

University of Bremen

TODO: Make sure
this is centered!

figures/unibremen.pdf

Faculty 3 — Mathematics & Computer Science

Master's Thesis

Building Code Cities using the Language Server Protocol

— *DRAFT* —

Falko Galperin

- | | |
|--------------------|---|
| 1. <i>Reviewer</i> | Prof. Dr. Rainer Koschke
<i>Working group Software Engineering</i> |
| 2. <i>Reviewer</i> | Prof. Dr. Ute Bormann
<i>Working group Computer Networks</i> |

November 17, 2024

Abstract

TODO!

Declaration

I hereby declare that I have completed this master's thesis independently and without any external assistance, unless explicitly stated otherwise. I have not used any sources or aids other than those specified. All passages that have been quoted verbatim or paraphrased from published sources are clearly identified as such.

Bremen, November 17, 2024

Falko Galperin

Acknowledgements

TODO!

Contents

1	Introduction	1
1.1	Format	1
1.2	Motivation	2
1.2.1	Code Cities & SEE	2
1.2.2	Language Server Protocol	3
1.2.3	Integration	6
1.3	Goals & Research Questions	8
1.4	Thesis Structure	9
2	Concepts	11
2.1	SEE	11
2.1.1	Basics	12
2.1.2	Project Graph	13
2.1.3	Other Relevant Features	14
2.2	Language Server Protocol	15
2.2.1	Basics	16
2.2.2	Planned Capabilities	18
2.2.3	Unplanned Capabilities	21
2.3	Interim Conclusion	23
3	Implementation	25
3.1	Preliminary Changes	25
3.1.1	Specification Cleanup	25
3.1.2	Preparing SEE	26
3.2	Generating Code Cities using LSP	27
3.2.1	Algorithm	27
3.2.2	Augmented Interval Trees	35
3.2.3	Usage in Practice	38
3.3	Integrating LSP Functionality into Code Cities	41
3.4	Integrating LSP Functionality into Code Windows	41
3.5	Technical Evaluation	41
3.6	Interim Conclusion	41

4	User Study	43
4.1	Plan	43
4.1.1	Existing Research	44
4.1.2	VSCode	44
4.1.3	Hypotheses	44
4.2	Structure	45
4.2.1	Questionnaire	45
4.2.2	Tasks	45
4.3	Results	45
4.3.1	Demographics	45
4.3.2	Correctness	48
4.3.3	Time	49
4.3.4	Usability	49
4.3.5	Comments	49
4.4	Threats to Validity	49
4.5	Interim Conclusion	49
5	Conclusion	51
5.1	Limitations	51
5.2	Future Work	51
5.3	The End	51
A	List of TODOs	53
B	Glossary	57
C	Acronyms	59
D	Attached Files	61
E	List of Figures	63
F	List of Tables	65
	List of Listings	67
G	List of Listings	67
H	Bibliography	69

TODO: Listings
appear twice in
TOC

Introduction



LINKING the Language Server Protocol with code cities, specifically SEE, is the core focus of this master's thesis. The main novel contribution will be a way of creating code cities using only information provided by the Language Server Protocol, while additional contributions consist of the integration of Language Server Protocol-based functionality into code cities as well as their integration into SEE's code windows. Finally, the penultimate chapter 4 describes a controlled experiment (with $n = 20$ participants) in which code cities are compared to traditional IDEs via a user study—this serves as an evaluation for this thesis and for code cities as IDE replacements in general.

TODO: Maybe
rewrite. Could
start by going over
languages in
general, giving
spell check as
example.

In this chapter, apart from explaining some formatting semantics, we will examine the motivation and basics behind each of the central concepts (i. e., code cities, the Language Server Protocol, and the integration of the two), name the goals and research question of the thesis, and finally describe the structure of upcoming chapters.

1.1 Format

This document uses many technical terms that not every reader may know. To remedy this, starting from the next section, whenever a technical term or acronym appears for the first time, it will be printed in *this color*, and an explanation of that term will appear in a box of the same color nearby. Terms that are already explained in the text itself will not receive such a box. The terms and acronyms are collected within the glossaries in appendices B and C—all mentions of such terms also link (in the digital version of this thesis, at least) to the corresponding part of the glossary. In total, the following colors are used to convey specific meanings:

- **Maroon** for the introduction of a glossary term or acronym,
- **Fuchsia** for internal links (e. g., to other sections),

New Term: *IDE* (Integrated Development Environment)

Editor for source code with features that are useful for development (e. g., highlighting errors). Examples are *Eclipse* or *JETBRAINS IntelliJ*.

- **Blue** for external links (e. g., to web pages),
- **Green** for cited literature, and
- **Cyan** for references to attached files (see appendix D).

1.2 Motivation

As mentioned at the beginning of this section, this master’s thesis is about *integrating* the *Language Server Protocol* into *Code Cities*. I will motivate each of these italicized central points individually in the following sections. Note that these will get more thorough explanations in chapter 2.

1.2.1 Code Cities & SEE

Visualization in general often helps facilitate the understanding of complex systems by representing them with a simplified visual model. This can be especially useful in the area of software engineering, where it is often hard to get an intuitive overview of large software systems when only equipped with standard tools, like *Integrated Development Environments* (IDE). One such software visualization—called *Software Engineering Experience* (SEE)—is being developed at the University of Bremen and will be introduced in the next section.

SEE is an interactive software visualization tool using the *code city* metaphor in 3D, developed in the Unity game engine. It features collaborative “multiplayer” functionality across multiple platforms¹, allowing multiple participants to view and interact with the same code city together.

In the code city metaphor, software components are visualized as buildings within a city. Various metrics from the original software can then be represented by different visual properties of each building—for example, the *Lines of Code* (LOC) within a file might correlate to the height of the corresponding building. Relationships between software components, such as where components are referenced, are instead represented by edges

¹Notably, besides usual desktop and touchscreen-based environments, virtual reality (e. g., via the *Valve Index*) is supported as well.

figures/SEE_screenshot.png

Figure 1.1: A code city visualized in SEE.

New Term: LOC (Lines of Code)

The number of lines in a source code file.

drawn between the respective buildings. The exception to this are part-of relations, that is, relations that describe which component belongs to which other component. These are instead visualized in SEE by buildings being nested within their corresponding “parent” building. In this way, the data model of SEE can be represented as a graph in which the software components are the nodes and the relationships are the edges.

For example, in fig. 1.1, we can see the source code of the SPOTBUGS project² rendered as a code city. A few very tall buildings—indicating that the respective component is very big and that a refactoring into smaller pieces may be in order—immediately jump out. Additionally, this visualization also makes the number of methods readily apparent: the redder a node, the higher its method count. fig. 1.2 instead visualizes the modeled architecture of a very small system, as compared to a city “empirically” generated by an implementation like in the previous example. Here, we can also see yellow edges between the components, in this case representing desired references that should be present between components.

1.2.2 Language Server Protocol

As stated on its website:

Adding features like auto complete, go to definition, or documentation on hover for a programming language takes significant effort. Traditionally[,] this work

²<https://github.com/spotbugs/spotbugs> (last access: 2024-09-11)

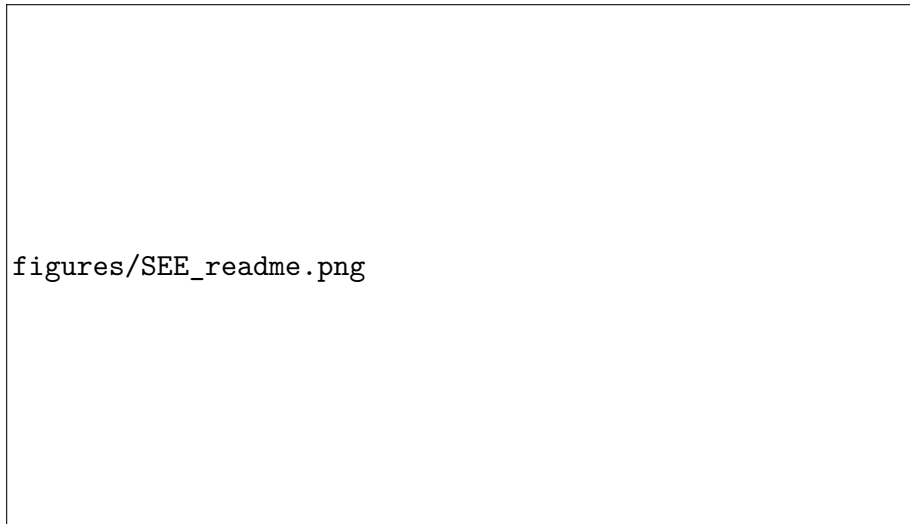


Figure 1.2: An example of what edges can look like in SEE. **TODO:** Use higher quality image.

had to be repeated for each development tool, as each tool provides different APIs for implementing the same feature.

A *Language Server* is meant to provide the language-specific smarts and communicate with development tools over a protocol that enables inter-process communication.

The idea behind the *Language Server Protocol* (LSP) is to standardize the protocol for how such servers and development tools communicate. This way, a single Language Server can be re-used in multiple development tools, which in turn can support multiple languages with minimal effort. ([LspSpec])

Since LSP³ is a central component of my master’s thesis, I have created a diagram in fig. 1.3 in the hope to strengthen intuitions around the motivation and use of the protocol. While the LSP specification has originally been created by Microsoft, it is by now an open-source project⁴, where changes can be actively proposed using issues or pull requests. Apart from the specification itself, a great number of open-source implementations of Language Servers for all kinds of programming languages from Ada to Zig exist. A partial overview of available implementations is listed at <https://microsoft.github.io/language-server-protocol/implementors/servers/> (last access: 2024-09-11).

The protocol introduces the concept of so-called *capabilities*, which define a specific set of features a given Language Server (and *Language Client*) support. These include

³Note that I will often refer to the Language Server Protocol as just “LSP” instead of “the LSP” (e. g., “IDEs use LSP”) from now on, as this is how the specification [LspSpec] does it as well.

⁴Available at <https://github.com/Microsoft/language-server-protocol> (last access: 2024-09-11).

content/../../figures/overview_lsp_without_lsp.pdf

(a) IDE development without LSP.

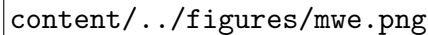
content/../../figures/overview_lsp_with_lsp.pdf

(b) IDE development with LSP.

Figure 1.3: An illustration of how LSP can help simplify IDE development.

New Term: *Language Client*

A development tool, such as an IDE, that supports LSP and can hence integrate language-specific features into itself using compatible Language Servers.



```
content/../../figures/mwe.png
```

Figure 1.4: An example of the `texlab` Language Server running in NEOVIM.

New Term: GXL (Graph eXchange Language)

A file format for graphs, used in SEE for representing dependency and hierarchy graphs of software projects.


navigational features, like the ability to jump to a variable’s declaration, and editing-related features, such as autocomplete. To give a specific example of what an LSP capability might look like in practice, the `texlab`⁵ Language Server for \LaTeX —which I am using while writing this document—provides a list of available packages when one starts typing text after “`\usepackage{`”. Additionally, for the currently hovered package, a short description of it is displayed. A screenshot of this behavior within the NEOVIM editor is provided in fig. 1.4.

The counterparts to Language Servers are the Language Clients: These are the IDEs and editors that incorporate the Language Server into themselves. Examples for IDEs that support acting as a Language Client in the LSP context include *Eclipse*, *Emacs*, JETBRAINS IntelliJ (NEO)VIM, and *Visual Studio Code*.

1.2.3 Integration

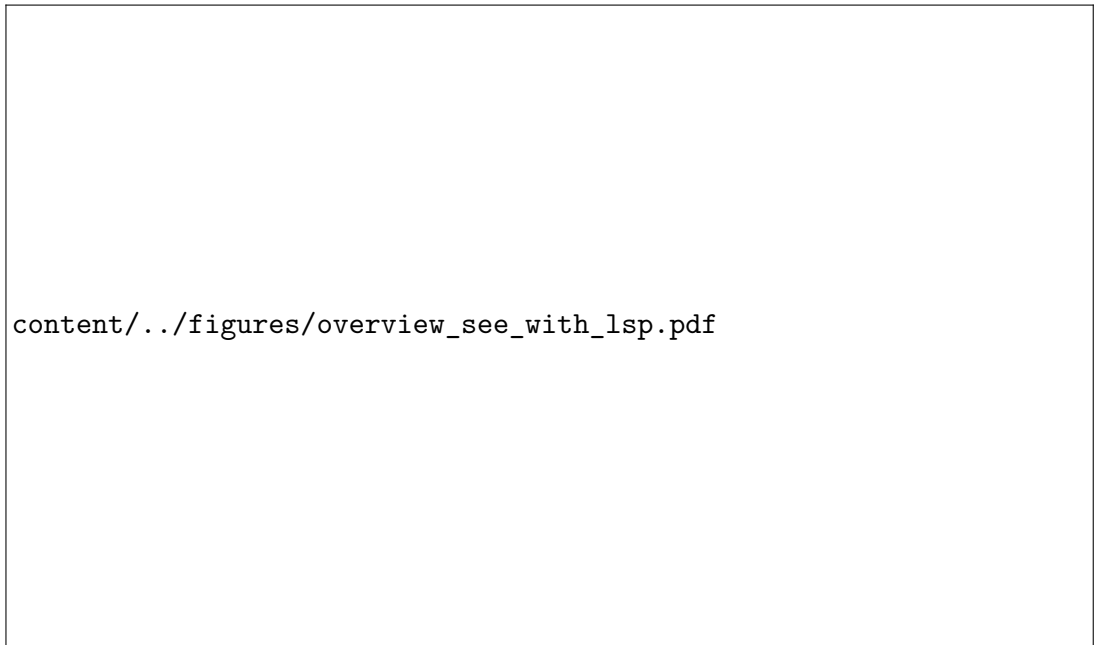
Currently, code cities in SEE are rendered by reading in pre-made *Graph eXchange Language (GXL)* files, which can be created by the proprietary Axivion Suite. This approach has the disadvantage of only supporting languages supported by the Axivion Suite, as well as making regenerating cities (e. g., if the source code changed) fairly cumbersome. Another current shortcoming of SEE is that information about the source code available to the user is limited when compared to an IDE—for example, quickly displaying documentation for a given component by hovering over it is not supported. This is where the Language Server Protocol can help.

