

Demystifying Dependency Bugs in Deep Learning Stack

Kaifeng Huang
Fudan University
China

Bihuan Chen*
Fudan University
China

Susheng Wu
Fudan University
China

Junmin Cao
Fudan University
China

Lei Ma
University of Alberta
Canada

Xin Peng
Fudan University
China

ABSTRACT

Recent breakthroughs in deep learning (DL) techniques have stimulated significant growth in developing DL-enabled applications. These DL applications, built upon a heterogeneous and complex DL stack (e.g., Nvidia GPU, Linux, CUDA driver, Python runtime, and TensorFlow), are subject to software and hardware dependencies across the DL stack. A persistent challenge in dependency management across the entire engineering lifecycle is posed by the asynchronous and radical evolution as well as the complex version constraints among dependencies. Developers might introduce dependency bugs (DBs) in selecting, using and maintaining dependencies. However, the characteristics of DBs in DL stack is still under-investigated, hindering practical solutions to dependency management in DL stack.

To fill this gap, this paper presents the first comprehensive study to characterize symptoms, root causes and fix patterns of DBs across the whole DL stack with 326 DBs collected from StackOverflow posts. For each DB, we first investigate the symptom as well as the lifecycle stage and dependency where the symptom is exposed. Then, we analyze the root cause as well as the lifecycle stage and dependency where the root cause is introduced. Finally, we explore the fix pattern as well as the knowledge sources that are used to fix it. Our findings from this study shed light on the implications on dependency management, e.g., constructing dependency knowledge graph for the entire DL stack, recommending dependencies, detecting, localizing and fixing dependency bugs, and upgrading and migrating dependencies.

1 INTRODUCTION

The significant breakthroughs in deep learning (DL) techniques have brought great success to various DL-enabled application domains, e.g., machine translation [25], medical diagnosis [40], voice assistants [18] and autonomous vehicles [7]. These DL applications are built upon a heterogeneous and complex DL stack, including hardware (e.g., Nvidia GPU), OS (e.g., Linux), drivers (e.g., CUDA and cuDNN), runtime (e.g., Python) and libraries (e.g., TensorFlow). In other words, engineering DL applications requires the software and hardware in the DL stack as prerequisite dependencies. One common challenge in engineering DL applications is dependency management across the DL stack [1, 8, 62], i.e., to properly manage versions and configurations of the software and hardware dependencies in the entire DL stack.

Motivation. Dependency management is challenging due to three main reasons. First, *software and hardware dependencies are complex, and evolve quickly in an asynchronous and radical way.* Dependency complexity originates from two sources, i.e., deep stack and rich vendors. For example, many vendors provide DL libraries, e.g., Google’s

TensorFlow, Facebook’s PyTorch and Microsoft’s CNTK. Besides, dependency evolution is performed at the vendor’s own pace and might introduce incompatible changes. For example, the micro-architecture of Nvidia GPU has evolved several generations over the years, from old versions such as Tesla and Fermi to new versions such as Turing and Ampere. In the meantime, CUDA has evolved from version 1.0 to 11.6 to support different GPUs distinguished by compute capability, which ranges from 1.0 to 9.0. Therefore, developers might miss some dependencies and build an incomplete stack, or have troubles in selecting, updating and migrating dependency versions.

Second, *software and hardware dependencies have complex version constraints to satisfy for working together properly.* For example, each TensorFlow version only works compatibly with certain cuDNN versions, CUDA versions and Nvidia GPU versions. A developer tried to set up an environment with TensorFlow gpu version 1.2.0rc0, Python 3.5.2, CUDA 8.0.61, cuDNN 8.0, and a GPU card with compute capability 2.1 on Windows 7¹. This setup failed to recognize a valid GPU card because this TensorFlow version required a GPU card with compute capability 3.0 or higher. Similarly, a developer observed that importing and showing GPU output via TensorFlow APIs was slow². It turned out that the used TensorFlow 2.3 did not support CUDA 11 or higher, but the minimal requirement of the used GPU card was CUDA 11. These version constraints are usually scattered across the documentations of software and hardware. Therefore, developers might build an incompatible stack, or introduce dependency incompatibilities when updating versions or deploying to a new environment.

Third, *each dependency version may contain bugs or need proper configuration.* While dependency version constraints are satisfied, there might be bugs in specific versions under certain circumstances. For example, a developer created a Seq2Seq model using TensorFlow 1.5 but encountered an error³. It turned out that it was caused by a bug only in TensorFlow 1.5, and could be alleviated by upgrading to 1.6 or downgrading to 1.4. In addition, there might be misconfigurations during the installation of dependencies. For example, some kernel modules are required to be signed on Secure Boot enabled systems when the Nvidia driver is installed. However, this may cause unknown errors raised from CUDA⁴, which could be fixed by disabling Secure Boot. Therefore, developers might use a buggy dependency version or misconfigure a dependency version.

In summary, developers may introduce various dependency management problems in selecting, using and maintaining dependencies in the DL stack during the entire lifecycle (i.e., environment setup,

¹<https://stackoverflow.com/questions/44269047/>

²<https://stackoverflow.com/questions/65300388/>

³<https://stackoverflow.com/questions/48713335/>

⁴<https://stackoverflow.com/questions/67045622/>

*B. Chen is the corresponding author.

development, deployment and maintenance) of engineering DL applications. We refer to these problems as dependency bugs (DBs).

Literature. On the one hand, a lot of advances have been made to investigate DBs in different ecosystems, e.g., Java [19, 55], C/C++ [27], JavaScript [41], Python [36, 54], Go [53], and Debian and Red Hat [3]. They only consider DBs at the homogeneous library layer. However, DBs in the DL ecosystem are different because they can occur across all the heterogeneous layers in the DL stack. On the other hand, a lot of efforts have also been made to explore characteristics (e.g. symptoms, root causes and fix patterns) of general bugs [21, 23, 24, 39, 64] and specific bugs [6, 9, 52, 61, 63] in DL applications. However, these studies are not specifically designed for DBs, and thus only uncover partial characteristics of DBs in DL stack; e.g., dependency-related questions are often the most difficult to answer among all questions of DL applications on StackOverflow [1, 62]. Therefore, although it is necessary to understand the characteristics of DBs in DL stack, no systematic study exists yet, hindering dependency management.

Our Study. To bridge this knowledge gap, we present the first comprehensive study to characterize DBs in DL stack. An overview of our study is presented in Fig. 1. After a brief introduction of the DL stack (see Sec. 2), we first collect 326 DBs from StackOverflow posts and then analyze these DBs to answer three research questions (see Sec. 3).

- **RQ1 Symptom:** What are the symptoms of DBs? At which lifecycle stages and dependencies are the symptoms exposed?
- **RQ2 Root Cause:** What are the root causes of DBs? At which lifecycle stages and dependencies are the root causes introduced?
- **RQ3 Fix Pattern:** What are the fix patterns of DBs? Which knowledge sources are used to fix DBs?

Through these research questions, we aim to provide useful findings for developers and researchers (see Sec. 4, 5 and 6). For example, 36.2% of the DBs manifest DL specific errors/anomalies in software and hardware dependencies, behavior, model and data, mostly leading to crashes. Violation of constraints among software and hardware dependencies in DL stack causes 85.3% of the DBs. Development is the most bug-affecting lifecycle stage, which exposes 61.3% of the DBs, while environment setup is the most bug-prone lifecycle stage, which introduces 87.4% of the DBs. 43.6% of the DBs are not introduced and exposed in the same dependency. Changing dependency version and fixing API usage at the application code level are the most common fix patterns, which are used to fix 78.5% and 12.6% of the DBs. 26.1% of the DBs can be fixed by combining multiple fix patterns. Source code, documentation, issue tracker and other online resource are the main knowledge sources of fixing DBs.

Our findings provide implications for developers and researchers on dependency management across the entire engineering lifecycle (see Sec. 7), e.g., automated techniques to construct dependency knowledge graph for the entire DL stack, recommend dependencies in the entire DL stack, detect, localize and fix dependency bugs, and upgrade and migrate dependencies.

In summary, our work makes the following contributions.

- We conducted the first comprehensive study to characterize symptoms, root causes and fix patterns of 326 DBs in DL stack.
- We provided practical implications for developers and researchers on dependency management in engineering DL applications.

2 DEEP LEARNING STACK

Developers need to set up a DL environment before developing or deploying DL applications. The setup process often involves the following steps. First, developers need to choose a physical machine with GPUs and operating system installed. Besides, developers can use a virtual machine on the physical machine, or choose a virtual machine on the cloud supported by cloud service providers (e.g., Amazon SageMaker). Second, to fully empower upper libraries and DL applications, developers need to install the corresponding GPU drivers and GPU-accelerated SDKs (e.g., CUDA and cuDNN). Third, developers need to select a runtime environment based on the programming language that DL applications are developed with (e.g., Python and Java). Forth, a number of libraries should be leveraged to boost the development of DL applications from difference perspectives. Finally, developers could develop and deploy DL applications on top of the software and hardware dependencies.

This setup process is complicated by involving a wide scope of software and hardware dependencies. To reduce the complexity and provide a complete solution, a DL stack is proposed by organizing dependencies into layers. For example, Patterson shows a generic program stack consisting of modeling code, framework, storage, driver, operating system and hardware [42]. By following the setup process and referencing the DL stack at Patterson Consulting [42], Intel [22], Huawei [20] and Nvidia [44], we summarize a DL stack in Fig. 1. The stack consists of five layers. From top to bottom, they are *Application*, *Library*, *Runtime* and *Driver*, *OS/Container*, and *Hardware*.

