

When Neural Model Meets NL2Code: A Survey

Daoguang Zan^{1,2,*}, Bei Chen³, Fengji Zhang³, Dianjie Lu⁴, Bingchao Wu¹,
Bei Guan⁵, Yongji Wang⁵, Jian-Guang Lou³

¹Cooperative Innovation Center, Institute of Software, Chinese Academy of Sciences

²University of Chinese Academy of Sciences; ³Microsoft Research Asia; ⁴Shandong Normal University

⁵Integrative Innovation Center, Institute of Software, Chinese Academy of Sciences

{daoguang@, bingchao2017, guanbei@, ywang@itechs.}iscas.ac.cn;

{beichen, v-fengjzhang, jlou}@microsoft.com; Ludianjie@sdnu.edu.cn

Abstract

Given a natural language that describes the user's demands, the NL2Code task aims to generate code that addresses the demands. This is a critical but challenging task that mirrors the capabilities of AI-powered programming. The NL2Code task is inherently versatile, diverse and complex. For example, a demand can be described in different languages, in different formats, and at different levels of granularity. This inspired us to do this survey for NL2Code. In this survey, we focus on how does neural network (NN) solves NL2Code. We first propose a comprehensive framework, which is able to cover all studies in this field. Then, we in-depth parse the existing studies into this framework. We create an online website to record the parsing results, which tracks existing and recent NL2Code progress. In addition, we summarize the current challenges of NL2Code as well as its future directions. We hope that this survey can foster the evolution of this field.

1 Introduction

How can novice programmers, even people without any programming background, build an application they desire just by describing their demands? It is a chronically crucial but challenging goal in software engineering (SE), programming language (PL), and artificial intelligence (AI) fields. Assuming we have achieved this goal successfully, it will entail an unprecedented impact on our lives, education, economy, labour market, etc. Let us imagine a scenario together: you have a strong and urgent demand for building an application. Once you describe what the application looks like, the machine can automatically write code to implement your demands, and quickly output an application for you. In the whole process above, instead of hiring a costly outsource to create your application,

* This work was done before October 2022 when the author, Daoguang Zan, was an intern at Microsoft Research Asia.

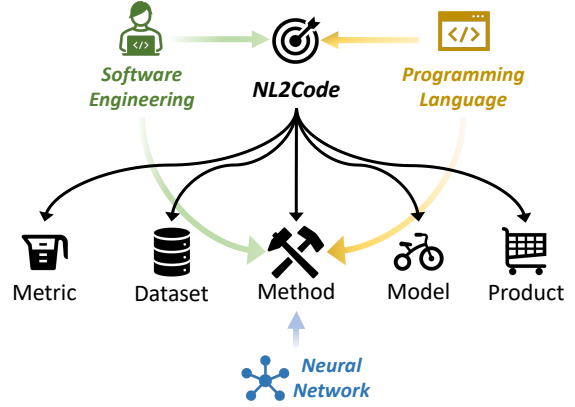


Figure 1: The framework of NL2Code.

you just need to interact with the machine. Thanks to such promising and fascinating prospects, AI-powered programming has drawn widespread interest in industry and academia. Therefore, the natural-language-to-code (NL2Code) task is proposed to evaluate the capability of AI-powered programming, which aims to generate code addressing the given natural language.

In practical programming, how do you describe your demand through NL of NL2Code? As is well known, describing one demand is not unique, which may be extremely diverse. For example, for a demand, NL can be described in various languages like English, French, and Chinese. Also, NL can be formatted as code comment, competition problem, and so on. Moreover, it can be described with various granularities, i.e. you can describe a demand in terms of its background knowledge, functionality, algorithmic process, etc. In addition, you should specify some programming constraints for the demand, such as programming language, code library, time-space complexity, etc. As we can see above, the NL2Code task is inherently highly versatile, diverse and complex, which would be detrimental to the evolution of the research direction if not combed. Such situation motivates us to

do the survey of NL2Code.

On top of the NL2Code task, many papers have been published. In order to do a well-organized and valuable survey on these works, we first summarize a holistic framework for the NL2Code task and then laboriously incorporate the existing papers into the framework. We will elaborate on this below. The NL2Code task was initially proposed in software engineering (SE) and programming language (PL) fields. Recently, many solutions were proposed to solve it. Among these solutions, neural networks (NN) has won wide attention in NL2Code, owing to their impressive code generation capabilities. We thus primarily focus on how NN solves NL2Code in this paper. We find that these NN-based methods usually use NN for modelling based on the insights of SE and PL fields. However, each of NN, SE, and PL is self-contained and covers a lot of studies themselves. So these methods, as a cross-cutting field between NN, SE, and PL, would be disorganized if not combed well. Given this, we provide an in-depth survey of these methods from a cross-cutting perspective. Besides the methods, NL2Code also covers many models, products, datasets and metrics. We also survey them respectively. Based on our survey findings, we summarize a framework for NL2Code, as shown in Figure 1. The framework is able to cover all NL2Code papers, which facilitates the researcher to quickly grasp a new NL2Code paper and locate its innovative ideas. Based on this framework, we parse the existing NL2Code papers. Moreover, we build an online website to display the parsed results, further facilitating the researchers. Consequently, we summarize the ongoing challenges and outline some future research directions.

Our contribution falls into the following aspects:

- To the best of our knowledge, the paper is the first to provide a comprehensive framework for NL2Code in terms of task inputs and outputs (I/Os), methods, datasets, evaluation metrics, models, and products.
- Based on the framework, we made an in-depth survey on a lot of NL2Code papers. Also, an online website was developed to show our survey findings, which can keep up-to-date progress of NL2Code. We also provide an easy-to-use interface for all developers to work together to maintain the website. The website can be found at <https://nl2code.github.io>.

- Based on the framework and our survey findings, we list the existing challenges, as well as draw out future research directions for NL2Code.

2 Framework

We propose a framework for organizing and understanding research on the NL2Code task. It aims to assist researchers in quickly identifying key points in NL2Code papers and to provide a resource for tracking existing and new research through a real-time updated website. This framework first introduces NL2Code (§3), a chronically challenging task in software engineering (SE) and programming language (PL). Under this task, many methods were produced. Our framework particularly focuses on these neural networks (NN)-based methods (§4), which are commonly used as the primary method for the NL2Code task. Therefore, the framework combines the insights from SE and PL to survey the NN-based methods. Besides the methods, our framework also includes the NL2Code’s models (§5), products (§6), datasets (§7), as well as evaluation metrics (§8). The correlation between each of the above modules is plotted in Figure 1, while Figure 2 provides more detailed information for each one. In the following, we will discuss the design of each module of the framework. The NL2Code task aims to produce the code based on a natural language (NL) description of a user’s demand. However, the NL description of a single demand can vary in language, format, and granularity; the inputs for the NL2Code task may include contextual information such as class, file, and repository, as is often the case when humans are programming. Apart from its input, there also exist many programming constraints on NL2Code’s output. For example, the code outputs can vary in language, domain, library, and level. We can see that the task is complicated. So, we conducted a systemic review based on our proposed framework to thoroughly analyze the input and output involved in the NL2Code task.

Our framework focuses on NN-based methods from the NN, SE and PL fields, respectively. For an in-depth survey, we first carefully categorize each of these three fields. Furthermore, we ensure that these categories are relatively well-rounded so that they can cover all papers in each of their fields. The category results of these three fields are placed in Figure 2. Based on these fine-grained categories,

we surveyed a lot of NL2Code methods. We put the findings of our survey on the website. Researchers can quickly learn a new NL2Code method via it.

Our framework also fine-grains each of these modules to assist the survey. For example, the pre-trained models include its supported natural language (NL) and programming language (PL), model architectures, training corpora, model sizes, and other relevant details of these models. Please see Figure 2 for the details of all framework modules. Based on these details we summarize, we also surveyed NL2Code’s models, products, datasets, and metrics separately. Our findings were also placed on the website for everyone to study.

3 What is NL2Code?

In this section, we will introduce the definition of NL2Code and detail its inputs and outputs.

Given a natural language (NL) to describe the user’s demands, NL2Code aims to generate the **Code** that tackles the demand. In this paper, we focus on the neural network (NN) method to face NL2Code. Therefore, this task can be formalized as: $\text{Code} = \mathcal{M}(\text{NL})$, where \mathcal{M} denotes the NN model. On top of the NL2Code definition, 1) how do users describe their demands using NL? Furthermore, 2) what are the user’s programming constraints for **Code** included? We present the answers to the two questions as shown in the upper part of Figure 2.

For the first question, as we all know, for one demand, the descriptions of different humans are different, and even the same human at different moments is different. We thus propose several dimensions to measure how to describe a demand. Concretely, these dimensions include the natural language used in NL (e.g., English, Chinese), the description format (e.g., code comment, competition problem), and the description granularity (e.g., functionality, algorithm process). Moreover, the inputs of the NL2Code task, besides the description of the demand, sometimes have additional information, such as code context, class, file, and repository. This mimics that humans may also refer to such additional information when programming.

For the second question, the programming constraints for **Code**, as a particular user demand, can also include many dimensions. In detail, we can use different programming languages (e.g., Python, Java) and code libraries (e.g., pandas, NumPy) to implement different levels (e.g., line, function) of

code oriented to different domains (e.g., general, data science).

Overall, although NL2Code is well-defined, its input and output have a wide range of expressions.

4 Methods

This section first gives an overview of all existing methods to solve NL2Code, followed by focusing on the neural network (NN)-based ones.

4.1 Overview of Existing Methods

At the beginning of the evolution of NL2Code, a large number of methods based on strong rules or expert systems emerged. For example, domain-specific language (DSL)-guided methods (Moura and Bjørner, 2008; Gulwani, 2010; Jha et al., 2010; Gvero et al., 2013; Evans and Grefenstette, 2018; Pu et al., 2018) require experts to instill in-depth domain knowledge like abstract syntax tree or grammar rule to produce the structured code we want, which drastically limits their flexibility. Besides, probabilistic grammar-based methods (Joshi and Rambow, 2003; Cohn et al., 2010; Allamanis and Sutton, 2014; Bielik et al., 2016a; Raychev et al., 2016) have been proposed to generate high-quality programs. However, such methods rely heavily on pre-defined rules, so they are less scalable as well.

In addition to the above rule-based methods, there are other proposed methods using static language models like n -gram (Hindle et al., 2016; Nguyen et al., 2013; Raychev et al., 2014; Bielik et al., 2016b; Hellendoorn and Devanbu, 2017) and Hidden Markov (Sutskever et al., 2008). Although they do not require domain knowledge and complex rules, they fail to model long-term dependencies, and their vector representations are sparse.

