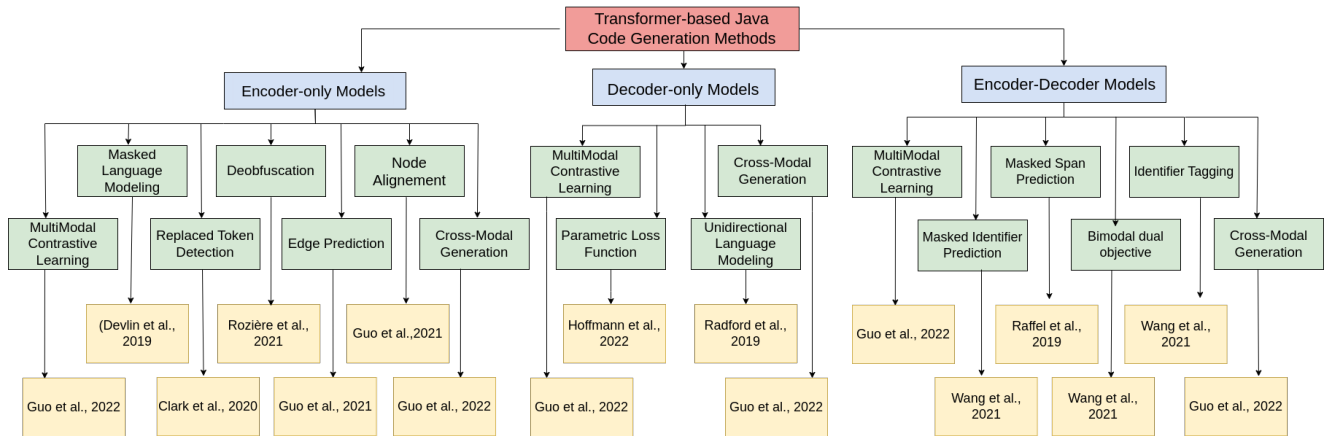# Graphical Abstract

**A Comprehensive Review of State-of-The-Art Methods for Java Code Generation from Natural Language Text**

Jessica López Espejel,Mahaman Sanoussi Yahaya Alassan,El Mehdi Chouham,Walid Dahhane,El Hassane Ettifouri

Transformer-based Java Code Generation Methods

Highlights

**A Comprehensive Review of State-of-The-Art Methods for Java Code Generation from Natural Language Text**

Jessica López Espejel,Mahaman Sanoussi Yahaya Alassan,El Mehdi Chouham,Walid Dahhane,El Hassane Ettifouri

- Code Generation is an important Natural Language Processing (NLP) task

- Initializing models from pretrained weights leads to better results than training them from scratch

- Decoder-only models achieve the best comprehension of the code syntax and semantics

- Combining multiple learning objectives lead to better code generation models.

- Improving Code Generation metrics is becoming a must in order to better compare and further improve state-of-the-art methods

# A Comprehensive Review of State-of-The-Art Methods for Java Code Generation from Natural Language Text[*]

Jessica López Espejel[a,*], Mahaman Sanoussi Yahaya Alassan[a], El Mehdi Chouham[a], Walid Dahhane[a] and El Hassane Ettifouri[a]

[a]*Novelis Research and Innovation Lab, 207 Rue de Bercy, Paris, 75012, , France*
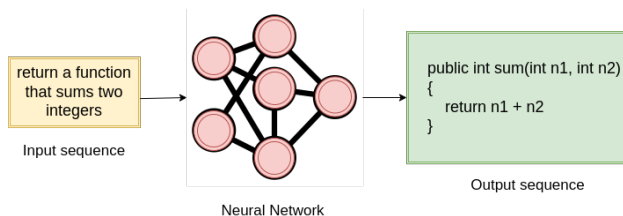
## ARTICLE INFO

## ABSTRACT

Java Code Generation consists in generating automatically Java code from a Natural Language Text. This NLP task helps in increasing programmers' productivity by providing them with immediate solutions to the simplest and most repetitive tasks. Code generation is a challenging task because of the hard syntactic rules and the necessity of a deep understanding of the semantic aspect of the programming language. Many works tried to tackle this task using either RNN-based, or Transformer-based models. The latter achieved remarkable advancement in the domain and they can be divided into three groups: (1) encoder-only models, (2) decoder-only models, and (3) encoder-decoder models. In this paper, we provide a comprehensive review of the evolution and progress of deep learning models in Java code generation task. We focus on the most important methods and present their merits and limitations, as well as the objective functions used by the community. In addition, we provide a detailed description of datasets and evaluation metrics used in the literature. Finally, we discuss results of different models on CONCODE dataset, then propose some future directions.

## 1. Introduction

In the last years, there has been a huge interest from the Natural Language Processing (NLP) community in the automation of software engineering to increase programmers' productivity (Ahmad et al., 2021). Code Generation Task (or *Program Synthesis Task*) helps considerably in reducing workload of programmers, by speeding up the implementation of simple functions, and letting them focus on the most complex tasks only.

In this paper, we are interested to show the progress over the years in the automatic Java source code generation from natural language. This task consists in taking as an input a natural language phrase, and producing an equivalent Java code as an output (Figure 1).



**Figure 1:** Java source code generation from natural language

In the Java Code Generation task as in many others, attention is put on both the model and the dataset. The latter are key ingredients to develop a powerful Artificial Intelligence (AI) system. Intuitively, most of code generation datasets are built from source code of the concerned Programming Language (PL) and documentation in Natural Language (NL). These datasets are collected from web sites such as Github[1] and Stack Overflow[2]. In the case of Java programming language, documentation is written in JavaDoc (Oracle, 2022), which is a format that describes all classes, member variables and methods. Therefore, JavaDoc documentation with its source code have been exploited to generate large datasets containing pairs of text-code lines (See Section 6 for more details). Along with datasets, several models have surged to tackle the Java Code Generation task. The first models used in literature are based on Recurrent Neural Networks (RNNs) such as LPN (Ling et al., 2016) and Seq2Seq (Yin and Neubig, 2017). However, the big success of Transformer-based pretrained language models (LMs) such as BERT (Devlin et al., 2019), GPT (Radford et al., 2019), and T5 (Raffel et al., 2020), in a wide range of NLP tasks, shifted the efforts of the community to adapt Transformers models (Vaswani et al., 2017) to generate Java code from natural language. The results of such efforts gave birth to more sophisticated Large Language Models (LLMs) such as ChatGPT (OpenIA, 2023), LLaMA (Touvron et al., 2023), and BARD (Google, 2023). Note that at the time of writing this article, the performance of the latter algorithms has not been tested on standard Java benchmarks. Nonetheless, initial indications suggest that their performance is exceptional.

Programming language generation presents more challenges than standard natural language generation. For instance, (1) the neural network should be able to correctly understand instructions in natural language to generate a corresponding correct source code, (2) the latter has lexical,

---

*Corresponding author

*Email addresses:* jlopezespejel@novelis.io (J. López Espejel); syahaya@novelis.io (M. Yahaya Alassan); elchouham@novelis.io (E.M. Chouham); wdahhane@novelis.io (W. Dahhane); eettifouri@novelis.io (E.H. Ettifouri)

ORCID(s): 0000-0001-6285-0770 (J. López Espejel); 0000-0001-5387-3380 (W. Dahhane); 0000-0001-5299-9053 (E.H. Ettifouri)

[1]https://github.com/
[2]https://stackoverflow.com/

grammatical, and semantic constraints (Wang et al., 2021; Scholak et al., 2021) that should be taken into account, and (3) one small mistake (such as a missing dot or colon) can change completely the semantic of the code and make the model's output incorrect.

In this paper, we provide a comprehensive review of state-of-the-art methods in Java code generation from natural language. To the best of our knowledge, this is the first review for Java Code Generation methods. Our goal is to provide a solid base for future researchers by highlighting what was already done, and what can be improved. These are the main findings of our paper:

- Initializing models from pretrained weights leads to better results than training them from scratch

- Decoder-only models have shown the best comprehension of the code syntax and semantics

- Combining multiple learning tasks lead to better code generation models.

- Improving Code Generation metrics is becoming a must in order to better compare and further improve state-of-the-art methods

The rest of the paper is organized as follows: Section 2 provides a formulation of the problem. Section 3 outlines the background of the studied task. Section 4 presents RNN-based methods to tackle Java Code Generation. Note that this type of methods are outdated nowadays and are not effective to tackle complex NLP tasks. Section 5 presents the most powerful Transformer-based methods. This section is the largest, and takes the biggest attention in this article. Section 6 provides a brief description of the most popular datasets used in Java Generation. Section 7 presents evaluation metrics used by the community to evaluate different models. In Section 8, we compare experimentally between the state-of-the-art methods and highlight their advantages and disadvantages. Finally, we conclude in Section 9 and provide some perspectives.

## 2. Problem Definition

We build on the notations used in Feng et al. (2020) and propose a unified notation for all methods described in this paper. The goal is to facilitate the comprehension and the comparison between state-of-the-art methods.

Given an input text $w = \{w_0, w_1, .., w_{n-1}\}$ of length $|w| = n$, the goal of a code generation model is to produce an equivalent source code sequence $c = \{c_0, c_1, .., c_{l-1}\}$ of length $|c| = l$, where $w_i$ and $c_i$ are input and output (target) tokens respectively. The model is trained on batches of size $B$ using a loss $\mathcal{L}(\theta)$, where $\theta$ are the model's parameters. The formula of loss function $\mathcal{L}(\theta)$ varies between code generation methods and depends heavily on the architecture used. We will provide a detailed explanation of these learning objectives in Section 5.

**Data type -** Code generation models can be trained on unimodal or bimodal data. On the one hand, unimodal data keep the source code without paired natural language, or keep the natural language without paired code. On the other hand, bimodal data have NL-PL pairs. The bimodal data have shown its usefulness in training the multi-modal pretrained models (Feng et al., 2020), which learn implicit alignment between input of different categories.
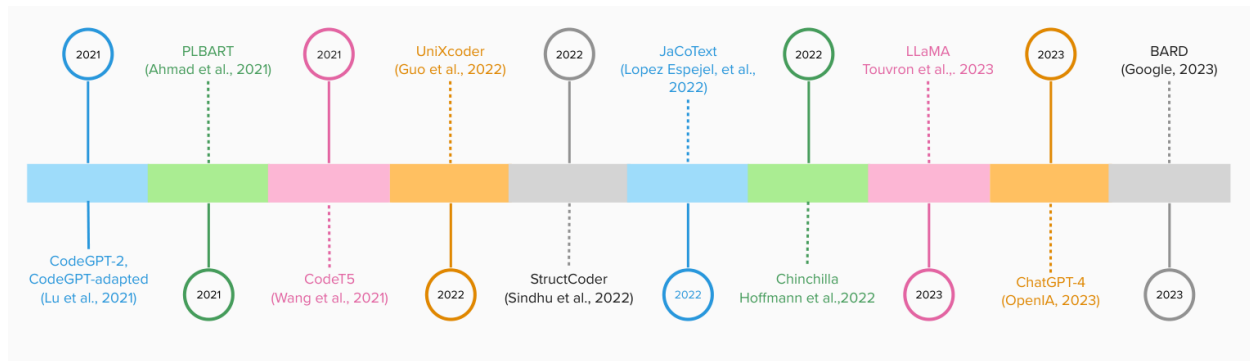
## 3. Background

The automatic source code generation task can be modeled as a translation task, where the input text is a natural language phrase and the output text is a programming source code. Early works in the field of language translation focused on regular expressions, logical forms, sequence of instructions, and agent-specific languages. Later, the efforts were shifted to use semantic parsing and SQL code generation. These developments paved the way for the generation of Java code for natural language. In the following paragraphs, we will provide a brief overview of remarkable methods proposed for each category.