⁵<https://github.com/latex-lsp/texlab> (last access: 2024-09-12)



content/../../figures/overview_see_without_lsp.pdf

(a) SEE as it exists right now. The analyzers on the right are developed by Axivion.



content/../../figures/overview_see_with_lsp.pdf

(b) SEE with implemented LSP integration. The Language Servers on the right are mostly developed by the open-source community.

Figure 1.5: A simplified illustration of how SEE currently gains and uses information about software projects, and how the integration of LSP would change that.

1.3 Goals & Research Questions

The goal of this master’s thesis—as outlined in section 1.2.3—is to integrate the Language Server Protocol into SEE by making it a Language Client, then evaluate this implementation by comparing it with traditional IDEs in a user study. To this end, the main contribution is a way of generating code cities using the Language Server, where all the information obtainable by relevant LSP capabilities should be manifested⁶ in the city in a suitable way. This is an unintended (or at the very least, unusual) use of LSP and may require some experimentation.

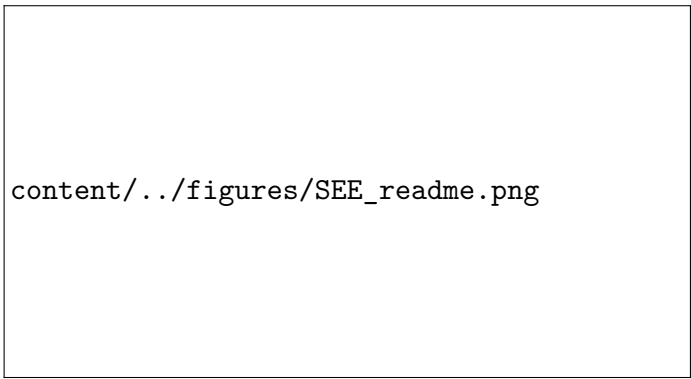


Figure 1.6: A code window.

Apart from the code cities, SEE also provides the so-called *code windows*, in which the source code of a specific component can be viewed in a similar way as in an IDE. This can be seen in fig. 1.6. An additional goal of this master’s thesis is to enhance the functionality of the code windows by implementing more IDE-like behavior into it (e. g., allowing users to go to

a variable’s declaration, or displaying diagnostics inline) using the Language Server.

TODO: Maybe remove the next part if already mentioned before.
TODO: In general, reword some stuff from “planning to do” ⇒ “done”

It should be noted that any LSP capabilities which involve modifying the underlying software project is out of scope for this master’s thesis. Rewriting the code windows to be editable (in a way that distributes the edits over the network) is complex enough to warrant its own thesis [see also Ble22], not even taking into account that this would also more than double the number of capabilities that would then become useful to implement. As such, only capabilities that are “read-only” (i. e., that do not modify the source code or project structure in any way) will be considered when planning which features to implement as part of the thesis. Consult section 2.2.2 for a full list of LSP capabilities which I plan to implement.

Additionally, I will not implement the C# interface to the LSP (i. e., translating between C# method calls and *JavaScript Object Notation—Remote Procedure Calls* (JSON-RPC)). There already exist well-made interfaces for this purpose⁷, and the focus of the thesis will

⁶Since there is a lot of diverse data available via LSP, it makes sense to only immediately display the most pertinent information and make the rest of it available upon request within SEE’s user interface.
⁷Such as OmniSharp’s implementation of LSP in C#, which I plan to use: <https://github.com/OmniSharp/csharp-language-server-protocol> (last access: 2024-28-01)

New Term: *JSON-RPC* (JavaScript Object Notation—Remote Procedure Call)

A remote procedure call protocol that uses JSON as its encoding, supporting (among other features) asynchronous calls and notifications. It is used as the base for LSP (even though LSP is technically not a remote protocol).

be on the integration of the protocol's *data* into SEE's code cities and code windows, not the integration of the protocol itself.

Hence, the goals for this thesis can be summarized as follows:

- Integrating an LSP framework into SEE and allowing users to manage Language Servers from within SEE
- Making SEE a Language Client, such that:
 1. Code cities can be generated directly from source code directories, using Language Servers that SEE interfaces with,
 2. code windows gain “read-only” IDE-like functionality, covering behavior of capabilities listed in ??, and
 3. Code cities gain similar functionality (where applicable), such as displaying relevant documentation when hovering over a node.
- Evaluating the above empirically in a controlled experiment via a user study.

The main research questions that I want to answer in this thesis are as follows:

RQ1 Is it feasible (i. e., realistically doable with usable results) to generate code cities using the Language Server Protocol?

RQ2 Are code cities a suitable means to present LSP information to developers as compared to IDEs + tables (on the dimensions of speed, accuracy, and usability)?

TODO: This is rather vague. Is that alright or do I need to operationalize here?

1.4 Thesis Structure

We begin by examining the Language Server Protocol and the concept of code cities along with SEE more closely in chapter 2. Next, in chapter 3 we take a look at my implementation of the LSP-based code city generation algorithm, alongside additional contributions to visualize this information in the cities and code windows of SEE. In the course of this, we will answer **RQ1**. To answer **RQ2**, I will carry out a user study comparing an LSP-enabled IDE (specifically, VSCode) with an LSP-enabled code city visualization (specifically, SEE) and report on its results in chapter 4. Finally, we wrap up in chapter 5 by summarizing the

thesis's results and providing an outlook on additional ideas for the implementation, while also listing some possible avenues of further research.

Concepts



OUTLINING the central concepts behind this thesis is a crucial first step before tackling the implementation, so that we can form a concrete idea of which parts of LSP are well-suited to being integrated into code cities. Thus, we will examine the concept of code cities, where we will use SEE as a concrete implementation (which we do in section 2.1), as well as the Language Server Protocol (which we do in section 2.2). For the former, we will go over some of the existing literature regarding code cities (as this is more of an academic topic than LSP) and describe the essentials that we need to work with in the implementation, such as the project graph. As for the latter, we will go over the available capabilities and take a look at both existing Language Servers and Language Clients to get an idea of what the Language Server Protocol offers and how it is most commonly used. For both topics, we will focus on the parts relevant to the implementation and evaluation—for example, we will only explore those capabilities in detail that actually end up being used in this thesis.

2.1 SEE

As explained in section 1.2.1, SEE is an interactive software visualization tool using the code city metaphor in 3D, with the aim to make it easier to work with large software projects on an architectural level. To name a few example scenarios in which using code cities could be useful, **compared to** using traditional IDEs:

- Senior developers may use it in its “multiplayer” mode to explain the structure of their software to newcomers
- Project planners could visualize *code smells* to find candidates for refactoring [Gal21]
- Software architects can find violations of the planned architecture using SEE’s *reflexion analysis*

TODO: Back this up with sources

TODO: Refer back here later, noting how LSP can generate code cities for each of these use cases

New Term: *Code Smell*

Certain structures in source code that suggest that a refactoring is in order, such as duplicated code, or a very long method [FB19, pp. 85–87].

New Term: *Reflexion Analysis*

The process of comparing the architecture and implementation of a software project and finding incongruencies between the two.

SEE is an open-source project, currently hosted at GitHub¹. It is a research project at the University of Bremen, where it has been in development **since 2019**, and where it is frequently offered as a bachelor or master project for the students. While it was at first a project for the *Unreal Engine*, the game engine has been switched to *Unity* relatively early, mainly because its editor eased development due to the ability to reload the *User Interface* (UI) without having to restart the whole engine.

TODO: is this correct? And put source here

After going over the basics of how SEE works, I will give an explanation and formalization of the project graph—as this is the central component in the LSP-based code city-generation algorithm—followed by a brief overview over other features that are relevant for this thesis.

2.1.1 Basics

As mentioned before, when analyzing a software project in SEE, its source code structure is reflected in the rendered buildings, where connections between the components, such as references, are displayed using edges rendered as splines (this will be elaborated in section 2.1.2). Metrics of the project, such as the McCabe complexity [McC76], can then be mapped onto visual properties of this rendering, such as its color, width, or height.

TODO: Needs a 3D SEE screenshot, esp.

So far, the kind of visualization described here would still be possible **in two dimensions**, but adding a third dimension gives us the benefit of having another axis to map metrics onto, as well as making it a more natural environment for humans to interact with. By giving users the ability to walk, move the camera around, and so on, as is possible in many video games, the code city can be navigated in a richer and more intuitive manner than if it were just a simple 2D image. This also makes it possible to offer SEE as a virtual reality application. Virtual reality gives users an even more immersive and natural way of interacting with the environment, leading to improved spatial memory of the city, as some recent studies have shown. It also makes it a more fitting environment for multiple people. Each player

for edges

TODO: add 2D screenshot

TODO: Source for 3D/game/VR being easier to navigate/understand

¹<https://github.com/uni-bremen-agst/SEE> (last access: 2024-09-18) (but note that, to clone this repository, you additionally need access to the Git LFS counterpart hosted at the University of Bremen's GitLab, as this is where paid plugins for SEE are hosted.)

New Term: *Polytree*

A directed acyclic graph which also has no undirected cycles.

can have their own avatar placed in the scene, interacting with one another and the city, communicating via voice chat.

I heavily recommend readers take a look at [this introductory video](#), as an interactive game-like project like this comes across much more clearly in a video than in mere texts and disconnected snapshots.

TODO: link the video

2.1.2 Project Graph

We will now take a look at the graph representing the source structure of the analyzed software projects in SEE. One of the main goals of this thesis is generating such a graph using LSP, as opposed to the currently available methods of generating this graph, which require access to the proprietary Axivion Suite to create the GXL files used in SEE.

These graphs can be formalized as $G = (V, E, a, s, t, \ell)$, with:

- V being a set of nodes and E being a set of edges.
- $a : (V \times \mathcal{A}_K) \rightarrow \mathcal{A}_V$ associating nodes and an attribute name ($\in \mathcal{A}_K$) with a value ($\in \mathcal{A}_V$). Note that this is a partial function.
- $s : E \rightarrow V$ denoting the source node of an edge.
- $t : E \rightarrow V$ denoting the target node of an edge.
- $\ell : E \rightarrow \Sigma$ providing a label for an edge over some alphabet Σ . Also, $\text{partOf} \in \Sigma$ such that the subgraph $(V, \{e \in E \mid \ell(e) = \text{partOf}\}, a, s, t, \ell)$ is a *polytree*.

Explained in natural language: The graph's nodes (V) can be connected by directed edges (E)—these are not necessarily unique for each possible tuple of nodes, meaning that there can be more than one edge between the same two nodes. A node represents a component in the source code (e. g., a class), while edges represent relationships between them (e. g., inheritance). This is an attributed graph. Specifically, each node can have multiple attributes (distinguished by attribute keys), while each edge always has one label that denotes the type of relationship it represents. The label `partOf` additionally has a special meaning: edges with this type induce the hierarchy of the source code (e. g., classes contained in files). Hence, if we look at the graph that removes all edges except for those with label `partOf`, and we make these edges undirected, we get a tree.

TODO: A lot of “represents” here

To illustrate, a few examples of attributes are:

- `Source.Path`: the path to the file the element is contained in.
- `Source.Name`: the name of the element within the source code.
- `Type`: the type of the element (e. g., “method”).

Another attribute that is also relevant for LSP is the `Source.Range`, which is often used in the upcoming section 2.2 and chapter 3. It describes a contiguous portion of source code. We formally define the domain of ranges \mathcal{R} as the Cartesian square of positions \mathcal{P}^2 , whereas the domain of positions \mathcal{P} is defined as the Cartesian square of natural numbers (including zero), that is, $\mathcal{P} = \mathbb{N}_0^2$. Hence, as a whole, $\mathcal{R} = \mathcal{P}^2 = \mathbb{N}_0^4$. Semantically, a position $(l, c) \in \mathcal{P}$ describes a zero-indexed line l and a zero-indexed character offset c , relative to the beginning of the line. A range $(b, e) \in \mathcal{R}$ can be understood as a beginning position b (inclusive) and an ending position e (exclusive). We will also occasionally “decompose” those positions and refer to a range r in decomposed form as $r = (b_l^r, b_c^r, e_l^r, e_c^r)$. For example, the interval of lines that the range r (partially or completely) covers is then given by $[b_l^r, e_l^r)$.

I should note that the model presented here is a simplification of SEE’s actual data model for source code graphs. For example, in reality, edges can have multiple attributes, there can be multiple edge types other than `partOf` that have the `polytree` property, and so on. This is just what is needed to understand the LSP-based city generation algorithm presented in section 3.2.

TODO: Give examples of software projects here to illustrate graph usage?

2.1.3 Other Relevant Features

Apart from the project graph, there are a few other features of SEE we need to go over, as they become relevant when integrating LSP into code cities (in section 3.3) and code windows (in section 3.4).

Context Menu When right-clicking any node or edge, a menu opens with several context-dependent options. These include the option to delete the element, to highlight it within the code city, to open its corresponding code window, and others. We will expand this menu with LSP-specific actions in the course of section 3.3.

TODO: Show screenshot of context menu (pre-LSP)

City Editor To actually generate a code city—assuming one has a GXL file for this purpose—a customized UI component within the *Unity Editor* exists in SEE. Here, a variety of options can be configured, such as the layout of the city, the mapping of metrics to visual attributes, or the *graph providers*, which create a project graph based on optional input parameters

TODO: Provide screenshot of editor

New Term: *Unity Editor*

The main UI of the Unity game engine, in which scenes can be set up, components can be configured, the game itself can be run, etc. Note that it is only used for development purposes, and hence not included within generated builds of a game.

