1. For the text below, summarise the problem and suggest a software system architecture solution. The main stakeholders in this case are: - Software system architects - DevOps engineers - Developers Problem, context, and related work Software engineering today is a challenging task. As the complexity of software keeps increasing, software engineers find it ever more challenging to cope with this complexity and produce high-quality source code. Furthermore, developers today spend a significant amount of their time solving similar problems as they struggle with code reuse (e.g. they are aware that the solution exists but it requires tedious changes so re-writing a new code seems easier, and they can misinterpret the requirements, or they are not aware if a solution is available and being used in a different context). Ideally, developers could profit from systems that assist them during the development process by suggesting already existing good (and tested) solutions. This could reduce development efforts. The developers could further benefit from systems that have the capabilities to understand the software requirements. This would enable the creation of assistance mechanisms that can index problem and solution domains and be aware if for an existing or similar problem there exist solutions. We are focusing on supporting developers when dealing with these challenges: 1. Avoidance of repetitive tasks (realizing code that already exists increases time efficiency during development and keeps developers motivated) 2. Avoiding bugs in new code (tested code is often far more valuable because potentially many bugs have been discovered and removed) These challenges are both a drain on motivation that developers have (e.g., influence their productivity and motivation to keep the existing work), and introduce a potential overhead due to all the activities that are related to testing newly implemented solutions. Even if solutions already exist, developers still would need to invest time in understanding them and reusing them. In particular, how well these solutions fit certain requirements. Because understanding source code is a time-consuming task, developers often find it quicker to write the new code. Ideally, there would exist solutions that can understand existing source code (and logic behind it) and relate it to the problems they are solving (requirements). AI is a technology that simulates human intelligence processes through computing systems. In an ideal case, AI is capable to rationalize and take actions towards a specific goal. AI, in this context, is interesting because it can recognize patterns in textual data. In software engineering, both requirements and source code are represented mostly in textual form. Therefore, if we can train AI models to understand requirements and to produce results in a restricted format that matches the semantics of programming languages, this could enable the creation of solutions that assist developers. Ideally, this could lead to the use of AI-enabled technologies as assistance tools. However, this is not an easy task. Gathering requirements for human is challenging as the requirement elicitation process need experience and constant communication between the stakeholders for understanding requirements. Expressing requirements is even hard for software engineers, as they imply many things (e.g. different use cases). They acquire customers' requirements, refine them into tasks and then use AI to assist them with the solutions. Another challenging task is to prioritize goals for development which are often not represented in the requirements and need an expert opinion. When developing software solutions, developers often can choose between various ways of implementation, different coding styles, various algorithms, and different programming languages. For example, if developers want to optimize for a reduced memory footprint, they might use processing algorithms that need more processing power but keep fewer states alive in memory. Finally, it is necessary to test the implementation. This is a challenge, as developers need to be skilled in the implemented technology stack. Gathering requirements can benefit from AI techniques, especially from NLP that can be used to automate this process by extracting valuable information from user requirements by semantic analysis. In an ideal scenario, AI can understand requirements and generate portions of the code to assist developers. These portions can be code snippets or function templates. The focus of this work is