**Regular expressions -** Angluin (1987) worked on the mapping of natural language to regular expressions by learning general rules from examples, and automatizing their generation. Alternatively, Ranta (1998) and Kushman and Barzilay (2013a) used rule-based techniques and Combinatory Categorical Grammar (CCG) to generate regular expressions from natural language, respectively. Similarly, Kushman and Barzilay (2013b) worked on generating regular expressions from natural language queries. While the method proved effective in handling a wide variety of natural language queries, including complex ones, it was not able to handle all of them, as it relies on semantic unification and could struggle with certain types of queries.

**Logical forms -** Authors of Zettlemoyer and Collins (2005a) mapped natural language sentences to logical forms. They introduced an approach based on probabilistic categorial grammars (PCGs). Even though the method handles ambiguity in natural language, and can generate multiple logical forms, the PCG parser used in the model is complex and computationally expensive, making it difficult to scale to large datasets.

**Agent-specific language -** Kate et al. (2005) proposed a method for automatically translating natural language sentences into agent-specific language expressions. The advantages of this method is that it can automate the translation process, saving time and effort for human translators, and can potentially improve the accuracy of formal language expression generation. However, the model could face difficulties in processing certain categories of complex sentences or in handling expressions that have multiple valid formal language translations.

**Sequence of instructions -** Branavan et al. (2010) presented a framework that learns to map natural language instructions to corresponding actions in an environment by training an agent with a combination of supervised learning

**Figure 2:** Timeline of some models in the Java Code Generation task

and reinforcement learning techniques. The method does not require hand-engineered features, making it highly adaptable to different tasks and environments. However, the framework struggles with instructions that are ambiguous or imprecise, leading to errors in action generation.

***Semantic Parsing -*** We have identified salient background approaches in the semantic parsing. One example is the work introduced by Zettlemoyer and Collins (2005a), who mapped NL to lambda–calculus encodings of their semantics. To check the syntax and the semantics, the authors used a log-linear model. Later, Wong and Mooney (2006), developed a statistical semantic parsing system called WASP. This system is designed to generate a formal representation of the meaning of a sentence. Similarly, Lu et al. (2008) proposed for parsing natural language, which involves utilizing a generative model to convert input sentences into representations of their meanings.

***SQL query generation -*** This is the closest type of approaches to Java generation task. One interesting work in this category was proposed by Miller et al. (1996) who used an interface that is based on trained statistical models. To train the model, authors used ATIS (Air Travel Information) dataset. This dataset also was used by Ramaswamy and Kleindienst (2000).

The aforementioned works offer an introduction to the methods utilized in the background and establish the groundwork for more advanced techniques based on RNN and Transformers, which will be discussed in subsequent sections of this paper. Unfortunately, most of the methods presented in this section are outdated and are unable to deal with complex Java generation programs.

## 4. RNN-based Code Generation Methods

The first neural-network-based approaches used to tackle source code generation from natural language were based on Recurrent Neural Networks (RNNs) such as Long Short-term Memory (LSTM) (Hochreiter and Schmidhuber, 1997) and Gated Recurrent Unit (GRU) (Cho et al., 2014a). For instance, Neelakantan et al. (2016) proposed an architecture known as the Neural Programmer, which integrates RNNs with a set of fundamental arithmetic and logic operations

to enhance program induction methods. One key feature of their approach is the integration of additional memory into neural networks (Graves et al., 2014; Kumar et al., 2016; Joulin and Mikolov, 2015), which allows for more advanced and sophisticated problem-solving capabilities. Moreover, Mou et al. (2015) used a method based on RNNs for generating computer programs directly from natural language input. This sequence-to-sequence model is an end-to-end approach that uses an LSTM-based encoder to create a fixed-length vector representation of the input, and a decoder that predicts each token based on the output sequence tokens and the input vector.

Similarly to Neelakantan et al. (2016), Reed and de Freitas (2016) introduced an approach to induce programs called NPI (Neural Programmer-Interpreter). As Mou et al. (2015), NPI is an end-to-end model. NPI does not need manual feature engineering or domain-specific knowledge. It can learn how to create and execute programs, but the method is limited to generate complex programs, as the size of the programs is restricted by the capacity of the neural network. In the same year, Ling et al. (2016) proposed Latent Predictor Network (LPN) that generates code from natural language descriptions by predicting the latent representation of code snippets. The LPN model consists of two components: a code generation model and a latent predictor model. The code generation model maps natural language descriptions to code representations, and the latent predictor model predicts the latent code representation given the previous generated code. However, LPN model requires a large amount of training data to perform well.

Another interesting work in this category is called Deep-API. It was proposed by Gu et al. (2016b), and consists in generating Java API sequences from natural language query. DeepAPI adapts a RNN encoder-decoder model (Cho et al., 2014b) to encode the input sequence to a fixed-length context vector. This helps in recognizing semantically related words. In addition, importance of individual APIs is used to enhance the performance of the model. DeepAPI is the first model from its kind, but unfortunately, it only focuses on JDK library.

Iyer et al. (2018) introduced a sequence-to-sequence model and evaluated its performance on their CONCODE

dataset. The model is composed of a bidirectional LSTM encoder, and an LSTM-based decoder. The latter uses a two-step attention mechanism, one of these steps is a supervised copy of Gu et al. (2016a). Alternatively, Iyer et al. (2018) tested seq2seq model introduced by Yin and Neubig (2017) on their CONCODE dataset. Seq2seq uses sequence-to-sequence LSTM-based architecture. However, the decoder uses supervised copying from the whole input sequence.

Despite the usefulness of RNNs in multiple NLP tasks, they are currently not very popular in the community since they were highly outperformed in 2017 by Transformer neural networks (Vaswani et al., 2017). Most of the works presented in this section did not focus on Java programming language. This is showed in Figure 2 that outlines some methods focus on Java code generation. We can clearly see that Java PL task is indeed tackled using Transformers. For this reason, we booked more room for them in the rest of the paper.
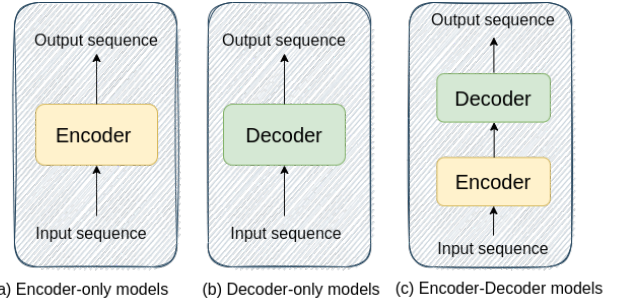
## 5. Transformer-based Java Code Generation Methods

Transformers-based models have shown outstanding scores on various NLP tasks such as translation (Vaswani et al., 2017), automatic summarization (Zhang et al., 2020) and question answering (Raffel et al., 2020). All of these models have two stages: (1) pre-training, and (2) fine-tuning. The pre-training consists of learning a language model or a learning objective task from unlabeled data, and fine-tuning is the step in which the model learns the knowledge that is related to a specific task using labeled data. Therefore, one of the main factors in the success of these models is the supervised learning objectives used. Code generation methods tried first to adapt existing standard text generation methods and their objectives to Code Generation (Phan et al., 2021; Ahmad et al., 2021). Later, other methods explored more specific techniques to improve the syntax and semantics of the generated code (Tipirneni et al., 2022; Wang et al., 2021). We study state-of-the-art methods based on their architecture and divide then into three main categories: encoder-only models, decoder-only models, and encoder-decoder models. Figure 3 shows a simplified illustration of the main difference in the architecture between the three types of models. Figure 4 shows the main learning objectives used in each category of methods, plus, an example of reference for each objective. We provide more details in the following three subsections.

### 5.1. Encoder-only models

These methods use an encoder only in the Transformer neural network architecture (Figure 3 (a)). They concatenate the natural language sequence $w$, and the programming language sequence $c$ and use them as an input $x = \{w, c\}$ to the encoder.

The main objectives used by this type of methods are as follows:



**Figure 3:** A simplified illustration of the different architecture types of Transformer neural networks

- *MLM: Masked Language Modeling* (Devlin et al., 2019) - It consists in selecting randomly a set of positions from the input tokens, then replacing tokens in these positions with a special [MASK] token. This objective aims to predict only the original tokens that are masked, and not to reconstruct the entire input. The token chosen at each position is selected from the model vocabulary. The MLM loss function is defined in Equation 1.

$$\mathcal{L}_{MLM}(\theta) = - \sum_{i \in m^w \cup m^c} log\, P_\theta^{D_1}(x_i | w^{masked}, c^{masked}) \quad (1)$$

where $\theta$ are the model's parameters, $w$ is a natural language input, $c$ is programming language output, $P_\theta^{D_1}$ is the discriminator that predicts a token from a large vocabulary, and $m^w$ and $m^c$ are random set of masked positions in $w$ and $c$, respectively.

- *RTD: Replaced Token Detection* (Clark et al., 2020) - It replaces (corrupts) some tokens in the input sequence with plausible tokens proposed from a generator network. The discriminator neural network is trained to determine if each token in the corrupt input is the original one or not. The RTD loss function is defined in Equation 2.

$$\mathcal{L}_{RTD}(\theta) = \sum_{i=0}^{n+l-1} (\delta(i)\, log\, P_\theta^{D_2}(x^{corrupt}, i) + (1 - \delta(i))(1 - log\, P_\theta^{D_2}(x^{corrupt}, i))) \quad (2)$$

where $P_\theta^{D_2}$ is the discriminator which predicts the probability of the $i$-th word being original, and $\delta(i)$ is an indicator function defined as:

$$\delta(i) = \begin{cases} 1 & if\ x_i^{corrupt} \equiv x_i, \\ 0 & otherwise. \end{cases}$$

- *DOBF: Deobfuscation* (Rozière et al., 2021) - It renames the class, function and variable names with uninformative labels. All instances of selected identifiers are replaced in the whole source code with the uninformative label. The model is trained to recover the original names by deobfuscating the code. It is the first approach to exploit the structure of programming languages.

