
A Common Graph Representation for Source Code and Developer Documentation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Semantic information plays a key role in the code search and synthesis settings. In
2 this work, we propose a graph-based representation for source code and natural
3 language which incorporates semantic and relational features from both domains.
4 We apply this graph to a parsing a corpus of code and developer documents, and
5 demonstrate the effectiveness of a common graph-based representation on three
6 downstream tasks: code search, document recommendation and link prediction.

1 Background and motivation

8 In addition to its syntactic structure, source code contains a rich denotational and operational
9 semantics [Henkel et al., 2018]. To effectively reason about code in semantically similar but
10 syntactically diverse settings requires models which incorporate features from the call graph [Gu
11 et al., 2016, Liu et al., 2019] and surrounding typing context [Allamanis et al., 2017]. Many semantic
12 features, such as data and control flow [Si et al., 2018] can be represented as a directed acyclic graph
13 (DAG), which admits linear-time solutions to a number of graph problems, including topological
14 sorting, single-source shortest path and reachability queries.

15 DAGs also have important applications in natural language parsing [Sagae and Tsujii, 2008, Quern-
16 heim and Knight, 2012]. Various attempts to build semantic representations for natural language have
17 been proposed, notably the pointer network architecture [Vinyals et al., 2015b,a]. Pointer networks
18 help to capture permutation-invariant semantic relations between natural language entities, and have
19 important applications in dependency parsing [Ma et al., 2018], named-entity recognition [Lample
20 et al., 2016], and other tasks where sequence representations fall short. Li et al. [2017] extend pointer
21 networks with a copy-mechanism to handle out-of-vocabulary code tokens.

22 Content recommendation for doc-to-doc (D2D) and code-to-code (C2C) is a relatively straightforward
23 application of existing link prediction [Zhang and Chen, 2018] and code embedding [Gu et al., 2018]
24 techniques, but cross-domain transfer remains largely unsolved. Robillard and Chhetri [2015],
25 Robillard et al. [2017] first explore the task of predicting reference API documentation from source
26 code using manual annotation. Prior work also studies the association between comments and code
27 entities [Panthaplackel et al., 2020] using machine learning, but only within source code.

28 Maintainers of widely-used software projects often publish web-based documentation, typically stored
29 in markup languages like HTML or Markdown. These files contain a collection of natural language
30 sentences, markup, and hyperlinks to other documents. Both the link graph and the document tree
31 contain important semantic information: the markup describes the text in relation to the other entities
32 in the document hierarchy [Yang et al., 2016], while the link graph describes the relationship between
33 the parent document and related documents or source code entities. Documents occasionally contain
34 hyperlinks to source code, but source code rarely contains links to developer documents.

35 Some programming languages allow users to specify which type of values will inhabit a given variable
36 at runtime. Types allow the compiler to reason about certain properties like nullity [Ekman and Hedin,
37 2007] and shape [Considine et al., 2019]. While types may not appear explicitly in source code,
38 they can often be inferred from the surrounding context using a dataflow graph (DFG). The Java
39 language recently introduced local variable type inference Liddell and Kim [2019], which allows
40 variable types to be omitted, and later inferred by the compiler.

41 **2 Proposed approach**

42 Given a single token in either source code or developer documentation and its surrounding context,
43 what are the most relevant source code or documentation entities related to the token in question? We
44 would like to infer which entities are relevant to a particular token, based on the semantic context.
45 To infer links across these two domains requires building a multi-relational graph, using features
46 extracted from both natural language and source code. Following Si et al. [2018], Gu et al. [2018], Liu
47 et al. [2019], we use a node embedding on the dataflow graph and type environment, and following
48 Yang et al. [2016], Zhang and Chen [2018], use the markup hierarchy and link graph to construct an
49 embedding for code-like tokens used within documentation.