Unlike the above static language models that require explicitly calculating the probability of each token, neural network (NN)-based methods (Schwenk and Gauvain, 2002; Mnih and Hinton, 2007; Maddison and Tarlow, 2014; Allamanis et al., 2015b,a) were proposed for solving NL2Code by first encoding each token into an intermediate vector and then decoding it. The entire process is end-to-end learning and does not rely on any feature engineering. As the computational capacity and the size of training corpus increase, the parameters of neural networks become more and more numerous, which makes the code generation capabilities progressively powerful (Chen et al., 2021; Wang et al., 2021; Ahmad et al., 2021;

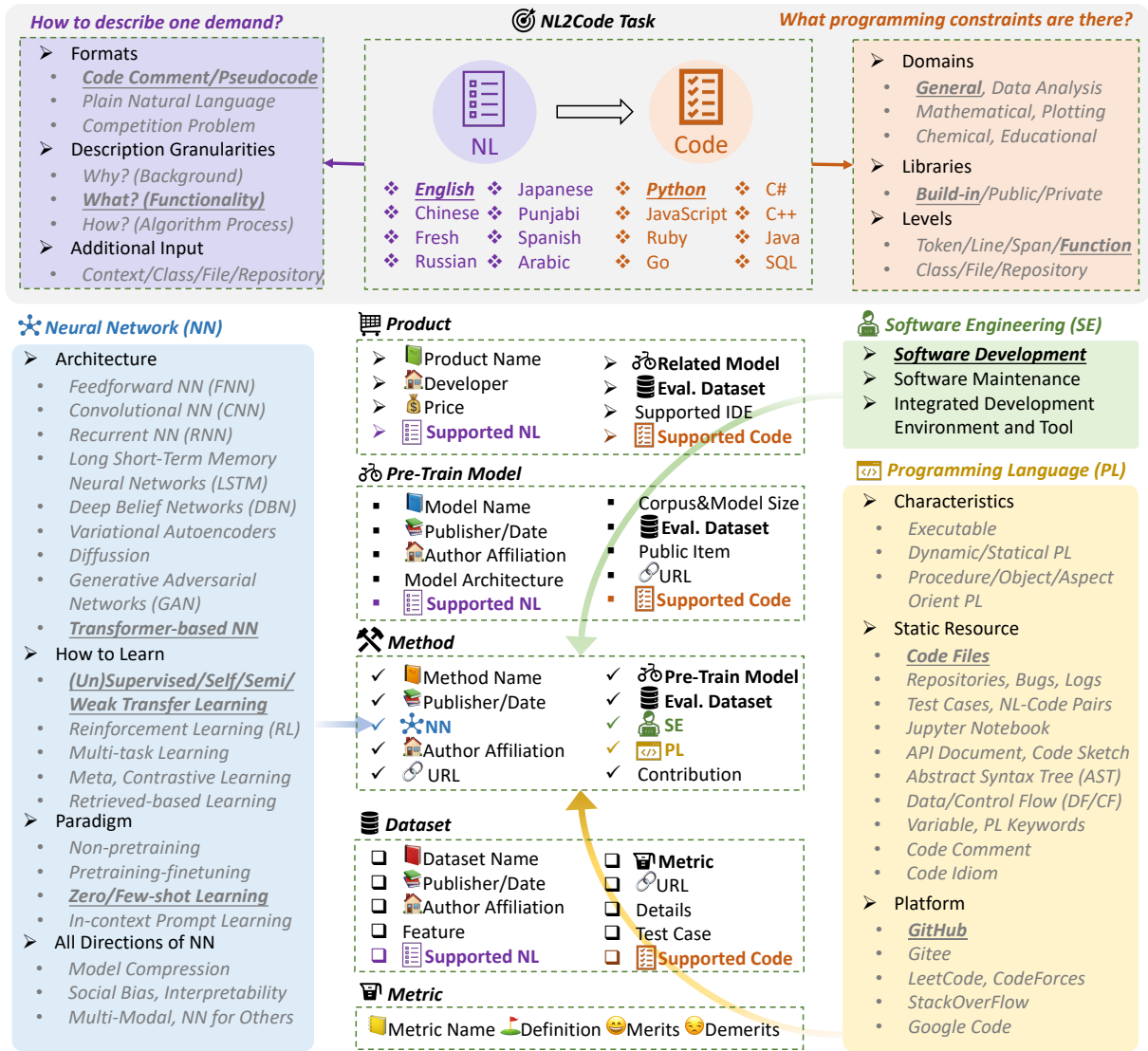


Figure 2: A fine-grained framework of NL2Code. Please refer to Figure 1 for the correlation between each module. Note that the most common settings are marked in underlined bold.

Li et al., 2022a; Nijkamp et al., 2022; Fried et al., 2022).

4.2 NN-based Methods

Compared to other types of methods, the neural network (NN) has exhibited surprising capabilities in the NL2Code task and can be deployed to facilitate coding efficiency in practice. So we specialize in how NN-based methods solve NL2Code. These methods belong to the intersection of NN, SE and PL fields; we thus present them from these perspectives separately.

4.2.1 NN Perspective

As is well known, many NN methods (Devlin et al., 2019; Brown et al., 2020; Lewis et al., 2020; Raffel et al., 2020) were first proposed in the natural

language processing (NLP) field and then applied to a wide variety of other fields such as NL2Code. As we can see, NN methods themselves are field-agnostic. Consequently, we need to first conduct an in-depth survey for NN if we would like to provide closer insights into how NN-based methods solve NL2Code. After a comprehensive survey, we summarized several common perspectives for NN and ensure that these perspectives can cover all NN methods. More details can be found in the bottom left of Figure 2. From these perspectives, we surveyed the existing NL2Code methods and came up with some insights.

- We first surveyed NL2Code from the perspective of NN architecture. Such perspective classifies NL2Code methods into the following as-

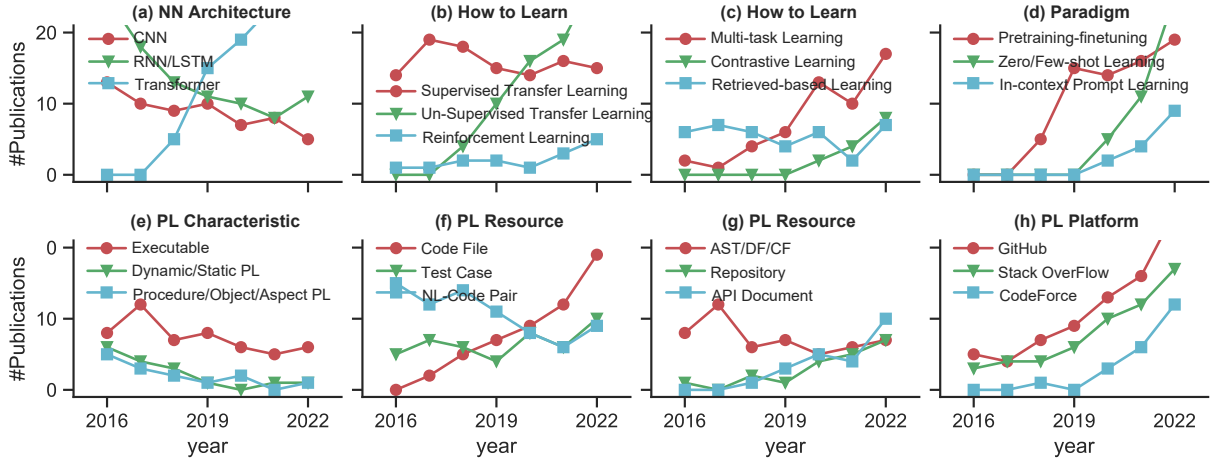


Figure 3: The change trend in the number of NL2Code publications from 2016 to 2022 in terms of multiple dimensions. All the above statistics come from multiple top conferences manually and may suffer from minor bias.

pects: CNN (Liu et al., 2016; Sun et al., 2019), RNN (Iyer et al., 2016; Wan et al., 2018), LSTM (Ling et al., 2016; Eriguchi et al., 2016; Yin and Neubig, 2017; Wei et al., 2019), GAN (Zhu et al., 2019; Wang et al., 2022b), Transformers (Chen et al., 2021; Nijkamp et al., 2022; Fried et al., 2022), etc. Also, we analyzed the trend of the number of NL2Code publications using CNN, RNN/LSTM, and Transformer at the top-tier conferences from 2016 to 2022, and the results are shown in (a) of Figure 3. As we can observe, Transformer was quickly applied to NL2Code since it was proposed in 2017, which is growing exponentially. However, both CNN and RNN/LSTM tend to decrease. This observation suggests that Transformer is superior to other architectures in modelling code.

- Besides the architectural perspective mentioned above, we also surveyed NL2Code from the perspective of how neural networks learn. The learning methods of NL2Code are categorized into supervised learning (Ling et al., 2016; Rabinovich et al., 2017; Iyer et al., 2018), unsupervised learning (Chen et al., 2021; Nijkamp et al., 2022; Fried et al., 2022), reinforcement learning (RL) (Le et al., 2022; Wang et al., 2022b), multi-task learning (MTL) (Wei et al., 2019; Phan et al., 2021; Liu et al., 2020; Ye et al., 2020), contrastive learning (CL) (Jain et al., 2021; Bui et al., 2021; Neelakantan et al., 2022a), retrieved-based learning (Hayati et al., 2018; Lu et al.,

2022; Zan et al., 2022c; Zhou et al., 2022b; Zan et al., 2022a), etc. We also plotted the trends in the number of publications for each of them, as shown in (b) and (c) of Figure 3. We can see the increasing trend of unsupervised learning to solve NL2Code compared to supervised one. Meanwhile, more and more efforts focused on leveraging human feedback (e.g., RL), additional knowledge and additional tasks (e.g., MTL, CL, retrieved-based learning), and so on to raise the reliability of the generated code.

