Towards Safe Automated Refactoring of Imperative Deep Learning Programs to Graph Execution

Raffi Khatchadourian*[†], Tatiana Castro Vélez[†], Mehdi Bagherzadeh[‡], Nan Jia[†], Anita Raja*[†]
*City University of New York (CUNY) Hunter College, [†]CUNY Graduate Center, [‡]Oakland University
Email: raffi.khatchadourian@hunter.cuny.edu, tcastrovelez@gradcenter.cuny.edu, mbagherzadeh@oakland.edu,
njia@gradcenter.cuny.edu, 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-supporting symbolic, graph-based Deep Neural Network (DNN) computation. While scalable, such development is error-prone, non-intuitive, and difficult to debug. Consequently, more natural, imperative DL frameworks encouraging eager execution have emerged at the expense of run-time performance. Though hybrid approaches aim for the "best of both worlds," using them effectively requires subtle considerations to make code amenable to safe, accurate, and efficient graph execution. We present our ongoing work on automated refactoring that assists developers in specifying whether and how their otherwise eagerly-executed imperative DL code could be reliably and efficiently executed as graphs while preserving semantics. The approach, based on a novel imperative tensor analysis, will automatically determine when it is safe and potentially advantageous to migrate imperative DL code to graph execution and modify decorator parameters or eagerly executing code already running as graphs. The approach is being implemented as a PvDev Eclipse IDE plug-in and uses the WALA Ariadne analysis framework. We discuss our ongoing work towards optimizing imperative DL code to its full potential.

Index Terms—deep learning, refactoring, graph execution

I. INTRODUCTION

Machine Learning (ML), including Deep Learning (DL), systems are pervasive. They use dynamic models, whose behavior is ultimately defined by input data. However, as datasets grow, efficiency becomes essential [1]. DL frameworks have traditionally embraced a *deferred* execution-style that supports symbolic, graph-based Deep Neural Network (DNN) computation [2], [3]. While scalable, development is error-prone, cumbersome, and produces programs that are difficult to debug [4]-[7]. Contrarily, more natural, less errorprone, and easier-to-debug *imperative* DL frameworks [8]– [10] encouraging eager execution have emerged. Though ubiquitous, such programs are less efficient and scalable as their deferred-execution counterparts [3], [9], [11]-[14]. Thus, hybrid approaches [11], [12], [15] execute imperative DL programs as static graphs at run-time. For example, in TensorFlow [16], AutoGraph [11] can enhance performance by decorating (annotating)—with optional yet influential decorator arguments—appropriate Python function(s) with @tf.

This material is based upon work supported by the National Science Foundation under Award Nos. CCF-22-00343 and CNS-22-13763.

function. Decorating functions with such hybridization APIs can increase code performance without explicit modification.

Though promising, hybridization necessitates non-trivial metadata [13] and exhibits limitations and known issues [17] with native program constructs. Subtle considerations are required to make code amenable to safe, accurate, and efficient graph execution [18], [19]. Alternative approaches [13] impose custom Python interpreters, which may be impractical for industry, and support only specific Python constructs. Thus, developers are burdened with making their code compatible with the underlying execution model conversion and *manually* specifying the functions to be converted. Manual analysis and refactoring (semantics-preserving, source-to-source transformation) can be overwhelming, error- and omission-prone [20], and complicated by Object-Orientation (OO) (e.g., *Keras* [10]) and dynamically-typed languages (e.g., Python).