New Term: *VSCode (Visual Studio Code)*

A proprietary, but free IDE developed by Microsoft with a plugin system from which LSP originated. See <https://code.visualstudio.com> (*last access: 2024-10-05*)

(such as the aforementioned GXL file). After implementing the city-generation algorithm for LSP, a new graph provider with its own Unity Editor-UI shall be implemented as well.

Code Windows As explained in section 1.3, there is also the option of opening code windows to view the source code of a component, an example of which is shown in fig. 1.6. Currently, these windows do little more than lexer-based syntax highlighting, so a goal is to include more IDE-like behavior by using functionality offered by LSP.

Erosion icons Following my bachelor's thesis [Gal21], SEE offers the possibility to indicate the number of code smells per component using the so-called *erosion icons*. These are essentially small icons that can be put above each node. The size and color of these erosion icons can then indicate the quantity of code smells for that given node. A controlled experiment has suggested that this gives developers a quicker, more intuitive overview over the distribution of code smells within a project, compared to traditional (i. e., tabular) ways of displaying this information [GKS22]. While LSP does not offer a standardized way of offering code smell information, we can use its diagnostics capability (see section 2.2.2) for the same purpose.

TODO: Show screenshot of erosion icons

2.2 Language Server Protocol

The very basic concepts behind LSP have already been explained in section 1.2.2. The protocol aims to make it easier for IDEs to support more programming languages—specifically, to support language-specific capabilities, which we will go over in section 2.2.2. It was originally developed for Microsoft's editor *Visual Studio Code* (VSCode) and was later converted into an open-source specification² (though there are still VSCode-specific extensions to LSP that are not in the official specification today). LSP has found widespread use: An overview

²Available at <https://github.com/microsoft/language-server-protocol> (*last access: 2024-10-05*).

page by Microsoft lists at least 269 Language Servers³ (i. e., servers offering support for some programming language) and 61 Language Clients⁴ (i. e., IDEs or development tools). The current (as of November 17, 2024) version of the protocol is 3.17, with version 3.18 being under active development.

2.2.1 Basics

Messages in the Language Server Protocol are built using JSON-RPC, which uses JSON to encode both requests (consisting of a method name and parameters) and responses (consisting of either a result object or an error object) for procedure calls. Requests include an ID that a response by the server can then reference to match it up with the request it is a reply to. There also so-called *notifications*, which are in essence requests without an ID, intended to send a message server that does not warrant a response [ISO13; Cro06]. In LSP, any message sent between the Language Client and Language Server consist of a header—describing the length and type of the content—followed by the content, which is always a JSON-RPC payload. While the specification does not mandate it, it lists some recommended communication channels on which the protocol messages can be sent, those being `stdio` (using standard input/output), `pipe` (using Windows's pipes), `socket` (using a socket), and `node-ipc` (IPC communication over Node.js⁵).

The specification lists all available types for requests and responses as TypeScript interfaces. A sample type definition for the hovering capability, taken from the documentation, can be seen in listing 2.1. From this example, we can also see that locations within documents⁶ are represented as a *Uniform Resource Identifier* (URI) representing the document along with a range (which we have formalized in section 2.1.2). The corresponding method name for this example is `textDocument/hover`.

³<https://microsoft.github.io/language-server-protocol/implementors/servers/> (last access: 2024-10-05)

⁴<https://microsoft.github.io/language-server-protocol/implementors/tools/> (last access: 2024-10-05)

⁵<https://nodejs.org> (last access: 2024-10-05)

⁶These documents must always be textual—there is no support for binary files.

Listing 2.1: Example specification of request and response objects for the Hover capability.

```
interface HoverParams {
    textDocument: string; /** The text document's URI in string form */
    position: { line: uinteger; character: uinteger; };
}

interface HoverResult {
    value: string;
}
```

The first request from the Language Client to the Language Server always has to be the `initialize` request, including the so-called client capabilities. These specify the capabilities that the Language Client supports. The response from the server will be a response object that includes the analogous server capabilities. The capability information being sent during this initial handshake is not just a pure list of the names of the corresponding procedure names, but also includes additional details on exactly which parts are supported (or, e. g., which encodings are used), specific to each capability. Afterwards, both parties will then restrict their usage of LSP to the subset of capabilities that both the client and server support.

The end of an LSP session is marked by the Language Client sending the Language Server a shutdown request. After the server confirms the success of the shutdown with a corresponding response, the client should send an exit notification that finally asks the server to quit their process.

For long-running operations, the specification also supports reporting progress on ongoing requests, and cancelling them. The progress reports not only allow indicating the status of the request to the user (e. g., by displaying it in the Language Client's UI), but also allows the server to return partial results to stream responses (e. g., showing the first few references to a variable while the rest are still loading).

The Language Client should also notify the Language Server whenever a document is opened or closed—that way, the server can, for example, start tracking diagnostics in the background as soon as a certain file is opened. There are also notifications related to modifications to the document that the Language Client should send, but these are irrelevant for us because modifying source code is out-of-scope for the LSP integration planned in this thesis.

Apart from describing fundamentals of the protocol like the above, the biggest part of the protocol's specification are the *capabilities*, that is, the JSON-RPC method names and types of the parameters and response objects for each feature that either the Language Client or Language Server can use [LspSpec].

2.2.2 Planned Capabilities

The capabilities I will make use of in SEE can roughly be grouped into the three categories *Navigation*, *Information*, and *Structure*. We will take a detailed look at each relevant capability below, along with how exactly they will be integrated into code cities. This level of detail is justified in the fact that the integration of the capabilities is the main focus of this whole thesis. Note that not all Language Servers support all capabilities—for example, for a language without a hierarchic type system, the capability *type hierarchy* cannot really be sensibly implemented.

Navigation This is the category containing the most capabilities that we can use for this thesis, since there are a lot of ways to navigate from one element within the source code to another. All of these take as input the position in a document, and return any number of locations, where the locations can contain a name, a file, and a range within the file.

All such available capabilities should appear in the context menu of SEE, either upon right-clicking a node, or right-clicking a code element in a code window. Selecting one of these navigation options from the menu should open a menu from which the user can select a single result. If there is only one result to begin with, this step should be skipped. Once a single result has been selected: If the request originates from a code window, the result should be opened and highlighted in that window. If the request instead originates from a code city, the node belonging to that result should be highlighted (e. g., by glowing and having a line pointed to it).

The other important part in SEE where this should be used is when building a city (i. e., when creating the project graph), as these capabilities provides us with the information we need to create (non-hierarchic) edges. Thus, we can create an edge e for each available navigation relation between two nodes, where $\ell(e)$ becomes the type of capability that was used. This is not as trivial as it sounds, since the ranges these capabilities return do not necessarily exactly match the actual ranges of the referenced nodes, so we need to implement some kind of matching algorithm (see section 3.2) .

TODO: More precise reference

The following navigation capabilities are available:

- **Call hierarchy:** Returns the incoming/outgoing calls for the symbol at the given location. Since incoming calls are already covered by the references capability, the context menu will only contain an option for showing outgoing calls.
- **Go to declaration:** Returns the declaration location for the symbol at the given location.
- **Go to definition:** Returns the definition location for the symbol at the given location.

- For this capability, another feature common in IDEs should be implemented in SEE’s code windows: Holding `Ctrl`, then clicking on a symbol, directly jumps to the definition of that symbol (or opens the corresponding selection menu, if there is more than one result).
- **Go to implementation:** Returns the definition location for the symbol at the given location.
- **Go to type definition:** Returns the location of the type definition for the symbol at the given location.
- **References:** Returns the location of the references to the symbol at the given location.
- **Type hierarchy:** Returns the sub/supertypes for the symbol at the given location. Since subtypes are already covered by the references capability, the context menu will only contain an option for showing supertypes.

Information Using these capabilities, the user can get information about either the project as a whole, or certain components of it. I have grouped the following capabilities into this category:

- **Diagnostics:** Returns diagnostics for a given file (e. g., warnings or errors). These will be integrated in exactly the same way as the Axivion Suite’s code smells in my bachelor’s thesis were [Gal21], that is:
 - In code cities, we will display erosion icons above affected nodes (see section 2.1.3).
 - In code windows, the corresponding parts of the source code should be highlighted, while hovering over the highlighted parts should reveal the diagnostic’s message and details.

Note that, instead of this being a proper request followed by a server response, this is only the case for the *pull diagnostics* capability, which has been added rather recently . The far more commonly used version is the *push diagnostics* capability, where the Language Server sends out diagnostics for the currently opened files at its own discretion, as notifications (see section 2.2.1)—this makes it difficult to collect this information during city construction, a topic we will explore in section 3.2.

- **Hover:** Returns hover information for a given location. The specification does not specify what exactly this “hover text” should be [LspSpec], but most implementations of Language Servers display the documentation of the hovered element, or the signature if it’s a method, or other helpful associated details. We can simply implement

TODO:
Inconsistent use
of should/will in
this section
TODO: When?

this part of the specification as intended: If the user hovers above an element in a code window, or a node in a code city, we should reveal the hover information in some kind of box near the mouse cursor, hiding it again once the cursor is moved away.

- **Semantic tokens:** Returns semantic tokens for the given file, which are intended for syntax highlighting. Similar to normal (e. g., lexer-based) syntax highlighting, requesting semantic tokens for a document yields a list of tokens containing their positions and a type, where the type can be one that is specified in the protocol (e. g., `enum`) or one that was previously announced as supported in the client capabilities. IDEs can then render each token type in a different color. An interesting addition to usual syntax highlighting is that each token can also be affected by any number of *token modifiers*, where each modifier may add an additional rendering effect on top of the type-based color. For example, tokens with the `static` modifier might be rendered in *italics*, while ones with the `deprecated` modifier could be rendered in ~~strikethrough~~.

SEE currently uses ANTLR⁷-based syntax highlighting, where we need to manually group each parser’s token into some categories to determine colors. The added value of the semantic tokens capability here would be ease of use (i. e., no need to manually configure each Language Server) on the one hand, and support for token modifiers (i. e., “extended” syntax highlighting) on the other hand.

Structure This category actually only comprises a single capability—one that we can use to build the hierarchy of the code city’s project graph (outlined in section 2.1.2), because it gives us information about the structure of the project. The capability I am talking about here is the *document symbols* capability. Given a document, it will return all symbols present within that file, along with some additional information for each symbol, such as its type or range.

There are actually two different possible kinds of symbols that this capability may return:

1. an array of `SymbolInformation`, which is “a flat list of all symbols” (`[LspSpec]`) that should not be used to infer a hierarchy. Because of this limitation, if a Language Server is only able to return symbols of this data type, we cannot use it to build code cities. This is an older data type, and instead modern Language Servers should rather return
2. an array of `DocumentSymbol`. This contains a field `children`, which stores `DocumentSymbols` that are contained in this one. Using this property, we can establish a hierarchy and

⁷<https://www.antlr.org> (last access: 2024-06-10)

build a code city by recursively enumerating all `DocumentSymbols` and their children for each file, then querying for all relevant information by using the other capabilities outlined above.

2.2.3 Unplanned Capabilities

There are a number of capabilities that I will not implement into SEE. These can be grouped into roughly three categories, based on the reasoning behind them being unused: The first concerns those capabilities that relate to editing only and provide no features related to simply viewing code—as explained in section 1.3, editing code is not part of the goals here and requires additional large-scale preparatory changes to SEE. The second concerns the complex capabilities, that is, ones whose implementation would take a lot of time and effort and thus go beyond the scope of this master’s thesis. The third concerns the niche capabilities that provide only a very marginal benefit, or do so only in rare situations. For these, I also have not deemed the effort worth it to implement them, at least not as part of this thesis.

I will quickly list the contents of all these groups here.

Editing capabilities

- **Code Actions:** Allows the programmer to apply refactoring actions to the code, such as importing a referenced library.
- **Completion:** Computes autocomplete items while the user is typing, and applies them when chosen.
- **Formatting:** Applies automatic formatting to a file, or range, of code.
- **Rename:** Executes a project-wide rename of a symbol, which also renames all references to that symbol.
- **Linked Editing Range:** Returns a list of ranges that will be edited upon executing a rename of a symbol, with the purpose of highlighting those ranges during the rename.
- **Signature Help:** Returns signature information at the given cursor information. This may seem relevant for our purposes, but it is actually intended to be shown while editing (e. g., highlighting the active parameter as one types), and its information is given by the hover capability anyway in almost all cases.

New Term: *LSIF* (Language Server Index Format)

A format which language servers can emit to persist LSP-based information about a software project.

New Term: *DAP* (Debug Adapter Protocol)

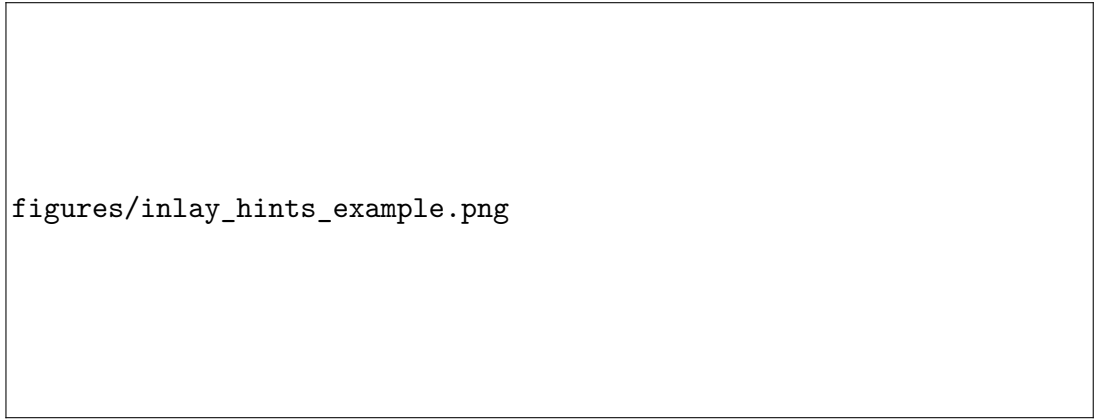
A protocol that can be viewed as the analogue to LSP for debuggers, with the goal to make it easier to integrate debuggers into development tools.