- In addition to the above two perspectives, we also surveyed NL2Code from the perspective of the modelling paradigm. They include non-pretraining (Dam et al., 2016; Lin et al., 2017), pretraining-finetuning (Wei et al., 2019; Xu et al., 2020; Yan et al., 2021), zero-shot learning (Chen et al., 2021; Li et al., 2022a), few-shot learning (Wang et al., 2022c; Madaan et al., 2022), prompt tuning (Wang et al., 2022a), in-context learning (Rajkumar et al., 2022; Trummer, 2022; Poesia et al., 2022), etc. Furthermore, we also analyzed the trend of the number of NL2Code publications for them. The results are shown in (d) of Figure 3. Our findings reveal that, as NN language models become powerful, zero-shot, few-shot, in-context learning receive ever-increasing attention. Moreover, they are gradually surpassing pretraining-finetuning as the dominant paradigm.
- Neural networks elicit many research direc-

tions, such as social bias, model compression, and model interpretability. In NL2Code, these directions also exist and incorporate the characteristics of code; for example, code bias (Mouselinos et al., 2022; Jones and Steinhart, 2022), code model compression (Shi et al., 2022; Zhang et al., 2022c), code model interpretability (Yan and Li, 2022; MacNeil et al., 2022b). We observe that an increasing interest is paid to these NN directions of NL2Code. Such observation also reflects that the goal of NL2Code is gradually focusing on applying it to practical scenarios rather than just toy ones.

4.2.2 SE&PL Perspectives

NN methods solve NL2Code based on the insights of SE and PL. We thus surveyed the NL2Code methods not only from the NN perspective but also from the SE and PL ones. Concretely, for the SE and PL fields, we also did in-depth surveys and sorted out some popular dimensions (bottom right of figure 2). Notably, we ensure that these dimensions are sufficient to outline the whole SE and PL fields. From these dimensions, we surveyed the NL2Code methods as well as drew some insights.

Software engineering (SE) is a discipline that studies the development and maintenance of software using engineering practices. So it mainly includes the following dimensions: software development, software maintenance, integrated development environment (IDE). We all know that the goal of NL2Code aims to improve the efficiency of software development and maintenance by designing a better IDE. Therefore, some methods (Alor-Hernández et al., 2015; Svyatkovskiy et al., 2020; Xu et al., 2022b) have been proposed to achieve this goal.

A programming language (PL) is a formal language used to define computer programs. Moreover, NL2Code lets NN methods learn how to use the PL to build the software users want. So, just like humans, NN methods need to learn what the PL involves if they wish to program correctly. We thus divided PL into multiple dimensions like PL characteristics, PL resources and code hosting platform. Followed by this, we surveyed the NL2Code methods from these dimensions.

- NN methods may use PL characteristics for learning, such as executable (Le et al., 2022), dynamic/statical (Han et al., 2019),

process/object/aspect-oriented (Syriani et al., 2018; Lattner et al., 2021). However, as the parameter numbers of the models exponentially increase, they are decreasingly dependent on these well-designed characteristics, as shown in (e) in Figure 3. This may be because the model can learn the PL characteristics just from a large scale of code files.

- Training a NL2Code model requires a large-scale corpus. The corpus can be code files, test cases, NL-code pairs, repositories, API documentation or even abstract syntax trees (AST), data flow (DF), or control flows (CF). As shown in (f) and (g) of Figure 3, code files, as a vital resource of PL, were increasingly leveraged to train the model in an unsupervised manner in recent years. Also, due to the unsupervised manner does not require manually labelled NL-code pairs, these pairs show a downward trend. Furthermore, recently proposed methods leveraged additional resources to improve the code quality. For example, Li et al. (2022a) and Chen et al. (2022a) used test cases to produce better code; Shrivastava et al. (2022) leveraged the repository information to improve the accuracy of API; Zhou et al. (2022b) and Zan et al. (2022a) used API documentation to generate private APIs; Paik and Wang (2021) and Guo et al. (2020) exploited AST, DF, CF to model more robust vectors.
- The massive code corpus exists on various platforms, for which we combed. Our findings suggest that GitHub, as the largest code hosting platform, contributes extensive code corpus for training NN models (Chen et al., 2021; Nijkamp et al., 2022). Also, we observe from (h) of Figure 3 that Q&A posts in StackOverflow are gradually being mined to train NN models (Orlanski and Gittens, 2021; Drain et al., 2021; Lai et al., 2022; Orlanski et al., 2022). In addition, the corpus on the programming competition website like CodeForces¹ and LeetCode² is used to train models with the expectation that they can automatically generate solutions for any programming problem (Hendrycks et al., 2021; Li et al., 2022a).

¹<https://codeforces.com>

²<https://leetcode.cn>

<i>Model</i>	<i>Arch.</i>	<i>Size</i>	<i>NL</i>	<i>PL</i>	<i>Public</i>
GPT-C (2020)	De	366M	English	Multilingual	None
PyMT5 (2020)	En&De	374M	English	Python	None
CodeGPT (2021)	De	124M	English	Multilingual	Model
PLBART (2021)	En&De	140M	Multilingual	Multilingual	Model, Code, Data
Codex (2021)	De	12M~175B	English	Multilingual	OpenAI API
CodeT5 (2021)	En&De	60M~770M	Multilingual	Multilingual	Model, Code
CodeParrot	De	110M, 1.5B	English	Python	Model, Code, Data
CodeClippy	De	125M, 1.3B	English	Python	Model, Code, Data
AlphaCode (2022a)	En&De	300M~41B	English	Multilingual	None
CodeGen (2022)	De	350M~16.1B	English	Multilingual	Model, Code
PyCodeGPT (2022b)	De	110M	English	Python	Model, Code
PanGu-Coder (2022)	De	317M, 2.6B	English	Python	None
CodeGeeX (b)	De	13B	Multilingual	Multilingual	Model, Code, Data
InCoder (2022)	De	1.3B, 6.7B	Multilingual	Multilingual	Model
FIM (2022)	De	50M~6.9B	English	Python	None
PolyCoder (2022a)	De	160M~2.7B	English	Multilingual	Model, Data
CodeRL (2022)	En&De	770M	English	Python	None
JuPyT5 (2022)	En&De	350M	English	Python	Data

Table 1: The pre-trained models of NL2Code. The two transformer architectures of encoder-decoder and decoder are abbreviated as En&De and De.

5 Pre-trained Models

Since transformer (Vaswani et al., 2017) was proposed in 2017, a large number of pre-trained language models have emerged in various fields. The pre-trained language model refreshes a range of downstream tasks with its powerful zero-shot capability, showing amazing performance. The NL2Code field is no exception, and many pre-trained models have been released for generating code automatically. However, these models vary in size, training corpus, supported natural and programming languages, open-source situations, and so on. So, we surveyed these models in this section from the above aspects. Table 1 lists some results of our survey.

The generative models consist of two main NN architectures: encoder-decoder like T5 (Raffel et al., 2020), BART (Lewis et al., 2020), and decoder like GPT (Radford et al., 2018, 2019; Brown et al., 2020). Afterwards, these architectures were applied to NL2Code; thus, PyMT5 (Clement et al., 2020), CodeT5 (Wang et al., 2021), PLBART (Ahmad et al., 2021), GPT-C (Svyatkovskiy et al., 2020), CodeGPT (Lu et al., 2021), and ERNIE-Code (Chai et al., 2022) were trained. Since these above models are relatively small (million-level parameters), they do not exhibit strong capabilities

in generating code. Codex (Chen et al., 2021), a large language model with billion-level parameters, broke the bottleneck. It shows surprising performance by training on a large-scale and high-quality code corpus. Unfortunately, Codex only provides access via OpenAI’s API ³ and does not release its data, code and model. DeepMind also soon came up with AlphaCode (Li et al., 2022a), which achieves comparable performance to Codex but is likewise unavailable. So, many efforts like CodeClippy ⁴, CodeParrot ⁵ and PolyCode (Xu et al., 2022a) were devoted to training a publicly available model. Although these models are available completely, even including their training data, their performance is still inferior to Codex and AlphaCode. Subsequently, CodeGen (Nijkamp et al., 2022), PyCodeGPT (Zan et al., 2022b), PanGu-Coder (Christopoulou et al., 2022), CodeRL (Le et al., 2022) and other models were proposed. Despite their impressive performance, they only release their models, and neither share their training data. The above is the roadmap for NL2Code pre-trained models. In addition, some models were proposed to address more scenarios in NL2Code. For example, CodeGeeX was trained to support

³<https://beta.openai.com/docs/models/codex>

⁴<https://github.com/CodedotAI/gpt-code-clippy>

⁵<https://huggingface.co/blog/codeparrot>

<i>Product</i>	<i>Model</i>	<i># PL</i>	<i># IDE</i>
Copilot (i)	Codex	-	4
tabnine (l)	-	25	20
Ghostwriter (h)	-	16	1
CodeWhisperer (d)	-	3	4
CodeGenX (c)	GPT-J	1	1
CodeGeeX (b)	CodeGeeX	>20	>12
aiXcoder (a)	-	8	8
FauxPilot (g)	CodeGen	6	-
Diffblue Cover (f)	-	1	-
IntelliCode (j)	-	9	2
cosy (e)	-	1	1
kite (offline) (k)	-	16	16

Table 2: The products of NL2Code. # PL and # IDE denote the number of supported PLs and IDEs.

multiple natural and programming languages; In-Coder (Fried et al., 2022) and FIM (Bavarian et al., 2022) support not only left-to-right code prediction but also infill arbitrary regions of code.

6 Products

Pre-trained models, as the underlying technology, are wrapped into products to improve the programming efficiency of users. In 2021, GitHub and OpenAI jointly launched Copilot⁶, a programming assistance that uses Codex (Chen et al., 2021) to suggest code and entire functions in real-time. It supports multiple programming languages and integrated development environments (IDEs). After many empirical experiments, we observe that it can assist programmers in using a new programming language or even solving bugs. We can see that AI-powered programming assistance has shown great potential. Table 2 lists the currently popular products. Most of the commercial products support more programming languages and IDEs one by one, which mirrors the fact that commercial competition remains intense. In addition, we find that these products are usually launched by large companies. This means that AI-powered programming products have extremely high thresholds in terms of computing capacity, core technology, etc.

Recent studies (Sobania et al., 2022; Pearce et al., 2022; Nguyen and Nadi, 2022; Vaithilingam et al., 2022; Imai, 2022; Ernst and Bavota, 2022; Barke et al., 2022) on the practicality of these products have been done thoroughly. These studies revealed that those products also introduced a series of draw-

backs, even though they could recommend code that users fail to write. For example, a large piece of code generated by a product may harbour some minor bugs, which makes it hard for users to detect. This will probably cause debugging code that may take more time than programming it from scratch. Therefore, one should thoroughly investigate the advantages and disadvantages when using these products. And the product developer should not only improve the performance and speed of the model but also consider how to design a good product, such as how to leverage user feedback for continuous learning.

7 Datasets

To evaluate the code generation capabilities, many datasets have been released. We surveyed them, and the results are listed in Table 3.