on generating full functions to assist developers. While research in AI and software engineering is potentially providing a variety of solutions in this domain, there still do not exist solutions that can outmatch humans in terms of building program logic and fully understanding the requirements. Even more, although some very impressive solutions can generate source code (e.g., GPT-3), there are no practical solutions that can integrate with development practices and assist developers with writing particularly interesting code and choosing optimal solutions based on the code context and requirements. Therefore, we define our problem statement as follows: Are AI-enabled solutions mature enough to integrate with development processes, index problem and solutions domains, and assist developers by suggesting solutions to common problems? If this is not the case, what is preventing the adoption of AI-based solutions in this way, in software engineering? Therefore, in this thesis, we aim to identify challenges that prevent the adoption of AI from being used as an assistance tool in software development by providing an elaborate list of identified issues and an AI-based approach that can aid software development by generating solutions for already solved problems using user requirements as input. This considers the understanding of how the machine learning algorithms associate source code with natural language and how it understands it. Although there have been several studies focused on source code generation tools that use AI ("Imam, Ayad Tareq and Rousan, Thamer and Aljawarneh, Shadi; An expert code generator using rule-based and frames knowledge representation techniques; 2014", "Soni, Naitik Chetan and Pawar, Dhruv Ashok and Tambe, Namita Sandeep and Deolekar, Rugved Vivek; Automated system for code generation from unstructured algorithm; 2016", "Shah, Meet and Shenoy, Rajat and Shankarmani, Radha; Natural Language to Python Source Code using Transformers; 2021", "Lu, Chao and Cao, Xin and Zhu, Yi and Huang, Tao and Pan, Zhaoming and Li, Xiu; Procedure2Command: an AI-based Nuclear Power Plant Control Command Code Generation Prototype System; 2020", "Liliane do Nascimento Vale and Marcelo de Almeida Maia; Towards a question answering assistant for software development using a transformer-based language model; 2021", "Zhang, Chen and Niu, Xinyi and Yu, Bin; A Method of Automatic Code Generation Based on AADL Model; 2018"), there haven't been significant contributions toward their classification or identification of problems that prevent them to be usable in practice. Requirements for the solution Functional requirements • RQ.1: An approach is able to use dataset for supervised learning and testing & Input • RQ.2: An approach is able to generate a training/testing dataset (For training an algorithm, we would need a specific formatted dataset or a structured dataset. The generated dataset will be the consolidated data collected by a design template (including fields like inputs, function requirements, source code, and output) for users to input, hence making our training data. The same logic can also be applied to gathering the testing data. However, there can be other means to gather datasets, for instance, crawling/scrapping/REST-API calls from open source projects or public/private repositories.) from taking user input via a design template. • The design template will include fields to structure the natural language input