure 2** Diagram showing the various tools in the SPIR-V ecosystem and how they interact.

162 directly by programmers. Instead, programmers write their shaders in a higher-level language
 163 like GLSL or HLSL, and use a tool to compile the shaders into SPIR-V. SPIR-V modules
 164 use static single assignment (SSA) form [12], including the use of Phi instructions [12], and
 165 functions contain blocks with branches.

166 2.2 The SPIR-V Tooling Ecosystem

167 Figure 2 summarizes various open source tools for analyzing and transforming SPIR-V
 168 shaders, and translating to and from SPIR-V.

169 As mentioned in §2.1, most shaders are written in high level languages such as GLSL
 170 and HLSL and translated to SPIR-V. For example, `glslang` [24] and `DXC` [37] can compile
 171 GLSL and HLSL, respectively, to SPIR-V. A binary SPIR-V shader can be loaded by the
 172 Vulkan API and executed as part of a graphics pipeline on a GPU device. Google provides
 173 a software implementation of Vulkan, `SwiftShader` [18], which allows Vulkan applications
 174 (including their SPIR-V shaders) to be executed in the absence of Vulkan-capable hardware.
 175 This is useful to bring Vulkan support to old devices, as a fall-back renderer if a GPU driver
 176 goes into an unstable state, and as a “second opinion” for GPU driver writers.

177 The code generated by front-ends such as `glslang` and `DXC` is not typically optimized. In
 178 fact `glslang` deliberately performs as straightforward a syntax-directed translation of a GLSL
 179 shader as possible. The `spirv-opt` tool, part of the Khronos SPIRV-Tools framework [27],
 180 implements many target-agnostic optimizations as SPIRV-V-to-SPIR-V passes.

181 The philosophy of the Vulkan API is to allow drivers to assume that the Vulkan workloads
 182 with which they are presented are valid, pushing the onus of validation to the application. In
 183 support of this, the `spirv-val` tool (also part of the SPIR-V tools framework), checks whether
 184 a SPIR-V shader obeys the (many) rules mandated by the SPIR-V specification [23]. The
 185 `spirv-dis` and `spirv-as` disassembler and assembler (again, part of SPIRV-Tools) allow a shader
 186 to be translated into text format and back, which is useful for debugging.

187 Finally, the `spirv-cross` tool [28] allows SPIR-V to be translated into various shading
 188 languages including GLSL, HLSL and Apple’s Metal shading language (MetalSL, not shown
 189 in the figure). Translation to these higher-level languages can help in understanding the
 190 intended behavior of a SPIR-V shader, and the SPIR-V-to-MetalSL pathway is used by the
 191 MoltenVK project, which provides an implementation of most of Vulkan on top of Apple’s
 192 Metal graphics API [26].

193 2.3 The Vulkan Conformance Test Suite

194 The Khronos Vulkan Conformance Test Suite (Vulkan CTS) [25] is a set of tests for the Vulkan
 195 API. In theory, every part of the Vulkan specification should have one or more corresponding
 196 tests in the Vulkan CTS. Each test should invoke the relevant Vulkan API functions to check
 197 that a Vulkan implementation conforms to the Vulkan specification. Indeed, the Vulkan
 198 CTS mostly consists of a set of functional tests (there are over 550,000 Vulkan CTS tests at
 199 the time of writing) that attempt to test features in isolation. The Vulkan CTS is part of
 200 the larger Khronos Conformance Test Suite called dEQP (drawElements Quality Program⁵)
 201 that additionally contains tests for OpenGL ES and EGL.

202 Any implementation of Vulkan (including any Vulkan graphics driver with its associated
 203 GPUs) must pass the Vulkan CTS (and upload the results to Khronos for peer review) before
 204 the Vulkan name or logo can be used in association with the implementation. Thus, the
 205 Vulkan CTS sets a minimum quality standard for every conformant Vulkan implementation.
 206 Of course, the test suite is also extremely useful during development of a Vulkan driver; as
 207 with most test suites, it can be used to identify bugs and regressions, and to measure progress
 208 towards becoming a conformant implementation. The OpenGL ES and EGL test suite is
 209 similarly used as part of the conformance process for those APIs, and as a useful aid during
 210 driver development. The dEQP test suite is included in the Android Compatibility Test Suite
 211 (Android CTS), which is an even larger test suite for Android devices. Original equipment
 212 manufacturers (OEMs) will typically customize the Android OS for a given device, but
 213 these Android implementations must still pass the Android CTS to be deemed “compatible”.
 214 Thus, the Vulkan CTS also sets a minimum quality standard for Vulkan on every compatible
 215 Android device, which can have a large impact on the Android ecosystem.

216 Vulkan CTS development is mostly done by Khronos members, although anybody can
 217 contribute. New tests are reviewed by GPU vendors before being accepted. Tests need to be
 218 deterministic, and clear enough to allow debugging of Vulkan implementations if a test fails.

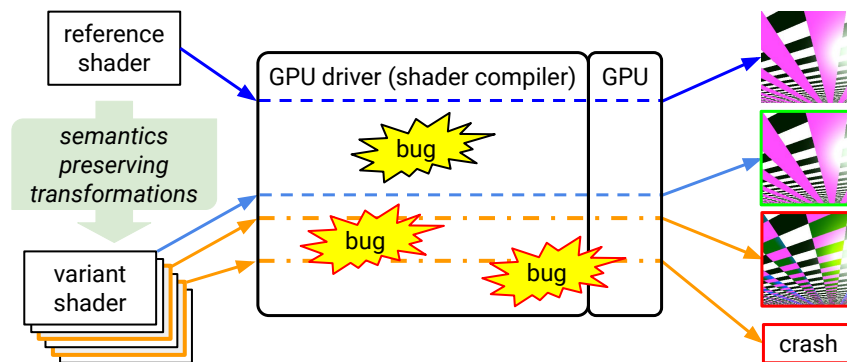
219 2.4 Metamorphic Compiler Testing Using GraphicsFuzz

220 The GraphicsFuzz tool originated from a research project at Imperial College London, and
 221 formed the basis of a spin-out company, GraphicsFuzz Ltd., founded by the authors of this
 222 paper, which Google acquired during 2018.

223 Figure 3 gives an overview of the GraphicsFuzz approach to testing shader compilers. It
 224 starts with an existing reference shader which, after being compiled by the shader compiler
 225 embedded in the GPU driver and executed on the GPU hardware, leads to a given reference
 226 image. A classic way to test a GPU driver would be to compare this resulting image to what
 227 a reference implementation of the graphics API would produce. However, graphics API are
 228 purposefully relaxed to let GPU vendors reach very high performance through aggressive
 229 optimizations, such that there are various images that can be deemed acceptable. It is
 230 currently not possible, and not desirable for GPU vendors, to agree on a strict reference
 231 implementation that would serve as a test oracle.

232 Inspired by the *equivalence modulo inputs* method for testing C compilers [31], GraphicsFuzz
 233 works around this lack of oracle by using *metamorphic testing* [9, 42]: here the GPU input
 234 (shader) is transformed in a way that should not change its output (image). In practice, the
 235 `gsl-fuzz` tool applies *semantics-preserving transformations* to the reference shader source code

⁵ dEQP was a commercial product developed by the drawElements company. Google acquired drawElements in 2014 and donated the dEQP test suite to Khronos, where it is now open source.



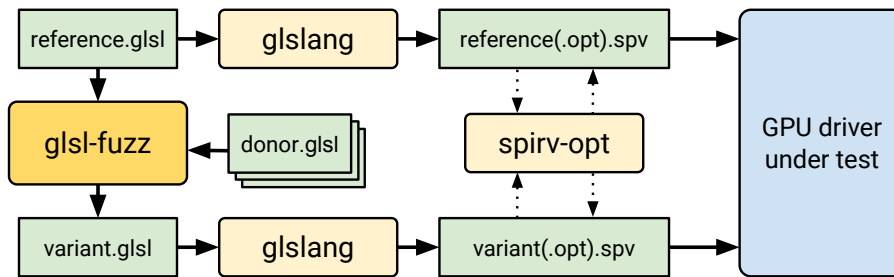
■ **Figure 3** Illustration of the metamorphic testing approach used by GraphicsFuzz

236 to obtain a family of variant shaders. As a very simple example, one can add zero to an existing
 237 integer operation, `int x = y + 0`: this source code change should not impact the program
 238 behavior. The `glsl-fuzz` tool contains many such semantics-preserving transformations [15],
 239 including: arithmetic operations (such as adding zero, multiplying by one, etc), boolean
 240 operations (e.g. `bool b = x && true`), dead code injection (adding valid yet unreachable
 241 code, e.g. wrapped inside an `if (false) { ... }`), live code injection (adding code that
 242 will be executed while making sure to save and restore all variables affected by it, e.g.
 243 `{t = x; x = foo(x); x = t;}`), control flow wrapping (e.g. wrapping existing code in a
 244 single-iteration loop `do { ... } while(false)`), packing scalar data into composite data
 245 types (such as structs and vectors), and outlining expressions into functions. Some of the
 246 values used in these transformations, such as zero, one, true and false, are obtained via
 247 program inputs whose value is guaranteed at execution time, but unknown at compilation
 248 time. This is to make sure compilers cannot trivially remove some transformations, e.g. by
 249 statically detecting dead code injections to be unreachable.

250 Care is required when applying these transformations to ensure that program semantics
 251 are preserved. For instance, one cannot wrap some code in a single-iteration loop if this code
 252 contains a top-level `break` statement: the `break` would now apply to the newly-introduced
 253 loop, rather than to the original loop or switch statement in which it originally appeared.

254 Each variant shader is syntactically distinct from the reference, yet has the same semantics
 255 (modulo floating-point error). It may thus exercise a different path in the shader compiler but
 256 should still lead to a visually similar image being rendered, so long as the reference shader
 257 is sufficiently numerically stable. This is illustrated by the top, blue variant shader line in
 258 Figure 3. However, some variants may lead to significantly different images, or a driver crash,
 259 which are symptoms of bugs, most likely in the shader compiler but also potentially in other
 260 parts of the driver or GPU hardware. These are illustrated by the lower two orange variant
 261 shader lines in Figure 3. For a given GPU, we cannot know what to expect as a reference
 262 image, but we do expect variants to lead to an extremely similar image.

263 Semantics-preserving transformations are used in other contexts, e.g. compiler optimiza-
 264 tions and code obfuscation tools modifying a program representation while keeping the
 265 same behavior. Code refactoring, when understood as improving the program structure
 266 while keeping the same functional features, can also be considered as a semantics-preserving
 267 transformation at a bigger scale than the transformations used in `glsl-fuzz`. For our testing
 268 purpose, we are interested in any kind of semantics-preserving transformation that may
 269 potentially have interesting effects on how the shader is processed by the GPU.



■ **Figure 4** Targeting SPIR-V shader compilers from GLSL.

270 Although a bug-inducing variant can be used as a starting point for debugging, its source
 271 code is often barely understandable by a human because of the hundreds of transformations
 272 that have been applied to it. To ease debugging, the `glsl-reduce` tool progressively shrinks
 273 the variant source code while making sure that the bug is still triggered.

274 There are two reduction modes:

275 **Semantics-preserving reduction** For shader miscompilation bugs leading to wrong
 276 images, `glsl-reduce` performs *semantics-preserving reduction* by removing the `glsl-fuzz`
 277 transformations in a way that still preserves semantics. This typically leads to a variant that
 278 differs from its reference only by a handful of transformations necessary to trigger the wrong
 279 image bug. The pair of semantically identical shaders is useful as a debugging starting point.

