

A qualitative study of neural language models of source code

Anonymous Author(s)

ABSTRACT

An increasing number of developers rely upon pretrained language models to help them write source code. But how reliable are they and what are the consequences of overreliance? Properly vetted, the risk of errors should be minimal, however complacency can lead to increased reviewer burden and introduce bugs into version control, a risk which users must remain vigilant towards when accepting their suggestions. Companies which host these models are unlikely to give users access to the pretrained models for evaluation. Nevertheless, we would like some assurance the model's suggestions are robust to naturally-occurring variance. In this work, we provide a systematic approach and empirical results for evaluating neural code completion models using automated source code transformations.

ACM Reference Format:

Anonymous Author(s). 2022. A qualitative study of neural language models of source code. In *Proceedings of The 44th International Conference on Software Engineering (ICSE 2022)*. ACM, New York, NY, USA, 4 pages. <https://doi.org/10.1145/1122445.1122456>

1 INTRODUCTION

Statistical language models play an increasingly synergetic role in software development, most prominently in code completion [1]. Yet from a user's perspective, the behavior of these models is opaque: partially completed source code written in an editor is sent to a remote server, which returns a real-time suggestion. This client-server architecture can be seen as a black-box or *extensional* function from a mathematical perspective. How can we evaluate the behavior of neural language models in this setting? The role of automated testing becomes apparent.

First conceived in the software testing literature, metamorphic testing [2] is a concept known in machine learning as *self-supervision*. In settings where labels are scarce but certain *invariants* or *metamorphic relations* are known, one can generate an effectively unlimited amount of data by selecting and recombining synthetic transformations. For example, computer vision models should be invariant to shift, scale and rotation, so given a small dataset of labeled images, we can apply these transformations to obtain far more data for training or evaluation. Do similar invariants exist for source code?

One would expect code completion models to exhibit invariance to certain refactorings. For example, code completion should be invariant to synonym variable renaming or reordering dataflow-independent statements.

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

ICSE 2022, May 21–29, 2022, Pittsburgh, PA, USA

© 2022 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00

<https://doi.org/10.1145/1122445.1122456>

Code transformations...

How to evaluate source code transformations:

Comparing different types of transformations

Combining different transformations or bounds?

Single vs. multiple transformation benefits (cite testing literature?)

Code characteristics: different categories of code based on... selection: slicing, tokenization.. language: java, kotlin size: method length complexity: halstead, cyclometric

2 METHOD

Do language models learn to recognize patterns in source code? If so, what kinds? A longstanding research area in software engineering tries to extract idioms from large software repositories. Likewise, a longstanding goal of neural program synthesis is compositional generalization. Both currently require designing a feature representation that “selects” for a specific kind of invariance — for example, structural similarity selects for name-invariance, and semantic similarity captures structural invariance.

In our work, we explore the extent to which recent language models trained on source code learn regularities they were not explicitly designed to represent. We analyze three different language models, and compare their zero-shot generalization on synthetically-generated patterns. We generate code snippet pairs corresponding to manually-designed variance, and measure the effect on various language models. Each model is evaluated on a set of code snippets exposed to synthetic transformations and predicts a masked token.

Our goal then, is to identify equivariant representations captured by the language model (LM). To do so, we generate program transformations (PTs) and measure vector equivariance in four rewriting categories:

- (1) Syntactic - can the LM detect syntactically invalid PTs? (e.g. syntax corruption, imbalanced parenthesis)
- (2) Structural - can the LM detect syntactically valid, but structurally altered PTs? (e.g. use before declaration, permuted argument order)
- (3) Semantic - can the LM detect structurally valid, but semantically altered PTs? (e.g. constant modification, operator substitution and order of operations)
- (4) Equivalence - can the LM detect semantically valid but rewritten PTs? (e.g. semantically valid rewrites)

To generate these SCTs, we implement a source code transformation tool which generates synthetic variants of naturally-occurring source code snippets from each of these three categories and evaluates the model for completion accuracy. We have implemented the following source code transformations (SCTs):

- (1) **Synonym renaming:** rename variables with synonyms
- (2) **Dead code:** introduce dead code to source code
- (3) **Statement reordering:** swap non-interfering statement order
- (4) **Loop bounds fuzzing:** change loop bounds conditions

- (5) **Permute argument order**: scramble method arguments
- (6) **Literal mutation**: mutate contents of primitive types

3 EXPERIMENTS

In this work, we attempt to understand the relationship between entities in a software project. Our research seeks to answer the following questions:

- (1) Which contexts in a software project share mutual information?
- (2) To what degree can we claim the model has learned to:
 - (a) Locate contextually relevant artifacts within a software project?

Model	IRA	QDL	P/R
CODEGPT [7]	X	X	X
GRAPHCODEBERT [4]	X	X	X
CODEBERT-SMALL [3]	X	X	X
ROBERTA-JAVA [6]	X	X	X
COPILLOT[1]	X	X	X

Table 1: Experiments table for comparing pretrained LM embeddings on source code snippets. IRA: Iterater agreement. QDL: query description length, P/R: Precision/recall

For each of these models (available on HuggingFace), we sample a set of code snippets from our synthetic dataset, and compare how well they agree on each task.