provided, such as function input (e.g. parameter integer A, parameter string B, etc.), function requirements (which can include function comments, description, use cases, etc.), function source code and function output (e.g. return type, etc.) • RQ.3: An approach can extract and generate a training/testing dataset from open source projects and public or private repositories (e.g. GitHub, BitBucket, etc.) • RQ.4: An approach is able to persist the consolidated training/testing dataset in a storage system (e.g. database, file storage, etc.) & Pre-Processing • RQ.5: An approach is able to pre-process the stored dataset for training and testing for the AI models (e.g. tokenization, word2vec, stop word removal, etc.) • RQ.6: An approach is able to transform the extracted repository data to a structured format, such that it aligns with the design template dataset • RQ.7: An approach is able to train on Python programming language source code & Training • RQ.8: An approach is able to train on the structured dataset that is prepared using the design template. • RQ.9: An
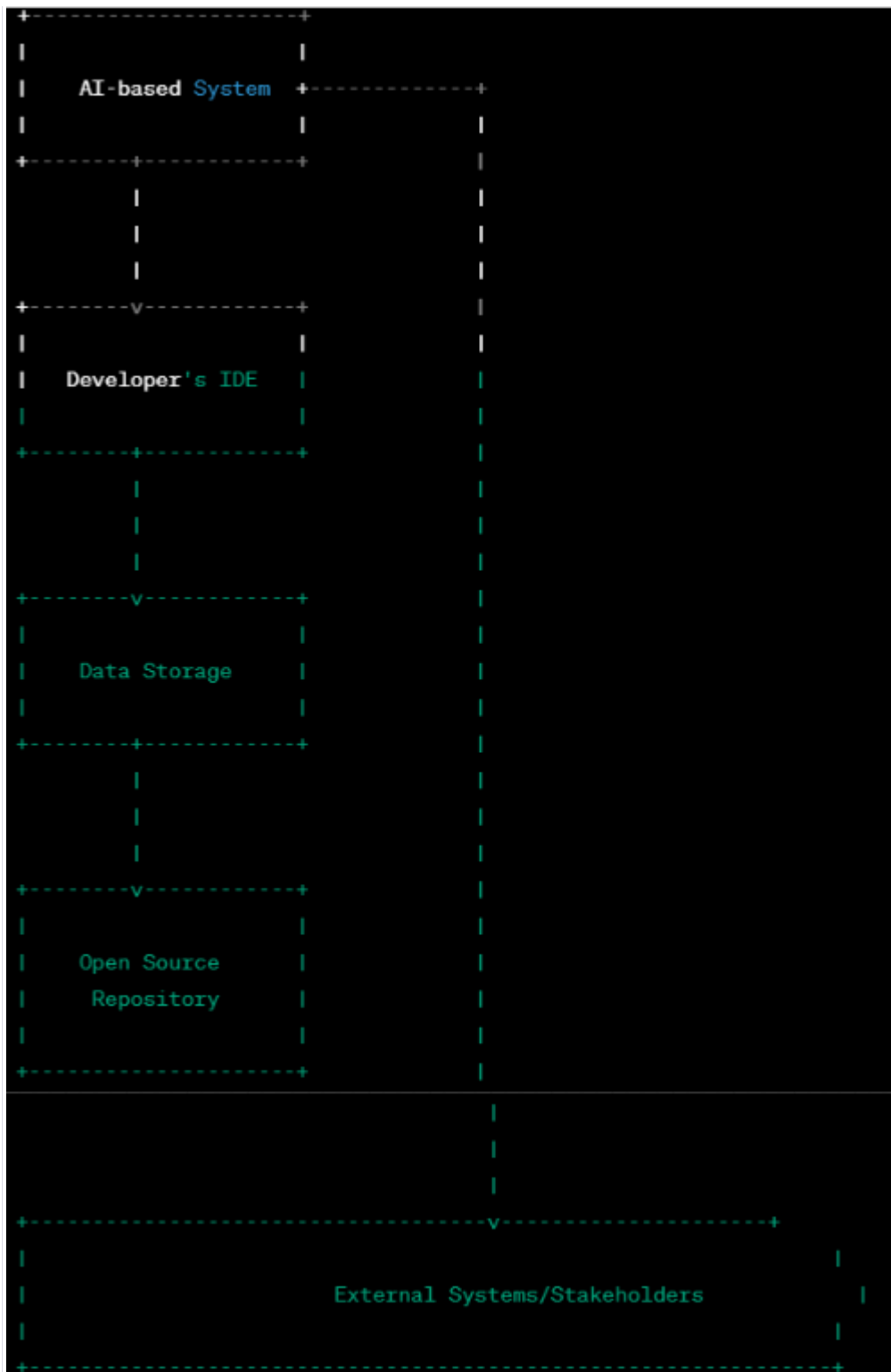
approach is able to train on deep neural networks by using existing models (e.g. GPT2, GPT3, etc.) • RQ.10: An approach is able to understand the requirements in the structured dataset to generate the source code • RQ.11: An approach is able to generate source code for Python programming language based on a set of input requirements (testing design template) & Output • RQ.12: An approach is able to generate multiple source code functions or snippets as a suggestion for the user to select • RQ.13: An approach is able to evaluate the generated source code accuracy based on syntax and logical correctness (by using the following metrics e.g. BLEU score, F-Score, WER, etc.) & Testing and Evaluation • RQ.14: The approach is able to integrate in an IDE (e.g. Eclipse, NetBeans, VS Code, etc.) by a plugin. This will allow the end user to take benefit of the deployed approach in a real-world use case when they are programming in on-going projects. & Deployment • RQ.15: The approach is able to save the trained model in a compressed format (e.g. ZIP, PKL, etc.) in a storage system for testing, evaluating and deploying Quality requirements • Q.1: - Environment: ♣ An approach generating a function source code based on a structured input template ♣ Number of syntax faults = E; Parameters = p - Stimulus: ♣ Re-training the model with hyperparameter tuning ♣ Parameters = q - Response: ♣ Better syntax accuracy ♣ Number of syntax faults = E/2 • Q.2: - Environment: ♣ An approach generating a function source code based on a structured input template ♣ Cyclomatic complexity = CE ♣ Parameters = p - Stimulus: ♣ Re-training the model with hyper-parameter tuning ♣ Parameters = q - Response: ♣ Better syntax accuracy ♣ Number of syntax faults = E/2 • Q.3: - Environment: ♣ An approach generating a function source code in an on-going project ♣ Lack of cohesion in methods (LCOM) = C - Stimulus: ♣ Existing source code in an on-going project which are logically equivalent to the generated source code ♣ Parameters = q ♣ Response Low Cohesion [6] LCOM = 0 • Q.4: - Environment: ♣ An approach generating a function source code in an on-going project ♣ Size of the training dataset and its trained model accuracy; Size = S; Accuracy = A% - Stimulus: ♣ Re-training the model with the increased data size ♣ Size = 2S - Response: ♣ Improved training model accuracy ♣ Accuracy = 2A%