280 **Semantics-changing reduction** For bugs leading to a driver crash, `glsl-reduce` performs
 281 *semantics-changing reduction* by removing source code, the only requirement being to keep
 282 it statically valid (which includes being syntactically valid and well-typed). No valid shader
 283 should not cause a driver crash, so there is no need to keep a semantic equivalence with
 284 the reference shader. Semantics-changing reductions can lead to very short crash-inducing
 285 shaders (e.g. Figure 1), which are useful for debugging and as regression test cases.

286 3 Integrating GraphicsFuzz Tests With Vulkan CTS

287 As described in §2.4, the `GraphicsFuzz` tool was originally designed to find bugs in OpenGL
 288 and OpenGL ES drivers by transforming shaders written in the GLSL shading language.
 289 However, our interest is in making shader compilers for the more modern Vulkan API as
 290 reliable as possible by improving the Vulkan CTS, and Vulkan uses SPIR-V as its shading
 291 language (see §2.1). We explain the process we used to allow GLSL-based fuzzing of SPIR-V
 292 shader compilers via translation (§3.1). We explain why we did *not* opt for embedding the
 293 fuzzer inside CTS, or directly contributing large numbers of fuzzer-generated tests, instead
 294 preferring to add tests that are known to expose shader compiler bugs (§3.2). We then
 295 describe how we paved the way for tests that expose crash and wrong image bugs to be
 296 added to CTS (§3.3 and §3.4, respectively).

297 3.1 Fuzzing SPIR-V Compilers via GLSL Shaders

298 In order to target SPIR-V shader compilers with a tool that operates on GLSL, we leverage
 299 the `glslang` translator, which takes GLSL as input and has a SPIR-V back-end. By design,
 300 `glslang` performs a very straightforward translation from GLSL to SPIR-V, performing no
 301 optimization beyond some basic constant folding and elimination of functions that are never
 302 invoked. As a result, the SPIR-V that `glslang` emits is rather basic (e.g. it rarely exhibits uses

303 of Phi instructions). While it is vital that SPIR-V shader compilers correctly handle this
304 “vanilla” SPIR-V, we are also interested in testing their support for more interesting SPIR-V
305 features. Towards this aim, we optionally invoke the `spirv-opt` tool on the SPIR-V that `glslang`
306 generates, with its `-O` flag (optimize for speed), its `-Os` flag (optimize for size), or a random
307 selection of its finer-grained flags (which include things like `--ssa-rewrite`, which changes
308 variable uses to register uses, and can add Phi instructions, and `--eliminate-dead-inserts`,
309 which avoids unnecessary insertions of data into composite structures).

310 This use of `glslang` and `spirv-opt` allows us to perform metamorphic fuzzing at the GLSL
311 level to generate a variant from a reference, send both the variant and reference through
312 `glslang` to turn them into SPIR-V, and then (optionally, and at random) transform the variant
313 using a configuration of `spirv-opt`. The resulting SPIR-V shaders can then be compiled and
314 executed on a Vulkan driver, and the results they compute can be compared. This process is
315 illustrated graphically in Figure 4. This translation-based approach allows us to also find
316 bugs in `glslang` and `spirv-opt`, which benefits the Vulkan ecosystem. However, as discussed
317 further in §6, it can be hard to determine—in the case of wrong image bugs—which of these
318 tools or the driver’s shader compiler has miscompiled.

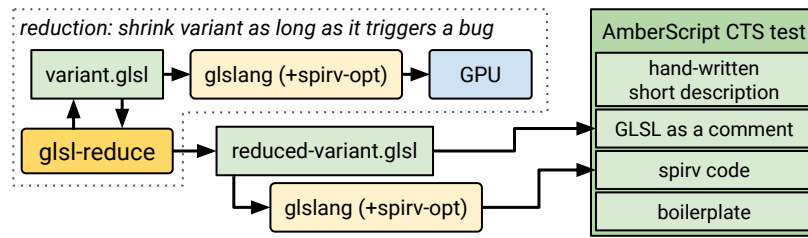
319 **Making shaders “Vulkan-friendly”** Unlike in GLSL, where a global variable of almost
320 any type can be declared as `uniform` (see §2.1), SPIR-V requires that every uniform is
321 declared as a field of a structure called a *uniform block*, with the whole structure being
322 declared to be `uniform`. The number of uniform blocks allowed in a SPIR-V module is
323 implementation-dependent. GLSL has been updated with “Vulkan-friendly features” to
324 allow uniforms to be presented in this way, and `glslang` will only compile Vulkan-friendly
325 shaders into SPIR-V. We thus wrote a simple pass to turn a standard GLSL shader into
326 Vulkan-friendly form. For simplicity of implementation we approached this by placing each
327 original uniform variable in its own (single-field) uniform block. Our pass limits the number
328 of such blocks to 10, as we have not encountered a Vulkan implementation that supports
329 fewer than 10 uniform blocks, and none of the reference shaders we currently use for testing
330 feature more than 10 uniforms. When `glsl-fuzz` generates a variant shader with more than 10
331 uniforms (due to injecting code from other shaders), our Vulkan preparation pass demotes
332 superfluous uniforms to standard global variables initialized to concrete values.

333 3.2 Argument for Not Running Fuzzing in CTS

334 We briefly considered pitching to Khronos the idea of running `GraphicsFuzz` as part of
335 running CTS, so that to pass CTS a driver would have to pass all of the regular tests, and
336 additionally survive running a certain number of `GraphicsFuzz`-generated tests unscathed.
337 We quickly dismissed this idea because it is important to GPU driver makers that qualifying
338 as Vulkan-conformant involves passing a fixed number of tests that run in a deterministic
339 fashion. However much enthusiasm driver makers have for randomized testing as a way to
340 discover bugs, it is understandable that there is little appetite for a conformance test suite
341 that exhibits randomization.

342 Another issue with embedding `GraphicsFuzz` in CTS is that inevitable defects in `Graphic-`
343 `sFuzz` (such as generating a variant shader that turns out *not* to be semantically equivalent
344 to the reference shader) would manifest as a driver failing to pass CTS.

345 An alternative to actually running `GraphicsFuzz` in CTS would be to generate a reasonably
346 large set of shaders—e.g. 1000 shaders—and contribute them as CTS tests. We also quickly
347 decided against this strategy for a few reasons. First, the intended behavior of a CTS test
348 should be feasible for a Vulkan expert to understand. The generated variant shaders are
349 large (in order to maximize the probability of finding a bug), and not feasible for humans



■ **Figure 5** Overview of the reduction process and the creation of a CTS test in AmberScript.

350 to realistically comprehend in isolation; the reducer, `glsl-reduce`, is essential in shrinking a
 351 bug-inducing variant to a comprehensible form. Furthermore, 1000 large randomized shaders
 352 would be a substantial addition to CTS in terms of the test suite’s runtime, but is not a large
 353 enough number of tests to run with the expectation of thoroughly testing a shader compiler.

354 We opted instead for setting up a continuous fuzzing process whereby we could use
 355 `GraphicsFuzz` to find bugs that affect current shader compilers, use `glsl-reduce` to shrink
 356 the associated tests down to small examples that reproduce said bugs, and contribute the
 357 resulting tests. We now explain the format we settled on for adding crash and wrong image
 358 tests to CTS. We detail our tooling for continuous fuzzing in §4.

359 3.3 Supporting Crash Tests

360 Around the same time we commenced our plan to add tests exposing shader compiler crash
 361 bugs to CTS, a new tool called Amber was launched [17]. Amber provides a simple domain-
 362 specific language, AmberScript, in which some aspects of a Vulkan graphics pipeline can be
 363 specified, including the SPIR-V shaders that should be executed, and input and output data
 364 (including uniform blocks and framebuffers) on which the shaders should operate. It also
 365 allows querying the results of running a shader, e.g. probing pixels in the output framebuffer.
 366 The motivation for Amber was to make it easy to write stand-alone shader compiler tests,
 367 hiding the (very substantial amount of) Vulkan API boilerplate required for even a simple
 368 graphics pipeline. Since early 2019, Amber has been integrated into Vulkan CTS and is now
 369 the preferred method for writing shader compiler tests.

370 We wrote a script that takes a reduced GLSL shader known to trigger a SPIR-V shader
 371 compiler crash (after translation to SPIR-V and possibly optimization using some specific
 372 `spirv-opt` flags) and produces an Amber test comprised of:

- 373 ■ A brief comment, supplied as an argument to the script, to describe the test and the
 374 reason why it should be expected to pass;
- 375 ■ A comment showing the original GLSL code for the reduced shader; this is useful because
 376 GLSL is much easier to read compared with SPIR-V;
- 377 ■ Assembly code for the SPIR-V fragment shader that was obtained from this GLSL by
 378 translation using `glslang` and (optional) optimization using `spirv-opt`;
- 379 ■ A comment listing the `spirv-opt` arguments that were applied (if any);
- 380 ■ Commands to create the target framebuffer and to populate the shader uniforms;
- 381 ■ A command, supplied as an argument to the script, to check some property of the image
 382 finally obtained in the framebuffer.

383 Figure 5 illustrates the process of Amber test creation following test case reduction.

```

1  # A test for a bug found by GraphicsFuzz.
2
3  # Short description: A fragment shader that uses pow
4
5  # We check that all pixels are red. The test passes because main does
6  # some computation and then writes red to _GLF_color.
7
8  SHADER vertex variant_vertex_shader PASSTHROUGH
9
10 # variant_fragment_shader is derived from the following GLSL:
11 # #version 310 es
12 #
13 # precision highp float;
14 # precision highp int;
15 #
16 # layout(location = 0) out vec4 _GLF_color;
17 #
18 # void main()
19 # {
20 #   vec2 a = vec2(1.0);
21 #   vec4 b = vec4(1.0);
22 #   pow(vec4(a, vec2(1.0)), b);
23 #   _GLF_color = vec4(1.0, 0.0, 0.0, 1.0);
24 # }
25 SHADER fragment variant_fragment_shader SPIRV-ASM
26 ; SPIR-V
27 ; Version: 1.0
28 ; Generator: Khronos Glslang Reference Front End; 7
29 ; Bound: 28
30 ; Schema: 0
31
32     OpCapability Shader
33     %1 = OpExtInstImport "GLSL.std.450"
34     OpMemoryModel Logical GLSL450
35     OpEntryPoint Fragment %main "main" %_GLF_color
36     OpExecutionMode %main OriginUpperLeft
37     OpSource ESSL 310
38     OpName %main "main"
39     OpName %a "a"
40     OpName %b "b"
41     OpName %_GLF_color "_GLF_color"
42     OpDecorate %_GLF_color Location 0
43
44     %void = OpTypeVoid
45     %3 = OpTypeFunction %void
46     %float = OpTypeFloat 32
47     %v2float = OpTypeVector %float 2
48     ...
49     %21 = OpCompositeConstruct %v4float %17 %18 %19 %20
50     %22 = OpLoad %v4float %b
51     %23 = OpExtInst %v4float %1 Pow %21 %22
52     OpStore %_GLF_color %27
53     OpReturn
54     OpFunctionEnd
55
56 END
57
58 BUFFER variant_framebuffer FORMAT B8G8R8A8_UNORM
59
60 PIPELINE graphics variant_pipeline
61   ATTACH variant_vertex_shader
62   ATTACH variant_fragment_shader
63   FRAMEBUFFER_SIZE 256 256
64   BIND BUFFER variant_framebuffer AS color LOCATION 0
65 END
66 CLEAR_COLOR variant_pipeline 0 0 0 255
67
68 CLEAR variant_pipeline
69 RUN variant_pipeline DRAW_RECT POS 0 0 SIZE 256 256
70
71 EXPECT variant_framebuffer IDX 0 0 SIZE 256 256 EQ_RGBA 255 0 0 255

```

■ **Figure 6** CTS test that exposed a shader compiler crash bug, in AmberScript form. Some of the SPIR-V assembly has been omitted.