• *Edge Prediction* (Guo et al., 2021) - It is used with Data Flow Graph (DFG) and consists in masking some direct edges in the data flow. Then, the model aims to predict these masked edges. Edge Prediction loss is defined in Equation 3.

$$\mathcal{L}_{EdgePred}(\theta) = -\sum_{e_{ij} \in E_c} [\delta(e_{ij} \in E^{masked}) \, log \, p_{e_{ij}} + (1 - \delta(e_{ij} \in E^{masked})) \, log \, (1 - p_{e_{ij}})]$$

(3)

where $E_c$ is the set of candidate edges, $E^{masked}$ is the set of masked edges, $p_{e_{ij}}$ is the probability that an edge exists between nodes $i$ and $j$, and $\delta(e_{ij} \in E)$ is defined as follows:

$$\delta(e_{ij} \in E) = \begin{cases} 1 & if \; \langle v_i, v_j \rangle \in E, \\ 0 & otherwise. \end{cases}$$

• *Node Alignement* (Guo et al., 2021) - It is similar to the Edge Prediction task, but instead of predicting edges between the nodes of DFG, it predicts edges between the source code tokens and data flow nodes.

**Methods -** CodeBERT (Feng et al., 2020) is based on BERT (Devlin et al., 2019), and is the first large bimodal (NL-PL) pretrained model on CodeSearchNet (Husain et al., 2019) dataset. In the pretraining stage, CodeBERT uses two objective functions: Masked Language Modeling (MLM), and Replaced Token Detection (RTD).

GraphCodeBERT (Guo et al., 2021) is based on BERT (Devlin et al., 2019), and is the first pretrained model that leverages the code structure via Data Flow Graph (DFG). In the pretraining stage, GraphCodeBERT uses MLM as a learning objective, and two additional tasks: prediction of code structure edges in the data flow, and variable-alignment over source code and data flow.

## 5.2. Decoder-only models

These methods use decoder(s) only in the Transformer neural network architecture (Figure 3 (b)).

• *PAR: Parametric Loss Function* (Hoffmann et al., 2022) - This is one of the most recent objectives proposed for decoder-only Chinchilla model. It is a parametric loss function, based on models parameters count and the number of already-seen tokens. The PAR loss is defined in Equation 4:

$$\mathcal{L}_{PAR}(\theta, N, D) \triangleq E + \frac{A}{N^\alpha} + \frac{B}{D^\beta}$$

(4)

where $E$ is the entropy of the natural language text and refers to ideal generative process, $\frac{A}{N^\alpha}$ refers to the negative gap in performance between a perfectly trained Transformer with N parameters and an ideal generative process, and $\frac{B}{D^\beta}$ refers to incomplete convergence of the Transformer model because of its limited number of optimization steps. The estimation of $\alpha, \beta, A, B, E$ is done using a Huber Loss (Hernandez et al., 2021) that uses the L-BFGS algorithm (Nocedal, 1980), as in Equation 5:

$$\min_{A,B,E,\alpha,\beta} \sum_{runs \, i} Huber_\delta \left( log \, \mathcal{L}_{PAR}(\theta, N_i, D_i) - log \, L_i \right) \quad (5)$$

where $\delta = 10^{-3}$ and $L$ is the pretraining loss. Huber loss is robust to outliers and provides an overall good performance.

• *ULM: Unidirectional Language Modeling* (Radford et al., 2019) - It aims to predict the next token $x_i$ given the previous tokens. The ULM loss function is defined in Equation 6:

$$\mathcal{L}_{ULM}(\theta) = -\sum_{i=0}^{l-1} log \, P_\theta(x_i | x_{<i})$$

(6)

where $x_{<i}$ refers to all the tokens that are before the $i^{th}$ one.

**Methods -** Lu et al. (2021) tested the performance of GPT-2 (Radford et al., 2019) model on code generation task using three configurations: GPT-2, CodeGPT-2 and CodeGPT-adapt. The latter is pretrained on text extracted from 45 million links on the web. Unlike GPT-2, CodeGPT-2 and CodeGPT-adapt that are pretrained on CodeSearch-Net dataset, CodeGPT-2 is pretrained from scratch, and CodeGPT-adapt is initialized from GPT-2 checkpoint.

**Large Language Models (LLMs) -** GPT-3 (Brown et al., 2020) marked the beginning of the era of Large Language Models (LLMs). It showed that, by scaling up language models, the few-shot learning is possible, and it can sometimes even reach competitive results in tasks such as machine translation, code generation, reading comprehension, etc. For all tasks, GPT-3 is applied without any gradient updates or fine-tuning, with tasks and few-shot demonstrations specified purely via text interaction with the model. This model follows GPT-2 architecture, and has 175B of parameters (10 times more than the previous one).

There is a series of GPT-3.5 models that are based on GPT-3 (Brown et al., 2020). We cite them hereafter:

- *Code-davinci-002*: is one of the Codex (Chen et al., 2021a) models, that are trained on natural language and programming languages including Java.

- *Text-davinci-002*: is an InstructGPT (Ouyang et al., 2022) model based on code-davinci-002.

- *Text-davinci-003*: is the most powerful GPT-3 model compared with previous GPT-based models. It can perform better in quality language generation, instruction-following, etc.

An upgrade of GPT-3.5 lead to the birth of Chat-GPT (OpenIA, 2023), an impressive language model that fueled a debate about the impact of AI on our daily life. Currently, ChatGPT is based on an even more powerful and sophisticated GPT-4 architecture. The latter is a large

multimodal model capable of processing both text and image inputs and producing text outputs. It surpasses existing models by a considerable margin in English and also exhibits strong performance in other languages. The loss function computation is one of the main modifications made compared to previous GPT models. The authors were able to predict GPT-4's final loss on their internal codebase by fitting a scaling law with an irreducible loss term. As for Java, ChatGPT is able to generate very complex code from natural language request.

Although GPT-4 is an advanced language model that outperforms its predecessors, it still has limitations similar to earlier GPT models. These limitations include generating false information, providing harmful advises, having a restricted context window, and lacking the ability to learn from experience. However, the authors conducted a thorough evaluation of GPT-4's performance on a variety of datasets, such as academic exams. It was found that GPT-4 achieved a human-level performance on most of these tests. Additionally, in their internal factuality evaluations designed to test the model's accuracy, GPT-4 scored 19% points higher than the latest GPT-3.5 model. Note that ChatGPT (and other large language models) is recent and was not tested yet on CONCODE dataset. For this reason, we did not include its results in our review.

At the same time than ChatGPT, Rae et al. (2021) built on GPT-2 and proposed Gopher, a model that contains 280 billion parameters. However, Gopher has two modifications compared with GPT-2: the authors used RMSNorm (Zhang and Sennrich, 2019) instead of LayerNorm (Ba et al., 2016), and also used the relative positional encoding (Dai et al., 2019) instead of the absolute positional encodings. The model was evaluated on 152 different tasks, including Java code generation. Even though Gopher outperforms the state of the art in 81% of the task, when the authors evaluated it in 124 tasks, there are still some challenges related to the distributional bias, and the training definition criteria.

After Gopher, Hoffmann et al. (2022) introduced a new language model called Chinchilla. This model has the same architecture as Gopher but with 70 billion parameters and additional data. Unlike Gopher, Chinchilla uses variations in the number of training steps, model sizes of different training FLOP counts, a different loss function ($PAR$ function from Equation 4), and the AdamW (Loshchilov and Hutter, 2019) optimizer instead of Adam (Kingma and Ba, 2015). Chinchilla outperforms various language models such as Gopher and GPT-3 in many tasks. However, the quality of code generation in Chinchilla is related to its size, and scaling training tokens can improve it. Although Chinchilla's authors recommend training a 10B model on 200B tokens, Touvron et al. (2023) has shown that a 7B model continues to improve even after 1T tokens.

Inspired by Chinchilla, Touvron et al. (2023) proposed the LLaMA family of LLMs, whose parameters range from 7B to 65B. The architecture of LLaMA incorporates the strengths of previous LLMs, such as input normalization inspired by GPT-3 (Brown et al., 2020), and SwiGLU activation function from PaLM (Driess et al., 2023). Additionally, it employs rotary positional embeddings instead of absolute positional embeddings used by GPTNeo (Black et al., 2022). LLaMA-13B, which is only one-tenth the size of GPT-3, performs better on most benchmarks and can be run on a single GPU. Furthermore, the 65B-parameter model is competitive with other LLMs like Chinchilla-70B and PaLM-540B.

## 5.3. Encoder-Decoder models

These methods use both encoder(s) and decoder(s) in the Transformer neural network architecture (Figure 3 (c)).

• *MSP: Masked Span Prediction* (Raffel et al., 2020)- It consists in randomly masking spans of text with arbitrary lengths. Later, the model predicts the masked spans. MSP loss is defined in Equation 7.

$$\mathcal{L}_{MSP}(\theta) = \sum_{i=0}^{k-1} -log\ P_\theta(w_i^{masked}|w^{\backslash masked}, w_{<i}^{masked})\ \ (7)$$

where $x^{masked}$ is the masked input of length $k$, $x^{\backslash masked}$ is the masked sequence to predict by the decoder, and $x_{<i}^{masked}$ is the span sequence generated so far.

• *IT: Identifier Tagging* (Wang et al., 2021) - It aims to determine if code tokens are identifiers (e.g., function or variable names). Therefore, authors compute a probability value $p_i$ for each output token.

To see if a node is an identifier or not, the authors convert the code sequence to an AST (Abstract Syntax Tree). The IT loss is computed using a binary cross entropy loss, where the binary label $y_i = 1$ means that the token is an identifier, and $y_i = 0$ means that it is not (Equation 8).

$$\mathcal{L}_{IT}(\theta_e) = \sum_{i=0}^{l-1} -[y_i\ log\ p_i + (1 - y_i)\ log\ (1 - p_i)]\ \ (8)$$

where $\theta_e$ are the encoder's parameters, and $p_i$ are the probabilities of the output sequence tokens.

• *MIP: Masked Identifier Prediction* (Wang et al., 2021) - It consists in masking all identifiers (function and variable names) in the PL, and replacing all recurrences of a specific identifier with a unique sentinel token. It is inspired by deobfuscation previously used in Rozière et al. (2021). The loss function is defined as:

$$\mathcal{L}_{MIP}(\theta) = \sum_{i=0}^{n-1} -log\ P_\theta(w_i|w^{masked}, w_{<i})\ \ \ \ (9)$$

where $w^{masked}$ is the masked input.

• *Bimodal dual objective* (Wang et al., 2021) - It aims to leverage bimodal data in order to achieve better NL-PL alignments. During the training process, if NL is the encoder input, then PL is the decoder input and vice versa.

• *MCL: MultiModal Contrastive Learning* (Guo et al., 2022) - This loss helps in learning the semantic embedding of mapped AST sequences. Authors forward the same input in the neural network and use different hidden dropout masks. The resulting hidden representations are used as positives, while the other representations from the same batch are used as negatives.

$$\mathcal{L}_{MCL}(\theta) = -\sum_{i=0}^{B-1} log \ \frac{e^{cos(\tilde{h}_i, \tilde{h}_i^+)}/\tau}{\sum_{j=0}^{b-1} e^{cos(\tilde{h}_j, \tilde{h}_j^+)}/\tau} \quad (10)$$

where $\tau$ is the temperature scalar, $\tilde{h}_i^+$ and $\tilde{h}_i$ are the positive and negative representations, respectively, and $cos()$ is the cosine similarity.

