

DOI: 10.15514/ISPRAS-2024-36(4)-6



GraphTyper: Neural Types Inference from Code Represented as Graph

G.A. Arutyunov, ORCID: 00000-0003-4537-4332 <gaarutyunov@edu.hse.ru>

S.M. Avdoshin, ORCID: 0000-0001-8473-8077 <savdoshin@hse.ru>

HSE University

20, Myasnitskaya st., Moscow, 101000, Russia.

Abstract. Although software development is mostly a creative process, there are many scrutiny tasks. As in other industries, there is a trend for automation of routine work. In many cases, machine learning and neural networks have become a useful assistant in that matter. Programming is not an exception: GitHub has stated that Copilot is already used to write up to 30% of code in the company. Copilot is based on Codex, a Transformer model trained on code as a sequence. However, a sequence is not a perfect representation for programming languages. In this work, we claim and demonstrate that by combining the advantages of Transformers and graph representations of code, it is possible to achieve excellent results even with comparably small models.

Keywords: neural networks; transformers; graphs; abstract syntax tree.

For citation: Arutyunov G.A., Avdoshin S.M. GraphTyper: Neural types inference from code represented as graph. Trudy ISP RAN/Proc. ISP RAS, vol. 36, issue 4, 2024. pp. 69-80. DOI: 10.15514/ISPRAS-2024-36(4)-6.

Acknowledgements. This research was supported in part through computational resources of HPC facilities at HSE University.

GraphType: Вывод типов из графовой репрезентации кода посредством нейронных сетей

Г.А. Арутюнов, ORCID: 0000-0003-4537-4332 <gaarutyunov@edu.hse.ru>

С.М. Авдошин, ORCID: 0000-0001-8473-8077 <savdoshin@hse.ru>

*Национальный исследовательский университет «Высшая школа экономики» (НИУ ВШЭ),
101000, Россия, г. Москва, ул. Мясницкая, д. 20*

Аннотация. Несмотря на то, что программирование – это творческий процесс, достаточно много времени уходит на решение рутинных задач. Как и в других индустриях в сфере информационных технологий стремятся автоматизировать рутинные задачи. Во многих случаях применяются нейронные сети. Программирование не является исключением: Github заверяют, что уже около 30% кода написано при помощи Copilot. Этот инструмент основан на модели Codex – трансформере, обученном на исходном коде программ. Однако представление кода в виде последовательности, как это сделано в Copilot, не так эффективно. В данной работе мы показали, что использование трансформеров и графового представления кода приводит к очень хорошим результатам даже для маленьких моделей.

Ключевые слова: нейронные сети; трансформеры; графы; абстрактное синтаксическое дерево.

Для цитирования: Арутюнов Г.А., Авдошин С.М. GraphType: Вывод типов из графовой репрезентации кода посредством нейронных сетей. Труды ИСП РАН, том 36, вып. 4, 2024 г., стр. 69–80 (на английском языке). DOI: 10.15514/ISPRAS–2024–36(4)–6.

Благодарности. Исследование выполнено с использованием суперкомпьютерного комплекса НИУ ВШЭ.

1. Introduction

Application of Transformers yet again has managed to break the deadlock: this time in the task of code generation [1–4]. Nevertheless, the versatile Transformer architecture has displayed good results on several benchmarks, in the recent work [5] it was shown that increasing the size of the model doesn't result in a better performance. Moreover, it is evident that context matters a lot to produce a working code. However, it is not practical to relentlessly increase the length of context sequence in a Transformer. Therefore, a different approach is needed to boost the efficiency in machine programming tasks [6].

First of all, an expressive code representation has to be selected. Several ways, including token-based, structured and graph-based approaches, have been reviewed [7]. For instance, graph representation using abstract syntax tree (AST), data-flow graph (DFG) and control-flow graph (CFG) yield good results in such tasks as variable misuse detection and correction [8]. Such graph representation can capture an extensive amount of information about the program's code.

Secondly, a versatile model architecture that supports learning on graphs must be used. Multiple models such as RNN [9], LSTM [10] and CNN [11] with flattened graphs have been used. However, graph-aware model architecture is more suitable for the graph representation of code. For this reason, Graph Neural Networks (GNN) are a more reasonable choice of architecture, namely message-passing neural networks [8].

Nonetheless, in this work we aim to make the most of both worlds: the advantages of Transformer architecture and graph representation of code. For instance, we will use Transformer architecture optimizations [12] and graph code representation created from AST and DFG. To make this possible, we will use Pure Transformers [13] instead of models that have some architectural alterations to support graph structure [14–16].