22:12 Putting Randomized Compiler Testing into Production

```
1 #version 310 es
2 precision highp float;
3
4 layout(location = 0) out vec4 _GLF_color;
5
6 vec3 GLF_live6mand()
7 {
8     return mix(uintBitsToFloat(ulevec3(38730u, 63193u, 63173u)),
9               floor(vec3(463.499, 4.7, 0.7)), vec3(1.0) + vec3(1.0));
10 }
11 void main()
12 {
13     GLF_live6mand();
14     _GLF_color = vec4(1.0, 0.0, 0.0, 1.0);
15 }
```

■ **Figure 7** Reduced shader that triggered a floating-point constant folding bug in LLVM

384 ► **Example 1.** The GLSL shader of Figure 1, which we used to illustrate the GLSL language
385 in §2.1, triggered a SPIR-V shader compiler crash in the GPU driver of a popular Android
386 device, after translation to SPIR-V (and without requiring any use of `spirv-opt`). This shader
387 was reduced from a much larger variant shader generated by `GraphicsFuzz`, which we edited
388 by making the variable names simpler, and by adding the final line of executable code, which
389 ensures that the colour red is written to the framebuffer. We believe the shader compiler
390 crash was due to an assertion failing in the lowering of the `pow` intrinsic to LLVM bytecode.
391 This is somewhat surprising given that the result of `pow` is not used, but was presumably
392 due to dead code elimination being executed after lowering.

393 An abbreviated version of the Amber test corresponding to this shader is shown in
394 Figure 6 (we omit some of the SPIR-V assembly). The test and its intent are described on
395 lines 1–6; line 8 indicates that a standard trivial *vertex* shader (not otherwise relevant in
396 this experience report) should be used in the test pipeline; lines 10–24 show the GLSL code
397 for the fragment shader, and match Figure 1; the corresponding SPIR-V shader (emitted
398 by `glslang`) is shown on lines 26–52 as SPIR-V assembly (notice the invocation of the `Pow`
399 intrinsic on line 49); line 55 declares a framebuffer, and lines 57–62 define a graphics pipeline
400 based on the vertex and fragment shaders, with the framebuffer attached; line 63 sets the
401 back buffer to black (so that any pixels not rendered to would remain black); lines 65–66 run
402 the pipeline, and line 68 asserts that the framebuffer ends up red at every pixel.

403 The purpose of adding this test to CTS was to expose the driver bug that it triggered,
404 so that future drivers cannot be Vulkan conformant unless the underlying bug is (at least
405 partially) fixed. We write red to the framebuffer and assert that the framebuffer indeed
406 ends up being red so that the test has at least some runtime oracle; it does a little more
407 than just checking that shader compilation succeeds. The test would be better if the shader
408 stored values into one or more components of `_GLF_color` using the result of the call to `pow`,
409 and then asserted a suitable framebuffer colour; as it stands the test would pass even if `pow`
410 were compiled incorrectly but without a compiler crash. We occasionally work to contribute
411 higher quality test oracles, but do not agonize over this since the main motivation for adding
412 these tests is to force the elimination of compiler crash bugs from conformant drivers.

413 ► **Example 2.** Figure 7 shows a reduced shader that triggered a bug in AMD’s LLVM-Based
414 Pipeline Compiler (LLPC) [19]: an assertion failed during constant folding:

```
415 amdllpc: external/llvm/lib/Support/APFloat.cpp:1521: llvm::lostFraction llvm::detail::IEEEFloat::
416 addOrSubtractSignificand(const llvm::detail::IEEEFloat &, bool): Assertion '!carry' failed.
```

417 We reported this bug,⁶ and the LLVM compiler developers traced its root cause to a bug
418 in LLVM’s floating-point emulation code related to handling of subnormal numbers, which
419 was promptly fixed.⁷ This demonstrates that shader compiler fuzzing can have positive
420 impact on common infrastructure (LLVM in this case) that is used by many compilers for
421 C-family languages. We contributed a Vulkan CTS test based on this bug, with a structure
422 similar to the example of Figure 6.⁸

423 3.4 Supporting Wrong Image Tests

424 Recall from §2.4 and Figure 3 that GraphicsFuzz finds miscompilation bugs via a variant
425 shader that renders a significantly different image compared to the image rendered by the
426 associated reference shader. In this case `gsl-reduce` reverses as many of the transformations
427 that were applied to the variant shader as possible while the difference persists. To create
428 Vulkan CTS tests suitable for exposing such bugs we worked with the Amber developers to
429 add AmberScript features related to comparing the outputs of multiple graphics pipelines. In
430 particular, we added the ability to compare framebuffers in a *fuzzy* manner. This allows us
431 to turn a GraphicsFuzz reference and reduced-variant shader pair into a single AmberScript
432 file that (a) creates and runs a separate pipeline for each shader, rendering to distinct
433 framebuffers, and (b) asserts fuzzy equality between these framebuffers.

434 A challenge associated with this is the selection of a suitable fuzzy comparison metric
435 for our purpose. We collected a corpus of image pairs that—based on our shader compiler
436 fuzzing experience—we would like to be deemed similar, and a set of pairs that we would like
437 to be deemed different. The corpus includes image pairs produced by graphics drivers during
438 our testing efforts, plus a few manually crafted image pairs that we believe could occur in
439 theory and that we thought may prove challenging for certain comparison algorithms. We
440 experimented with various image comparison algorithms provided by the `scikit-image` [47]
441 Python library, including MSE, NRMSE, SSIM, and PSNR. We also tried several custom
442 image comparison algorithms based on obtaining and comparing image histograms. We
443 found that image histogram comparison was very effective at correctly classifying image pairs
444 in our corpus, except for some manually crafted image pairs where one image was a rotation
445 or mirror of the other. Indeed, the key weakness of image histogram comparison is that all
446 spatial information is lost. A key advantage is it is very resilient to minor differences that
447 other algorithms flag as important, but which we would typically like to be ignored. We
448 chose to initially proceed with an image histogram comparison algorithm for the following
449 reasons: it correctly classifies image pairs in our corpus as well as or better than most other
450 algorithms; it is very simple to understand and implement (which is important because
451 we don’t want GPU vendors to struggle to understand why a Vulkan CTS test has failed
452 and have to debug the image comparison algorithm itself); it has fairly low performance
453 requirements;⁹ with a high tolerance value, it is fairly forgiving of minor differences, and—to
454 achieve a low false alarm rate—we would prefer to incorrectly classify an image pair as
455 similar than incorrectly classify the pair as different (most image differences we encounter in
456 practice are easily detected with a forgiving algorithm/tolerance).

⁶ <https://github.com/GPUOpen-Drivers/llpc/issues/211>

⁷ <https://reviews.llvm.org/D69772>

⁸ <https://github.com/KhronosGroup/VK-GL-CTS/blob/master/external/vulkancts/data/vulkan/amber/graphicsfuzz/mix-floor-add.amber>

⁹ When running the Vulkan CTS on Android, the image comparison is done on the Android device using the CPU, which has some overhead, especially when using a simulated (software) CPU, as is commonly done when testing next-generation hardware.

22:14 Putting Randomized Compiler Testing into Production

457 We implemented our image comparison algorithm, *Histogram EMD* (where EMD stands
458 for Earth Mover’s Distance [29]), in the Amber code base, and added a command to
459 AmberScript of the form:

```
460 EXPECT buffer_1 EQ_HISTOGRAM_EMD_BUFFER buffer_2 TOLERANCE value
```

461 where `buffer_1` and `buffer_2` are framebuffers containing the images we wish to compare
462 and `value` is the tolerance value. The test fails if the difference value returned by the
463 algorithm exceeds the tolerance value.

464 These extensions to Amber provide a pathway for landing tests that expose wrong image
465 bugs in CTS, and we have implemented the necessary scripts to directly generate such tests.
466 We recently put several such tests up for Vulkan CTS code review, and a reviewer quickly
467 found that the validity of one of the tests was questionable due to floating-point precision
468 issues. We discuss this as Example 3 in §3.5. To err on the side of caution, we retracted the
469 other wrong image tests we had put forward and manually simplified each one to double-check
470 that it really did correspond to a driver bug rather than a floating-point precision issue.
471 After sufficient manual simplification, we were able to add an Amber test for each of these
472 bugs, consisting of a *single* shader (and a single pipeline) with a straightforward assertion to
473 check that the single output image is red.

474 In order to be able to add wrong image tests with a pair of shaders to CTS with confidence,
475 we are working on a corpus of reference shaders that are highly numerically stable.

476 3.5 Avoiding Invalid Tests

477 We are anxious not to waste Vulkan CTS reviewer time by proposing tests that turn out
478 to be invalid and get rejected, or—worse—that get accepted (due to the invalidity being
479 subtle, and not leading to failures on current drivers) and subsequently found to be invalid
480 (necessitating their removal from every CTS release they have made it into). We discuss our
481 main concerns related to possible invalid tests.

482 **Preserving semantics during generation and reduction** As explained in §2.4, `GraphicsFuzz`
483 produces a variant shader by having `gsl-fuzz` repeatedly apply semantics-preserving
484 transformations to a reference, and upon finding a potential wrong image bug, invokes
485 `gsl-reduce` to reduce the test case by repeatedly attempting to reverse or simplify trans-
486 formations. For wrong image bugs, it is critical that all transformations preserve semantics
487 both when applied and reversed/simplified. The way `GraphicsFuzz` has been designed, all
488 information about the transformations that have been applied is recorded by `gsl-fuzz` via
489 syntactic markers in the generated shaders. Examples of syntactic markers include using
490 special preprocessor macros, and giving variables and functions special names or name
491 prefixes. The `gsl-reduce` tool then needs to understand these markers and use them to
492 reverse and simplify certain transformations without spoiling the syntactic markers that
493 represent other transformations. In practice we have encountered several hard-to-diagnose
494 bugs where `gsl-reduce` has erroneously changed the semantics of a shader, usually due to
495 reversal of one transformation having messed up the syntactic markers associated with
496 another transformation, which as a result gets incorrectly reversed.¹⁰

497 We guard against this in practice via a degree of manual inspection of the final reduced
498 shader emitted by `gsl-reduce`, and as `gsl-fuzz` and `gsl-reduce` continue to become more stable
499 this issue becomes less relevant. However, based on our experience, we regard having a

¹⁰ See <https://github.com/google/graphicsfuzz/pull/599> as an example pull request that fixes such an issue.

500 separate generator and reducer that must understand one another in an intricate manner
501 to be a serious pitfall of the `GraphicsFuzz` approach. Recent research on *internal* test case
502 reduction has the potential to avoid the need for a separate generator and reducer [34], and
503 could thus be useful in our domain.

504 **Loop limiters** Recall that the *live code injection* transformation performed by `gsl-fuzz`
505 (see §2.4) injects code from a donor shader into the shader under transformation in a manner
506 such that the injected code really gets executed at runtime. A problem here is that the
507 injected code may contain loops, and these loops may run for potentially large numbers of
508 iterations. In particular, if the declarations of variables that control loop execution are not
509 themselves injected, `gsl-fuzz` creates declarations for such variables and initializes them to
510 randomized expressions, which can lead to infinite loops. Programs that risk containing
511 infinite loops are used for compiler testing by tools such as `Csmith` [50], with the philosophy
512 that it is better to accept that some programs will not terminate, and to use a timeout to
513 bound the runtime of any individual test, than to put in place draconian measures to ensure
514 that all loops terminate. Unfortunately, in the world of GPU shader compilers, long-running
515 shaders cause display freezes, leading to the operating system’s GPU watchdog killing the
516 executing shader. This can lead to the shader rendering what appears to be an incorrect
517 image when in fact the image was simply incomplete.

