



Challenges in Migrating Imperative Deep Learning Programs to Graph Execution: An Empirical Study

Tatiana Castro Vélez

City University of New York (CUNY) Graduate Center
New York, NY, USA
tcastrovelez@gradcenter.cuny.edu

Mehdi Bagherzadeh

Oakland University
Rochester, MI, USA
mbagherzadeh@oakland.edu

Raffi Khatchadourian

City University of New York (CUNY) Hunter College
New York, NY, USA
raffi.khatchadourian@hunter.cuny.edu

Anita Raja

City University of New York (CUNY) Hunter College
New York, NY, USA
anita.raja@hunter.cuny.edu

ABSTRACT

Efficiency is essential to support responsiveness w.r.t. ever-growing datasets, especially for Deep Learning (DL) systems. DL frameworks have traditionally embraced *deferred* execution-style DL code that supports symbolic, graph-based Deep Neural Network (DNN) computation. While scalable, such development tends to produce DL code that is error-prone, non-intuitive, and difficult to debug. Consequently, more natural, less error-prone imperative DL frameworks encouraging *eager* execution have emerged at the expense of run-time performance. While hybrid approaches aim for the “best of both worlds,” the challenges in applying them in the real world are largely unknown. We conduct a data-driven analysis of challenges—and resultant bugs—involved in writing reliable yet performant imperative DL code by studying 250 open-source projects, consisting of 19.7 MLOC, along with 470 and 446 manually examined code patches and bug reports, respectively. The results indicate that hybridization: (i) is prone to API misuse, (ii) can result in performance *degradation*—the opposite of its intention, and (iii) has limited application due to execution mode incompatibility. We put forth several recommendations, best practices, and anti-patterns for effectively hybridizing imperative DL code, potentially benefiting DL practitioners, API designers, tool developers, and educators.

CCS CONCEPTS

• **General and reference** → **Empirical studies**; • **Computing methodologies** → **Machine learning**; • **Software and its engineering** → *Language features*; *Software evolution*.

KEYWORDS

empirical studies, deep learning, imperative programs, hybrid programming paradigms, graph-based execution, software evolution

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

MSR '22, May 23–24, 2022, Pittsburgh, PA, USA

© 2022 Association for Computing Machinery.

ACM ISBN 978-1-4503-9303-4/22/05...\$15.00

<https://doi.org/10.1145/3524842.3528455>

ACM Reference Format:

Tatiana Castro Vélez, Raffi Khatchadourian, Mehdi Bagherzadeh, and Anita Raja. 2022. Challenges in Migrating Imperative Deep Learning Programs to Graph Execution: An Empirical Study. In *19th International Conference on Mining Software Repositories (MSR '22)*, May 23–24, 2022, Pittsburgh, PA, USA. ACM, New York, NY, USA, 13 pages. <https://doi.org/10.1145/3524842.3528455>

1 INTRODUCTION

Machine Learning (ML), including Deep Learning (DL), systems are pervasive in society. Central to such systems are dynamic models, whose behavior is ultimately defined by input data. However, as datasets grow, efficiency becomes essential to support responsiveness [103]. For industrial applications, DL frameworks—pillars of DL systems [56,58,68,100]—must quickly execute complex computations on large datasets while supporting easy-to-use programming paradigms [60]. For efficiency, DL frameworks have traditionally embraced a *deferred* execution-style that supports symbolic, graph-based Deep Neural Network (DNN) computation [25,46]. While scalable, development is error-prone, cumbersome, and produces programs that are difficult to debug [56,57,99,100]. Furthermore, because graph computation executes statements in a non-imperative order, traditional Software Engineering (SE) tools cannot help troubleshoot bugs [9]. Contrarily, more natural, less error-prone, and easier-to-debug *imperative* DL frameworks [3,27,79] encouraging *eager* execution have emerged. Though ubiquitous, eagerly-executed imperative DL programs are less efficient and scalable as their deferred-execution counterparts [25,37,43,60,72,79]. Executing (imperative) DL programs eagerly “makes tensor [matrix-like data structures central to DL] evaluation trivial but at the cost of lower performance” [30].¹ Thus, hybrid approaches [6,37,72]—integrated into mainstream DL frameworks—execute imperative DL programs as static graphs at run-time. For example, in *TensorFlow* [1]—a popular [54,100] DL framework—*AutoGraph* [72] can potentially enhance performance by decorating (annotating)—with optional yet influential decorator arguments—appropriate Python function(s) with `@tf.` function. Decorating functions with such hybridization Application Programming Interfaces (APIs) can increase imperative DL code performance without explicit modification.

Though promising, hybrid approaches necessitate non-trivial specialized metadata [60] and exhibit limitations and known issues [42] with native program constructs. Subtle considerations

¹Performance in this paper refers to *run-time* performance (speed), not model accuracy.

are required to make code amenable to safe, accurate, and efficient graph execution—avoiding performance bottlenecks and semantically inequivalent results. Therefore, developers are burdened with making their code compatible with the underlying execution model conversion, as well *manually* specifying which functions should be converted. While alternatives [60] exist, they impose custom Python interpreters, which may be impractical for industry, and support only specific Python constructs. Thus, there is a knowledge gap in how hybridization is used in real-world DL applications, leading to the challenges in successfully applying it underexplored. Without such insight, DL systems may be inefficient, fallible, and difficult to maintain. Moreover, advances in DL are likely to be futile if they cannot be effectively used.

To fill this gap, we conduct an empirical study on common development challenges in migrating imperative DL code to graph execution using hybridization in open-source DL systems. Particularly, we aim to answer the following research questions: (RQ1) what bug patterns and corresponding challenges are involved in writing *reliable* yet *performant* imperative DL code, and (RQ2) which best practices and anti-patterns can be extracted from (RQ1)? Such knowledge can help drive new automated migration techniques, IDE code completion, and automated (data science-specific [11,32,33]) refactoring mining approaches [95]. The results: (i) advance knowledge of this emerging yet pervasive hybrid paradigm, (ii) provide feedback to language and API designers for future API versions, (iii) help tool designers comprehend difficulties with writing performant imperative DL code, (iv) include preliminary *recommendations*, *best practices*, and *anti-patterns* for practitioners in using hybridization effectively, and (v) assist educators in teaching hybridization APIs.

Our study involves analyzing occurrences of `tf.function` in 250 projects, consisting of 19.7 MLOC, along with 470 and 446 manually examined code patches (Git commits) and bug reports (GitHub issues), respectively. Challenges—along with their causes, symptoms, and fix patterns—are taxonomized using manual processes aided by automated software repository mining. Due to its popularity and extensive analysis by previous work [26,55,56,58,68,74,99,100], we focus on hybridization in *TensorFlow*. Our study indicates that: (i) `tf.function` is widely used, (ii) misusing `tf.function` was a major theme in migrating imperative DL programs to graph execution, (iii) subtle bugs in using `tf.function` can result in performance *degradation*—the opposite of its intention, and (iv) `tf.function` is commonly incompatible in a given context—limiting its application.

Our contributions can be summarized as follows:

Hybridization bug hierarchical taxonomy From 470 and 446 patches and bug reports, respectively, of 250 projects manually examined, we build a rich hierarchical taxonomy of common hybridization challenges.

Recommendations, best practices, & anti-patterns We propose preliminary recommendations, best practices, and anti-patterns for effectively hybridizing imperative DL code from our statistical results, as well as an in-depth analysis.

Complete results of our study are available in our dataset [24].

2 MOTIVATING EXAMPLES & BACKGROUND

Popular DL frameworks have historically embraced *deferred* execution-style (low-level) APIs, making DNNs straight-forward to

```

1 # Build a graph.
2 a = tf.constant(5.0)
3 b = tf.constant(6.0)
4 c = a * b

5 # Launch graph in a session.
6 sess = tf.Session()
7 # Evaluate the tensor 'c'.
8 print(sess.run(c)) # prints 30.0

```

Listing 1: *TensorFlow* deferred execution-style code [48].

```

1 class SequentialModel(tf.keras.Model):
2     def __init__(self, **kwargs):
3         super(SequentialModel, self).__init__(...)
4         self.flatten = layers.Flatten(
5             input_shape=(28, 28))
6         num_layers = 100 # Add many small layers.
7         self.layers = [layers.Dense(64, activation =
8             "relu") for n in range(num_layers)]
9         self.dropout = tf.keras.layers.Dropout(0.2)
10        self.dense_2 = tf.keras.layers.Dense(10)
11 @tf.function(...) # Executes
12 # the model as a graph
13 # (with optional args).
14 def __call__(self, x):
15     x = self.flatten(x)
16     for layer in self.layers:
17         x = layer(x)
18     x = self.dropout(x)
19     x = self.dense_2(x)
20     return x

```

Listing 2: *TensorFlow* imperative (OO) DL model code [43].

```

1 @tf.function
2 def f(x):
3     print("Input: ", x)
4     f(1)
5     f(1)
6     f(2)

```

Output (expecting 1, 1, 2):

```

Input: 1
Input: 2

```

Listing 3: Imperative *TensorFlow* code with Python side-effects [42].

execute as symbolic graphs that enable various run-time optimizations. For example, during graph building (lines 2–4 of listing 1), line 4 does not execute until the `Session` created on line 6 is run on line 8. While efficient, legacy code using such APIs are cumbersome, error-prone, and difficult to debug and maintain [56,57,99,100]. Such APIs also do not natively support common imperative program constructs, e.g., iteration [5]. Contrarily, *eager* execution-style DL APIs [3,79] facilitating higher-level, imperative, and Object-Oriented (OO) [27] (Python) programs that are easier-to-debug, less error-prone, and more extensible have emerged. For instance, with *eager* execution, line 4 of listing 1 would execute and immediately evaluate tensor `c`, foregoing the need of a session. In many DL frameworks, *eager* execution is now the default.

Despite the benefits, executing (imperative) DL programs eagerly comes at the cost of run-time performance [30]. Thus, hybridization approaches [6,37,72] that execute imperative DL programs as graphs at run-time have been integrated into mainstream DL frameworks. For example, listing 2 portrays *TensorFlow* imperative (OO) DL code representing a modestly-sized model for classifying images. On line 11, *AutoGraph* [72] is used to potentially improve performance by decorating the model’s `call()` method with `@tf.function`, possibly providing optional yet influential decorator arguments. At run-time, `call()`’s execution will be “traced” and an equivalent graph will be generated [42]. In this case, a speedup ($\text{runtime}_{old}/\text{runtime}_{new}$) of ~ 9.22 , averaged over five runs, ensues [63].

As noted in Section 1, while promising, hybridization presents unique challenges [42,60] in ensuring that programs run reliably *and* efficiently. If used incorrectly, hybridization may yield programs that result in unexpected run-time behavior. Decorating the right functions, supplying the correct decorator arguments, using the appropriate API, and properly structuring imperative DL code so that it is amenable to graph execution can be daunting, especially for developers (data scientists) lacking SE expertise.