50 To compensate for the sparsity of hyperlinks between code and documentation, we must design a
51 heuristic to connect the documentation graph and source code entities. One heuristic which developers
52 often use to discover relevant documents is plaintext search on a salient lexical string. Co-occurrence
53 of an infrequent token indicates the two entities are likely related, even though they may not share
54 an explicit grammatical link. If we can recover this relationship without observing the lexical token
55 itself, only using dataflow and type-related information, this indicates our representation is providing
56 useful information.

57 **3 Data availability and computational requirements**

58 Java, one of the most popular programming languages on GitHub, is a statically typed language with
59 an extensive amount of API documentation on the web. It has a variety of tools for parsing and
60 analyzing both code [Kovalenko et al., 2019] and natural language [Manning et al., 2014, Grella and
61 Cangialosi, 2018], making it a suitable candidate both as a dataset and implementation language. Our
62 dataset consists of Java repositories on GitHub, and their accompanying docs on the Zeal software
63 documentation aggregator. All projects in our dataset have a collection of source code files and
64 multiple related repositories on GitHub.

65 We construct two datasets consisting of naturally-occurring links between developer documentation
66 and source code, and a surrogate set of links constructed by matching lexical tokens available in
67 both domains. Our target is the recovery of ground truth links in our test set and surrogate links
68 in the lexical matching graph. By adding weighted edges between source code and documentation
69 and learning the relations using, e.g. a pointer network architecture, we evaluate our approach on
70 reconstructing synthetic links between tokens contained in code-like fragments and markup entities
71 which refer to the selected token. In addition, we evaluate our approach both on D2D and C2C link
72 retrieval, as well as precision and recall on the surrogate link graph.

73 To perform our experiments, we require a large number of CPUs for semantic parsing, link extraction
74 and graph preprocessing, and a single P100 GPU for training a graph neural network. We have
75 applied and received access to the Niagra CC cluster.

76 **References**

- 77 Miltiadis Allamanis, Marc Brockschmidt, and Mahmoud Khademi. Learning to represent programs
78 with graphs. *arXiv preprint arXiv:1711.00740*, 2017. URL <https://arxiv.org/pdf/1711.00740.pdf>.
79
- 80 Breandan Considine, Michalis Famelis, and Liam Paull. Kotlin ∇ : A shape-safe eDSL for differen-
81 tiable programming. <https://github.com/breandan/kotlingrad>, 2019.