Niche capabilities

- **Document Color:** Lists all color references in the code (e. g., symbolic references like `Colors.red`) along with their color value in the RGB format.
- **Color Presentation:** Allows users to modify color references in the code by using a color picker.
- **Document Link:** Returns the location of links in the document.
- **Code Lens:** Returns commands that can be shown next to source code, such as the number of implementers of an abstract method.
- **Monikers:** This is a description of what symbols a project imports and which ones it exports, and is intended to make relations between multiple projects possible. As LSP usually only deals with a single project at a time, this is more useful in the *Language Server Index Format* (*LSIF*), whose specification it also originates from [*LsifSpec*].

Complex capabilities

- **Folding Range:** Returns ranges that can be collapsed in the code viewer. For example, the contents of a function could be collapsed, leaving only its signature visible. Due to the way code windows are implemented in SEE, this would increase the complexity of the implementation quite a bit.
- **Inline Value:** In debugging contexts, this supplies the contents of a variable with the purpose of displaying them inline in the Language Client, next to the variable itself. It uses the *Debug Adapter Protocol* (*DAP*) [*DapSpec*], which has previously already been integrated into SEE [*Roh22*], but it would take a lot of refactoring work to make the two implementations compatible with each other.
- **Inlay Hint:** Returns textual hints that can be rendered within the source code. An example would be parameter names that are intended to be shown at the call site, such as in the screenshot in fig. 2.1.



`figures/inlay_hints_example.png`

Figure 2.1: An example of inlay hints in the JetBrains IntelliJ IDE. (From <https://www.jetbrains.com/help/idea/inlay-hints.html> (*last access: 2024-10-04*))

- **Notebook-related capabilities:** These are intended for interactive notebook systems such as Jupyter⁸, which SEE does not currently support.

2.3 Interim Conclusion

In this chapter, we have taken a detailed look at the concepts behind the two topics central to this thesis—namely, the Language Server Protocol and code cities. We have also motivated and laid out the specific ways in which the existing LSP capabilities could be integrated into SEE, including a formalization of the project graph that will become central to section 3.2.

We are now ready to tackle the actual implementation in the next chapter. As a quick overview and recap before then, there is a table of planned capabilities along with their intended use in SEE in table 2.1.

⁸See <https://jupyter.org/> (*last access: 2024-10-03*)

Table 2.1: LSP capabilities that will be integrated into SEE as part of this thesis.

Capability	Code Windows	Code Cities
<i>Call hierarchy</i>	Show incoming/outgoing calls and allow jumping to caller	Generate corresponding edges
<i>Diagnostics</i>	Highlight corresponding code ranges and display details on hover	Display code smell icons [see Gal21]
<i>References</i>	Show references and allow jumping to usage	Generate corresponding edges
<i>Document symbols</i>	—	Generate corresponding nodes and hierarchy
<i>Go to location*</i>	Show locations and allow jumping to them	Generate corresponding edges
<i>Hover</i>	Show hover information when hovering above item	Show hover information when hovering above node
<i>Semantic tokens</i>	Extended (“semantic”) syntax highlighting	—
<i>Type hierarchy</i>	Show sub-/supertypes and allow jumping to them	Generate corresponding edges

*This includes the *Go to declaration / definition / implementation / type definition* capabilities.

Implementation



ELYING on the foundations of SEE and the Language Server Protocol established in the previous chapter, we can now turn to the core part of this thesis: The integration of LSP into SEE, with a special focus on how to build code cities using LSP’s capabilities. We will start by briefly going over some preliminary changes to both SEE and the LSP specification. Then, we will spend the majority of this chapter specifying and explaining the algorithm which “converts” LSP information into code cities, before looking into how additional capabilities can be integrated into SEE’s code cities and code windows specifically. Finally, we will conduct a brief technical evaluation, with a more thorough user study following in the next chapter.

3.1 Preliminary Changes

As promised in the preceding paragraph, we will first quickly list some preparations.

3.1.1 Specification Cleanup

While familiarizing myself with the LSP specification, I noticed and fixed a few small issues along the way. Most of these were of a formal nature (e. g., spelling, grammar, formatting, consistent usage of terms), some were fixing incorrect TypeScript syntax in the definition of LSP’s data models. The rest of the changes were related to the so-called snippet grammar.

In the context of LSP, snippets are in essence string templates that are inserted on certain completions (see section 2.2.3), with some designed parts being filled in by the programmer on insertion. There are also parts that can be filled in by certain values (e. g., the file name), which can themselves be transformed using regular expressions.¹ The complexity of the

¹I am skipping over some additional features and details here because this is not that relevant a capability for us—to get the full picture, see https://microsoft.github.io/language-server-protocol/specifications/lsp/3.18/specification/#snippet_syntax (last access: 2024-10-10).

New Term: *EBNF* (Extended Backus–Naur Form)

A syntax in which context-free grammars can be formally expressed.

combinations of all these features increase the possibility of misunderstandings, which is why the snippet's grammar has been formally specified in *Extended Backus–Naur Form* (**EBNF**). However, as it was written down in the specification, the grammar had a few problems that I have fixed. Three notable examples are:

- Some alternatives were incorrectly grouped, contradicting the explanatory text above them. Also, the rules on how and when control characters had to be escaped were inconsistent with the surrounding text.²
- The grammar contained some string transformations that were unexplained in the text. Since the LSP specification is based on VSCode, I added explanations to the text based on what these transformations did in VSCode's source code.
- Finally, there were ambiguities present in the grammar that led to FIRST/FOLLOW conflicts. I have rewritten the grammar to eliminate these, and it should now be *LL*(1)-parseable [Aho+07, pp. 222–224].

I have submitted these fixes as a pull request³. After addressing the resulting code review, it has been merged, and the changes will be incorporated in the upcoming 3.18 release of the specification.

3.1.2 Preparing SEE

There was not much I had to do in terms of getting SEE ready, so this section will be very short:

- I have integrated the OmniSharp LSP C# library⁴ into SEE, which we will leverage in the subsequent sections so that we can use LSP without needing to worry about JSON-RPC encoding, data models, and so on.
- Code windows have previously been made editable by Blecker [Ble22] in his bachelor thesis, also enabling collaborative editing over the internet. I unfortunately had to remove these changes because they did not work anymore in the current version of

²This has lead to confusion in some projects making use of snippets. See, for example, <https://github.com/neovim/neovim/issues/30495> (last access: 2024-10-10).

³<https://github.com/microsoft/language-server-protocol/pull/1886> (last access: 2024-10-10)

⁴Available at <https://github.com/OmniSharp/csharp-language-server-protocol> (last access: 2024-10-11)

New Term: *interval tree*

A data structure meant to store intervals/ranges in such a way that overlapping or contained intervals can be found efficiently.

SEE, and additionally caused a lot of complexity overhead in the code window implementation that would have made the LSP integration much harder to accomplish.

- Finally, the attribute space \mathcal{A} in SEE did not allow for ranges of the form LSP needs, so I had to replace the existing attributes (which track the line and column, but not a full range) with a proper set of range attributes. In section 2.1.2, we have introduced this as a single `Source.Range` attribute, but in reality, there are four range attributes—one per member of the decomposed form. We will ignore this reality for the rest of this thesis and act like the range is a single attribute, that is, for all project graphs with nodes V and attributes a , it holds that $\{a(v, \text{Source.Range}) \mid v \in V \wedge (v, \text{Source.Range}) \in \text{dom}(a)\} \subseteq \mathcal{R}$.

All of these changes have been made across two pull requests to SEE.⁵

3.2 Generating Code Cities using LSP

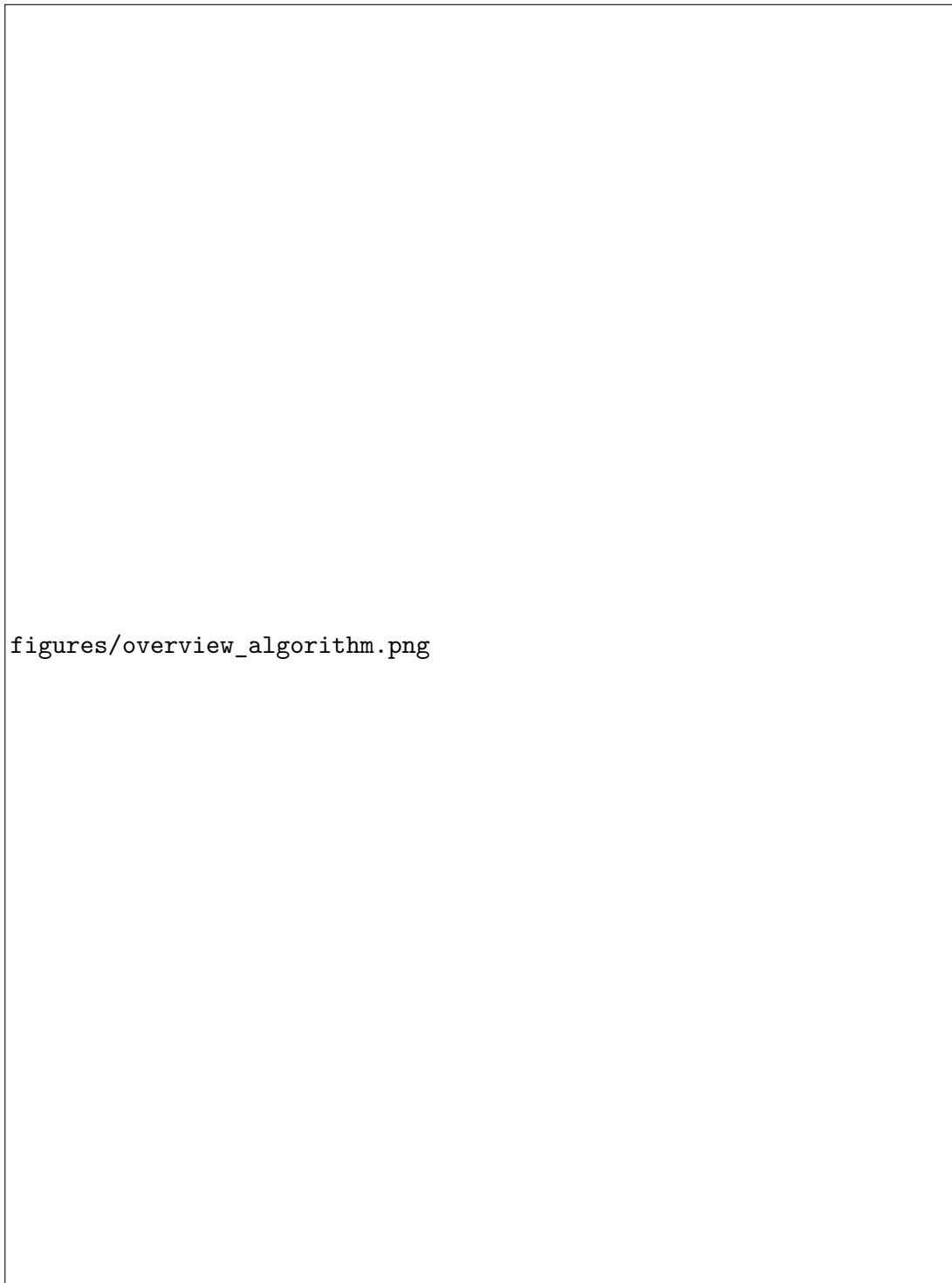
In this section, we will examine the centerpiece of this thesis: The algorithm with which code cities can be generated using the Language Server Protocol. While going over how the algorithm works, we will take a quick look at *interval trees* and how they relate to the algorithm, before finally taking a look at what the import process looks like in practice in SEE.

3.2.1 Algorithm

We will take a look at the algorithm in a generalized and programming language independent form here. In this form, the algorithm takes as input a set of source code documents, as well as a family of LSP functions belonging to a specific instantiation of a Language Server. These functions will be used to analyze the documents and extract the required information from them. The output of the algorithm, then, is a graph representing the given software project.

TODO: Replace diagram sketch with TikZ picture.

⁵<https://github.com/uni-bremen-agst/SEE/pull/687> and <https://github.com/uni-bremen-agst/SEE/pull/715> (last access: 2024-10-11).



figures/overview_algorithm.png

Figure 3.1: A high-level overview of the basic steps of the algorithm.

Overview Before diving into the specifics of *how* the algorithm works, it may help to take a look at the diagram in fig. 3.1, which gives us a high-level overview of *what* it does. To summarize, the steps can be broken down into three major parts:

I Node synthesis: Here, we create the graph's nodes and combine them together into a hierarchy.

1. We recreate the parts of the filesystem hierarchy that are relevant to the given documents (i. e., directories the documents are contained in and their relation to each other).
2. For each code symbol within that document, a node will be created as a child to the document. Any symbols contained within that symbol are recursively added in the same manner as a child to their parent nodes.
3. Finally, we will pull diagnostics for the document⁶ and attach their counts to the nodes they correspond to.

II Edge synthesis: Here, we connect the nodes by creating edges between them. To do this, we go over each node, using LSP functions to check for definition locations, references, and so on, and then determine which of our existing nodes best corresponds to that location. This turns out to be the most difficult and complex part of the algorithm to implement efficiently, as we will later see.

III Aggregation: Finally, we want to aggregate LOC and diagnostic count metrics upwards in the hierarchy. This way, we can, for example, see how many diagnostics are contained as a whole in a class, or even a directory.

Specification Now that we know what it is the algorithm does, we can take a look at its detailed specification. The specification is given in algorithm 1 and follows a few special formatting rules that I will briefly list here:

- Text in SMALL CAPS refers to functions. Those starting with LSP specifically refer to functions provisioned by the Language Server.
- Sentences in *gray italics* are comments.
- **Bold** text represents keywords.
- A normal font represents strings (i. e., text that represents itself).

⁶In the actual algorithm, we cannot rely on pulling diagnostics alone, since only few Language Servers support this. Instead, we will collect pushed diagnostics in the background and handle them all at once at the very end.