Python Side-effects. Side-effect producing, native Python statements, e.g., printing, list appending, global variable mutation [42],

```

1 class Model(tf.Module):
2     def __init__(self):
3         self.v = tf.Variable(0)
4         self.counter = 0
5
6     @tf.function
7     def __call__(self):
8         if self.counter == 0:
9             self.counter += 1
10            self.v.assign_add(1)
11            return self.v

```

Output (expecting 1, 1, 1):

```

1
2
3

```

Listing 4: Imperative TensorFlow code using a counter [42].

```

1 model = SequentialModel()
2 res1 = model(tf.constant([1, 2, 3]))
3 res2 = model(tf.constant([1, 2, 3, 4, 5]))

```

WARNING: 5 of the last 5 calls triggered ... retracing. Tracing is expensive.

Listing 5: DL model (listing 2) client code using varying datasets [42].

are problematic for `tf.function`-decorated functions.² Because they are traced, a function’s behavior is “etched” into its corresponding graph and thus can have unexpected results, executing side-effects multiple times or not at all. Side-effects occur when `tf.functions` are called the first time; subsequent calls with similar arguments execute the graph instead. For example, on line 3 of listing 3, `f()` outputs `x`. On line 1, `f()` is decorated with `@tf.function`, which migrates it to a graph at run-time. Then, `f()` is invoked three times, the first two with the argument 1 and the last with 2. In the output on the right, the first invocation of `f()` on line 4 results in a graph being built (through tracing) that—due to a similar argument—is later used on line 5. Consequently, the side-effecting code on line 3 is *not* exercised. In contrast, line 3 is exercised as a result of the call on line 6 due to a different argument being supplied.

Although listing 3 is simple, unexpected behavior can generally be difficult to notice. Consider listing 4, where a model uses a counter to safeguard a variable incrementation. The initial value of counter, however, is captured during tracing upon the first model invocation (line 14). The overall effect is that the value of `v` is incremented *unconditionally* (line 10) each time the model is invoked. Such problems are common in migrating deferred-execution-style DL code (e.g., listing 1) to an imperative style (e.g., listing 2). Worse yet, developers only realize such errors after observing suspicious numerical results or significantly lower performance than expected (e.g., when guarded operations are costly) [42].

When To Use Hybridization? Besides ensuring that DL code is amenable hybridization [36], developers must also know *when* and *where* to use it to avoid performance bottlenecks and other undesired behavior. For example, confusion exists on how often `@tf.function` should be applied [87], and calling `tf.functions` recursively could cause infinite loops [42]. Even if a recursion seems to work, the `tf.function` will be traced *multiple* times (“retracing”), potentially impacting performance. Also, using `@tf.function` on small computations can be dominated by graph creation overhead [43].

Using Hybridization Parameters. Decorating the *correct* function but with *incorrect* decorator arguments may result in performance degradation. For instance, retracing helps ensure that the correct graphs are generated for each set of inputs; however, excessive retracing may cause code to run more slowly had `tf.function` *not* been used [42,80,81]. Listing 5 depicts code that invokes the model

²Herein, “`tf.function`-decorated” functions will be referred to as “`tf.functions`.”

Table 1: Studied subjects.

	subj	KLOC	studied periods	cmts/iss	kws	exe
fixes	122	10,879	2015-11-06 to 2021-01-14	199,140	470	470
reports	167	17,378	2012-05-07 to 2021-08-11	237,232	704	446
Total	250*	19,677*	2012-05-07 to 2021-08-11	436,372	1,174	916

* Represents unique totals due to subject overlap between the study portions.

declared in listing 2 multiple times using different (hypothetical) datasets, producing the warning on the right. To limit retracing, an `input_signature` can be specified on line 11, listing 2 as follows:

```
@tf.function(input_signature=(tf.TensorSpec(shape=[None], dtype=tf.int32),))
```

A `[None]` dimension in the `tf.TensorSpec` allows for flexibility in trace (graph) reuse. Since tensors are matched on their shape, a `None` wild card allows `tf.functions` to reuse traces for variably-sized input—occurring when sequences or images are of different lengths or sizes, respectively. Since each call no longer produces a trace, the warning disappears—averting any performance bottlenecks.

These simplified examples demonstrate that effectively using hybridization is not always straight-forward, potentially requiring complex analyses and a thorough understanding of API intricacies—a compounding problem in more extensive programs. As imperative DL programming becomes more widespread, statistical insight into how such programs are best written efficiently and how to avoid common bugs would be extremely valuable to developers.

3 METHODOLOGY

Subjects. We examined Git commit changesets (code patches; row **fixes**, Table 1) representing bug fixes involving `tf.function` and GitHub issues (row **reports**) mentioning `tf.function`. Our study encompassed 250 open-source DL systems (column **subj**), comprising ~19.7 million lines of source code (column **KLOC**), 199,140 Git commits (column **cmts** for **commits**), 237,232 GitHub issues (column **iss** for bug **reports**), and 460.21 years of combined project history, averaging 1.86 years per subject. Subject details may be found in our dataset [24]; subjects sources are publicly available on GitHub. While we focus `tf.function` client usages, we include *TensorFlow* as developers often file GitHub issues against it to discuss `tf.function` usage challenges and potential bugs. Subjects include those used in previous studies [26,32,55,56,58,59,68,99,100] and appearing in data science-specific datasets [17]. To determine if a project represents a DL system, i.e., one with at least one DL module, we searched repositories for specific keywords, e.g., “keras,” “layer,” “net,” “neural network,” “deep learning.” The keywords have also been used in related work [59] for a similar purpose; the keywords were only used to ensure that subjects were DL systems, not for finding hybridization bugs. We then verified the code to ensure that the keywords represented DL contexts.

For changesets (bug fixes), subject criteria consists of having at least one commit whose changeset contains `tf.function`. For issues, subjects must have at least one GitHub issue mentioning “`tf.function`.” Subjects were mostly written in Python, which is popular for DL [16]. While the subjects include popular open-source repositories from well-known and reputable organizations, e.g., Apache [7], Apple [8], Google [45], NVIDIA [75], they also include lesser-known repositories to understand hybridization challenges

facing the DL community-at-large. Furthermore, hybridization is relatively new—`tf.function` was released on September 30, 2019.

Mining. To find changesets (patches) representing hybridization bug fixes, we mined repositories for commits referencing `tf.function` using `gitcproc` [23], a tool for classifying Git commits used by previous work [12,65,92,94]. Row **fixes**, column **kws** of Table 1 is the commits containing `tf.function` in their changesets. We manually examined all 470 commits, portrayed by row **fixes**, column **exe**. To find issues related to hybridization, we mined repositories for GitHub issues mentioning “`tf.function`” by first filtering out issues containing only irrelevant discussion (e.g., “social conversation”) using a pre-trained classification model [10] used by previous work [76,98,102]. We then invoked the *GitHub Search API* [41] to select (open and closed) issues that included “`tf.function`” using several different criteria, e.g., “best match,” “most commented.” To reduce false positives, since the API ignores punctuation, we further filtered the results to ensure that they included the period. Row **reports**, column **kws** of Table 1 is the issues³ containing “`tf.function`” in either their title or body (description and conversations). We randomly selected a subset of these to examine manually (details below), portrayed by row **reports**, column **exe**. The aforementioned tools [10,23,41] were only used to narrow the search space and not for classification, which was done manually. The GitHub search was performed in a (standard) manor consistent with previous work.

Identification. We used a `gitcproc` feature that leverages heuristics based on log messages to identify bug fix commits. Natural language processing (NLP) is internally used by `gitcproc` to determine the commits that fall into this category. Doing so helps us to focus on likely bug fix commits for further manual examination. Random matching issues—with ones containing code being favored—were chosen for manual inspection. Next, the authors manually examined the commits and issues to ascertain if they indeed relate to hybridization bugs. Two authors are SE and PL professors with extensive expertise in software evolution, system performance, and empirical SE. Another author is a data mining and ML professor with substantial proficiency in AI and SE. Three authors have several years of industrial SE experience.

Although the researchers did not converse during the initial identification and classification process to avoid bias, this mix of expertise is effective in studying SE tasks in DL systems. The researchers convened regularly during the study, as well as at the end for finalization, to solidify the results. Cohen’s Kappa coefficients [96] for identification and classification were 0.80 and 0.57, respectively.⁴ As the authors did not always have detailed knowledge of the particular systems, only changes where a bug fix was extremely likely were marked as such. The authors also used commit comments and referenced bug databases to ascertain whether a change was a bug fix. GitHub issues tags were also considered.

Classification. For commits, once bug fixes were identified, the authors studied the code changes to determine the category of bug fixes and whether the category relates to hybridization. For issues, the authors examined issue descriptions and discussions, paying attention to the `tf.function` challenges being described and their possible solutions and workarounds. Particular attention was paid

³Also includes pull (patch) requests as these are treated similarly in GitHub.

⁴Moderate agreement is expected; the team has mixed ML/SE expertise.

Table 2: Discovered top-level problem categories.

problem	abbr	cmts	iss	total
Performance	PRF	74	37	111
API misuse	APM	23	30	53
Incompatibility	INC	16	33	49
TensorFlow bug	TFB	4	18	22
Other	OTH	14	2	16
Unknown	UKN	10	0	10
Test	TST	8	0	8
Debuggability	DBG	4	2	6
Exposed variable state	EVS	1	1	2
Compilation error	CMP	1	0	1
Numerical errors	NME	1	0	1
Segmentation fault	SEG	1	0	1
Total		157	123	280

to code snippets. No scripts were involved in the classification—only manual examination. Categories were then formed into a hierarchy, in part by using the *TensorFlow* documentation [42]. On several occasions, developers were contacted for clarification using the GitHub line comment mechanism and via email.

4 RESULTS

This section answers (RQ1) by summarizing our results, noting trends, exceptions, and unexpected outcomes. Contrarily, Section 5 consolidates, comments on, and connects the main findings. Related discussion in Section 5 is referenced where appropriate.

4.1 Quantitative Analysis

From the 470 commits and 446 GitHub issues (totaling 916) manually examined (column **exe**, Table 1), we found 157 and 123 (totaling 280) `tf.function` bug fixes and developer challenges depicted in columns **cmts** (commits) and **iss** (GitHub issues) of Table 2, respectively. Finding these bugs and understanding their relevance required a significant amount of manual labor that may not be feasible in more large-scale, automated studies. Python, being a dynamic language, can be difficult to analyze, particularly w.r.t. inheritance relationships; subclassing Keras models is a common way to write imperative DL code in *TensorFlow* (cf. line 1, listing 2). Furthermore, our number of findings (280) is comparable with previous studies involving manual inspection (e.g., Tang et al. [91] found 285, Zhang et al. [100] found 175, Khatchadourian et al. [65] found 61). Nevertheless, as `tf.function` becomes more popular, we expect its usage and number of related bugs to grow.

4.1.1 Problem Categories. We group bug fixes and GitHub issues into common problem categories, shown in Fig. 1 and Table 2 (column **abbr** is the category abbreviation). The former includes combined data (commits and issues), while the latter separates the two. Figure 1 presents a hierarchical categorization—with varying levels of detail—of the 280 discovered `tf.function`-related challenges in our subjects. Challenges are represented by their problem category name and are followed by their counts. Categories without instances are *abstract*, i.e., they *only* group together other categories. Table 2 portrays a nonhierarchical, top-level view of Fig. 1; the innermost (top) layers of Fig. 1 represent the rows of Table 2.

