

Karl Oskar Magnus Holm

GeoGPT: An LLM-Based Geographic Information System

Making GIS More Accessible by Unlocking Geospatial Insights Through Natural Language

Master Thesis in Computer Science and Geomatics, June 2024

Supervisor at NTNU: Hongchao Fan

External supervisors from Norkart: Alexander Salveson Nossum and Arild Nomeland

Department of Geomatics
Faculty of Engineering
Norwegian University of Science and Technology



KARTAN

Oppgåve med omfang som kan tilpassast både prosjekt og masteroppgåve

LLMs - GIS-analysens død

(kan justerast seinare)

BAKGRUNN

Nyere modeller for kunstig intelligens har demonstrert spesielt gode evner til å kunne lære av store mengder ustrukturert og semi-strukturert informasjon. ChatGPT fra OpenAI tok verden med storm – og chat-baserte systemer florerer. Kan chat-baserte modeller skapes for å hente ut GIS-data effektivt? Norkart har en stor dataplattform hvor brukere utvikler mot API'er som i stor grad har GIS/Geografiske data i bunn. GeoNorge er en stor datakatalog hvor brukere slår opp, eller søker kategorisert for å finne data. QGIS, Python, PostGIS, FME og andre verktøy brukes ofte til å gjennomføre GIS-analyser – hvor en GIS-analytiker/data-scientist gjennomfører dette.

«*Finn alle bygninger innenfor 100-meters-belte som er over 100 kvm og har brygger*»

Er dette mulig å få til med dagens tilgjengelige chat-modeller?

OPPGAVEBESKRIVELSE

Oppgaven har som hovedmål å undersøke hvordan nyere språkmodeller kan benyttes for å gjennomføre klassiske GIS-analyser ved å bruke standard GIS-teknologi som PostGIS/SQL og datakataloger (OGC API Records fks). Hva finnes av tilgjengelig chat-løsninger? Hvordan spesialtilpasses til GIS-anvendelser? Hvor presise kan en GIS-Chat bli?

Relevante delmål for oppgaven:

1. Kartlegge state-of-the-art
2. Utvikle proof-of-concepts
3. Analysere begrensninger og kvalitet

Oppgaven vil med fordel deles i prosjektoppgave og masteroppgave

- Prosjektoppgave
 - State-of-the-art: Ai-modeller og multi-modal maskinlæring
 - Innhente og utvikle datagrunnlag og API-tilgjengelighet
- Masteroppgave
 - Utvikle proof-of-concepts med tilgjengelige åpne modeller/teknologi
 - Gjennomføre eksperimenter for analyse av kvalitet

Detaljert oppgavebeskrivelse utvikles i samarbeid med studenten.

ADMINISTRATIVT/VEILEDNING

Ekstern veileder: (en eller flere)

Mathilde Ørstavik, Norkart

Rune Aasgaard, Norkart

Alexander Nossum, Norkart

Aktuelle vegleiarar og ansvarleg professor ve NTNU (den som har fagansvar nærmast oppgåva):

Terje Midtbø (GIS, kartografi, visualisering)

Hongchao Fan (3D modellering, fotogrammetri, laser)

Abstract

Recent technological developments within Natural Language Processing have lead to the development of powerful Large Language Models (LLMs) like those powering ChatGPT — a AI-based chat interface created by OpenAI. LLMs have shown to be versatile, and the overarching goal for this master’s thesis is to utilize the logical reasoning and code generation abilities of modern LLMs to develop a GIS application that the user interacts with through a chat interface using natural language. The application, named *GeoGPT*, features three different agent types which, as demonstrated in experiments, can perform common GIS analyses on data from OpenStreetMap (OSM) with little to no help from the user. By utilize function calling, the agents are able to order invocations of predefined functions/tools with parameters specified by the LLM. The agents differ in the set of tool that they have been provided, and in the way they access the data — which is identical across the agents. One agent has accesses the data through a PostGIS database, another through an OGC API Features endpoint, and the third by having direct access to the shapefiles stored locally in its environment. The results from the experiments show that the PostGIS agent solves most tasks correctly, solves tasks in the shortest time on average, and is cheapest in terms of token usage. Results from an experiment conducted to evaluate the importance of the initial prompt from the user show that a more detailed step-by-step prompt — resembling that of a person with great experience with GIS — significantly improves the chance of GeoGPT producing a successful response. Overall, the work done in this thesis shows that an LLM-based GIS like GeoGPT can improve productivity by solving common geospatial tasks from simple natural language prompts, but that GIS expertise is still required as tasks become more difficult.

Sammendrag

Nyere teknologiske utviklinger innenfor språkprosessering har ledet til utviklingen av kraftige store språkmodeller, for eksempel de som driver ChatGPT — en AI-basert chattegrensesnitt skapt av OpenAI. Store språkmodeller har vist seg å være allsidige, og denne masteroppgavens vil teste denne allsidigheten ved utnytte språkmodellenes evner til logisk resonnering og kodegenerering til å utvikle en GIS-applikasjon som brukeren interagerer med gjennom et chattegrensesnitt ved bruk av naturlig språk. Applikasjonen har fått navnet *GeoGPT* og inneholder tre forskjellige agenttyper som kan utføre vanlige GIS-analyser på data fra OpenStreetMap (OSM) med minimal hjelp fra brukeren. Ved å utnytte såkalt *function calling*, er agentene i stand til å bestille kall av forhåndsdefinerte funksjoner/verktøy med parametere spesifisert av språkmodellen selv. Agentene varierer i hvilket sett med verktøy de har blitt gitt, og i måten de får tilgang til dataene på — som er identisk på tvers av agentene. Den ene agenten har tilgang til datagrunnlaget gjennom en PostGIS-database, en annen gjennom et OGC API Features-endepunkt, og den tredje ved å ha direkte tilgang til shapefiler lagret lokalt i miljøet den kjører i. Resultat fra eksperimenter viser at PostGIS-agenten løser de fleste oppgavene korrekt, har kortest mediantid per oppgave, og er billigst hva gjelder bruk av *tokens*. Resultater fra et eksperiment utført for å evaluere betydningen av den innledende meldingen fra brukeren av GeoGPT viser at et mer detaljert trinn-for-trinn-melding — som ligner det fra en person med stor erfaring med GIS kunne skrevet — betydelig forbedrer sjansen for at GeoGPT produserer et vellykket svar. Samlet sett viser arbeidet i denne masteroppgaven at et språkmodellbasert GIS som GeoGPT kan øke produktivitet ved å løse vanlige GIS-oppgaver basert på “spørninger” formulert ved naturlig språk, men at GIS-ekspertise fortsatt er nødvendig ettersom oppgavene blir mer utfordrende.

Preface

Karl Oskar Magnus Holm
Trondheim, 23rd May 2024

Contents

Abstract	iii
Sammendrag	iv
Preface	v
Acronyms	ix
List of Figures	xii
List of Tables	xiii
List of Code Snippets	xiv
1. Introduction	1
1.1. Background and Motivation	1
1.2. Goals and Research Questions	1
1.3. Research Method	2
1.4. Contributions	2
1.5. Thesis Structure	2
2. Background Theory	4
2.1. Large Language Models	4
2.1.1. Tokens and the Context Window	4
2.1.2. Attention and the Transformer Architecture	5
2.1.3. State-of-the-Art Decoder-Only Models	6
The GPT Family	6
The Gemini Family	8
The Claude Family	8
Open-Source Alternatives	8
2.1.4. Prompt Engineering	9
2.1.5. Function Calling LLMs	10
2.2. LangChain	11
2.3. Geospatial Databases and Data Catalogues	12
2.3.1. PostGIS	12
2.3.2. OGC API Features	13
2.3.3. SpatioTemporal Asset Catalogs	14

Contents

3. Related Work	15
3.1. Using LLMs for Geospatial Purposes	15
3.2. Agent Patterns	16
4. Datasets	18
4.1. Data Sources	18
4.2. Data Access	20
5. Architecture	22
5.1. High-Level Application Architecture	22
5.1.1. LangChain Server	22
5.1.2. Redis for Conversations	25
5.1.3. PostGIS and OGC API Features	25
5.1.4. Web UI	26
5.2. Agent Architecture	27
5.2.1. LangGraph Agent Implementation	27
5.2.2. Tools	29
5.2.3. Prompt Templating	32
6. Experiments	35
6.1. Experimental Setup	35
6.1.1. GIS Benchmark	35
Outcome Evaluation	35
Cost and Duration	35
Repeatability	36
6.1.2. Prompt Quality Experiment	36
6.1.3. Hardware and Model Version	37
6.2. Experimental Results	37
6.2.1. GIS Benchmark — Results	38
Outcome Evaluation	38
Cost and Duration	38
Repeatability	43
Successful Responses	43
Unsuccessful Responses	48
6.2.2. Prompt Quality Experiment — Results	52
7. Discussion	56
7.1. GeoGPT’s Place in the Field of GIS	56
7.2. Why Does the SQL Agent Outperform the Others?	57
7.2.1. Likely Higher Prevalence of PostGIS Examples During Pre-Training	57
7.2.2. Limitations with OGC API Features	58
7.3. Where GeoGPT Struggles	59
7.3.1. Walking Into Dead Ends	59
7.3.2. Self-Verification	59

Contents

7.4. Multi-Agent Architectures	60
8. Future Work	63
8.0.1. Ability to Answer Questions with no Clear Answer	63
8.0.2. Comparing Different Models	63
8.0.3. Automated Data Access	63
9. Conclusion	65
Bibliography	67
Appendices	72
A. Experiments	73
A.1. Quantitative	73
A.2. Experience Levels	81

Acronyms

AI Artificial Intelligence.