518 We found that this problem confounded our test results, requiring significant manual
519 inspection of final shaders to check for long-running loops. To overcome this we decided to
520 go ahead and put a relatively draconian measure in place: every loop in every live-injected
521 shader is truncated via a *loop limiter*. This is an additional counter variable specific to a
522 loop. It is initialized to zero immediately before the loop. A conditional statement at the
523 start of the loop body breaks from the loop if the counter exceeds a small positive value
524 (randomly chosen at generation time), and increments the counter otherwise.

525 With reference to our discussion above about keeping the generator and reducer in
526 sync: loop limiters are given special names when inserted by `gsl-fuzz`, and when simplifying
527 live-injected code `gsl-reduce` checks for these names and takes care not to remove loop limiters
528 unless removing the entire associated loop. Again, this coupling between generator and
529 reducer is fragile and can be hard to maintain.

530 When reducing a compiler crash bug `gsl-reduce` aggressively shrinks a shader. In this case
531 we allow it to remove loop limiters, which can mean that finally-reduced shaders may contain
532 infinite loops. While the resulting shaders are good enough to reproduce a compiler crash,
533 they are not suitable for addition to CTS, as all CTS tests should be runnable. We therefore
534 inspect shaders manually and edit them to avoid any infinite loops—while preserving the
535 compiler crash—before submitting them for CTS review.

536 **Array bounds clamping** Live-injected code may also contain access into arrays and
537 vector/matrix types, which have the potential to be out-of-bounds if their indexing expressions
538 depend on variables that `gsl-fuzz` initializes to randomized expressions. SPIR-V for Vulkan
539 requires that all accesses are in-bounds. Fortunately, array and vector/matrix sizes are
540 always known statically in GLSL and there are no pointers in the language. We therefore
541 rewrite every array index expression e that appears in live-injected code as `clamp(e, 0, N - 1)`,
542 where N is the size of the array or vector/matrix being accessed, and `clamp(a, b, c)` is the
543 GLSL built-in that clamps a into the range $[b, c]$. An exception to this is when e is a literal
544 that is already in-bounds. As with loop limiters, `gsl-reduce` is responsible for preserving
545 these in-bounds clamping expressions during test case reduction.

546 **Floating-point stability** We use an example to illustrate the risk of submitting invalid
547 CTS tests posed by floating-point instability.

► **Example 3.** A transformation that `GraphicsFuzz` may try to apply is to replace a floating-point expression e with an expression e/ONE , where ONE is an expression guaranteed to evaluate to 1.0 at runtime. `GraphicsFuzz` has many possible ways of synthesizing an expression that is expected to evaluate to 1.0, one method being to generate an expression of the form `length(normalize(v))`, where v is some non-zero vector. The `normalize` GLSL built-in yields a unit vector (when applied to a non-zero vector), and `length` yields the length of a vector, so the expression intuitively should evaluate to 1.0. However, it turns out that the floating-point precision requirements on SPIR-V instructions mean that the result might not *quite* evaluate to 1.0; some round-off error is allowed [20, pp. 1754–1759].

We thought we had found a wrong image bug in `SwiftShader` upon finding a major image difference to be caused by transforming the following code snippet:

```

559 // 'ref' and 's' are 'float' variables; 'ref' has value 32.0 at runtime
560 for (int i = 1; i < 800; i++) {
561     // 'mod' is the floating-point modulus operation
562     if (mod(float(i), ref) <= 0.01) {
563         s += 0.2;
564     }
565     ...
566 }
567
568

```

This code snippet causes `s` to increase by 0.2 every time `i` is a multiple of 32, since this is the only scenario where `mod(float(i), ref)` will be sufficiently small for the `if` condition to evaluate to true. `GraphicsFuzz` replaced `ref` with `ref / length(normalize(vec3(...)))`, where the `...` is a placeholder for a non-trivial but sensible expression that evaluates to 1.0 (so that the resulting vector is (1.0, 1.0, 1.0)).

What we assumed was a bug in `SwiftShader` turned out to be a false alarm. After some manual analysis we found that the divisor `length(normalize(vec3(...)))` evaluated to a value *slightly larger* than 1.0, so that the second argument to the floating-point `mod` built-in was *slightly smaller* than 32.0 (due to `ref` being exactly 32.0). As a result, the statement `s += 0.2` became *unreachable*, even for loop iterations where `i` is a multiple of 32 since the modulus of a multiple of 32 with a value v slightly smaller than 32.0 leads to the value v .

Floating-point precision issues like this hammer home the importance of using numerically stable shaders when searching for wrong image bugs using `GraphicsFuzz`—the code snippet in Example 3 demonstrates that the shader in question was *not* numerically stable. It is also important to maximize the extent to which the transformations that `GraphicsFuzz` applies actually preserve floating-point semantics. The deliberately ambiguous approach that graphics shading languages take to floating-point (in order to accommodate many disparate GPUs) means that we can never be certain that a program transformation will completely preserve semantics (since it can affect the optimizations the shader compiler performs, and those optimizations are permitted to have small effects on floating-point results). However, where possible we try to take measures to avoid floating-point error; for instance we have changed the representation of 1.0 discussed in Example 3 from `length(normalize(v))` to `round(length(normalize(v)))`, where the `round` GLSL built-in rounds its floating-point argument to the nearest integer value; this ensures that the result will indeed be 1.0.

4 gfacto

`gfacto` (short for `GraphicsFuzz auto`) is a set of tools for using `GraphicsFuzz` in a “push-button” fashion with minimal interaction, geared towards generation of tests that can be added to CTS using the pathways described in §3. Pre-`gfacto`, performing a fuzzing run required

597 manually generating a set of variant shaders offline from a set of reference shaders, followed
598 by a number of manual steps to run the reference and variant shaders on target devices,
599 waiting for the shaders to finish, and then manually triggering reductions of interesting
600 variant shaders. This approach is unnecessarily inefficient when the main objective is to
601 find as many interesting variants (i.e. those that expose bugs) for a given device as possible
602 within a fuzzing run. In contrast, the high-level, *automatic* workflow of `gfauto` is: generate
603 a variant shader from a reference shader; run the shaders on the target device; reduce the
604 variant if it is interesting, otherwise discard it; repeat. This process can run continuously for
605 long periods of time, without interaction, which maximizes the number of interesting variant
606 shaders, and thus the potential number of new CTS tests. Using `gfauto` greatly decreases
607 the length of time needed to perform a fuzzing run and submit a number of CTS tests from
608 that run; we estimate the time period has gone from about 1-2 days (pre-`gfauto`) down to 1-2
609 hours (when using `gfauto`).

610 We detail three key features of `gfauto`: creation and replay of self-contained tests (§4.1),
611 bug de-duplication and prioritization (§4.2), and automatic Vulkan CTS test export (§4.3).

612 4.1 Creation and replay of self-contained tests

613 Pre-`gfauto`, the output of a fuzzing run was a directory of images and log files from running
614 reference and variant shaders; the reference and variant shaders themselves were stored
615 in a different directory. Collecting the shaders and output files needed to reproduce and
616 investigate a bug required copying files from different directories, and these files were stored
617 in an ad-hoc format. Furthermore, the versions of the tools required in order to run the
618 test (such as `glslang` and `spirv-opt`) were not captured. Details about the target device were
619 available but were again outputted in yet another directory and were typically archived in
620 an ad-hoc format, if at all. Thus, reproducing and investigating a bug was difficult and
621 time-consuming, and useful information was often lost.

622 `gfauto` generates a self-contained test from the start. The generated test directory contains
623 a `test.json` metadata file and the reference and variant GLSL shaders. The metadata file
624 contains all information needed to run the test, including a list of required tools and their
625 versions (which are downloaded on-the-fly), an error signature for the test (described below,
626 and initially empty until a crash or wrong image is observed), details of the device on which
627 the test should be run (including the driver version), and the steps needed to run the test
628 (e.g. running `spirv-opt` with a given series of optimization passes). In particular, `gfauto` runs
629 the test for the first time using the test metadata file, and is restricted to the tools specified
630 in the metadata; this ensures that no tool dependency can be missed. A test directory can
631 thus be replayed with a single command. In the case of Android, the `test.json` file even
632 captures the Android device serial number so that the test can be automatically replayed on
633 the target device, with no interaction, as long as it is connected to the host machine.

634 4.2 Bug de-duplication and prioritization

635 A fuzzing run pre-`gfauto` would often find a large number of crash-inducing variant shaders,
636 but upon inspection of crash stack traces it would become apparent that many variants
637 were exposing the same bug. Clearly, we would like to prioritize the *unique* bugs found. We
638 wrote several ad-hoc scripts to classify variants that caused crashes into unique “buckets”,
639 so that each bucket represents a unique bug (based on the top function name in the stack
640 trace). However, this process was still tedious (as it involved several manual steps) and
641 unreliable (as the scripts were typically hand-tuned for a given fuzzing run). Furthermore,

642 this classification was never made permanent, so the information would typically be absent
 643 in future fuzzing runs. Thus, we would often re-find bugs that had already been found in
 644 previous runs and we would have to manually avoid investigating these.

645 In `gfauto`, generated tests that expose bugs are stored in buckets in the file system, where
 646 a bucket is a directory named using the “signature” (usually the top function name in the
 647 stack trace). This makes it trivial to identify tests that expose unique bugs (pick one test
 648 from each bucket). The signature is also stored in the test metadata, ensuring the information
 649 is never lost, even if the test is moved. A Python function `get_signature` takes the log
 650 contents as its only input and outputs the signature string; we update this function as needed
 651 to get an accurate bug signature in a number of scenarios. For example, if a stack trace
 652 is present (in one of several different formats), the top function name of the stack trace is
 653 used, if available, falling back to the hex offset of the function otherwise. Alternatively, if
 654 a recognized error message or assertion failure pattern is seen, the error message itself can
 655 be used as the signature. This approach ensures we reliably classify tests in most cases.
 656 A configurable threshold ensures only a small number of tests are stored in each bucket;
 657 subsequent tests are discarded and, crucially, do not need to be reduced, which is expensive.
 658 Additionally, `gfauto` supports downloading and running our Vulkan CTS tests on the target
 659 device, capturing the signatures (if an error occurs), and ignoring these signatures during the
 660 next fuzzing run. This allows `gfauto` to ignore bugs that can already be found by existing
 661 tests, even if the signatures change between fuzzing runs; this might happen due to a graphics
 662 driver update on the target device or due to changes in `gfauto`’s `get_signature` function. In
 663 particular this allows unfixed bugs found in previous fuzzing runs to be ignored, assuming
 664 appropriate CTS tests were created.

665 **Bug de-duplication challenges** The above approach works well most of the time, but
 666 some issues remain. Some bugs are nondeterministic in nature. In particular, some of our
 667 tests appear to trigger memory leaks in certain shader compilers, which can cause an abort to
 668 occur at arbitrary places. Our `gsl-reduce` tool runs the test up to five times initially (before
 669 commencing reduction) in order to validate that the the originally-observed crash signature
 670 can be reproduced. Highly nondeterministic tests will often fail this validation step, as the
 671 signature will be different every time.