Challenges are grouped into several (top-level) problem categories. Categories include performance (PRF, 90; further discussed

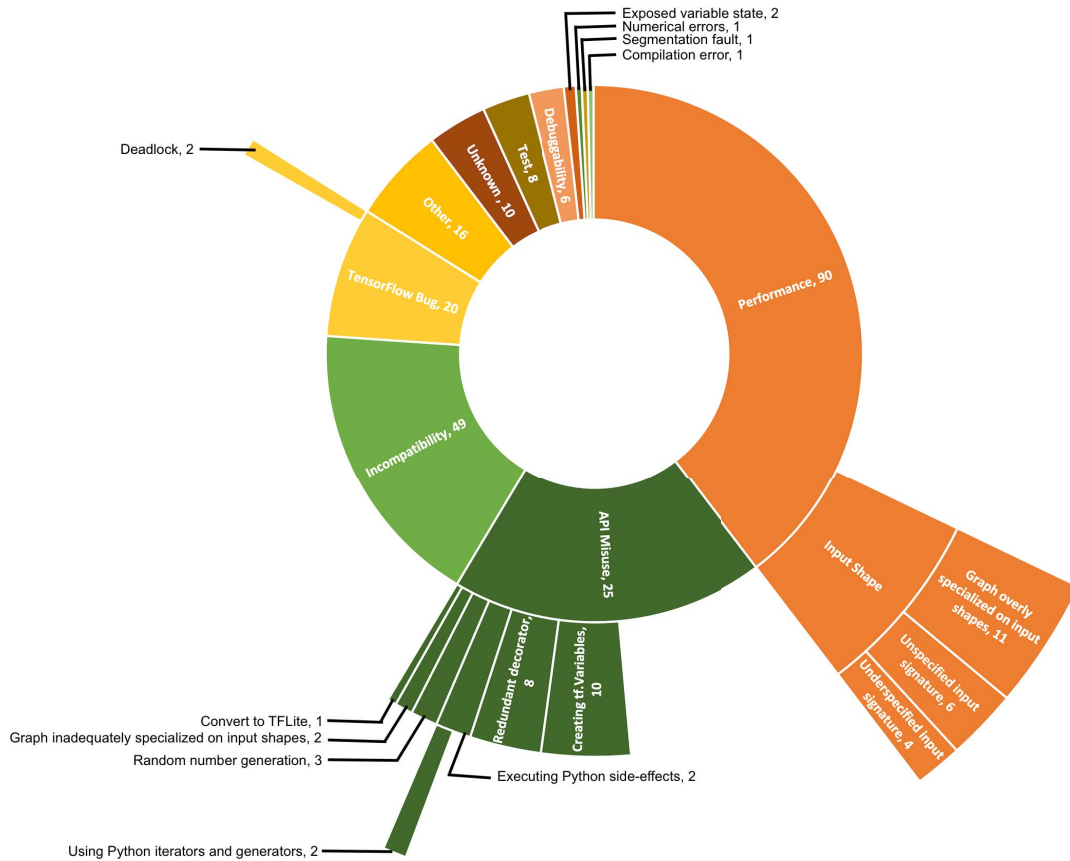


Figure 1: Discovered problem categories (hierarchical).

later), API misuse (APM, 25; further discussed later), and incompatibility between execution modes, i.e., eager and deferred, where `tf.function` is used in a context not amenable to graph conversion (INC, 48; further discussed later). An example of the latter is where particular loss functions cannot be used in graph mode or there is an *AutoGraph* limitation that prevents graph conversion. Other problem categories include dealing with or working around open bugs related to `tf.function` in *TensorFlow* (TFB, 20; further discussed later) and “other” (OTH, 16), which involves syntactic corrections, general cleanup, and refactorings—a category similar to that used by previous work [65,94]. “Unknown” (UKN, 10) represents situations where the problem category was indeterminable without further domain knowledge or developer input. Only 3.57% of problems had unknown categories. Code changes involving `tf.function` appearing in tests were categorized as “Test” (TST, 8).

Debuggability. Debuggability (DBG, 6) represent situations where using `tf.function` to improve performance of DL code may, in turn, reduce a developer’s ability to easily debug it. “In general, debugging code is easier in eager mode than inside `tf.function`” [42]. In such situations, developers may not understand that using `tf.function` is the reason why they are not able to debug their code, e.g., intermediate variable values may be missing. Or, `tf.function` may temporarily be removed (via a commit) to facilitate debugging, but

developers inadvertently neglect to replace it (cf. Section 4.2.5). This latter situation is unfortunate as, to assist in the debugging process, a flag can be used to globally (temporarily) toggle `tf.function` [42].

Other Categories. Other (top-level) categories were more minor in terms of their counts, yet have potentially significant consequences. For example, exposed variable state (EVS, 2) occurs when saving (exposed) program state (variables) is problematic during `tf.function` conversion at run-time, e.g., variables becoming undefined [71]. Numerical errors (NME, 1) involve possible numeric overflow. *Autograph* compilation errors (CMP, 1) surface when `tf.functions` are compiled and subsequently result in compilation errors. This problem may arise when certain dynamic Python features, e.g., lexical scoping, are utilized (cf. Section 4.2.2). Segmentation fault (SEG, 1) is when using `tf.function` causes a program crash. While compilation and numerical errors and segmentation faults may be considered symptoms, we focus on `tf.function` client usage; these categories represent problems from a client perspective. Their underlying causes are bugs within the framework.

Performance. As the main purpose to hybridization is to improve the performance of imperative style DL code by building a bridge to graph-based execution, it was not surprising that performance—at 39.64% (111/280)—was the largest category:

Table 3: Performance fixes.

fix category	count
Add <code>tf.function</code> decorator	61
Change <code>tf.function</code> argument	20
Add <code>input_signature</code> argument to <code>tf.function</code>	9
Remove <code>tf.function</code> decorator	8
Upgrade to new library version	4
Relocate <code>tf.function</code> (use on different function)	5
Re-add <code>tf.function</code> decorator	2
Unsolved (open)	2
Total	111

Finding 1: At 39.64% (111/280), performance was the largest problem category encompassing `tf.function` usage.

Performance problems represent a spectrum of situations, stemming from using `tf.function` to solve a DL code performance bug to not observing the expected speedup from using `tf.function` to exhibiting *worse* performance that *not* using `tf.function`. Table 3 portrays the various fixes used to solve performance problems. Though the majority of times it was used to enhance performance of imperative DL code, we found that in 7.21% (8/111) of cases, `tf.function` was *removed* to alleviate performance problems:

Finding 2: Despite intent to improve performance, `tf.function` caused performance degradation in 7.21% (8/111) of cases.

Moreover, only 54.95% of imperative DL code performance problems were fixed by *adding* `tf.function`. Thus, the remaining 45.05% of cases were due to *existing* hybridization:

Finding 3: Only 54.95% (61/111) of imperative DL code performance problems were fixed by *adding* `tf.function`. The remaining 45.05% were due to *using* `tf.function`.

In fact, 25.23% of performance fixes involved altering `tf.function` arguments:

Finding 4: Performance fixes entailed altering developer-supplied `tf.function(...)` arguments at a rate of 25.23%.

Performance problems are further categorized into those related to “input shapes,” which make up 18.92% of all such problems:

Finding 5: Performance problems involved incorrect input tensor shape specifications at a rate of 18.92%.

Tensors are heavily used DL programs, and accurately matching tensor shapes (dimensions) is often required to write reliable DL code. In hybridization, since `tf.functions` are being traced and thus converted into graphs, the underlying framework (by default) attempts to build specialized graphs for each kind of input. However, when tensors are involved, graphs may be specialized to particular input shapes, creating a situation where function retracing is excessive. Retracing can lead to significant performance degradation [43].

To curb this problem, an (`input_signature`) argument may be supplied to `tf.function` that specifies an *expected* range of shapes. In effect, developers provide contextual information to the framework about how `tf.functions` will be used. For instance, setting `experimental_relax_shapes` to `True` may cause `tf.functions` to generate fewer graphs that are less specialized on input shapes. However, this may not match reality, especially when dealing with dynamic shapes. As such, we further divide “input shape” challenges: **Graph overly specified on input shapes (11)** Generated graphs are too specific for the context where a `tf.function` is being used, which can occur when either:

Table 4: API misuse causes.

cause	count
API confusion	20
Use of graph mode	14
Decorated outer function calls unnecessarily decorated inner function	8
Incorrect <code>tf.function</code> argument	7
Use of eager mode	2
Lost variable state due to graph conversion	1
Lack of static shape specifications	1
Total	53

- (i) `experimental_relax_shapes` is incorrectly set to `False`.
- (ii) `input_signature` is unnecessarily specified. Either it should be either removed or set to `None` (the default).

Underspecified input signature (4) The `input_signature` parameter lacks proper arguments to avoid excessive retracing.

Unspecified input signature (6) The `input_signature` is missing in contexts that are advantageous to graph specialization.

API Misuse. API Misuse—the second largest problem category at 18.93%—involves situations where `tf.function` is not used in a way recommended by the API documentation:

Finding 6: At 18.93%, API misuse—using `tf.function` inconsistent to documentation—was the 2nd largest category.

Misusing APIs typically results in either run-time errors or unexpected behavior. Violating DL API constraints may lead to crashes and poor performance [56,58]. In high-level, e.g., imperative DL, code, bugs are commonly due to misunderstandings of the guarantees offered and obligations imposed by increasingly layered software, e.g., those written against the *TensorFlow* API [66]. *TensorFlow* documentation contains a prominent sections regarding `tf.function` and *AutoGraph* usage constraints and limitations. If such constraints, e.g., w.r.t. control-flow, side-effects, global variables, are violated, *AutoGraph* will not properly generate graphs from Python code. Despite the vast documentation, at 37.74%, API confusion was the largest cause of API misuse (Table 4):

Finding 7: API misuse was caused by developers not understanding hybridization APIs at a rate of 37.74% (20/53).

Regarding potential category overlap, recall that API misuse is defined above as a violation of intended API usage per the documentation. Consider changing a `tf.function` argument. Performance degradation can occur when parameter usage is *consistent* with the documentation; it can be a tuning issue, e.g., shape-related. In such a case, according to the earlier definition, shape mismatches would not be considered API misuse as they are dependent on context.

We found that the most common way (28.30% or 15/53) to fix API misuse was to *remove* `@tf.function`. Of these, in 46.67% of cases (7/15), the problem cause was that `@tf.function` was used to decorate an *inner* function called by an *already* decorated *outer* function. As `tf.function` applies to the decorated function *and* all other functions it calls and since the inner function cannot be called from any other function besides the outer function, the inner function decorator is unnecessary and can thus be safely removed [43]. However, another 46.67% (7/15) of cases were caused by API confusion. Thus, in these cases, unfortunately, developers *abandoned* `@tf.function`—along with its potential to enhance performance—due to their confusion over how to use it. Most likely, developers were doing so to avoid run-time errors, which occurred in 62.50%

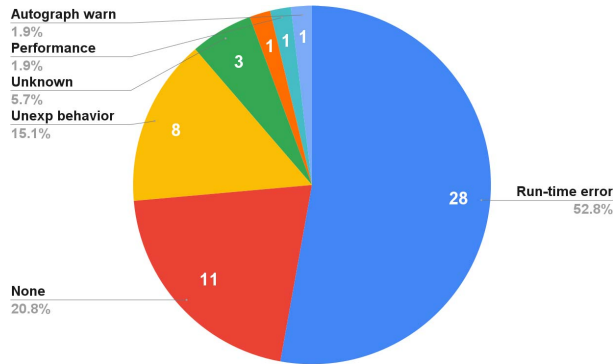


Figure 2: API misuse symptoms.

