# A Literature Study of Embeddings on Source Code

**Zimin Chen & Martin Monperrus**
School of Electrical Engineering and Computer Science
KTH Royal Institute of Technology
zimin@kth.se, martin.monperrus@csc.kth.se

## ABSTRACT

Natural language processing has improved tremendously after the success of word embedding techniques such as word2vec. Recently, the same idea has been applied on source code with encouraging results. In this survey, we aim to collect and discuss the usage of word embedding techniques on programs and source code. The articles in this survey have been collected by asking authors of related work and with an extensive search on Google Scholar. Each article is categorized into five categories: 1. embedding of tokens 2. embedding of functions or methods 3. embedding of sequences or sets of method calls 4. embedding of binary code 5. other embeddings. We also provide links to experimental data and show some remarkable visualization of code embeddings. In summary, word embedding has been successfully applied on different granularities of source code. With access to countless open-source repositories, we see a great potential of applying other data-driven natural language processing techniques on source code in the future.

## 1 INTRODUCTION

With recent success in deep learning, it has become more and more popular to apply it on images, natural language, audios and *etc.* One key aspect of this research is to compute a numerical vector representing the object of interest: this vector is called an "embedding". For words in natural language, researchers have come up with a set of techniques, called word embedding, that maps words or phrases to a vector of real numbers.

Moving from natural language words and sentences to programming language tokens and statements is natural. It is also meaningful because Hindle et al. (2012) showed that programming languages are like natural languages, with comparable repetitiveness and predictability . In this line of thought, authors have envisioned potentially powerful applications of word embeddings on source code, and experimental results are encouraging. This paper presents this new research area on embeddings on source code.

Much like how natural languages have characters, words, sentences and paragraphs, programming languages also have different granularities, such variables, expressions, statements and methods. Therefore, word embeddings can be applied to source code on different granularities.

In this paper, we collect and discuss the usage of word embeddings in programs. Our methodology is to identify articles that compute and use an embedding on source code. We do that with asking authors of related work directly and extensive search with Google Scholar.

Our contributions are:

- A survey of word embeddings on source code in five categories: embeddings of tokens, embeddings of expressions, embeddings of APIs, embeddings of methods, and other miscellaneous embeddings.

- A collection of visual representations of embeddings, that add an aesthetic dimension to the power of code embeddings.

- A curated list of publicly available code embedding data.

## 2 BACKGROUND

### 2.1 EMBEDDINGS

An embedding is mapping from objects to vectors of real numbers. Word embedding refers to all natural language processing approaches where words (or sequence of words) are mapped onto vectors of real numbers. It makes it possible to work with textual data in a mathematical model. It also has the advantage that fundamentally discrete data (words) is transformed into continuous vectors space. With recent progress in word embedding, especially word2vec (Mikolov et al. (2013)), the embedding vectors preserve the semantic. Word embedding such as word2vec requires large unlabeled corpora to train on.

### 2.2 VISUALIZATION

Embedding vectors usually have more than 3 dimensions, up to hundreds and even thousands of dimensions. This means we have to reduce the dimensionality in order to visualize them in a 2D plane or 3D space. The most commonly used techniques are Principal Component Analysis (PCA, Pearson (1901)) and t-Distributed Stochastic Neighbour Embedding (t-SNE, Maaten & Hinton (2008)).

### 2.3 SIMILARITY METRICS IN THE EMBEDDING SPACE

Similarity between words can be measured by calculating a similarity metric between their embeddings. The most used similarity metric for word embeddings is cosine similarity. The cosine similarity measures the angle between two vectors, which is independent of their magnitude. Other similarity metrics exists such as Euclidean distance, but in high dimensional space where the data points are sparse, Beyer et al. (1999) showed that the ratio of distance between the nearest and farthest point is close to 1. And Aggarwal et al. (2001) proved that between $L_k$ norms, smaller k is more preferable in high dimensional space, *i.e.,* $L_1$ norm (Manhattan distance) is more preferable than $L_2$ norm (Euclidean distance).