API Application Programming Interface.

BFCL Berkeley Function-Calling Leaderboard.

CQL Common Query Language.

CRS Coordinate Reference System.

DAG Directed Acyclic Graph.

DBMS DataBase Management System.

GIS Geographic Information System.

GML Geography Markup Language.

GPT Generative Pre-trained Transformer.

GQA Grouped-Query Attention.

HTML HyperText Markup Language.

HTTP Hypertext Transfer Protocol.

JSON JavaScript Object Notation.

Llama Large Language Model Meta AI.

LLM Large Language Model.

MBPP Mostly Basic Python Programming.

NLP Natural Language Processing.

OGC Open Geospatial Consortium.

OSM OpenStreetMap.

Acronyms

PoSE Positional Skip-wisE.

PPO Proximal Policy Optimization.

RAG Retrieval Augmented Generation.

RAM Random Access Memory.

REPL Read-Eval-Print Loop.

RLHF Reinforcement Learning from Human Feedback.

RNN Recurrent Neural Network.

RoPE Rotary Position Embedding.

SMoE Sparse Mixture-of-Experts.

SQL Structured Query Language.

SSD Solid-State-Disk.

STAC SpatioTemporal Asset Catalog.

SWA Sliding Window Attention.

UI User Interface.

List of Figures

2.1.	Tokenization example for an English sentence	5
2.2.	Tokenization of the word “revolution” with different suffixes	5
2.3.	How multiple choice tasks are inlined to work with fine-tuning of GPT models (Radford et al., 2018, p. 4)	7
2.4.	Collections, items, and features in OGC API Features specification. Retrieved from https://features.developer.ogc.org/ on April 29, 2024.	13
4.1.	A plot of four selected datasets constrained to a polygon of Trondheim . .	21
5.1.	Architecture overview	23
5.2.	Sequence diagram showing the information flow as the user loads and sends a message to GeoGPT	24
5.3.	Web UI	26
5.4.	Generic tool agent graph	28
5.5.	Example of a chat history	30
5.6.	Conversation showcasing the effect of giving an LLM a persona through the system message	33
5.7.	Chat template	34
6.1.	Outcome distribution between different agent types	38
6.2.	Task Durations for the Agent Types	39
6.3.	Cost and token usage	41
6.4.	Correlation matrices for metrics for the three agent types	42
6.5.	Successful response from GeoGPT’s SQL agent when asked how many counties the Glomma river runs through	44
6.6.	Partially successful response from GeoGPT’s Python agent when asked how many trees there are along Munkegata in Trondheim	46
6.7.	Unsuccessful attempt by GeoGPT’s Python agent to retrieve high-speed roads in Oslo	48
6.8.	Unsuccessful attempt by GeoGPT’s OGC API Features agent to create a geodesic line between Oslo Airport Gardermoen and Bergen Airport Flesland	50
6.9.	Outcome distribution for different levels of GIS experience	53
6.10.	Comparison between novice- and expert-level prompting of GeoGPT’s OGC API Features agent for calculation of the number of trees along Munkegata in Trondheim	55

List of Figures

7.1. ChatGPT multi-modal GPT-4 correctly identifying that the map layer intended to show high-speed roads in Oslo in fact “extends far beyond Oslo”	61
7.2. Illustration of how an agent supervisor takes in a user message and selects which sub-agent is to solve the next sub-task	62

List of Tables

4.1. Datasets used in experiments	18
5.1. Summary of Server Endpoints	23
5.2. Overview of the agent types and their tools	29
6.1. Description of Success	36
6.2. Encoding for Test Outcome	37
6.3. Standard Deviation by Agent Type	43
A.1. Questions for quantitative tests	73
A.2. Test results for quantitative tests	75
A.3. Questions for experience level tests	81
A.4. Test results for experience level tests	83

List of Code Snippets

2.1.	Example of a tool definition	10
2.2.	PostGIS example code for retrieving building outlines within a specified bounding box	13
2.3.	CQL examples	14
5.1.	Conversion from CQL to SQL	25
5.2.	Prevalences of common values for the <i>fclass</i> attribute in the <i>osm_natural_- points</i> collection	31
6.1.	GeoGPT-generated Python code aimed at computing the difference be- tween the Oslo outline and water features within it	40
6.2.	SQL code generated by GeoGPT to retrieve the counties that the Glomma river runs through	45
6.3.	Python code to exclude roads named Munkegata <i>outside</i> of Trondheim . .	47
6.4.	Python code to calculate the number of trees along Munkegata in Trondheim	47
6.5.	Python code that saves high-speed roads with a misleading filename . . .	49
6.6.	Tool invocation of <code>query_collection</code> with no results	51
6.7.	“Desperate” Python code that picks the first feature of the downloaded collection and <i>hopes</i> that they correspond to Gardermoen and Flesland .	52
7.1.	Multi-collection CQL query using the <i>Search</i> extension	58

1. Introduction

The introductory chapter will explain the motivation behind the thesis, its goals, and the research questions it will attempt to answer. Section 1.4 will list the main contributions of the thesis, and section 1.5 will give a high-level overview over the thesis.

1.1. Background and Motivation

The release of OpenAI's ChatGPT in November, 2022 (OpenAI, 2022) generated a hype within the general population and chat-based systems are now flourishing. Furthermore, significant advancements have been made within code generation, which makes LLMs useful even for technical tasks, enabling individuals with little to no prior programming experience to carry out computational tasks that require the code execution.

Geographic Information System (GIS) analysis has traditionally been reserved for GIS experts. Furthermore, GIS professional are commonly required to know their way around one or more GISs, and to be proficient in programming languages suitable to data science task, such as Python or R. Extensive domain knowledge is often also necessary when tackling GIS tasks, like knowing which data to use for a particular tasks and where to find them. All of these points, and more, are barriers to entry for people that wish to make use of powerful GIS tools for their particular purposes, but lack the technical know-how required to use them correctly. This challenge serves as the overall motivation behind this master's thesis, which will mitigate these issues by utilizing the vast background knowledge and code generation abilities of modern Large Language Models (LLMs).

1.2. Goals and Research Questions

Deriving from the motivation described in the section above, the overarching goal of this master's thesis is to investigate the possibilities of utilizing LLMs to create a natural language interface with a system that is capable of solving GIS-related tasks. The thesis' hypothesis is that modern LLMs are embedded with an understanding of common GIS workflows, and that their code generation abilities are now of such a level that they can solve a variety of such tasks.

Based on the overarching goal, three research questions have been constructed and are listed below:

1. Can an LLM-based system perform common GIS tasks?
2. What are core challenges in developing LLM-based GISs?

1. Introduction

3. Can an autonomous LLM-based GIS agent replace GIS professionals?

1.3. Research Method

Prior to this master's thesis, a specialization project on the same topic was conducted. The specialization project — as detailed in (Holm, 2023) — was of a theoretical character predominantly served as a literature study leading up to the master's thesis. The research questions listed in the above section, however, call for a different and more practical approach. Therefore, this master's thesis will revolve around the development of a “proof of concept”. The usefulnesss of the “proof of concept” will be evaluated through a series of tests.

Is this section necessary?