• *CMG: Cross-Modal Generation* (Guo et al., 2022) - This loss is complementary to the previous one. It pushes the model to provide comments to describe the code's function. This helps the model to better understand the semantics of the code and unify its representation across programming languages. Its is defined in Equation 11.

$$\mathcal{L}_{CMG}(\theta) = -\sum_{i=0}^{l-1} log \ P_\theta(d_i|X, d_{t<i}) \quad (11)$$

where X is the flattened AST sequence, and $d = \{d_0, d_1, .., d_{n-1}\}$ is the code comments.

**<u>Methods</u>** - Ahmad et al. (2021) introduced PLBART (Program and Language BART), a model based on $BART_{base}$ (Lewis et al., 2020). The latter was originally proposed to tackle text generation and comprehension tasks, such as translation, question answering, and automatic summarization. The only difference between the two model architectures is that PLBART added a normalization layer on top of both the encoder and the decoder. It is the first model to exploit the possibility to generate source code from NL with a pretrained decoder.

Posterior to PLBART, Phan et al. (2021) introduced CoTexT (Code and Text Transfer Transformer), which builds on $T5_{base}$ (Raffel et al., 2020) model. The authors studied the CoTexT performance based on two criteria: (1) combining corpora (**C**odeSearchNet and **G**itHub repositories), and (2) pretraining with unimodal and bimodal data. Therefore, the three resulting models are as follows: CoTexT (1-CC), CoTexT (2-CC), and CoTexT (1-CCG). Since all of the models are initialized from $T5_{base}$ weights, and the latter was originally pretrained on the C4 dataset, the comparison is fair.

Most recently, Lopez Espejel et al. (2023) introduced JaCoText, a model that is initialized from CotexT weights.
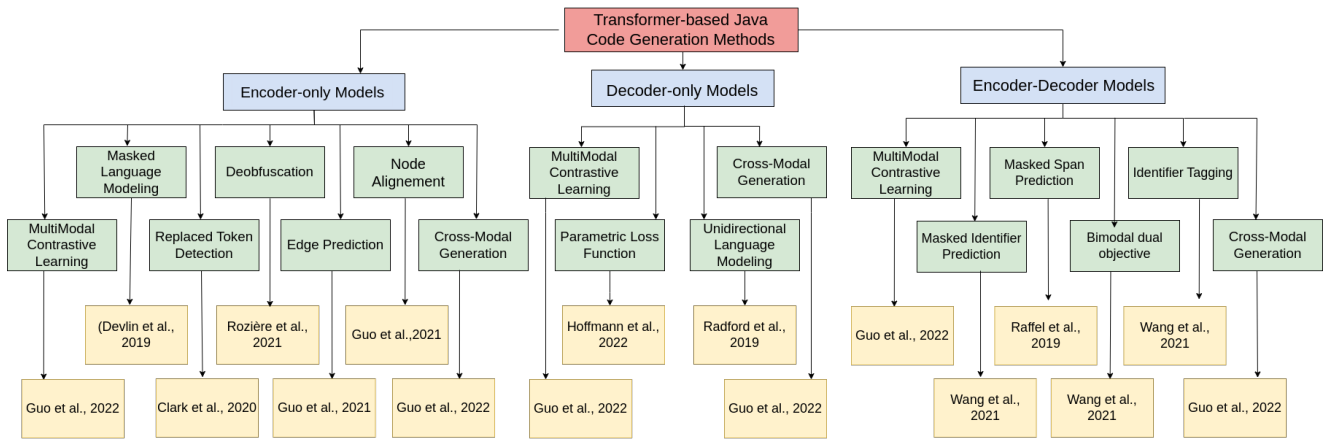
It explores additional pretraining using Java code only. In addition, the authors explore the impact of input and output sequence length during fine-tuning.

Parvez et al. (2021) introduced REDCODER (Retrieval augmentED CODe gEneration and summaRization framework). It is composed of a retrieval module called SCODER-R (Summary and CODE Retriever), and a generator module called SCODER-G, which is based on PLBLART (Ahmad et al., 2021). SCODER-R is initialized from GraphCode-BERT (Guo et al., 2021) weights, it receives the NL text, and returns the top-$k$ source codes. Later, the NL sequence is concatenated with the top-$k$ code sequences to have the augmented input sequence. This sequence is fed to PLBART to generate Java code.

CodeT5 (Wang et al., 2021) is based on T5 (Raffel et al., 2020) neural network architecture. Authors convert PL segments into an Abstract Syntax Tree (AST) and introduce two pretraining tasks: Identifier-aware denoising objective, and bimodal dual generation objective. On the one hand, Identifier-aware denoising objective uses Masked Span Prediction (MSP) objective similar to T5 (Raffel et al., 2020), and introduced two supplementary tasks: Identifier Tagging (IT), and Masked Identifier Prediction (MIP). Identifier-aware objective optimizes these three losses (MSP, IT, MIP) equally. On the other hand, bimodal dual objective optimizes the model simultaneously using bimodal data with different order (NL-PL, then PL-NL). Surprisingly, experiments lead by the authors have shown that the order between the code and the natural language test matters in the model's performance.

Later, Tipirneni et al. (2022) proposed StructCoder model, which is initialized from CodeT5 (Wang et al., 2021) weights, and pretrained using DAE (Denoising Autoencoding) task. StructCoder also uses an abstract syntax tree and propose two learning objectives to train both the encoder and the decoder: DFG prediction, and AST paths prediction. Specifically, the decoder is trained to predict the node types on all the root-leaf paths in the target AST, and to predict the data flow edges.

Very recently, UniXcoder (Guo et al., 2022) surged to tackle code-related understanding and generation tasks. UniXcoder leverages the AST information through transforming AST to sequence text. This method is functional in three modes: encoder-only, decoder-only, and encoder-decoder. Authors used three learning objective tasks to pretrained the model: MLM (Devlin et al., 2019), Unidirectional Language Modeling (ULM) (Radford et al., 2019), and Denoising Objective (Raffel et al., 2020). In addition, they introduced two more objectives: Multi-modal Contrastive Learning (MCL), and Cross-Modal Generation (CMG) that target the semantic understanding of code and its unification across different programming languages.

**Figure 4:** Overview of learning objectives in Transformer-based Code Generation Methods, and some examples of such methods.

## 6. Datasets for Java Code Generation

In this section, we present the most used datasets in Java Code Generation. A brief summary of them, as well as their merits and limitations, are presented in Table 1.

***CONCODE*** (Iyer et al., 2018) was collected from Java projects on Github (approx. 33,000 repositories). It is composed of tuples of class environment (member variables and member methods), natural language (JavaDoc comments), and source code. It aims to generate Java code from NL documentation and the class environment (used as a context). Java code belongs to several domains. It has been used to fine-tune the models. It contains $100K$, $2K$, and $2K$ lines of NL-PL in the train, validation, and test sets, respectively.

***CoNaLa*** (Yin et al., 2018) consists of query-code pairs collected from the website Stack Overflow via a mining method. It contains two programming languages: Python and Java. It contains 11125, 1237, and 500 lines of NL-PL in the train, validation, and test sets, respectively.

***MTG corpus*** (Ling et al., 2016) aims to generate source code given the card description in Trading Card Games (TCGs). The authors collected the data from two open source implementations of TCGs: Magic the Gathering (MTG), and Hearthtone (HS). While the latter digital implementation is in python. we focus in this work on MTG dataset that contains Java digital implementation. It contains 11969, 664, and 664 lines of NL-PL in the train, validation, and test sets, respectively.

***CodeSearchNet*** (General Language Understanding Evaluation benchmark for CODE) was introduced by Lu et al. (2021). It contains six programming languages (Ruby, JavaScript, Go, Python, Java, and PHP) with 6.4M of unimodal (only PL) data, and 2.1M bimodal (PL-NL) datapoints. The data was collected from Github code repositories, and includes 10 tasks over 14 different datasets.

***CodeNet*** was proposed by Puri et al. (2021b). It contains a set of programming problems collected from two websites: AIZU (Watanobe, 2022) and AtCoder (Inc., 2022). It contains 13,916,868 submissions, divided into 4053 problems. Submissions are in 55 different programming languages, but 95% are coded in C++, Python, Java, C, Ruby, and C#. Only the other 5% are in Java.

***Java dataset*** (Ahmad et al., 2021) was collected from Github repositories associated with Java language available on Google BigQuery. Authors follow Roziere et al. (2020) to preprocess the data.

***AlphaCode Pre-training*** was introduced by Li et al. (2022) to pretrained their model. The authors scrawled the data from public Github repositories. The pretrained dataset contains 715.1 GB of code, including C++, C#, Go, Java, JavaScript, Lua, PHP, Python, Ruby, Rust, Scala, and TypeScript. From these programming languages, 113.8 GB consist of Java data.

***CodeContests*** was introduced by Li et al. (2022) to fine-tuned their model. The authors combines scraped data from Codeforces (Mirzayanov, 2020) with existing data from Description2Code (Caballero et al., 2016), and CodeNet (Puri et al., 2021a). It contains C++, Python, and Java.

***BIGQUERY CODEGEN*** (Pang et al., 2022) contains six programming languages: C, C++, Go, Java, JavaScript, and Python. In total, it contains, 342.1 GB of data, including around 120.2 GB of Java data. It was used to pretrain the CODEGEN models.