- 82 Torbjörn Ekman and Görel Hedin. Pluggable checking and inferencing of nonnull types for Java.
83 *Journal of Object Technology*, 6(9):455–475, 2007. URL [http://www.jot.fm/issues/issue_](http://www.jot.fm/issues/issue_2007_10/paper23.pdf)
84 [2007_10/paper23.pdf](http://www.jot.fm/issues/issue_2007_10/paper23.pdf).
- 85 Matteo Grella and Simone Cangialosi. Non-projective dependency parsing via latent heads repre-
86 sentation LHR. *arXiv preprint arXiv:1802.02116*, 2018. URL [https://arxiv.org/pdf/1802.](https://arxiv.org/pdf/1802.02116.pdf)
87 [02116.pdf](https://arxiv.org/pdf/1802.02116.pdf).
- 88 Xiaodong Gu, Hongyu Zhang, Dongmei Zhang, and Sunghun Kim. Deep API learning. In *Pro-*
89 *ceedings of the 2016 24th ACM SIGSOFT International Symposium on Foundations of Software*
90 *Engineering*, pages 631–642, 2016. URL <https://arxiv.org/pdf/1605.08535.pdf>.
- 91 Xiaodong Gu, Hongyu Zhang, and Sunghun Kim. Deep code search. In *2018 IEEE/ACM 40th*
92 *International Conference on Software Engineering (ICSE)*, pages 933–944. IEEE, 2018. URL
93 <https://guxd.github.io/papers/deepcs.pdf>.
- 94 Jordan Henkel, Shuvendu K Lahiri, Ben Liblit, and Thomas Reps. Code vectors: Understanding
95 programs through embedded abstracted symbolic traces. In *Proceedings of the 2018 26th ACM*
96 *Joint Meeting on European Software Engineering Conference and Symposium on the Foundations*
97 *of Software Engineering*, pages 163–174, 2018. URL [https://arxiv.org/pdf/1803.06686.](https://arxiv.org/pdf/1803.06686.pdf)
98 [pdf](https://arxiv.org/pdf/1803.06686.pdf).
- 99 Vladimir Kovalenko, Egor Bogomolov, Timofey Bryksin, and Alberto Bacchelli. PathMiner: a
100 library for mining of path-based representations of code. In *Proceedings of the 16th International*
101 *Conference on Mining Software Repositories*, pages 13–17. IEEE Press, 2019. URL [https:](https://github.com/JetBrains-Research/astminer)
102 [//github.com/JetBrains-Research/astminer](https://github.com/JetBrains-Research/astminer).
- 103 Guillaume Lample, Miguel Ballesteros, Sandeep Subramanian, Kazuya Kawakami, and Chris Dyer.
104 Neural architectures for named entity recognition. *arXiv preprint arXiv:1603.01360*, 2016. URL
105 <https://arxiv.org/pdf/1603.01360.pdf>.
- 106 Jian Li, Yue Wang, Michael R Lyu, and Irwin King. Code completion with neural attention and
107 pointer networks. *arXiv preprint arXiv:1711.09573*, 2017. URL [https://www.ijcai.org/](https://www.ijcai.org/Proceedings/2018/0578.pdf)
108 [Proceedings/2018/0578.pdf](https://www.ijcai.org/Proceedings/2018/0578.pdf).
- 109 Clayton Liddell and Donghoon Kim. Analyzing the adoption rate of local variable type infer-
110 ence in open-source Java 10 projects. *Journal of the Arkansas Academy of Science*, 73(1):
111 51–54, 2019. URL [https://scholarworks.uark.edu/cgi/viewcontent.cgi?article=](https://scholarworks.uark.edu/cgi/viewcontent.cgi?article=3346&context=jaas)
112 [3346&context=jaas](https://scholarworks.uark.edu/cgi/viewcontent.cgi?article=3346&context=jaas).
- 113 Bohong Liu, Tao Wang, Xunhui Zhang, Qiang Fan, Gang Yin, and Jinsheng Deng. A neural-network
114 based code summarization approach by using source code and its call dependencies. In *Proceedings*
115 *of the 11th Asia-Pacific Symposium on Internetworking*, Internetworking ’19, New York, NY, USA, 2019.
116 Association for Computing Machinery. ISBN 9781450377010. doi: 10.1145/3361242.3362774.
117 URL <https://doi.org/10.1145/3361242.3362774>.
- 118 Xuezhe Ma, Zecong Hu, Jingzhou Liu, Nanyun Peng, Graham Neubig, and Eduard Hovy. Stack-
119 pointer networks for dependency parsing. *arXiv preprint arXiv:1805.01087*, 2018. URL [https:](https://arxiv.org/pdf/1805.01087.pdf)
120 [//arxiv.org/pdf/1805.01087.pdf](https://arxiv.org/pdf/1805.01087.pdf).
- 121 Christopher D Manning, Mihai Surdeanu, John Bauer, Jenny Rose Finkel, Steven Bethard, and David
122 McClosky. The stanford coreNLP natural language processing toolkit. In *Proceedings of 52nd*
123 *annual meeting of the association for computational linguistics: system demonstrations*, pages
124 55–60, 2014. URL <https://nlp.stanford.edu/pubs/StanfordCoreNlp2014.pdf>.
- 125 Sheena Panthaplackel, Milos Gligoric, Raymond J. Mooney, and Junyi Jessy Li. Associating natural
126 language comment and source code entities. In *AAAI*, 2020. URL [https://arxiv.org/pdf/](https://arxiv.org/pdf/1912.06728.pdf)
127 [1912.06728.pdf](https://arxiv.org/pdf/1912.06728.pdf).
- 128 Daniel Quernheim and Kevin Knight. Dagger: A toolkit for automata on directed acyclic graphs. In
129 *Proceedings of the 10th International Workshop on Finite State Methods and Natural Language*
130 *Processing*, pages 40–44, 2012. URL <https://www.aclweb.org/anthology/W12-6207.pdf>.