Should I mention the specialization project at all

1.4. Contributions

Below is a brief description of the contributions of this master's thesis:

1. A chat-based GIS named *GeoGPT*, powered by LLMs, that can solve tasks commonly solved using GIS software.
2. A new benchmark that will give insight into the ability of a system like GeoGPT to solve common GIS tasks when provided only with a natural language problem formulation.
3. Experimental results that highlight the importance of the initial problem formulation from the user, suggesting that GIS experience will not become redundant in the foreseeable future.

1.5. Thesis Structure

Below is an outline of the thesis' structure:

- Chapter 2 introduces the theory and tools necessary for the reader to be familiar with in order to understand the rest of the work.
- Chapter 3 gives insight into the work that has been done in regard to autonomous, LLM-based systems, in particular those within the field of geomatics.
- Chapter 4 presents the data provided for GeoGPT in the experiments.
- Chapter 5 will lay out the GeoGPT's architecture, providing both a high-level overview and details on important parts of the system.
- Chapter 6 presents the experimental setup, the datasets utilized, and the results and evaluations obtained from these experiments.

1. Introduction

- Chapter 7 will elaborate upon points of discussion that arise from the experimental results.
- Chapter 8 will suggest areas of improvement that are suitable for future research.
- Chapter 9 will conclude the master's thesis, reiterating the main contributions of the work.

2. Background Theory

NB! Parts of the Background Theory chapter is reused material from the specialization project (Holm, 2023) that preceded this master's thesis. Below are the sections in question, accompanied with descriptions of the extent to which, and how, the materials is reused:

- *Subsection 2.1.2: Reused, with minor adjustments.*
- *Subsection 2.1.3: GPT part reused without modifications.*

Chapter 2 will lay a theoretical basis for the work done in this master thesis, providing the user with the necessary knowledge to understand the contributions of the work. Section 2.1 will focus on Large Language Model (LLM). First, the section will introduce the terms “tokens” and “context window”. Thereafter, subsection 2.1.2 will present the component that most modern LLMs are based upon — namely the Transformer — and the attention mechanism that drives it. Continuing, subsection 2.1.3 section will present some of the leading LLMs as of 23rd May 2024, both proprietary and open-source ones. Subsection 2.1.4 will present the concept of *prompt engineering* before subsection 2.1.5 concludes section 2.1. Section 2.2 will give a brief intro to LangChain, a Python/JavaScript library that is used extensively throughout the code base of GeoGPT. Section 2.3 will conclude the Background Theory chapter, presenting the geospatial technologies that is used within GeoGPT.

2.1. Large Language Models

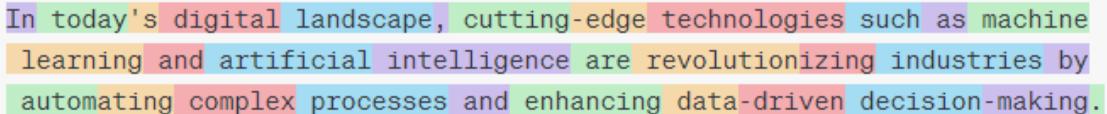
This section will lay the theoretical groundwork required to gain an overview of the inner workings of LLMs. LLMs are a type of neural networks that excel at language processing. They can be developed for different Natural Language Processing (NLP) tasks, such as text classification, masked language modelling, and text generation. While they all have their use cases, only text generation will be relevant for this thesis. Generative LLMs are designed to generate some output sequence based upon some input sequence, and are the types of models behind technologies like OpenAI’s *ChatGPT* (OpenAI, 2022). This input sequence is also known as the “context window”, and all sequences are made up from “tokens”, as subsection 2.1.1 will explain.

2.1.1. Tokens and the Context Window

While humans understand sentences as sequences of words, LLMs perceive them as sequences of **tokens**. An LLM possesses a fixed set of unique tokens in its vocabulary,

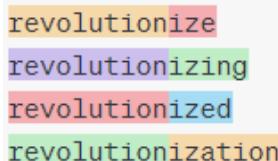
2. Background Theory

from which it constructs words and sentences. Tokens can be entire words, short character sequences, or single characters. Figure 2.1 illustrates how the OpenAI’s GPT-3.5 and GPT-4 models *tokenize* an English sentence. Notice how some words are deconstructed into more than one token. For example, the word “revolutionizing” is split into its root, “revolution”, and its suffix, “izing”. Figure 2.2 shows how this process applies to other suffixes as well. This way, the model only has to learn the meaning of “revolution”, and can append different suffixes to modify its function within a sentence, as opposed to learning an entirely new token for each version of the word.



In today's digital landscape, cutting-edge technologies such as machine learning and artificial intelligence are revolutionizing industries by automating complex processes and enhancing data-driven decision-making.

Figure 2.1.: Tokenization example for an English sentence



revolutionize
revolutionizing
revolutionized
revolutionization

Figure 2.2.: Tokenization of the word “revolution” with different suffixes

The **context window** of an LLM is the range of tokens that an LLM is able to process. When an LLM generates text, it does so by generating a new token based on the tokens it sees in the span of the context window. A larger context window will allow an LLM to take more and longer documents as context, and generate answers based on these. Leveraging sophisticated techniques like Rotary Position Embedding (RoPE) (Su et al., 2024) and Positional Skip-wisE (PoSE) (Zhu et al., 2024), researchers have been able to efficiently extend the context window of open-source models like, for instance, Meta AI’s Llama models.

2.1.2. Attention and the Transformer Architecture

Vaswani et al. (2017) managed to achieve new state-of-the-art results for machine translation tasks with their introduction of the Transformer architecture. The Transformer has later been proved effective for numerous downstream tasks, and for a variety of modalities. Titling their paper *Attention Is All You Need*, Vaswani et al. suggest that their attention-based Transformer architecture renders network architectures like Recurrent Neural Networks (RNNs) redundant, due to its superior parallelization abilities and the shorter path between combinations of position input and output sequences, making it easier for the model to learn long-range dependencies (Vaswani et al., 2017, p. 6).

The Transformer employs self-attention, which enables the model to draw connections between arbitrary parts of a given sequence, bypassing the long-range dependency issue

2. Background Theory

commonly found with RNNs. An attention function maps a query and a set of key-value pairs to an output, calculating the compatibility between a query and a corresponding key (Vaswani et al., 2017, p. 3). Looking at Vaswani et al.’s proposed attention function (2.1), we observe that it takes the dot product between the query Q and the keys K , where Q is the token that we want to compare all the keys to. Keys similar to Q will get a higher score, i.e., be *more attended to*. These differences in attention are further emphasized by applying the softmax function. The final matrix multiplication with the values V (the initial embeddings of the input tokens) will yield a new embedding in which all individual tokens have some context from all other tokens. The attention mechanism is improved by multiplying queries, keys, and values with *learned* weight matrices that are obtained through backpropagation. Self-attention is a special kind of attention in which queries, keys, and values are all the same sequence.

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V \quad (2.1)$$

Attention blocks can be found in three places in the Transformer architecture (Vaswani et al., 2017, p. 5). Below is an example using translation from Norwegian to German, showing how the attention blocks are used:

1. In the encoder block, to perform self-attention on the input sequence (which is in Norwegian)
2. In the decoder block, to perform self-attention on the output sequence (which is in German)
3. In the decoder block, to perform cross-attention (also known as encoder-decoder attention) where each position in the decoder attends to all positions in the encoder

The Transformer represented a breakthrough in the field of NLP, and is the fundamental building block of most modern LLMs.

2.1.3. State-of-the-Art Decoder-Only Large Language Models

While the work of Vaswani et al. is still considered perhaps the greatest breakthrough in NLP, most moderns LLM do not apply this exact encoder-decoder architecture. The evolution following the Transformer has favoured generative decoder-only models, focusing entirely on the generative component of the Transformer. The goal is to create models that can produce coherent and context-aware text.