2. Problem Statement

In this work, we test the ability of Pure Transformers to add types to Python source code based on its graph structure. This task was selected as a starting point for future research due to its practical relevance.

Firstly, dynamically typed languages, such as Python and JavaScript, have gained quite some traction during the last years [17]. However, it doesn't mean they're easier [18–20] or less error-prone than statically typed languages [21]. Moreover, lack of type hints in libraries might lead to expensive errors in fields such as Data Science [22].

There are some tools outside the neural networks domain that perform static type checking and inferencing type annotations [23,24]. Nonetheless, these utilities do not work without type hints in the source code of the dependencies, which is pretty common. To alleviate this, there are proposals about Domain-Specific Languages (DSL) for Data Science [22]. However, it wouldn't work on existing code base and massive adoption is not very likely.

On the other hand, absence of type hints is not a restriction for neural networks (see. Section 5.2). In addition, they don't only find erroneous types in existing codebase [25] but can also be used during development to annotate code on the fly [26].

Most importantly, inferring types requires a model to learn a lot about the source code. Therefore, developing a model with versatile architecture to infer types allows it to be later applied for other tasks.

2.1 Metrics

To test the model, we use two metrics from the Typilus paper [25]:

- Exact Match: Predicted and ground truth types match exactly.
- Match up to Parametric Type: Exact match when ignoring all type parameters.

3. Previous Work

3.1 Graph Representation of Code

AST and DFG have already been used with Transformers in the code generation and summarization tasks [27–29]. In addition, some joint graph structure representations that include different code graphs have been developed, namely code property graph (CPG) [30], that incorporates AST, CFG and PDG (program dependency graph). This graph representation has already been used for vulnerability detection [30] and similarity detection [31].

3.2 Graph Transformers

Graph Transformers is a novel architecture that has been developing in the past few years. They have been applied for several tasks, mostly in the field of molecule generation, node classification and node feature regression [13–16]. Apart from models with alterations to Transformer base architecture [15,16] researchers have recently developed simpler models [13] that are compatible with many popular techniques developed for standard Transformers [12].

3.3 Type Inference with Neural Networks

The task of type inference has been also extensively covered in recent research. Many different architectures have been used for this task: GNNs [25], RNNs [26,32] and Transformers [33,34] among others. Moreover, graph representation of code has been used for the task of type inference in dynamically typed programming languages such as Python [25] and Javascript [35].

However, the power of Graph Transformers and Graph Representation of code hasn't been combined yet to solve the task of type inference in source code. This is the gap our model aims to fill. The results of our model compared to previous work [25,26,32–34] are displayed in Table 1.

Table 1. Quantitative evaluation of models measuring their ability to predict ground truth type annotations. EM – exact match, UTPT – Match up to parametric type.

Top-n	Top-1		Top-3		Top-5	
	EM	UTPT	EM	UTPT	EM	UTPT
GraphTyper	34.7	36.42	45.47	55.01	50.69	64.58
TypeBERT	45.4	48.1	51.4	53.5	54.1	56.5
TypeWriter	56.1	58.3	63.7	67.3	65.9	70.4
Typilus	66.1	74.2	71.6	79.8	72.7	80.9
Type4Py	75.8	80.6	78.1	83.8	78.7	84.7
TypeGen	79.2	87.3	85.6	91	87	91.7

4. Proposed Solution

4.1 Dataset

To train and test the model we gathered 600 Python repositories from GitHub containing type annotations from Typilus [25]. We clone these repositories and use pytype [24] for static analysis, augmenting the corpus with inferred type annotations. The top 175 most downloaded libraries are added to the Python environment for type inference. Through deduplication, we remove over 133 thousand code duplicates to prevent bias.

The resulting dataset comprises 118,440 files with 5,997,459 symbols, of which 252,470 have non-Any non-None type annotations. The annotations exhibit diversity with a heavy tailed distribution, where the top 10 types cover half of the dataset, primarily including str, bool, and int. Only 158 types have over 100 annotations, while the majority of types are used fewer than 100 times each, forming 32% of the dataset. This distribution underscores the importance of accurately predicting annotations, especially for less common types. The long-tail of types consists of user-defined and generic types with various type arguments.

The source files are processed to generate graphs that contain AST, DFG, as well as lexical and syntactical information. An example of such a graph is shown on Fig. 1.

In addition to extracting graphs from source code AST, we split them by setting a maximum node and edges number in one graph. For this, we prune the graphs around nodes that have annotations that are later used as targets during training and testing. Finally, we split the data into train-validation-test sets with proportions of 70-10-20, respectively.