- Finally, text in a typewriter font have two purposes: They represent attribute keys ($\in \mathcal{A}_K$) as well as properties of LSP-returned objects, the latter of which are prefixed by a dot.

Page 31 contains the main algorithm. We can map fig. 3.1 onto it as follows: Lines 1 to 12 contain part I (node synthesis), lines 12 to 34 contain part II (edge synthesis), and finally, lines 34 to 38 contain part III (aggregation) as well as handling of pushed diagnostics.

The following pages 32 to 33 contain the functions referenced in the main algorithm. Here, we represent the “types” of each parameter by specifying the function domain, that is, by noting the sets each parameter must be in. Most of these sets are already defined in chapter 2 or the algorithm itself, but two exceptions to this are the set of LSP code symbols \mathcal{S} and the set of LSP diagnostics \mathcal{D} (not to be confused with the set of input documents \mathcal{D}).

Simplifications As mentioned before, algorithm 1 is a generalized version of the actual C# algorithm that was implemented into SEE. Hence, a few simplifications⁷ were made to not make this section even more technical and longer than it already is. Some noteworthy omissions are:

- Details on adding nodes, such as the definition of the `NEWNODE` function, or the unique IDs assigned to each node.
- The configuration of the algorithm. We present it as only having two input parameters, but in reality, there are many additional explicit and implicit parameters (e. g., importing only certain types of nodes). We will take a look at some of these configuration parameters that are relevant for the user in section 3.2.3.
 - Actually, even the two input parameters that we did include are simplifications. In truth, the user selects a directory (instead of a set of documents), with the option to exclude some subdirectories, and we then automatically include every file with an extension that is relevant to the selected Language Server.
- Progress reporting. A progress bar noting the approximate progress (in percent) appears in both the Unity Editor’s UI as well as in-game while the code city is constructed, so the user has an idea of how long the conversion is going to take.
- Asynchronicity. Specifically, this refers to the way the algorithm is executed in Unity. If we were to implement it as a normal synchronous function, the UI would become irresponsive to input and freeze until the whole algorithm is done. Instead, we make use of C#’s task-based `async` capabilities [Wag23] in combination with the `UniTask`⁸

⁷Apart from the obvious simplifications that occur naturally due to the difference between declarative mathematical notation and imperative programming syntax.

⁸<https://github.com/Cysharp/UniTask> (last access: 2024-10-25)

Algorithm 1 How code city graphs can be generated from LSP information.

Input: Family of LSP functions provided by the Language Server, set of documents D

Output: Graph G representing the underlying software project

```
1  $V, E, a, s, t, \ell, C \leftarrow \emptyset$  ▷ Initialize empty graph components.
2 for all  $d \in D$  do
3    $\text{LSPOPENDOCUMENT}(d)$  ▷ Document needs to be opened for all capabilities to work.
4    $v_d \leftarrow \text{ADDDOCUMENTNODE}(d)$  ▷ Each document becomes a node...
5   for all  $x \in \text{LSPDOCUMENTSYMBOLS}(d)$  do
6      $\text{MAKECHILD}(\text{ADDSYMBOLNODE}(x), v_d)$  ▷ ... with its symbols as children.
7   if  $\text{LSPLANGUAGESEVERSUPPORTSPULLDIAGNOSTICS}()$  then
8      $\text{HANDLEDIAGNOSTICS}(\text{LSPPULLDOCUMENTDIAGNOSTICS}(d))$ 
9   else
10    ▷ We will save any incoming diagnostics in the background and handle them at the end. ◀
11     $\text{LSPREGISTERPUSHDIAGNOSTICSCALLBACK}(d, (c) \mapsto (C \leftarrow C \cup \{c\}))$ 
12     $\text{LSPCLOSEDOCUMENT}(d)$ 

13 for all  $v \in V : (v, \text{Source.Range}) \in \text{dom}(a)$  do
14   ▷ First, connect nodes to each other based on LSP relations. ◀
15    $\text{CONNECTNODEVIA}(\text{LSPGoToDefinition}, \text{Definition}, v)$ 
16    $\text{CONNECTNODEVIA}(\text{LSPGoToDeclaration}, \text{Declaration}, v)$ 
17    $\text{CONNECTNODEVIA}(\text{LSPGoToTypeDefinition}, \text{TypeDefinition}, v)$ 
18    $\text{CONNECTNODEVIA}(\text{LSPGoToImplementation}, \text{Implementation}, v)$ 
19    $\text{CONNECTNODEVIA}(\text{LSPReferences}, \text{Reference}, v)$ 

20   if  $a(v, \text{Type}) = \text{Method}$  then ▷ We need to integrate the call hierarchy into the graph.
21      $I \leftarrow \text{LSPPREPARECALLHIERARCHY}(a(v, \text{Source.Path}), a(v, \text{Source.Range}))$ 
22     ▷ GETMATCHINGITEM returns the item in I with the same name and location as v. ◀
23      $i \leftarrow \text{GETMATCHINGITEM}(I, v)$ 
24      $R \leftarrow \text{LSPCALLHIERARCHYOUTGOINGCALLS}(i)$ 
25      $V' \leftarrow \bigcup_{r \in R} \text{FINDNODESBYLOCATION}(r.\text{path}, r.\text{range})$ 
26     for all  $v' \in V'$  do
27        $\text{ADDEDGE}(v, v', \text{Call})$ 
28   else if  $a(v, \text{Type}) = \text{Type}$  then ▷ We need to integrate the type hierarchy into the graph.
29      $I \leftarrow \text{LSPPREPARETYPEHIERARCHY}(a(v, \text{Source.Path}), a(v, \text{Source.Range}))$ 
30      $i \leftarrow \text{GETMATCHINGITEM}(I, v)$ 
31      $R \leftarrow \text{LSPTYPEHIERARCHYSUPERTYPES}(i)$ 
32      $V' \leftarrow \bigcup_{r \in R} \text{FINDNODESBYLOCATION}(r.\text{path}, r.\text{range})$ 
33     for all  $v' \in V'$  do
34        $\text{ADDEDGE}(v, v', \text{Extend})$ 

35  $\text{HANDLEDIAGNOSTICS}(C)$  ▷ Handle diagnostics that were collected in the background.
36  $\text{AGGREGATEMETRICS}(\{\text{Metric.LOC}\})$ 
37  $\text{AGGREGATEMETRICS}(\{\text{ErrorCount}, \text{WarningCount}, \text{InformationCount}, \text{HintCount}\})$ 
38 return  $(V, E, a, s, t, \ell)$ 
```

```

39 function ADDDOCUMENTNODE( $d \in D$ )
40    $v_d \leftarrow \text{NEWNODE}()$ 
41    $a' \leftarrow \emptyset$ 
42    $a'(v_d, \text{Type}) \leftarrow \text{File}$ 
43    $a'(v_d, \text{Source.Path}) \leftarrow d$ 
44    $a'(v_d, \text{Metric.LOC}) \leftarrow |\text{READLINES}(d)|$   $\triangleright \text{READLINES returns the set of lines in the file.}$ 
45    $V \leftarrow V \cup \{v_d\}$ 
46    $a \leftarrow a \cup a'$ 
47    $\text{MAKECHILD}(v_d, \text{ADDDIRECTORYNODE}(d.\text{directory}))$ 
48   return  $v_d$ 

49 function ADDSYMBOLNODE( $x \in S$ )
50    $v \leftarrow \text{NEWNODE}()$ 
51    $a' \leftarrow \emptyset$ 
52    $a'(v, \text{Source.Name}) \leftarrow x.\text{name}$ 
53    $a'(v, \text{Source.Path}) \leftarrow d$ 
54    $a'(v, \text{Type}) \leftarrow x.\text{kind}$ 
55    $a'(v, \text{Deprecated}) \leftarrow (\text{deprecated} \in x.\text{tags})$ 
56    $a'(v, \text{Source.Range}) \leftarrow x.\text{range}$ 
57    $a'(v, \text{Metric.LOC}) \leftarrow e_l^{x.\text{range}} - b_l^{x.\text{range}}$ 
58    $\triangleright \text{Several other similar attributes omitted here...}$   $\triangleleft$ 
59    $a'(v, \text{HoverInfo}) \leftarrow \text{LSPHOVER}(d, x.\text{range})$ 
60   if  $a' \not\subseteq a$  then  $\triangleright \text{If an isomorphic node does not already exist...}$ 
61      $V \leftarrow V \cup \{v\}$   $\triangleright \text{...add it and handle its children.}$ 
62      $a \leftarrow a \cup a'$ 
63     for all  $x' \in x.\text{children}$  do
64        $\text{MAKECHILD}(\text{ADDSYMBOLNODE}(x'), v)$   $\triangleright \text{Recurse.}$ 
65   return  $v$ 

66 function MAKECHILD( $v_c \in V, v_p \in V$ )
67    $\triangleright \text{The partOf edges must induce a tree structure. Hence, if a node already is a part of another}$   $\triangleleft$ 
68    $\triangleright \text{node, we must not add another partOf edge.}$ 
69   if  $\exists e \in E : \ell(e) = \text{partOf} \wedge s(e) = v_c$  then
70     output Warning: Hierarchy is cyclic. Some children will be omitted.
71   else
72      $\text{ADDEDGE}(v_c, v_p, \text{partOf})$ 

73 function CONNECTNODEVIA( $\text{LSPFUN} \in (D \times \mathcal{R})^{D \times \mathcal{R}}, l \in \Sigma, v \in V$ )
74    $\triangleright \text{Function LSPFUN only returns locations, so we need to find the relevant nodes first.}$   $\triangleleft$ 
75   for all  $(d, r) \in \text{LSPFUN}(a(v, \text{Source.Path}), a(v, \text{Source.Range}))$  do
76     for all  $v' \in \text{FINDNODESBYLOCATION}(d, r)$  do
77        $\text{ADDEDGE}(v, v', l)$ 

78 function ADDEDGE( $v_s \in V, v_t \in V, l \in \Sigma$ )
79    $e' \leftarrow \text{NEWEDGE}()$ 
80    $E \leftarrow E \cup \{e'\}$ 
81    $s(e') \leftarrow v_s$ 
82    $t(e') \leftarrow v_t$ 
83    $\ell(e') \leftarrow l$ 

```

```

1 function FINDNODESBYLOCATION( $d \in D, r \in \mathcal{R}$ )
2    $\triangleright$  We pick the nodes with the most specific range containing the given location.  $\triangleleft$ 
3    $\text{getR}(v) = a(v, \text{Source.Path})$ 
4    $N \leftarrow \{v \in V \mid a(v, \text{Source.Path}) = d \wedge b_l^r, e_l^r \in [b_l^{\text{getR}(v)}, e_l^{\text{getR}(v)}]$ 
5      $\wedge b_c^r \geq b_c^{\text{getR}(v)} \wedge e_c^r \leq e_c^{\text{getR}(v)}\}$ 
6    $N \leftarrow \arg \min_{v \in N} e_l^{\text{getR}(v)} - b_l^{\text{getR}(v)}$ 
7   return  $\arg \min_{v \in N} e_c^{\text{getR}(v)} - b_c^{\text{getR}(v)}$ 

7 function ADDDIRECTORYNODE( $p \in \mathcal{A}_V$ )
8   if  $\exists! v \in V : a(v, \text{Source.Path}) = p$  then  $\triangleright$  Check if node exists already.
9     return  $v$   $\triangleright$  If so, just pick that one.
10   $v_p \leftarrow \text{NEWNODE}()$ 
11   $a' \leftarrow \emptyset$ 
12   $a'(v_p, \text{Type}) \leftarrow \text{Directory}$ 
13   $a'(v_p, \text{Source.Path}) \leftarrow p$ 
14   $V \leftarrow V \cup \{v_p\}$ 
15   $a \leftarrow a \cup a'$ 
16   $v_p^* \leftarrow \text{ADDIRECTORYNODE}(\text{GETPARENTDIRECTORY}(p))$   $\triangleright$  Recurse to add parent directories.
17  if  $v_p \neq v_p^*$  then
18     $\text{MAKECHILD}(v_p, v_p^*)$ 
19  return  $v_p$ 

20 function HANDLEDIAGNOSTICS( $d \subset \mathcal{D}$ )
21   for all  $c \in d$  do
22      $V_c \leftarrow \text{FINDNODESBYLOCATION}(c.\text{path}, c.\text{range})$ 
23     for all  $v \in V_c$  do  $\triangleright$  Save diagnostics count (grouped by severity) in all affected nodes.
24        $n \leftarrow c.\text{severity} + \text{Count}$   $\triangleright$  Concatenate Count to attribute name.
25       if  $(v, n) \in \text{dom}(a)$  then
26          $a(v, n) \leftarrow a(v, n) + 1$ 
27       else
28          $a(v, n) \leftarrow 1$ 

29 function AGGREGATEMETRICS( $M \subset \mathcal{A}_K$ )
30   for all  $v \in V : \nexists e \in E : t(e) = v \wedge \ell(e) = \text{partOf}$  do  $\triangleright$  Aggregate from each root node.
31     for all  $m \in M$  do
32        $\text{AGGREGATEMETRICFROMROOT}(m, v)$ 

33 function AGGREGATEMETRICFROMROOT( $m \in \mathcal{A}_K, v_r \in V$ )
34    $V_c \leftarrow \{v \in V \mid \exists e \in E : t(e) = v_r \wedge s(e) = v \wedge \ell(e) = \text{partOf}\}$   $\triangleright$  Immediate children.
35   for all  $v \in V_c$  do
36      $\text{AGGREGATEMETRICFROMROOT}(m, v)$ 
37    $\triangleright$  After the recursion above, immediate children now definitely have attribute  $m$ .  $\triangleleft$ 
38   if  $(v_r, m) \notin \text{dom}(a)$  then  $\triangleright$  We don't want to overwrite existing metrics.
39      $a(v_r, m) \leftarrow \sum_{v \in V_c} a(v, m)$ 

```

framework: We yield control to the Unity event loop when progress is suspended (e. g., while we wait for an answer from the Language Server), allowing frames to be rendered while the algorithm is running in the background. This also allows us to implement cancellation support, making it possible for the user to cancel the algorithm at any time⁹.