At first, some datasets (Oda et al., 2015; Ling et al., 2016; Yin et al., 2018; Lin et al., 2018; Zavershynskyi et al., 2018; Iyer et al., 2018; Husain et al., 2019; Kulal et al., 2019) provided both the training and test sets, where the model needs to be learned on the training set and then verified on the test set. As the model becomes more powerful, it can work well in a zero-shot manner without any additional training set. Therefore, a number of follow-up datasets (Chen et al., 2021; Austin et al., 2021; Hendrycks et al., 2021; Schuster et al., 2021) provided only the test set without the training set.

Earlier proposed datasets such as Django (Oda et al., 2015), CoNaLa (Yin et al., 2018), CON-CODE (Iyer et al., 2018) and CodeSearchNet (Husain et al., 2019) tend to be evaluated by hard metrics like BLEU and accuracy rather than by executing on test cases. This is because earlier models were not able to generate code that passed any of the test cases, so they had to be evaluated with the former. In contrast, recent datasets such as HumanEval (Chen et al., 2021), MBPP (Austin et al., 2021), APPS (Hendrycks et al., 2021), P3 (Schuster et al., 2021), CodeContests (Li et al., 2022a) and DS-1000 (Lai et al., 2022) were usually equipped with multiple test cases per programming problem. These datasets equipped with test cases usually contain relatively few instances compared to those without test cases, since it is extremely challenging to write test cases for programming problems.

We also observe that the natural language of the majority of the existing dataset is English, and the programming language is Python. To test the

⁶<https://github.com/features/copilot>

<i>Benchmark</i>	<i># Num</i>	<i>NL</i>	<i>PL</i>	<i>Metric</i>	<i>TS</i>	<i>Feature</i>
Django (2015)	1,805	Multilingual	Python	BLEU	No	None
MTG (2016)	664	English	Java	Acc, BLEU	No	Card Game
HS (2016)	66	English	Python	Acc, BLEU	No	Card Game
CoNaLa (2018)	500	English	Python, Java	BLEU	No	None
NL2Bash (2018)	9,305	English	Shell	Acc_F^k, Acc_T^k	No	None
NAPS (2018)	485	English	UAST	Acc	Yes	Competition
CONCODE (2018)	2,000	English	Java	EM, BLEU	No	None
CodeSearchNet (2019)	100K	English	Multilingual	MRR	No	None
SPoC (2019)	18,356	English	C++	Acc	Yes	None
HumanEval (2021)	164	English	Python	Pass@k	Yes	Standalone
MBPP (2021)	974	English	Python	Pass@k	Yes	Standalone
APPS (2021)	10,000	English	Python	Pass@k	Part	Competition
P3 (2021)	397	English	Python	Acc	Yes	Puzzle
CodeContests (2022a)	165	English	Multilingual	n@k	Yes	Competition
DS-1000 (2022)	1,000	English	Python	Pass@k	Yes	Data Science
MCoNaLa (2022d)	896	Multilingual	Python, Java	BLEU	No	Multilingual
MBXP (2022)	974	English	Multilingual	Pass@k	Yes	Multilingual
HumanEval-X	164	English	Multilingual	Pass@k	Yes	Multilingual
MultiPL-E (2022)	164	English	Multilingual	Pass@k	Yes	Multilingual
GSM8K-Python (2022)	8,500	English	Python	Pass@k	Yes	Mathematics
PandasEval (2022b)	101	English	Python	Pass@k	Yes	Public Library
NumpyEval (2022b)	101	English	Python	Pass@k	Yes	Public Library
TorchDataEval (2022a)	50	English	Python	Pass@k	Yes	Private Library
MonkeyEval (2022a)	101	English	Python	Pass@k	Yes	Private Library
BeatNumEval (2022a)	101	English	Python	Pass@k	Yes	Private Library
MTPB (2022)	115	English	Python	Pass@k	Yes	Multi-turn
SecurityEval (2022)	130	English	Python	Acc	No	Secure

Table 3: The datasets of NL2Code. # Num and TS denote the number of test sets and the availability of test cases.

multilingual-oriented code generation capability of the model, MCoNaLa (Wang et al., 2022d) extended CoNaLa’s English (Yin et al., 2018) into multiple natural languages like Spanish, Japanese and Russian; MBXP (Athiwaratkun et al., 2022) extended MBPP (Austin et al., 2021) that uses python in multiple programming languages such as Java, JavaScript, TypeScript, C++, and C#; HumanEval (Chen et al., 2021) had also been extended from the python version to HumanEval-X⁷ and MultiPL-E (Cassano et al., 2022) in multiple programming languages.

Another trend in these datasets is that they are becoming more and more relevant to real-world scenarios rather than just toy ones. HumanEval (Chen et al., 2021) and MBPP (Austin et al., 2021) were hand-written and ensured that they were not seen during training. Such setting is in line with actual practice. In addition, APPS (Hendrycks et al., 2021) and CodeContests (Li et al., 2022a) were pro-

posed to evaluate actual programming competition scenarios. Subsequently, domain-oriented datasets, such as DS-1000 (Lai et al., 2022) for data science and GSM8K-Python (Chowdhery et al., 2022) and MathQA-Python (Austin et al., 2021) for mathematics, were proposed. Moreover, PandasEval and NumpyEval (Zan et al., 2022b) were proposed to evaluate library-oriented code generation, as we use not only built-in libraries for coding in practice but also often third-party libraries. In addition to the built-in and public libraries mentioned above, we also encounter private ones in our daily programming. So Zan et al. (2022a) proposed three private library datasets named TorchDataEval, MonkeyEval and BeatNumEval.

All datasets mentioned above input each programming problem into the model in one step, while MTPB (Nijkamp et al., 2022) explored whether the programming problem can be factorized into multiple sub-problems to feed the model in turn. This is also called multi-turn program syn-

⁷<http://keg.cs.tsinghua.edu.cn/codegeex>

thesis. While most datasets focus on evaluating the accuracy of generated code, their security is likewise crucial in practical programming. So SecurityEval (Siddiq and Santos, 2022) was proposed to evaluate the code security.

Overall, the NL2Code datasets support diverse settings, providing valuable resources for researchers to evaluate and improve their methods, models and products.

8 Evaluation Metrics

Evaluation metrics can measure the quality of the generated code, which mirrors the NL2Code performance (Evtikhiev et al., 2022; Takaichi et al., 2022). In this section, we surveyed these evaluation metrics and classified them into four categories.

8.1 Human Evaluation

We all know that manual evaluation of the generated code is the gold standard (Yang et al., 2022). Such way enables us to evaluate the generated code all round, including its accuracy, syntactic, semantic, complexity, security. However, it is also impractical to manually evaluate each programming problem, as doing so is very time-consuming and costly. Therefore, this way be only suitable for small-scale empirical experiments, not large-scale ones.

8.2 Off-the-shelf Metrics in Machine Translation

Technically speaking, all existing metrics in machine translation can also be employed to evaluate code by regarding code as natural language. These metrics include BLEU, ROUGE, METEOR, chrF, EM, and EM@N. We introduce their definitions separately below.

- BLEU (Papineni et al., 2002), referring to BiLingual Evaluation Understudy, measures the similarity between the predicted values and the reference ones. In detail, it evaluates the quality of the predicted output by counting the number of occurrences in the reference value for each gram of the output. Usually, a higher BLEU score means a better quality.
- ROUGE (Lin, 2004) refers to Recall-Oriented Understudy for Gisting Evaluation, which can also measure the similarity between the predicted values and the reference ones by counting their overlap. A higher ROUGE indicates that the predicted value is more reliable.

- METEOR (Denkowski and Lavie, 2014) is known as Metric for Evaluation of Translation with Explicit Ordering. Similar to BLEU and ROUGE, it is also a metric that evaluates the similarity between two sentences. It considers the order information between different grams, while BLEU and ROUGE do not. It also solves BLEU's shortcoming of not considering the recall rate. However, it is still sensitive to the sentence length, as BLEU is.
- chrF (Popović, 2015) uses the character n-gram F-score for evaluating the similarity. Unlike the above metrics that consider word-level information, it considers character-level one. Such a feature makes it more suitable for scenarios with complex language structures.
- Exact match (EM) measures the complete consistency between the predicted value and the reference one. That is, EM is set to 1 if they are exactly equal and 0 otherwise. While the above is sampled one prediction for each instance, we can also sample N ones. EM@N (Rajpurkar et al., 2016; Clement et al., 2021) means that it is set to 1 if at least one of these N predictions is correct. The accuracy of the dataset refers to the percentage of instances with EM = 1.

8.3 Code-oriented Modified Metrics

Although machine translation metrics like BLEU and ROUGE have been widely used for evaluating code, these metrics do not evaluate code well, as they only consider the characteristics of the natural language rather than the programming language. Therefore, Tran et al. (2019) proposed a new metric, RUBY, which considers the program dependency graphs and abstract syntax tree. Besides, Ren et al. (2020) also proposed a composite metric named CodeBLEU, which not only maximizes the advantages of BLEU but also considers code syntax and semantics by AST and data-flow.

8.4 Execution-based Metrics

For one demand, there exist various solutions. Under such situation, the machine translation metrics and its modified versions cannot correctly evaluate the generated code. Many execution-based metrics like pass@k, n@k, success rate and test case average were thus proposed. We will elaborate on their definitions in the following.

- Pass@ k was proposed by [Chen et al. \(2021\)](#). For each instance, we sample n candidate programs and then randomly pick k of them. If any of the k ones passes the given test cases, the instance can be regarded as solved. So pass@ k is the proportion of solved instances in the dataset.
- $n@k$ ([Li et al., 2022a](#)) is similar to pass@ k . It selects k programs from n ones by a specific strategy, not at random.
- Test case average was proposed by [Hendrycks et al. \(2021\)](#). Literally, it counts the average percentage of test cases that pass.

9 An Online Website

To keep tracking the progress of NL2Code, we developed an on-line real-time update website, which can be found at <https://nl2code.github.io>. This website is designed precisely following our proposed framework in Section 2. Moreover, as much as possible, we add all existing and latest papers to this website. So, a researcher can quickly learn some key points of a new NL2Code paper via this website. This would be beneficial to the NL2Code field.

We built this website using the jekyll technology⁸. The website supports several features. First, it supports everyone in adding a new NL2Code paper by pulling requests on GitHub. Second, it supports specifying the tag of this paper when adding a new paper. These tags can be one or more of methods, models, products, datasets or evaluation metrics. In addition, it supports adding the published date and publisher for each paper. Moreover, it also supports the fuzzy search of content, which facilitates researchers to find the papers they want.

10 Challenges and Opportunities

In this section, we will present the challenges NL2Code is encountering and give our insights on its future opportunities.