Specifically, the *Application* layer contains DL applications from a variety of domains, e.g., facial recognition and autonomous driving. The *Library* layer contains the dependencies the upper-layer DL applications directly/transitively depend on. It covers a wide range of libraries, including frameworks (e.g., TensorFlow, PyTorch and CNTK) which provide abstraction and generic functionality implementation for DL algorithms, front-end libraries providing high-level abstraction or language bindings (e.g., Keras, ktrain and NeuPy), and other libraries in the ecosystem. The *Runtime* layer includes interpreters for dynamically typed languages (e.g., Python and JavaScript) and virtual machines for statically typed languages (e.g., Java and .Net). The *Driver* layer contains the dependencies for interacting with GPUs, including GPU drivers, computing platforms and GPU-accelerated SDKs (e.g., Nvidia GPU driver, CUDA and cuDNN). The *Library* layer can directly interact with the *Runtime* and *Driver* layer, and thus they are put at the same layer. The *OS/Container* layer contains operating systems, containers and other virtual environments (e.g., Ubuntu, Windows, macOS, Docker, and Amazon SageMaker). The *Hardware* layer contains fundamental hardware like CPU, GPU, mobile chips, and vendor-specific chips (e.g., Google's TPU).

3 EMPIRICAL STUDY METHODOLOGY

We first introduce the design of our empirical study, and then present our process of data collection and data labeling.

3.1 Study Design

The goal of our study is to characterize DBs in DL stack. To this end, we propose three research questions, which are introduced in Sec. 1.

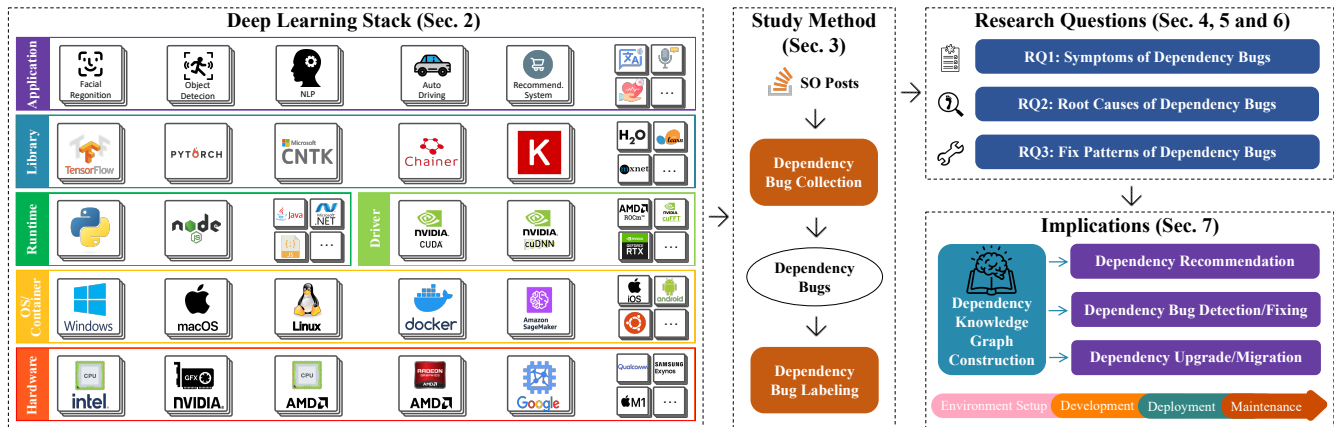


Figure 1: An Overview of Our Empirical Study on Dependency Bugs in DL Stack

Our *symptom analysis* in **RQ1** aims to characterize the observable consequences of DBs, which is helpful to assess impacts and provide insights for DB diagnosis and detection. Moreover, it aims to identify the lifecycle stage and dependency where the symptom is exposed, which is helpful to guide both developers and researchers to focus more effort on these bug-affecting stages and dependencies so as to achieve the most benefit for DB diagnosis and detection.

Our *root cause analysis* in **RQ2** seeks to understand the fundamental nature of DBs, which is helpful to provide insights for DB detection and localization. Further, it seeks to locate the lifecycle stage and dependency where the root cause is introduced, which is helpful to guide both developers and researchers to spend more effort on these bug-prone stages and dependencies in order to achieve the most benefit for DB avoidance, detection and localization.

Our *fix pattern analysis* in **RQ3** attempts to characterize the fixes of DBs, which is helpful to provide insights for DB fixing. In addition, it attempts to explore the distribution of fix patterns for root causes as well as the knowledge sources that are leveraged to fix DBs, which is helpful for both developers and researchers to achieve DB fixing in a more automated and effective fashion.

3.2 Data Collection

We collected DBs from a well-known Q&A site StackOverflow to obtain a comprehensive understanding of DBs. We selected StackOverflow because i) it is a popular site, and contains a wide range of questions raised by world-wide developers in real-life development activities; and ii) it has a high potential to contain questions about the dependencies in the entire DL stack due to its diversity. Our DB collection process consists of the following three steps.

Step 1: Dependency Tag Selection. Developers often attach several tags to a post to indicate the topics or concepts related to the question. Therefore, tags can be used to select the posts that are relevant to dependency problems in DL stack, and we need to determine a set of tags that have a high coverage of the dependencies in DL stack. To this end, we first collected all the 21,978,327 posts from Stack Exchange Data Dump⁵ on December 20, 2021. Then, for each post with an accepted answer, we iterated its tag list, and searched for tags that co-occurred with the tag “deep learning” or “neural network”. In

this way, we obtained an initial set of 1,576 tags. Here we did not directly use the tag “deep learning” or “neural network” to select posts because it may miss posts that were not tagged with “deep learning” and “neural network” but with other dependency related tags.

Next, two of the authors independently determined whether each of the 1,576 tags was related to the dependencies in DL stack by reading the excerpt provided by StackOverflow and online materials obtained by search engines. We used Cohen’s Kappa coefficient to measure agreement, and it reached 0.906. A third author was involved to resolve disagreements. Finally, we obtained 57 *Library* tags, 3 *Driver* tags, 59 *Runtime* tags, 23 *OS/Container* tags and 14 *Hardware* tags.

Step 2: Dependency Post Selection. We selected dependency-related posts in two steps. First, we selected from the 21,978,327 posts the ones whose tags contained one of the 57 *Library* tags and 3 *Driver* tags, or contained the tag “deep learning” or “neural network” as well as one of the 59 *Runtime* tags, 23 *OS/Container* tags and 14 *Hardware* tags. As *Runtime*, *OS/Container* and *Hardware* tags often have a weaker correlation with DL than *Library* and *Driver* tags, here we enforced their co-occurrence with either “deep learning” or “neural network” to reduce noisy posts. This led to 66,422 posts.

Second, to focus on high-quality posts, we removed the posts that did not have an accepted answer which contained dependency version information. The information of dependency versions was considered as important to determine root causes and fix patterns of DBs. We used regular expression matching to check the existence of version information. This restricted our selection to 3,814 posts.

Step 3: DB Identification. We manually verified the 3,814 posts to reduce noise that was not about DBs in DL stack. For example, some posts happened to have dependency-related tags but did not discuss dependencies; and some posts discussed general problems of dependencies (e.g., API usages). In particular, two of the authors independently investigated each post to identify DBs. The Cohen’s Kappa coefficient was 0.899. A third author was involved to resolve disagreements. Finally, we identified a total of 326 DBs.

3.3 Data Labeling

To answer the three research questions, we manually labeled each of the 326 DBs with respect to eight aspects, i.e., symptom, exposing stage, exposing dependency, root cause, introducing stage, introducing dependency, fix pattern, and knowledge source for fixing.

⁵<https://archive.org/download/stackexchange/>

In particular, two of the authors first randomly sampled 100 posts for a pilot labeling, following an open coding procedure [46]. They separately read all contents of a post (including title, question description, comments, answers, and reference links mentioned during discussion) and relied on search engines to carefully label DBs. Basically, the symptom of a DB was determined by analyzing the question description. The root cause, fix pattern and knowledge source for fixing of a DB were inferred from the question description and the accepted answer. The exposing stage and dependency of a DB were determined by analyzing where its symptom was exhibited, while the introducing stage and dependency of a DB were determined by analyzing where its root cause was located. A group discussion was conducted to exchange experience and summarize the initial taxonomies.

Then, two of the authors independently labeled all the 326 posts based on the initial taxonomies, and finally reached Cohen's Kappa coefficients of 0.905, 0.822, 0.806, 0.767, 0.993, 0.798, 0.961 and 0.856 for the eight labeling aspects. A third author was involved to resolve disagreements in both the pilot and final labeling. It is worth mentioning that the manual effort, involved in our data collection and labeling procedure, required six person-months.

4 RQ1: SYMPTOM ANALYSIS

We first present the taxonomy of DB symptoms, and then explore the stages and dependencies where DB symptoms are exposed.

4.1 Symptom Taxonomy

The taxonomy of DB symptoms is presented in Fig. 2. It is organized into five inner categories (i.e. *Syntactic Error*, *Performance Anomaly*, *DL Specific Error/Anomaly*, *Termination* and *Warning*) and 14 leaf categories. For each category, the number in parentheses refers to the number of DBs exhibiting the corresponding symptom.

Syntactic Error. 176 (54.0%) of the DBs exhibit general syntactic errors that are common and similar to those in traditional programs. It is the most common symptom. In particular, 92 (28.2%) of the DBs manifest *Element Not Found* errors; i.e., the used syntactic elements like module, class, function, key and attribute cannot be retrieved. In detail, *Module Not Found*, *Class Not Found*, *Function Not Found*, *Key Not Found* and *Attribute Not Found* errors respectively occur in 11, 6, 11, 7 and 57 of the DBs. Further, 32 (9.8%) of the DBs exhibit *Type Mismatch* errors; i.e., the variable type is inconsistent with the one that is expected. In addition, 21 (6.4%) and 16 (4.9%) of the DBs result in *Illegal Value* and *Illegal Argument* errors respectively, where a variable receives an illegal value, and a function call receives an illegal argument. Moreover, 7 (2.1%) of the DBs report *Undefined Variable* errors, denoting that the variable is not defined or initialized. Besides, some infrequent errors (e.g., compilation errors) are included in the *Others* category, which account for 8 (2.5%) of the DBs.