We present our ongoing work on a fully automated, semantics-preserving refactoring approach that transforms otherwise eagerly-executed imperative (Python) DL code for enhanced performance by specifying whether and how such code could be reliably and efficiently executed as graphs at run-time. The approach—based on a novel tensor analysis specifically for imperative DL code-will infer when it is safe and potentially advantageous to migrate imperative DL code to graph execution and modify decorator parameters or eagerly executing code already running as graphs. It will also discover possible side-effects in Python functions to safely transform imperative DL code to either execute eagerly or as a graph at run-time. While LLMs [21] and big datadriven refactorings [22] have emerged, obtaining a (correct) dataset large enough to automatically extract the proposed refactorings is challenging as developers struggle with (manually) migrating DL code to graph execution [18]. Also, while developers generally underuse automated refactorings [23], [24], since data scientists and engineers may not be classically trained software engineers, they may be more open to using automated (refactoring) tools. Furthermore, our approach will be fully automated with minimal barrier to entry. Our refactoring approach is being implemented as an open-source PyDev Eclipse Integrated Development Environment (IDE) plug-in [25] that integrates analyses from the WALA Ariadne analysis framework [26]. Moreover, while the refactorings will operate on imperative DL code that is easier-to-debug than its deferred-execution counterparts, the refactorings themselves

```
1 import tensorflow as tf
2 class SequentialModel (Model) :
                                       2 class SequentialModel (Model):
                                       3 def init (self, **kwargs):
   def init (self, **kwargs):
    super (SequentialModel, self)
                                          super(SequentialModel, self)
       __init__(...)
f.flatten = layers.Flatten(
                                            __init_
                                             <u>__init__</u>(...)
lf.flatten = layers.Flatten(
     input_shape=(28, 28))
                                           input_shape=(28, 28))
    num lavers = 100 # Add
                                          num lavers = 100 # Add lavers.
    self.layers = [layers
                                           self.layers = [layers
      .Dense(64,activation="relu")
                                            .Dense(64,activation="relu")
11
      for n in range (num lavers) ]
                                      11
                                             for n in range(num lavers)]
    self.dropout = Dropout(0.2)
                                           self.dropout = Dropout(0.2)
13
    self.dense_2 = layers.Dense(10)13
                                           self.dense_2 = layers.Dense(10)
14
                                         @tf.function
                                             __call__ (self, x):
= self.flatten(x)
16
   def __call__(self, x):
                                      16
17
         self.flatten(x)
                                      17
    for layer in self.layers:
                                          for layer in self.layers:
19
     x = layer(x)
                                      19
                                           x = layer(x)
    x = self.dropout(x)
                                      20
                                          x = self.dropout(x)
20
    x = self.dense_2(x)
                                          x = self.dense_2(x)
                                      21
    return x
                                      22
                                          return ×
```

(a) Code snippet before refactoring.(b) Improved code via refactoring.Listing 1: TensorFlow imperative (OO) DL model code [14].

will not improve debuggability but instead enable developers to have *performant* easily-debuggable (imperative) DL code.

II. MOTIVATING EXAMPLES

Lst. 1a portrays TensorFlow imperative (OO) DL code representing a modestly-sized model for classifying images. By default, this code runs eagerly; however, it may be possible to enhance performance by executing it as a graph at run-time. Lst. 1b, lines 1 and 15 display the refactoring with the imperative DL code executed as a graph at runtime (added code is underlined). AutoGraph [11] is now used to potentially improve performance by decorating—with optional yet influential decorator arguments—call() with Otf.function. At run-time, call()'s execution will be "traced" and an equivalent graph will be generated [17]. In this case, a speedup $(runtime_{old}/runtime_{new})$ of ~ 9.22 ensues [27]. Though promising, using hybridization reliably and efficiently is challenging [13], [17]. For instance, side-effect producing, native Python statements are problematic for tf.functiondecorated functions [17]. Because their executions are traced, a function's behavior is "etched" (frozen) into its corresponding graph and thus can have unexpected results.