672 Another issue is when a driver returns a “shader compile error” or “shader link error”
 673 message, even though the provided shaders are valid. The driver often provides no additional
 674 information, and so there is no way to further distinguish the shader compiler bug. Thus, if
 675 we find hundreds of “shader compile error” bugs, we may have found hundreds of distinct
 676 compiler bugs, or just one, or any number in between. The same issue applies for tests that
 677 expose wrong image bugs, which are simply given a signature of “`wrong_image`”. In future
 678 work, we hope to identify tests that likely expose distinct wrong image bugs by comparing the
 679 semantics-preserving transformations that remain in the fully-reduced variant shaders. Tests
 680 that contain very distinct transformations are perhaps more likely to be triggering different
 681 shader compiler bugs than tests that contain similar transformations. For compile/link
 682 errors (where reduction need not be semantics-preserving) we may be able to use a similarity
 683 measure on fully-reduced shaders for de-duplication purposes, drawing on ideas for “taming”
 684 compiler fuzzers [10].

685 4.3 Vulkan CTS test export

686 Creating a Vulkan CTS test from a bug found by `GraphicsFuzz` (using `gfauto` or otherwise)
 687 requires some manual iteration on the test. As explained in §4.1, reproducing and investigating
 688 a bug pre-`gfauto` was time-consuming, and information was liable to be lost. The self-contained

689 nature of a `gfauto` test greatly improves the experience. Iterating on a CTS test typically
690 requires tweaking the original GLSL shaders and re-generating the SPIR-V (using the correct
691 versions of `glslang` and `spirv-opt`) again and again. As already stated above, we believe this
692 has greatly decreased the time needed to get from a fuzzing run to a number of submitted
693 CTS tests, from about 1-2 days (pre-`gfauto`) down to 1-2 hours (when using `gfauto`).

694 We took the push-button nature of `gfauto` further by automating the end-to-end process
695 of adding a Vulkan CTS test (after manually tweaking the GLSL shaders). Alongside each
696 `gfauto` test, we store a Python script that generates the final `.amber` file for the Vulkan
697 CTS test. The `.amber` file is similar to the one generated when running the test, but with a
698 copyright header and, as illustrated in Figure 6, a short description and a comment explaining
699 *why* the test passes; note that the short description and comment are manually written by
700 us. The Python script includes the name of the output `.amber` file, the contents of these
701 comments, and some optional tweaks, such as additional AmberScript commands that we
702 might want to add to provide an oracle for the test. Another utility tool then takes this
703 `.amber` file and inserts it into the Vulkan CTS source tree, taking care of updating various
704 index files based on the `.amber` file name and the short description comment. This yields a
705 patch that can be directly put up for Vulkan CTS code review.

706 **5 Finding Test Coverage Gaps Using GraphicsFuzz and gfauto**

707 **5.1 Absolute Code Coverage and its Limitations**

708 Line coverage is a widely-used metric for assessing the adequacy of a test suite at a basic
709 level. While many more thorough notions of coverage have been proposed [1], line coverage
710 is appealing because it is both simple to compute and *actionable* [4]—a lack of line coverage
711 can typically be addressed by crafting appropriate tests. A simple idea for growing Vulkan
712 CTS is therefore to run CTS on an open source Vulkan driver and then attempt to write
713 tests to cover parts of the driver that are not reached.

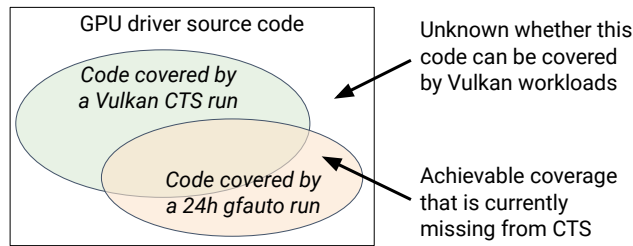
714 This simple idea suffers from two key problems:

- 715 1. It might be *inappropriate* for a Vulkan CTS test to reach certain driver code;
- 716 2. For code that could be covered in principle, it is likely very labour-intensive to manually
717 write tests that cover it in practice.

718 To illustrate problem 1, a recent run of Vulkan CTS on the open source Mesa driver with
719 an AMD back-end [11] identified much uncovered code, but a lot of this code turned out to
720 be (a) debug code (such as routines for dumping data structures in text format), (b) code
721 specific to APIs other than Vulkan (such as OpenCL), and (c) code specific to non-AMD
722 GPUs. It is perfectly legitimate for this code to remain uncovered.

723 **5.2 Differential Code Coverage**

724 To partially solve problem 1 of §5.1 we appeal to *differential* code coverage. Suppose we
725 know which lines of an open source Vulkan driver are covered during a CTS run; call these
726 lines A . Suppose further that we know which lines of the driver are covered by running some
727 other valid Vulkan workload, such as a Vulkan-based game, or a 24-hour run of `gfauto`; call
728 these lines B . For any line $l \in B \setminus A$, we know that l *can* be exercised by valid use of the
729 Vulkan API, so the fact that a CTS run does not exercise l demonstrates a coverage gap in



■ **Figure 8** Illustration of differential code coverage. Driver code covered during a `gfauto` run but not during a CTS run is code that *can* be exercised but for which CTS test coverage is lacking. The tests that `gfauto` generated to achieve such coverage provide a basis for new CTS tests.

730 CTS that can certainly be plugged in principle.¹¹

731 The idea of using `gfauto` and differential coverage to identify code that CTS could in
 732 principle cover is illustrated in Figure 8. One might also imagine that differential coverage
 733 analysis could be used to drive improvements in `GraphicsFuzz`: code that CTS can cover but
 734 that `gfauto` cannot might indicate that `gfauto` should be seeded with a richer set of reference
 735 shaders, or that `GraphicsFuzz` should implement more adventurous transformations. However,
 736 the scope of `GraphicsFuzz` is limited to *shader compilers*, while CTS tests the whole of the
 737 Vulkan API, so some knowledge of which parts of the driver relate to shader compilation
 738 specifically would be required.

739 Although the idea of differential coverage analysis is not new (e.g., continuous integration
 740 systems often provide facilities for visualizing the coverage trajectory of a project), we could
 741 not find a suitable open source project that provides it, so we implemented our own tooling
 742 for differential coverage, which we describe in §5.5.

743 5.3 Using Test Case Reduction to Synthesize Small Tests

744 While differential coverage helps with problem 1 of §5.1, it does not help directly with
 745 problem 2: just knowing that a line is coverable in principle does not yield a suitable CTS
 746 test that covers the line. If workload B (see §5.2) were an interactive game, it might be very
 747 difficult to reverse-engineer a stand-alone test that provides coverage of a particular line.

748 However, if workload B is a `gfauto` run, we can at least obtain a number of `GraphicsFuzz`-
 749 generated variant shaders that provide new coverage. Adding these tests to CTS would serve
 750 to fill the coverage gap, but recall from §3.2 that generated tests are very large, and virtually
 751 impossible for humans to understand in practice, thus unsuitable for direct addition to CTS.

752 To overcome this problem, we appeal to *test case reduction* in the following manner.
 753 Having performed a `gfauto` run for, say, 24 hours, and identified a set of driver source code
 754 lines $B \setminus A$ that were reached by `gfauto` but not by CTS, we manually choose one such line
 755 and prefix it with `assert(false)`—i.e., we pretend that it is *erroneous* to reach the line.
 756 We recompile the driver without coverage instrumentation and run `gfauto` again using the
 757 same parameters (random seed and corpus of shaders) as in the original run. `gfauto` will,
 758 once again, reach the line, this time leading to an assertion failure. `gfauto` will treat the
 759 assertion failure as a shader compiler crash and invoke `gsl-reduce` to shrink the shader to a

¹¹ It is theoretically possible that, e.g. due to concurrency, reachability of line l might be nondeterministic, but we have not encountered this in practice.

760 minimal form that still covers the line. The minimized shader is an excellent candidate for
 761 being added to CTS since it is small enough to be human-readable. The process is repeated
 762 by choosing another line from $B \setminus A$, avoiding lines that we believe are likely to already be
 763 covered by the candidate CTS tests found so far. We periodically re-run CTS after adding
 764 the new tests to update workload A , thus ensuring we don't miss any coverage gaps.

765 At present we have been adopting this approach by collecting differential line coverage of
 766 `SwiftShader`, which incorporates `spirv-opt` and large parts of LLVM internally. The fact that
 767 `SwiftShader` is open source simplifies the process considerably, although it should be possible
 768 to apply a similar process to closed-source binary drivers using instruction coverage instead
 769 of line coverage, and by overwriting an instruction with an interrupt instruction (or some
 770 invalid instruction) instead of prefixing a line with `assert(false)`.

771 **Custom interestingness test** Like many reducers, `gsl-reduce` supports a custom “inter-
 772 estingness test” script that signals to the reducer whether the shader that is being reduced is
 773 still “interesting” (e.g. still crashes the driver). Thus, instead of modifying the driver source
 774 code, we could simply provide an interestingness test that runs the shader using the driver
 775 with coverage instrumentation, processes the coverage data, and checks if the line of interest
 776 was covered. Unfortunately, processing the coverage data is slow, and thus usually done
 777 offline. `gsl-reduce` will typically run the interestingness test hundreds or thousands of times
 778 before finding a minimal shader. Thus, making the driver crash via an assertion failure is
 779 a much faster approach and, conveniently, already triggers a reduction in `gfauto` without
 780 requiring a customized interestingness test.

781 **Less coverage after reduction** A potential downside of our approach is that, after
 782 reduction, a shader may cover fewer lines of interest than before. For example, an unreduced
 783 shader might cover three seemingly unrelated functions, `f`, `g`, and `h`, that are all not covered
 784 by CTS, while the reduced shader might cover just one of the functions, `f`, because we only
 785 added an assertion in `f` and thus the reducer did not try to preserve coverage of `g` and `h`.
 786 Although this may seem undesirable, focusing on just one function at a time typically allows
 787 the reducer to go further (potentially *much* further) in minimizing the shader. We would
 788 much rather have three *simple* and *small* CTS tests, each covering a different function, than
 789 one *complex* and *large* CTS test that covers all three functions.

790 5.4 Manually Tweaking Tests to Improve Oracles

791 Although the reduced tests could be added to the CTS directly, this is almost never appro-
 792 priate. As with crash bugs, we need to add an oracle to the test, else a driver that does
 793 nothing could pass the test. Again as with crash bugs, we typically add code to the shader
 794 to make it render the colour red and add a check to the test to ensure all rendered pixels are
 795 indeed red. However, the test can be made much more useful if the newly covered lines affect
 796 the output colour value so that if a bug was introduced in the newly covered lines, the test
 797 would fail. Also note that a coverage gap test will fill a coverage gap for the Vulkan driver
 798 that we were running (e.g. `SwiftShader`), but the hope is that it may also be a meaningful
 799 test for other drivers, especially if it relates to a feature for which test coverage is generally
 800 lacking. The test should be written with this in mind, as we will see below.