Performance Anomaly. 17 (5.2%) of the DBs manifest abnormal performance with respect to execution time, memory usage and processor usage. Specifically, 8 (2.5%) of the DBs exhibit *Long Execution Time*; i.e., a program takes a long time to initialize or execute DL tasks, or even hangs in the middle of the execution. Further, 8 (2.5%) of the DBs cause *Memory Anomaly*, including abnormal memory utilization, memory leak, or even out of memory errors. Besides, one DB results in *Processor Anomaly* (i.e., high GPU utilization).

DL Specific Error/Anomaly. 118 (36.2%) of the DBs exhibit DL specific errors or anomalies. It is the second most common symptom, and is divided into five leaf categories. *Software Error/Anomaly* means errors or anomalies raised by software dependencies, accounting for 64 (19.6%) of the DBs. There are four cases. (1) 14 (4.3%) of the DBs exhibit software internal errors, which are indicated by an error message that contains the software name, e.g., `CUDA_ERROR_UNKNOWN` and `CUDNN_STATUS_INTERNAL_ERROR`. (2) 32 (9.8%) of the DBs report that required software dependencies cannot be found. (3) 11 (3.4%) of the DBs manifest dependency initialization failures, indicating that required dependencies are not successfully set up. (4) 7 (2.1%) of the DBs report that required software dependency versions do not match or are not available from the repository.

Moreover, *Hardware Error/Anomaly* denotes errors or anomalies raised by hardware dependencies; e.g., the GPU card is not correctly connected. It accounts for 10 (3.1%) of the DBs. Further, 20 (6.1%) of the DBs manifest *Behavior Anomaly*, e.g., abnormal accuracy metrics and unexpected return values of APIs. In addition, 15 (4.6%) of the DBs exhibit *Model Error*, which is indicated by an error message that contains model elements, e.g., computation operator missing, model save/load failure, tensor conversion error, and layer unrecognized. Besides, 9 (2.8%) of the DBs manifest *Data Anomaly*, reporting that input data has abnormal values or mismatched property (e.g., size).

Termination. 11 (3.4%) of the DBs caused the program directly terminated without any informative error code or error message. For example, it only reports a segmentation fault, or it simply reports that the task is killed or canceled.

Warning. 4 (1.2%) of the DBs show warning messages, including warnings about function change, version compatibility, and semantic mismatch in API arguments. For example, a version compatibility warning reveals that the installed version violates the working version requirements. These warnings forecast the potential DBs due to using versions with changed elements.

Summary. General syntactic errors and DL specific errors and anomalies are the most common symptoms, which account for 90.2% of the DBs and mostly cause crashes. Besides, 5.2% of the DBs slow executions down or consume high resources. These severe consequences of DBs motivate the significance of DBs.

4.2 Exposing Stage and Dependency

We identify the stage and dependency where the symptom of each DB is exposed, and analyze DB distribution over these two dimensions.

Exposing Stage Analysis. We classify the entire lifecycle of engineering DL applications into four stages, i.e., environment setup, development, deployment, and maintenance. Our stage taxonomy is different from the one in Islam et al.'s work [23]. They mainly focus on the bugs in the development of DL applications, and thus classify the development pipeline into six stages, i.e., data preparation, model building, training, evaluation, hyper parameter tuning, and prediction. However, as indicated by our manual labeling, DBs can be introduced and exposed beyond the development process. Thus, we add the environment setup stage, where software and hardware dependencies are installed and configured to build a DL environment. We add the deployment stage, where a DL application is deployed

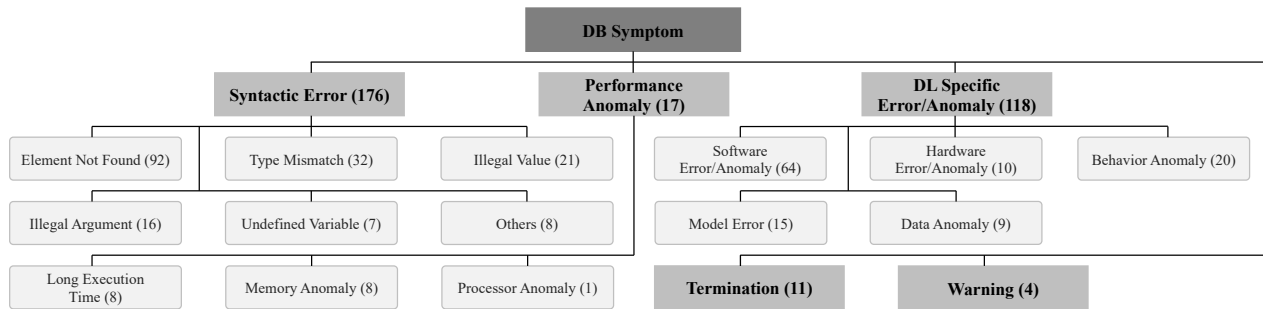


Figure 2: Taxonomy of DB Symptoms

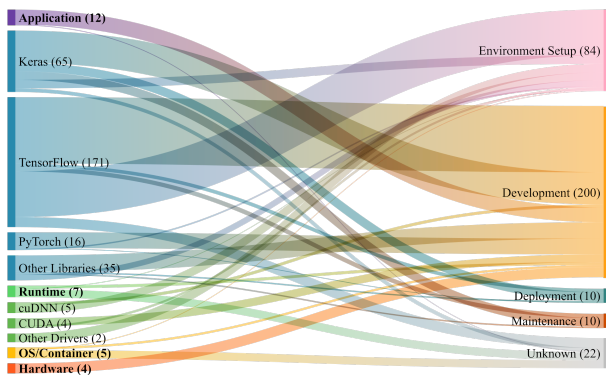


Figure 3: Exposing Dependency vs. Exposing Stage

and used for clients. We also add the maintenance stage, where developers conduct regular maintenance tasks like dependency upgrade, dependency migration, or DL application migration.

We report the DB distribution over the exposing stages in the right part of Fig. 3. Development is the most bug-affecting stage, where 200 (61.4%) of the DBs are exposed. This indicates that although the setup process of DL stack is presumably finished, a majority of DBs will not occur until DL application development. Environment setup is the second most bug-affecting stage, where 84 (25.8%) of the DBs are exposed. It indicates that the setup of a feasible DL stack is a non-trivial task. Apart from the two dominating stages, both deployment and maintenance expose 10 (3.1%) of the DBs, which are relatively smaller than in environment setup and development. One potential explanation is that developers have experienced DBs and built a feasible DL stack in environment setup and development, and hence most of the DBs are already fixed. The remaining 22 (6.7%) DBs have no clear indication about the exposing stage in the posts, and thus are included in the *Unknown* category.

Summary. The most bug-affecting stages are environment setup and development, which expose 87.1% of the DBs.

Exposing Dependency Analysis. We show the DB distribution over the exposing dependencies in the left part of Fig. 3, which is organized by the layer hierarchy in DL stack (see Fig. 2) with dominating dependencies separately highlighted. The *Library* layer is the most bug-affecting layer, where 287 (88.0%) of the DBs are exposed. Specifically, Keras, TensorFlow and PyTorch in the *Library* layer expose 65 (19.9%), 171 (52.5%) and 16 (4.9%) of the DBs respectively, which are

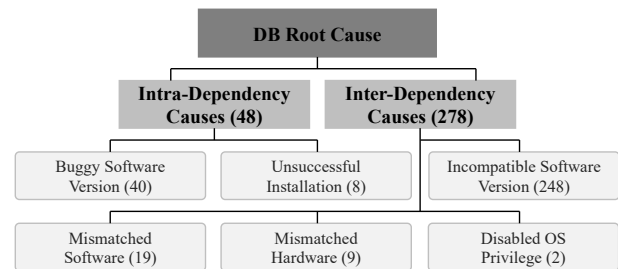


Figure 4: Taxonomy of DB Root Causes

the most bug-affecting libraries. This is reasonable as they are currently the most popular DL frameworks. The *Application* layer exposes 12 (3.7%) of the DBs, while the *Driver* layer exposes 11 (3.4%) of the DBs. CUDA and cuDNN are the most bug-affecting drivers, and respectively expose 4 (1.2%) and 5 (1.5%) of the DBs. Besides, there are at most 7 (2.1%) of the DBs that are exposed at the dependencies at the *Runtime*, *OS/Container* or *Hardware* layer.

Besides, the Sankey diagram in Fig. 3 illustrates where the DBs exposed in a dependency are exposed across the lifecycle stages. The width of the flow is proportional to the number of DBs. Generally, a DB can be exposed at any dependency at any layer in DL stack at any stage of the engineering lifecycle. It indicates the complexity of DBs.

Summary. *Library*, *Application* and *Driver* are the most bug-affecting stack layers, exposing 95.1% of the DBs. Keras, TensorFlow and PyTorch are the most bug-affecting libraries, and CUDA and cuDNN are the most bug-affecting drivers.

5 RQ2: ROOT CAUSE ANALYSIS

We first present the taxonomy of DB root causes, and then elaborate the stages and dependencies where DB root causes are introduced.

5.1 Root Cause Taxonomy

The taxonomy of DB root causes is shown Fig. 4. We first classify the root causes based on the criterion that whether a DB is caused solely by one dependency itself (i.e., *Intra-Dependency Cause*) or caused by constraints among dependencies across DL stack (i.e., *Inter-Dependency Cause*). Then, we summarize six leaf categories.