(5/8) of `tf.function` removals not caused by unnecessary inner function decoration and 52.83% (28/53) overall (Fig. 2):

Finding 8: To fix API misuse, `tf.function` was removed 28.30% of the time. In 46.67% of these, hybridization was abandoned due to API confusion, with 62.50% causing run-time errors.

API misuse is further divided into several categories, the largest of which involves creating `tf.Variables` within `tf.functions` (10). A `tf.Variable` represents a tensor whose value is mutable [44]. Currently, `tf.function` only supports singleton `tf.Variables`; creating multiple `tf.Variables` within the scope of a `tf.function` results in a run-time exception [42]. Redundant decoration (8) is where multiple functions on a call path are unnecessarily decorated with `@tf.function`; all functions called from a `tf.function` are also automatically migrated to graphs. Accurately approximating such paths statically—especially in the context of a dynamic language such as Python—may be difficult, and there is ample confusion among developers on where to apply `@tf.function` [87] (cf. Section 2).

Executing Python side-effects (2) refers to the situation where `tf.functions` contain side-effect producing Python statements. As described in Section 2, executing such statements within migrated graphs can have unexpected results, sometimes executing twice or not all. A specific pattern of side-effects were those involving the use of iterators and generators (2), a common looping mechanism in Python code. Random number generation (RNG, 3) problems occur when developers do not use RNG facilities consistently with the documentation, commonly resulting in unexpected behavior under graph mode. For example, RNG creation inside a `tf.function` can only happen during the first run of the function [47]. Seeding may also not work as expected in graph mode (e.g., [89])—“when [a] global seed is set but [TensorFlow] operation seeds are not, the sequence of random numbers are the same for each `tf.function`” [51]. “Graph inadequately specialized on input shapes” (2) involves an API misuse that is opposite to the “graph overly specified on input shapes” performance problem category described earlier. Such problems may be fixed by setting `experimental_relax_shapes` to `False` (the default). In other words, the shape specification is too general, which may result in a situation that is not amenable to graph migration [38]. For example, an `input_signature` may be supplied using a wild card shape to improve performance (q.v. Section 2) but results in a run-time error due to a tensor dimension mismatch [42,83].

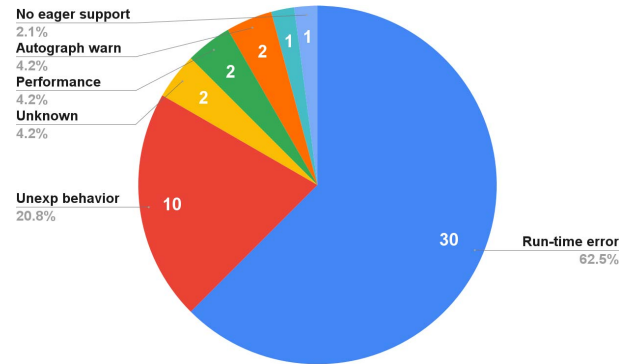


Figure 3: Incompatibility problem symptoms.

Conversion to TFLite (1) represents problems with an alternate use case of `tf.function` to convert a DL model to a portable format.

Execution Mode Incompatibility. At 17.50%, incompatibility is the third largest problem category:

Finding 9: Execution mode incompatibility, at 17.50% (49/280), was the 3rd largest problem category, meaning that seamlessly using similar constructs in different modes was problematic.

Developers seemingly struggle with seamlessly using imperative DL program constructs, e.g., particular loss functions, across execution modes. Ideally, developers could toggle between eager and graph execution modes—with *AutoGraph* simply enhancing performance—without making code changes. In other words, incompatibility problems prevents developers from focusing on the correctness of their DL code—thinking of performance as an afterthought. Instead, to use hybridization effectively, developers must be cognizant of its internal structure, i.e., *how* their DL code is being migrated to graphs. Moreover, developers must (manually) be aware of which constructs are amenable to graph conversion, how best to write code that works in either mode, and how to interact with code that may be executed in a different mode.

Execution mode incompatibility problems have dire consequences. As shown in Fig. 3, 81.63% of symptoms resulting from incompatibility involve run-time errors or unexpected behavior. Such problems that only occur at run-time are difficult to uncover and, if found, may be found after deployment:

Finding 10: Incompatibility problems led to run-time errors or unexpected results, which do not surface until *after* running the code, 81.63% of the time.

TensorFlow Bugs. TensorFlow bugs (TFB) made up 7.86% of bugs:

Finding 11: TensorFlow bugs, where developers were offered workarounds or awaited new framework versions, made up 7.86% of problems. Of these, 9.09% involve deadlocks.

Such bugs involve dealing with or working around open TensorFlow bugs related to `tf.function`. As hybridization is relatively new, the `tf.function` API is under active development. Thus, it was not uncommon for such bugs to be reported to TensorFlow by filing issues against its GitHub repository; 81.82% of TFBs appear as GitHub issues (see Table 2 and Fig. 4). We categorize bugs as TFB if they were in fact real bugs with TensorFlow that required a workaround—often suggested by TensorFlow contributors—or a new TensorFlow library

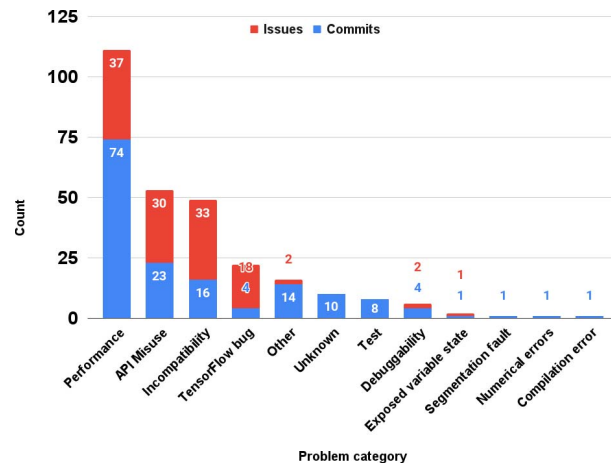


Figure 4: Top-level problem category comparison.

```

1 + @tf.function
2 def pm(linear):
3     state = lpt_init(linear, a0=0.1, order=1)
4     final_state = nbody(state, stages, nc)
5     tfinal_field = cic_paint(tf.zeros_like(linear), final_state[0])
6     return tfinal_field

```

Listing 6: Commit af1664e7 in galference: bug boxsize=nc

version to solve. If the reported bugs were not resolved to be the result of problems with *TensorFlow*, such bugs were not categorized as TFB but perhaps other categories.

TFB is further categorized into deadlock (2). Situations leading to the execution of a `tf.function` being deadlocked include using tensors as stopping condition of a recursive `tf.function` [29]. Deadlock may also occur as a result of other, specific `tf.function` code patterns—causing the *TensorFlow* run time to deadlock. For example, deadlock may occur when calling a `tf.function` from within a `tf.py_function` [21], which executes native Python functions as graph operations eagerly [50].

4.1.2 Commits vs. GitHub Issues. Figure 4 compares the different sources—commits (bottom/blue bars) and GitHub issues (top/red bars)—of problem categories. Performance—the largest problem area—was $\frac{2}{3}$ more likely to appear in *commits* vs. *issues*. In contrast, “incompatibility” was $\frac{2}{3}$ more likely to appear in *issues* vs. *commits*. Notably, 80% of TFB problems were found in GitHub issues compared to commits. Lastly, all UKN and TST bugs were found in commits. Section 5 discusses possible reasons for these differences.

4.2 Qualitative Analysis

This section answers (RQ2) by highlighting bug patterns with examples, summarizing causes, symptoms, and fixes, and proposing preliminary best practices and anti-patterns.

4.2.1 Performance. In listing 6, `pm()` is decorated with `@tf.function` (line 1). Using `tf.function` “ensure[s] that the graph for [a] function is compiled once and not every time it is called, thus gaining in speed and performance” [18], leading to best practice 1, Fig. 5.

- (1) Favor `@tf.function` on Python functions containing imperative, otherwise eagerly-executed, DL code to improve performance.
- (2) If possible, supply an `input_signature` argument to `tf.function` with the intended shape and types of any input tensors to avert retracing—a practice similar to that of providing type annotations to variables in dynamic languages to assist with type inferring.
- (3) When an operation is deemed incompatible with hybridization, check the documentation to see if additional steps are required to make the imperative DL code more amenable to graph conversion.
- (4) Framework limitations may impede performance enhancements. Check for potential workarounds of (unresolved) *TensorFlow* bugs.
- (5) Use `tf.config.run_functions_eagerly(True)` to temporarily disable `tf.function` to facilitate debugging.

Figure 5: Preliminary hybridization best practices.

- (1) Hybridizing nested functions may cause performance degradation. Sacrifice some modularity by either hybridizing the top-level function or refactoring the nested function to a top-level function [84].
- (2) Since shared variables must be singleton, using `tf.Variable`s in `tf.functions`, either directly or indirectly, may cause run-time exceptions. Either rewrite the function or do not hybridize it.
- (3) Since `tf.functions` are compiled, using dynamic language features, e.g., lexical scoping, either directly or indirectly, may lead to run-time exceptions. Avoid such features in `tf.functions` where possible.

Figure 6: Preliminary hybridization anti-patterns.

```

1 - @tf.function
2 + @tf.function(input_signature=[
3 +     tf.TensorSpec(shape=(None, self.num_states), dtype=tf.float32),
4 +     tf.TensorSpec(shape=(None, self.num_actions), dtype=tf.float32),
5 +     tf.TensorSpec(shape=(None, 1), dtype=tf.float32),
6 +     tf.TensorSpec(shape=(None, self.num_states), dtype=tf.float32),])
7 def update_weights(s, a, r, sn): # ...

```

Listing 7: Commit 02a3f297 in DDPG-tf2: Fixed all...this should work

While hybridization can enhance the performance of their imperative, otherwise eagerly-executed, DL code, we found that developers struggled to use it correctly. Some distrusted it, stating, e.g., that “it does far too much hidden magic” [2]. Others [82] struggled with uncontrolled retracing (q.v. Sections 2 and 4.1.1), which actually results in worse performance—speedup of 0.13 in this case—by using `tf.function` than not using it: “`tfa.image.equalize()` uses an internal `scale_channel()` function[,] which triggers excessive retracing ...” The problem is related to hybridizing *inner* functions: “I... tried using `@tf.function` at the `scale_channel()` and... `equalize_image()` level[s], but the further I moved it ‘inside,’ the slower `equalize()` became” [85]. The fix involved “using `@tf.function` at the top-level of `equalize()`, which made it run ~25%–40% faster ...” The root cause is that, “using [embedded] functions [i.e.,] defining functions inside function) will retrace the graph multiple times [as] the [ir] scope is not [publicly] visible, and the graphs cannot be cached” [84]. As a modularity mechanism, embedding (nesting) function definitions is a common idiom in Python, yet, currently, *TensorFlow* documentation does not mention this problem. Developers are left to consider the internals of *AutoGraph* in writing performant imperative DL code, leading to anti-pattern 1, Fig. 6.