801 ► **Example 4.** The following is a generated, reduced fragment shader that covers constant
 802 folding code in `SwiftShader` that replaces a dot product call with zero if one of the operands
 803 is a vector of zeros:

```
804 void main() {
805     if(1.0 >= dot(vec2(1.0, 0.0), vec2(0.0, 0.0))) {
806
```

22:22 Putting Randomized Compiler Testing into Production

```
803     _GLF_color = vec4(1.0);
804 }
805 }
810 }
```

811 To transform the shader into a form suitable for the Vulkan CTS, we could simply add
812 code to the end of `main` that assigns the colour red to `_GLF_color`, but this has two key
813 disadvantages: (1) any driver that *incorrectly* constant folds the dot function call will still
814 always pass the test, and; (2) a driver might eliminate everything above our final write to
815 `_GLF_color`, and thus, on this hypothetical driver, the test would not even cover code related
816 to the dot product operation. A simple fix is to use the output of the dot function call in the
817 output colour value, as follows:

```
818 void main() {
819     float zero = dot(vec2(1.0, 0.0), vec2(0.0));
820     _GLF_color = vec4(1.0, zero, 0.0, 1.0); // we expect red
821 }
824 }
```

824 However, even this is not ideal; if a driver incorrectly replaced the dot call with a negative
825 float value, the output colour would still be red, as the output colour components are clamped
826 by the driver to be between 0 and 1, as required by the Vulkan specification. In our final
827 version of the test,¹² we check that the result of the dot call is exactly 0, and we only output
828 red in this case:

```
829 void main() {
830     if(dot(vec2(1.0, 0.0), vec2(0.0)) == 0.0) // precise check
831         _GLF_color = vec4(1.0, 0.0, 0.0, 1.0); // we expect red
832     else
833         _GLF_color = vec4(0.0);
834 }
836 }
```

837 The process of manually changing the test to be useful is non-trivial and probably cannot
838 be automated as it requires some creativity, but the reduced test synthesized by `gfauto` is an
839 excellent starting point and is usually not that different to the final version of the test.

840 5.5 Implementing Differential Code Coverage

841 The GCC compiler supports compiling an application with *coverage instrumentation*, causing
842 coverage data to be output when the application runs, which can then be processed with the
843 `gcov` tool. For example, to capture line coverage of `SwiftShader` when running the Vulkan
844 CTS, we could perform the following steps:

- 845 ■ **Compile `SwiftShader` with GCC, adding the `--coverage` flag.** This builds the
846 `SwiftShader` Vulkan library with coverage instrumentation. For each `.o` file that was
847 written, the compiler also writes a `.gcno` file at the same location. The `.gcno` (`gcov`
848 notes) files describe the control flow graph of the corresponding `.o` files, and include
849 mappings to source code file paths and line numbers.
- 850 ■ **Run the Vulkan CTS using `SwiftShader`.** Due to the coverage instrumentation in
851 the `SwiftShader` library, `.gda` files are output alongside the corresponding `.o` and `.gcno`
852 files. The `.gda` (`gcov` data) files contain the control flow graph block and edge execution
853 counts (i.e. the number of times each block and edge was executed).

¹²<https://github.com/KhronosGroup/VK-GL-CTS/blob/master/external/vulkancts/data/vulkan/amber/graphicsfuzz/cov-const-folding-dot-determinant.amber>

854 ■ **Run gcov to process the .gcno and .gcda files to get line coverage information.**
 855 Manually executing gcov on every .gcno file from the required directory (or directories)
 856 while avoiding output filename clashes is a tedious process. Additionally, the line coverage
 857 output from gcov is fairly primitive. Thus, there are third-party tools, such as lcov and
 858 gcovr, that invoke gcov automatically, yielding information in an intermediate data format,
 859 and then further process this data to generate, say, an HTML report that shows every
 860 source file annotated with the execution count of each line.

861 Unfortunately, we could not find any tools capable of obtaining *differential* line coverage
 862 as described in §5.2. Thus, we created our own set of tools¹³ that are similar in spirit to lcov
 863 and gcovr, but support obtaining differential line coverage.

864 `cov_from_gcov` processes .gcno and .gcda files into a single .cov output file that contains
 865 the set of source file lines that were executed. The .gcno and .gcda files are processed by
 866 invoking gcov in each required directory to output intermediate data files that are then
 867 processed further to produce the .cov output file.

868 `cov_new` takes A.cov and B.cov as inputs, and outputs the differential coverage to
 869 `new.cov`. The `new.cov` file is simply A.cov minus B.cov (i.e. $A \setminus B$ from §5.2).

870 `cov_to_source` takes `new.cov` as its input, and outputs two parallel directory structures
 871 `zero/` and `new/`. Both directories contain copies of the original source files with a prefix
 872 added to every line: the `zero/` directory has a “0” prefix added to every line; the `new/`
 873 directory has a “1” prefix added to every line present in `new.cov`, and a “0” prefix for every
 874 other line. The two directory structures can be compared using a diff tool; lines that are
 875 different are the lines of interest (i.e. the lines in $A \setminus B$ from §5.2), and will be highlighted.
 876 The approach of generating parallel directory structures means we avoid generating large
 877 HTML reports, which can be slow to open and navigate, and instead allows the use of any
 878 existing diff tool that is already optimized for this type of task.

879 **6 Fuzzing the SPIR-V Tooling Ecosystem**

880 Recall the rich ecosystem of SPIR-V-related tools shown in Figure 2. Because `gfauto` uses
 881 many of these tools during a fuzzing run, we have the potential to find bugs in them as
 882 well as finding bugs in vendor shader compilers. Furthermore we have also conducted some
 883 fuzzing runs that include the `spirv-cross` tool (not part of the default `gfauto` workflow).

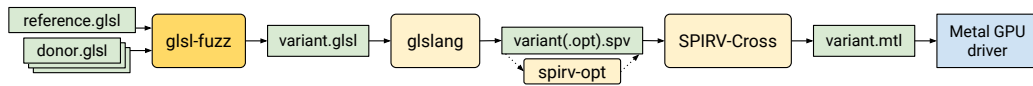
884 We illustrate, via examples, both the strength of being able to find bugs in multiple tools
 885 and the challenge associated with determining which tool is to blame when a problem arises.

886 ► **Example 5.** Running `spirv-opt` on a shader generated by `gfauto` led to a non-zero exit code.
 887 We reported this bug to the SPIRV-Tools project,¹⁴ assuming it to be a bug in `spirv-opt`.
 888 The `spirv-opt` authors investigated and determined that in fact the shader contained *invalid*
 889 SPIR-V that `spirv-val` (which `gfauto` runs at every transformation stage) had missed. This
 890 identified a validator bug (in the form of a validator omission) but raised the question of
 891 what had created the invalid SPIR-V. It turned out that `gfauto` had run `spirv-opt` as part of
 892 generating the shader, and its block merging pass had been too aggressive: loops in SPIR-V
 893 assembly have designated *merge* and *continue* blocks, the *merge* block denoting the loop’s
 894 exit, and the *continue* block denoting the start of a region of code that must be traversed in
 895 order to return to the loop head. The block merging pass was allowing the merge block of

¹³<https://github.com/google/graphicsfuzz/blob/master/gfauto/docs/coverage.md>

¹⁴<https://github.com/KhronosGroup/SPIRV-Tools/issues/3031>

22:24 Putting Randomized Compiler Testing into Production



■ **Figure 9** Fuzzing Metal drivers from GLSL.

896 one loop and the continue block of another loop to be merged. This turned out to be illegal
897 according to the SPIR-V specification, though the wording of the specification did not spell
898 the rule out very clearly. The SPIRV-Tools team enhanced `spirv-val` to detect this kind of
899 invalidity, and fixed the bug in `spirv-opt`'s block merging pass.¹⁵

900 This resolved the combined validator and optimizer bug that we had found with `gfauto`.
901 Unfortunately, it turned out that 55 existing Vulkan CTS tests contained SPIR-V that was
902 invalid for the same reason—in many cases the SPIR-V in question had been processed by
903 `spirv-opt`'s previously buggy block merging pass. This necessitated fixing these tests in the
904 master branch of CTS, as well as in multiple release branches.

905 The SPIR-V working group are discussing how to clarify the specification with respect to
906 its rules about the structure of loops and other control flow constructs, in part due to this
907 (and other) reports from our fuzzing efforts.

908 ► **Example 6.** An early assertion failure that we triggered in `spirv-opt`,¹⁶ by fuzzing using a
909 random combination of optimizer flags, turned out to be due to a “merge return” optimization
910 pass being applied to a SPIR-V control flow graph that it was known not to be able to handle.
911 The SPIRV-Tools team hardened `spirv-opt` by having the “merge return” pass explicitly
912 check for unsupported control flow graph idioms and, on encountering an unsupported idiom,
913 gracefully exit with an error informing the user that they should run the “eliminate dead
914 code” pass first.

915 Subsequently, we ensured that `gfauto` only generates lists of `spirv-opt` optimization passes
916 in which “merge return” (if present) runs after “eliminate dead code”.

917 ► **Example 7.** We ran some experiments testing MoltenVK [26], an implementation of most of
918 Vulkan on top of Apple's Metal graphics API, on a MacBook Pro. MoltenVK uses `spirv-cross`
919 to translate SPIR-V to the Metal shading language (MetalSL) so it can be sent to the Metal
920 shader compiler within the Metal driver on a Mac or iOS device. The full translation pipeline
921 is illustrated in Figure 9. When a wrong image is produced in this setup, the bug could be
922 in the Metal driver or in any of the tools that come before (shown as rounded-rectangles
923 in Figure 9). Differential testing can come to the rescue here: if the bad image is also
924 produced by the variant shader in a “vanilla” setup, e.g. by using `glslang` and rendering the
925 resulting SPIR-V using some other Vulkan driver, the bug is very likely in `glsl-fuzz` or `glslang`.
926 Otherwise, if the bug manifests only when adding the same `spirv-opt` passes (before running
927 on this other Vulkan driver), the problem is likely in `spirv-opt`. Otherwise, the bug is likely
928 in `spirv-cross` or the Metal driver.

929 We found such a wrong image bug in our testing, and used differential testing to conclude
930 that the bug was likely in `spirv-cross` or the Metal driver. After inspecting the MetalSL code,
931 we found it to be incorrect and so ascertained that the bug was in `spirv-cross`. We submitted
932 a bug report¹⁷ and the bug was promptly fixed.

¹⁵ <https://github.com/KhronosGroup/SPIRV-Tools/pull/3068>

¹⁶ <https://github.com/KhronosGroup/SPIRV-Tools/issues/1962>

¹⁷ <https://github.com/KhronosGroup/SPIRV-Cross/issues/1091>

7 Related Work

933

934 **Randomized and metamorphic compiler testing techniques** Randomized testing of
935 compilers has a long history; see e.g. [21] for a very early example, and multiple surveys [8,
936 3, 30]. Random differential testing (RDT) of C compilers was investigated to some extent by
937 McKeeman [36], and the Csmith project from the University of Utah [50] has triggered a lot
938 of interest in the topic over the last decade. An early approach to metamorphic compiler
939 testing involved generating equivalent programs from scratch [46]. More recent work on
940 *equivalence modulo inputs* testing (EMI) [31, 44], a form of metamorphic testing, showed
941 that approaches based on transforming programs in a manner that preserves semantics at
942 least for certain inputs can be an effective way of triggering wrong code bugs. The RDT
943 and EMI approaches have been extended to allow testing of OpenCL compilers [33], and
944 the EMI approach was the inspiration for the approach to metamorphic testing employed
945 by GraphicsFuzz [16, 15], on which the work described in this experience report was built.
946 Randomized differential testing has also been applied to other program processing tools,
947 such as refactoring engines [14] and static analyzers [13], and there is scope for applying
948 techniques from compiler testing to program analyzers more generally [6].

949 **Experience reports related to compiler testing** A short report on work at the UK’s
950 National Physical Laboratory describes some experiences testing compilers for Pascal, Ada
951 and Haskell using random program generators, mainly focusing on the relative difficulty
952 of constructing program generators for each of these languages [48]. McKeeman’s seminal
953 paper on differential testing includes a section on randomized compiler testing that is written
954 in the style of an experience report [36]. In comparison to our paper, these reports do
955 not discuss the challenges of setting up a pathway from randomly-generated test cases to
956 test cases suitable for incorporation into a standard compiler test suite. An edited volume
957 on validation of Pascal compilers provides a number of experience reports related to the
958 testing and validation process [49]. These reports discuss issues related to constructing
959 compiler conformance tests, but do not mention randomized testing. Furthermore, none of
960 the aforementioned experience reports discuss the challenges of testing *graphics* compilers.