In contrast with classical code completion models which only require a file-local context, our method is designed to navigate an entire project. In the following experiments, we compare completion accuracy with a vanilla sequence prediction model, as well as an AST-structured sequence prediction model trained from a corpus of Java projects on the same task.

We hypothesize that by jointly learning to choose locations in a project over which to attend while solving a downstream task, such as masked sequence completion, our model will produce a feature representation capable of locating and extracting information from semantically relevant contexts. We evaluate our hypothesis both qualitatively and quantitatively.

In our first set of experiments, we try to understand what is shared between sequences mapped to similar latent vectors. Do similar sequences share salient keywords? Are those keywords relevant to the task?

In our second experiment, we try to measure the information gain from including and excluding various filetypes through ablation. For graphs containing filetypes which include Markdown or Java, what kinds of information do these resources provide and which categories are most salient?

In our third and final set of experiments, we compare performance across hyperparameters. Does contextual expansion lead to better task performance for a given sequence prediction task? By relaxing edge-construction criteria and increasing hyperparameters such as beam search budget, we would expect corresponding task performance to increase.

If our hypothesis is correct, the virtual knowledge graph will span both natural language and source code artifacts. If so, this

For each SCT, we mask various tokens in the code snippet and compare the model’s ability to fill in the correct token before and after the transformation is applied. Our goal is to measure how sensitive the pretrained model is to each type of SCT.

would provide evidence to support the broader hypothesis [5] that documentation is a useful source of information. In addition to being useful for the prediction task itself, we anticipate our model could also be used for knowledge graph extraction and suggest semantically relevant code snippets to developers.

4 EVALUATION

REFERENCES

- [1] Mark Chen et al. Evaluating large language models trained on code, 2021.
- [2] TY Chen, J Feng, and TH Tse. Metamorphic testing of programs on partial differential. *Information and Computation*, 121(1):93–102, 1995.
- [3] Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, et al. CodeBERT: A pre-trained model for programming and natural languages. *arXiv preprint arXiv:2002.08155*, 2020.
- [4] Daya Guo, Shuo Ren, Shuai Lu, Zhangyin Feng, Duyu Tang, Shujie Liu, Long Zhou, Nan Duan, Alexey Svyatkovskiy, Shengyu Fu, Michele Tufano, Shao Kun Deng, Colin Clement, Dawn Drain, Neel Sundaresan, Jian Yin, Daxin Jiang, and Ming Zhou. Graphcodebert: Pre-training code representations with data flow, 2021.
- [5] Jin Guo, Jinghui Cheng, and Jane Cleland-Huang. Semantically enhanced software traceability using deep learning techniques. In *2017 IEEE/ACM 39th International Conference on Software Engineering (ICSE)*, pages 3–14. IEEE, 2017.
- [6] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach, 2019.
- [7] Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin Clement, Dawn Drain, Daxin Jiang, Duyu Tang, Ge Li, Lidong Zhou, Linjun Shou, Long Zhou, Michele Tufano, Ming Gong, Ming Zhou, Nan Duan, Neel Sundaresan, Shao Kun Deng, Shengyu Fu, and Shujie Liu. Codexglue: A machine learning benchmark dataset for code understanding and generation, 2021.

A ANALYSIS OF INTERNAL STRUCTURE

Our preliminary results compare distance metrics (Fig. 1), explore embedding quality (Fig. 2) and visualize the synthetic knowledge graphs (Fig. 3).

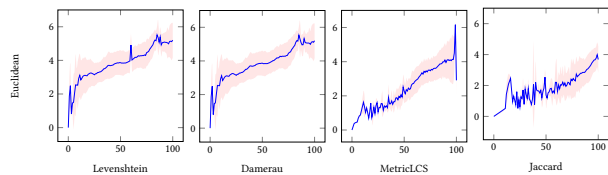


Figure 1: CodeBERT latent space distance correlates with string edit distance.

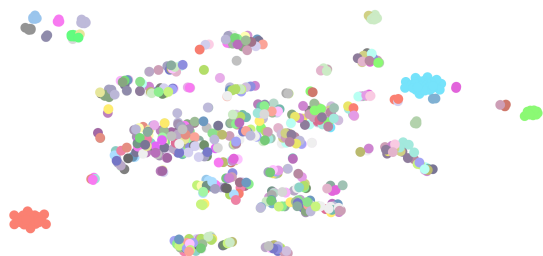


Figure 2: Reduced dimensionality TNSE embeddings colored by line length.

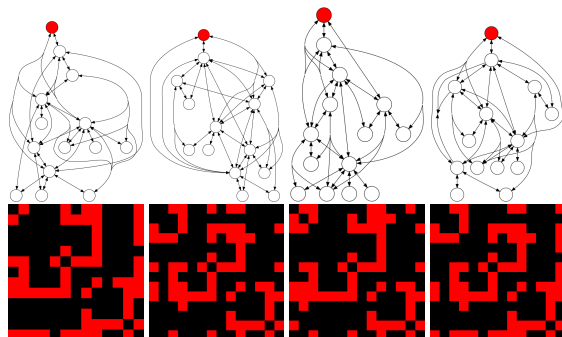


Figure 3: Virtual knowledge graph constructed by visiting nearest-neighbors in latent space.