Input Signatures. Arguments to `tf.function()`, particularly involving input tensor shapes, may also influence performance (q.v. Section 2). Listing 7 portrays an underspecified input signature (q.v. Section 4.1.1)—one of the most used `tf.function` parameters


```

1 def ndiagquad(funcs, H: int, Fmu, Fvar, logspace: bool = False, **Ys):
2     # Computes N Gaussian expectation integrals of one or more functions ...
3     def unify(f_list): # Stack a list of means/vars into a full block.
4         return tf.reshape(tf.concat([tf.reshape(f, (-1, 1)) for f in f_list],
5             axis=1), (-1, 1, Din))
6     if isinstance(Fmu, (tuple, list)):
7         Din = len(Fmu)
8     def unify(f_list): # Stack a list of means/vars into a full block.
9         return tf.reshape(tf.concat([tf.reshape(f, (-1, 1)) for f in f_list],
10            axis=1), (-1, 1, Din))
11     Fmu, Fvar = map(unify, [Fmu, Fvar]) # both [N, 1, Din]

```

Listing 8: Commit b65848a2 in GPFflow: fix compilation issue...

```

1 @tf.function
2 def interpolate_bilinear(grid, query_points, indexing="ij", name=None): # ...
3     tf.debugging.assert_equal(query_shape[2], 2, message=
4         "Query points must be size 2 in dim 2.")

```

Listing 9: Commit 8bab3226 in tensorflow/addons: remove tf.func

that we observed. On lines 2–6, a performance regression was fixed by adding an `input_signature` to a weight distribution `tf.function` to “make sure it does not recreate graph, which will slow down training significantly” [40]. The sequence of `tf.TensorSpecs` specifies the intended tensor shapes and data types (`dtypes`) that will be supplied to `update_weights()`. Otherwise, a separate (concrete) function (graph) is instantiated for each *inferred* input signature, which may result in retracing, leading to best practice 2, Fig. 5.

4.2.2 Compilation Errors. Consider `unify()` originally defined on lines 3–5, listing 8 that accesses `Din` on line 5. This variable, however, it is defined *after* the function definition on line 7—legal due to Python’s lexical scoping rules. In other words, the value of `Din` will come from the *calling* context. In this case, `Din` on line 5 is replaced with the value defined on line 7 due to `unify()` being accessed on line 11. The code, though, results in the following (runtime) `NameError` on line 5: free variable ‘Din’ referenced before assignment in enclosing scope [19]. The problem is that, while it itself is not a `tf.function`, `ndiagquad()` is *called* by a `tf.function` elsewhere—it will also be compiled into a (static) graph (cf. Section 4.1.1). Thus, dynamic language features like lexical scoping are not available in static contexts. As a result, `unify()` is moved to line 8, where `Din` is in its declaration scope. Although Python is a dynamic language, developers must be aware that certain code will be compiled to static graphs, leading to anti-pattern 3, Fig. 6.

4.2.3 API Misuse. On line 1 in listing 9, `@tf.function` is removed to fix a bug that is causing flaky tests [31]. The problem is deemed to be that `@tf.function` and the `assert` statement on lines 3–4 is incompatible. The developers express that “removing the decorator is not ideal, but stability is more important than the [speedup] we [would] get with [it]” [39]. However, this code likely causes a race condition because of a missing control dependency following the assertion. To use the assertion within a `tf.function`, a control dependency is required “to block follow-up computation[s] until the check has executed” as a result of the function being converted to a (static) graph [49]. This leads to best practice 3, Fig. 5.

4.2.4 TensorFlow Bugs. On line 1, listing 10, `@tf.function` is once again removed. The problem is that—with `@tf.function`—the application “can only process one image before” needing to restarted [20], terminating with the message: `ValueError: tf.function-decorated function tried to create variables on non-first call`. Recall from Section 4.1.1 that shared variables inside a `tf.function` must

```

1 @tf.function
2 def train_step(image):
3     with tf.GradientTape() as tape:
4         outputs = extractor(image)
5         loss = style_content_loss(outputs)
6         loss += total_variation_weight * tf.image.total_variation(image)
7         grad = tape.gradient(loss, image)
8         opt.apply_gradients([(grad, image)])
9         image.assign(clip_0_1(image))

```

Listing 10: Commit 8bab3226 in neuro-art: Multiple request bugfix...

```

1 @tf.function
2 def call(self, inputs):
3     """Call `Layer`"""
4     if not self.initialized:
5         self._data_dep_init(inputs)
6     if not self._initialized:
7         self._initialize_weights(inputs)
8     self._compute_weights() # Recompute weights for each forward pass ...

```

Listing 11: Commit 16ee6c59 in tensorflow/addons: tf.func for debug

- (1) More tool-support for assisting with using `tf.function` may help produce reliable yet performant imperative DL code.
- (2) Modernize and reformulate existing tensor shape mismatch detectors for imperative DL code and `tf.function(...)` input shapes.
- (3) More formal specification in a design-by-contract (DbC) style may be helpful for new tool-support aimed to alleviate API misuse.
- (4) Testing (dynamic analysis) focused on hybridized (imperative) DL code that runs under *multiple* execution modes may localize bugs.

Figure 7: Preliminary hybridization recommendations.

be singleton; a run-time exception ensues otherwise [42]. However, it is not obvious from listing 10 where the variable creation occurs—there are no explicit `tf.Variables`. The developer expresses that “removing the ... decorator is a viable workaround but not [a] best practice,” and that the root cause is an (unresolved) *TensorFlow* bug [88]. In terms of listing 10, the problematic line is 8, as “calling `apply_gradients()` on an optimizer for the first time will create its internal variables” [78]. In terms of the framework, it transpires to be related to software layering, as, “sadly[,] there [is] currently no public API to just initialize the optimizer state but not [apply it].” While several developers found workarounds for their particular situations, imperative DL code such as that in listing 10 have foregone any potential performance gains from using `@tf.function`, leading to best practice 4 and anti-pattern 2 in Figs. 5 and 6, respectively.

4.2.5 Debuggability. To improve debuggability, `@tf.function` is removed on line 1, listing 11 (cf. Section 4.1.1). However, in the latest file version, `@tf.function` has not been replaced. Thus, the developer may have inadvertently sacrificed permanent performance gains for temporary debuggability, leading to best practice 5, Fig. 5.

5 DISCUSSION

We summarize and comment on our main findings while connecting them to other research. To help solve the problems, we put forth preliminary recommendations for practitioners, tool developers, and researchers. Though our hope is that the findings will shed light on future tool challenges and that the aforementioned descriptions and real-world examples will provide sufficient, generalizable, and actionable contexts, we nonetheless outline potential solutions.

Performance. It is not surprising that performance is our largest category (finding 1) since hybridization is centrally related to performance enhancement. The volume of performance problems is a

testament to the struggles developers have in writing performant, imperative DL code. However, 45.05% of performance problems (finding 3) were due to *existing* `tf.function` usages, suggesting that developers also struggle with using hybridization *effectively* to achieve the performance they desire. A feasible explanation is that developers must *manually* decide: (i) where and when to use `tf.function`, (ii) the arguments to supply `tf.function` for their code to perform optimally, and (iii) which code is amenable to (efficient) graph conversion and which is not, all of which can be error-prone.

The overarching goal is reliable *and* performant DL code; reliability stems from being able to write DL code in a less error-prone imperative-style, while performance is achieved in migrating that code to graph execution at run-time. *AutoGraph*, as well as other hybridization technologies, attempt to achieve this goal by automating the migration process as much as possible—frequently requiring contextual information from developers as to their intentions and imposing limitations of where the technology can be used. The end result is a trade-off—one of many typically made by DL frameworks [57]. As discussed in Section 1, others [60] attempt to automate the entire migration process—not requiring any contextual metadata—but impose new trade-offs, such as necessitating custom Python interpreters that may not be practical for industrial applications and support only specific Python constructs.

As *AutoGraph* and other hybridization technologies are pervasively used, as well as being integrated into official distributions of popular DL frameworks, our suggestion is to retain (and continually improve upon) hybridization platforms, while simultaneously posing this problem as one of (API) *usability*. Our perspective is that what is needed is tool-support that will guide developers in using this technology correctly given a particular context, as well as automated refactoring and other source code transformation tools that can detect and repair hybridization problems. Such techniques would alleviate hybridization issues well-before they are seen beyond (production) deployment or after long training sessions, leading to recommendation 1, Fig. 7.

Cao et al. [22] and Zhang et al. [100] also study performance of DL code and found that a modest portion of *non-imperative TensorFlow* program bugs involved performance problems. However, these problems were caused by confusion with the underlying computation model, which essentially requires developers to build graphs *manually*. In our case, graphs are built *automatically*. Since imperative DL code runs eagerly by default, it is understandable that our study would uncover more performance problems. In fact, Tambon et al. [90] also observe performance degradation of imperative DL code. Wan et al. [97] found that performance bugs took the longest average time to fix in blockchain systems.

Per finding 4, developer-supplied arguments to `tf.function()` played a major role in performance problems, comprising 25.23% of performance fixes. Furthermore, per finding 5, a significant percentage (18.92%) of performance problems involved parameters representing input shape specification—one of the most frequently used `tf.function` parameters (q.v. Section 4.2.1). Input shape problems are a central focus of related work [56,66,100] on DL programs; related studies [55,56,58,100] also found shape problems. A feasible explanation is that developers are challenged to determine tensor shapes from all possible call sites statically. We again advocate for more tool-support in this area, e.g., an adaptation of Lagouvardos

et al. [66] for imperative DL programs focused on hybridization parameters, leading to recommendation 2 in Fig. 7.

API Misuse. Per finding 6, using the `tf.function` API inconsistently with its documentation was a major theme. Feasible explanations include: (i) DL APIs—along with their documentation [53]—are particularly vast and complex [57], (ii) often, documentation consumers (developers) are not software experts [53], (iii) although developers are writing imperative DL code, there exist situations where they must nevertheless be cognizant of hybridization limitations, and (iv) error messages may not be helpful. Due to Item (i), learning how to use DL APIs effectively necessitates a steep learning curve, especially considering that hybridization is relatively new. As ML systems have a quick time-to-market [86], developers may not have the luxury of time to thoroughly understand the documentation. This is especially evident in finding 7, with 37.74% of misuses caused by API confusion. Item (ii) has been recognized by other ML/DL software studies (e.g., [91]). We conjecture that Item (iii) can also be alleviated with more tool-support, however, such tool-support in this context may require (e.g., design-by-contract) formalization of DL API specifications (e.g., modeling operation limitations in particular contexts), leading to recommendation 3, Fig. 7. A potential downside to recommendation 3 is the rapid change of ML APIs [32]. For Item (iv), developers often expressed frustration with error messages, e.g., “the main complexity in [*TensorFlow*] 2 is in `@tf.function[:]` ... error messages should be as clear as possible, especially for common problems” [28].