***PolyCoder Dataset*** (Xu et al., 2022) was collected from Github repositories that have at least 50 stars. It contains 12 PLs (C, C#, C++, Go, Java, JavaScript, PHP, Python, Ruby, Rust, Scala, TypeScript), where each PL has a maximum of 25K repositories. The raw data size is around 631 GB (including 60 GB for Java). To preprocess the data, the authors follow Codex process (Chen et al., 2021b). After

| Dataset | Programming language (PL) | Size | Source | Proposed for | merits | limitations |
|---|---|---|---|---|---|---|
| CONCODE | Java | 100K, 2K, 2K lines in train/val/test | Github | code generation | large scale | query nl incorporates both member variables and member methods. However, in practice, is not an optimal approach. |
| CoNaLa | Python, Java | 11125, 1237, and 500 lines in train/val/test | Stack Overflow | classification | manual annotation | - If there are several valid source code for a question, the annotators might not identify all of them. - The annotators can fail when the code source is complex. |
| MTG corpus | Python, Java | 11969, 664, and 664 lines in train/val/test | trading card games Magic the Gathering and Hearthtone | code generation | new area of focus on code generation for card games. | - The programming of card games is a niche industry that caters to a specific demographic within society. |
| CodeSearchNet | Ruby, JavaScript, Go, Python, Java, PHP | 6.4M of PL data, 2.1M of PL-NL data | Github | code search | - large scale - expert annotators | The authors create proxy dataset of lower quality |
| BIGQUERY CODEGEN | C, C++, Go, Java, JavaScript, Python | 342.1 GB of total code, 120.3 GB of Java code | - | Multi-Turn program synthesis | - large scale - multi-lingual dataset | - The code in the dataset exhibits certain vulnerabilities and safety issues. |
| CodeNet | 55 different programming languages | 14 million of code samples and approx. 500 million lines of code | AIZU and AtCoder | - code similarity and classification code translation - code performance | - large scale - high-quality annotations | - The code samples may lack extensive comments and could be in multiple languages. - The code samples offer solutions for programming problems that are aimed at beginners in college and high school. |
| AlphaCode Pre-training dataset | C++, C#, Go, Java, JavaScript, Lua, PHP, Python, Ruby, Rust, Scala, TypeScript | 715.1 GB of total code, 113.8 GB of Java code | Github | code generation | - large scale - multiple PLs | The dataset may contain unsafe code |
| CodeContests | C++, Python, Java | 13328, 117, 165 in lines train/valid/test | Codeforces Description2Code CodeNet | code generation | - includes problems, solutions and test cases - contains correct and incorrect human submissions | - the raw datasets of CodeContests contain a high rate of false positive samples. |
| PolyCoder Dataset | C, C#, C++, Go, Java, JavaScript, PHP, Python, Ruby, Rust , Scala, TypeScript | 631 GB of total code, 60 GB of Java code | Github | code completion code synthesis | - large scale - multiple PLs | The code source may contain vulnerabilities and safety issues. |
| The Stack Dataset | 358 languages including Java | 6 TB of total code | Github | - code completion - documentation generation - auto-completion of code snippets | - large scale - multiple PLs | it contains unsafe code (malicious files) |
| Java dataset | Java | 352 GB | Github and StackOverflow | - code generation | - specialized in Java | - the code may contain vulnerabilities |

**Table 1**
Summary of datasets used in literature

filtering the data, the final data size is 243.6 GB, with 41 GB for Java language.

***The Stack Dataset*** was proposed by Kocetkov et al. (2022). It is the largest code dataset so far. It contains 6 TB of data from 358 programming languages. The code was extracted from Github repositories with permissive licenses, such as MIT, BSD-3-Clause, ISC, ECL and Apache 2.0.

## 7. Evaluation metrics

The following metrics are used to evaluate Java Code Generation methods. A brief summary of these metrics is presented in Table 2.

***Exact Match (EM)*** measures the accuracy of the model. It determines if the prediction is the same as the ground truth (also known as reference or gold standard). Even though this metric is easy and fast to compute, and provides straightforward understanding of model's accuracy, it is too rigid that it does not take into account the semantic similarity between reference and prediction codes. It penalizes the model even for white spaces and the variable names. Naturally, this is the most difficult metric to satisfy.

***BLEU*** (Papineni et al., 2002) was originally introduced to evaluate machine translation task. It measures the percentage of overlapped n-grams between the ground truth and the prediction generated by the model. It considers some criteria to penalize the score, such as the sentence length. Overall, this metric is less rigid than EM, and accepts small variations between two predicted codes. However, since the Ngrams order is not considered, it also ignores the semantic of the generated code. Also, it can overestimate the model's score.

***CodeBLEU*** (Ren et al., 2020) uses the n-gram match score of the BLEU metric between the ground truth and the generated sequence. Furthermore, it takes into account the code syntax, and code semantics by matching the AST and DFG, respectively. This metric is the most accurate one so far and is better adapted for code generation. Unlike BLEU metric, it can be too strict and underestimate the model's score.

## 8. Results and Discussion

***Main Results -*** Table 3 displays the hyper-parameters utilized by state-of-the-art techniques during the pretraining and fine-tuning stages to evaluate the CONCODE dataset. Meanwhile, Table 4 presents the outcomes of diverse approaches on the same dataset. Since CONCODE was released in 2018, it has been the main benchmark to evaluate Java code generation task. The first two lines of the table

| Metric | Proposed for | Merits | Limitations |
|---|---|---|---|
| Exact Match (EM) | - code completion<br>- code summarization | - easy to compute<br>- fastest in inference<br>- easy to interpret | - rigid metric<br>- semantic not considered<br>- penalizes even white spaces |
| BLEU (Papineni et al., 2002) | machine translation | - less rigid that EM<br>- easy to compute<br>- accepts multiple correct answers | - Ngram order not considered<br>- semantic not considered<br>- can be too optimistic |
| CodeBLEU (Ren et al., 2020) | code generation | - more accurate<br>- better adapted for code generation<br>- better handling of semantic | - slower in inference<br>- more difficult to compute<br>- can be too strict |

**Table 2**
Summary of metrics used in literature

| Model | Input / Output Length | LR | Optimizer | Dropout | Batch Size (B) | Training steps | Vocab Size |
|---|---|---|---|---|---|---|---|
| GPT-2 pretraining | 1024 | 5e-5 | Adam | 0.1 | 32 | - | 50,000 |
| PLBART pretraining | 512 /512 | 1e-6 | Adam | 0.1, 0.5 | 2048 | 100K | 50,004 |
| PLBART fine-tuning | 512 /512 | 3e-5 | Adam | 0.1 | 32 | 100K | 50,004 |
| CoTexT pretraining | 1024 /1024 | 1e-3 | Adam | - | 128 | 200K | 32,000 |
| CoTexT-base fine-tuning | 256 / 256 | 1e-3 | Adam | - | 128 | 45K | 32,000 |
| JaCoText-base fine-tuning | 379 / 379 | 1e-3 | Adam | - | 128 | 60K | 32,000 |
| CodeBERT fine-tuning | 512 | 5e-5 | Adam | 0.1 | 32 | | 50,265 |
| GraphCodeBERT fine-tuning | 256 | 1e-4 | Adam | - | 64 | | - |
| CodeT5-base pretraining | 512 / 256 | 2e-4 | Adam | - | 1024 | 150 epochs | 32,100 |
| CodeT5-base fine-tuning | 320 / 150 | 5e-5 | Adam | - | 32 | 30 epochs | 32,100 |
| StructCoder pretraining | 400 / 400 | 5e-5 | AdamW | - | 32 | 12K | 32,100 |
| StructCoder fine-tuning | 325 / 155 | 5e-5 | AdamW | - | 32 | 100K | 32,100 |

**Table 3**
Hyper-parameters used by state-of-the-art methods during pretraining and fine tuning

correspond to results of RNN-based methods. The latter obtain the lowest scores. This finding is not surprising because recurrent neural networks are not as large as Transformers, and have limited generalization capability. In addition, they suffer often from vanishing and exploding gradient problems (Hochreiter, 1998; Squartini et al., 2003; Fadziso, 2020).

As for the Transformer-based methods, most of encoder-decoder methods are based on T5 (namely CodeT5, CoTexT, JaCoText, and StructCoder) or BART (namely PLBART) models. In the case of T5-base, the scores are obtained after fine-tuning directly from the pretrained model. However, PLBART, CoTexT and JaCoText were pretrained using source code datasets, and are based on previously trained models.

Encoder-only methods are generally based on BERT. Authors of CodeBERT and GraphCodeBERT did not originally used them for Java Code Generation on CONCODE dataset. Instead, they used them in other tasks such as code search and code documentation generation. Results presented here were reported by Parvez et al. (2021) who used the models for code generation. Note that CodeBERT (Feng et al., 2020) is the first model to be pretrained on various programming languages.

Decoder-only models presented by Lu et al. (2021) are versions based on GPT-2 (Radford et al., 2019). The first model is directly fine-tuned from GPT-2 weights, while the second and third models were trained on the CodeSearch-Net dataset. The main difference is that CodeGPT-adapted initializes the training process from the GPT-2 checkpoint, and CodeGPT-2 performs a training from scratch.

Results show that performance of the Transformer-based models varies significantly. The highest overall results are obtained with the very recent state-of-the-art methods that are based on T5 (Raffel et al., 2020) and GPT-2 (Radford et al., 2019) models. The best model in terms of BLEU and CodeBLEU metrics is CodeT5-large. This is intuitive insofar as it has the largest number of parameters (770 M). However, this boost in performance comes at an expense of longer training time due to the large complexity of the model. The best model in terms of Exact Match metric is REDCODER (Parvez et al., 2021), indicating that it generates code that is more similar syntactically to the ground-truth code. It is noteworthy to mention that, on the one hand, results are very comparable between CodeT5-large and REDCODER, while on the other hand, there is a huge difference between their number of parameters. In fact, CodeT5-large is x5.5 larger than REDCODER. Depending on the application domain and on the available computational resources, a compromise should be done in order to prioritize the gain in memory or the gain in performance.

CodeBERT (Feng et al., 2020) and GraphCodeBERT (Guo et al., 2021) have also achieved relatively high scores, but not as high as the CodeT5-large and REDCODER. This is likely due to the fact that they have fewer number of parameters (125M, 110M versus 770M and 140M, respectively). Interestingly, REDCODER outperforms CodeGPT-2 even thought it has less than half of its parameters only (140 M vs 345 M). This proves the effectiveness of retrieval and generation models SCODER-R and SCODER-G, the main building blocks of REDCODER.

StructCoder (Tipirneni et al., 2022) and JaCoText (Lopez Espejel et al., 2023) have very comparable EM (22.35 vs 22.15) and BLEU (40.91 vs 39.07) results, with StructCoder being slightly better than JaCoText. The latter improves the results of CoText by an average of 2.5 points. JaCotext increments CotexT by adding a pretraining phase in which large input sequence were used. Finally, UniXcoder, CodeT5-small, and $T5_{base}$ achieve the lowest results in Transformer-based methods. Overall Code-T5-small is the model that best balances the compromise between performance and number of parameters, and is preferable is stricter application domains.

Last but not least, results show that using pretraining leads to a better performance than training from scratch. This can be seen when comparing the performances of CodeGPT-adapted and CodeGPT-2. As mentioned previously, the former is initialized from GPT-1 checkpoints, and the latter is trained from scratch. This finding was also confirmed by Feng et al. (2020) and Rozière et al. (2021). Therefore, we recommend using pretrained models whose representation is already stable and can facilitate the training of downstreaming tasks.

Finally, since the highest achieved score on CONCODE dataset is 45.08, there is still a large gap between the best Java code generation methods and the ground truth. This means that more effort is needed to tackle this task.

***General Discussion -*** Over the years, generating code through natural language has been a challenging task. The main goal of this task is to increase programmers' productivity by automating tedious and repetitive coding tasks. In this paper, we present advancements in code generation within the Java programming language. However, in order to fully comprehend the subject matter, it is crucial to have a grasp of the historical development of code generation in other fields.

Early attempts to generate code automatically involved translating natural language requests into regular expressions (Angluin, 1987; Ranta, 1998), logical forms based on probabilistic categorial grammars (Zettlemoyer and Collins, 2005b), agent-specific language (Kate et al., 2005). These approaches laid the groundwork for more sophisticated techniques, such as semantic parsing (Wong and Mooney, 2006; Lu et al., 2008) and SQL query generation (Miller et al., 1996; Ramaswamy and Kleindienst, 2000), which were discussed in Section 3.