Intra-Dependency Cause. 48 (14.7%) of the DBs are caused solely by one dependency itself, and are divided into two leaf categories. Particularly, 40 (12.3%) of the DBs are caused by *Buggy Software Version*; i.e., a DB is caused by triggering bugs in software dependencies

in DL stack. For example, `tf.data.Dataset.from_generator()` did not work with strings on Python 3.x⁶, which was actually caused by a bug that affected Python 3.x and was fixed after TensorFlow 1.4.

Moreover, 8 (2.5%) of the DBs are caused by *Unsuccessful Installation* of dependencies. There are two cases. (1) Dependency installation does not complete. For example, a developer found that there was no file named `cuda64_6.dll` in the directory, which was caused by the missed installation of cuDNN on his/her machine⁷. (2) Dependency installation completes, but lacks proper path configuration (e.g., missing path configuration or configuring incorrect path).

Inter-Dependency Cause. 278 (85.3%) of the DBs are caused by constraints among software and hardware dependencies across DL stack; i.e., multiple dependencies have to be considered together to have a feasible DL stack, otherwise, DBs might be introduced. It is divided into four leaf categories. Specifically, 248 (76.1%) of the DBs are caused by *Incompatible Software Version*, which is the most common root cause. An incompatible software version is introduced if it violates the version constraint that has to be satisfied for it to work with other dependencies. For 56 of the 248 DBs, detailed API-level incompatibility information is provided in the posts, and we classify the incompatibility from the perspective of API changes [5], i.e., API removal, API addition, API replacement, API movement, API parameter list change, API renaming, and API behavior change. API addition, API behavior change and API removal are the most common root causes of API incompatibility, which respectively account for at least 19 (5.8%), 13 (4.0%) and 10 (3.1%) of the DBs. API replacement, API parameter list change, API movement and API renaming respectively cause at least 4, 4, 3 and 3 DBs. For the 138 of the 248 DBs, we can only distinguish whether they are caused by backward incompatibility (accounting for 80 (24.5%) of the DBs) or forward incompatibility (accounting for 58 (17.8%) of the DBs). For the remaining 54 of the 248 DBs, we can only determine they are caused by incompatibility due to the limited information in the posts.

19 (5.8%) of the DBs are caused by *Mismatched Software*; i.e., while different software can provide similar functionalities, only some of them can work with the other dependencies in DL stack, but others are regarded as mismatched. Specifically, 11 of the 19 DBs are caused by selecting wrong software as dependency. For example, the DL framework Keras and the `tf.keras` module introduced in TensorFlow 1.10 implement and provide similar APIs, but Keras does not support TensorFlow 2.0. In that sense, if TensorFlow 2.0 is used in DL stack, Keras would be mismatched and thus cannot be used⁸. Moreover, 5 of the 19 DBs are caused by choosing wrong software distribution. For example, the official pre-built TensorFlow 2.0 requires CUDA Toolkit 10.0. Developers have to re-build TensorFlow 2.0 with CUDA Toolkit 10.1 to work with CUDA Toolkit 10.1. Hence, using pre-built TensorFlow 2.0 with CUDA Toolkit 10.1 would cause a DB⁹. Further, 3 of the 19 DBs are caused by selecting multiple conflicting software. For example, loading both TensorFlow and TensorFlow-gpu¹⁰ or loading both Keras and `tf.keras`¹¹ would cause a DB.

⁶<https://stackoverflow.com/questions/47705684/>

⁷<https://stackoverflow.com/questions/43721690/>

⁸<https://stackoverflow.com/questions/56851895/>

⁹<https://stackoverflow.com/questions/58610020/>

¹⁰<https://stackoverflow.com/questions/42473052/>

¹¹<https://stackoverflow.com/questions/56851895/>

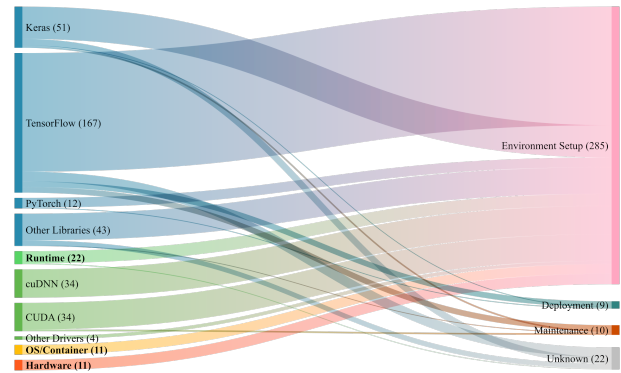


Figure 5: Introducing Dependency vs. Introducing Stage

9 (2.8%) of the DBs are caused by *Mismatched Hardware*; i.e., the hardware is incapable of meeting the requirements of dependencies in upper stack layers. For example, TensorFlow 1.6 started using AVX (Advanced Vector Extensions) feature of modern CPUs, which is supported by Sandy Bridge or newer CPU architectures. Hence, using TensorFlow with non-AVX CPUs would cause a DB¹².

2 (0.6%) of the DBs are caused by *Disabled OS Privilege*; i.e., the permissions required by software dependencies are not allowed from the OS or container. For example, the System Integrity Protection (SIP) is enabled on MacOS 10.11 to prevent the execution of unauthorized code. However, SIP prevents a path variable from being overridden, causing many dependencies not found¹³.

Summary. Violation of constraints among software and hardware dependencies in DL stack causes 85.3% of the DBs, where incompatible software version is the dominating root cause. Besides, bugs in software dependencies cause 12.3% of the DBs.

5.2 Introducing Stage and Dependency

We locate the stage and dependency where the root cause of each DB is introduced, and analyze DB distribution over these two dimensions.

Introducing Stage Analysis. The taxonomy of stages is the same to the one in Sec. 4.2. We present the DB distribution over the introducing stages in the right part of Fig. 5. Environment setup is the most bug-prone stage, where 285 (87.4%) of the DBs are introduced, while no DB is introduced in development because the DL stack is already determined in environment setup. It indicates that the setup of a feasible DL stack is important but challenging. Besides, deployment and maintenance introduce 9 (2.8%) and 10 (3.0%) of the DBs. The remaining 22 (6.7%) DBs have no clear indication about the introducing stage in the posts, and thus are put in the *Unknown* category.

Summary. The most bug-prone stage is environment setup, which introduces 87.4% of the DBs.

Introducing Dependency Analysis. For the DBs caused by inter-dependency causes, their root causes can be introduced by any of the involved dependencies. For example, a DB is caused by version constraint violation between TensorFlow and CUDA, and then both TensorFlow and CUDA can be the introducing dependency of this DB. If there is no clear indication about the introducing dependency in the

¹²<https://stackoverflow.com/questions/51599488/>

¹³<https://stackoverflow.com/questions/44354694/>

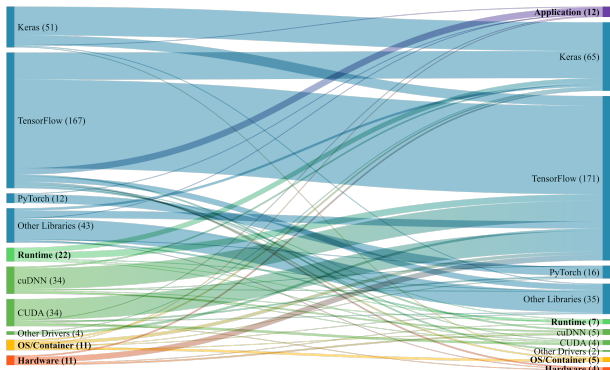


Figure 6: Introducing Dependency vs. Exposing Dependency

posts, we consider all involved dependencies as the introducing dependencies; otherwise, we use the introducing dependency that is determined in the posts. Specifically, of the 278 DBs that are caused by inter-dependency causes, 257 DBs have their introducing dependencies clearly indicated in the posts.

We show the DB distribution over the introducing dependencies in the left part of Fig. 5, which is organized in the same way in Fig. 3. No DB is introduced in the *Application* layer as it is the client of dependencies. The *Library* and *Driver* layers are the most bug-prone layers, respectively introducing 246 (75.5%) and 59 (18.1%) of the DBs. The number is larger than the summation of DBs in all libraries or drivers because a DB can have multiple introducing dependencies. In particular, Keras, TensorFlow and PyTorch in the *Library* layer are the most bug-prone libraries, introducing 51 (15.6%), 167 (51.2%) and 12 (3.7%) of the DBs respectively. CUDA and cuDNN are the most bug-prone drivers, respectively introducing 34 (10.4%) of the DBs. There are at most 22 (6.7%) of the DBs introduced at the dependencies at the *Runtime*, *OS/Container* or *Hardware* layer.

Besides, the Sankey diagram in Fig. 5 illustrates where the DBs introduced in a dependency are introduced across the lifecycle stages. Generally, a DB can be introduced at any dependency at any stack layer (except for *Application*) at any lifecycle stage (except for development). It reveals the complexity of DB localization.

Summary. *Library* and *Driver* are the most bug-prone stack layers, which introduce 90.5% of the DBs. Keras, TensorFlow and PyTorch are the most bug-prone libraries, while CUDA and cuDNN are the most bug-prone drivers.

5.3 Introducing and Exposing Dependency

We further analyze where the DBs introduced in a dependency are exposed across the dependencies in DL stack, which is illustrated in the Sankey diagram in Fig. 6. Overall, 142 (43.6%) of the DBs are not introduced and exposed in the same dependency. For example, 32 (9.8%) of the DBs introduced in TensorFlow are exposed in Keras, and 41 (12.6%) of the DBs introduced in CUDA and cuDNN are exposed in TensorFlow. At the stack layer level, 97 (29.8%) of the DBs are not introduced and exposed at the same stack layer. For example, 8 DBs introduced at the *Hardware* layer are exposed at the *Library* layer. These results indicate that DB localization is difficult, and need systematic knowledge over the entire DL stack.