The GPT Family

Generative Pre-trained Transformer (GPT) is an LLM that was introduced by OpenAI in 2018 (Radford et al., 2018). Specifically designed for text generation, a GPT is essentially a stack of Transformer *decoders*. It demonstrates through its vast pre-training on unlabelled data that such unsupervised training can help a language model learn

2. Background Theory

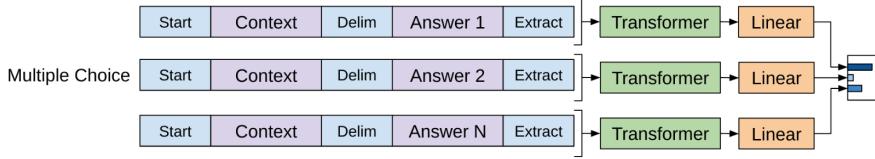


Figure 2.3.: How multiple choice tasks are inlined to work with fine-tuning of GPT models (Radford et al., 2018, p. 4)

good representations, providing a significant performance boost while alleviating the dependence on supervised learning. While the original Transformer architecture as described by Vaswani et al. (2017) was intended for machine translation — thus having encoders to learn the representation of the origin language representation of a given input sequence and decoders to learn the representation in the target language and perform cross-attention between the two — the GPT is designed only to *imitate* language. This is why there are no encoders to be found in the GPT architecture, only decoders. The model employs masked multi-head attention (running the input sequence through multiple attention heads in parallel), and is restricted to only see the last k tokens — with k being the size of the context window — and tasked to predict the next one.

Training consists of two stages: unsupervised pre-training and supervised fine-tuning. The former is used to find a good initialization point, essentially teaching the model to imitate the corpora upon which it is trained. This results in a model that will ramble on uncontrollably, simply attempting to elaborate upon the input sequence it is given, to the best of its abilities. This will naturally produce undefined behaviour, and it is therefore necessary to fine-tune the model on target tasks in a *supervised* manner. Radford et al. (2018, p. 4) explain how the model can be fine-tuned directly on tasks like text classification, but also how other tasks require the conversion of structured inputs into ordered sequences, as is the case when fine-tuning for tasks like multiple choice (see Figure 2.3). This *inlining* is required because the pre-trained model was trained on contiguous sequences of text. In the case of ChatGPT, OpenAI used Reinforcement Learning from Human Feedback (RLHF) (Christiano et al., 2023) by employing a three-step strategy: first training using a supervised policy, then using trained reward models to rank alternative completions produced by ChatGPT models, before fine-tuning the model using Proximal Policy Optimization (PPO) (Schulman et al., 2017), which is a way of training AI policies. This pipeline is then performed for several iterations until the model produces the desired behaviour (OpenAI, 2022).

OpenAI's API currently features three flagship models, which are all proprietary¹: the GPT-3.5 Turbo model, which is fast and inexpensive; the GPT-4 Turbo model, which as of 23rd May 2024 is described by OpenAI as their “previous high-intelligence model”; and lastly, GPT-4o, their “fastest and most affordable flagship model”.¹ The latter model has unique multi-modal abilities and can reason across audio, vision, and text in real

¹<https://openai.com/api/>

2. Background Theory

time, according to OpenAI.²

The Gemini Family

The suite of models known as Gemini is Google’s latest response OpenAI’s GPT models. The Gemini 1.0 suite (Gemini Team et al., 2024a), which is the first suite of Gemini models, includes three different models: Ultra, Pro, and Nano. These are listed in descending order in terms of size (number of model parameters). The models of the Gemini 1.0 suite are multi-modal, supporting text, image, audio, and video. Gemini 1.0 Ultra displayed new state-of-the-art performance on most major benchmarks, but performed significantly worse on the HellaSwag benchmark (which measures a model’s common-sense understanding) compared to the latest GPT-4 model at the time. The models scored 87.8% and 95.3%, respectively. Supporting a context length up to 1M tokens, Gemini 1.0 Ultra surpassed Claude 2.1’s context window of 200k (see the next paragraph for more on the Claude models) with a wide margin, and with the release of Gemini 1.5 Pro came also the possibility of utilizing a context window of up to 10M tokens, though in production this number is currently limited to 1M (Gemini Team et al., 2024b; Pichai and Hassabis, 2024). Furthermore, Gemini 1.5 Pro outperforms Gemini 1.0 Ultra in some capabilities despite using significantly less training compute (Gemini Team et al., 2024b, p. 31).

The Claude Family

Developed at Anthropic, Claude is the third major suite of proprietary LLM. Anthropic is one of the actors in the LLM market that has helped push in the direction of long-context LLMs, with their Claude 2 model (released in November 2023) being the first to support up to 200k tokens (Anthropic, 2023, p. 9). Latest in line is Claude 3, a family of language models that — much like the Gemini family — features three models of different sizes: Opus, Sonnet, and Haiku. Again, these are listed in descending order in terms of number of parameters. Opus offers the most advanced capabilities, and outperforms GPT-4 and Gemini 1.0 Ultra on most benchmarks (Anthropic, 2024, p. 6). Haiku is Anthropic’s fast and economical option, while Sonnet serves as a balance between Opus’ complexity and Haiku’s speed.

Open-Source Alternatives

As mentioned, OpenAI’s GPT models, Google’s Gemini models, and Anthropic’s Claude models are all proprietary. This prevents developers from downloading these models and making improvements and customizations to them through fine-tuning. This, and several other reasons, have led to the emergence of a number of *open-source* LLMs.

The **Llama** family of LLMs, developed at Meta AI, is perhaps the most famous open-source option to the proprietary LLMs. At the time of writing, the last in line is the Llama 3 model (Meta AI, 2024), which comes in two sizes: 8B and 70B parameters.

²<https://openai.com/index/hello-gpt-4o/>

2. Background Theory

Both models display state-of-the-art performance on most major benchmarks compared to comparable open-source alternatives, and the 70B model even surpasses proprietary models like Gemini 1.5 Pro and Claude 3 Sonnet on certain benchmarks.

Mistral AI is one of the most prominent actors in the world of open-source LLMs. Their debut model, Mistral 7B, outperformed Llama 2 13B (which was the best open-source LLM at the time) across all the benchmarks they evaluated (Jiang et al., 2023). Mistral AI has also gained fame for their Sparse Mixture-of-Experts (SMoE) architecture, which was introduced with the Mixtral 8x7B model (Jiang et al., 2024). It shares the same architecture as Mistral 7B, except that each layer of the model is composed of 8 feed-forward blocks. Using a router at each layer, Mixtral 8x7B is able to use only 13B out of a total of 47B parameters during inference, keeping cost and latency low.

Along with their Gemini models, Google released a family of open-source models called **Gemma**, which are based on the same research conducted for their Gemini models (Gemma Team et al., 2024). Gemma comes in two sizes: 2B and 7B parameters. At its release, the Gemma 7B surpassed Llama 2 13B and Mistral 7B LLMs in 11 out of 18 benchmarks. Note, however, that Llama 3 8B has improved upon its predecessor and now also performs better than Gemma 7B, overall.

2.1.4. Prompt Engineering

Prompt engineering refers to the process of constructing a query that will be used as input to an LLM. In a chat-based context like that of ChatGPT, the prompt generally consists of a series of messages. These messages can hold one of three different roles³, each listed below:

- **System:** Generally used to set the behaviour of the LLM assistant, giving the assistant a specific personality or information on how to answer the user.
- **User:** A message from the human user that the assistant should respond to.
- **Assistant:** A message generated by the AI/LLM.

White et al. (2023) stress the importance of prompt engineering to efficiently converse with LLMs. They provide a catalogue of prompt patterns that aim to help enforce certain qualities of the output generated by the LLM. These patterns are organized into six distinct categories (White et al., 2023, p. 4):

- **Input Semantics:** Clarifying what information is fed into the LLM and how this input should be used to generate responses.
- **Output Customization:** Strategies to guide how the LLM should format and structure its responses.
- **Error Identification:** Methods to identifying and resolving errors in the outputs produced by the LLM.

³<https://platform.openai.com/docs/guides/text-generation/chat-completions-api>

2. Background Theory

- **Prompt Improvement:** Techniques to improve the quality of both the input provided to the LLM and the output it generates.
- **Interaction:** Strategies for enhancing interactions between the user and the LLM.
- **Context Control:** Controlling the contextual domain within which the LLM operates.

Patterns that turned out useful to GeoGPT — the main contribution of this thesis — include the *Template* pattern (Output Customization), the *Reflection* pattern (Error Identification), and the *Infinite Generation* pattern (Interaction Strategy). The *Template* pattern allows the user or the system to define a template for the LLM to fill out. This is closely related to function calling, which is discussed in subsection 2.1.5. Also related to function calling is the *Reflection* pattern, which allows the LLM to inspect its own output in order to identify and correct errors. The *Infinite Generation* pattern lets the LLM generate output indefinitely without requiring the user to re-enter the conversation after each generated message. This pattern is important when developing agentic behaviours for the LLM.