961 **Empirical studies related to compiler bugs and compiler testing** There have been
962 three recent empirical studies related to compiler bugs and randomized compiler testing: a
963 study on the relative effectiveness of compiler testing based on RDT vs. EMI testing [7],
964 a study on the characteristics of bugs in the GCC and LLVM compilers (not specifically
965 focusing on bugs found via randomized testing) [45], and a study that aims to assess the
966 relative impact on end-user software of fuzzer-found compiler bugs compared with compiler
967 bugs encountered and reported “in the wild” by users [35]. Unlike our work these studies all
968 focus on C/C++ compilers, not compilers for graphics shading languages. A main finding
969 from [35]—that bugs found by fuzzers appear to have at least as much practical impact
970 as bugs reported by users—supports our belief that adding fuzzer-found compiler bugs to
971 compiler regression test suites is a worthwhile endeavour.

972 **Compiler test case reduction and bug de-duplication** The *semantics-changing* reduction
973 mode of `gsl-reduce` is similar in nature to the approach taken by the C-Reduce tool [40],
974 following the well-known delta debugging method [51]. A difference between `gsl-reduce` and
975 C-Reduce is that `gsl-reduce` exclusively uses valid abstract syntax tree transformations to
976 reduce a shader, whereas C-Reduce uses a combination of methods, including language-aware
977 transformations built on top of the Clang framework, language-agnostic transformations
978 based on line and token deletion, and methods in-between that assume only basic language
979 properties, such as that the language is block-structured, with blocks delimited by braces.

980 As a result, C-Reduce can be applied to programs from a variety of languages (e.g. it has
 981 been successfully applied to OpenCL C [39]), while `gsl-reduce` is specific to GLSL. It would
 982 be interesting to investigate how well C-Reduce works for the reduction of GLSL programs
 983 that induce shader compiler crashes.

984 The semantics-preserving mode of `gsl-reduce` is intimately tied to the semantics-preserving
 985 transformations applied during metamorphic testing. As discussed in §3.5, the tight coupling
 986 between `gsl-fuzz` and `gsl-reduce` related to this mode has made it hard to maintain the pair
 987 of tools. Recent work proposes leveraging a test-case generator to provide test case reduction
 988 “for free”, by repeatedly re-generation to search for smaller tests that still trigger a bug [34].
 989 An approach along these lines may be effective in avoiding the need for a tightly coupled
 990 generator and reducer in our domain.

991 Work on automated ranking of compiler bug reports proposes several metrics that can
 992 be used to order bug-inducing tests, with the aim of presenting a diverse selection of test
 993 cases exposing distinct bugs first [10]. Our de-duplication of crash bugs based on crash
 994 signatures (see §4.2) has not yet required this level of sophistication, but we believe such
 995 techniques could be brought to bear for de-duplication of wrong image bugs for which there
 996 is no analogue to a crash signature.

997 **8 Conclusions and Future Work**

998 We have described our experience rolling out graphics shader compiler fuzzing, based on
 999 the `GraphicsFuzz` tool chain, in a production environment with the goal of improving the
 1000 Vulkan Conformance Test Suite via new tests that expose shader compiler bugs or provide
 1001 additional coverage of shader processing tools. We hope the various insights in this report
 1002 will be useful to researchers interested in testing programming language implementations.

1003 We identify several directions for future practical work in this area.

1004 **Direct fuzzing for SPIR-V** We have gotten significant mileage from testing SPIR-V
 1005 shader compilers via GLSL shaders, but the SPIR-V features this flow will exercise are
 1006 inevitably limited, motivating the need for a fuzzer that works at the SPIR-V level.

1007 **Stability tests** We discuss the avoidance of invalid tests in §3.5. However, unintentionally
 1008 invalid tests (due to bugs in the `GraphicsFuzz` tooling) have sometimes led to the discovery of
 1009 serious driver stability issues, e.g. Android devices rebooting after accessing invalid memory,
 1010 or failing to gracefully recover from long-running shaders [15]. It would be valuable to put in
 1011 place a suite of tests, distinct from Vulkan CTS, to check that invalid shaders cannot derail
 1012 an operating system.

1013 **Higher confidence in wrong image bugs** As discussed in §3.4 we are presently exercising
 1014 caution regarding adding wrong image tests to CTS. A corpus of highly numerically stable
 1015 shaders would allow us to proceed with greater confidence here, as would a more detailed
 1016 analysis of the possible floating-point effects of the transformations that `gsl-fuzz` employs.

1017 **Acknowledgements**

1018 We are grateful to David Neto and to the anonymous ECOOP 2020 reviewers for their
 1019 feedback on an earlier draft of this work.

1020 **References**

-
- 1021 1 Paul Ammann and Jeff Offutt. *Introduction to Software Testing*. Cambridge University Press,
 1022 2 edition, 2017.

- 1023 2 Apple. About the security content of ios 10.3, 2017. [https://support.apple.com/en-gb/](https://support.apple.com/en-gb/HT207617)
1024 HT207617, see “Processing maliciously crafted web content may result in the disclosure of
1025 process memory”.
- 1026 3 Abdulazeez S. Boujarwah and Kassem Saleh. Compiler test case generation methods: a
1027 survey and assessment. *Information & Software Technology*, 39(9):617–625, 1997. URL: [https://doi.org/10.1016/S0950-5849\(97\)00017-7](https://doi.org/10.1016/S0950-5849(97)00017-7), doi:10.1016/S0950-5849(97)00017-7.
1028
- 1029 4 Arkady Bron, Eitan Farchi, Yonit Magid, Yarden Nir, and Shmuel Ur. Applications of
1030 synchronization coverage. In Keshav Pingali, Katherine A. Yelick, and Andrew S. Grimshaw,
1031 editors, *Proceedings of the ACM SIGPLAN Symposium on Principles and Practice of Parallel*
1032 *Programming, PPOPP 2005, June 15-17, 2005, Chicago, IL, USA*, pages 206–212. ACM, 2005.
1033 URL: <https://doi.org/10.1145/1065944.1065972>, doi:10.1145/1065944.1065972.
- 1034 5 bugs.chromium.org. Issue 675658: Security: Malicious WebGL page can capture and upload
1035 contents of other tabs, 2016. [https://bugs.chromium.org/p/chromium/issues/detail?id=](https://bugs.chromium.org/p/chromium/issues/detail?id=675658)
1036 675658.
- 1037 6 Cristian Cadar and Alastair F. Donaldson. Analysing the program analyser. In Laura K. Dillon,
1038 Willem Visser, and Laurie Williams, editors, *Proceedings of the 38th International Conference*
1039 *on Software Engineering, ICSE 2016, Austin, TX, USA, May 14-22, 2016 - Companion*
1040 *Volume*, pages 765–768. ACM, 2016. URL: <https://doi.org/10.1145/2889160.2889206>,
1041 doi:10.1145/2889160.2889206.
- 1042 7 Junjie Chen, Wenxiang Hu, Dan Hao, Yingfei Xiong, Hongyu Zhang, Lu Zhang, and Bing Xie.
1043 An empirical comparison of compiler testing techniques. In Laura K. Dillon, Willem Visser,
1044 and Laurie Williams, editors, *Proceedings of the 38th International Conference on Software*
1045 *Engineering, ICSE 2016, Austin, TX, USA, May 14-22, 2016*, pages 180–190. ACM, 2016.
1046 URL: <https://doi.org/10.1145/2884781.2884878>, doi:10.1145/2884781.2884878.
- 1047 8 Junjie Chen, Jibesh Patra, Michael Pradel, Yingfei Xiong, Hongyu Zhang, Dan Hao, and
1048 Lu Zhang. A survey of compiler testing techniques. *ACM Computing Surveys*, 2020. To
1049 appear.
- 1050 9 T.Y. Chen, S.C. Cheung, and S.M. Yiu. Metamorphic testing: a new approach for generating
1051 next test cases. Technical Report HKUST-CS98-01, Department of Computer Science, The
1052 Hong Kong University of Science and Technology, 1998.
- 1053 10 Yang Chen, Alex Groce, Chaoqiang Zhang, Weng-Keen Wong, Xiaoli Z. Fern, Eric Eide,
1054 and John Regehr. Taming compiler fuzzers. In Hans-Juergen Boehm and Cormac Flanagan,
1055 editors, *ACM SIGPLAN Conference on Programming Language Design and Implementation,*
1056 *PLDI '13, Seattle, WA, USA, June 16-19, 2013*, pages 197–208. ACM, 2013. URL: <https://doi.org/10.1145/2491956.2462173>, doi:10.1145/2491956.2462173.
1057
- 1058 11 Igalia / codecov.io. Coverage report for vulkan cts on open source mesa driver with amd
1059 bck-end, 2020. <https://codecov.io/gh/Igalia/mesa/>.
- 1060 12 Keith Cooper and Linda Torczon. *Engineering a Compiler*. Morgan Kaufmann, 2002.
- 1061 13 Pascal Cuoq, Benjamin Monate, Anne Pacalet, Virgile Prevosto, John Regehr, Boris
1062 Yakobowski, and Xuejun Yang. Testing static analyzers with randomly generated pro-
1063 grams. In Alwyn Goodloe and Suzette Person, editors, *NASA Formal Methods - 4th In-*
1064 *ternational Symposium, NFM 2012, Norfolk, VA, USA, April 3-5, 2012. Proceedings*, volume
1065 7226 of *Lecture Notes in Computer Science*, pages 120–125. Springer, 2012. URL: https://doi.org/10.1007/978-3-642-28891-3_12, doi:10.1007/978-3-642-28891-3_12.
1066
- 1067 14 Brett Daniel, Danny Dig, Kely Garcia, and Darko Marinov. Automated testing of refactoring
1068 engines. In Ivica Crnkovic and Antonia Bertolino, editors, *Proceedings of the 6th joint meeting*
1069 *of the European Software Engineering Conference and the ACM SIGSOFT International*
1070 *Symposium on Foundations of Software Engineering, 2007, Dubrovnik, Croatia, September*
1071 *3-7, 2007*, pages 185–194. ACM, 2007. URL: <https://doi.org/10.1145/1287624.1287651>,
1072 doi:10.1145/1287624.1287651.