## 3 EMBEDDING FOR SOURCE CODE

Here we present source code embeddings, each subsection focusing a specific granularity. Those different granularities are needed depending on the downstream task, *e.g.,* we need token embedding in code completion and we need function embedding in function clone detection. Table 1 shows links to publicly available experimental data.

### 3.1 EMBEDDING OF TOKENS

Harer et al. (2018) use word2vec to generate word embedding for C/C++ tokens for software vulnerability prediction. The token embedding is used to initialize a TextCNN model for classification.

White et al. (2017) generate Java token embedding using word2vec for automatic program repair. They used the embedding to initialize a recursive encoder of abstract syntax trees.

Azcona et al. (2019) explore Python token embeddings for profiling student submissions. Each student's submission is transformed into a flattened list of token embeddings, and then each student is represented as one matrix where each row represents one submission.

Chen & Monperrus (2018) use word2vec to generate java token embedding for finding the correct ingredient in automated program repair . They use cosine similarity on the embeddings to compute a distance between pieces of code.

### 3.2 EMBEDDING OF FUNCTIONS OR METHODS

Alon et al. (2019) compute Java method embeddings for predicting method names. Abstract syntax tree is used to construct a path representation between two leaf nodes. Bag of path representations are all aggregated into one single embedding for a method. In a subsequent paper Alon et al. (2018),
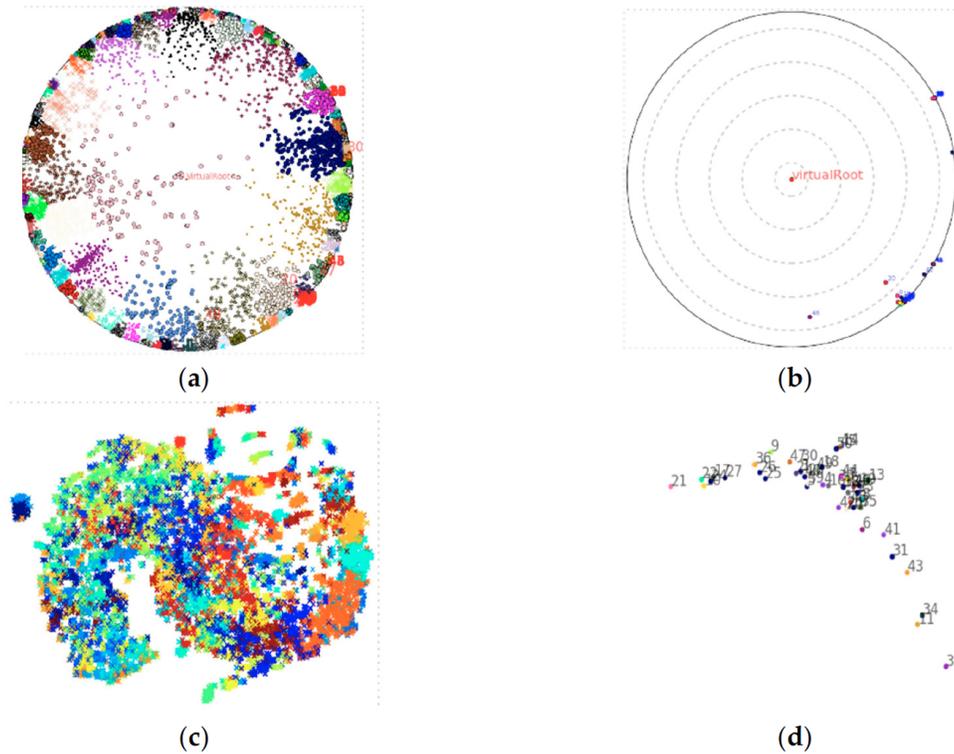
Figure 1: Function embedding from Lu et al. (2019), reproduced with permission.

the path representation is computed using Long Short Term Memory network (LSTM) instead of a single layer neural network.

Allamanis et al. (2015) use a logbilinear context model to generate an embedding for method names. They define a local context which captures tokens around the current token, and a global context which is composed of features from the method code. Their technique is able to generate new method names never seen in the training data.