Summary. 142 (43.6%) of the DBs are not introduced and exposed in the same dependency, while 97 (29.8%) of the DBs are not introduced and exposed at the same stack layer.

6 RQ3: FIX PATTERN ANALYSIS

We first show the taxonomy of DB fix patterns, and then report their distribution for root causes and the knowledge source used to fix DBs.

6.1 Fix Pattern Taxonomy

The taxonomy of DB fix patterns is reported in Fig. 7. It is organized into four inner categories (i.e., *Change Application Code*, *Change Dependency*, *Change DL Stack* and *Change Environment*) and 15 leaf categories. A DB can be fixed by applying multiple fix patterns. Hence, the summation of the number of DBs in Fig. 7 is larger than 326.

Change Application Code. 59 (18.1%) of the DBs are fixed via changing the application code although their root causes are not introduced by the application. Specifically, *Fixing API Usage* is used to fix 41 (12.6%) of the DBs; i.e., the library API usage has to be changed with the incompatible library version evolution. Moreover, *Adding Missing Code Logic* is utilized to fix 7 (2.1%) of the DBs. In such cases, some library APIs are removed or the behavior of some library APIs is changed, and hence developers have to implement the code logic of these library APIs by themselves at the application code level. Further, *Reformatting Data* is leveraged to fix 7 (2.1%) of the DBs in order to make the data format compatible with the changed library APIs. Besides, *Changing Hyper-Parameter* (e.g., batch size and learning rate) is used to fix 4 (1.2%) of the DBs, because the constraints on hyper-parameters are changed with library version evolution.

Change Dependency. This is the most common fix pattern, which is used to fix 304 (93.3%) of the DBs. It is divided into five leaf categories. In particular, *Changing Dependency Version* is leveraged to fix 256 (78.5%) of the DBs, indicating that it is the most common pattern to fix DBs. Of these 256 DBs, upgrading dependency version is used in the fix of 153 DBs, and downgrading dependency version is used in the fix of 104 DBs. In 17 of the DBs, dependency version is changed but there is no clear indication in the posts to determine upgrade or downgrade. Further, *Adding Dependency* is used to fix 23 (7.1%) of the DBs where some required dependencies are missing or not successfully installed. Moreover, *Re-building Dependency* is used to fix 19 (5.8%) of the DBs. In such cases, the source code of dependencies is re-built with other required dependencies to properly work with them, or the source code of dependencies is first changed (e.g., to fix bugs or to remove incompatibilities) and then re-built, potentially because of the huge maintenance effort in changing dependency versions. In addition, *Changing Dependency Configuration* is leveraged to fix 4 (1.2%) of the DBs, e.g., disabling SIP in MacOS. Besides, *Removing Dependency* is applied to fix 2 (0.6%) of the DBs in order to remove conflicted dependencies.

Change DL Stack. 22 (6.7%) of the DBs are fixed by changing the DL stack; i.e., some dependencies are switched to alternatives, and hence the DL stack becomes fundamentally different. Specifically, it is divided into three leaf categories, i.e., *Switching Software* (including libraries, drivers and runtimes), which accounts for 11 (3.4%) of the DBs, *Switching Hardware*, which accounts for 6 (1.8%) of the DBs, and *Switching OS*, which accounts for 5 (1.5%) of the DBs. For

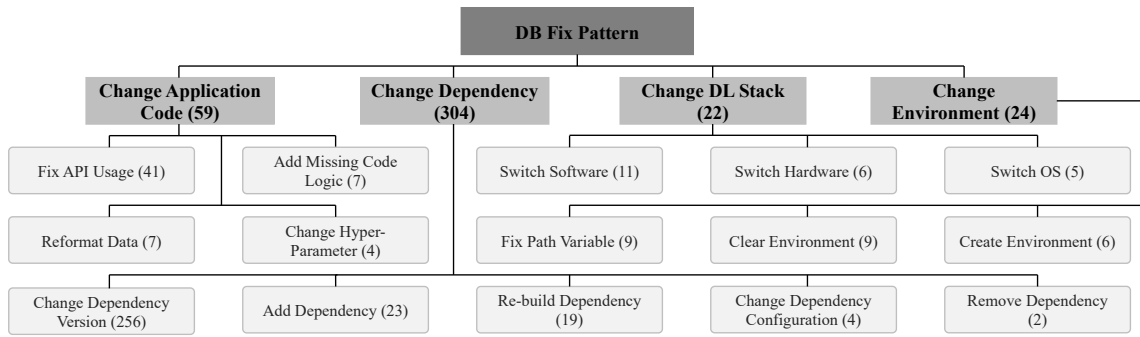


Figure 7: Taxonomy of DB Fix Patterns

example, a developer switched the OS from Windows to Ubuntu to successfully support distributed TensorFlow¹⁴.

Change Environment. 24 (7.4%) of the DBs are fixed by changing the environment where libraries, drivers and runtimes can be found. Specifically, *Fixing Path Variable* is used to fix 9 (2.8%) of the DBs; i.e., the path variable is fixed to point to the correct directory that contains the required dependencies. Besides, *Clearing Environment* and *Creating Environment* are used to respectively fix 9 (2.8%) and 6 (1.8%) of the DBs. In these cases, the virtual environment (i.e., a directory that contains a specific collection of installed packages) of package managers (e.g., pip and conda) is cleared or created.

It is worth mentioning that 279 (85.6%) of the DBs can be fixed by applying one fix pattern, while 64 (19.6%), 18 (5.5%) and 3 (0.9%) of the DBs can be fixed by combining two, three and four fix patterns at the same time. The summation here is larger than 326 because 36 DBs can be fixed by several different combinations of fix patterns.

Summary. The most common fix pattern is to change dependency versions, which is used to fix 78.5% of the DBs. Fixing API usage at the application code level is the second most common pattern, which is leveraged to fix 12.6% of the DBs. 26.1% of the DBs can be fixed by combining multiple fix patterns.

6.2 Distribution of Fix Patterns for Root Causes

We report the distribution of fix patterns for root causes in Fig. 8, where each cell denotes the number of DBs that are caused by a particular root cause and fixed by a particular fix pattern. Specifically, except for *Removing Dependency* and *Switching Software*, all fix patterns are utilized in fixing DBs that are caused by *Incompatible Software Version* for at least once. It is reasonable because this root cause results from the version constraint violation among dependencies. In other words, the involved dependencies are needed, but their versions are not satisfiable. Therefore, it is unlikely to fix it by removing dependencies or switching software. Besides, while *Fixing API Usage* and *Changing Dependency Version* are the two dominating fix patterns, there still exist diverse ways to fix the most common root cause *Incompatible Software Version*. The key challenge is to determine which fix pattern to use given a specific DB context.

Further, *Changing Dependency Version* is leveraged in mitigating four root causes, and is involved in the fix for 256 (78.6%) of the DBs.

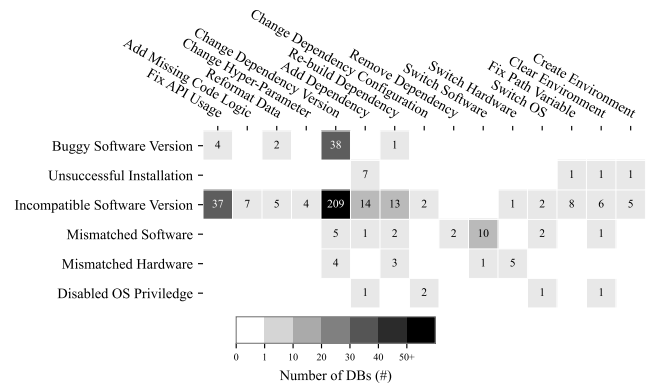


Figure 8: Distribution of Fix Patterns for Root Causes

It has a strong correlation to the root causes *Buggy Software Version* and *Incompatible Software Version*. While the fix pattern itself is very simple, the key challenge is to determine which dependency version to use for addressing a DB. Besides, *Adding Dependency*, *Re-building Dependency* and *Clearing Environment* are the other three fix patterns spanning four root causes. It is worth mentioning that *Adding Dependency* is the accompanied fix pattern for fixing DBs caused by *Incompatible Software Version*. For example, upgrading dependency version solves an *Incompatible Software Version*, but this upgraded dependency version may further depend on a new dependency.

Summary. DBs caused by *Incompatible Software Version* can be fixed by diverse fix patterns, whereas *Fixing API Usage* and *Changing Dependency Version* are the two dominating fix patterns. *Changing Dependency Version* is also the dominating fix pattern for DBs cause by *Buggy Software Version*.

6.3 Knowledge Source of DB Fixing

To fix a DB, developers often rely on certain knowledge about DL stack, e.g., dependency version constraints and dependency bugs. To characterize how fixes of DBs are derived, we investigate the knowledge sources that are used to fix DBs. We identify five knowledge sources. Multiple knowledge sources can be used in fixing one DB, and thus the summation of the number of DBs below is larger than 326.

Library Source Code. 51 (15.6%) of the DBs are fixed after digging into the source code of libraries. The source code of libraries is a good knowledge source to know library version evolution, e.g., how

¹⁴<https://stackoverflow.com/questions/52041077>

a library API is renamed or moved, how a library API's code logic is changed, and to what extent a library API is broken.

Dependency Documentation. 62 (19.0%) of the DBs are fixed after looking into the dependency documentation. Documentation of libraries, drivers and hardware often provide informative knowledge about dependency's installation requirements and version constraints. For example, to use TensorFlow with GPU support, TensorFlow documentation lists both the hardware requirements and software requirements¹⁵. Similarly, CUDA Toolkit documentation lists the required driver versions on Linux and Windows¹⁶.