The previously mentioned methods utilize basic techniques such as rule-based systems and simple statistical methods to overcome the different tasks. However, as time has progressed and technology has advanced, the demand for more sophisticated techniques has increased to tackle the most complex code generation obstacles. One of the techniques that has made significant strides in this field is the utilization of RNNs, such as LSTMs (Hochreiter and Schmidhuber, 1997), and later GRUs (Cho et al., 2014a). These advanced techniques have greatly enhanced the capability to address the challenges of code generation. For instance, Iyer et al. (2018) introduced CONCODE, a dataset with the purpose of generating Java code through natural language. The authors used a sequence-to-sequence model based on RNNs, which was similar to other models used in the field (Neelakantan et al., 2016; Ling et al., 2016). At the time of the development of this model, the integration of the attention mechanism (Gu et al., 2016a) was a significant breakthrough in code and text generation.

Although models based on Recurrent Neural Networks (RNNs) were initially promising, they were hindered by several drawbacks. One of these was the vanishing gradient problem (Hochreiter, 1998; Squartini et al., 2003; Fadziso, 2020), and also they were limited in their ability to effectively process large amounts of input data. As a result, RNN-based models often truncated natural language requests and code sequences, leading to significant information loss and incomplete generated code. Furthermore, due to Java's strict programming rules, the code generated by these models was often riddled with syntax errors and could not be executed.

After the big success of Transformers Vaswani et al. (2017) in NLP tasks, many Java Code Generation researchers shifted their efforts to this neural network architecture, and proposed the above mentioned methods. The most succesful models are based on T5 and GPT-2 models. CodeT5 models use a learning objective task that was specifically designed to take advantage of the programming language structure. In Wang et al. (2021), CodeT5-small and CodeT5-base achieved remarkable improvements in the performance of code generation task using bimodal dual learning objective. Authors leverage the token type of identifiers from the AST such as function names and variables to enrich the semantic features. Wang et al. (2021) studied three objectives: Masked Span Prediction (MSP), Identifier Tagging (IT), and Masked Identifier Prediction (MIP) . It turns out that Masked Span Prediction is the most crucial objective for learning the syntactic information in the generation tasks. However, the simultaneous optimization of the three of them achieve the best performance. Alternatively, Le et al. (2022) show that CodeT5-large obtains higher scores than CodeT5-small and CodeT5-base. This finding is intuitive since the model architecture is more complex.

After the success of CodeT5, StructCoder (Tipirneni et al., 2022) surged. It combines T5-base architecture with ASTs. Moreover, authors use DFGs to make aware both the encoder and decoder of the syntax and the data flow. This system achieved better results than CodeT5 using T5-base. Therefore, AST and DFG objectives are beneficial to learn correctly the semantic and syntax of Java programming

| | Category | Model | Score (↑) | | | # params (↓) |
|---|---|---|---|---|---|---|
| | | | EM | BLEU | CodeBLEU | |
| | RNN-based | Seq2Seq (Yin and Neubig, 2017) | 6.65 | 21.29 | - | - |
| | | Iyer et al. (2018) | 8.60 | 22.11 | - | - |
| | Encoder-only | CodeBERT (Feng et al., 2020) | 18.00 | 28.70 | 31.40 | 125M |
| | | GraphCodeBERT (Guo et al., 2021) | 18.70 | 33.40 | 35.90 | 110M |
| | Decoder-only | GPT-2 (Radford et al., 2019) | 17.35 | 25.37 | 29.69 | 345M |
| | | CodeGPT-2 (Lu et al., 2021) | 18.35 | 28.69 | 32.71 | 124M |
| | | CodeGPT-adapted (Lu et al., 2021) | 20.10 | 32.79 | 35.98 | 124M |
| Transformer-based Methods | | $T5_{base}$ (Raffel et al., 2020) | 18.65 | 32.74 | 35.95 | 220 M |
| | | CodeT5-small (Wang et al., 2021) | 21.55 | 38.13 | 41.39 | |
| | | +dual-gen | 19.95 | 39.02 | 42.21 | 60 M |
| | | +multi-task | 20.15 | 35.89 | 38.83 | |
| | | CodeT5-base (Wang et al., 2021) | 22.30 | 40.73 | 43.20 | |
| | | +dual-gen | 22.70 | 41.48 | 44.10 | 220 M |
| | Encoder-decoder | +multi-task | 21.15 | 37.54 | 40.01 | |
| | | CodeT5-large (Le et al., 2022) | 22.65 | **42.66** | **45.08** | 770 M |
| | | PLBART (Ahmad et al., 2021) | 18.75 | 36.69 | 38.52 | 140 M |
| | | CoTexT (Phan et al., 2021) | 20.10 | 37.40 | 40.14 | 220 M |
| | | JaCoText (Lopez Espejel et al., 2023) | 22.15 | 39.07 | 41.53 | 220 M |
| | | REDCODER (Parvez et al., 2021) | **23.40** | 41.60 | 43.40 | 140 M |
| | | REDCODER-EXT (Parvez et al., 2021) | 23.30 | 42.50 | 43.40 | 140 M |
| | | StructCoder (Tipirneni et al., 2022) | 22.35 | 40.91 | 44.76 | 220 M |
| | | UniXcoder (Guo et al., 2022) | 22.60 | 38.23 | - | 125M |

**Table 4**
Results of Java code generation task on the CONCODE dataset. We present EM, BLEU, and CodeBLEU scores, as well as the number of parameters for each method. Best results are in bold.

language. Another interesting models is REDCODER, the latter is the first model to use two encoders. The first encoder is used for the retrieval module to search for relevant code in the dataset, the second one is part of the generator module (encoder-decoder model) to generate the Java source code. REDCODER achieves the best Exact Match (EM) score.

GPT-3 model in an improvement of GPT-2, introduced by Brown et al. (2020). It is a breakthrough in Large Language Models (LLMs), as it demonstrates the feasibility of scaling the model and achieving competitive scores in various tasks through few-shot learning. This model was the starting point to many variants among which OpenAI's ChatGPT chatbot (OpenIA, 2023). ChatGPT (OpenIA, 2023) is based on the GPT-3.5 model, which has shown remarkable performance in various tasks, such as machine translation in multiple languages, automatic summarization, parts of speech tagging, and even code generation through natural language requests. In particular, ChatGPT has demonstrated impressive proficiency in writing complex Java source code, even for applications.

Despite the impressive results showed by ChatGPT, the latter has many limitations, out of which, we mention the following. (1) ChatGPT is not able to fully understand natural language requests that are not well-formulated (such as missing commas), and (2) ChatGPT is sensitive to the prompts used by users to interact with it. The new version of ChatGPT (called ChatGPT4) was proposed precisely to augment the understanding capacity of it. ChatGPT4 has made significant advancements in Java Code generation from complex queries, as well as in working on multimodal inputs such as text, images, and video.

ChatGPT4 is not the only attempt to improve ChatGPT, other decoder-only competitors are trying hard to outperform it. Two examples of such systems are Google's BARD (Google, 2023) that is based on LaMDA model (Thoppilan et al., 2022), and LLaMA (Touvron et al., 2023) that is based on Chinchilla model. It is noteworthy to mention that all very recent LLMs have not been yet fully benchmarked on any of state-of-the-art datasets for Java code generation, at the time of writing this article.

## 9. Conclusions and Perspectives

In this paper, we presented a comparison review of state-of-the-art methods in Java Code Generation using CONCODE dataset. Many methods surged during the last years and we categorized them based on their architecture: (1) RNN-based methods, and (2) Transformer-based methods. The latter are divided into three subcategories: encoder-based methods, decoder-based methods, and encoder-decoder-based methods. We propose the following future directions:

- It would be interesting to use beam search algorithm to detect syntax errors, similarly to what is done in the semantic parser PICARD in SQL code generation. (Scholak et al., 2021)

- One big drawback of current code generation metrics is their high dependency on syntax similarity between reference and predicted codes. It is really important and right time to start focusing on handling the semantic part of it. Even though that two codes are

syntactically different, they still can perform the same task.

- One straightforward line of work is to exploit more Java datasets such as AlphaCode (Li et al., 2022), PolyCoder (Xu et al., 2022), and AlphaCode (Li et al., 2022) in the generation process, especially at the fine tuning step. This is because current methods only use Java data to pretrain the model, while they fine-tune models on Python benchmarks.

- Despite the fact that large language models such as ChatGPT, Bard and LLaMA provide extraordinary performances in code generation, their number of parameters is huge. Future research would focus on the compression of such language models, or the proposal of smaller-scale efficient language models. This will be very useful for researcher and companies that do not have enough computational resources or memory (Tao et al., 2022).

- Alternatively, another perspective would be to work on training time optimization, especially for large models such as ChatGPT. For instance, continual learning algorithms can be used in order to integrate new data on-the-fly, without retraining the model from scratch at each update (Gao et al., 2023).

- Current code generation models lack of reasoning abilities. Another perspective would be to develop more the common sense of language models to improve their generation semantic, thus, their performance (Huang and Chang, 2022).

# References

Ahmad, W., Chakraborty, S., Ray, B., Chang, K.W., 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.

Angluin, D., 1987. Learning regular sets from queries and counterexamples. Information and Computation 75, 87–106. URL: https://www.sciencedirect.com/science/article/pii/0890540187900526, doi:https://doi.org/10.1016/0890-5401(87)90052-6.

Ba, J.L., Kiros, J.R., Hinton, G.E., 2016. Layer normalization. URL: http://arxiv.org/abs/1607.06450. cite arxiv:1607.06450.

Black, S., Biderman, S., Hallahan, E., Anthony, Q., Gao, L., Golding, L., He, H., Leahy, C., McDonell, K., Phang, J., Pieler, M., Prashanth, U.S., Purohit, S., Reynolds, L., Tow, J., Wang, B., Weinbach, S., 2022. GPT-NeoX-20B: An open-source autoregressive language model, in: Proceedings of BigScience Episode #5 – Workshop on Challenges & Perspectives in Creating Large Language Models.

Branavan, S., Zettlemoyer, L., Barzilay, R., 2010. Reading between the lines: Learning to map high-level instructions to commands, in: Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics, Association for Computational Linguistics, Uppsala, Sweden. pp. 1268–1277. URL: https://aclanthology.org/P10-1129.

Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., Amodei, D., 2020. Language models are few-shot learners, in:

Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., Lin, H. (Eds.), Advances in Neural Information Processing Systems, Curran Associates, Inc.. pp. 1877–1901. URL: https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf.

Caballero, E., OpenAI, ., Sutskever, I., 2016. Description2Code Dataset. URL: https://github.com/ethancaballero/description2code, doi:10.5281/zenodo.5665051.

Chen, M., Tworek, J., Jun, H., Yuan, Q., Ponde, H., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., Ray, A., Puri, R., Krueger, G., Petrov, M., Khlaaf, H., Sastry, G., Mishkin, P., Chan, B., Gray, S., Ryder, N., Pavlov, M., Power, A., Kaiser, L., Bavarian, M., Winter, C., Tillet, P., Such, F.P., Cummings, D.W., Plappert, M., Chantzis, F., Barnes, E., Herbert-Voss, A., Guss, W.H., Nichol, A., Babuschkin, I., Balaji, S.A., Jain, S., Carr, A., Leike, J., Achiam, J., Misra, V., Morikawa, E., Radford, A., Knight, M.M., Brundage, M., Murati, M., Mayer, K., Welinder, P., McGrew, B., Amodei, D., McCandlish, S., Sutskever, I., Zaremba, W., 2021a. Evaluating large language models trained on code. ArXiv abs/2107.03374.