III. OPTIMIZATION APPROACH

We work towards two new refactorings, namely, CON-VERT EAGER FUNCTION TO HYBRID and OPTIMIZE HYBRID FUNCTION. The former transforms otherwise eagerly-executed imperative (Python) DL code for enhanced performance, automatically specifying whether and how such code could be reliably and efficiently executed as graphs at run-time. It infers when it is safe and potentially advantageous to migrate imperative DL code to graph execution. The latter either modifies existing decorator parameters or the structure of imperative DL code already running as graphs. While the DL code portrayed in lst. 1b is sequentially executed, hybrid functions share some commonality with concurrent programs. For example, to avoid unexpected behavior, such functions should avoid side-effects. In our refactoring formulation, we will approximate aspects like side-effects in deciding which transformations to perform to ensure that they are safe, i.e., that the original program semantics are preserved. To ensure that the transformations are advantageous, we will involve (imperative) tensor analysis to avoid function "retracing" so that newly hybridized functions have tensor parameters whose shapes are sufficiently general. Otherwise, the transformed function would be traced *every* time it called, potentially *degrading* performance [28]. Furthermore, DL code interacts with many third-party libraries [5], [6], [29]–[31]. Our approach will operate on a *closed-world* assumption that assumes access to all source code that could possibly affect or be affected by the refactorings. We will then relax this assumption in our implementation by conservatively failing refactoring preconditions for functions defined elsewhere.

Challenges include a lack of static type information, which is necessary to determine candidate functions (must have at least one parameter of type Tensor). Our current approach is to use Python 3 type hints if present. We are also augmenting Ariadne [26] to analyze imperative Python code (TensorFlow 2). Also, unlike, e.g., Java, Python has no restrictions on decorator (annotation) arguments. Thus, we utilize Ariadne for dataflow analysis to determine configuration values. Furthermore, tf. function may be used as a first-class function instead of a decorator. To this end, we are working towards building a fluent API typestate analysis for imperative DL code by adapting the work of Khatchadourian et al. [32]. Existing work for determining tensor shapes only works for *procedural* TensorFlow (TF v1) code. Even finding a fully qualified name (FQN) of a program entity (e.g., tf.function) statically is difficult in Python. Import statements can appear anywhere in the code. Moreover, there is also import aliasing (e.g., import tensorflow as tf) to handle.

Complex static analysis can be expensive and not scalable. However, such analyses may be useful for future approaches by the community. Faster *speculative* analysis [1] uses contextual (ML) keywords to make assumptions. Here, assumptions are explicitly presented to developers. Developers then decide if assumptions are valid and may *reject* the refactoring. However, it involves more developer input, and DL frameworks evolve constantly, potentially changing "keywords." We are currently considering devising a hybrid analysis [13] that runs DL code for several epochs to collect type information. Such an approach can be fast but relies on particular datasets and thus may be less generalizable.

IV. CONCLUSION & FUTURE WORK

Imperative DL code is easier to debug, write, and maintain than traditional DL code that runs in a deferred execution. However, it comes at the expense of (run-time) performance. Hybrid approaches bridge the gap between eager and graph execution. Using hybrid techniques to achieve optimal performance and semantics preservation is difficult. Our in progress work aims to automate client-side analyses and transformations to *use* hybridization APIs correctly and optimally. In the future, we plan to evaluate our approach by curating a DL Python project dataset from Castro Vélez *et al.* [18] that is ready to be analyzed, with dependency and build infrastructure. We will also open-source our refactoring tools.