2.1.5. Function Calling LLMs

Function calling — also known as *tool calling* — was first introduced by OpenAI in April 2023 (Eleti et al., 2023). Function calling allows developers to provide function definitions to an LLM and have the LLM output a JSON object containing the name of one or more of the functions provided, as well as suitable parameters to these functions. Code Snippet 2.1 shows a description of a function that performs a mathematical division. Function calling LLMs are able to detect when a certain function should be called based on the name, description, and input parameters specified in such JSON objects. Function calling makes it possible to give an LLM *hooks* into the real world, and provides a more reliable way for developers to integrate LLMs into applications.

```
1 {
2   "type": "function",
3   "function": {
4     "name": "divide",
5     "description": "Performs a mathematical division.",
6     "parameters": {
7       "type": "object",
8       "properties": {
9         "dividend": {
10           "description": "The number to be divided (numerator).",
11           "type": "number"
12         },
13         "divisor": {
14           "description": "The number by which the dividend is divided (denominator).",
15           "type": "number"
16         }
17       }
18     }
19   }
20 }
```

2. Background Theory

```
16      }
17      },
18      "required": ["dividend", "divisor"]
19    }
20  }
21 }
```

Code Snippet 2.1: Example of a tool definition

Function calling can be used to provide correct and up-to-date information that would otherwise require extensive training and fine-tuning. Having the LLM use function calling for information retrieval also make them more transparent, making it possible to trace a claim back to its source, something that is normally hard to do with LLMs. Another use case is code execution. Take this rather simple function signature:

```
execute_python_code(code: string) -> string
```

Such a function could take some Python code as a string parameter and return the standard output that results from executing that code. This is the principle behind what was previously known as ChatGPT’s “Code Interpreter” mode, where ChatGPT serves as a code executing agent that can generate, execute, and self-correct its own code. Similar functions could be constructed for SQL, making it possible for LLMs to work against relational databases. Furthermore, as Eleti et al. (2023) describes, function calling can also be used to extract structured data from text.

An initiative amongst researchers at Berkeley (Yan, Fanjia et al., 2024) lead to the creation of a benchmark that aims to evaluate the LLM’s ability to call functions and tools. The benchmark, named Berkeley Function-Calling Leaderboard (BFCL), includes four different test scenarios: *single function accuracy*, where the LLM is provided with a single function definition; *multiple function accuracy*, where 2 to 4 functions are passed and the model must select the appropriate function; *parallel function*, where the model needs to determine how many functions should be called; and *parallel multiple function*, a combination of parallel function and multiple function. Some models support different levels of function calling natively, while others have to be carefully prompted to help accommodate function calling abilities. At the time of writing (23rd May 2024), a prompted version of GPT-4-0125-Preview tops the leaderboard.⁴

2.2. LangChain

LangChain (LangChain AI, 2022) is an open-source project that provides tooling which simplifies the way developers interface with LLMs. As the name suggests, LangChain revolves around *chains*. Chains are Directed Acyclic Graphs (DAGs), which can be constructed from so-called *runnables*. The “Runnable” protocol is a standard interface, which makes enforces chains to implement **stream**, **invoke**, and **batch** methods, as well

⁴<https://gorilla.cs.berkeley.edu/leaderboard.html>

2. Background Theory

as their asynchronous counterparts. LangChain is supported by a large community of developers, and has components that simplify interaction with LLMs, like prompt templates, output parsers, toolkits for function calling purposes, and off-the-shelf agents.

Common use cases for LangChain are:

- Building chatbots for question answering that use semantic retrieval from document store
- Creating agents with access to external tools by leveraging function calling (see subsection 2.1.5)
- Creating code executing agents for Python, SQL, or other programming languages

In January 2024, LangChain AI rolled out a framework called LangGraph, which is built on top of the LangChain ecosystem. While *chains* are well-suited for DAG workflows, they are not easily adapted to cyclic graphs. LangGraph, however, is specifically designed to simplify development of cyclic graphs for LLM applications, which are commonly used to create agent-like behaviours (LangChain AI, 2024). A LangGraph *graph* is a set of nodes that pass around and modify a state dictionary. The nodes are connected by edges that define what node can succeed another node. Edges can be conditional, meaning the state produced from one node decides which edge is followed. This allows for complex logic and simplifies implementation of advanced agent patterns, some of which are discussed in section 3.2.

2.3. Geospatial Databases and Data Catalogues

This section will discuss the geospatial technologies that were used or considered for use in this master's thesis.

2.3.1. PostGIS

PostGIS (*PostGIS* 2001) is an open-source extension for the PostgreSQL DBMS. The PostGIS extension adds support for storing, indexing, and querying geospatial data. Data can be stored in both two and three dimensions, and they can have types like points, lines, polygons. Geospatial types can be stored along with a spatial index which can significantly reduce search time for their geometries. GiST (Generalized Search Tree)⁵ is commonly used in PostgreSQL/PostGIS to take advantage of various tree-based search algorithms that are developed to retrieve spatial features quickly.

PostGIS also comes with a plethora of spatial database functions that are used to analyse and process geospatial data. These function's names are prefixed with `ST_-` (short for “spatio-temporal”), and some examples are `ST_DWITHIN`, `ST_BUFFER`, and `ST_TRANSFORM`. Code Snippet 2.2 shows a typical query for retrieving building outlines within a bounding box.

⁵<https://en.wikipedia.org/wiki/GiST>

2. Background Theory

```
1 SELECT *
2 FROM osm_buildings_polygons
3 WHERE type = 'house'
4 AND ST_Intersects(geom, ST_MakeEnvelope(min_lon, min_lat, max_lon,
max_lat, 4326));
```

Code Snippet 2.2: PostGIS example code for retrieving building outlines within a specified bounding box

2.3.2. OGC API Features

OGC API Features is an API specification that defines modular API building blocks for interacting with features, which are real-world objects (Open Geospatial Consortium, 2022). This includes blocks for creating, modifying, and querying features on the Web. A typical implementation of OGC API Features implements these building blocks for HTML, GeoJSON, and GML. These are called *requirement classes*, though none of them are strictly required. The HTML requirement class gives the user of the API a visualization of the features, whereas the GeoJSON and GML requirements classes are typically meant for use in other applications.

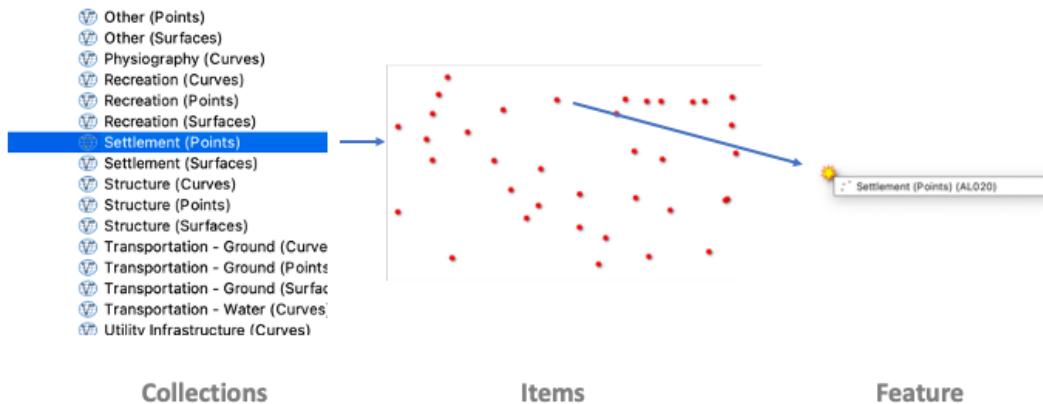


Figure 2.4.: Collections, items, and features in OGC API Features specification. Retrieved from <https://features.developer.ogc.org/> on April 29, 2024.