Issue Tracker. 18 (5.5%) of the DBs are fixed after being aware of the dependency bugs. Such bugs are tracked on issue trackers with their symptoms and affected versions described in detail.

Other Online Resource. 19 (5.8%) of the DBs are fixed by referencing other online resources, e.g., mailing lists, StackOverflow posts and technical blogs.

Unknown. For 200 (61.3%) of the DBs, there is no clear indication about the leveraged knowledge source in the posts, and hence we include them into the *Unknown* category. However, such StackOverflow posts themselves actually become the knowledge source.

Summary. Library source code, dependency documentation, issue tracker, and other online resource are the main knowledge sources that are directly used in the posts to fix DBs.

7 IMPLICATIONS AND THREATS

We discuss the implications for developers and researchers from the findings of our study, and analyze the threats to our study.

7.1 Implications

Developers. Our study uncover the common symptoms of DBs that developers should be aware of when engineering DL applications for detecting potential DBs as early as possible. Our study also identifies the common root causes and fix patterns of DBs that could be useful for developers to diagnose, localize and fix DBs. Our study also shows the most bug-prone and bug-affecting dependencies where developers should pay more attention when installing, using or maintaining them so that most DBs could be avoided or detected at the first place. Moreover, our findings provide some engineering suggestions. Developers should be trained to have a comprehensive understanding of the DL stack, as our study reports that a DB could be introduced or exposed across the entire DL stack and engineering lifecycle. In this way, developers are equipped with the sufficient knowledge to deal with DBs. Developers should carefully look into the dependency documentation to learn version constraints, and be aware of the bugs and API changes in library version evolution. In this way, DBs caused by the most common root causes (i.e., *Buggy Software Version* and *Incompatible Software Version*) might be effectively reduced.

Researchers. Our findings provide future research implications in four directions. First, *a dependency knowledge graph for the entire DL stack is needed to provide fundamental knowledge for the ease of dependency management*. As uncovered by our root cause analysis, a diversity of dependency knowledge is involved in DBs, e.g., version constraints among software and hardware dependencies, bugs in

dependencies, and API changes in version evolution. However, such knowledge is scattered across different sources, e.g., documentation, issue tracker and source code, as revealed by our investigation of the knowledge source of DB fixing. Besides, online resources like StackOverflow posts also provide practical solutions to fix DBs. Hence, the main challenges to construct the knowledge graph are that i) designing a high-level schema to fuse various knowledge into a graph, ii) leveraging various techniques like natural language processing and program analysis to automatically extract knowledge from different sources and keep them up-to-date; and iii) developing graph analysis techniques for various dependency management tasks. As shown in Fig. 1, this knowledge graph serves as the foundation of the following three research directions.

Second, *dependency recommendation techniques are needed*. Our introducing stage analysis reveals that environment setup is the most bug-prone stage which introduces 87.4% of the DBs. Therefore, developers often face difficulties in setting up a feasible DL stack. Further, our root cause analysis shows that 76.1% of the DBs are caused by *Incompatible Software Version*, although dependency documentation provides prerequisite information about setting up dependencies and their version constraints. Therefore, developers might not always refer to the documentation. In that sense, dependency recommendation techniques become useful for developers to ease the setup of a feasible DL stack; i.e., given some dependencies installed, they recommend other dependencies to form a complete DL stack. For example, given the available hardware and OS, they suggest the required dependencies in the *Driver*, *Runtime* and *Library* layers.

Third, *DB detection, localization and fixing techniques are needed*. Our study indicates that 87.4% of the DBs are introduced in environment setup, while only 25.8% of the DBs are exposed in environment setup. Hence, many DBs stay undetected until later lifecycle stages. To detect or localize DBs as early as possible, one possible remedy is to first identify the dependencies currently adopted in the DL stack, and then check against our dependency knowledge graph to detect potential dependency constraint violations. Here the challenging is to automatically identify all heterogeneous dependencies as well as their versions across the entire DL stack. Along this direction, Tan et al. [49] proposed a technique to identify homogeneous dependencies at the *Application* and *Library* layer. Moreover, as many DBs are caused by software bugs or API incompatibilities, fine-grained call graph analysis is needed to accurately detect and localize DBs, i.e., to determine whether such bugs or incompatible APIs are in the execution path and thus can be triggered. Once a DB is localized, automated fixing techniques can use the fix patterns derived from our study to fix it. However, the challenge is to determine which fix pattern or combination of fix patterns is applicable and how a fix pattern is instantiated. One potential way is to use a search-based method by applying fix patterns to generate potential fixes and using the dependency knowledge graph to decide the fitness of fixes.

Fourth, *dependency upgrading and migration techniques are needed*. Our introducing stage analysis uncovers that some DBs are introduced in deployment and maintenance. More specifically, the DL stack in deployment environment can be different from the one in development environment. Therefore, dependency migration techniques are needed to check whether dependencies in development environment can be replaced with the one in deployment environment. Besides, in maintenance, dependency versions can be upgraded for

¹⁵<https://www.tensorflow.org/install/gpu>

¹⁶<https://docs.nvidia.com/cuda/cuda-toolkit-release-notes/index.html>

the benefit of fixed bugs and improved features. However, it may also introduce incompatibilities. Therefore, dependency upgrading techniques are needed to analyze the API changes and assess the risk in terms of potential DBs and the effort in terms of potential code adaptation in order to safely upgrade dependency versions.

7.2 Threats

One threat to our study is about the generality of DL stack. We cannot claim that the DL stack can be generalizable. However, we derive it by referencing the DL stack at various industrial companies (as introduced in Sec. 2). In fact, we also continuously refine the DL stack during our data labeling procedure such that every identified dependency from the posts has a proper position in the DL stack. Therefore, we believe the DL stack is reasonable and representative.

Another threat to our study is about the collection source of DBs. We collect DBs from StackOverflow posts due to the wide topics of questions. We also attempt to collect DBs from the commit history of GitHub projects. However, GitHub projects mainly contain DBs at the *Application* and *Library* layer as GitHub only hosts the source code. Therefore, DBs at the lower layers are often not contained in GitHub projects. Therefore, to avoid potential DB distribution biases, we decide to only use StackOverflow as the source of DBs.

The other threat to our study is about the manual analysis involving in our data labeling procedure, which might introduce biases. To mitigate it, three authors make great effort (i.e., around six person-months) to carefully carry out our analysis by following an open coding procedure, and use the Cohen's Kappa coefficient to measure the labeling agreement (which turns out to be satisfiable).

8 RELATED WORK

We review and discuss the relevant work in three dimensions, i.e., dependency bugs, deep learning bugs, and empirical studies about DL.

8.1 Dependency Bugs

Dependency bugs have been widely investigated for different ecosystems. Artho et al. [3] categorized dependency conflicts in Debian and Red Hat, and discussed potential solutions to common categories. Patra et al. [41] investigated dependency conflicts caused by the lack of namespaces in JavaScript, and proposed CONFLICTJS to combine dynamic analysis and targeted test synthesis to validate them. Wang et al. [54–57] studied the manifestation and fixing patterns of dependency conflicts in Java and Python, and developed automated techniques to detect, test or monitor them. Instead of detecting conflicts among dependencies, Jia et al. [27] identified dependency incompatibilities caused by incompatible changes within dependencies in C/C++. Similarly, Mukherjee et al. [36] focused on dependency incompatibilities in Python, and proposed to detect and fix them to ensure build reproducibility. Huang et al. [19] studied inconsistent dependency versions (i.e., multiple versions of a dependency are used in a project) in Java, and proposed to harmonize inconsistent versions. Wang et al. [53] analyzed missing or wrong dependencies caused by two mixed dependency management mechanisms in Go, and proposed to detect them and suggest fixes. Macho et al. [34] proposed a rule-based approach to repair Maven builds that were broken due to dependency bugs. In make-based build systems, several advances

have been made to detect unspecified dependencies [4], missing dependencies [11, 12, 29, 48] and redundant dependencies [11, 48] in build scripts, which could lead to build failures, incorrect builds or poor performance. Similarly, Ghorbani et al. [13] detected inconsistent module dependencies specified in Java-9 applications.

Our work differs from the prior work in two aspects. First, to the best of our knowledge, our work is the first to systematically investigate dependency bugs in the DL ecosystem, which is different from previous ecosystems. Second, our work analyzes dependency bugs across different layers in the DL stack, which is a unique characteristic that is absent in previous ecosystems.

8.2 Deep Learning Bugs

Increasing interest has been gained in understanding characteristics, e.g., symptoms, types, root causes, impacts, pipeline stages and fix patterns, of bugs in DL systems through analyzing StackOverflow posts, inspecting GitHub commits or conducting developer interviews [21, 23, 24, 39, 64]. These studies are focused on a general scope of bugs in DL systems, while recent studies have started to explore more specific bugs in DL systems, e.g., DL job failures [61], deployment faults on mobile devices [9], training problems [63], cloud API misuses [52] and performance problems [6]. Moreover, Jia et al. [26] and Shen et al. [47] analyzed symptoms and root causes of bugs in DL library and compiler.

These studies uncover partial characteristics of dependency bugs in DL stack. Some studies [23, 24, 26, 47, 61, 64] recognized incompatibilities due to dependency API changes or version updates as a common bug root cause. Some studies [6, 9, 61] identified buggy, missing or mismatched dependency versions as the bug root cause. Despite these efforts, there still lacks a comprehensive study to characterize dependency bugs in DL stack, and our work aims to bridge this gap and raise the awareness of dependency bugs in DL stack.