The full C# implementation in SEE is available at <https://github.com/uni-bremen-agst/SEE/blob/c4e3de908a022d65723bf82d3b350dade8b5f01a/Assets/SEE/DataModel/DG/I0/LSPImporter.cs> (last access: 2024-10-25).

TODO: Convert all software mentions to biblatex @software type.

Performance Considerations When we analyze this algorithm in terms of complexity, we can quickly see that the most relevant portions are in matching locations to nodes, that is, `FINDNODESBYLOCATION`¹⁰. This is only meant to give a quick motivation on why we need the optimization described in section 3.2.2—a full complexity analysis of the algorithm is outside the scope of this thesis. Part I as a whole (ignoring diagnostics, where `FINDNODESBYLOCATION` is called) can be considered to fall in $\Theta(|V|)$, since there is a constant amount of work per added node. Part III (again ignoring diagnostics), as written, yields a worst-case runtime in $O(|V|^2 \cdot |E|)$, because we potentially need to search through all edges for each node to identify child nodes, and this happens once per node while aggregating metrics upwards. However, in the actual implementation, child nodes are saved alongside their parent and can be accessed in $O(1)$, so this reduces to a runtime of $\Theta(|V|)$.

Part II is where things get interesting: We are calling `CONNECTNODEVIA` a few times for each node (ignoring the handling of the type/call hierarchy, where very similar things happen), so let us look at what happens in here. `CONNECTNODEVIA` first retrieves all target locations from the given LSP function, and then calls `FINDNODESBYLOCATION` for each of those locations to identify the node in our graph to which the connecting edge should be drawn. This effectively means `CONNECTNODEVIA` is called once per added (non-partOf) edge. In this function, each node's range is compared to the given location to check whether it is contained therein. We then take the minimum over these nodes twice to determine the nodes with the most specific fit, so in total, the runtime of this function is in $\Theta(|V|)$. This function is called once per added edge, and also once per diagnostics (because diagnostics also need to be associated to nodes), so part II's runtime can be said to be in $\Theta((|E| + n_d) \cdot |V|)$, where n_d refers to the total number of diagnostics for the project.

Hence, the runtime for algorithm 1 as a whole lies in $\Theta(|V| + |E| \cdot |V| + n_d \cdot |V| + |V|)$. In practice for LSP-built code cities, assuming the default configuration, $n_d < |V|$, but

⁹Internally, we do this by checking every so often if a so-called *cancellation token* has been revoked by the user. If it has, we halt execution.

¹⁰At least, this is the case when we assume the runtime complexity of externally supplied LSP functions is constant. This is an oversimplification, but the claim that this function is the most expensive part of the algorithm holds up to analyses of real-world test runs of the C# algorithm.

New Term: *k-d tree*

A data structure (using a binary tree as a basis) with which certain spatial data in k dimensions can be efficiently stored and retrieved.

$|E| \gg |V|$. Hence, part II with $\Theta(|E| \cdot |V|)$ easily dwarfs the rest of the algorithm's runtime due to the high cost of `FINDNODESBYLOCATION`, which searches through every node for every added edge. For this reason, it would be nice to optimize that function somehow.

TODO: Put in examples or more concrete estimation of how much more edges there are than nodes.

TODO: Do actual comparison between two ways of finding target nodes. Also mention example times from eval section at end.
TODO: Rewrite “the above” to something else, may not be above in printed version

3.2.2 Augmented Interval Trees

To restate the problem outlined in the performance considerations **above**: The locations returned by LSP functions such as “show references” or “go to location” need to be matched to the nodes in our constructed project graph. We cannot simply create a lookup table from locations to nodes, as the locations are not necessarily equal to the location of the nodes, even when they describe the same logical element. Instead, we need to match the location to the node with the “tightest fitting range”, that is: from the nodes whose range completely contains the given location, we pick the one whose range is the smallest (to find the most specific fitting element). The naive solution used in algorithm 1—simply going over all nodes each time to find the best fit—leads to an unacceptable runtime.

There are a few possible ways to solve this problem (which is in essence a variant of the stabbing problem) more efficiently, such as segment trees or range trees, but we will use augmented interval trees with *k-d trees* as a base, as this configuration best fits our circumstances—there is no need to update/re-balance the tree (so the high cost associated with that is fine), the membership query should not be in $\Omega(n)$, we need to represent two dimensions (which is possible with a 2d-tree), and so on.

A detailed explanation of augmented interval trees can be found in Cormen et al. [Cor+22, section 17.3], while *k-d trees* are described in a paper by Bentley [Ben75]. In our case, we can create the tree by constructing a 2d-tree¹¹ out of all elements (as explained in the previous sources), where the key is the starting position of the range. Afterwards, we save in each node the maximum line number and maximum character offset, respectively, for the subtree rooted by that node, thereby turning this *k-d tree* into an interval tree.

An updated `FINDNODESBYLOCATIONEFFICIENT` is given in algorithm 2, with these details:

- An augmented interval tree is modeled here as a quintuple $T = (v, l, c, \lambda, \rho)$. The set of all such trees is defined as \mathcal{T} , so $T \in \mathcal{T}$. The first element $v \in V$ is the node in

¹¹Note that, in the actual implementation, the construction of the *k-d tree* (excluding its augmentation to an interval tree) is handled by the external `Supercluster.KDTree` library (<https://github.com/ericreg/Supercluster.KDTree> (last access: 2024-10-31)).

the project graph represented by the node in this tree. The second element $l \in \mathbb{N}_0$ refers to the maximum line number across all nodes within this tree, whereas the third element $c \in \mathbb{N}_0$ analogously refers to the maximum character offset. The fourth element $\lambda \in \mathcal{T}$ refers to the left subtree rooted by this node, while the fifth element $\rho \in \mathcal{T}$ conversely refers to the right subtree. Taking a single element x from the tuple T is written as x_T .

- We construct such an augmented interval tree for each document (directly after part I in algorithm 1) and save them all in a hash table H (modeled as a function $H : D \rightarrow \mathcal{T}$), where each document maps to its interval tree.
- Checking whether a range r_1 is contained in another range r_2 is written as $r_1 \subseteq r_2$.
- We need a way to compare two ranges against each other to check which one is “more specific”. The difficulty here is that the character offsets are hard to compare to each other, since the lines can be of different lengths—handling different line lengths correctly would be both algorithmically more expensive and harder to implement, so we have given up transitivity: We define a homogeneous relation \lesssim on \mathcal{R} that is anti-reflexive and asymmetric (but not necessarily transitive), where $r_1 \lesssim r_2$ implies that r_1 is more specific than r_2 . Similarly, \approx is a homogeneous relation on \mathcal{R} that is reflexive and symmetric (but also not necessarily transitive), where $r_1 \approx r_2$ if they are both “equally as specific”.¹²

Algorithm 2 Efficiently associating LSP-returned ranges to existing elements using an already constructed augmented interval tree.

```

1 function FINDNODESBYLOCATIONEFFICIENTEFFICIENT( $d \in D, r \in \mathcal{R}$ )
2   return QUERYTREE( $H(d), r, \emptyset$ )

3 function QUERYTREE( $T \in \mathcal{T}, q \in \mathcal{R}, R \subseteq V$ )
4    $r \leftarrow a(v_T, \text{Source.Range})$ 
5   if  $b_l^q > l_T \vee (b_l^q = l_T \wedge b_c^q \geq c_T)$  then           ▷ Range is to the right of all nodes in this subtree.
6     return  $R$ 
7   if  $q \subseteq r$  then                                         ▷ Range is contained in this node, but we want only minimal fits.
8      $m \leftarrow \{v \in R \mid r \lesssim a(v, \text{Source.Range})\}$ 
9     if  $|m| > 0$  then                                       ▷ This range is smaller than other results.
10       $R \leftarrow (R \setminus m) \cup \{v_T\}$ 
11     else if  $\forall v \in R : a(v, \text{Source.Range}) \approx r$  then      ▷ Other ranges are equally minimal.
12        $R \leftarrow R \cup \{v_T\}$ 
13     ▷ Otherwise, r is not minimal, so we don't add the node. ◀
14    $R \leftarrow \text{QUERYTREE}(\lambda_T, q, R)$ 
15   if  $e_l^q \geq b_l^r \wedge (e_l^q \neq b_l^r \vee e_c^q > b_c^r)$  then       ▷ Range could be contained in right subtree.
16      $R \leftarrow \text{QUERYTREE}(\rho_T, q, R)$ 
17   return  $R$ 

```

TODO: So how much does it help in practice? Reference tech eval here

¹²My actual implementation of the CompareTo function can be found here: <https://github.com/uni-bremen-agst/SEE/blob/c4e3de908a022d65723bf82d3b350dade8b5f01a/Assets/SEE/DataModel/DG/Range.cs> (last access: 2024-10-31).

Listing 3.1: Example C# source code with demarcated symbol ranges.

```

0 +\textcolor{red}{|}|+public class Example {
1   +\textcolor{Emerald}{|}|+const char someValue = 'A'+\textcolor{Emerald}{|}|+;
2
3   +\textcolor{Peach}{|}|+public static long pow(+\textcolor{OliveGreen}{|}|+int
   ↪ num+\textcolor{OliveGreen}{|}|+) {
4     +\textcolor{Brown}{|}|+long result = num * num+\textcolor{Brown}{|}|+;
5     return result;
6   }+\textcolor{Peach}{|}|+
7 }+\textcolor{red}{|}|+

```

As a concrete example on how this works, take a look at the example code in listing 3.1. Here, we have the following elements:

- The Example class with range (0, 0, 7, 1).
- The someValue field with range (1, 2, 1, 28).
- The pow method with range (3, 2, 6, 3).
- The num parameter with range (3, 25, 3, 32).
- The result variable with range (4, 4, 4, 27).

Now, if we were to convert this into an augmented 2-d interval tree, it might look like fig. 3.2. Here, the key refers to the starting position of each element and the max value refers to the maximum line and maximum character offset in the subtree rooted by each node. Let us say we want to find out what the range (1, 13, 1, 22) (comprising just the name of someValue) belongs to. We will start by checking the root node pow:

1. $b_l^q < l_{\text{pow}}$ because $1 < 7$, so the range is not to the right of all nodes.
2. It is not contained by pow's range.
3. We query the left subtree $\lambda_{\text{pow}} = \text{someValue}$:
 - a) $b_l^q < l_{\text{someValue}}$ because $1 < 7$, so the range is not to the right of all nodes.
 - b) It is contained by someValue's range, so now $R = \{v_{\text{someValue}}\}$.
 - c) We query the left subtree $\lambda_{\text{someValue}} = \text{Example}$.
 - i. $b_l^q < l_{\text{Example}}$ because $1 < 7$, so the range is not to the right of all nodes.
 - ii. It is contained by Example's range, but $r_{\text{Example}} \gtrsim r_{\text{someValue}}$, so someValue is still the tightest fit and R stays as it is.
 - iii. There are no subtrees.

figures/example_kd.pdf

Figure 3.2: An augmented interval tree using a k -d tree, representing the elements in listing 3.1.

- d) There is no right subtree.
- 4. $e_l^q \not\geq b_l^r$ because $1 < 3$, so there is no need to check the right subtree.
- 5. Our final result is $R = \{v_{\text{someValue}}\}$, which is intuitively the right answer.

TODO: Convert to TikZ diagram and use colors from listing!

This new and improved `FINDNODESBYLOCATION` implementation in algorithm 2 has a runtime that is in both $\Omega(1)$ and $O(k + \log |V|)$, where k is the number of ranges that a node is contained in—however, its worst-case runtime is still in $O(|V|)$, because in the worst case, the queried range is contained in every element in the tree, meaning $k = |V|$. Still, this case is almost impossible in real-world generated graphs. In actuality, a given range is only going to be contained by very few elements compared to the number of total elements within a given document tree. Taking into account that in almost all situations, $k \ll |V|$, and especially $k < \log |V|$, our average-case runtime reduces to an upper bound of $O(\log |V|)$. Hence, while the theoretical worst-case runtime complexity stays the same, the average-case runtime of algorithm 1 reduces from $\Theta(|V| + n_d \cdot |V| + |E| \cdot |V|)$ to $O(|V| + n_d \cdot \log |V| + |E| \cdot \log |V|)$. Since it is easy to make mistakes here due to the inherent complexity of this part, unit tests have also been implemented to make sure both normal and edge cases are handled correctly.

3.2.3 Usage in Practice

Now, we will have a quick look at how to actually use the algorithm in SEE. As explained before in section 2.1.3, the LSP importer—referring to the component responsible for using

algorithm 1 to generate a code city—is implemented as a graph provider. The UI of that graph provider in the Unity Editor is shown in fig. 3.3.

In that figure, we can see a configuration for the `dcaf-rs` project¹³, where the Language Server is the Rust Analyzer. At the top, we can see that the path to the source project has been set, and that the *Rust Analyzer*¹⁴ has been selected as a Language Server. Below that, there is a field for “Source Paths”. This refers to the directories that shall be scanned for relevant documents, whereas a document’s relevance is determined by its file extension. This is distinct from the “Project Path” as that is instead used by the Language Server to, for example, scan for project configuration files that describe dependencies. Conversely, “Excluded Source Paths” can be used for those paths within the configured source paths that should be ignored (e. g., generated files, tests).

The selectable Language Servers are only those that I have explicitly tested and confirmed to work with the algorithm—many servers are unusable, for example, due to required capabilities missing (such as the document symbols one). All in all, there are working¹⁵ Language Servers for the programming languages C, C++, C#, Dart, Go, Haskell, Java, JavaScript, Kotlin, Lua, MATLAB, PHP, Python, Ruby, Rust, TypeScript, and Zig. There are also functioning Language Server configurations for the miscellaneous languages JSON, \LaTeX , Markdown, and XML. A code city of a markup language like \LaTeX , for example, would have nodes for each section, figure, and so on, with edges representing references between these elements.