The development of the Features standard is divided into several parts that are meant to build on top of one another. Here are the five parts that are listed on OGC's websites:⁶

- Features - Part 1: Core⁷

⁶<https://ogcapi.ogc.org/features/>

⁷<https://docs.opengeospatial.org/is/17-069r4/17-069r4.html>

2. Background Theory

- Features - Part 2: Coordinate Reference Systems by Reference⁸
- Features - Part 3: Filtering⁹
- Features - Part 4: Create, Replace, Update and Delete¹⁰
- Features - Part 5: Schemas¹¹

Part 1 specifies core capabilities that are described in the first paragraph of this section, while parts 2-4 specify additional capabilities. Part 2 allows for retrieval of features in Coordinate Reference Systems (CRSs) different to the default WGS 84 reference system. Part 3 enables filtering of features using Common Query Language (CQL). CQL is a language similar to SQL. This allows for filtering of collections, so that users of the API can retrieve subsets of a given collection. Below are two examples of CQL queries:

```
1 \\\ Example 1
2 county in ('Akershus', 'Buskerud', 'Ostfold')
3
4 \\\ Example 2
5 DWITHIN(the_geom, Point(63.4265, 10.3960), 1, kilometers)
```

Code Snippet 2.3: CQL examples

Part 4 defines how an API that implements the specification should handle addition, replacement, modification, and removal of a collection's features. Part 5 describes how features can be described by a logical schema and how these schemas are published. In addition to parts 1 through 5, several other *proposed* extensions have emerged, such as the *Search* extension, which would allow for multi-collection queries, or the *Geometry Simplification* extension which proposes the use of simplification algorithms for retrieving simplified versions of a collection.¹²

2.3.3. SpatioTemporal Asset Catalog

The SpatioTemporal Asset Catalog (STAC) specification¹³ is closely related to OGC API Features, and Chris Holmes, former board member of Open Geospatial Consortium, stated in a blog post that “STAC API implements and extends the OGC API — Features standard, and our shared goal is for STAC API to become a full OGC standard” (Holmes, 2021). The main difference between OGC API Features and STAC is the latter’s requirement that all items/features should have a temporal component, thus making it *spatio-temporal*.

⁸<https://docs.ogc.org/is/18-058r1/18-058r1.html>

⁹<https://docs.ogc.org/DRAFTS/19-079r1.html>

¹⁰<https://docs.ogc.org/DRAFTS/20-002.html>

¹¹<https://docs.ogc.org/DRAFTS/23-058r1.html>

¹²<https://github.com/opengeospatial/ogcapi-features/tree/master/proposals>

¹³<https://stacspec.org/en>

Move
to Back-
ground?

3. Related Work

NB! Subsection 3.1 in this chapter reused material (with some minor modifications) from the specialization project (Holm, 2023) that preceded this master's thesis.

This chapter provides an overview of previous research that shares objectives similar to those of this thesis. The related work is divided into two main parts: section 3.1, which presents research investigating potential use cases of LLMs in the field of geospatial information technology; and section 3.2, which presents examples of three different types of patterns commonly used in LLM-based agents.

3.1. Using LLMs for Geospatial Purposes

Roberts et al. (2023) investigated the extent of GPT-4's geospatial awareness through a set of case studies with increasing difficulty, starting with general factual tasks and finishing with complex questions such as generating country outlines and travel networks. The authors found that GPT-4 is "skillful at solving a variety of application-centric tasks", almost having the ability to "see", despite being a language model and therefore only able to interface with the world through sequenced, textual input (note that multi-modal models were not particularly widespread when they wrote their paper). Examples include its ability to serve as a travel assistant in providing itinerary suggestions for a trip when provided with requirements, and its ability to provide generally correct start and end locations of bird migration paths. While it quickly became obvious that a lot of geospatial context have been embedded within the model during the vast pre-training, the question of whether this is memorization or reasoning is a central one. The authors suggest that the variability of tasks in their experiments deems it unlikely that it is all memorization, but they say that some things appear to be memorized.

Mooney et al. (2023) examined the performance of ChatGPT in a Geographic Information System (GIS) exam, aiming to assess its ability to grasp various geospatial concepts, highlighting its capabilities and limitations. Experiments were conducted on GPT-3.5 and GPT-4, which delivered performances equivalent to grades of D and B+, respectively. Additional experiments were conducted for more specialized areas of GIS, including True/False questions about spatial analysis, and simple tasks in applied GIS workflows. Experiments on the latter showed that GPT-4 was able to correctly answer a relatively complex GIS task involving seven different datasets, requiring seven steps in order to obtain a perfect score. Generally, GPT-4 outperformed GPT-3.5 in all tasks. While clearly powerful, the authors highlight a range of challenges, among which the multi-modal nature of GIS.

3. Related Work

Li and Ning (2023) state that “autonomous GIS will need to achieve five autonomous goals: self-generating, self-organizing, self-verifying, self-executing, and self-growing.”. They provide a “divide-and-conquer”-based method to address some of these goals. Furthermore, they propose a simple trial-and-error approach to address the self-verifying goal. They also highlight the need for a memory system in a mature LLM-based GIS system, referring to the use of vector databases in autonomous agents made with frameworks like AutoGPT (Richard, 2023). Even with its shortages, the solution that Li and Ning (2023) provide, called LLM-Geo, is able to produce good solutions in various case studies by providing executable assemblies in a Python environment when provided with URLs to relevant data sets, along with a user-specified query.

Zhang et al. (2023) use the LangChain framework in order to combine different GIS tools in a sequence to solve various sub-goals, focusing on using the semantic understanding and reasoning abilities of LLMs to call externally defined tools, employing the LLM as an agent or controller. The authors take great inspiration from the AutoGPT framework (Richard, 2023). The externally defined tools are described by their names and descriptions. These descriptions contain information about the input parameters and output types of the tools/functions. Tools are defined for geospatial data collection, data processing and analysis, and data visualization. The effectiveness of the system is showcased through case studies.

3.2. Agent Patterns

LLM-based agents can be implemented in many different ways, and researchers have developed a plethora of *agent patterns* that seek to improve upon areas where LLMs tend to be less effective. This section will shortly explain three different types of patterns that were considered during the development of GeoGPT.

The **multi-agent** pattern takes inspiration from human collaboration in that it is made up from multiple specialized agents that work together to achieve some objective. MetaGPT (Hong et al., 2023) is a LLM-based multi-agent system consisting of agents with human-level domain expertise. Using an assembly line paradigm, where the overall goal is divided into subtasks, Hong et al. showed that MetaGPT could generate more coherent solutions compared to the previous state-of-the-art multi-agent systems. At the time of release, MetaGPT set a new state-of-the-art performance on the HumanEval and MBPP benchmarks (Hong et al., 2023, p. 7), demonstrating the potential of the multi-agent pattern. AutoGen (Wu et al., 2023) and crewAI (Moura, 2024) are other examples of frameworks geared towards the multi-agent pattern.

Patterns that employ **self-reflection** are commonly used with autonomous LLM-based agents. *Reflextion* (Shinn et al., 2023) is a framework that reinforces agents through linguistic feedback, essentially allowing the agent to reflect upon the outcomes of its actions. The framework utilizes three distinct models: an *Actor* model responsible for generating text and actions; an *Evaluator* model which assesses the quality of the outputs from the Actor; and a *Self-Reflection* model that generates reinforcement cues for the Actor based on the output and quality assessment of the other two models. Together,

3. Related Work

these three models form a loop that will run until the Evaluator deems the output from the Actor as correct.

Step-by-step reasoning is another pattern that has proved to be efficient in helping LLMs produce correct responses. Wei et al. (2023) demonstrated that so-called *chain-of-thought prompting* can be used for enhancing reasoning in LLMs by providing it with examples of how to reason for tasks *similar* to that which it is trying to solve. By helping the LLM with decomposing multi-step problems into intermediate steps, Wei et al. managed to achieve state-of-the-art accuracy on the GSM8K benchmark of math word problems.

4. Datasets

Section 4.1 provides a description of the datasets used in the experiments. Furthermore, it was decided to explore different ways for GeoGPT to discover this data. Section 4.2 elaborates on this.

4.1. Data Sources

A total of eighteen datasets were used in the experiments, all downloaded from Geofabrik's website.¹ Geofabrik, which is German for "geo factory", is a company that "extract, select, and process free geodata". They have gathered data from OpenStreetMap and published them as a collection of shapefiles, dividing them into categories such as "places of worship", "points of interest", and "traffic". Data can be downloaded for different regions of the world, and for experiments conducted in this thesis, data for Norway was used. Table 4.1 in ?? lists all datasets, along with short descriptions of their contents. Common for all datasets are their *fclass* attribute, which is short for *feature class*. Some datasets have additional attributes, such as the *maxspeed* attribute in the road data and the *type* attribute in the building data. Figure 4.1 shows a plot containing four selected datasets, constrained to a bounding box of Trondheim. This plot was created by GeoGPT.

Table 4.1.: Datasets used in experiments

Dataset	Data Type	Description
Buildings	Polygon	Contains building outlines. Its <i>type</i> attribute can have values like <i>house</i> , <i>university</i> , and <i>restaurant</i> .
Land Use	Polygon	Represents areas designated to different purposes and activities. Its <i>fclass</i> attribute can have values like <i>forest</i> , <i>farmland</i> , and <i>residential</i> .