Besides, several advances have been made to detect DL bugs. For example, Zhang et al. [65], Wardat et al. [58] and Yan et al. [60] identified numerical bugs by static analysis and dynamic analysis. Lagouvardos et al. [28], Verma and Su [51], Liu et al. [30] and Wu et al. [59] detected shape bugs by static analysis, dynamic analysis, constraint solving and machine learning. Nikanjam et al. [37] detected faults expressed as rules by graph-based verification. However, little attention has been received to detecting dependency bugs in DL stack, and our work sheds light on this direction.

8.3 Empirical Studies about DL

Many studies have empirically investigated different aspects in developing, deploying and maintaining DL systems, e.g., software engineering for DL systems [2, 14, 35], challenges in developing DL systems [1, 62] and deploying DL systems [8], pain-points in using cloud services of computer vision [10], accuracy variance in training DL systems [43], performance variance in deploying DL models to different mobile devices and web browsers [15, 33], discussion topics about DL frameworks [17], API evolution in DL frameworks [66], technical debt in DL frameworks [31, 32] and DL systems [38, 45, 50], and dependency supply chain of DL libraries [16, 49].

Although these studies are not designed for dependency bugs, they motivate the importance of dependency management. For example, Han et al. [17] pinpointed the need of solving version problems and improving the compatibility of DL libraries. Liu et al. [32] found that

the removal of dependency compatibility debt was the slowest. Chen et al. [8], Zhang et al. [62] and Alshangiti et al. [1] recognized incompatible dependency installation or environment setup as a common challenge in developing and deploying DL systems, due to complex software and hardware dependencies and dependency version incompatibilities. Han et al. [16] investigated the dependency of open-source projects on three DL libraries (i.e., TENSORFLOW, PYTORCH and THEANO), while Tan et al. [49] built a supply chain for TENSORFLOW and PYTORCH via identifying packages and downstream projects that directly/transitively depended on TENSORFLOW and PYTORCH. However, these two studies do not consider dependencies across the whole DL stack. Our work is inspired by these studies to systematically characterize dependency bugs across the whole DL stack.

9 CONCLUSIONS

In this paper, we have conducted the first comprehensive study to understand the characteristics (e.g., symptoms, root causes and fix patterns) of DBs across the entire DL stack with a collection of 326 DBs from StackOverflow posts. We provide useful findings to raise the awareness of DBs in DL stack in the DL community and uncover the nature of DBs. We also provide actionable implications for both developers and researchers on dependency management across the engineering lifecycle. In future, we plan to develop techniques to automatically construct the dependency knowledge graph for the entire DL stack and build solutions to dependency management tasks upon this graph. The data of our study is available at our replication site <https://dl-dep.github.io> to foster future research.

REFERENCES

- [1] Moayad Alshangiti, Hitesh Sapkota, Pradeep K Murukannaiah, Xumin Liu, and Qi Yu. 2019. Why is developing machine learning applications challenging? a study on stack overflow posts. In *Proceedings of the ACM/IEEE International Symposium on Empirical Software Engineering and Measurement*. 1–11.
- [2] Salema 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 *Proceedings of the IEEE/ACM 41st International Conference on Software Engineering: Software Engineering in Practice*. 291–300.
- [3] Cyrille Artho, Kuniyasu Suzuki, Roberto Di Cosmo, Ralf Treinen, and Stefano Zacchiroli. 2012. Why do software packages conflict?. In *Proceedings of the 9th IEEE Working Conference on Mining Software Repositories*. 141–150.
- [4] Cor-Paul Bezemer, Shane McIntosh, Bram Adams, Daniel M German, and Ahmed E Hassan. 2017. An empirical study of unspecified dependencies in make-based build systems. *Empirical Software Engineering* 22, 6 (2017), 3117–3148.
- [5] Aline Brito, Laerte Xavier, Andre Hora, and Marco Tulio Valente. 2018. APIDiff: Detecting API breaking changes. In *Proceedings of the IEEE 25th International Conference on Software Analysis, Evolution and Reengineering*. 507–511.
- [6] Junming Cao, Bihuan Chen, Chao Sun, Longjie Hu, and Xin Peng. 2021. Characterizing Performance Bugs in Deep Learning Systems. *CoRR* abs/2112.01771 (2021).
- [7] Chenyi Chen, Ari Seff, Alain Kornhauser, and Jianxiong Xiao. 2015. Deepdriving: Learning affordance for direct perception in autonomous driving. In *Proceedings of the IEEE international conference on computer vision*. 2722–2730.
- [8] Zhenpeng Chen, Yanbin Cao, Yuanqiang Liu, Haoyu Wang, Tao Xie, and Xuanzhe Liu. 2020. A comprehensive study on challenges in deploying deep learning based software. In *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 750–762.
- [9] 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 *Proceedings of the IEEE/ACM 43rd International Conference on Software Engineering*. 674–685.
- [10] Alex Cummaudo, Rajesh Vasa, Scott Barnett, John Grundy, and Mohamed Abdelrazek. 2020. Interpreting Cloud Computer Vision Pain-Points: A Mining Study of Stack Overflow. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1584–1596.
- [11] Gang Fan, Chengpeng Wang, Rongxin Wu, Xiao Xiao, Qingkai Shi, and Charles Zhang. 2020. Escaping dependency hell: finding build dependency errors with the unified dependency graph. In *Proceedings of the 29th ACM SIGSOFT International Symposium on Software Testing and Analysis*. 463–474.
- [12] Anders Fischer-Nielsen, Zhoulai Fu, Ting Su, and Andrzej Wąsowski. 2020. The forgotten case of the dependency bugs: On the example of the robot operating system. In *Proceedings of the IEEE/ACM 42nd International Conference on Software Engineering: Software Engineering in Practice*. 21–30.
- [13] Negar Ghorbani, Joshua Garcia, and Sam Malek. 2019. Detection and repair of architectural inconsistencies in java. In *Proceedings of the IEEE/ACM 41st International Conference on Software Engineering*. 560–571.
- [14] Görkem Giray. 2021. A software engineering perspective on engineering machine learning systems: State of the art and challenges. *Journal of Systems and Software* 180 (2021), 111031.
- [15] Qianyu Guo, Sen Chen, Xiaofei Xie, Lei Ma, Qiang Hu, Hongtao Liu, Yang Liu, Jianjun Zhao, and Xiaohong Li. 2019. An empirical study towards characterizing deep learning development and deployment across different frameworks and platforms. In *Proceedings of the 34th IEEE/ACM International Conference on Automated Software Engineering*. 810–822.
- [16] Junxiao Han, Shuiguang Deng, David Lo, Chen Zhi, Jianwei Yin, and Xin Xia. 2020. An empirical study of the dependency networks of deep learning libraries. In *Proceedings of the IEEE International Conference on Software Maintenance and Evolution*. 868–878.
- [17] Junxiao Han, Emad Shihab, Zhiyuan Wan, Shuiguang Deng, and Xin Xia. 2020. What do programmers discuss about deep learning frameworks. *Empirical Software Engineering* 25, 4 (2020), 2694–2747.
- [18] Johann Hauswald, Michael A Laurenzano, Yunqi Zhang, Cheng Li, Austin Rovinski, Arjun Khurana, Ronald G Dreslinski, Trevor Mudge, Vinicius Petrucci, Lingjia Tang, et al. 2015. Sirius: An open end-to-end voice and vision personal assistant and its implications for future warehouse scale computers. In *Proceedings of the Twentieth International Conference on Architectural Support for Programming Languages and Operating Systems*. 223–238.
- [19] Kaifeng Huang, Bihuan Chen, Bowen Shi, Ying Wang, Congying Xu, and Xin Peng. 2020. Interactive, effort-aware library version harmonization. In *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 518–529.
- [20] Huawei. 2020. *Overview of the Ascend AI Software Stack*. https://support.huaweicloud.com/intl/en-us/access-g-Atlas200dkappc32/atlaspd_19_0001.html#atlaspd_19_0001_fig10978124614814
- [21] Nargiz Humbatova, Gunel Jahangirova, Gabriele Bavota, Vincenzo Riccio, Andrea Stocco, and Paolo Tonella. 2020. Taxonomy of real faults in deep learning systems. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1110–1121.
- [22] Intel. 2020. *Deep Learning Reference Stack*. <https://intel.github.io/stacks/dlrs/index.html>
- [23] Md Johirul Islam, Giang Nguyen, Rangeet Pan, and Hriday Rajan. 2019. A comprehensive study on deep learning bug characteristics. In *Proceedings of the 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 510–520.
- [24] Md Johirul Islam, Rangeet Pan, Giang Nguyen, and Hriday Rajan. 2020. Repairing deep neural networks: Fix patterns and challenges. In *Proceedings of the IEEE/ACM 42nd International Conference on Software Engineering*. 1135–1146.
- [25] Sébastien Jean, KyungHyun Cho, Roland Memisevic, and Yoshua Bengio. 2015. On Using Very Large Target Vocabulary for Neural Machine Translation. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing of the Asian Federation of Natural Language Processing*. 1–10.
- [26] Li Jia, Hao Zhong, Xiaoyin Wang, Linpeng Huang, and Xuansheng Lu. 2020. An Empirical Study on Bugs Inside TensorFlow. In *Proceedings of the International Conference on Database Systems for Advanced Applications*. 604–620.
- [27] Zhouyang Jia, Shanshan Li, Tingting Yu, Chen Zeng, Erci Xu, Xiaodong Liu, Ji Wang, and Xiangke Liao. 2021. DepOwl: Detecting Dependency Bugs to Prevent Compatibility Failures. In *Proceedings of the IEEE/ACM 43rd International Conference on Software Engineering*. 86–98.
- [28] Sifis Lagouvardos, Julian Dolby, Neville Grech, Anastasios Antoniadis, and Yanis Smaragdakis. 2020. Static analysis of shape in TensorFlow programs. In *Proceedings of the 34th European Conference on Object-Oriented Programming*. 1–29.
- [29] Nándor Licker and Andrew Rice. 2019. Detecting incorrect build rules. In *Proceedings of the IEEE/ACM 41st International Conference on Software Engineering*. 1234–1244.
- [30] 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 *Proceedings of the 36th IEEE/ACM International Conference on Automated Software Engineering*. 55–66.
- [31] 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 *Proceedings of the ACM/IEEE 42nd International*