DeFreez et al. (2018) generate function embeddings for C code using control-flow graphs. They perform a random walk on interprocedural paths in the program, and used the paths to generate function embeddings. The embedding is used for detecting function clones.

Murali et al. (2017) compute a function embedding for program sketches. They use the embedding to generate source code given a specification that lists API calls and API types to be used in the generated code. An encoder-decoder model is used to generate program sketches from the specification. The final hidden state of the encoder is seen as embedding for the function.

Lu et al. (2019) propose a new function embedding method that learns embeddings in a hyperbolic space. They first construct function call graphs where the weight of edges are replaced with a Ricci curvature. The embeddings are then learned by using a Poincaré model. The function embedding in is visualized in Figure 1.

Devlin et al. (2017) use function embedding as input to repair variable misuse in Python. They encode the function AST by doing a depth first traversal and create an embedding by concatenating the absolute position of the node, the type of node, the relationship between the node and its parent and the string label of the node.

Büch & Andrzejak (2019) generate embedding for methods to do code clone detection, by jointly considering embeddings of AST node type and of node contents. The AST to be embedded is traversed from leaf nodes to the root, and the obtained sequences are fed to a LSTM for learning sequences. The similarity metric between methods is computed using a Siamese network.
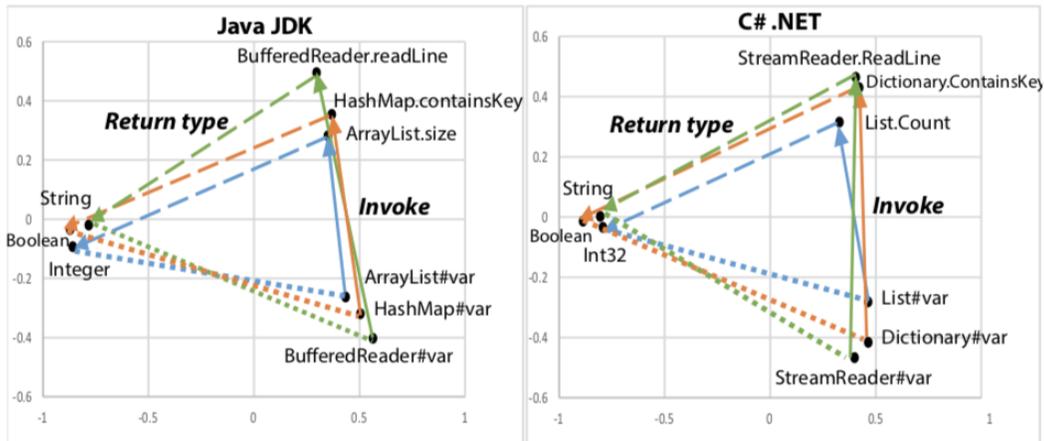
3

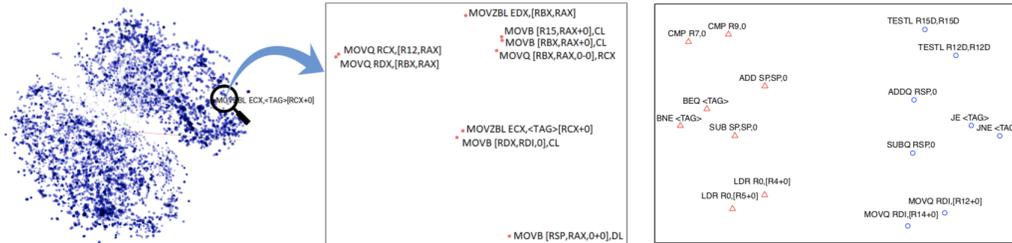Figure 2: API embedding from Nguyen et al. (2016), reproduced with permission.



Figure 3: Binary instruction embedding from Zuo et al. (2018), reproduced with permission.