Continued on next page

¹<https://download.geofabrik.de/europe/norway.html>

4. Datasets

Table 4.1 continued from previous page

Dataset	Data Type	Description
Natural	Point	Contains outlines of various objects found in nature. Its <i>fclass</i> attribute can have values like <i>beach</i> , <i>glacier</i> , and <i>cave_entrance</i> .
Natural	Polygon	Similar to the point data equivalent.
Places of Worship	Point	Common values for <i>fclass</i> attribute: <i>christian</i> , <i>buddhist</i> , and <i>muslim</i> .
Places of Worship	Polygon	Similar to the point data equivalent.
Places	Point	Common values for <i>fclass</i> attribute: <i>farm</i> , <i>village</i> , and <i>island</i> . Repeated entries trimmed for brevity.
Places	Polygon	Similar to the point data equivalent.
Points of Interest	Point	Common values for <i>fclass</i> attribute: <i>tourist_info</i> , <i>bench</i> , and <i>kindergarten</i> .
Points of Interest	Polygon	Similar to the point data equivalent.
Railways	Lines	Common values for <i>fclass</i> attribute: <i>rail</i> , <i>subway</i> , and <i>tram</i> . Also has True/False attributes saying if a given line segment is a bridge or a tunnel.
Roads	Lines	Common values for <i>fclass</i> attribute: <i>rail</i> , <i>subway</i> , and <i>tram</i> . Has additional attributes <i>oneway</i> , <i>maxspeed</i> , <i>bridge</i> , and <i>tunnel</i> .
Traffic	Point	Common values for <i>fclass</i> attribute: <i>crossing</i> , <i>street_lamp</i> , and <i>parking</i> .
Traffic	Polygon	Common values for <i>fclass</i> attribute: <i>parking</i> , <i>pier</i> , and <i>dam</i> .
Transport	Point	Common values for <i>fclass</i> attribute: <i>bus_stop</i> , <i>ferry_terminal</i> , and <i>railway_station</i> .
Transport	Polygon	Similar to the point data equivalent.
Water	Polygon	Common values for <i>fclass</i> attribute: <i>water</i> , <i>wetland</i> , and <i>river_bank</i> .

Continued on next page

4. Datasets

Table 4.1 continued from previous page

Dataset	Data Type	Description
Waterways	Lines	Common values for <i>fclass</i> attribute: <i>stream</i> , <i>river</i> , and <i>canal</i> .

4.2. Data Access

In order to investigate how an LLM-based GIS agent most comfortably accesses geospatial data, a decision was made to include three different methods for data access in the experiments.

The first method for data access is to have the files described in section 4.1 remain untouched. The files would be stored locally on the computer on which the experiments are conducted, and GeoGPT, which runs on the same computer, would be able to interact with the data through the code it generates. These files are made available in a working directory assigned to GeoGPT.

The second method used is to load the data into a spatial SQL database and provide the model with database schemas that can be used to generate queries. The datasets were uploaded to a Dockerized PostGIS database using QGIS’s DB Manager plugin.

The third method for data access is to use the OGC API Features standard. This method exposes the data stored in the PostGIS database through a web API that allows consumers of the API to download the data over HTTP as GeoJSON.

4. Datasets

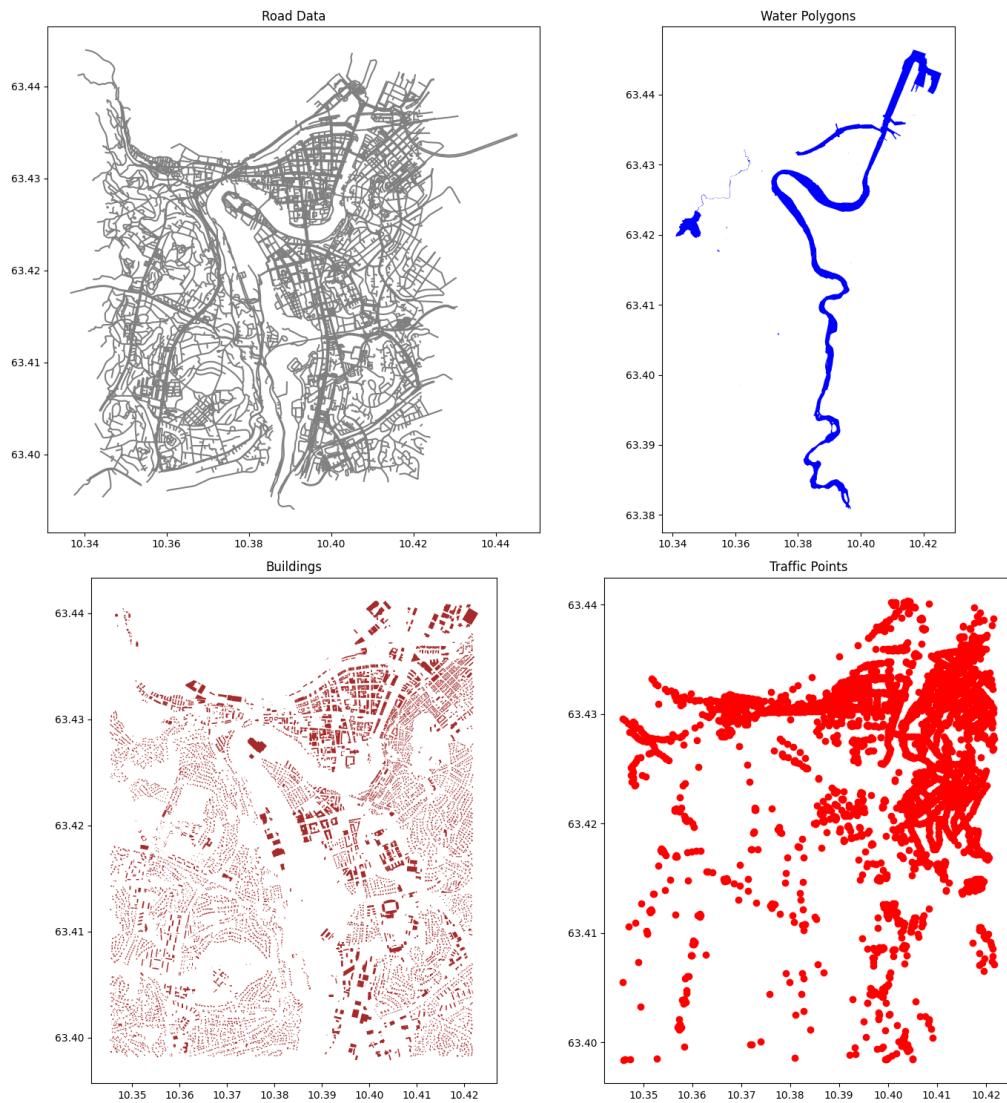


Figure 4.1.: A plot of four selected datasets constrained to a polygon of Trondheim

5. Architecture

The Architecture chapter will give a detailed description of the inner workings of GeoGPT, the thesis' proposed solution to an LLM-based GIS. GeoGPT is a webpage that features a chat interface and a web map, where users can ask questions and receive answers as text and/or map geometries. Section 5.1 will present a high-level overview of the different modules that make up GeoGPT, and the way that they interact with each other. Section 5.2 will delve into the architecture of the generic tool agent that is implemented using LangChain and LangGraph (see section 2.2 for detailed information), as well as the *three* different agent types that are implemented using this generic agent.

5.1. High-Level Application Architecture

A microservice architecture was adopted to simplify development and separate concerns among the system's different functionalities. The services are deployed as Docker Containers, and are orchestrated using Docker Compose. Figure 5.1 shows how the application is divided into five distinct services, and shows the general direction of information flow between these. Figure 5.2 goes into greater detail, showing how the different services interact with each other when the user first loads GeoGPT and enters a question into its chat interface.

5.1.1. LangChain Server

The *LangChain Server* service is the heart of the application, and is where the Large Language Model (LLM)-related logic is situated. It is responsible for taking requests from the *Web UI* and returning suitable responses in what becomes a client-server architecture between the two services. Table 5.1 shows the endpoints exposed by the server and how they can be used by a client.