Next down the list in fig. 3.3, we have the import settings, which can be used to further refine the LSP-generated data that is used for the project graph. Specifically, we can select what kinds of nodes and edges we want to import, and which diagnostics we want to include. Additionally, for the edges, there is the option to exclude loops and edges (excluding `partOf`) to a node’s direct parent. The purpose behind this option is that otherwise, there are going to be a lot of definition/declaration edges from elements to their immediate parents, which happens because the locations returned by these LSP functions often extend slightly outside the node’s actual locations. For example, a quirk of many Language Servers is that a returned location may begin at `int value;`, while the actual node for this variable starts at `int value;`—in that case, the returned location would instead resolve to the outer container, such as the function it is contained in.

Finally, there is an option to enable the LSP functions for code windows that are detailed in section 3.4, a setting to generate log files for the transferred JSON-RPC messages, and a slider to adjust the maximum time we should wait for a Language Server’s response to

¹³This is one of the sample projects we use in section 3.5’s technical evaluation.

¹⁴The menu is grouped by language when opening it, hence the `Rust/` in front in fig. 3.3.

¹⁵There may well be more, I could not test some of the available Language Servers either due to setup problems, or (like in the case of Wolfram Mathematica) because I do not have a license for the language in question.



Figure 3.3: Unity Editor UI for the LSP graph provider

a request. With all of these options configured, the graph can be loaded and drawn (and optionally exported to a GXL file), where a bar displays the approximate progress of the process. At this point, we will once again refer back to the explanatory video **TODO**, which goes over the implemented functionality in a bit more detail and may be easier to follow along than this textual description.

TODO: Link to video again here

3.3 Integrating LSP Functionality into Code Cities

TODO!

3.4 Integrating LSP Functionality into Code Windows

TODO!

3.5 Technical Evaluation

TODO!

3.6 Interim Conclusion

TODO!

Table 3.1: All submitted pull requests done as part of this thesis.

Summary	Pull Request URI	Δ LOC	
		+	-
Cleanup of LSP specification	microsoft/language-server-protocol#1886	475	453
Preparing SEE for LSP integration	uni-bremen-agst/SEE#687	282	2773
Introducing Source . Range	uni-bremen-agst/SEE#715	392	313
Generating code cities using LSP	uni-bremen-agst/SEE#727	3475	180
LSP functions in code cities	uni-bremen-agst/SEE#747	1139	432
LSP functions in code windows	uni-bremen-agst/SEE#751	4080	2024
Preparing SEE for user evaluation	uni-bremen-agst/SEE#772	541	150

Only C# line changes have been counted in SEE pull requests.
GitHub pull requests are specified in the format namespace/repository#PR_number.

4

User Study



ESPITE having already done a technical evaluation in section 3.5, it is usually also a good idea to compare new approaches with existing state-of-the-art tools in a user study. In our case, we want to compare SEE's LSP-generated code cities with the capabilities a normal LSP-enabled IDE offers. For the latter, the Microsoft-developed IDE VSCode is a good fit, given that the Language Server Protocol itself originates here (see section 2.2) and that it still has a deep integration with it.

First, we will outline the general aim of this study by going over some existing related research, explaining important aspects of VSCode, and then enumerating our hypotheses. Next, we will explain the details of the design of the study itself, before analyzing its results with the 20 participants in detail. Finally, we describe some relevant threats to validity.

4.1 Plan

Our main aim here is to answer our second research question that we defined in section 1.3:

Are code cities a suitable means to present LSP information to developers as compared to IDEs + tables (on the dimensions of speed, accuracy, and usability)?

To empirically evaluate this research question, we will devise a series of short software engineering related tasks. Participants then get randomly assigned to either use SEE (along with our implementation from chapter 3) or VSCode (with an active Language Server). However, evaluating the supported capabilities (see table 2.1) in this way turns out to be quite difficult—for example, how would one evaluate the *Hover* capability, let alone features like semantic tokens which are almost identically implemented across SEE and VSCode? For this reason, we will abstain from incorporating the code window-related capabilities

from section 3.4. Limiting ourselves, then, to the code city-related changes from section 3.3, we have:

1. *Diagnostics* being displayed as erosion icons above corresponding nodes.
 - This feature was essentially already evaluated in my bachelor's thesis, albeit with the Axivion Dashboard as a data source instead of LSP [Gal21; GKS22].
2. *Hover* details being displayed when the user hovers the mouse above a node.
 - Since this is used here almost identically as in code windows, it does not make much sense to compare it against VSCode.
3. *Go to location*, *references*, and *call/type hierarchy* being used for rendered edges and context menu actions.
 - The context menu actions are not interesting to evaluate for the same reasons as above, though this does not apply to the generated edges.

It appears that the only capability reasonably evaluable in a user study of this form are actually the ones used in the generation of the code city in section 3.2. Besides, the bulk of the implementation pertains to the generation of code cities, so it makes sense to focus on them here. As a result, the user study is now actually of a form (directly comparing code cities with IDEs) that has been researched in previous literature before, so let us take a look at that research first before designing our own study.

4.1.1 Existing Research

TODO!

4.1.2 VSCode

TODO!

4.1.3 Hypotheses

TODO!

4.2 Structure

TODO!

4.2.1 Questionnaire

TODO!

4.2.2 Tasks

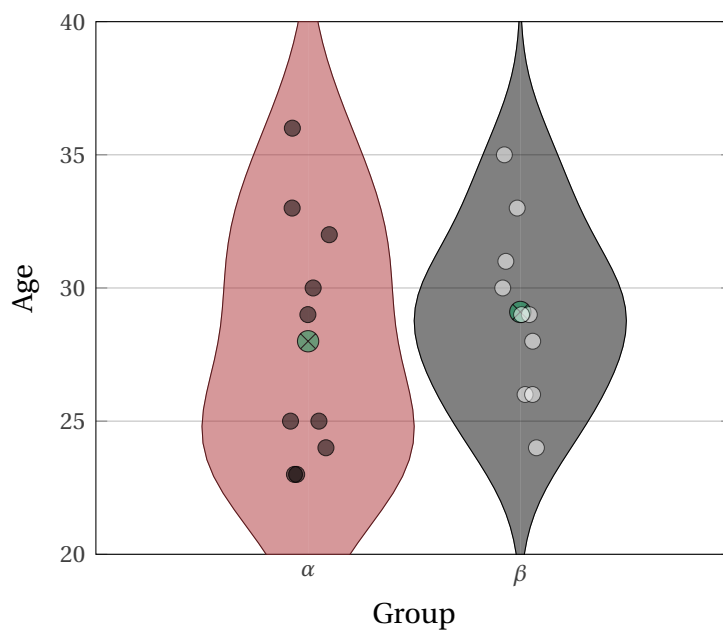
TODO!

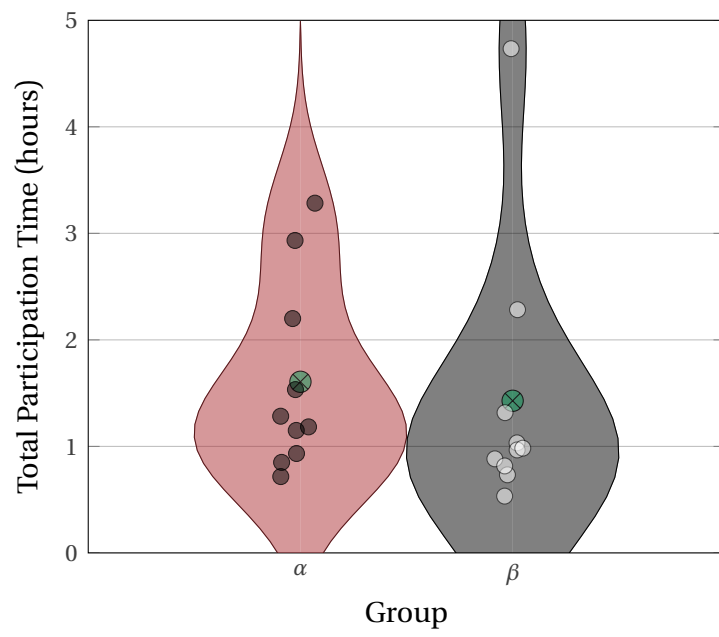
4.3 Results

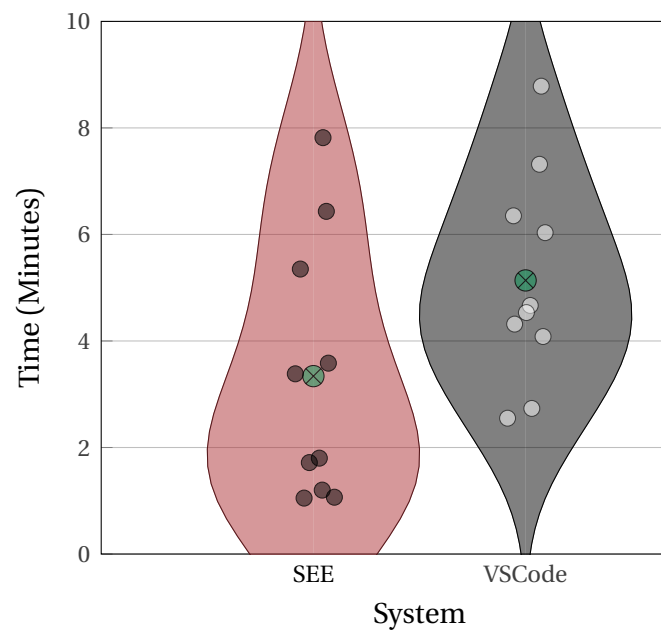
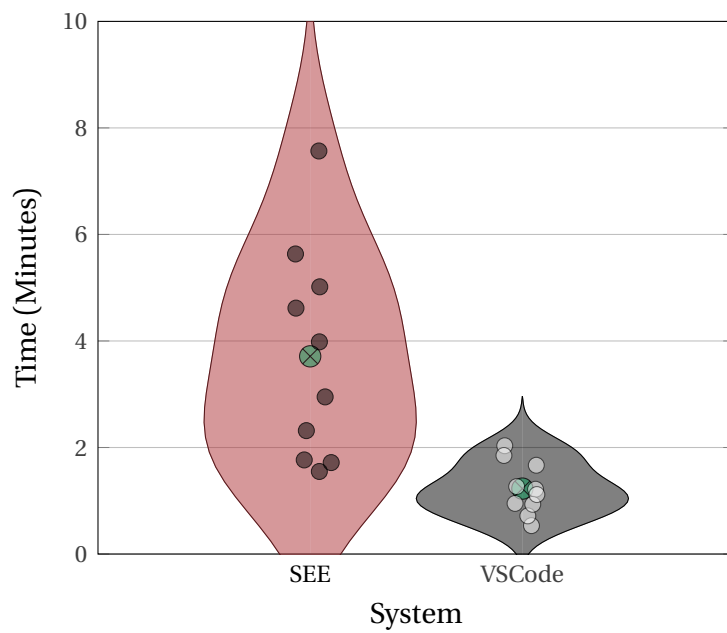
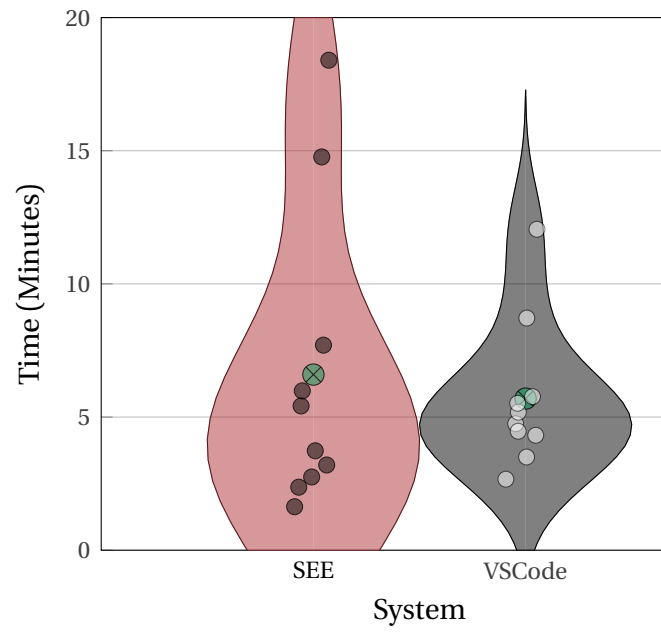
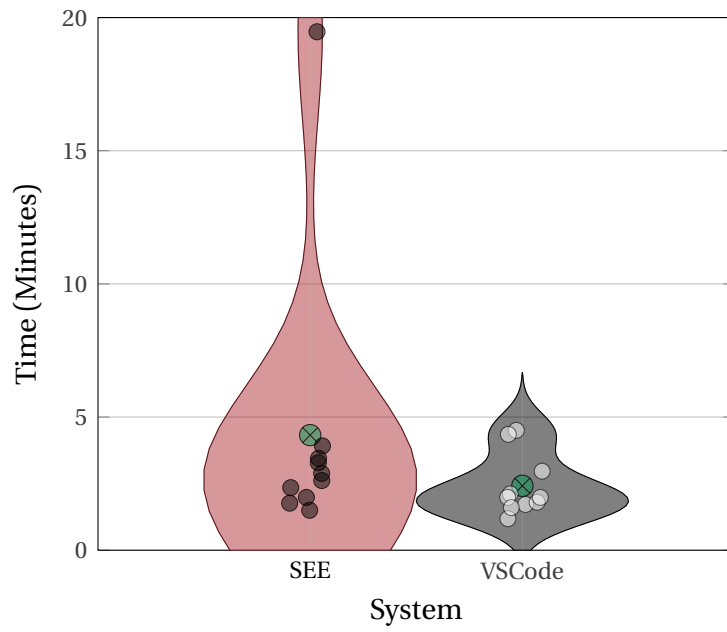
TODO!