Chen, M., Tworek, J., Jun, H., Yuan, Q., Ponde, H., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., Ray, A., Puri, R., Krueger, G., Petrov, M., Khlaaf, H., Sastry, G., Mishkin, P., Chan, B., Gray, S., Ryder, N., Pavlov, M., Power, A., Kaiser, L., Bavarian, M., Winter, C., Tillet, P., Such, F.P., Cummings, D.W., Plappert, M., Chantzis, F., Barnes, E., Herbert-Voss, A., Guss, W.H., Nichol, A., Babuschkin, I., Balaji, S.A., Jain, S., Carr, A., Leike, J., Achiam, J., Misra, V., Morikawa, E., Radford, A., Knight, M.M., Brundage, M., Murati, M., Mayer, K., Welinder, P., McGrew, B., Amodei, D., McCandlish, S., Sutskever, I., Zaremba, W., 2021b. Evaluating large language models trained on code. ArXiv abs/2107.03374.

Cho, K., van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., Bengio, Y., 2014a. Learning phrase representations using RNN encoder–decoder for statistical machine translation, in: Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 1724–1734.

Cho, K., van Merrienboer, B., Gülçehre, Ç., Bahdanau, D., Bougares, F., Schwenk, H., Bengio, Y., 2014b. Learning phrase representations using RNN encoder-decoder for statistical machine translation, in: Moschitti, A., Pang, B., Daelemans, W. (Eds.), Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL, ACL. pp. 1724–1734.

Clark, K., Luong, M., Le, Q.V., Manning, C.D., 2020. ELECTRA: pre-training text encoders as discriminators rather than generators. 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020 URL: https://openreview.net/forum?id=r1xMH1BtvB.

Dai, Z., Yang, Z., Yang, Y., Carbonell, J.G., Le, Q.V., Salakhutdinov, R., 2019. Transformer-xl: Attentive language models beyond a fixed-length context, in: Korhonen, A., Traum, D.R., Màrquez, L. (Eds.), ACL (1), Association for Computational Linguistics. pp. 2978–2988.

Devlin, J., Chang, M., Lee, K., Toutanova, K., 2019. BERT: pre-training of deep bidirectional transformers for language understanding, in: Burstein, J., Doran, C., Solorio, T. (Eds.), Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers), Association for Computational Linguistics. pp. 4171–4186.

Driess, D., Xia, F., Sajjadi, M.S.M., Lynch, C., Chowdhery, A., Ichter, B., Wahid, A., Tompson, J., Vuong, Q., Yu, T., Huang, W., Chebotar, Y., Sermanet, P., Duckworth, D., Levine, S., Vanhoucke, V., Hausman, K., Toussaint, M., Greff, K., Zeng, A., Mordatch, I., Florence, P., 2023. Palm-e: An embodied multimodal language model. arXiv preprint arXiv:2303.03378 .

Fadziso, T., 2020. Overcoming the vanishing gradient problem during learning recurrent neural nets (rnn). Asian Journal of Applied Science and Engineering 9, 207–218.

Feng, Z., Guo, D., Tang, D., Duan, N., Feng, X., Gong, M., Shou, L., Qin, B., Liu, T., Jiang, D., Zhou, M., 2020. Codebert: A pre-trained model

for programming and natural languages. URL: http://arxiv.org/abs/2002.08155. cite arxiv:2002.08155Comment: Accepted to Findings of EMNLP 2020. 12 pages.

Gao, S., Zhang, H., Gao, C., Wang, C., 2023. Keeping pace with ever-increasing data: Towards continual learning of code intelligence models. CoRR abs/2302.03482.

Google, 2023. An overview of bard: an early experiment with generative ai. URL: https://ai.google/static/documents/google-about-bard.pdf.

Graves, A., Wayne, G., Danihelka, I., 2014. Neural turing machines. arxiv Cite arxiv:1410.5401.

Gu, J., Lu, Z., Li, H., Li, V.O.K., 2016a. Incorporating copying mechanism in sequence-to-sequence learning. CoRR abs/1603.06393.

Gu, X., Zhang, H., Zhang, D., Kim, S., 2016b. Deep api learning, in: Proceedings of the 2016 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering, Association for Computing Machinery, New York, NY, USA. p. 631–642. URL: https://doi.org/10.1145/2950290.2950334, doi:10.1145/2950290.2950334.

Guo, D., Lu, S., Duan, N., Wang, Y., Zhou, M., Yin, J., 2022. Unixcoder: Unified cross-modal pre-training for code representation. arxiv .

Guo, D., Ren, S., Lu, S., Feng, Z., Tang, D., Liu, S., Zhou, L., Duan, N., Svyatkovskiy, A., Fu, S., Tufano, M., Deng, S.K., Clement, C.B., Drain, D., Sundaresan, N., Yin, J., Jiang, D., Zhou, M., 2021. Graphcodebert: Pre-training code representations with data flow. 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021 URL: https://openreview.net/forum?id=jLoC4ez43PZ.

Hernandez, D., Kaplan, J., Henighan, T., McCandlish, S., 2021. Scaling laws for transfer. arXiv preprint arXiv:2102.01293 .

Hochreiter, S., 1998. The vanishing gradient problem during learning recurrent neural nets and problem solutions. International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems 6, 107–116.

Hochreiter, S., Schmidhuber, J., 1997. Long short-term memory. Neural computation 9, 1735–1780.

Hoffmann, J., Borgeaud, S., Mensch, A., Buchatskaya, E., Cai, T., Rutherford, E., de Las Casas, D., Hendricks, L.A., Welbl, J., Clark, A., Hennigan, T., Noland, E., Millican, K., van den Driessche, G., Damoc, B., Guy, A., Osindero, S., Simonyan, K., Elsen, E., Rae, J.W., Vinyals, O., Sifre, L., 2022. Training Compute-Optimal Large Language Models. arXiv e-prints .

Huang, J., Chang, K.C.C., 2022. Towards reasoning in large language models: A survey. arXiv preprint arXiv:2212.10403 .

Husain, H., Wu, H., Gazit, T., Allamanis, M., Brockschmidt, M., 2019. Codesearchnet challenge: Evaluating the state of semantic code search. CoRR abs/1909.09436. URL: http://arxiv.org/abs/1909.09436, arXiv:1909.09436.

Inc., A., 2022. Atcoder. https://atcoder.jp/. Accessed: 2022-10-28.

Iyer, S., Konstas, I., Cheung, A., Zettlemoyer, L., 2018. Mapping language to code in programmatic context, in: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Brussels, Belgium. pp. 1643–1652. URL: https://aclanthology.org/D18-1192, doi:10.18653/v1/D18-1192.

Joulin, A., Mikolov, T., 2015. Inferring algorithmic patterns with stack-augmented recurrent nets. Advances in Neural Information Processing Systems 28. URL: https://proceedings.neurips.cc/paper_files/paper/2015/file/26657d5ff9020d2abefe558796b99584-Paper.pdf.

Kate, R.J., Wong, Y.W., Mooney, R.J., 2005. Learning to transform natural to formal languages, in: Proceedings of the 20th National Conference on Artificial Intelligence - Volume 3, AAAI Press. p. 1062–1068.

Kingma, D.P., Ba, J., 2015. Adam: A method for stochastic optimization. 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings URL: http://arxiv.org/abs/1412.6980.

Kocetkov, D., Li, R., Ben Allal, L., Li, J., Mou, C., Muños Ferrandis, C., Hughes, S., Wolf, T., Bahdanau, D., von Werra, L., de Vries, H., 2022. The stack: 3 tb of permissively licensed source code. arxiv .

Kumar, A., Irsoy, O., Ondruska, P., Iyyer, M., Bradbury, J., Gulrajani, I., Zhong, V., Paulus, R., Socher, R., 2016. Ask me anything: Dynamic memory networks for natural language processing, in: Balcan, M.F.,

Weinberger, K.Q. (Eds.), Proceedings of The 33rd International Conference on Machine Learning, PMLR, New York, New York, USA. pp. 1378–1387. URL: https://proceedings.mlr.press/v48/kumar16.html.

Kushman, N., Barzilay, R., 2013a. Using semantic unification to generate regular expressions from natural language, in: Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Association for Computational Linguistics, Atlanta, Georgia. pp. 826–836. URL: https://aclanthology.org/N13-1103.

Kushman, N., Barzilay, R., 2013b. Using semantic unification to generate regular expressions from natural language, in: Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Association for Computational Linguistics, Atlanta, Georgia. pp. 826–836. URL: https://aclanthology.org/N13-1103.

Le, H., Wang, Y., Gotmare, A.D., Savarese, S., Hoi, S.C.H., 2022. Coderl: Mastering code generation through pretrained models and deep reinforcement learning. arXiv preprint arXiv:2207.01780 .

Lewis, M., Liu, Y., Goyal, N., Ghazvininejad, M., Mohamed, A., Levy, O., Stoyanov, V., Zettlemoyer, L., 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, Association for Computational Linguistics, Online. pp. 7871–7880. URL: https://aclanthology.org/2020.acl-main.703, doi:10.18653/v1/2020.acl-main.703.

Li, Y., Choi, D., Chung, J., Kushman, N., Schrittwieser, J., Leblond, R., Eccles, T., Keeling, J., Gimeno, F., Dal Lago, A., Hubert, T., Choy, P., de Masson d'Autume, C., Babuschkin, I., Chen, X., Huang, P.S., Welbl, J., Gowal, S., Cherepanov, A., Molloy, J., Mankowitz, D., Sutherland Robson, E., Kohli, P., de Freitas, N., Kavukcuoglu, K., Vinyals, O., 2022. Competition-level code generation with alphacode. arXiv preprint arXiv:2203.07814 .

Ling, W., Blunsom, P., Grefenstette, E., Hermann, K.M., Kočiský, T., Wang, F., Senior, A., 2016. Latent predictor networks for code generation, in: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 599–609.

Lopez Espejel, J., Alassan, M.S.Y., Dahhane, W., Ettifouri, E.H., 2023. Jacotext: A pretrained model for java code-text generation. International Journal of Computer and Systems Engineering 17, 100 – 105.

Loshchilov, I., Hutter, F., 2019. Decoupled weight decay regularization. 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019 URL: https://openreview.net/forum?id=Bkg6RiCqY7.

Lu, S., Guo, D., Ren, S., Huang, J., Svyatkovskiy, A., Blanco, A., Clement, C.B., Drain, D., Jiang, D., Tang, D., Li, G., Zhou, L., Shou, L., Zhou, L., Tufano, M., Gong, M., Zhou, M., Duan, N., Sundaresan, N., Deng, S.K., Fu, S., Liu, S., 2021. Codexglue: A machine learning benchmark dataset for code understanding and generation. ArXiv abs/2102.04664.