## 3.3 EMBEDDING OF SEQUENCES OR SETS OF METHOD CALLS

Henkel et al. (2018) generate embedding for abstracted symbolic traces from C projects. The traces are collected from symbolic executions, and then names in symbolic traces are renamed to reduce the vocabulary size. The embedding is learned from the abstracted symbolic traces using word2vec.

Nguyen et al. (2016) explore embeddings for Java and C# APIs, and used it to find similar API usage between the languages. They use word2vec for generating API element embeddings, and based on known API mappings, they can find cluster with similar usage across two languages. The API embedding is visualized in Figure 2. In a subsequent paper Nguyen et al. (2017), they translate source code to API sequences, and learned API embedding from abstracted API sequences using word2vec.

Gu et al. (2016) consider the problem of translating natural language query to a sequence of API calls. They use an encoder-decoder model with custom a loss function, where the input is the natural language query and output is the sequence of API calls. The final hidden state of the encoder can be seen as an embedding for the natural language query.

Pradel & Sen (2018) propose to use code snippet embeddings to identify potential buggy code. To generate the training data, they collect code examples from a code corpus and apply simple code transformation to insert an artificial bug. Embedding for positive and negative code examples are generated and they are used to train a classifier.

## 3.4 EMBEDDING OF BINARY CODE

Ben-Nun et al. (2018) use the LLVM Intermediate Representation (IR) to build a contextual flow graph, which is used to generate binary code embedding for C/C++ statements. The source code is
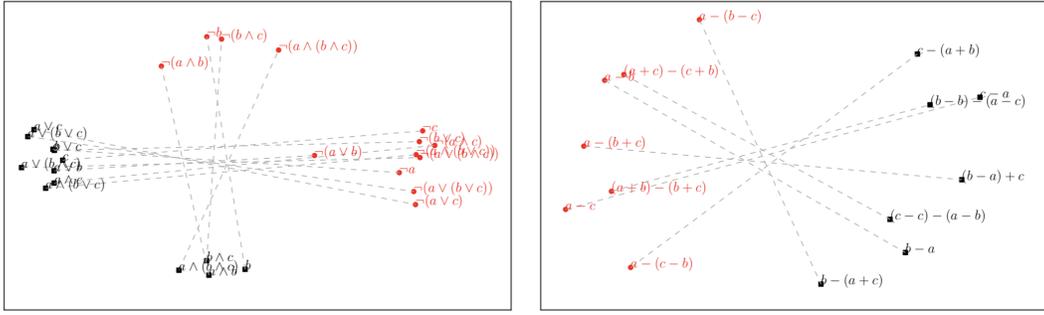
Figure 4: Expression embedding from Allamanis et al. (2017b), reproduced with permission.

first compiled to LLVM IR, where the identifier and literals are renamed. The abstracted representation is used to build a contextual flow graph, where paths are extracted for training the embeddings.

Zuo et al. (2018) compute binary instruction embedding to calculate the similarity between two binary blocks. Binary block embeddings are trained on x86-64 and ARM instructions, and instructions in a basic block are combined using LSTM. A Siamese architecture is adopted to determine the similarity between two basic blocks. The binary instruction embedding is visualized in Figure 3.

Redmond et al. (2018) explore binary instruction embedding across architectures. Binary instruction embeddings are learned by cluster similar instructions in the same architecture, and preserve the semantic relationship between binary instruction in different architectures.

Xu et al. (2017) generate cross-platform binary function embedding for similarity detection. The binary function is converted to control flow graph and embedding is generated using Structure2vec. Then, Siamese network is trained for similarity detection.

## 3.5   OTHER EMBEDDINGS

Allamanis et al. (2017b) use TreeNN to generate embedding for algebraic and logical expressions. Their technique recursively encodes non-leafs node in the abstract syntax tree structure by combining the children nodes using a multi-layer perceptron. The expression embedding is visualized in Figure 4.

Yin et al. (2018) introduce the problem of learning embedding for edits over source code, *i.e.,* patches. The edits are represented as sequence or graph, and a neural network is used to generate the embedding.