REFERENCES

- [1] W. Zhou, Y. Zhao, G. Zhang, and X. Shen, "HARP: Holistic analysis for refactoring Python-based analytics programs," in *International Conference on Software Engineering*, 2020, pp. 506–517. DOI: 10.1145/3377811.3380434.
- [2] Google LLC. "Migrate your TensorFlow 1 code to TensorFlow 2, Automatic conversion script," TensorFlow Core. (May 27, 2021), [Online]. Available: https://tensorflow.org/guide/migrate#automatic_conversion_script (visited on 05/27/2021).
- [3] T. Chen, M. Li, Y. Li, M. Lin, N. Wang, M. Wang, T. Xiao, B. Xu, C. Zhang, and Z. Zhang, "MXNet: A flexible and efficient Machine Learning library for heterogeneous distributed systems," in Workshop on Machine Learning Systems at NIPS, 2015. arXiv: 1512.01274 [CS.DC].
- [4] Y. Zhang, Y. Chen, S.-C. Cheung, Y. Xiong, and L. Zhang, "An empirical study on TensorFlow program bugs," in *International Symposium on Software Testing and Analysis*, 2018. DOI: 10.1145/3213846.3213866
- [5] M. J. Islam, G. Nguyen, R. Pan, and H. Rajan, "A comprehensive study on Deep Learning bug characteristics," in *Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, Aug. 2019. DOI: 10.1145/3338906.3338955.
- [6] M. J. Islam, H. A. Nguyen, R. Pan, and H. Rajan, What do developers ask about ML libraries? a large-scale study using Stack Overflow, 2019. arXiv: 1906.11940 [cs.SE].
- [7] T. Zhang, C. Gao, L. Ma, M. Lyu, and M. Kim, "An empirical study of common challenges in developing Deep Learning applications," in *International Symposium on Software Reliability Engineering*, Oct. 2019. DOI: 10.1109/ISSRE.2019.00020.
- [8] A. Agrawal, A. N. Modi, A. Passos, A. Lavoie, A. Agarwal, A. Shankar, I. Ganichev, J. Levenberg, M. Hong, R. Monga, and S. Cai, TensorFlow Eager: A multi-stage, Python-embedded DSL for Machine Learning, 2019. arXiv: 1903.01855 [cs.PL].
- [9] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Köpf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, *PyTorch: An imperative style, high-performance Deep Learning library*, Dec. 3, 2019. arXiv: 1912.01703 [CS.LG].
- [10] F. Chollet, Deep Learning with Python, 2nd ed. Manning, 2020.
- [11] D. Moldovan, J. M. Decker, F. Wang, A. A. Johnson, B. K. Lee, Z. Nado, D. Sculley, T. Rompf, and A. B. Wiltschko, *AutoGraph: Imperative-style coding with graph-based performance*, 2019. arXiv: 1810.08061 [cs.PL].
- [12] Facebook Inc. "PyTorch documentation, TorchScript." en. (2019), [Online]. Available: https://pytorch.org/docs/stable/jit.html (visited on 02/19/2021).
- [13] E. Jeong, S. Cho, G.-I. Yu, J. S. Jeong, D.-J. Shin, T. Kim, and B.-G. Chun, "Speculative symbolic graph execution of imperative Deep Learning programs," SIGOPS Oper. Syst. Rev., vol. 53, no. 1, pp. 26–33, Jul. 2019, ISSN: 0163-5980. DOI: 10.1145/3352020. 3352025.
- [14] Google LLC. "Introduction to graphs and tf.function." (Jan. 19, 2022), [Online]. Available: https://tensorflow.org/guide/intro_to_graphs (visited on 01/20/2022).
- [15] Apache. "Hybridize, Apache MXNet documentation." (Apr. 8, 2021), [Online]. Available: https://mxnet.apache.org/versions/1.8.0/api/ python/docs/tutorials/packages/gluon/blocks/hybridize.html (visited on 04/08/2021).
- [16] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, M. Kudlur, J. Levenberg, R. Monga, S. Moore, D. G. Murray, B. Steiner, P. Tucker, V. Vasudevan, P. Warden, M. Wicke, Y. Yu, and X. Zheng, "TensorFlow: A system for large-scale Machine Learning," in Symposium on Operating Systems Design and Implementation, 2016.
- [17] Google LLC. "Better performance with tf.function." (Feb. 4, 2021), [Online]. Available: https://tensorflow.org/guide/function (visited on 02/19/2021).
- [18] T. Castro Vélez, R. Khatchadourian, M. Bagherzadeh, and A. Raja, "Challenges in migrating imperative Deep Learning programs to graph execution: An empirical study," in *International Conference on Mining Software Repositories*, ser. MSR '22, ACM/IEEE, ACM, May 2022. DOI: 10.1145/3524842.3528455. arXiv: 2201.09953 [cs.SE].