- Conference on Software Engineering: Software Engineering in Society*. 1–10.
- [32] Jiakun Liu, Qiao Huang, Xin Xia, Emad Shihab, David Lo, and Shanping Li. 2021. An exploratory study on the introduction and removal of different types of technical debt in deep learning frameworks. *Empirical Software Engineering* 26, 2 (2021), 1–36.
- [33] Yun Ma, Dongwei Xiang, Shuyu Zheng, Deyu Tian, and Xuanzhe Liu. 2019. Moving deep learning into web browser: How far can we go?. In *Proceedings of the World Wide Web Conference*. 1234–1244.
- [34] Christian Macho, Shane McIntosh, and Martin Pinzger. 2018. Automatically repairing dependency-related build breakage. In *Proceedings of the IEEE 25th International Conference on Software Analysis, Evolution and Reengineering*. 106–117.
- [35] Silverio Martínez-Fernández, Justus Bogner, Xavier Franch, Marc Oriol, Julien Siebert, Adam Trendowicz, Anna Maria Vollmer, and Stefan Wagner. 2022. Software Engineering for AI-Based Systems: A Survey. *ACM Transactions on Software Engineering and Methodology* 31, 2 (2022), 1–59.
- [36] Suchita Mukherjee, Abigail Almanza, and Cindy Rubio-González. 2021. Fixing dependency errors for Python build reproducibility. In *Proceedings of the 30th ACM SIGSOFT International Symposium on Software Testing and Analysis*. 439–451.
- [37] Amin Nikanjam, Housseem Ben Braiek, Mohammad Mehdi Morovati, and Foutse Khomh. 2021. Automatic fault detection for deep learning programs using graph transformations. *ACM Transactions on Software Engineering and Methodology (TOSEM)* 31, 1 (2021), 1–27.
- [38] Amin Nikanjam and Foutse Khomh. 2021. Design Smells in Deep Learning Programs: An Empirical Study. In *Proceedings of the IEEE International Conference on Software Maintenance and Evolution*. 332–342.
- [39] Amin Nikanjam, Mohammad Mehdi Morovati, Foutse Khomh, and Housseem Ben Braiek. 2022. Faults in deep reinforcement learning programs: a taxonomy and a detection approach. *Automated Software Engineering* 29, 1 (2022), 1–32.
- [40] Ziad Obermeyer and Ezekiel J Emanuel. 2016. Predicting the future—big data, machine learning, and clinical medicine. *The New England journal of medicine* 375, 13 (2016), 1216.
- [41] Jibesh Patra, Pooja N Dixit, and Michael Pradel. 2018. Conflictjs: finding and understanding conflicts between javascript libraries. In *Proceedings of the 40th International Conference on Software Engineering*. 741–751.
- [42] Josh Patterson. 2019. *A Practical Guide for Data Scientists Using GPUs with TensorFlow*. http://www.pattersonconsultingtm.com/blog/datascience_guide_tensorflow_gpus.html
- [43] Hung Viet Pham, Shangshu Qian, Jiannan Wang, Thibaud Lutellier, Jonathan Rosenthal, Lin Tan, Yaoliang Yu, and Nachiappan Nagappan. 2020. Problems and opportunities in training deep learning software systems: an analysis of variance. In *Proceedings of the 35th IEEE/ACM International Conference on Automated Software Engineering*. 771–783.
- [44] Felix Abecassis Ryan Olson, Jonathan Calmels and Phil Rogers. 2016. *NVIDIA Docker: GPU Server Application Deployment Made Easy*. <https://developer.nvidia.com/blog/nvidia-docker-gpu-server-application-deployment-made-easy/>
- [45] D. Sculley, Gary Holt, Daniel Golovin, Eugene Davydov, Todd Phillips, Dietmar Ebner, Vinay Chaudhary, Michael Young, Jean-Francois Crespo, and Dan Dennison. 2015. Hidden Technical Debt in Machine Learning Systems. In *Proceedings of the 28th International Conference on Neural Information Processing Systems*. 2503–2511.
- [46] Carolyn B. Seaman. 1999. Qualitative methods in empirical studies of software engineering. *IEEE Transactions on software engineering* 25, 4 (1999), 557–572.
- [47] Qingchao Shen, Haoyang Ma, Junjie Chen, Yongqiang Tian, Shing-Chi Cheung, and Xiang Chen. 2021. A comprehensive study of deep learning compiler bugs. In *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 968–980.
- [48] Sarah Spall, Neil Mitchell, and Sam Tobin-Hochstadt. 2020. Build Scripts with Perfect Dependencies. *Proc. ACM Program. Lang.* 4, OOPSLA (2020).
- [49] Xin Tan, Kai Gao, Minghui Zhou, and Li Zhang. 2022. An Exploratory Study of Deep Learning Supply Chain. In *Proceedings of the IEEE/ACM 44th International Conference on Software Engineering*.
- [50] 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 *Proceedings of the IEEE/ACM 43rd International Conference on Software Engineering*. 238–250.
- [51] Sahil Verma and Zhendong Su. 2020. ShapeFlow: Dynamic Shape Interpreter for TensorFlow. *CoRR* abs/2011.13452 (2020).
- [52] Chengcheng Wan, Shicheng Liu, Henry Hoffmann, Michael Maire, and Shan Lu. 2021. Are Machine Learning Cloud APIs Used Correctly?. In *Proceedings of the IEEE/ACM 43rd International Conference on Software Engineering*. 125–137.
- [53] Ying Wang, Liang Qiao, Chang Xu, Yepang Liu, Shing-Chi Cheung, Na Meng, Hai Yu, and Zhiliang Zhu. 2021. HERO: On the Chaos When PATH Meets Modules. In *Proceedings of the IEEE/ACM 43rd International Conference on Software Engineering*. 99–111.
- [54] Ying Wang, Ming Wen, Yepang Liu, Yibo Wang, Zhenming Li, Chao Wang, Hai Yu, Shing-Chi Cheung, Chang Xu, and Zhiliang Zhu. 2020. Watchman: Monitoring dependency conflicts for python library ecosystem. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 125–135.
- [55] Ying Wang, Ming Wen, Zhenwei Liu, Rongxin Wu, Rui Wang, Bo Yang, Hai Yu, Zhiliang Zhu, and Shing-Chi Cheung. 2018. Do the dependency conflicts in my project matter?. In *Proceedings of the 26th ACM joint meeting on european software engineering conference and symposium on the foundations of software engineering*. 319–330.
- [56] Ying Wang, Ming Wen, Rongxin Wu, Zhenwei Liu, Shin Hwei Tan, Zhiliang Zhu, Hai Yu, and Shing-Chi Cheung. 2019. Could i have a stack trace to examine the dependency conflict issue?. In *Proceedings of the IEEE/ACM 41st International Conference on Software Engineering*. 572–583.
- [57] Ying Wang, Rongxin Wu, Chao Wang, Ming Wen, Yepang Liu, Shing-Chi Cheung, Hai Yu, Chang Xu, and Zhi-liang Zhu. 2021. Will Dependency Conflicts Affect My Program’s Semantics. *IEEE Transactions on Software Engineering* (2021).
- [58] Mohammad Wardat, Wei Le, and Hriday Rajan. 2021. DeepLocalize: Fault Localization for Deep Neural Networks. In *Proceedings of the IEEE/ACM 43rd International Conference on Software Engineering*. 251–262.
- [59] Dangwei Wu, Beijun Shen, Yuting Chen, He Jiang, and Lei Qiao. 2021. Tensfa: Detecting and Repairing Tensor Shape Faults in Deep Learning Systems. In *Proceedings of the IEEE 32nd International Symposium on Software Reliability Engineering*. 11–21.
- [60] Ming Yan, Junjie Chen, Xiangyu Zhang, Lin Tan, Gan Wang, and Zan Wang. 2021. Exposing numerical bugs in deep learning via gradient back-propagation. In *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 627–638.
- [61] Ru Zhang, Wencong Xiao, Hongyu Zhang, Yu Liu, Haoxiang Lin, and Mao Yang. 2020. An empirical study on program failures of deep learning jobs. In *Proceedings of the IEEE/ACM 42nd International Conference on Software Engineering*. 1159–1170.
- [62] Tianyi Zhang, Cuiyun Gao, Lei Ma, Michael Lyu, and Miryung Kim. 2019. An empirical study of common challenges in developing deep learning applications. In *Proceedings of the IEEE 30th International Symposium on Software Reliability Engineering*. 104–115.
- [63] Xiaoyu Zhang, Juan Zhai, Shiqing Ma, and Chao Shen. 2021. AUTOTRAINER: An Automatic DNN Training Problem Detection and Repair System. In *Proceedings of the IEEE/ACM 43rd International Conference on Software Engineering*. 359–371.
- [64] Yuhao Zhang, Yifan Chen, Shing-Chi Cheung, Yingfei Xiong, and Lu Zhang. 2018. An empirical study on TensorFlow program bugs. In *Proceedings of the 27th ACM SIGSOFT International Symposium on Software Testing and Analysis*. 129–140.
- [65] Yuhao Zhang, Luyao Ren, Liqian Chen, Yingfei Xiong, Shing-Chi Cheung, and Tao Xie. 2020. Detecting numerical bugs in neural network architectures. In *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 826–837.
- [66] 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 *Proceedings of the IEEE/ACM 43rd International Conference on Software Engineering: Software Engineering in Practice*. 238–247.