Challenge 1: how should users describe their demands? As described in Section 3, we have shown that describing one demand is versatile, diverse, and complex, which will aggravate the instability of the code generation model. So how should one demand be described to maximize the performance of the model? This is an meaningful

research question. To this end, one potential idea is to first use ChatGPT⁹ or other dialogue models to sufficiently communicate the demands with the user via multi-round interactions, and then input the post-communication demands into the model to return the final code. This idea may be a possible future research direction because a prerequisite for an excellent model is that it can sufficiently understand the user’s demands.

Challenge 2: how to enable models to learn like humans? As we all know, the current powerful code generation models are trained via large-scale, high-quality code files without using any other information. However, this process differs from humans while humans can learn a wide variety of knowledge, including low-quality noise code files. Therefore, it is challenging to empower the model to accept a wide range of knowledge. Moreover, there exists another difference between models and humans. When models are trained directly on code files, humans first learn the programming basics and then learn and practice with code files. Therefore, how to enable models like humans learn progressively from simple to complex? This is also a research question to be explored.

Challenge 3: how to edit the knowledge inside a large language model (LLM)? Knowledge is often changing. For example, some APIs may be deprecated, and some are newly generated as the version is updated. So, how to dynamically update or edit such knowledge inside LLM? If we fine-tune alone on such knowledge, it would be disruptive to previous knowledge. In contrast, if we re-train on the whole knowledge, and this way would be uneconomical. The existing works ([De Cao et al., 2021](#); [Meng et al., 2022](#)) edit the knowledge of the model by modifying its parameters, but these works fail to achieve remarkable results. Also, several works ([Prasad et al., 2022](#); [Chowdhury et al., 2022](#)) directly design prompts to guide the knowledge of the model, and they do not work well. So, to efficiently edit the knowledge inside LLM, we still need to explore more solutions in the future.

Challenge 4: how to enable the model to know which programming problems it can answer correctly? which ones can not? Once you enter one programming problem, the model will definitely return one solution to solve it, even if it has no answer. Obviously, this is unreasonable. Therefore, the existing models do not have the capability to

⁸<https://jekyllrb.com>

⁹<https://chat.openai.com>

self-judgment. Recently, some methods (Ouyang et al., 2022) have been based on reinforcement learning to improve its self-judgment capability by leveraging user feedback. However, collecting a lot of high-quality user feedback is tricky. We can see that the self-judgment capability of the model is vital. Therefore, we should explore this direction in the future as well.

Challenge 5: how to interpret the generated code? It would be extremely user-friendly if the code is generated along with the corresponding explanation. However, NN models are black boxes and cannot do this inherently. So, we use specific prompts to let the model produce the explanation, such as “please comment on each line of generated code”. Unfortunately, this way often fails to generate constructive explanations. In addition, we also infill code comments for the generated code by InCoder (Fried et al., 2022). Such a way also fails to work well. Therefore, it is a very challenging and meaningful direction we need to explore in the future.

Challenge 6: how does the language model combine with the knowledge graph in NL2Code? In natural language processing, language models combined with knowledge graphs have been widely used and have demonstrated their superiority (Fei et al., 2021; Yu et al., 2022). But this approach has yet to be explored in the code. To do this, we list below some possible feasible methods. We first need to define what the knowledge graph can store about the code. It may store the programming language’s grammar, code, repository structure, etc. Followed by this, we need to explore how the code language model can be combined with the knowledge graph. Potential solutions are to train them jointly or to train language models based on the results of retrieval and reasoning in the knowledge graph. Overall, there exist many possibilities that are worth exploring.

Challenge 7: how does the model generalize to other programming languages? It would be fascinating if we trained a model on one programming language and it could be quickly migrated to another one. However, the current models are not able to do this well, even with over a hundred billion parameters (Chowdhery et al., 2022). The primary reason is that the model fails to abstract the high-level knowledge of the programming languages and thereby fails to learn their commonalities. In conclusion, the model does not yet have the capa-

bility to migrate as quickly as humans. Therefore, improving the model’s generalizability across multiple programming languages is a challenge and opportunity.

Challenge 8: how to accelerate training and inference speed? Either for code or natural language, only if the model has a sufficiently large number of parameters, it can show the emergency capability. However, a larger number of parameters makes the training and inference of the model slower. Recent works (Polino et al., 2018; Gou et al., 2021) leveraged techniques such as distillation, quantification, or pruning to tackle it, but that still leaves room for improvement. In order to better deploy models online, we still need to explore more efficient lightweight methods in the future to improve their training and inference speed.

Challenge 9: is it reasonable to model code as plain text? The transformer-based model refreshes a wide variety of downstream tasks, including NL2Code. It models code as plain text outright. As we all know, code and text are fundamentally different. So, is it reasonable to model code as plain text? Furthermore, is there a better method to model the code? Although these scientific questions are still unanswered, they are well worth exploring.

Challenge 10: how to better evaluate the generated code? Our findings show that execution-based metrics can quickly and accurately evaluate whether the generated code is correct. However, the generated codes are not always runnable, so they cannot be evaluated using test cases for execution. For those unexecutable codes, they must be evaluated only using human evaluation or hard metrics like CodeBLEU. Nevertheless, as shown in Section 8, human evaluation is not suitable for large-scale evaluation, and hard metrics do not support precise evaluation. As we can see, the code evaluation is still imperfect, which needs further exploration.

Challenge 11: how does the model address long code? Some scenarios, such as file- or repository-level code generation, need to model longer code. But the present model can only support fixed-length codes due to its computational complexity. So, Clement et al. (2021) proposed eWASH that leverages the syntax hierarchy of code to extend the modelling window. Although eWASH can handle the file level, it still fails to handle the repository level. Thus, we should make more ef-

forts to explore this direction for a larger scope of modelling.

Challenge 12: what scenarios can LLMs of code solve? While LLMs of code like Codex have many challenges, they still exhibit amazing performance, where they are trained on a large-scale code corpus and can store rich knowledge. So, these models are capable of solving a wide range of scenarios. Consequently, exploring more downstream scenarios for LLMs has become an interesting scientific question. We will next list some common scenarios.

- While LLMs default to the general domain for generating code, they can also be applied to other domains. For example, some works (Davis, 2022; Lewis, 2022; Sarsa et al., 2022; MacNeil et al., 2022a; Hocky and White, 2022; Zhang et al., 2022b) have applied LLMs to education, which can automatically generate exam questions and even review the code students write. Also, LLMs can be applied to chemistry (Krenn et al., 2022; Hocky and White, 2022), mathematics (Drori et al., 2022; Drori and Verma, 2021; Tang et al., 2021), numerical computation (Zan et al., 2022b; Chandel et al., 2022), and other domains. Practically, LLMs of code can solve much more than the above domains. Therefore, we can explore more potential domains in future research.
- LLMs are trained on public code corpus, so they excel at the seen code libraries but not the unseen private ones. So, Zan et al. (2022a) empowered LLMs to tackle private libraries using retriever-based techniques. However, this approach assumes that public libraries are prohibited when using private ones. Such an assumption is unreasonable since we often use both public and private libraries in practical programming. Furthermore, how do LLMs solve this scenario? It has not been explored yet. As we can see, the code library perspective can induce many meaningful scenarios. We should therefore explore how to address these scenarios with LLMs in future.
- Code levels include a token, line, span, function, class, file, and repository. Currently, LLMs have the capability to generate a token, line, span, or even a function or class. But LLMs fail to generate a file or repository at

once. Therefore, exploring how LLMs generate file- or repository-level code is necessary. In practical programming, users may prefer different code levels at different moments. We can see that the code level demands of users are changing, so how should LLMs solve this problem? This is also challenging research. In fact, the code level perspective can also induce many practical scenarios, and these scenarios may not be well solved by LLMs now, so this points to the future direction for NL2Code.

- LLMs like Codex not only can be applied to NL2Code scenarios, but also to other scenarios that can be formalized by the code style. For example, Rajkumar et al. (2022), Trummer (2022) and Poesia et al. (2022) formatted the text2sql task into the code style for the LLMs to solve, and the results proved their effectiveness. Some works (Joshi et al., 2022; Prenner and Robbes, 2021) also transformed the code repair task into the code style that LLMs can recognize, allowing LLMs to repair the code automatically. LLMs can also be applied to the reasoning tasks by converting them to code style, which can model the reasoning process as code logic (Chen et al., 2022b; Wei et al., 2022; Zhou et al., 2022a). Besides, LLMs can seamlessly solve bug location and fixing (Zhang et al., 2022a), code review (Li et al., 2022b), code retrieval (Nee-lakantan et al., 2022b) and other scenarios. Overall, LLM may yield a decent or even surprising result as long as the task can be described in code style.

11 Conclusion

This paper provides a framework for NL2Code, which includes the task I/Os, methods, models, products, datasets, and evaluation metrics. We ensure that this framework can cover all NL2Code papers. So, we added as many NL2Code papers as possible to this framework to guide our survey. We built an online website to record our survey findings. Meanwhile, we summarized the ongoing challenges and future opportunities for NL2Code based on our framework and survey findings. We hope that this paper will facilitate the research work of researchers in NL2Code. Furthermore, we also hope everyone can contribute to the website by adding up-to-date NL2Code papers.