Zhao et al. (2018) use an embedding approach to determine if two Android components have an Inter-Component Communication link. An intent is a message sent between different Android components, and an intent filter describes what types of intent the component is willing to receive. The intent and filter embeddings are used to train a classifier for possible Inter-Component Communication link.

Wang et al. (2017) generate program embedding using execution traces. The program is represented as multiple variable traces, and they are interleaved to simulate dependencies between traces. A full program embedding is the average over all embeddings of the traces.

Chistyakov et al. (2018) extract behaviour patterns from a sequence of system events, and combine them to an embedding that represents the program behaviour. First, they generate an embedding using an autoencoder, and then all behaviour patterns are combined using min, max and mean functions over elements.

Allamanis et al. (2017a) use a graph neural network to generate embedding for each node in an abstract syntax tree. The graph is constructed from nodes and edges from the AST, and then enriched with additional edges indicating data flow and type hierarchies.

| | |
|---|---|
| White et al. (2017) | `https://sites.google.com/view/deeprepair/home` |
| Azcona et al. (2019) | `https://github.com/dazcona/user2code2vec` |
| Chen & Monperrus (2018) | `https://github.com/kth-tcs/3sFix-experiments` |
| Alon et al. (2019) | `https://github.com/tech-srl/code2vec` |
| Alon et al. (2018) | `http://github.com/tech-srl/code2seq` |
| Allamanis et al. (2015) | `http://groups.inf.ed.ac.uk/cup/naturalize/` |
| DeFreez et al. (2018) | `https://github.com/defreez-ucd/func2vec-fse2018-artifact` |
| Murali et al. (2017) | `https://github.com/capergroup/bayou` |
| Devlin et al. (2017) | `https://iclr2018anon.github.io/semantic_code_repair/index.html.` |
| Henkel et al. (2018) | `https://github.com/jjhenkel/code-vectors-artifact` |
| Nguyen et al. (2017) | `http://home.eng.iastate.edu/~trong/projects/jv2cs/` |
| Gu et al. (2016) | `https://guxd.github.io/deepapi/` |
| Pradel & Sen (2018) | `https://github.com/michaelpradel/DeepBugs` |
| Zuo et al. (2018) | `https://nmt4binaries.github.io/` |
| Ben-Nun et al. (2018) | `https://github.com/spcl/ncc` |
| Redmond et al. (2018) | `https://github.com/nlp-code-analysis/cross-arch-instr-model` |
| Xu et al. (2017) | `https://github.com/xiaojunxu/dnn-binary-code-similarity` |
| Allamanis et al. (2017b) | `http://groups.inf.ed.ac.uk/cup/semvec/` |
| Yin et al. (2018) | `https://github.com/Microsoft/msrc-dpu-learning-to-represent-edits` |
| Wang et al. (2017) | `https://github.com/keowang/dynamic-program-embedding` |
| Theeten et al. (2019) | `https://zenodo.org/record/2546488#.XKHb2S97GL4` |

Table 1: Publicly available embeddings on code.

Theeten et al. (2019) explore software library embedding for Python, Java and JavaScript. They collected library usage in open source repositories and ignored the import order. The library embedding is trained by consider all imported libraries in the same file as context.

## 4 FUTURE RESEARCH DIRECTIONS

Now, we speculate about potentially interesting future research directions on embedding for code.

### 4.1 EMBEDDING VERSUS DOWNSTREAM TASKS

Code embeddings should capture the semantics, but it is a difficult task to assess its quality in a general manner. One way is to use analogies ( $a$ is to $b$, what $x$ is to $y$ ). Another way could be to list top-n closest elements in an embedding space. Yet, the state-of-the-art in NLP suggests to evaluate code embedding techniques on downstream tasks rather than in a generic way. The existing embeddings can be used as initialization for the downstream tasks.