4.2 Model Architecture

We base our model architecture on TokenGT [13]. The main advantage of this model is that standard Transformer architecture is not altered to support graph data. It allows us to use some advantages developed specifically for Transformers. For instance, Performer [12] is used to speed up training by using linear time as space complexity.

The main idea of the authors is that combining appropriate token-wise embeddings and self-attention over the node and edge tokens is expressive enough to accurately encode graph structure to make

graph and node-wise predictions. The embeddings in the model are composed of orthonormal node identifiers, namely Laplacian eigenvectors obtained from eigendecomposition of graph Laplacian matrix. In addition, type identifiers are used to encode types of tokens (nodes or edges).

In our model, we use node and edge types extracted from code as token features. Node ground truth annotations are added to the features and randomly masked during training. The overall architecture of the model is displayed at Fig. 2.

Predicting type annotations in graph domain is a node classification task. However, since we are using a Pure Transformer with graphs represented as a sequence of tokens, the task can be reduced to token classification. In the Natural Language Processing (NLP) domain, this is a ubiquitous task, also known as Named Entity Recognition (NER).

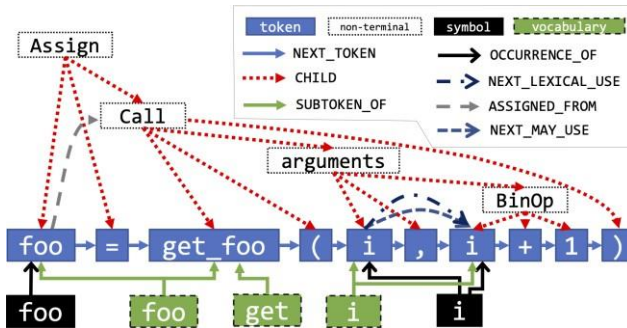


Fig. 1. Sample graph for `foo=get foo(i, i+1)` showing different node and edge types implemented by Allamanis et al. [25].

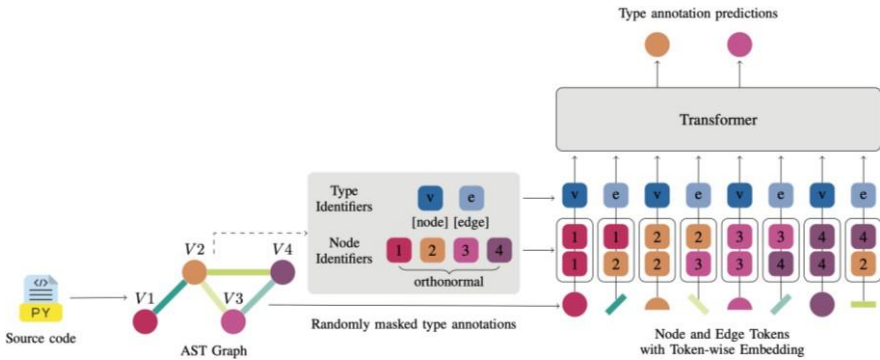


Fig. 2. GraphTyper Architecture. The source code is first transformed into AST graph, then type annotations are randomly masked. The graph is enriched by token type identifiers (node or edge) and orthonormal node identifiers obtained from eigendecomposition of Laplacian matrix. The resulting graph is fed through a Transformer Encoder to obtain type annotations for masked nodes.

Encoder-only architecture has been widely used for the NER task, namely BERT is one of the most popular models [36,37]. We adapt similar architecture by randomly masking type annotations. We then apply an MLP layer to the output of TokenGT [13] to get logits of type annotations.

Masked model architecture is very versatile, and the pre-trained model can be later easily fine-tuned for other tasks, similar to the approaches from the NLP-domain [36]. For example, error [38] and vulnerability [39] data can be added to the code graph to detect and fix them [40–44].

5. Experiment and Ablation Results

To select the final model architecture, we test different models. For our experiments and ablation analysis, we train and test the models using one sample repository. We also limit the number of types

in the vocabulary to one hundred to speed up training and use less resources. To test the models, we calculate Top-n predictions similar to the previous work [26]. Table 2 depicts the results of the experiments and ablation. The model was trained on 1 NVIDIA Tesla V100 32 GB NVLink [45] for 10 epochs with the batch size of 32 graphs.

Table 2. Experiment results of top-n predictions for different model variants.
EM – exact match, UTPT – Match up to parametric type.