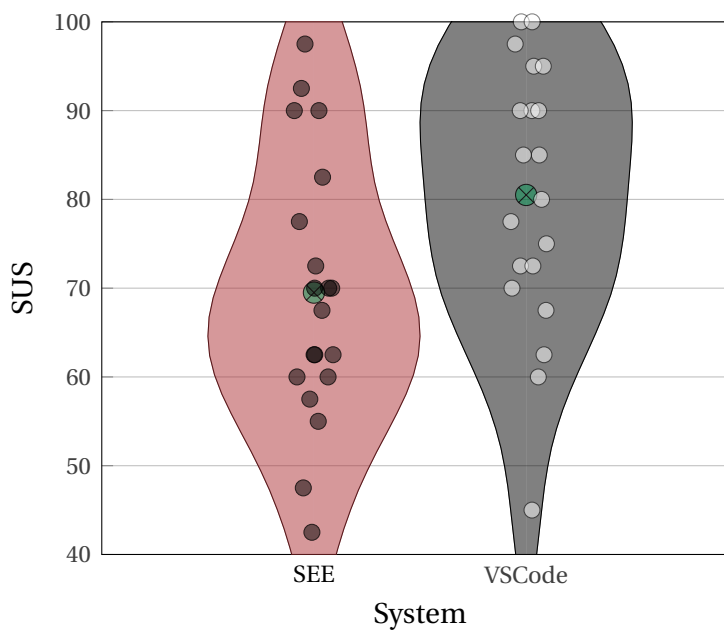
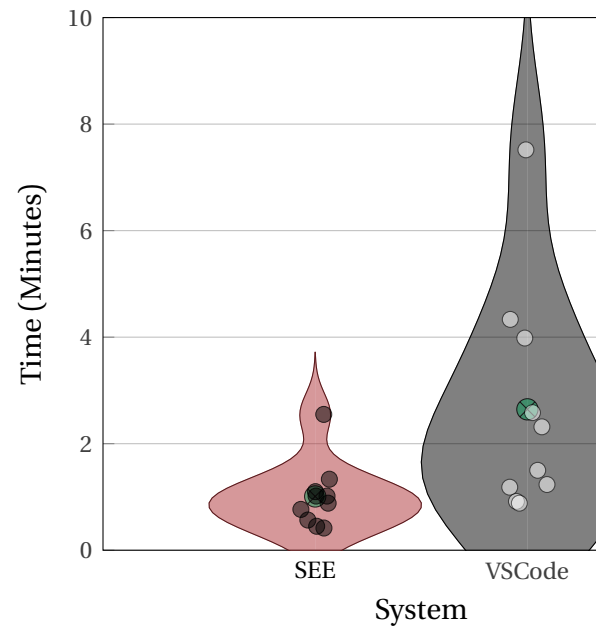
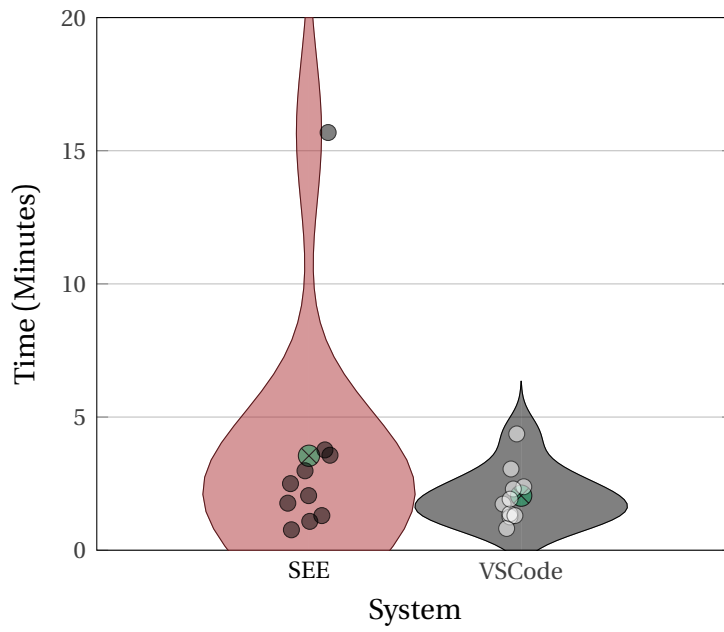
4.3.1 Demographics

TODO!









4.3.2 Correctness

TODO!

4.3.3 Time

TODO!

4.3.4 Usability

TODO!

4.3.5 Comments

TODO!

TODO: Also a section on the effect of experience (and others)

4.4 Threats to Validity

TODO!

4.5 Interim Conclusion

TODO!

5

Conclusion



PON... **TODO!**

5.1 Limitations

TODO!

5.2 Future Work

TODO!

5.3 The End

TODO!

TODO: Find
better title here



List of TODOs



ANY open tasks/notes/mistakes for this master's thesis are collected within this appendix. If you see any mistake or empty section not covered by such a note, please tell me. Note that this appendix will only appear in draft versions.

Make sure this is centered!	i
Unfinished section: ()	iii
Unfinished section: ()	v
Listings appear twice in TOC	viii
Convert diagrams to TikZ	1
Maybe rewrite. Could start by going over languages in general, giving spell check as example.	1
Use higher quality image.	4
Maybe remove the next part if already mentioned before.	8
In general, reword some stuff from “planning to do” ⇒ “done”	8
This is rather vague. Is that alright or do I need to operationalize here?	9
Back this up with sources	11
Refer back here later, noting how LSP can generate code cities for each of these use cases	12
is this correct? And put source here	12
Needs a 3D SEE screenshot, esp. for edges	12
add 2D screenshot	12
Source for 3D/game/VR being easier to navigate/understand	12

link the video	13
A lot of “represents” here	13
Give examples of software projects here to illustrate graph usage?	14
Show screenshot of context menu (pre-LSP)	14
Provide screenshot of editor	14
Show screenshot of erosion icons	15
More precise reference	18
Inconsistent use of should/will in this section	19
When?	19
Replace diagram sketch with TikZ picture.	27
Convert all software mentions to biblatex @software type.	34
Put in examples or more concrete estimation of how much more edges there are than nodes.	35
Do actual comparison between two ways of finding target nodes. Also mention example times from eval section at end.	35
Rewrite “the above” to something else, may not be above in printed version . . .	35
So how much does it help in practice? Reference tech eval here	36
Convert to TikZ diagram and use colors from listing!	38
Link to video again here	41
Unfinished section: Integrating LSP Functionality into Code Cities (3.3)	41
Unfinished section: Integrating LSP Functionality into Code Windows (3.4) . . .	41
Unfinished section: Technical Evaluation (3.5)	41
Unfinished section: Interim Conclusion (3.6)	41
Unfinished section: Existing Research (4.1.1)	44
Unfinished section: VSCode (4.1.2)	44
Unfinished section: Hypotheses (4.1.3)	44
Unfinished section: Structure (4.2)	45

Unfinished section: Questionnaire (4.2.1)	45
Unfinished section: Tasks (4.2.2)	45
Unfinished section: Results (4.3)	45
Unfinished section: Demographics (4.3.1)	45
Unfinished section: Correctness (4.3.2)	48
Unfinished section: Time (4.3.3)	49
Unfinished section: Usability (4.3.4)	49
Also a section on the effect of experience (and others)	49
Unfinished section: Comments (4.3.5)	49
Unfinished section: Threats to Validity (4.4)	49
Unfinished section: Interim Conclusion (4.5)	49
Unfinished section: Conclusion (5)	51
Unfinished section: Limitations (5.1)	51
Unfinished section: Future Work (5.2)	51
Unfinished section: The End (5.3)	51
Find better title here	51
Use higher quality image.	63
Bibliography doesn't work	70



Glossary

This document is incomplete. The external file associated with the glossary ‘main’ (which should be called `thesis.gls`) hasn’t been created.

Check the contents of the file `thesis.gls`. If it’s empty, that means you haven’t indexed any of your entries in this glossary (using commands like `\gls` or `\glsadd`) so this list can’t be generated. If the file isn’t empty, the document build process hasn’t been completed.

If you don’t want this glossary, add `nomain` to your package option list when you load `glossaries-extra.sty`. For example:

```
\usepackage[nomain]{glossaries-extra}
```

Try one of the following:

- Add `automake` to your package option list when you load `glossaries-extra.sty`. For example:

```
\usepackage[automake]{glossaries-extra}
```

- Run the external (Lua) application:

```
makeglossaries-lite.lua "thesis"
```

- Run the external (Perl) application:

```
makeglossaries "thesis"
```

Then rerun \LaTeX on this document.

This message will be removed once the problem has been fixed.



Acronyms

This document is incomplete. The external file associated with the glossary ‘abbreviations’ (which should be called `thesis.gls-abr`) hasn’t been created.

Check the contents of the file `thesis.gls-abr`. If it’s empty, that means you haven’t indexed any of your entries in this glossary (using commands like `\gls` or `\glsadd`) so this list can’t be generated. If the file isn’t empty, the document build process hasn’t been completed.

Try one of the following:

- Add `automake` to your package option list when you load `glossaries-extra.sty`.
For example:

```
\usepackage[automake]{glossaries-extra}
```

- Run the external (Lua) application:

```
makeglossaries-lite.lua "thesis"
```

- Run the external (Perl) application:

```
makeglossaries "thesis"
```

Then rerun \LaTeX on this document.

This message will be removed once the problem has been fixed.



Attached Files

Notation	Description
SEE.zip	The build of SEE used for the evaluation.



List of Figures

1.1	A Code City visualized in SEE.	3
1.2	An example of what edges can look like in SEE. TODO: Use higher quality image.	4
1.3	An illustration of how LSP can help simplify IDE development.	5
1.4	An example of the <code>texlab</code> language server running in <code>NEOVIM</code>	6
1.5	A simplified illustration of how SEE currently gains and uses information about software projects, and how the integration of LSP would change that. .	7
1.6	A code window.	8
2.1	An example of inlay hints in the JetBrains IntelliJ IDE. (From https://www.jetbrains.com/help/idea/inlay-hints.html (<i>last access: 2024-10-04</i>)) .	23
3.1	A high-level overview of the basic steps of the algorithm.	28
3.2	An augmented interval tree using a k -d tree, representing the elements in listing 3.1.	38
3.3	Unity Editor UI for the LSP Graph provider	40



List of Tables

2.1	LSP capabilities that will be integrated into SEE as part of this thesis.	24
3.1	All submitted pull requests done as part of this thesis.	41



List of Listings

2.1.	Example specification of request and response objects for the Hover capability.	17
3.1.	Example C# source code with demarcated symbol ranges.	37



Bibliography

- [Aho+07] Alfred V. Aho, Monica S. Lam, Ravi Sethi, and Jeffrey D. Ullman, eds. *Compilers: Principles, Techniques, & Tools*. 2. ed., Pearson internat. ed. Boston Munich: Pearson Addison-Wesley, 2007. 1009 pp. (cit. on p. 26).
- [Ben75] Jon Louis Bentley. “Multidimensional Binary Search Trees Used for Associative Searching”. In: *Communications of the ACM* 18.9 (Sept. 1975), pp. 509–517 (cit. on p. 35).
- [Ble22] Moritz Blecker. “Developing a Multi-user Source Code Editor for SEE”. BA thesis. University of Bremen, Jan. 31, 2022 (cit. on pp. 8, 26).
- [Cor+22] Thomas H. Cormen, Charles Eric Leiserson, Ronald Linn Rivest, and Clifford Stein. *Introduction to Algorithms*. Fourth edition. Cambridge, Massachusetts London: The MIT Press, 2022. 1 p. (cit. on p. 35).
- [Cro06] Douglas Crockford. *The Application/Json Media Type for JavaScript Object Notation (JSON)*. Request for Comments RFC 4627. Internet Engineering Task Force, July 2006. 10 pp. (cit. on p. 16).
- [DapSpec] *Debug Adapter Protocol*. Microsoft, May 2022 (cit. on p. 22).
- [FB19] Martin Fowler and Kent Beck. *Refactoring: Improving the Design of Existing Code*. Second edition. The Addison-Wesley Signature Series. Boston Columbus New York San Francisco Amsterdam Cape Town Dubai London Munich: Addison-Wesley, 2019. 418 pp. (cit. on p. 12).
- [Gal21] Falko Galperin. “Visualizing Code Smells in Code Cities”. BA thesis. University of Bremen, Sept. 27, 2021 (cit. on pp. 11, 15, 19, 24, 44).
- [GKS22] Falko Galperin, Rainer Koschke, and Marcel Steinbeck. “Visualizing Code Smells: Tables or Code Cities? A Controlled Experiment”. In: *2022 Working Conference on Software Visualization (VISSOFT)*. 2022 Working Conference on Software Visualization (VISSOFT). Limassol, Cyprus: IEEE, Oct. 2022, pp. 51–62 (cit. on pp. 15, 44).

- [JSO13] JSON-RPC Working Group. *JSON-RPC 2.0 Specification*. Jan. 4, 2013 (cit. on p. 16).
- [LsifSpec] *Language Server Index Format*. Microsoft, Sept. 2021 (cit. on p. 22).
- [LspSpec] *Language Server Protocol*. Microsoft, Oct. 2022 (cit. on pp. 4, 17, 19, 20).
- [McC76] T.J. McCabe. “A Complexity Measure”. In: *IEEE Transactions on Software Engineering* SE-2 (4 Dec. 1976), pp. 308–320 (cit. on p. 12).
- [Roh22] Ferdinand Rohlfing. “Debug Adapter Protocol in SEE”. BA thesis. University of Bremen, Aug. 2022 (cit. on p. 22).
- [Wag23] Bill Wagner. *The Task Asynchronous Programming (TAP) Model with Async and Await’ - C#*. Feb. 13, 2023. URL: <https://learn.microsoft.com/en-us/dotnet/csharp/asynchronous-programming/task-asynchronous-programming-model> (visited on Oct. 25, 2024) (cit. on p. 30).

TODO:
Bibliography
doesn't work