The problem addressed in this text is the difficulty that developers face in coping with the increasing complexity of software and producing high-quality source code, as well as the challenges of code reuse and understanding software requirements. The proposed solution is to use AI models to understand requirements and generate full functions to assist developers. The functional requirements for the solution include the ability to use a structured dataset for supervised learning and testing, generate and extract training/testing datasets, pre-process the stored dataset, train on Python programming language source code, and generate multiple source code functions as suggestions for the user to select. The approach should also be able to integrate into an IDE, save the trained model in a compressed format, and evaluate the generated source code's accuracy. The quality requirements include better syntax accuracy, reduced cyclomatic complexity, low cohesion in methods, and improved training model accuracy.

2. Suggest architecture of a conceptual solution for this problem. We need a C4 system context diagram (https://c4model.com/#SystemContextDiagram). Please use text (arrows, dots, different shapes, and different colours) to visualise the solution. Consider the following: - Identify relevant external systems (if any) and add them to the diagram. - Identify relevant external stakeholders (if any) and add them to the diagram. - Identify high level abstract data types transferred between components and stakeholders. - Summarise data types transmitted between components. Add this information to the connections and interfaces between components.

Here's a suggested C4 system context diagram for the conceptual solution:

```
+--------------------+
|                    |
|   AI-based System  +----------------+
|                    |                |
+--------+-----------+                |
         |                            |
         |                            |
         |                            |
+--------v-----------+                |
|                    |                |
|   Developer's IDE  |                |
|                    |                |
+--------+-----------+                |
         |                            |
         |                            |
         |                            |
+--------v-----------+                |
|                    |                |