Top-n	Top-1		Top-3		Top-5	
Metric	EM	UTPT	EM	UTPT	EM	UTPT
Plane Transformer	10.15	19.46	15.06	29.40	16.81	37.91
+ Node & Type Identifiers	30.88	36.55	40.33	50.37	42.82	56.01
+ Type Annotations	33.36	42.28	41.71	52.90	43.62	57.00
+ Decoder (Autoencoder)	15.90	16.65	28.26	32.81	44.17	56.33
or Longer Context	38.49	39.80	53.14	57.41	58.80	67.38
or More Parameters	29.39	31.82	44.85	49.72	49.74	56.14

5.1 Validating the necessity of node and type identifiers that encode graph structure

First of all, we remove the node and type identifiers introduced by Kim et. al [13]. Our ablation analysis demonstrates that indeed, the graph structure embeddings play a key role in model quality. By removing them from the model, we are left with a simple Transformer that makes predictions only based on AST nodes and edges types without any information about graph structure. Such a model outputs the worst results among all the experiments.

5.2 Using the model without node type annotations

In addition, we try to remove the type annotations from the model completely. This alteration turns our training into a masked NER task. Surprisingly, our model performs well in such conditions. This means that the selected graph representation of code contains a lot of necessary information to infer types.

5.3 Increasing the number of parameters

As we can see, increasing the number of parameters also increases the predictive power of the model. However, increasing the parameters indefinitely is not very practical and requires a lot of computational resources [6]. Moreover, keeping the low number of parameters allows us to use longer context length (more node and edges in graph) during inference with same resource capabilities. Therefore, we don't change the parameter number of the final model, so it remains compact.

5.4 Testing different context length

As for the context length, i.e., maximum number of nodes in graph (512 vs. 1024), our findings are aligned with the conclusions from previous work [6]: longer context increases the performance of the model. However, the AST representation of source code is very bloated and even having a lot of nodes in the graph might not capture enough useful information to make quality predictions. In addition, increasing the context length drastically slows down the training process. Thus, in future research, we will be working on finding a better and more compact graph representation of code.

5.5 Testing different Transformer architectures

Recently, Masked Graph Autoencoders have been applied for the tasks of link prediction and node classification [46], as well as feature reconstruction [47,48]. To validate the robustness of the Encoder-only Model, we also implement a Masked Autoencoder Model. For this, we adapt the approach of Hou et. al [48] for our model. We introduce a learnable mask token and a decoder based on the encoder layers. We reconstruct the type annotations by re-masking the target nodes before feeding them into the decoder. However, we do not observe as good results as with a simple Encoder-only model.

6. Known Limitations

6.1 Size of Type Vocabulary

Since we define our task as node (token) classification, we feed our transformer output into a classifier linear head. Therefore, our type vocabulary is limited. Because of the computational resources' constraints, we limit it to one thousand types.

This issue is addressable by formulating the task as Deep Similarity Learning Problem [49,50]. In this way, the model will output vector representations of types that can be grouped into cluster of similar types. After that, an algorithm such as KNN [51] is used to transform vector representation into a probability of each type [25,26]. Illustration for such an approach is depicted on Fig. 3. The approach follows the methodology described in Typilus paper [25].

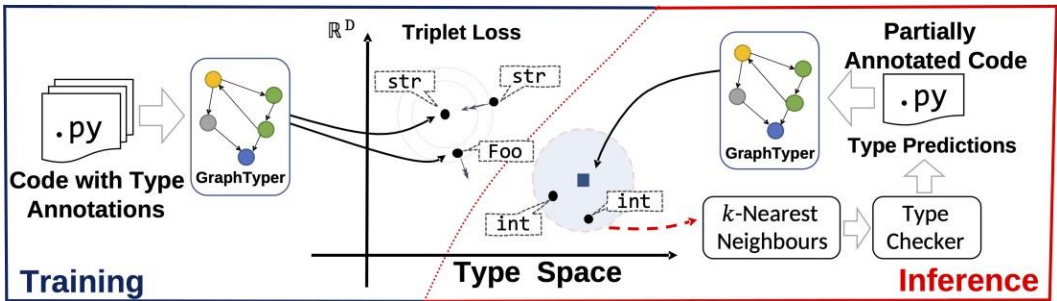


Fig. 3. Solution to the problem using Deep Similarity Learning [25].

6.2 Absence of Natural Language information

In our work, we use only categorical features of nodes and edges of code graph, e.g., AST node types and Python type annotations. Therefore, it would be challenging to apply it directly for tasks such as code generation, because the representation doesn't encode any information about variable names.

There are several approaches that would help address this issue. Firstly, it is possible to use the model output as graph encoding that would be later fed into another model along with tokenized code [52]. This approach could also address the issue from the previous section, since types would be treated as a set of text tokens [34]. Secondly, it is possible to use neural networks to infer variables' names from the context they are used in [53].