Zhang et al. [99] likewise observed broader API misuse in DL systems. Nadi et al. [73] also found API misuse despite ample documentation in the context of cryptography—developers prefer higher-level documentation. Current hybridization documentation tends to focus on lower-level details—future research may explore whether a similar concept will work for DL APIs. Furthermore, our findings coincide with Jin et al. [61] that many performance bugs are due performance implication misunderstandings of certain functions.

Incompatibility. Execution incompatibility of particular Python constructs was also a major theme (q.v. finding 9). Zhang et al. [99] found a similar problem in DL systems w.r.t. CPU/GPU compatibility. We again advocate for more automation to circumvent such problems. To use hybridization effectively, developers must understand which constructs are amenable to *both* eager and graph execution and make appropriate considerations. Tool-support, e.g., IDE recommendations, may be helpful here. To alleviate run-time errors and unexpected results, we also advocate for more testing (dynamic analysis) of (imperative) DL code that runs the same code under *multiple* execution modes. Testing of DL systems is an emerging yet promising area, and testing focusing on (imperative) DL code hybridization may help to shed light on: (i) where developers struggle to write performant yet reliable (imperative) DL code and (ii) potential areas of where hybridization technologies can be improved. This leads to recommendation 4, Fig. 7.

Commits vs. GitHub Issues. Performance bugs appeared more in commits than GitHub issues. The reason may be that enhancing performance typically requires a code change, which can be benchmarked. Contrarily, “incompatibility” is more difficult to quantify, often resulting in unexpected behavior or run-time errors (q.v. Fig. 3). Therefore, developers may be more likely to seek external assistance. Developers commonly file GitHub issues against

TensorFlow; 93.75% of TFB issues are against the TensorFlow subject. That all UKN and TST bugs appeared in commits may be due to GitHub issues being easier to categorize than changesets and DL testing remains an emerging area, respectively.

6 THREATS TO VALIDITY

Subjects may not be representative of DL systems. To mitigate this, subjects encompass diverse domains and sizes, have been used in previous studies, and are from a data science-specific dataset (q.v. Section 3). Various GitHub metrics and DL-related keywords were used in choosing subjects. Also, hybridization is relatively new; we expect a larger selection of subjects as it grows in popularity.

Our study involved many hours of manual validation to understand and categorize bugs. To mitigate bias, we investigated referenced resources and comments made by developers to help more fully understand the challenges faced. The NLP of `gitcproc` may have missed bug fix changesets. Nevertheless, using it, we were still able to find 157 bugs (280 overall) that contributed to a rich bug categorization, best practices, and anti-patterns. Furthermore, `gitcproc` has been used previously in other studies (q.v. Section 3).

Hybridization in comparable DL frameworks may have yielded different challenges. Nevertheless, focusing on *TensorFlow* enables us to more thoroughly understand the intricacies involved in using hybridization effectively. Moreover, *TensorFlow* is a widely-studied and popular (industrial) DL framework (q.v. Section 1).

7 RELATED WORK

Cao et al. [22] characterizing performance bugs in DL systems. During their analysis of general performance bugs, they also find that developers often struggle with knowing where to add `@tf.function` and how to implement decorated functions for optimal performance. Beyond performance bugs, our study includes a rich, hierarchical taxonomy of varying hybridization bug types, including input shape mismatches, API misuse, and construct incompatibility, whose results include run-time errors, unexpected behavior, and deadlock. Tambon et al. [90] examine (silent) behavioral bugs *within* DL frameworks and their impact on client code. Their work is reminiscent of our TFB problem category (q.v. Section 4.1.1) and also note that performance degradation may lead to significant problems at run-time. While they do not explicitly mention hybridization performance bugs, some of their performance bugs in imperative DL code may be alleviated by using `@tf.function`. Zhang et al. [101] study API change trends in *TensorFlow* and for which reasons; our focus is on *client* code modifications involving hybridization. Baker et al. [15] extract 11 common *TensorFlow* API misuse patterns. Only one of the patterns (and corresponding fix suggestion) involves (a specific use case of) `tf.function`. In contrast, our study goes beyond API misuse and entails 12 top-level problem categories—24 overall—encompassing hybridization challenges.

Zhang et al. [99] present a large-scale empirical study of general DL questions on Stack Overflow. Particularly, their “CPU/GPU incompatibility” problem category resembles our execution mode incompatibility category. Concerning hybridization, whether the migrated graph executes on a GPU is typically decided by the underlying DL framework; our focus is on conversion itself. Islam

et al. [56] and Zhang et al. [100] study general DL bug characteristics and present anti-patterns to avoid bugs. Islam et al. [58] study patterns in which such bugs are fixed. Chen et al. [26] explore faults in deploying DL models to mobile applications. Nikanjam and Khomh [74] catalog various design smells in DL systems and recommend suitable refactorings. Jebnoun et al. [59] correlate code smells with bugs in DL code. Liu et al. [68] characterize technical debt in DL frameworks, while Humbatova et al. [55] taxonomize (functional) faults in DL systems. Arpteg et al. [9] categorize (general) SE challenges in DL systems into three areas—development, production, and organizational. Liu et al. [67] study failed *TensorFlow* industrial jobs and propose a constraint-based approach for detecting shape-related errors. Amershi et al. [4] conduct a study at Microsoft, observing software teams as they developed AI applications. Lwakatere et al. [69] also classify SE challenges for ML systems at six different companies, focusing mainly on deployment issues. Thung et al. [93] examine bugs in three general ML systems, finding that nonfunctional bugs, of which performance problems may be categorized, require the most involved fixes. Dilhara et al. [32] study ML library evolution and its resulting client-code modifications. And, Dilhara et al. [33] and Tang et al. [91] analyze repetitive code changes and refactorings made in ML systems, respectively. While valuable, these studies do not deal with challenges faced in migrating imperative DL code to graph execution.

Several studies involve performance in other contexts. Han and Yu [52] study configurability and performance. Future work entails correlating their findings with `tf.function` arguments. Jin et al. [62] study performance slowdowns caused by system side inefficiencies. Bagherzadeh et al. [13] investigate performance in Actor-based systems. Others study language features. Parnin et al. [77] study Java generics adoption. Dyer et al. [34] study language feature evolution. Khatchadourian and Masuhara [64] empirically assess default methods. There are many general empirical studies. Makhshari and Mesbah [70] taxonomize development challenges of IoT systems. Bagherzadeh and Khatchadourian [14] investigate common questions asked by big data developers, and Khatchadourian et al. [65] examine the use and misuse of Java streams. Engler et al. [35] and Tian and Ray [94] study errors in systems code.

8 CONCLUSION & FUTURE WORK

This study advances knowledge of the development challenges involved in migrating imperative DL code to graph execution via hybridization. A hierarchical taxonomy of common hybridization challenges was formulated and preliminary recommendations, best practices, and anti-patterns were proposed. In the future, we will explore analyzing alternative developer resources, e.g., Stack Overflow, and integrating our results into automated bug finders and refactoring detection approaches [11,95].

ACKNOWLEDGMENTS

We thank Manal Zneit, Ye Paing, and Jack Cruse-Mulhall, developers providing feedback, and the anonymous reviewers for thorough and insightful comments. Support for this project was provided by PSC-CUNY Award #638010051, jointly funded by The Professional Staff Congress and The City University of New York.