|   Data Storage     |                |
|                    |                |
+--------+-----------+                |
         |                            |
         |                            |
         |                            |
+--------v-----------+                |
|                    |                |
|   Open Source      |                |
|   Repository       |                |
|                    |                |
+--------------------+                |
                                      |
                                      |
                                      |
+-------------------------------------v--------------------+
|                                                          |
|            External Systems/Stakeholders                 |
|                                                          |
+----------------------------------------------------------+
```
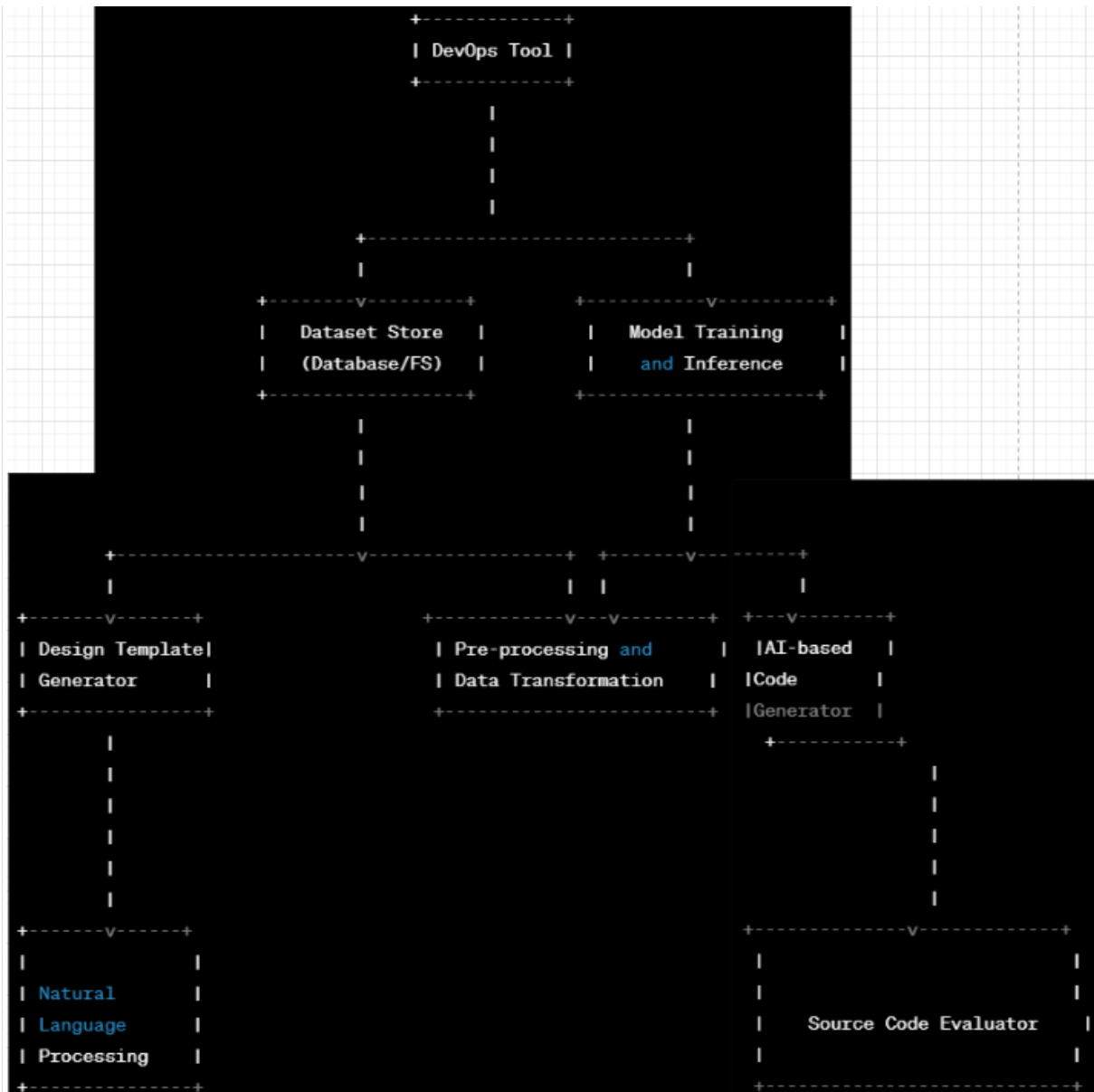
The AI-based system is the main component in the system context diagram. It interacts with the Developer's IDE, Data Storage, and Open Source Repository to perform its functions. The Developer's IDE is where developers write code and can receive suggestions from the AI-based system. The Data Storage component stores the structured datasets and trained models. The Open Source Repository component provides additional data for training and testing.

The main data types transmitted between components are structured datasets, trained models, and source code functions. The AI-based system can extract and generate structured datasets from user input or from the Open Source Repository. It then pre-processes the stored dataset for training and testing by performing operations such as tokenization, word2vec, and stop word removal. The trained model is then used to generate source code functions based on the input requirements. The generated source code functions can be evaluated for accuracy based on syntax and logical correctness.

The external systems and stakeholders are not specified in the text, so they are represented by a single box in the system context diagram. They may include stakeholders such as end-users who benefit from the software developed with the assistance of the AI-based system or external systems that interact with the AI-based system, such as a Continuous Integration/Continuous Deployment (CI/CD) pipeline.