References

- 2014k. Kite. <https://www.kite.com>.
- 2018a. aiXcoder. <https://aixcoder.com>.
- 2018l. TabNine. <https://www.tabnine.com>.
- 2019j. IntelliCode. <https://github.com/MicrosoftDocs/intellicode>.
- 2020f. Diffblue Cover. <https://www.diffblue.com>.
- 2021i. GitHub Copilot. <https://github.com/features/copilot>.
- 2022b. CodeGeeX: A Multilingual Code Generation Model. <https://keg.cs.tsinghua.edu.cn/codegeex>.
- 2022c. CodeGenX. <https://docs.deepgenx.com>.
- 2022d. CodeWhisperer. <https://aws.amazon.com/cn/codewhisperer>.
- 2022e. Cosy. <https://github.com/alibaba-cloud-toolkit/cosy>.
- 2022g. FauxPilot. <https://github.com/moyix/fauxpilot>.
- 2022h. Ghostwriter. <https://replit.com/site/ghostwriter>.
- Wasi Ahmad, Saikat Chakraborty, Baishakhi Ray, and Kai-Wei Chang. 2021. Unified pre-training for program understanding and generation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2655–2668.
- Miltiadis Allamanis, Earl T Barr, Christian Bird, and Charles Sutton. 2015a. Suggesting accurate method and class names. In *Proceedings of the 2015 10th joint meeting on foundations of software engineering*, pages 38–49.
- Miltiadis Allamanis and Charles Sutton. 2014. Mining idioms from source code. In *Proceedings of the 22nd acm sigsoft international symposium on foundations of software engineering*, pages 472–483.
- Miltos Allamanis, Daniel Tarlow, Andrew Gordon, and Yi Wei. 2015b. Bimodal modelling of source code and natural language. In *International conference on machine learning*, pages 2123–2132. PMLR.
- Giner Alor-Hernández, Viviana Yarel Rosales-Morales, Jorge Luis García Alcaráz, Ramón Zatarain Cabada, and María Lucía Barrón Estrada. 2015. An analysis of tools for automatic software development and automatic code generation. *Revista Facultad de Ingeniería Universidad de Antioquia*, pages 75–87.
- Ben Athiwaratkun, Sanjay Krishna Gouda, Zijian Wang, Xiaopeng Li, Yuchen Tian, Ming Tan, Wasi Uddin Ahmad, Shiqi Wang, Qing Sun, Mingyue Shang, et al. 2022. Multi-lingual evaluation of code generation models. *arXiv preprint arXiv:2210.14868*.
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. 2021. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*.
- Shraddha Barke, Michael B James, and Nadia Polikarpova. 2022. Grounded Copilot: How programmers interact with code-generating models. *arXiv preprint arXiv:2206.15000*.
- Mohammad Bavarian, Heewoo Jun, Nikolas Tezak, John Schulman, Christine McLeavey, Jerry Tworek, and Mark Chen. 2022. Efficient training of language models to fill in the middle. *arXiv preprint arXiv:2207.14255*.
- Pavol Bielik, Veselin Raychev, and Martin Vechev. 2016a. Phog: probabilistic model for code. In *International Conference on Machine Learning*, pages 2933–2942. PMLR.
- Pavol Bielik, Veselin Raychev, and Martin Vechev. 2016b. Program synthesis for character level language modeling. In *ICLR*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Nghi DQ Bui, Yijun Yu, and Lingxiao Jiang. 2021. Self-supervised contrastive learning for code retrieval and summarization via semantic-preserving transformations. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 511–521.
- Federico Cassano, John Gouwar, Daniel Nguyen, Sydney Nguyen, Luna Phipps-Costin, Donald Pinckney, Ming Ho Yee, Yangtian Zi, Carolyn Jane Anderson, Molly Q Feldman, et al. 2022. A scalable and extensible approach to benchmarking nl2code for 18 programming languages. *arXiv preprint arXiv:2208.08227*.
- Yekun Chai, Shuohuan Wang, Chao Pang, Yu Sun, Hao Tian, and Hua Wu. 2022. ERNIE-Code: Beyond english-centric cross-lingual pretraining for programming languages. *arXiv preprint arXiv:2212.06742*.
- Shubham Chandel, Colin B Clement, Guillermo Serato, and Neel Sundaresan. 2022. Training and evaluating a jupyter notebook data science assistant. *arXiv preprint arXiv:2201.12901*.

- Bei Chen, Fengji Zhang, Anh Nguyen, Daoguang Zan, Zeqi Lin, Jian-Guang Lou, and Weizhu Chen. 2022a. Codet: Code generation with generated tests. *arXiv preprint arXiv:2207.10397*.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W Cohen. 2022b. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *arXiv preprint arXiv:2211.12588*.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. PaLM: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*.
- Jishnu Ray Chowdhury, Yong Zhuang, and Shuyi Wang. 2022. Novelty controlled paraphrase generation with retrieval augmented conditional prompt tuning. *AAAI*.
- Fenia Christopoulou, Gerasimos Lampouras, Milan Gritta, Guchun Zhang, Yinpeng Guo, Zhongqi Li, Qi Zhang, Meng Xiao, Bo Shen, Lin Li, et al. 2022. Pangu-coder: Program synthesis with function-level language modeling. *arXiv preprint arXiv:2207.11280*.
- Colin Clement, Dawn Drain, Jonathan Timcheck, Alexey Svyatkovskiy, and Neel Sundaresan. 2020. PyMT5: multi-mode translation of natural language and python code with transformers. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9052–9065.
- Colin B Clement, Shuai Lu, Xiaoyu Liu, Michele Tufano, Dawn Drain, Nan Duan, Neel Sundaresan, and Alexey Svyatkovskiy. 2021. Long-range modeling of source code files with ewash: Extended window access by syntax hierarchy. *arXiv preprint arXiv:2109.08780*.
- Trevor Cohn, Phil Blunsom, and Sharon Goldwater. 2010. Inducing tree-substitution grammars. *The Journal of Machine Learning Research*, 11:3053–3096.
- Hoa Khanh Dam, Truyen Tran, and Trang Pham. 2016. A deep language model for software code. *arXiv preprint arXiv:1608.02715*.
- Ernest Davis. 2022. Limits of an ai program for solving college math problems. *arXiv preprint arXiv:2208.06906*.
- Nicola De Cao, Wilker Aziz, and Ivan Titov. 2021. Editing factual knowledge in language models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6491–6506.
- Michael Denkowski and Alon Lavie. 2014. Meteor universal: Language specific translation evaluation for any target language. In *Proceedings of the ninth workshop on statistical machine translation*, pages 376–380.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186.
- Dawn Drain, Changran Hu, Chen Wu, Mikhail Breslav, and Neel Sundaresan. 2021. Generating code with the help of retrieved template functions and stack overflow answers. *arXiv preprint arXiv:2104.05310*.
- Iddo Drori and Nakul Verma. 2021. Solving linear algebra by program synthesis. *arXiv preprint arXiv:2111.08171*.
- Iddo Drori, Sarah Zhang, Reece Shuttleworth, Leonard Tang, Albert Lu, Elizabeth Ke, Kevin Liu, Linda Chen, Sunny Tran, Newman Cheng, et al. 2022. A neural network solves, explains, and generates university math problems by program synthesis and few-shot learning at human level. *Proceedings of the National Academy of Sciences*, 119(32):e2123433119.
- Akiko Eriguchi, Kazuma Hashimoto, and Yoshimasa Tsuruoka. 2016. Tree-to-sequence attentional neural machine translation. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 823–833.
- Neil A Ernst and Gabriele Bavota. 2022. Ai-driven development is here: Should you worry? *IEEE Software*, 39(2):106–110.
- Richard Evans and Edward Grefenstette. 2018. Learning explanatory rules from noisy data. *Journal of Artificial Intelligence Research*, 61:1–64.
- Mikhail Evtikhiev, Egor Bogomolov, Yaroslav Sokolov, and Timofey Bryksin. 2022. Out of the bleu: how should we assess quality of the code generation models? *arXiv preprint arXiv:2208.03133*.
- Hao Fei, Yafeng Ren, Yue Zhang, Donghong Ji, and Xiaohui Liang. 2021. Enriching contextualized language model from knowledge graph for biomedical information extraction. *Briefings in bioinformatics*, 22(3):bbaa110.

- Daniel Fried, Armen Aghajanyan, Jessy Lin, Sida Wang, Eric Wallace, Freda Shi, Ruiqi Zhong, Wen-tau Yih, Luke Zettlemoyer, and Mike Lewis. 2022. Incoder: A generative model for code infilling and synthesis. *arXiv preprint arXiv:2204.05999*.
- Jianping Gou, Baosheng Yu, Stephen J Maybank, and Dacheng Tao. 2021. Knowledge distillation: A survey. *International Journal of Computer Vision*, 129(6):1789–1819.
- Sumit Gulwani. 2010. Dimensions in program synthesis. In *Proceedings of the 12th international ACM SIGPLAN symposium on Principles and practice of declarative programming*, pages 13–24.
- Daya Guo, Shuo Ren, Shuai Lu, Zhangyin Feng, Duyu Tang, Shujie Liu, Long Zhou, Nan Duan, Alexey Svyatkovskiy, Shengyu Fu, et al. 2020. GraphCodeBERT: Pre-training code representations with data flow. *arXiv preprint arXiv:2009.08366*.
- Tihomir Gvero, Viktor Kuncak, Ivan Kuraj, and Ruzica Piskac. 2013. Complete completion using types and weights. In *Proceedings of the 34th ACM SIGPLAN conference on Programming language design and implementation*, pages 27–38.
- HyungSeok Han, DongHyeon Oh, and Sang Kil Cha. 2019. CodeAlchemist: Semantics-aware code generation to find vulnerabilities in javascript engines. In *NDSS*.
- Shirley Anugrah Hayati, Raphael Olivier, Pravalika Avvaru, Pengcheng Yin, Anthony Tomasic, and Graham Neubig. 2018. Retrieval-based neural code generation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 925–930.
- Vincent J Hellendoorn and Premkumar Devanbu. 2017. Are deep neural networks the best choice for modeling source code? In *Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering*, pages 763–773.
- Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Song, et al. 2021. Measuring coding challenge competence with apps. *NeurIPS*.
- Abram Hindle, Earl T Barr, Mark Gabel, Zhendong Su, and Premkumar Devanbu. 2016. On the naturalness of software. *Communications of the ACM*, 59(5):122–131.
- Glen M Hocky and Andrew D White. 2022. Natural language processing models that automate programming will transform chemistry research and teaching. *Digital discovery*, 1(2):79–83.
- Hamel Husain, Ho-Hsiang Wu, Tiferet Gazit, Miltiadis Allamanis, and Marc Brockschmidt. 2019. CodeSearchNet Challenge: Evaluating the state of semantic code search. *arXiv preprint arXiv:1909.09436*.
- Saki Imai. 2022. Is github copilot a substitute for human pair-programming? an empirical study. In *Proceedings of the ACM/IEEE 44th International Conference on Software Engineering: Companion Proceedings*, pages 319–321.
- Srinivasan Iyer, Ioannis Konstas, Alvin Cheung, and Luke Zettlemoyer. 2016. Summarizing source code using a neural attention model. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2073–2083.
- Srinivasan Iyer, Ioannis Konstas, Alvin Cheung, and Luke Zettlemoyer. 2018. Mapping language to code in programmatic context. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1643–1652.
- Paras Jain, Ajay Jain, Tianjun Zhang, Pieter Abbeel, Joseph Gonzalez, and Ion Stoica. 2021. Contrastive code representation learning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5954–5971.
- Susmit Jha, Sumit Gulwani, Sanjit A Seshia, and Ashish Tiwari. 2010. Oracle-guided component-based program synthesis. In *2010 ACM/IEEE 32nd International Conference on Software Engineering*, volume 1, pages 215–224. IEEE.
- Erik Jones and Jacob Steinhardt. 2022. Capturing failures of large language models via human cognitive biases. *arXiv preprint arXiv:2202.12299*.
- Aravind Joshi and Owen Rambow. 2003. A formalism for dependency grammar based on tree adjoining grammar. In *Proceedings of the Conference on Meaning-text Theory*, pages 207–216. MTT Paris, France.
- Harshit Joshi, José Cambronero, Sumit Gulwani, Vu Le, Ivan Radicek, and Gust Verbruggen. 2022. Repair is nearly generation: Multilingual program repair with llms. *arXiv preprint arXiv:2208.11640*.
- Mario Krenn, Qianxiang Ai, Senja Barthel, Nessa Carson, Angelo Frei, Nathan C Frey, Pascal Friederich, Théophile Gaudin, Alberto Alexander Gayle, Kevin Maik Jablonka, et al. 2022. Selfies and the future of molecular string representations. *Patterns*, 3(10):100588.
- Sumith Kulal, Panupong Pasupat, Kartik Chandra, Mina Lee, Oded Padon, Alex Aiken, and Percy S Liang. 2019. SPoC: Search-based pseudocode to code. *Advances in Neural Information Processing Systems*, 32.
- Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang, Ruiqi Zhong, Luke Zettlemoyer, Scott Wen-tau Yih, Daniel Fried, Sida Wang, and Tao Yu. 2022. DS-1000: A natural and reliable benchmark for data science code generation. *arXiv preprint arXiv:2211.11501*.