- 1073 15 Alastair F. Donaldson, Hugues Evrard, Andrei Lascu, and Paul Thomson. Automated
1074 testing of graphics shader compilers. *PACMPL*, 1(OOPSLA):93:1–93:29, 2017. URL: <https://doi.org/10.1145/3133917>, doi:10.1145/3133917.
- 1076 16 Alastair F. Donaldson and Andrei Lascu. Metamorphic testing for (graphics) compilers. In
1077 *Proceedings of the 1st International Workshop on Metamorphic Testing, MET@ICSE 2016,*
1078 *Austin, Texas, USA, May 16, 2016*, pages 44–47. ACM, 2016. URL: <https://doi.org/10.1145/2896971.2896978>, doi:10.1145/2896971.2896978.
- 1080 17 Google. Amber GitHub repository, 2020. <https://github.com/google/amber>.
- 1081 18 Google. SwiftShader GitHub repository, 2020. <https://github.com/google/SwiftShader>.
- 1082 19 GPUOpen Drivers. LLVM-based pipeline compiler GitHub repository, 2020. <https://github.com/GPUOpen-Driver-LLPC>.
- 1084 20 The Khronos Vulkan Working Group. *Vulkan 1.1.141 - A Specification (with all registered*
1085 *Vulkan extensions)*. The Khronos Group, 2019. [https://www.khronos.org/registry/vulkan/](https://www.khronos.org/registry/vulkan/specs/1.1-extensions/pdf/vkspec.pdf)
1086 [specs/1.1-extensions/pdf/vkspec.pdf](https://www.khronos.org/registry/vulkan/specs/1.1-extensions/pdf/vkspec.pdf).
- 1087 21 K. V. Hanford. Automatic generation of test cases. *IBM Systems Journal*, 9:242–257, 1970.
- 1088 22 John Kessenich, editor. *The OpenGL Shading Language Version 4.60.7*. The Khronos Group,
1089 2019. <https://www.khronos.org/registry/OpenGL/specs/gl/GLSLangSpec.4.60.pdf>.
- 1090 23 John Kessenich, Boaz Ouriel, and Raun Krisch, editors. *SPIR-V Specification, Version*
1091 *1.5, Revision 2, Unified*. The Khronos Group, 2019. [https://www.khronos.org/registry/](https://www.khronos.org/registry/spir-v/specs/unified1/SPIRV.pdf)
1092 [spir-v/specs/unified1/SPIRV.pdf](https://www.khronos.org/registry/spir-v/specs/unified1/SPIRV.pdf).
- 1093 24 Khronos Group. glslang GitHub repository, 2020. [https://github.com/KhronosGroup/](https://github.com/KhronosGroup/glslang)
1094 [glslang](https://github.com/KhronosGroup/glslang).
- 1095 25 Khronos Group. Khronos Vulkan, OpenGL, and OpenGL ES conformance tests GitHub
1096 repository, 2020. <https://github.com/KhronosGroup/VK-GL-CTS>.
- 1097 26 Khronos Group. MoltenVk GitHub repository, 2020. [https://github.com/KhronosGroup/](https://github.com/KhronosGroup/MoltenVK)
1098 [MoltenVK](https://github.com/KhronosGroup/MoltenVK).
- 1099 27 Khronos Group. SPIR-V Tools GitHub repository, 2020. [https://github.com/KhronosGroup/](https://github.com/KhronosGroup/SPIRV-Tools)
1100 [SPIRV-Tools](https://github.com/KhronosGroup/SPIRV-Tools).
- 1101 28 Khronos Group. SPIRV-Cross GitHub repository, 2020. [https://github.com/KhronosGroup/](https://github.com/KhronosGroup/SPIRV-Cross)
1102 [SPIRV-Cross](https://github.com/KhronosGroup/SPIRV-Cross).
- 1103 29 Jeffery Kline. Properties of the d-dimensional earth mover’s problem. *Discrete Applied*
1104 *Mathematics*, 265:128 – 141, 2019. URL: [http://www.sciencedirect.com/science/article/](http://www.sciencedirect.com/science/article/pii/S0166218X19301441)
1105 [pii/S0166218X19301441](http://www.sciencedirect.com/science/article/pii/S0166218X19301441), doi:<https://doi.org/10.1016/j.dam.2019.02.042>.
- 1106 30 Alexander S. Kossatchev and Mikhail Posypkin. Survey of compiler testing methods. *Pro-*
1107 *gramming and Computer Software*, 31(1):10–19, 2005. URL: [https://doi.org/10.1007/](https://doi.org/10.1007/s11086-005-0008-6)
1108 [s11086-005-0008-6](https://doi.org/10.1007/s11086-005-0008-6), doi:10.1007/s11086-005-0008-6.
- 1109 31 Vu Le, Mehrdad Afshari, and Zhendong Su. Compiler validation via equivalence modulo
1110 inputs. In Michael F. P. O’Boyle and Keshav Pingali, editors, *ACM SIGPLAN Conference on*
1111 *Programming Language Design and Implementation, PLDI ’14, Edinburgh, United Kingdom -*
1112 *June 09 - 11, 2014*, pages 216–226. ACM, 2014. URL: [https://doi.org/10.1145/2594291.](https://doi.org/10.1145/2594291.2594334)
1113 [2594334](https://doi.org/10.1145/2594291.2594334), doi:10.1145/2594291.2594334.
- 1114 32 Jon Leech, editor. *OpenGL ES Version 3.2*. The Khronos Group, 2019. [https://www.khronos.](https://www.khronos.org/registry/OpenGL/specs/es/3.2/es_spec_3.2.pdf)
1115 [org/registry/OpenGL/specs/es/3.2/es_spec_3.2.pdf](https://www.khronos.org/registry/OpenGL/specs/es/3.2/es_spec_3.2.pdf).
- 1116 33 Christopher Lidbury, Andrei Lascu, Nathan Chong, and Alastair F. Donaldson. Many-core
1117 compiler fuzzing. In David Grove and Steve Blackburn, editors, *Proceedings of the 36th ACM*
1118 *SIGPLAN Conference on Programming Language Design and Implementation, Portland, OR,*
1119 *USA, June 15-17, 2015*, pages 65–76. ACM, 2015. URL: [https://doi.org/10.1145/2737924.](https://doi.org/10.1145/2737924.2737986)
1120 [2737986](https://doi.org/10.1145/2737924.2737986), doi:10.1145/2737924.2737986.
- 1121 34 David R. MacIver and Alastair F. Donaldson. Test-case reduction via test-case generation:
1122 Insights from the hypothesis reducer. In *34th European Conference on Object-Oriented*
1123 *Programming, ECOOP 2020*, volume 166 of *LIPICs*, pages 13:1–13:28. Schloss Dagstuhl -
1124 Leibniz-Zentrum für Informatik, 2020.

- 1125 35 Michaël Marcozzi, Qiyi Tang, Alastair F. Donaldson, and Cristian Cadar. Compiler fuzzing:
1126 how much does it matter? *PACMPL*, 3(OOPSLA):155:1–155:29, 2019. URL: <https://doi.org/10.1145/3360581>, doi:10.1145/3360581.
- 1128 36 William M. McKeeman. Differential testing for software. *Digital Technical Journal*, 10(1):100–
1129 107, 1998. URL: <http://www.hpl.hp.com/hpjournal/dtj/vol10num1/vol10num1art9.pdf>.
- 1130 37 Microsoft. DirectX shader compiler GitHub repository, 2020. <https://github.com/microsoft/DirectXShaderCompiler>.
- 1132 38 NVIDIA. Security bulletin: Nvidia gpu display driver contains multiple vulnerabilities in
1133 the kernel mode layer handler, 2018. https://nvidia.custhelp.com/app/answers/detail/a_id/4525/, see “NVIDIA GPU Display Driver contains a vulnerability in the kernel mode
1134 layer handler where an incorrect detection and recovery from an invalid state produced by
1135 specific user actions may lead to a denial of service”.
- 1137 39 Moritz Pflanzner, Alastair F. Donaldson, and Andrei Lascu. Automatic test case reduction for
1138 opencl. In *Proceedings of the 4th International Workshop on OpenCL, IWOCL 2016, Vienna, Austria, April 19-21, 2016*, pages 1:1–1:12. ACM, 2016. URL: <https://doi.org/10.1145/2909437.2909439>, doi:10.1145/2909437.2909439.
- 1141 40 John Regehr, Yang Chen, Pascal Cuoq, Eric Eide, Chucky Ellison, and Xuejun Yang. Test-case
1142 reduction for C compiler bugs. In Jan Vitek, Haibo Lin, and Frank Tip, editors, *ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI '12, Beijing, China - June 11 - 16, 2012*, pages 335–346. ACM, 2012. URL: <https://doi.org/10.1145/2254064.2254104>, doi:10.1145/2254064.2254104.
- 1146 41 Mark Segal and Kurt Akeley, editors. *The OpenGL Graphics System: A Specification Version 4.6 (Core Profile)*. The Khronos Group, 2019. <https://www.khronos.org/registry/OpenGL/specs/gl/glspec46.core.pdf>.
- 1149 42 Sergio Segura, Gordon Fraser, Ana B. Sánchez, and Antonio Ruiz Cortés. A survey on
1150 metamorphic testing. *IEEE Trans. Software Eng.*, 42(9):805–824, 2016. URL: <https://doi.org/10.1109/TSE.2016.2532875>, doi:10.1109/TSE.2016.2532875.
- 1152 43 Robert J. Simpson and John Kessenich, editors. *The OpenGL ES Shading Language Version 3.20.6*. The Khronos Group, 2019. https://www.khronos.org/registry/OpenGL/specs/es/3.2/GLSL_ES_Specification_3.20.pdf.
- 1155 44 Chengnian Sun, Vu Le, and Zhendong Su. Finding compiler bugs via live code mutation. In
1156 Eelco Visser and Yannis Smaragdakis, editors, *Proceedings of the 2016 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications, OOPSLA 2016, part of SPLASH 2016, Amsterdam, The Netherlands, October 30 - November 4, 2016*, pages 849–863. ACM, 2016. URL: <https://doi.org/10.1145/2983990.2984038>, doi:10.1145/2983990.2984038.
- 1161 45 Chengnian Sun, Vu Le, Qirun Zhang, and Zhendong Su. Toward understanding compiler bugs
1162 in GCC and LLVM. In Andreas Zeller and Abhik Roychoudhury, editors, *Proceedings of the 25th International Symposium on Software Testing and Analysis, ISSTA 2016, Saarbrücken, Germany, July 18-20, 2016*, pages 294–305. ACM, 2016. URL: <https://doi.org/10.1145/2931037.2931074>, doi:10.1145/2931037.2931074.
- 1166 46 Qiuming Tao, Wei Wu, Chen Zhao, and Wuwei Shen. An automatic testing approach
1167 for compiler based on metamorphic testing technique. In Jun Han and Tran Dan Thu, editors, *17th Asia Pacific Software Engineering Conference, APSEC 2010, Sydney, Australia, November 30 - December 3, 2010*, pages 270–279. IEEE Computer Society, 2010. URL: <https://doi.org/10.1109/APSEC.2010.39>, doi:10.1109/APSEC.2010.39.
- 1171 47 Stéfan van der Walt, Johannes L. Schönberger, Juan Nunez-Iglesias, François Boulogne,
1172 Joshua D. Warner, Neil Yager, Emmanuelle Gouillart, Tony Yu, and the scikit-image
1173 contributors. scikit-image: image processing in Python. *PeerJ*, 2:e453, 6 2014. URL:
1174 <https://doi.org/10.7717/peerj.453>, doi:10.7717/peerj.453.

22:30 Putting Randomized Compiler Testing into Production

- 1175 **48** Brian A. Wichmann. Some remarks about random testing, 1998. Available online at <https://www.semanticscholar.org/paper/Some-Remarks-about-Random-Testing-Wichmann/2ad3c4c2e1b0b5867a1aa3e7c2de4a17d9facead>.
1176
1177
- 1178 **49** Brian A. Wichmann and Z. J. Ciechanowicz, editors. *Pascal Compiler Validation*. John Wiley
1179 & Sons, Inc., 1983.
- 1180 **50** Xuejun Yang, Yang Chen, Eric Eide, and John Regehr. Finding and understanding bugs in
1181 C compilers. In Mary W. Hall and David A. Padua, editors, *Proceedings of the 32nd ACM*
1182 *SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2011,*
1183 *San Jose, CA, USA, June 4-8, 2011*, pages 283–294. ACM, 2011. URL: <https://doi.org/10.1145/1993498.1993532>, doi:10.1145/1993498.1993532.
1184
- 1185 **51** Andreas Zeller and Ralf Hildebrandt. Simplifying and isolating failure-inducing input. *IEEE*
1186 *Trans. Software Eng.*, 28(2):183–200, 2002. URL: <https://doi.org/10.1109/32.988498>,
1187 doi:10.1109/32.988498.