3. Suggest architectural components that we need to solve this problem. We need a C4 container solution diagram (https://c4model.com/#ContainerDiagram). Please use text (arrows, dots, different shapes, and different colours) to visualise the solution. Consider the following: - Define how the main stakeholder will use the solution. Add it to the diagram.. - Which quality properties are important? Suggest adequate architectural solutions to address them (architectural style, design patterns, design tactics, components and their relationships). Structure the components according to it. - Define data types transferred between components. - Summarise data types transmitted between components. Add this information to the connections and interfaces between components. - Define interfaces between components. Add information about required and provided interfaces to the components, along with the data types. - Document which anti-patterns we need to avoid when implementing this solution.

Based on the problem description, here is a possible C4 container diagram:

```
                          +-------------+
                          | DevOps Tool |
                          +-------------+
                                 |
                                 |
                                 |
                                 |
                    +----------------------------+
                    |                          |
            +--------v--------+         +----------v----------+
            |  Dataset Store  |         |   Model Training    |
            |  (Database/FS)  |         |    and Inference    |
            +-----------------+         +---------------------+
                    |                          |
                    |                          |
                    |                          |
                    |                          |
        +---------------------v---------------+ +-------v----------+
        |                          |   |        |
+--------v-------+     +-------------v---v--------+ +---v--------+
| Design Template|     | Pre-processing and      |  |AI-based   |
| Generator      |     | Data Transformation     |  |Code       |
+----------------+     +-------------------------+  |Generator  |
        |                                            +-----------+
        |                                                  |
        |                                                  |
        |                                                  |
        |                                                  |
        |                                                  |
+-------v------+                           +-------------v------------+
|            |                             |                        |
| Natural    |                             |                        |
| Language   |                             |   Source Code Evaluator |
| Processing |                             |                        |
+------------+                             +------------------------+
```

- The main stakeholders will use the solution by providing input to the Design Template Generator, which will generate a structured dataset that can be used for training an AI-based Code Generator. The trained model can then be used to generate code snippets or full functions based on the input requirements provided by the user. The Source Code Evaluator will check the correctness of the generated code.
- The important quality properties are accuracy, reliability, and scalability. To address these, we can use a microservices architecture, containerization, and automated testing. We can also use the Façade pattern to hide the complexity of the system and present a simplified interface to the end-users.
- The Design Template Generator generates a structured dataset that includes the following data types: function input, function requirements, function source code, and function output. The pre-processing and data transformation component will apply tokenization, word2vec, stop word removal, etc. to prepare the dataset for training. The AI-based Code Generator will generate source code based on the structured dataset and input requirements. The Source Code Evaluator will evaluate the generated code based on syntax and logical correctness.

- The data types transmitted between components include structured datasets, training/testing data, source code, and evaluation metrics. The Design Template Generator provides the structured dataset to the pre-processing and data transformation component, which generates training/testing data for the AI-based Code Generator. The AI-based Code Generator generates source code based on the training data and input requirements, which is then evaluated by the Source Code Evaluator. Evaluation metrics are provided back to the AI-based Code Generator for further training.
- The interfaces between components include the following:
  - Design Template Generator provides a REST API for the end-users to provide input requirements.
  - Pre-processing and data transformation component provides a REST API for the AI-based Code Generator to access the prepared training/testing data.
  - AI-based Code Generator provides a REST API for the end-users to receive generated code and for the Source Code Evaluator to evaluate the generated code.
  - Source Code Evaluator provides a REST API for the AI-based Code Generator to receive evaluation metrics.
- Anti-patterns to avoid include: spaghetti code, tight coupling, and monolithic architecture. These can be addressed by using a modular, microservices-based architecture with clear separation of concerns and loose coupling between components.

4. Suggest architectural components that we need to solve this problem. We need a C4 component solution diagram (https://c4model.com/#ComponentDiagram). Please use text (arrows, dots, different shapes, and different colours) to visualise the solution. Consider following: - Define how the main stakeholder will use the solution. Add it to the diagram. - Which quality properties are important? Suggest adequate architectural solutions to address them (architectural style, design patterns, design tactics, components and their relationships). Structure the components according to it. - Define data types transferred between components. - Summarise data types transmitted between components. Add this information to the connections and interfaces between components. - Define interfaces between components. Add information about required and provided interfaces to the components, along with the data types. - Document which anti-patterns we need to avoid when implementing this solution.