7. Future Work

In this work, we explored the application of Graph Transformers for type inference. The versatile architecture of the proposed solution lets us explore other tasks.

7.1 Universal code graph representation

If a universal version of graph code representation is used, similar to CPG [30], we can train the model for multiple programming languages [29]. However, because of the differences of type systems, separate models would be trained for each programming language for better results.

7.2 Detecting duplicates

It is crucial to address the issue of duplicates in source code to train neural networks for code. Several architectures have already been used for such task: Transformers [54], GNNs [55] and RNNs [56]. We believe that the graph representation obtained with our model can be successfully used for code clone detection.

7.3 Code and docstring generation

Firstly, we can train the model using a technique similar to generative pretrained models [57,58] or masked language models [52] to generate code. Secondly, our model can be used to generate code summarization or docstring generation [59,60]. This could only be possible if we adapt some of the approach discussed in the previous section.

7.4 Vulnerability and error detection

Another useful task is to detect errors and generate fixes [61,62]. This is possible by simply adding features that contain error indication or types. Similar approach can be used to scan for vulnerabilities [40,41,44]. Fixing bugs and vulnerabilities, however, would imply that the graph structure could change. Therefore, solving this task would require the model to be modified for graph generation [63].

7.5 Refactoring

Finally, we can extend our model with information about changes to analyze them and propose refactoring possibilities [64]. This goal could be achieved by using the model from the previous section.

8. Conclusion

As for the conclusion, we were able to create a universal model based on TokenGT [13] and code represented as graphs. One of the most important advantages of this model is that it uses the code graph directly. Secondly, the model can be modified to fit other tasks, such as code generation and summarization, docstring generation, refactoring and many more. The code graph can also be extended by different features and node types, since the representation does not differ depending on graph structure. The source code is available at this <https> URL [65].

References / Список литературы

- [1]. M. Chen, J. Tworek, H. Jun, Q. Yuan, H.P. de O. Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman, Evaluating large language models trained on code, ArXiv Prepr. ArXiv210703374 (2021).
- [2]. D. Hendrycks, S. Basart, S. Kadavath, M. Mazeika, A. Arora, E. Guo, C. Burns, S. Puranik, H. He, D. Song, Measuring coding challenge competence with apps, ArXiv Prepr. ArXiv210509938 (2021).
- [3]. Y. Li, D. Choi, J. Chung, N. Kushman, J. Schrittwieser, R. Leblond, J. Keeling, F. Gimeno, A.D. Lago, T. Hubert, P. Choy, C. de, Competition-Level Code Generation with AlphaCode, (n.d.) 74.
- [4]. E. Nijkamp, B. Pang, H. Hayashi, L. Tu, H. Wang, Y. Zhou, S. Savarese, C. Xiong, A Conversational Paradigm for Program Synthesis, ArXiv Prepr. ArXiv220313474 (2022).
- [5]. F.F. Xu, U. Alon, G. Neubig, V.J. Hellendoorn, A Systematic Evaluation of Large Language Models of Code, ArXiv Prepr. ArXiv220213169 (2022).