REFERENCES

- [1] Martin Abadi et al. 2016. TensorFlow: a system for large-scale Machine Learning. In *Symposium on Operating Systems Design and Implementation*.
- [2] 2020. Added jitted ncon. Pull request #623. google/TensorNetwork. Xanadu. (May 26, 2020). Retrieved 01/10/2022 from <https://git.io/J9cMx>.
- [3] Akshay Agrawal et al. 2019. TensorFlow Eager: a multi-stage, Python-embedded DSL for Machine Learning. (2019). arXiv: 1903.01855 [cs. PL].
- [4] Saleema Amershi, Andrew Begel, Christian Bird, Robert DeLine, Harald Gall, Ece Kamar, Nachiappan Nagappan, Besmira Nushi, and Thomas Zimmermann. 2019. Software Engineering for Machine Learning: a case study. In *International Conference on Software Engineering*. *Software Engineering in Practice*. IEEE. IEEE, (May 2019), 291–300. doi: 10.1109/ICSE-SEIP.2019.00042.
- [5] Apache. 2018. Customer layers (beginners). Apache MXNet documentation. Retrieved 07/23/2021 from https://mxnet.apache.org/versions/1.7/api/python/docs/tutorials/packages/gluon/blocks/custom_layer_beginners.html.
- [6] Apache. 2021. Hybridize. Apache MXNet documentation. (April 8, 2021). Retrieved 04/08/2021 from <https://mxnet.apache.org/versions/1.8.0/api/python/docs/tutorials/packages/gluon/blocks/hybridize.html>.
- [7] Apache Software Foundation. 2021. Open Deep Learning compiler stack. (December 1, 2021). Retrieved 12/01/2021 from <https://git.io/JMr6c>.
- [8] Apple Inc. 2021. Core ML tools. (December 1, 2021). Retrieved 12/01/2021 from <https://git.io/JMr61>.
- [9] A. Arpteg, B. Brinne, L. Crnkovic-Friis, and J. Bosch. 2018. Software Engineering challenges of Deep Learning. In *Euromicro Conference on Software Engineering and Advanced Applications*. IEEE, 50–59. doi: 10.1109/SEAA.2018.00018.
- [10] Deeksha Arya, Wenting Wang, Jin L. C. Guo, and Jinghui Cheng. 2019. Analysis and detection of information types of open source software issue discussions. In *International Conference on Software Engineering*. (May 2019), 454–464. doi: 10.1109/ICSE.2019.00058.
- [11] Hassan Atwi, Bin Lin, Nikolaos Tsantalis, Yutaro Kashiwa, Yasutaka Kamei, Naoyasu Ubayashi, Gabriele Bavota, and Michele Lanza. 2021. PyRef: refactoring detection in Python projects. In *International Working Conference on Source Code Analysis and Manipulation*. doi: 10.1109/SCAM52516.2021.00025.
- [12] Hlib Babii, Julian Aron Prenner, Laurin Stricker, Anjan Karmakar, Andrea Janes, and Romain Robbes. 2021. Mining software repositories with a collaborative heuristic repository. In *International Conference on Software Engineering (ICSE-NIER '21)*. IEEE, 106–110. doi: 10.1109/ICSE-NIER52604.2021.00030.
- [13] Mehdi Bagherzadeh, Nicholas Fireman, Anas Shawesh, and Raffi Khatchadourian. 2020. Actor concurrency bugs: a comprehensive study on symptoms, root causes, API usages, and differences. *Proc. ACM Program. Lang.*, 4, OOPSLA, Article 214, (November 2020), 1–32. doi: 10.1145/3428282.
- [14] Mehdi Bagherzadeh and Raffi Khatchadourian. 2019. Going big: a large-scale study on what big data developers ask. In *Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. ACM, 432–442. doi: 10.1145/3338906.3338939.
- [15] Wilson Baker, Michael O'Connor, Seyed Reza Shahamiri, and Valerio Terragni. 2022. Detect, fix, and verify TensorFlow API misuses. In *International Conference on Software Analysis, Evolution and Reengineering*, 1–5.
- [16] Houssein Ben Braiek, Foutse Khomh, and Bram Adams. 2018. The open-closed principle of modern Machine Learning frameworks. In *Mining Software Repositories*, 353–363. ISBN: 978-1-4503-5716-6. doi: 10.1145/3196398.3196445.
- [17] S. Biswas, M. J. Islam, Y. Huang, and H. Rajan. 2019. Boa meets Python: a Boa dataset of Data Science software in Python language. In *Mining Software Repositories*, 577–581. doi: 10.1109/MSR.2019.00086.
- [18] 2021. bug boxsize=nc. modichirag/galference. af1664e. UC Berkeley. (April 16, 2021). Retrieved 01/10/2022 from <https://git.io/J9ciM>.
- [19] 2020. Bug fix: compilation issue quadrature. Pull request #1418. GPflow/GPflow. Cambridge University. (April 8, 2020). Retrieved 01/14/2022 from <https://github.com/GPflow/GPflow/pull/1418#issue-596552141>.
- [20] 2020. Bug in TensorFlow only allows it to run once. Issue #13. MLH-Fellowship/neuro-art. MLH Fellowship. (November 10, 2020). Retrieved 01/13/2022 from <https://github.com/MLH-Fellowship/neuro-art/issues/13>.
- [21] 2019. Calling tf.function from tf.py_function in dataset.map hangs. (September 11, 2019). Retrieved 01/05/2022 from <https://git.io/JSSBw>.
- [22] Junming Cao, Bihuan Chen, Chao Sun, Longjie Hu, and Xin Peng. 2021. Characterizing performance bugs in Deep Learning systems. (December 3, 2021). arXiv: 2112.01771 [cs. SE].
- [23] Casey Casalnuovo, Yagnik Suchak, Baishakhi Ray, and Cindy Rubio-González. 2017. GitProc: a tool for processing and classifying GitHub commits. In *International Symposium on Software Testing and Analysis (ISSTA '17)*. ACM, 396–399. doi: 10.1145/3092703.3098230.
- [24] Tatiana Castro Vélez, Raffi Khatchadourian, Mehdi Bagherzadeh, and Anita Raja. Challenges in migrating imperative DL programs to graph execution: an empirical study. Zenodo, (March 31, 2022). doi: 10.5281/zenodo.5601987.
- [25] Tianqi Chen, Mu Li, Yutian Li, Min Lin, Naiyan Wang, Minjie Wang, Tianjun Xiao, Bing Xu, Chiyuan Zhang, and Zheng Zhang. 2015. MXNet: a flexible and efficient Machine Learning library for heterogeneous distributed systems. In *Workshop on Machine Learning Systems at NIPS*. arXiv: 1512.01274 [cs. DC].
- [26] Zhenpeng Chen, Huihan Yao, Yiling Lou, Yanbin Cao, Yuanqiang Liu, Haoyu Wang, and Xuanzhe Liu. 2021. An empirical study on deployment faults of Deep Learning based mobile applications. In *International Conference on Software Engineering*. ACM/IEEE. IEEE. doi: 10.1109/icse43902.2021.00068.
- [27] François Chollet. 2020. *Deep Learning with Python*. (2nd edition). Manning.
- [28] 2019. Cryptic error message when assigning to a variable in a tf.function. Issue #30768. tensorflow/tensorflow. (July 21, 2019). Retrieved 01/17/2022 from <https://github.com/tensorflow/tensorflow/issues/30768>.
- [29] 2021. Deadlock on recursive tf.function-decorated function. Issue #35540. (October 8, 2021). Retrieved 01/05/2022 from <https://git.io/JSS4P>.
- [30] 2021. Deep Learning examples. NVIDIA. (March 1, 2021). Retrieved 05/05/2021 from <https://git.io/J2vFG>.
- [31] 2019. Dense image warp tests are flaky. AWS. (April 3, 2019). Retrieved 01/13/2022 from <https://github.com/tensorflow/addons/issues/138#issue-428951400>.
- [32] Malinda Dilhara, Ameya Ketkar, and Danny Dig. 2021. Understanding software-2.0: a study of Machine Learning library usage and evolution. *ACM Transactions on Software Engineering and Methodology*. doi: 10.1145/3453478.
- [33] Malinda Dilhara, Ameya Ketkar, Nikhith Sannidhi, and Danny Dig. 2022. Discovering repetitive code changes in Python ML systems. In *International Conference on Software Engineering (ICSE '22)*. To appear.
- [34] Robert Dyer, Hridesh Rajan, Hoan Anh Nguyen, and Tien N. Nguyen. 2014. Mining billions of AST nodes to study actual and potential usage of Java language features. In *International Conference on Software Engineering*, 779–790. ISBN: 978-1-4503-2756-5.
- [35] Dawson Engler, David Yu Chen, Seth Hallem, Andy Chou, and Benjamin Chelf. 2001. Bugs as deviant behavior: a general approach to inferring errors in systems code. In *Symposium on Operating Systems Principles (SOSP '01)*. ACM, Banff, Alberta, Canada, 57–72. doi: 10.1145/502034.502041.
- [36] 2020. Ensure compatibility with tf.function. secondmind-labs/trieste. Issue #90. Secondmind Labs. (December 2, 2020). Retrieved 11/08/2021 from <https://github.com/secondmind-labs/trieste/issues/90>.
- [37] Facebook Inc. 2019. PyTorch documentation. TorchScript. en. Retrieved 02/19/2021 from <https://pytorch.org/docs/stable/jit.html>.
- [38] 2021. Fit bug in Blatt-Weisskopf; update Dalitz decomposition. Apoluekt/AmpliTF.d02db12. CERN. (July 19, 2021). Retrieved 01/03/2022 from <https://git.io/JSi9f>.
- [39] 2019. FIX: dense image warp bug. (April 17, 2019). Retrieved 01/13/2022 from <https://github.com/tensorflow/addons/pull/187>.
- [40] 2021. Fixed all ... this should work. samuelmat19/DDPG-tf2. 02a3f29. ML6. (February 26, 2021). Retrieved 01/12/2022 from <https://github.com/samuelmat19/DDPG-tf2/commit/02a3f297#47584455>.
- [41] GitHub, Inc. 2021. Search. REST API Reference. GitHub Docs. Retrieved 12/02/2021 from <https://docs.github.com/en/rest/reference/search>.
- [42] Google LLC. 2021. Better performance with tf.function. (February 4, 2021). Retrieved 02/19/2021 from <https://tensorflow.org/guide/function>.
- [43] Google LLC. 2022. Introduction to graphs and tf.function. (January 19, 2022). Retrieved 01/20/2022 from https://tensorflow.org/guide/intro_to_graphs.
- [44] Google LLC. 2021. Introduction to variables. TensorFlow core. (November 11, 2021). Retrieved 01/03/2022 from <https://www.tensorflow.org/guide/variable>.
- [45] Google LLC. 2021. IREE: intermediate representation execution environment. (December 1, 2021). Retrieved 12/01/2021 from <https://git.io/JMrPT>.
- [46] Google LLC. 2021. Migrate your TensorFlow 1 code to TensorFlow 2. Automatic conversion script. TensorFlow Core. (May 27, 2021). Retrieved 05/27/2021 from https://tensorflow.org/guide/migrate#automatic_conversion_script.
- [47] Google LLC. 2021. Random number generation. Interaction with tf.function. TensorFlow Core. (November 16, 2021). Retrieved 01/07/2022 from https://tensorflow.org/guide/random_numbers#interaction_with_tffunction.
- [48] Google LLC. 2021. tf.compat.v1.Session. (May 14, 2021). Retrieved 07/06/2021 from https://tensorflow.org/api_docs/python/tf/compat/v1/Session#run.
- [49] Google LLC. 2021. tf.debugging.assert_equal. Retrieved 01/13/2022 from https://tensorflow.org/api_docs/python/tf/debugging/assert_equal#returns.
- [50] Google LLC. 2021. tf.py_function. Version 2.7.0. (November 5, 2021). Retrieved 01/05/2022 from https://tensorflow.org/api_docs/python/tf/py_function.
- [51] Google LLC. 2021. tf.random.set_seed. (December 4, 2021). Retrieved 01/03/2022 from https://tensorflow.org/api_docs/python/tf/random/set_seed.
- [52] Xue Han and Tingting Yu. 2016. An empirical study on performance bugs for highly configurable software systems. In *International Symposium on Empirical Software Engineering and Measurement*. doi: 10.1145/2961111.2962602.
- [53] Y. Hashemi, M. Nayebi, and G. Antoniol. 2020. Documentation of Machine Learning software. In *International Conference on Software Analysis, Evolution and Reengineering*. (February 2020). doi: 10.1109/SANER48275.2020.9054844.
- [54] Horace He. 2019. The state of Machine Learning frameworks in 2019. Retrieved 04/01/2021 from <https://thegradient.pub/state-of-ml-frameworks-2019-pytorch-dominates-research-tensorflow-dominates-industry>.
- [55] Nargiz Humbatova, Gunel Jahangirova, Gabriele Bavota, Vincenzo Riccio, Andrea Stocco, and Paolo Tonella. 2020. Taxonomy of real faults in Deep