As a language model, I cannot access visualisation tools to create a diagram. However, I can provide you with the information that can be included in a C4 component solution diagram based on the previous question. Here are the answers:

- Define how the main stakeholder will use the solution: The main stakeholder, i.e., developers, will use the solution by inputting the structured data through a web-based user interface, which includes function input, function requirements, function source code, and function output. The system will then generate the suggested source code functions and display them on the same interface. The developer can then select the most appropriate function for their use case.
- Which quality properties are important? Suggest adequate architectural solutions to address them (architectural style, design patterns, design tactics, components and their relationships). Structure the components according to it:
  - Usability: The system should be user-friendly, and the web-based user interface should be intuitive to use. To achieve this, the design should focus on simplicity

and consistency. Also, the user interface should be responsive and work well on different devices.

- Performance: The system should generate source code functions quickly. To achieve this, the system should use efficient algorithms and caching to reduce the time required to generate source code functions.
- Accuracy: The generated source code functions should be syntactically and logically correct. To achieve this, the system should include a testing component to evaluate the generated source code.
- Maintainability: The system should be easy to maintain and extend. To achieve this, the system should follow best practices for code organization and design patterns. The system should also have clear documentation and include tools for debugging and testing.

- Define data types transferred between components: The data types transferred between components include structured data for input (function input, function requirements, function source code, and function output), the training dataset, and the generated source code functions.
- Summarize data types transmitted between components: The structured data for input will be transferred from the web-based user interface to the machine learning model for generating source code functions. The training dataset will be used to train the machine learning model. The generated source code functions will be displayed on the web-based user interface for the developer to select.
- Define interfaces between components: The web-based user interface will require a REST API to communicate with the machine learning model. The machine learning model will require a database to store the training dataset. The machine learning model will also require a testing component to evaluate the generated source code functions.
- Document which anti-patterns we need to avoid when implementing this solution: To avoid anti-patterns, the system should follow best practices for code organization and design patterns. Also, the system should use efficient algorithms and caching to reduce the time required to generate source code functions. The system should also have clear documentation and include tools for debugging and testing.

5. Please suggest a toolchain to implement the solution. Consider which tools we can reuse. For the components we need to develop, please suggest a programming language. Please use text (arrows and dots) to visualise the solution. Define data types transferred between components. Summarise data types transmitted between components. Add this information to the connections and interfaces between components.

To implement the proposed solution, we suggest the following toolchain:

- For data extraction from open-source repositories, we can use Git or similar tools.
- For pre-processing and data cleaning, we can use Python-based libraries such as Pandas, Numpy, and Scikit-Learn.
- For building deep learning models, we can use Python-based libraries such as TensorFlow or PyTorch.
- For building the web application, we can use Django or Flask, which are Python-based web frameworks.
- For integrating the solution into an IDE, we can use the corresponding plugin development kit provided by the IDE vendor.

The programming language of choice for this solution is Python, as it has strong support for data manipulation, machine learning, and web development.

The data types transferred between components include natural language text, tokenized text, source code snippets, and trained machine learning models.

The data flows are as follows:

- The data extraction component retrieves data from open-source repositories and stores it in a structured format.
- The pre-processing component cleans and transforms the data into a suitable format for machine learning.
- The machine learning component trains a deep learning model on the pre-processed data.
- The web application component provides a user interface for developers to input natural language requirements and receive generated source code snippets.
- The plugin component integrates the solution into an IDE for seamless usage.

The interfaces between components include HTTP APIs for the web application and plugin components to communicate with the machine learning and data components. The data types transmitted between components are primarily JSON-encoded data, such as input/output data for the machine learning component and generated source code snippets for the web application and plugin components.