- [6]. G.A. Arutyunov, S.M. Avdoshin, Big Transformers for Code Generation, Proc. Inst. Syst. Program. RAS 34 (2022) 79–88. [https://doi.org/10.15514/ispras-2022-34\(4\)-6](https://doi.org/10.15514/ispras-2022-34(4)-6).
- [7]. S.M. Avdoshin, G.A. Arutyunov, Code Analysis and Generation Methods Using Neural Networks: an Overview, Inf. Technol. 28 (2022) 378–391. <https://doi.org/10.17587/it.28.378-391>.
- [8]. M. Allamanis, M. Brockschmidt, M. Khademi, Learning to represent programs with graphs, ArXiv Prepr. ArXiv171100740 (2017).
- [9]. M. White, M. Tufano, C. Vendome, D. Poshyvanyk, Deep learning code fragments for code clone detection, in: 2016 31st IEEEACM Int. Conf. Autom. Softw. Eng. ASE, IEEE, 2016: pp. 87–98.
- [10]. H. Wei, M. Li, Supervised Deep Features for Software Functional Clone Detection by Exploiting Lexical and Syntactical Information in Source Code., in: IJCAI, 2017: pp. 3034–3040.
- [11]. L. Mou, G. Li, L. Zhang, T. Wang, Z. Jin, Convolutional neural networks over tree structures for programming language processing, in: Thirtieth AAAI Conf. Artif. Intell., 2016.
- [12]. K. Choromanski, V. Likhoshesterov, D. Dohan, X. Song, A. Kane, T. Sarlos, P. Hawkins, J. Davis, A. Mohiuddin, L. Kaiser, D. Belanger, L. Colwell, A. Weller, Rethinking Attention with Performers, (2020).
- [13]. J. Kim, T.D. Nguyen, S. Min, S. Cho, M. Lee, H. Lee, S. Hong, Pure Transformers are Powerful Graph Learners, (2022). <https://doi.org/10.48550/arXiv.2207.02505>.
- [14]. V.P. Dwivedi, X. Bresson, A Generalization of Transformer Networks to Graphs, (2021). <https://doi.org/10.48550/arXiv.2012.09699>.
- [15]. D. Kreuzer, D. Beaini, W.L. Hamilton, V. Létourneau, P. Tossou, Rethinking Graph Transformers with Spectral Attention, (2021). <https://doi.org/10.48550/arXiv.2106.03893>.
- [16]. C. Ying, T. Cai, S. Luo, S. Zheng, G. Ke, D. He, Y. Shen, T.-Y. Liu, Do Transformers Really Perform Bad for Graph Representation?, (2021). <https://doi.org/10.48550/arXiv.2106.05234>.
- [17]. P.M. Julia Elliott, 2021 Kaggle Machine Learning & Data Science Survey, (2021). <https://kaggle.com/competitions/kaggle-survey-2021>.
- [18]. M.P. Robillard, What Makes APIs Hard to Learn? Answers from Developers, IEEE Softw. 26 (2009) 27–34. <https://doi.org/10.1109/MS.2009.193>.
- [19]. M.P. Robillard, R. Deline, A field study of API learning obstacles, Empir. Softw Engg 16 (2011) 703–732. <https://doi.org/10.1007/s10664-010-9150-8>.
- [20]. M.F. Zibran, F.Z. Eishita, C.K. Roy, Useful, But Usable? Factors Affecting the Usability of APIs, in: 2011 18th Work. Conf. Reverse Eng., 2011: pp. 151–155. <https://doi.org/10.1109/WCRE.2011.26>.
- [21]. N. Alzahrani, F. Vahid, A. Edgcomb, K. Nguyen, R. Lysecky, Python Versus C++: An Analysis of Student Struggle on Small Coding Exercises in Introductory Programming Courses, in: Proc. 49th ACM Tech. Symp. Comput. Sci. Educ., Association for Computing Machinery, New York, NY, USA, 2018: pp. 86–91. <https://doi.org/10.1145/3159450.3160586>.
- [22]. L. Reimann, G. Kniesel-Wünsche, Safe-DS: A Domain Specific Language to Make Data Science Safe, (2023).
- [23]. Pyre: A performant type-checker for Python 3, (n.d.). <https://pyre-check.org> (accessed May 12, 2024).
- [24]. PyType: A static type analyzer for Python code, (n.d.). <https://github.com/google/pytype> (accessed May 12, 2024).
- [25]. M. Allamanis, E.T. Barr, S. Ducousso, Z. Gao, Typilus: Neural Type Hints, in: PLDI, 2020.
- [26]. A.M. Mir, E. Latoskinas, S. Proksch, G. Gousios, Type4py: Deep similarity learning-based type inference for python, ArXiv Prepr. ArXiv210104470 (2021).
- [27]. Z. Sun, Q. Zhu, Y. Xiong, Y. Sun, L. Mou, L. Zhang, Treegen: A tree-based transformer architecture for code generation, in: Proc. AAAI Conf. Artif. Intell., 2020: pp. 8984–8991.
- [28]. Z. Tang, C. Li, J. Ge, X. Shen, Z. Zhu, B. Luo, AST-Transformer: Encoding Abstract Syntax Trees Efficiently for Code Summarization, (2021). <https://doi.org/10.48550/arXiv.2112.01184>.
- [29]. K. Wang, M. Yan, H. Zhang, H. Hu, Unified Abstract Syntax Tree Representation Learning for Cross-Language Program Classification, in: Proc. 30th IEEEACM Int. Conf. Program Comprehension, 2022: pp. 390–400. <https://doi.org/10.1145/3524610.3527915>.
- [30]. F. Yamaguchi, N. Golde, D. Arp, K. Rieck, Modeling and Discovering Vulnerabilities with Code Property Graphs, in: 2014 IEEE Symp. Secur. Priv., 2014: pp. 590–604. <https://doi.org/10.1109/SP.2014.44>.
- [31]. J. Liu, J. Zeng, X. Wang, Z. Liang, Learning Graph-based Code Representations for Source-level Functional Similarity Detection, in: 2023 IEEEACM 45th Int. Conf. Softw. Eng. ICSE, 2023: pp. 345–357. <https://doi.org/10.1109/ICSE48619.2023.00040>.
- [32]. M. Pradel, G. Gousios, J. Liu, S. Chandra, TypeWriter: neural type prediction with search-based validation, in: Proc. 28th ACM Jt. Meet. Eur. Softw. Eng. Conf. Symp. Found. Softw. Eng., Association