- Learning systems. In *International Conference on Software Engineering*. doi: 10.1145/3377811.3380395.
- [56] Md Johirul Islam, Giang Nguyen, Rangeet Pan, and Hridesh Rajan. 2019. A comprehensive study on Deep Learning bug characteristics. In *Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. (August 2019). doi: 10.1145/3338906.3338955.
- [57] Md Johirul Islam, Hoan Anh Nguyen, Rangeet Pan, and Hridesh Rajan. 2019. What do developers ask about ML libraries? a large-scale study using Stack Overflow. (2019). arXiv: 1906.11940 [cs. SE].
- [58] Md Johirul Islam, Rangeet Pan, Giang Nguyen, and Hridesh Rajan. 2020. Repairing Deep Neural Networks: fix patterns and challenges. In *International Conference on Software Engineering*. doi: 10.1145/3377811.3380378.
- [59] Hadhemi Jebnoun, Houssein Ben Braiek, Mohammad Masudur Rahman, and Foutse Khomh. 2020. The scent of Deep Learning code: an empirical study. In *Mining Software Repositories*. doi: 10.1145/3379597.3387479.
- [60] Eunji Jeong, Sungwoo Cho, Gyeong-In Yu, Joo Seong Jeong, Dong-Jin Shin, Taebum Kim, and Byung-Gon Chun. 2019. Speculative symbolic graph execution of imperative Deep Learning programs. *SIGOPS Oper. Syst. Rev.*, 53, 1, (July 2019), 26–33. issn: 0163-5980. doi: 10.1145/3352020.3352025.
- [61] Guoliang Jin, Linhai Song, Xiaoming Shi, Joel Scherpelz, and Shan Lu. 2012. Understanding and detecting real-world performance bugs. In *Programming Language Design and Implementation*. doi: 10.1145/2254064.2254075.
- [62] Hui Jin, Kan Qiao, Xian-He Sun, and Ying Li. 2011. Performance under failures of MapReduce applications. In *International Symposium on Cluster, Cloud and Grid Computing*. doi: 10.1109/ccgrid.2011.84.
- [63] Raffi Khatchadourian. 2021. graph_execution_time_comparison.ipynb. (February 23, 2021). Retrieved 11/03/2021 from <https://bit.ly/3bwrhVt>.
- [64] Raffi Khatchadourian and Hidehiko Masuhara. 2018. Proactive empirical assessment of new language feature adoption via automated refactoring: the case of Java 8 default methods. *The Art, Science, and Engineering of Programming*, 2, 6, 6:1–6:30, 3. doi: 10.22152/programming-journal.org/2018/2/6.
- [65] Raffi Khatchadourian, Yiming Tang, Mehdi Bagherzadeh, and Baishakhi Ray. 2020. An empirical study on the use and misuse of Java 8 streams. In *International Conference on Fundamental Aspects of Software Engineering*. ETAPS. (April 2020), 97–118. doi: 10.1007/978-3-030-45234-6_5.
- [66] Sifis Lagouvardos, Julian Dolby, Neville Grech, Anastasios Antoniadis, and Yannis Smaragdakis. 2020. Static analysis of shape in TensorFlow programs. In *European Conference on Object-Oriented Programming*. Volume 166, 15:1–15:29. doi: 10.4230/LIPIcs.ECOOP.2020.15.
- [67] Chen Liu, Jie Lu, Guangwei Li, Ting Yuan, Lian Li, Feng Tan, Jun Yang, Liang You, and Jingling Xue. 2021. Detecting TensorFlow program bugs in real-world industrial environment. In *International Conference on Automated Software Engineering*. IEEE. doi: 10.1109/ase51524.2021.9678891.
- [68] Jiakun Liu, Qiao Huang, Xin Xia, Emad Shihab, David Lo, and Shanping Li. 2020. Is using Deep Learning frameworks free? Characterizing technical debt in Deep Learning frameworks. In *International Conference on Software Engineering* (ICSE-SEIS '20). doi: 10.1145/3377815.3381377.
- [69] Lucy Ellen Lwakatare, Aiswarya Raj, Jan Bosch, Helena Holmström Olsson, and Ivica Crnkovic. 2019. A taxonomy of software engineering challenges for Machine Learning systems: an empirical investigation. In *Agile Processes in Software Engineering and Extreme Programming*, 227–243. doi: 10.1007/978-3-030-19034-7_14.
- [70] Amir Makhshari and Ali Mesbah. 2021. IoT bugs and development challenges. In *International Conference on Software Engineering*. ACM/IEEE. IEEE, (May 2021). doi: 10.1109/icse43902.2021.00051.
- [71] 2020. Migrate to tf.module and add support for SavedModels. Pull request #603. onnx/onnx-tensorflow. IBM. (July 2, 2020). Retrieved 12/14/2021 from <https://git.io/JDEoD>.
- [72] Dan Moldovan, James M. Decker, Fei Wang, Andrew A. Johnson, Brian K. Lee, Zachary Nado, D. Sculley, Tiark Rompf, and Alexander B. Wiltschko. 2019. AutoGraph: imperative-style coding with graph-based performance. (2019). arXiv: 1810.08061 [cs. PL].
- [73] Sarah Nadi, Stefan Krüger, Mira Mezini, and Eric Bodden. 2016. “Jumping through hoops”: why do Java developers struggle with cryptography APIs? In *International Conference on Software Engineering*. ACM, (May 2016), 935–946. doi: 10.1145/2884781.2884790.
- [74] Amin Nikanjam and Foutse Khomh. 2021. Design smells in Deep Learning programs: an empirical study. In *International Conference on Software Maintenance and Evolution*. IEEE, 332–342. doi: 10.1109/ICSM52107.2021.00036.
- [75] NVIDIA Corporation. 2021. TensorRT open source software. (December 1, 2021). Retrieved 12/01/2021 from <https://git.io/fjVoO>.
- [76] Shengyi Pan, Lingfeng Bao, Xiaoxue Ren, Xin Xia, David Lo, and Shanping Li. 2021. Automating developer chat mining. In *International Conference on Automated Software Engineering*. (November 2021), 854–866. doi: 10.1109/ASE51524.2021.9678923.
- [77] Chris Parnin, Christian Bird, and Emerson Murphy-Hill. 2013. Adoption and use of Java generics. *Empirical Softw. Engg.*, 18, 6, (December 2013), 1047–1089. issn: 1382-3256. doi: 10.1007/s10664-012-9236-6.
- [78] Alexandre Passos. 2019. tf.function-decorated function tried to create variables on non-first call. (June 11, 2019). Retrieved 01/13/2022 from <https://github.com/tensorflow/tensorflow/issues/27120#issuecomment-500975337>.
- [79] Adam Paszke et al. 2019. PyTorch: an imperative style, high-performance Deep Learning library. (December 3, 2019). arXiv: 1912.01703 [cs. LG].
- [80] 2021. Performance bottleneck due to tf.function retracing. Issue #74. TensorFlow. q-optimize/c3. (March 18, 2021). Retrieved 11/08/2021 from <https://github.com/q-optimize/c3/issues/74>.
- [81] 2020. Reduce tf.function retracing. Pull Request #100. (August 10, 2020). Retrieved 11/08/2021 from <https://github.com/keiohta/tf2rl/pull/100>.
- [82] 2020. Remove @tf.function in tf.image.equalize. Issue #2263. eomii. (December 3, 2020). Retrieved 01/10/2022 from <https://git.io/J9cAb>.
- [83] 2022. Remove `experimental_relax_shapes=true` to avoid none shape. University of Chinese Academy of Sciences. (January 4, 2022). Retrieved 01/10/2022 from <http://github.com/jiangyi15/tf-pwa/commit/0db1#62889998>.
- [84] 2020. Remove tf.function decorator in tf.image.equalize. (December 5, 2020). Retrieved 01/10/2022 from <https://git.io/J9cFt>.
- [85] 2020. Remove tf.function decorator in tf.image.equalize. eomii. (December 8, 2020). Retrieved 01/10/2022 from <https://git.io/J9cHq>.
- [86] D. Sculley, Gary Holt, Daniel Golovin, Eugene Davydog, Todd Phillips, Diemar Ebner, Vinay Chaudhary, Michael Young, Jean-François Crespo, and Dan Dennison. 2015. Hidden technical debt in Machine Learning systems. In *Advances in Neural Information Processing Systems*, 2503–2511.
- [87] Stack Exchange Inc. 2020. Should I use @tf.function for all functions? Stack Overflow. (January 21, 2020). Retrieved 11/08/2021 from <https://stackoverflow.com/questions/59847045/should-i-use-tf-function-for-all-functions>.
- [88] Eric Stavarache. 2019. tf.function-decorated function tried to create variables on non-first call. Issue #27120. ETH Zürich. (March 25, 2019). Retrieved 01/13/2022 from <https://github.com/tensorflow/tensorflow/issues/27120>.
- [89] 2019. Surprising random seed behavior when using @tf.function. Issue #33297. (October 13, 2019). Retrieved 01/03/2022 from <https://git.io/J5iac>.
- [90] Florian Tambon, Amin Nikanjam, Le An, Foutse Khomh, and Giuliano Antoniol. 2021. Silent bugs in Deep Learning frameworks: an empirical study of Keras and TensorFlow. (2021). arXiv: 2112.13314 [cs. SE].
- [91] Yiming Tang, Raffi Khatchadourian, Mehdi Bagherzadeh, Rhia Singh, Ajani Stewart, and Anita Raja. 2021. An empirical study of refactorings and technical debt in Machine Learning systems. In *International Conference on Software Engineering*, 238–250. doi: 10.1109/ICSE43902.2021.00033.
- [92] Yiming Tang, Allan Spektor, Raffi Khatchadourian, and Mehdi Bagherzadeh. 2022. Automated evolution of feature logging statement levels using Git histories and degree of interest. *Science of Computer Programming*, 214, (February 1, 2022), 102724. doi: 10.1016/j.scico.2021.102724.
- [93] Ferdian Thung, Shaowei Wang, David Lo, and Lingxiao Jiang. 2012. An empirical study of bugs in Machine Learning systems. In *International Symposium on Software Reliability Engineering*. doi: 10.1109/ISSRE.2012.22.
- [94] Yuchi Tian and Baishakhi Ray. 2017. Automatically diagnosing and repairing error handling bugs in C. In *Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. ACM, 752–762. doi: 10.1145/3106237.3106300.
- [95] Nikolaos Tsantalis, Ameya Ketkar, and Danny Dig. 2020. RefactoringMiner 2.0. *IEEE Trans. Softw. Eng.* doi: 10.1109/TSE.2020.3007722.
- [96] Anthony J. Viera and Joanne M. Garrett. 2005. Understanding interobserver agreement: the kappa statistic. *Family medicine*, 37, 360–363, 5.
- [97] Zhiyuan Wan, David Lo, Xin Xia, and Liang Cai. 2017. Bug characteristics in blockchain systems: a large-scale empirical study. In *Mining Software Repositories*. IEEE, (May 2017), 413–424. doi: 10.1109/MSR.2017.59.
- [98] Wenting Wang, Deeksha Arya, Nicole Novielli, Jinghui Cheng, and Jin L.C. Guo. 2020. ArguLens: anatomy of community opinions on usability issues using argumentation models. In *Conference on Human Factors in Computing Systems* (CHI '20). ACM, (April 2020). doi: 10.1145/3313831.3376218.
- [99] Tianyi Zhang, Cuiyun Gao, Lei Ma, Michael Lyu, and Miryung Kim. 2019. An empirical study of common challenges in developing Deep Learning applications. In *International Symposium on Software Reliability Engineering*. (October 2019). doi: 10.1109/ISSRE.2019.00020.
- [100] Yuhao Zhang, Yifan Chen, Shing-Chi Cheung, Yingfei Xiong, and Lu Zhang. 2018. An empirical study on TensorFlow program bugs. In *International Symposium on Software Testing and Analysis*. doi: 10.1145/3213846.3213866.
- [101] Zejun Zhang, Yanming Yang, Xin Xia, David Lo, Xiaoxue Ren, and John Grundy. 2021. Unveiling the mystery of API evolution in Deep Learning frameworks: a case study of TensorFlow 2. In *International Conference on Software Engineering* (ICSE-SEIP). doi: 10.1109/ICSE-SEIP52600.2021.00033.
- [102] Jiayuan Zhou, Shaowei Wang, Cor-Paul Bezemer, Ying Zou, and Ahmed E. Hassan. 2021. Studying the association between bounty issues and the issue-addressing likelihood of GitHub issue reports. *IEEE Trans. Softw. Eng.*, 47, 12, (December 2021), 2919–2933. doi: 10.1109/TSE.2020.2974469.
- [103] Weijie Zhou, Yue Zhao, Guoqiang Zhang, and Xipeng Shen. 2020. HARP: holistic analysis for refactoring Python-based analytics programs. In *International Conference on Software Engineering*, 506–517. doi: 10.1145/3377811.3380434.