- [19] J. Cao, B. Chen, C. Sun, L. Hu, and X. Peng, Characterizing performance bugs in Deep Learning systems, Dec. 3, 2021. arXiv: 2112.01771 [CS.SE].
- [20] D. Dig, J. Marrero, and M. D. Ernst, "Refactoring sequential Java code for concurrency via concurrent libraries," in *International Conference* on Software Engineering, IEEE, 2009, pp. 397–407. DOI: 10.1109/ ICSE.2009.5070539.
- [21] OpenAI, Inc. "ChatGPT." (Aug. 18, 2023), [Online]. Available: https://chat.openai.com (visited on 08/18/2023).
- [22] M. Dilhara, A. Ketkar, N. Sannidhi, and D. Dig, "Discovering repetitive code changes in Python ML systems," in *International Conference on Software Engineering*, ser. ICSE '22, To appear., 2022.
- [23] S. Negara, N. Chen, M. Vakilian, R. E. Johnson, and D. Dig, "A comparative study of manual and automated refactorings," in European Conference on Object-Oriented Programming, G. Castagna, Ed., Berlin, Heidelberg: Springer Berlin Heidelberg, 2013, pp. 552–576, ISBN: 978-3-642-39038-8.
- [24] M. Kim, T. Zimmermann, and N. Nagappan, "A field study of refactoring challenges and benefits," in *Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, ser. FSE '12, Cary, North Carolina: ACM, Nov. 2012, ISBN: 9781450316149. DOI: 10.1145/2393596.2393655.
- [25] F. Zadrozny. "Pydev." (Apr. 15, 2023), [Online]. Available: https://www.pydev.org (visited on 05/31/2023).
- [26] J. Dolby, A. Shinnar, A. Allain, and J. Reinen, "Ariadne: Analysis for Machine Learning programs," in *International Workshop on Machine Learning and Programming Languages*, ser. MAPL 2018, ACM SIGPLAN, Philadelphia, PA, USA: Association for Computing Machinery, 2018, pp. 1–10, ISBN: 9781450358347. DOI: 10.1145/3211346.3211349.
- [27] R. Khatchadourian. "graph_execution_time_comparison.ipynb." (Feb. 23, 2021), [Online]. Available: https://bit.ly/3bwrhVt (visited on 11/03/2021).
- [28] Google LLC. "Better performance with tf.function, Controlling retracing," TensorFlow Core, TensorFlow. (Nov. 11, 2021), [Online]. Available: https://www.tensorflow.org/guide/function#controlling_retracing (visited on 01/10/2022).
- [29] M. Dilhara, A. Ketkar, and D. Dig, "Understanding software-2.0: A study of Machine Learning library usage and evolution," ACM Transactions on Software Engineering and Methodology, 2021. DOI: 10.1145/3453478.
- [30] M. J. Islam, R. Pan, G. Nguyen, and H. Rajan, "Repairing Deep Neural Networks: Fix patterns and challenges," in *International Conference on Software Engineering*, 2020. DOI: 10.1145/3377811.3380378.
- [31] Z. Zhang, Y. Yang, X. Xia, D. Lo, X. Ren, and J. Grundy, "Unveiling the mystery of API evolution in Deep Learning frameworks: A case study of TensorFlow 2," in *International Conference on Software Engineering*, ser. ICSE-SEIP, 2021. DOI: 10.1109/ICSE-SEIP52600. 2021.00033.
- [32] R. Khatchadourian, Y. Tang, M. Bagherzadeh, and S. Ahmed, "Safe automated refactoring for intelligent parallelization of Java 8 streams," in *International Conference on Software Engineering*, ser. ICSE '19, IEEE Press, 2019, pp. 619–630. DOI: 10.1109/ICSE.2019.00072.