- for Computing Machinery, New York, NY, USA, 2020: pp. 209–220. <https://doi.org/10.1145/3368089.3409715>.
- [33]. K. Jesse, P.T. Devanbu, T. Ahmed, Learning type annotation: is big data enough?, in: Proc. 29th ACM Jt. Meet. Eur. Softw. Eng. Conf. Symp. Found. Softw. Eng., Association for Computing Machinery, New York, NY, USA, 2021: pp. 1483–1486. <https://doi.org/10.1145/3468264.3473135>.
- [34]. Y. Peng, C. Wang, W. Wang, C. Gao, M.R. Lyu, Generative Type Inference for Python, (2023).
- [35]. J. Schrouff, K. Wohlfahrt, B. Marnette, L. Atkinson, Inferring javascript types using graph neural networks, ArXiv Prepr. ArXiv190506707 (2019).
- [36]. Z. Liu, F. Jiang, Y. Hu, C. Shi, P. Fung, NER-BERT: A Pre-trained Model for Low-Resource Entity Tagging, (2021).
- [37]. H. Darji, J. Mitrović, M. Granitzer, German BERT Model for Legal Named Entity Recognition, in: Proc. 15th Int. Conf. Agents Artif. Intell., SCITEPRESS - Science and Technology Publications, 2023. <https://doi.org/10.5220/0011749400003393>.
- [38]. D. Bieber, R. Goel, D. Zheng, H. Larochelle, D. Tarlow, Static Prediction of Runtime Errors by Learning to Execute Programs with External Resource Descriptions, (2022).
- [39]. S. Sun, S. Wang, X. Wang, Y. Xing, E. Zhang, K. Sun, Exploring Security Commits in Python, (2023).
- [40]. V.-A. Nguyen, D.Q. Nguyen, V. Nguyen, T. Le, Q.H. Tran, D. Phung, ReGVD: Revisiting Graph Neural Networks for Vulnerability Detection, ArXiv Prepr. ArXiv211007317 (2021).
- [41]. Z. Li, D. Zou, S. Xu, X. Ou, H. Jin, S. Wang, Z. Deng, Y. Zhong, Vuldeepecker: A deep learning-based system for vulnerability detection, ArXiv Prepr. ArXiv180101681 (2018).
- [42]. S. Cao, X. Sun, L. Bo, Y. Wei, B. Li, Bgnn4vd: constructing bidirectional graph neural-network for vulnerability detection, Inf. Softw. Technol. 136 (2021) 106576.
- [43]. Z. Li, D. Zou, S. Xu, H. Jin, Y. Zhu, Z. Chen, SySeVR: A Framework for Using Deep Learning to Detect Software Vulnerabilities, IEEE Trans. Dependable Secure Comput. (2021) 1–1. <https://doi.org/10.1109/TDSC.2021.3051525>.
- [44]. R. Russell, L. Kim, L. Hamilton, T. Lazovich, J. Harer, O. Ozdemir, P. Ellingwood, M. McConley, Automated vulnerability detection in source code using deep representation learning, in: 2018 17th IEEE Int. Conf. Mach. Learn. Appl. ICMLA, IEEE, 2018: pp. 757–762.
- [45]. P.S. Kostenetskiy, R.A. Chulkevich, V.I. Kozyrev, HPC Resources of the Higher School of Economics, J. Phys. Conf. Ser. 1740 (2021) 012050. <https://doi.org/10.1088/1742-6596/1740/1/012050>.
- [46]. Q. Tan, N. Liu, X. Huang, R. Chen, S.-H. Choi, X. Hu, MGAE: Masked Autoencoders for Self-Supervised Learning on Graphs, (2022).
- [47]. S. Zhang, H. Chen, H. Yang, X. Sun, P.S. Yu, G. Xu, Graph Masked Autoencoders with Transformers, (2022).
- [48]. Z. Hou, X. Liu, Y. Cen, Y. Dong, H. Yang, C. Wang, J. Tang, GraphMAE: Self-Supervised Masked Graph Autoencoders, (2022).
- [49]. S. Chopra, R. Hadsell, Y. LeCun, Learning a similarity metric discriminatively, with application to face verification, in: 2005 IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit. CVPR05, 2005: pp. 539–546 vol. 1. <https://doi.org/10.1109/CVPR.2005.202>.
- [50]. W. Liao, M.Y. Yang, N. Zhan, B. Rosenhahn, Triplet-based Deep Similarity Learning for Person Re-Identification, (2018).
- [51]. T. Cover, P. Hart, Nearest neighbor pattern classification, IEEE Trans. Inf. Theory 13 (1967) 21–27. <https://doi.org/10.1109/TIT.1967.1053964>.
- [52]. S. Tipirneni, M. Zhu, C.K. Reddy, StructCoder: Structure-Aware Transformer for Code Generation, (2022). <https://doi.org/10.48550/arXiv.2206.05239>.
- [53]. R. Bavishi, M. Pradel, K. Sen, Context2Name: A Deep Learning-Based Approach to Infer Natural Variable Names from Usage Contexts, (2018).
- [54]. Zhang, L. Fang, C. Ge, P. Li, Z. Liu, Efficient transformer with code token learner for code clone detection, J Syst Softw 197 (2023). <https://doi.org/10.1016/j.jss.2022.111557>.
- [55]. W. Wang, G. Li, B. Ma, X. Xia, Z. Jin, Detecting code clones with graph neural network and flow-augmented abstract syntax tree, in: 2020 IEEE 27th Int. Conf. Softw. Anal. Evol. Reengineering SANER, IEEE, 2020: pp. 261–271.
- [56]. J. Yasarwi, S. Purini, C.V. Jawahar, Plagiarism Detection in Programming Assignments Using Deep Features, in: 2017 4th IAPR Asian Conf. Pattern Recognit. ACPR, 2017: pp. 652–657. <https://doi.org/10.1109/ACPR.2017.146>.
- [57]. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, Language models are unsupervised multitask learners, OpenAI Blog 1 (2019) 9.