### 4.2 CONTEXTUAL WORD EMBEDDING

Word embedding techniques have the drawback of having a one-to-one mapping between elements and their embeddings. Consider the Java keyword 'static', which have different meaning depending on put on a field or a method: it should also have two different embeddings for those two different contexts. Contextual word embedding is a technique that takes the context around the word into account, approach like BERT by Devlin et al. (2018) has achieved state-of-the-art results in several NLP tasks. Considering how powerful contextual word embedding is, and how the meaning of code elements can be contextual, we believe that it is a promising direction machine learning on code. We note that Allamanis et al. (2017a) can be seen as an implicit contextual word embedding, even if it does not use this terminology.

## 5 CONCLUSION

We have collected and discussed articles that define or use embedding on source code. This literature study shows how diverse embedding techniques are, and they are applied to different downstream tasks. We think that the field of embedding on code is only at the beginning and we expect that many papers will use this key and fascinating concept in the future.
The authors would like to thank Jordan Henkel, Miltos Allamanis and Hugo Mougard for valuable pointers and discussions.

REFERENCES

Charu C Aggarwal, Alexander Hinneburg, and Daniel A Keim. On the surprising behavior of distance metrics in high dimensional space. In *International conference on database theory*, pp. 420–434. Springer, 2001.

Miltiadis Allamanis, Earl T Barr, Christian Bird, and Charles Sutton. Suggesting accurate method and class names. In *Proceedings of the 2015 10th Joint Meeting on Foundations of Software Engineering*, pp. 38–49. ACM, 2015.

Miltiadis Allamanis, Marc Brockschmidt, and Mahmoud Khademi. Learning to represent programs with graphs. *arXiv preprint arXiv:1711.00740*, 2017a.

Miltiadis Allamanis, Pankajan Chanthirasegaran, Pushmeet Kohli, and Charles Sutton. Learning continuous semantic representations of symbolic expressions. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 80–88. JMLR. org, 2017b.

Uri Alon, Omer Levy, and Eran Yahav. code2seq: Generating sequences from structured representations of code. *arXiv preprint arXiv:1808.01400*, 2018.

Uri Alon, Meital Zilberstein, Omer Levy, and Eran Yahav. code2vec: Learning distributed representations of code. *Proceedings of the ACM on Programming Languages*, 3(POPL):40, 2019.

David Azcona, Piyush Arora, I-Han Hsiao, and Alan Smeaton. user2code2vec: Embeddings for profiling students based on distributional representations of source code. In *Proceedings of the 9th International Conference on Learning Analytics & Knowledge*, pp. 86–95. ACM, 2019.

Tal Ben-Nun, Alice Shoshana Jakobovits, and Torsten Hoefler. Neural code comprehension: A learnable representation of code semantics. In *Advances in Neural Information Processing Systems*, pp. 3589–3601, 2018.

Kevin Beyer, Jonathan Goldstein, Raghu Ramakrishnan, and Uri Shaft. When is "nearest neighbor" meaningful? In *International conference on database theory*, pp. 217–235. Springer, 1999.

L. Büch and A. Andrzejak. Learning-based recursive aggregation of abstract syntax trees for code clone detection. In *2019 IEEE 26th International Conference on Software Analysis, Evolution and Reengineering (SANER)*, pp. 95–104, Feb 2019. doi: 10.1109/SANER.2019.8668039.

Zimin Chen and Martin Monperrus. The remarkable role of similarity in redundancy-based program repair. *arXiv preprint arXiv:1811.05703*, 2018.

Alexander Chistyakov, Ekaterina Lobacheva, Arseny Kuznetsov, and Alexey Romanenko. Semantic embeddings for program behavior patterns. *arXiv preprint arXiv:1804.03635*, 2018.

Daniel DeFreez, Aditya V Thakur, and Cindy Rubio-González. Path-based function embedding and its application to specification mining. *arXiv preprint arXiv:1802.07779*, 2018.

Jacob Devlin, Jonathan Uesato, Rishabh Singh, and Pushmeet Kohli. Semantic code repair using neuro-symbolic transformation networks. *arXiv preprint arXiv:1710.11054*, 2017.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.

Xiaodong Gu, Hongyu Zhang, Dongmei Zhang, and Sunghun Kim. Deep api learning. In *Proceedings of the 2016 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering*, pp. 631–642. ACM, 2016.