- Chris Lattner, Mehdi Amini, Uday Bondhugula, Albert Cohen, Andy Davis, Jacques Pienaar, River Riddle, Tatiana Shpeisman, Nicolas Vasilache, and Oleksandr Zinenko. 2021. Mlir: Scaling compiler infrastructure for domain specific computation. In *2021 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*, pages 2–14. IEEE.
- Hung Le, Yue Wang, Akhilesh Deepak Gotmare, Silvio Savarese, and Steven CH Hoi. 2022. CodeRL: Mastering code generation through pretrained models and deep reinforcement learning. *arXiv preprint arXiv:2207.01780*.
- Clayton Lewis. 2022. Automatic programming and education. In *Proceedings of the 6th International Conference on the Art, Science, and Engineering of Programming*, pages 70–80.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880.
- Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. 2022a. Competition-level code generation with alphacode. *arXiv preprint arXiv:2203.07814*.
- Zhiyu Li, Shuai Lu, Daya Guo, Nan Duan, Shailesh Jannu, Grant Jenks, Deep Majumder, Jared Green, Alexey Svyatkovskiy, Shengyu Fu, et al. 2022b. Automating code review activities by large-scale pre-training. In *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, pages 1035–1047.
- Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Xi Victoria Lin, Chenglong Wang, Deric Pang, Kevin Vu, and Michael D Ernst. 2017. Program synthesis from natural language using recurrent neural networks. *University of Washington Department of Computer Science and Engineering, Seattle, WA, USA, Tech. Rep. UW-CSE-17-03-01*.
- Xi Victoria Lin, Chenglong Wang, Luke Zettlemoyer, and Michael D Ernst. 2018. NL2Bash: A corpus and semantic parser for natural language interface to the linux operating system. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*.
- Wang Ling, Phil Blunsom, Edward Grefenstette, Karl Moritz Hermann, Tomáš Kočíšký, Fumin Wang, and Andrew Senior. 2016. Latent predictor networks for code generation. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*, pages 599–609.
- Fang Liu, Ge Li, Yunfei Zhao, and Zhi Jin. 2020. Multi-task learning based pre-trained language model for code completion. In *Proceedings of the 35th IEEE/ACM International Conference on Automated Software Engineering*, pages 473–485.
- Zhiqiang Liu, Yong Dou, Jingfei Jiang, and Jinwei Xu. 2016. Automatic code generation of convolutional neural networks in fpga implementation. In *2016 International conference on field-programmable technology (FPT)*, pages 61–68. IEEE.
- Shuai Lu, Nan Duan, Hojae Han, Daya Guo, Seungwon Hwang, and Alexey Svyatkovskiy. 2022. Reacc: A retrieval-augmented code completion framework. *arXiv preprint arXiv:2203.07722*.
- Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin Clement, Dawn Drain, Daxin Jiang, Duyu Tang, et al. 2021. CodeXGLUE: A machine learning benchmark dataset for code understanding and generation. *arXiv preprint arXiv:2102.04664*.
- Stephen MacNeil, Andrew Tran, Juho Leinonen, Paul Denny, Joanne Kim, Arto Hellas, Seth Bernstein, and Sami Sarsa. 2022a. Automatically generating cs learning materials with large language models. *arXiv preprint arXiv:2212.05113*.
- Stephen MacNeil, Andrew Tran, Dan Mogil, Seth Bernstein, Erin Ross, and Ziheng Huang. 2022b. Generating diverse code explanations using the gpt-3 large language model. In *Proceedings of the 2022 ACM Conference on International Computing Education Research-Volume 2*, pages 37–39.
- Aman Madaan, Shuyan Zhou, Uri Alon, Yiming Yang, and Graham Neubig. 2022. Language models of code are few-shot commonsense learners. *arXiv preprint arXiv:2210.07128*.
- Chris Maddison and Daniel Tarlow. 2014. Structured generative models of natural source code. In *International Conference on Machine Learning*, pages 649–657. PMLR.
- Kevin Meng, David Bau, Alex J Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in gpt. In *Advances in Neural Information Processing Systems*.
- Andriy Mnih and Geoffrey Hinton. 2007. Three new graphical models for statistical language modelling. In *Proceedings of the 24th international conference on Machine learning*, pages 641–648.
- Leonardo de Moura and Nikolaj Bjørner. 2008. Z3: An efficient smt solver. In *International conference on Tools and Algorithms for the Construction and Analysis of Systems*, pages 337–340. Springer.

- Spyridon Mouselinos, Mateusz Malinowski, and Henryk Michalewski. 2022. A simple, yet effective approach to finding biases in code generation. *arXiv preprint arXiv:2211.00609*.
- Arvind Neelakantan, Tao Xu, Raul Puri, Alec Radford, Jesse Michael Han, Jerry Tworek, Qiming Yuan, Nikolas Tezak, Jong Wook Kim, Chris Hallacy, et al. 2022a. Text and code embeddings by contrastive pre-training. *arXiv preprint arXiv:2201.10005*.
- Arvind Neelakantan, Tao Xu, Raul Puri, Alec Radford, Jesse Michael Han, Jerry Tworek, Qiming Yuan, Nikolas Tezak, Jong Wook Kim, Chris Hallacy, et al. 2022b. Text and code embeddings by contrastive pre-training. *arXiv preprint arXiv:2201.10005*.
- Nhan Nguyen and Sarah Nadi. 2022. An empirical evaluation of github copilot’s code suggestions. In *Proceedings of the 19th International Conference on Mining Software Repositories*, pages 1–5.
- Tung Thanh Nguyen, Anh Tuan Nguyen, Hoan Anh Nguyen, and Tien N Nguyen. 2013. A statistical semantic language model for source code. In *Proceedings of the 2013 9th Joint Meeting on Foundations of Software Engineering*, pages 532–542.
- Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese, and Caiming Xiong. 2022. A conversational paradigm for program synthesis. *arXiv preprint arXiv:2203.13474*.
- Yusuke Oda, Hiroyuki Fudaba, Graham Neubig, Hideaki Hata, Sakriani Sakti, Tomoki Toda, and Satoshi Nakamura. 2015. Learning to generate pseudo-code from source code using statistical machine translation. In *2015 30th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, pages 574–584. IEEE.
- Gabriel Orlanski and Alex Gittens. 2021. Reading stackoverflow encourages cheating: Adding question text improves extractive code generation. In *Proceedings of the 1st Workshop on Natural Language Processing for Programming (NLP4Prog 2021)*, pages 65–76.
- Gabriel Orlanski, Seonhye Yang, and Michael Healy. 2022. Evaluating how fine-tuning on bimodal data effects code generation. *arXiv preprint arXiv:2211.07842*.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*.
- Incheon Paik and Jun-Wei Wang. 2021. Improving text-to-code generation with features of code graph on gpt-2. *Electronics*, 10(21):2706.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.
- Hammond Pearce, Baleegh Ahmad, Benjamin Tan, Brendan Dolan-Gavitt, and Ramesh Karri. 2022. Asleep at the keyboard? assessing the security of github copilot’s code contributions. In *2022 IEEE Symposium on Security and Privacy (SP)*, pages 754–768. IEEE.
- Long Phan, Hieu Tran, Daniel Le, Hieu Nguyen, James Annibal, Alec Peltekian, and Yanfang Ye. 2021. CoTexT: Multi-task learning with code-text transformer. In *Proceedings of the 1st Workshop on Natural Language Processing for Programming (NLP4Prog 2021)*, pages 40–47.
- Gabriel Poesia, Oleksandr Polozov, Vu Le, Ashish Tiwari, Gustavo Soares, Christopher Meek, and Sumit Gulwani. 2022. Synchromesh: Reliable code generation from pre-trained language models. *arXiv preprint arXiv:2201.11227*.
- Antonio Polino, Razvan Pascanu, and Dan Alistarh. 2018. Model compression via distillation and quantization. In *International Conference on Learning Representations*.
- Maja Popović. 2015. chrF: character n-gram f-score for automatic mt evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395.
- Archiki Prasad, Peter Hase, Xiang Zhou, and Mohit Bansal. 2022. Grips: Gradient-free, edit-based instruction search for prompting large language models. *arXiv preprint arXiv:2203.07281*.
- Julian Aron Prenner and Romain Robbes. 2021. Automatic program repair with openai’s codex: Evaluating quixbugs. *arXiv preprint arXiv:2111.03922*.
- Yewen Pu, Zachery Miranda, Armando Solar-Lezama, and Leslie Pack Kaelbling. 2018. Learning to select examples for program synthesis. *arXiv preprint arXiv:1711.03243*.
- Maxim Rabinovich, Mitchell Stern, and Dan Klein. 2017. Abstract syntax networks for code generation and semantic parsing. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1139–1149.
- Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving language understanding by generative pre-training. *OpenAI blog*.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, Peter J Liu, et al. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21(140):1–67.
- Nitarshan Rajkumar, Raymond Li, and Dzmitry Bahdanau. 2022. Evaluating the text-to-sql capabilities of large language models. *arXiv preprint arXiv:2204.00498*.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392.
- Veselin Raychev, Pavol Bielik, and Martin Vechev. 2016. Probabilistic model for code with decision trees. *ACM SIGPLAN Notices*, 51(10):731–747.
- Veselin Raychev, Martin Vechev, and Eran Yahav. 2014. Code completion with statistical language models. In *Proceedings of the 35th ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 419–428.
- Shuo Ren, Daya Guo, Shuai Lu, Long Zhou, Shujie Liu, Duyu Tang, Neel Sundaresan, Ming Zhou, Ambrosio Blanco, and Shuai Ma. 2020. CodeBLEU: a method for automatic evaluation of code synthesis. *arXiv preprint arXiv:2009.10297*.
- Sami Sarsa, Paul Denny, Arto Hellas, and Juho Leinonen. 2022. Automatic generation of programming exercises and code explanations using large language models. In *Proceedings of the 2022 ACM Conference on International Computing Education Research-Volume 1*, pages 27–43.
- Tal Schuster, Ashwin Kalyan, Oleksandr Polozov, and Adam Tauman Kalai. 2021. Programming puzzles. *arXiv preprint arXiv:2106.05784*.
- Holger Schwenk and Jean-Luc Gauvain. 2002. Connectionist language modeling for large vocabulary continuous speech recognition. In *2002 IEEE International Conference on Acoustics, Speech, and Signal Processing*, volume 1, pages I–765. IEEE.
- Jieke Shi, Zhou Yang, Bowen Xu, Hong Jin Kang, and David Lo. 2022. Compressing pre-trained models of code into 3 mb. *arXiv preprint arXiv:2208.07120*.
- Disha Shrivastava, Hugo Larochelle, and Daniel Tarlow. 2022. Repository-level prompt generation for large language models of code. In *ICML 2022 Workshop on Knowledge Retrieval and Language Models*.
- Mohammed Latif Siddiq and Joanna CS Santos. 2022. SecurityEval dataset: mining vulnerability examples to evaluate machine learning-based code generation techniques. In *Proceedings of the 1st International Workshop on Mining Software Repositories Applications for Privacy and Security*, pages 29–33.
- Dominik Sobania, Martin Briesch, and Franz Rothlauf. 2022. Choose your programming copilot: a comparison of the program synthesis performance of github copilot and genetic programming. In *Proceedings of the Genetic and Evolutionary Computation Conference*, pages 1019–1027.
- Zeyu Sun, Qihao Zhu, Lili Mou, Yingfei Xiong, Ge Li, and Lu Zhang. 2019. A grammar-based structural cnn decoder for code generation. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 7055–7062.
- Ilya Sutskever, Geoffrey E Hinton, and Graham W Taylor. 2008. The recurrent temporal restricted boltzmann machine. *Advances in neural information processing systems*, 21.
- Alexey Svyatkovskiy, Shao Kun Deng, Shengyu Fu, and Neel Sundaresan. 2020. Intellicode compose: Code generation using transformer. In *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, pages 1433–1443.
- Eugene Syriani, Lechanceux Luhunu, and Houari Sahraoui. 2018. Systematic mapping study of template-based code generation. *Computer Languages, Systems & Structures*, 52:43–62.
- Riku Takaichi, Yoshiki Higo, Shinsuke Matsumoto, Shinji Kusumoto, Toshiyuki Kurabayashi, Hiroyuki Kirinuki, and Haruto Tanno. 2022. Are nlp metrics suitable for evaluating generated code? In *International Conference on Product-Focused Software Process Improvement*, pages 531–537. Springer.
- Leonard Tang, Elizabeth Ke, Nikhil Singh, Nakul Verma, and Iddo Drori. 2021. Solving probability and statistics problems by program synthesis. *arXiv preprint arXiv:2111.08267*.
- Ngoc Tran, Hieu Tran, Son Nguyen, Hoan Nguyen, and Tien Nguyen. 2019. Does bleu score work for code migration? In *2019 IEEE/ACM 27th International Conference on Program Comprehension (ICPC)*, pages 165–176. IEEE.
- Immanuel Trummer. 2022. CodexDB: Synthesizing code for query processing from natural language instructions using gpt-3 codex. *Proceedings of the VLDB Endowment*, 15(11):2921–2928.
- Priyan Vaithilingam, Tianyi Zhang, and Elena L Glassman. 2022. Expectation vs. experience: Evaluating the usability of code generation tools powered by large language models. In *CHI Conference on Human Factors in Computing Systems Extended Abstracts*, pages 1–7.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