Lu, W., Ng, H.T., Lee, W.S., Zettlemoyer, L.S., 2008. A generative model for parsing natural language to meaning representations, in: Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Honolulu, Hawaii. pp. 783–792. URL: https://aclanthology.org/D08-1082.

Miller, S., Stallard, D., Bobrow, R., Schwartz, R., 1996. A fully statistical approach to natural language interfaces, in: 34th Annual Meeting of the Association for Computational Linguistics, Association for Computational Linguistics, Santa Cruz, California, USA. pp. 55–61. URL: https://aclanthology.org/P96-1008, doi:10.3115/981863.981871.

Mirzayanov, M., 2020. Codeforces: Results of 2020. URL: https://codeforces.com/blog/entry/89502.

Mou, L., Men, R., Li, G., Zhang, L., Jin, Z., 2015. On end-to-end program generation from user intention by deep neural networks. CoRR abs/1510.07211. URL: http://dblp.uni-trier.de/db/journals/corr/corr1510.html#MouMLZJ15.

Neelakantan, A., Le, Q.V., Sutskever, I., 2016. Neural programmer: Inducing latent programs with gradient descent. 4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings URL: http://arxiv.org/abs/

1511.04834.

Nocedal, J., 1980. Updating quasi-newton matrices with limited storage. Mathematics of computation 35, 773–782.

OpenIA, 2023. Gpt-4 technical report. arxiv URL: https://arxiv.org/pdf/2303.08774.pdf.

Oracle, 2022. javadoc. https://docs.oracle.com/javase/8/docs/technotes/tools/windows/javadoc.html. Accessed: 2022-09-05.

Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C.L., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., Schulman, J., Hilton, J., Kelton, F., Miller, L.E., Simens, M., Askell, A., Welinder, P., Christiano, P.F., Leike, J., Lowe, R.J., 2022. Training language models to follow instructions with human feedback. ArXiv abs/2203.02155.

Pang, B., Hayashi, H., Tu, L., Wang, H., Zhou, Y., Savarese, Silvio Xiong, C., 2022. Codegen: An open large language model for code with multi-turn program synthesis. ArXiv .

Papineni, K., Roukos, S., Ward, T., Zhu, W.J., 2002. Bleu: a method for automatic evaluation of machine translation, in: Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, Association for Computational Linguistics, Philadelphia, Pennsylvania, USA. pp. 311–318. URL: https://aclanthology.org/P02-1040, doi:10.3115/1073083.1073135.

Parvez, M.R., Ahmad, W.U., Chakraborty, S., Ray, B., Chang, K., 2021. Retrieval augmented code generation and summarization. CoRR abs/2108.11601.

Phan, L., Tran, H., Le, D., Nguyen, H., Annibal, J., Peltekian, A., Ye, Y., 2021. CoTexT: Multi-task learning with code-text transformer, in: Proceedings of the 1st Workshop on Natural Language Processing for Programming (NLP4Prog 2021).

Puri, R., Kung, D., Janssen, G., Zhang, W., Domeniconi, G., Zolotov, V., Dolby, J., Chen, J., Choudhury, M., Decker, L., Thost, V., Buratti, L., Pujar, S., Finkler, U., 2021a. Project codenet: A large-scale ai for code dataset for learning a diversity of coding tasks. arxiv .

Puri, R., Kung, D.S., Janssen, G., Zhang, W., Domeniconi, G., Zolotov, V., Dolby, J., Chen, J., Choudhury, M., Decker, L., Thost, V., Buratti, L., Pujar, S., Ramji, S., Finkler, U., Malaika, S., Reiss, F., 2021b. Project codenet: A large-scale ai for code dataset for learning a diversity of coding tasks. ArXiv arXiv:2105.12655.

Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., 2019. Language models are unsupervised multitask learners. arxiv .

Rae, J.W., Borgeaud, S., Cai, T., Millican, K., Hoffmann, J., Song, F., Aslanides, J., Henderson, S., Ring, R., Young, S., Rutherford, E., Hennigan, T., Menick, J., Cassirer, A., Powell, R., van den Driessche, G., Hendricks, L.A., Rauh, M., Huang, P.S., Glaese, A., Welbl, J., Dathathri, S., Huang, S., Uesato, J., Mellor, J., Higgins, I., Creswell, A., McAleese, N., Wu, A., Elsen, E., Jayakumar, S., Buchatskaya, E., Budden, D., Sutherland, E., Simonyan, K., Paganini, M., Sifre, L., Martens, L., Li, X.L., Kuncoro, A., Nematzadeh, A., Gribovskaya, E., Donato, D., Lazaridou, A., Mensch, A., Lespiau, J.B., Tsimpoukelli, M., Grigorev, N., Fritz, D., Sottiaux, T., Pajarskas, M., Pohlen, T., Gong, Z., Toyama, D., de Masson d'Autume, C., Li, Y., Terzi, T., Mikulik, V., Babuschkin, I., Clark, A., de Las Casas, D., Guy, A., Jones, C., Bradbury, J., Johnson, M., Hechtman, B., Weidinger, L., Gabriel, I., Isaac, W., Lockhart, E., Osindero, S., Rimell, L., Dyer, C., Vinyals, O., Ayoub, K., Stanway, J., Bennett, L., Hassabis, D., Kavukcuoglu, K., Irving, G., 2021. Scaling Language Models: Methods, Analysis & Insights from Training Gopher. arXiv e-prints , arXiv:2112.11446doi:10.48550/arXiv.2112.11446, arXiv:2112.11446.

Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., Liu, P.J., 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. Journal of Machine Learning Research 21, 1–67. URL: http://jmlr.org/papers/v21/20-074.html.

Ramaswamy, G.N., Kleindienst, J., 2000. Hierarchical feature-based translation for scalable natural language understanding. Sixth International Conference on Spoken Language Processing .

Ranta, A., 1998. A multilingual natural-language interface to regular expressions, in: Proceedings of the International Workshop on Finite State Methods in Natural Language Processing, Association for Computational Linguistics, USA. p. 79–90.

Reed, S.E., de Freitas, N., 2016. Neural programmer-interpreters. 4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings URL: http://arxiv.org/abs/1511.06279.

Ren, S., Guo, D., Lu, S., Zhou, L., Liu, S., Tang, D., Zhou, M., Blanco, A., Ma, S., 2020. Codebleu: a method for automatic evaluation of code synthesis. ArXiv abs/2009.10297.

Rozière, B., Lachaux, M., Szafraniec, M., Lample, G., 2021. DOBF: A deobfuscation pre-training objective for programming languages. CoRR abs/2102.07492.

Roziere, B., Lachaux, M.A., Chanussot, L., Lample, G., 2020. Unsupervised translation of programming languages, in: Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., Lin, H. (Eds.), Advances in Neural Information Processing Systems, Curran Associates, Inc.. pp. 20601–20611. URL: https://proceedings.neurips.cc/paper_files/paper/2020/file/ed23fbf18c2cd35f8c7f8de44f85c08d-Paper.pdf.

Scholak, T., Schucher, N., Bahdanau, D., 2021. PICARD: Parsing incrementally for constrained auto-regressive decoding from language models, in: Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics. pp. 9895–9901. URL: https://aclanthology.org/2021.emnlp-main.779.

Squartini, S., Hussain, A., Piazza, F., 2003. Preprocessing based solution for the vanishing gradient problem in recurrent neural networks, in: Proceedings of the 2003 International Symposium on Circuits and Systems, 2003. ISCAS'03., IEEE. pp. V–V.

Tao, C., Hou, L., Zhang, W., Shang, L., Jiang, X., Liu, Q., Luo, P., Wong, N., 2022. Compression of generative pre-trained language models via quantization. CoRR abs/2203.10705.

Thoppilan, R., Freitas, D., Hall, J., Shazeer, N., Kulshreshtha, A., Cheng, H.T., Jin, A., Bos, T., Baker, L., Du, Y., Li, Y., Lee, H., Zheng, H., Ghafouri, A., Menegali, M., Huang, Y., Krikun, M., Lepikhin, D., Qin, J., Le, Q., 2022. Lamda: Language models for dialog applications. arXiv .

Tipirneni, S., Zhu, M., Reddy, C.K., 2022. Structcoder: Structure-aware transformer for code generation. CoRR abs/2206.05239. URL: https://doi.org/10.48550/arXiv.2206.05239, doi:10.48550/arXiv.2206.05239, arXiv:2206.05239.

Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al., 2023. Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971 .

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I., 2017. Attention is all you need, in: Advances in Neural Information Processing Systems, pp. 5998–6008.

Wang, Y., Wang, W., Joty, S., Hoi, S.C., 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.

Watanobe, Y., 2022. Aizu online judge. https://onlinejudge.u-aizu.ac.jp. Accessed: 2022-10-28.

Wong, Y.W., Mooney, R., 2006. Learning for semantic parsing with statistical machine translation, in: Proceedings of the Human Language Technology Conference of the NAACL, Main Conference, Association for Computational Linguistics, New York City, USA. pp. 439–446. URL: https://aclanthology.org/N06-1056.

Xu, F.F., Alon, U., Neubig, G., Hellendoorn, V.J., 2022. A systematic evaluation of large language models of code, in: Proceedings of the 6th ACM SIGPLAN International Symposium on Machine Programming, Association for Computing Machinery, New York, NY, USA. p. 1–10.

Yin, P., Deng, B., Chen, E., Vasilescu, B., Neubig, G., 2018. Learning to mine aligned code and natural language pairs from stack overflow, in: Proceedings of the 15th International Conference on Mining Software Repositories, Association for Computing Machinery, New York, NY, USA. p. 476–486. URL: https://doi.org/10.1145/3196398.3196408, doi:10.1145/3196398.3196408.

Yin, P., Neubig, G., 2017. A syntactic neural model for general-purpose code generation, in: Proceedings of the 55th Annual Meeting of the

Association for Computational Linguistics (Volume 1: Long Papers), pp. 440–450.

Zettlemoyer, L.S., Collins, M., 2005a. Learning to map sentences to logical form: Structured classification with probabilistic categorial grammars, in: Proceedings of the Twenty-First Conference on Uncertainty in Artificial Intelligence, AUAI Press, Arlington, Virginia, USA. p. 658–666.

Zettlemoyer, L.S., Collins, M., 2005b. Learning to map sentences to logical form: Structured classification with probabilistic categorial grammars, in: Proceedings of the Twenty-First Conference on Uncertainty in Artificial Intelligence, AUAI Press, Arlington, Virginia, USA. p. 658–666.

Zhang, B., Sennrich, R., 2019. Root Mean Square Layer Normalization. Advances in Neural Information Processing Systems 32 URL: https://openreview.net/references/pdf?id=S1qBAf6rr.

Zhang, J., Zhao, Y., Saleh, M., Liu, P.J., 2020. Pegasus: Pre-training with extracted gap-sentences for abstractive summarization, in: Proceedings of the 37th International Conference on Machine Learning, pp. 11328–11339.