Jacob A Harer, Louis Y Kim, Rebecca L Russell, Onur Ozdemir, Leonard R Kosta, Akshay Rangamani, Lei H Hamilton, Gabriel I Centeno, Jonathan R Key, Paul M Ellingwood, et al. Automated software vulnerability detection with machine learning. *arXiv preprint arXiv:1803.04497*, 2018.

Jordan Henkel, Shuvendu K Lahiri, Ben Liblit, and Thomas Reps. Code vectors: understanding programs through embedded abstracted symbolic traces. In *Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, pp. 163–174. ACM, 2018.

Abram Hindle, Earl T Barr, Zhendong Su, Mark Gabel, and Premkumar Devanbu. On the naturalness of software. In *2012 34th International Conference on Software Engineering (ICSE)*, pp. 837–847. IEEE, 2012.

Mingming Lu, Yan Liu, Haifeng Li, Dingwu Tan, Xiaoxian He, Wenjie Bi, and Wendbo Li. Hyperbolic function embedding: Learning hierarchical representation for functions of source code in hyperbolic space. *Symmetry*, 11(2):254, 2019.

Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605, 2008.

Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pp. 3111–3119, 2013.

Vijayaraghavan Murali, Letao Qi, Swarat Chaudhuri, and Chris Jermaine. Neural sketch learning for conditional program generation. *arXiv preprint arXiv:1703.05698*, 2017.

Trong Duc Nguyen, Anh Tuan Nguyen, and Tien N Nguyen. Mapping api elements for code migration with vector representations. In *2016 IEEE/ACM 38th International Conference on Software Engineering Companion (ICSE-C)*, pp. 756–758. IEEE, 2016.

Trong Duc Nguyen, Anh Tuan Nguyen, Hung Dang Phan, and Tien N Nguyen. Exploring api embedding for api usages and applications. In *2017 IEEE/ACM 39th International Conference on Software Engineering (ICSE)*, pp. 438–449. IEEE, 2017.

Karl Pearson. Liii. on lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, 2(11):559–572, 1901.

Michael Pradel and Koushik Sen. Deepbugs: a learning approach to name-based bug detection. *Proceedings of the ACM on Programming Languages*, 2(OOPSLA):147, 2018.

Kimberly Redmond, Lannan Luo, and Qiang Zeng. A cross-architecture instruction embedding model for natural language processing-inspired binary code analysis. *arXiv preprint arXiv:1812.09652*, 2018.

Bart Theeten, Vandeputte Frederik, and Tom Van Cutsem. Import2vec: learning embeddings for software libraries. In *2019 IEEE/ACM 16th International Conference on Mining Software Repositories (MSR)*. IEEE, 2019.

Ke Wang, Rishabh Singh, and Zhendong Su. Dynamic neural program embedding for program repair. *arXiv preprint arXiv:1711.07163*, 2017.

Martin White, Michele Tufano, Matias Martinez, Martin Monperrus, and Denys Poshyvanyk. Sorting and transforming program repair ingredients via deep learning code similarities. *arXiv preprint arXiv:1707.04742*, 2017.

Xiaojun Xu, Chang Liu, Qian Feng, Heng Yin, Le Song, and Dawn Song. Neural network-based graph embedding for cross-platform binary code similarity detection. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, pp. 363–376. ACM, 2017.

Pengcheng Yin, Graham Neubig, Miltiadis Allamanis, Marc Brockschmidt, and Alexander L Gaunt. Learning to represent edits. *arXiv preprint arXiv:1810.13337*, 2018.

Jinman Zhao, Aws Albarghouthi, Vaibhav Rastogi, Somesh Jha, and Damien Octeau. Neural-augmented static analysis of android communication. In *Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, pp. 342–353. ACM, 2018.

Fei Zuo, Xiaopeng Li, Zhexin Zhang, Patrick Young, Lannan Luo, and Qiang Zeng. Neural machine translation inspired binary code similarity comparison beyond function pairs. *arXiv preprint arXiv:1808.04706*, 2018.