- [58]. T. Brown, B. Mann, N. Ryder, M. Subbiah, J.D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, Language models are few-shot learners, *Adv. Neural Inf. Process. Syst.* 33 (2020) 1877–1901.
- [59]. A.V.M. Barone, R. Sennrich, A parallel corpus of python functions and documentation strings for automated code documentation and code generation, *ArXiv Prepr. ArXiv170702275* (2017).
- [60]. X. Liu, D. Wang, A. Wang, Y. Hou, L. Wu, HACONVGGNN: Hierarchical attention based convolutional graph neural network for code documentation generation in jupyter notebooks, *ArXiv Prepr. ArXiv210401002* (2021).
- [61]. S. Bhatia, R. Singh, Automated correction for syntax errors in programming assignments using recurrent neural networks, *ArXiv Prepr. ArXiv160306129* (2016).
- [62]. Marginean, J. Bader, S. Chandra, M. Harman, Y. Jia, K. Mao, A. Mols, A. Scott, Sapfix: Automated end-to-end repair at scale, in: 2019 IEEEACM 41st Int. Conf. Softw. Eng. Softw. Eng. Pract. ICSE-SEIP, IEEE, 2019: pp. 269–278.
- [63]. Khajenezhad, S.A. Osia, M. Karimian, H. Beigy, Gransformer: Transformer-based Graph Generation, (2022).
- [64]. S. Cabrera Lozoya, A. Baumann, A. Sabetta, M. Bezzi, Commit2vec: Learning distributed representations of code changes, *SN Comput. Sci.* 2 (2021) 1–16.
- [65]. G. Arutyunov, gaarutyunov/graph-typer, (2024). <https://github.com/gaarutyunov/graph-typer> (accessed August 12, 2024).

Информация об авторах / Information about authors

Герман Аренович АРУТЮНОВ – магистр факультета компьютерных наук НИУ ВШЭ. Сфера научных интересов: генерация и анализ языков программирования посредством машинного обучения и глубоких нейронных сетей.

German Arsenovich ARUTYUNOV – Master’s graduate at the Faculty of Computer Science at HSE University. Research interests include programming language generation and programming language understanding using machine learning and deep neural networks.

Сергей Михайлович АВДОШИН – кандидат технических наук, профессор департамента компьютерной инженерии Московского института электроники и математики им. А.Н. Тихонова НИУ ВШЭ. Сфера научных интересов: разработка и анализ компьютерных алгоритмов, имитация и моделирование, параллельные и распределенные процессы, машинное обучение.

Sergey Mikchailovitch AVDOSHTIN – Cand. Sci. (Tech.), Professor of the School of Computer Engineering at Tikhonov Moscow Institute of Electronics and Mathematics HSE University. Research interests include design and analysis of computer algorithms, simulation and modeling, parallel and distributed processing, machine learning.