131 Martin P Robillard and Yam B Chhetri. Recommending reference API documentation. *Empirical*
132 *Software Engineering*, 20(6):1558–1586, 2015. URL [https://www.cs.mcgill.ca/~martin/](https://www.cs.mcgill.ca/~martin/papers/cr2014a.pdf)
133 [papers/cr2014a.pdf](https://www.cs.mcgill.ca/~martin/papers/cr2014a.pdf).

134 Martin P Robillard, Andrian Marcus, Christoph Treude, Gabriele Bavota, Oscar Chaparro, Neil
135 Ernst, Marco Aurélio Gerosa, Michael Godfrey, Michele Lanza, Mario Linares-Vásquez, et al.
136 On-demand developer documentation. In *2017 IEEE International Conference on Software*
137 *Maintenance and Evolution (ICSME)*, pages 479–483. IEEE, 2017.

138 Kenji Sagae and Jun’ichi Tsujii. Shift-reduce dependency DAG parsing. In *Proceedings of the 22nd*
139 *International Conference on Computational Linguistics-Volume 1*, pages 753–760. Association
140 for Computational Linguistics, 2008. URL [https://www.aclweb.org/anthology/C08-1095.](https://www.aclweb.org/anthology/C08-1095.pdf)
141 [pdf](https://www.aclweb.org/anthology/C08-1095.pdf).

142 Xujie Si, Hanjun Dai, Mukund Raghothaman, Mayur Naik, and Le Song. Learning
143 loop invariants for program verification. In *Advances in Neural Information Pro-*
144 *cessing Systems*, pages 7751–7762, 2018. URL [https://papers.nips.cc/paper/](https://papers.nips.cc/paper/8001-learning-loop-invariants-for-program-verification.pdf)
145 [8001-learning-loop-invariants-for-program-verification.pdf](https://papers.nips.cc/paper/8001-learning-loop-invariants-for-program-verification.pdf).

146 Oriol Vinyals, Samy Bengio, and Manjunath Kudlur. Order matters: Sequence to sequence for sets.
147 *arXiv preprint arXiv:1511.06391*, 2015a. URL <https://arxiv.org/pdf/1511.06391.pdf>.

148 Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. Pointer networks. In *Advances in neural*
149 *information processing systems*, pages 2692–2700, 2015b. URL [https://arxiv.org/pdf/](https://arxiv.org/pdf/1506.03134.pdf)
150 [1506.03134.pdf](https://arxiv.org/pdf/1506.03134.pdf).

151 Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He, Alex Smola, and Eduard Hovy. Hierarchi-
152 cal attention networks for document classification. In *Proceedings of the 2016 conference of*
153 *the North American chapter of the association for computational linguistics: human language*
154 *technologies*, pages 1480–1489, 2016. URL [https://www.cs.cmu.edu/~hovy/papers/](https://www.cs.cmu.edu/~hovy/papers/16HLT-hierarchical-attention-networks.pdf)
155 [16HLT-hierarchical-attention-networks.pdf](https://www.cs.cmu.edu/~hovy/papers/16HLT-hierarchical-attention-networks.pdf).

156 Muhan Zhang and Yixin Chen. Link prediction based on graph neural networks. In *Advances in*
157 *Neural Information Processing Systems*, pages 5165–5175, 2018. URL [https://papers.nips.](https://papers.nips.cc/paper/7763-link-prediction-based-on-graph-neural-networks.pdf)
158 [cc/paper/7763-link-prediction-based-on-graph-neural-networks.pdf](https://papers.nips.cc/paper/7763-link-prediction-based-on-graph-neural-networks.pdf).