- Yao Wan, Zhou Zhao, Min Yang, Guandong Xu, Haochao Ying, Jian Wu, and Philip S Yu. 2018. Improving automatic source code summarization via deep reinforcement learning. In *Proceedings of the 33rd ACM/IEEE international conference on automated software engineering*, pages 397–407.
- Chaozheng Wang, Yuanhang Yang, Cuiyun Gao, Yun Peng, Hongyu Zhang, and Michael R Lyu. 2022a. No more fine-tuning? an experimental evaluation of prompt tuning in code intelligence. In *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, pages 382–394.
- Xin Wang, Yasheng Wang, Yao Wan, Fei Mi, Yitong Li, Pingyi Zhou, Jin Liu, Hao Wu, Xin Jiang, and Qun Liu. 2022b. Compilable neural code generation with compiler feedback. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 9–19.
- Xingyao Wang, Sha Li, and Heng Ji. 2022c. Code4struct: Code generation for few-shot structured prediction from natural language. *arXiv preprint arXiv:2210.12810*.
- Yue Wang, Weishi Wang, Shafiq Joty, and Steven CH Hoi. 2021. Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 8696–8708.
- Zhiruo Wang, Grace Cuenca, Shuyan Zhou, Frank F Xu, and Graham Neubig. 2022d. MCoNaLa: A benchmark for code generation from multiple natural languages. *arXiv preprint arXiv:2203.08388*.
- Bolin Wei, Ge Li, Xin Xia, Zhiyi Fu, and Zhi Jin. 2019. Code generation as a dual task of code summarization. *Advances in neural information processing systems*, 32.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903*.
- Frank F Xu, Uri Alon, Graham Neubig, and Vincent Josua Hellendoorn. 2022a. A systematic evaluation of large language models of code. In *Proceedings of the 6th ACM SIGPLAN International Symposium on Machine Programming*, pages 1–10.
- Frank F Xu, Zhengbao Jiang, Pengcheng Yin, Bogdan Vasilescu, and Graham Neubig. 2020. Incorporating external knowledge through pre-training for natural language to code generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6045–6052.
- Frank F Xu, Bogdan Vasilescu, and Graham Neubig. 2022b. In-ide code generation from natural language: Promise and challenges. *ACM Transactions on Software Engineering and Methodology (TOSEM)*, 31(2):1–47.
- Weixiang Yan and Yuanchun Li. 2022. Why-gen: Explaining ml-powered code generation by referring to training examples. *arXiv preprint arXiv:2204.07940*.
- Yuanmeng Yan, Rumei Li, Sirui Wang, Hongzhi Zhang, Zan Daoguang, Fuzheng Zhang, Wei Wu, and Weiran Xu. 2021. Large-scale relation learning for question answering over knowledge bases with pre-trained language models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3653–3660.
- Guang Yang, Yu Zhou, Xiang Chen, Xiangyu Zhang, Tingting Han, and Taolue Chen. 2022. ExploitGen: Template-augmented exploit code generation based on codebert. *Journal of Systems and Software*, page 111577.
- Wei Ye, Rui Xie, Jinglei Zhang, Tianxiang Hu, Xiaoyin Wang, and Shikun Zhang. 2020. Leveraging code generation to improve code retrieval and summarization via dual learning. In *Proceedings of The Web Conference 2020*, pages 2309–2319.
- Pengcheng Yin, Bowen Deng, Edgar Chen, Bogdan Vasilescu, and Graham Neubig. 2018. Learning to mine aligned code and natural language pairs from stack overflow. In *2018 IEEE/ACM 15th international conference on mining software repositories (MSR)*, pages 476–486. IEEE.
- Pengcheng Yin and Graham Neubig. 2017. A syntactic neural model for general-purpose code generation. *arXiv preprint arXiv:1704.01696*.
- Donghan Yu, Chenguang Zhu, Yiming Yang, and Michael Zeng. 2022. Jacket: Joint pre-training of knowledge graph and language understanding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 10, pages 11630–11638.
- Daoguang Zan, Bei Chen, Zeqi Lin, Bei Guan, Yongji Wang, and Jian-Guang Lou. 2022a. When language model meets private library. *EMNLP Findings*.
- Daoguang Zan, Bei Chen, Dejian Yang, Zeqi Lin, Minsu Kim, Bei Guan, Yongji Wang, Weizhu Chen, and Jian-Guang Lou. 2022b. CERT: Continual pre-training on sketches for library-oriented code generation. *IJCAI*.
- Daoguang Zan, Wang Sirui, Zhang Hongzhi, Yan Yuanmeng, Wu Wei, Guan Bei, and Wang Yongji. 2022c. S²QL: Retrieval augmented zero-shot question answering over knowledge graph. In *PAKDD 2022*, volume 13282 of *Lecture Notes in Computer Science*, pages 223–236. Springer.
- Maksym Zavershynskiy, Alex Skidanov, and Illia Polosukhin. 2018. NAPS: Natural program synthesis dataset. *arXiv preprint arXiv:1807.03168*.

- Jialu Zhang, José Cambronero, Sumit Gulwani, Vu Le, Ruzica Piskac, Gustavo Soares, and Gust Verbruggen. 2022a. Repairing bugs in python assignments using large language models. *arXiv preprint arXiv:2209.14876*.
- Sarah Zhang, Reece Shuttleworth, Derek Austin, Yann Hicke, Leonard Tang, Sathwik Karnik, Darnell Granberry, and Iddo Drori. 2022b. A dataset and benchmark for automatically answering and generating machine learning final exams. *arXiv preprint arXiv:2206.05442*.
- Zhaowei Zhang, Hongyu Zhang, Beijun Shen, and Xiaodong Gu. 2022c. Diet code is healthy: Simplifying programs for pre-trained models of code. In *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, pages 1073–1084.
- Hattie Zhou, Azade Nova, Hugo Larochelle, Aaron Courville, Behnam Neyshabur, and Hanie Sedghi. 2022a. Teaching algorithmic reasoning via in-context learning. *arXiv preprint arXiv:2211.09066*.
- Shuyan Zhou, Uri Alon, Frank F Xu, Zhengbao JIang, and Graham Neubig. 2022b. Doccoder: Generating code by retrieving and reading docs. *arXiv preprint arXiv:2207.05987*.
- Yabing Zhu, Yanfeng Zhang, Huili Yang, and Fangjing Wang. 2019. Gancoder: An automatic natural language-to-programming language translation approach based on gan. In *CCF International Conference on Natural Language Processing and Chinese Computing*, pages 529–539. Springer.