5. Architecture

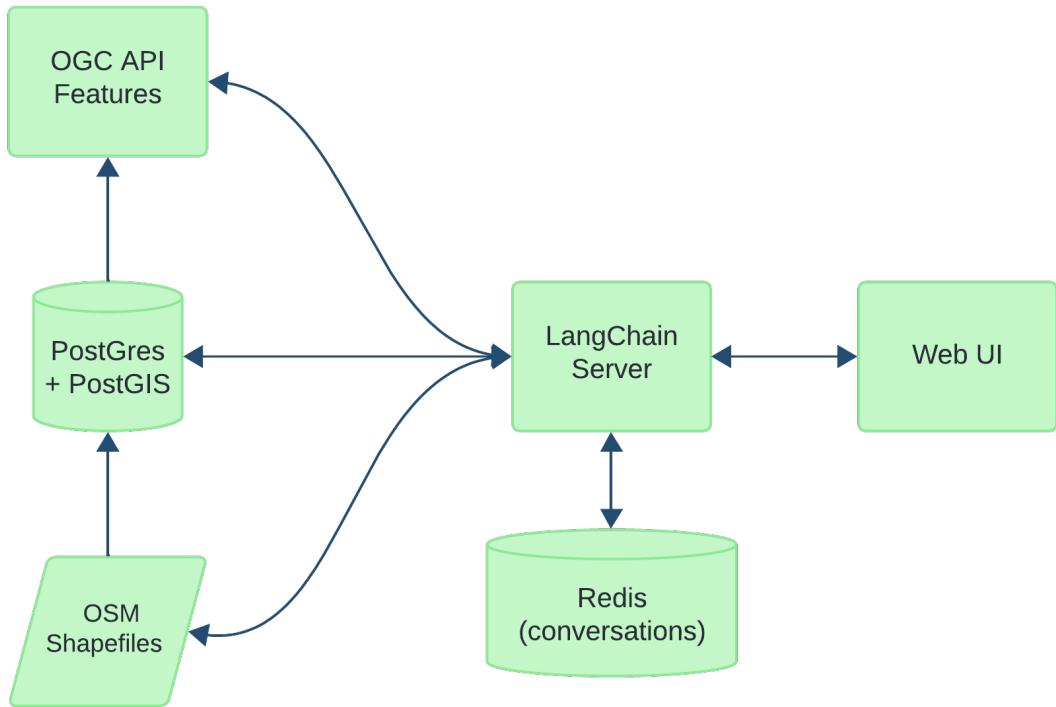


Figure 5.1.: Architecture overview

Table 5.1.: Summary of Server Endpoints

Endpoint	Method	Description
/session	GET	Takes a <code>session_id</code> as a query parameter, allowing the client to continue on a pre-existing session.
/session	POST	Creates a new session with an empty conversation.
/streaming-chat	GET	Endpoint for chatting the LLM. Takes a <code>message</code> as a query parameter and returns an event stream, allowing for token streaming from server to client.
/update-map-state	POST	Send the state of the client map to the server. Keeps the server updated on what layers are present in the map, their color, etc.
/geojson	GET	Takes a <code>geojson_path</code> as a query parameter. Allows the client to retrieve a given GeoJSON file that is stored in the working directory on the server.
/upload	POST	Allows the client to upload one or more files to the working directory on the server. 23

5. Architecture

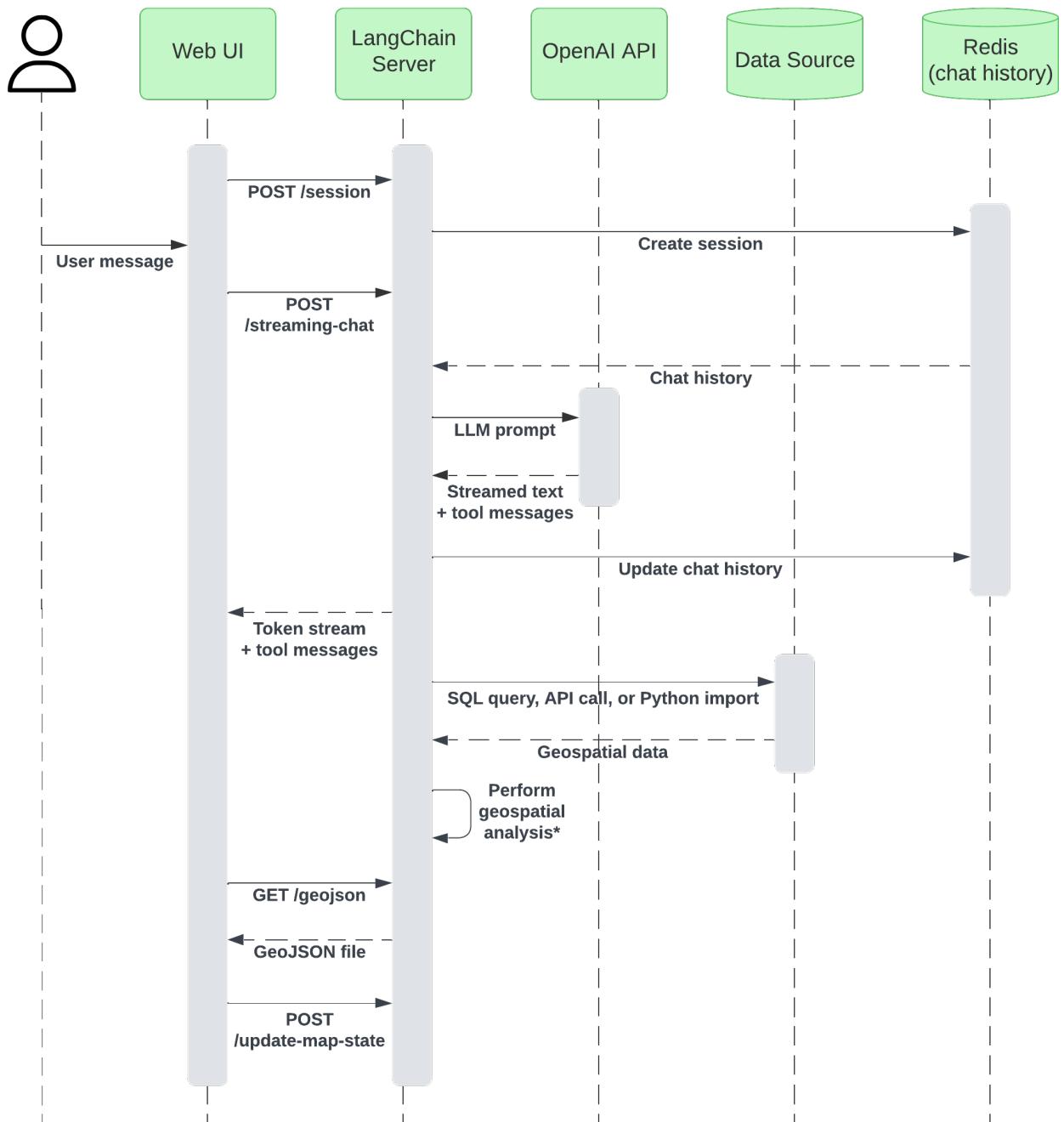


Figure 5.2.: Sequence diagram showing the information flow as the user loads and sends a message to GeoGPT

5. Architecture

5.1.2. Redis for Conversations

Redis (Sanfilippo, 2009) is a fast in-memory database that is often applied as a caching database that sits on top of some persistent database. It can also be used for vector-based storage, or as a simple NoSQL database. The latter option is the way it is used in GeoGPT’s architecture, and its sole purpose is to store conversations between GeoGPT and its users. Whenever a user starts a conversation with GeoGPT, a session object is stored in the Redis database, as can be seen in Figure 5.2. This object holds an array that represents the conversation. This array is written to every time the human, or GeoGPT, produces a message.

Storing messages, either in memory as a simple array or in a database like Redis, is crucial to enable multi-message conversations. In order for an LLM to act as a conversational agent, a chat history needs to be included in the prompt. In the case of GeoGPT, the entire chat history is included. This has the advantage of providing the LLM with the complete context of the chat history, but the disadvantage of potentially bloating the context window. Therefore, as the chat becomes longer each new token will be both more expensive and take longer to get generated. A long chat history could also make the resulting prompt exceed the token limit of the LLM. These issues were not considered much for this project and are left for future work.

5.1.3. PostGIS OGC API Features

A PostgreSQL database with the PostGIS extension, containing the datasets in ??, was deployed using Docker. On top of this database is an OGC API Features web server. It is realized through the `pramsey/pg_featureserv` Docker image¹ (CrunchyData, 2024), which is simply passed a connection string to the PostGIS database one wishes to expose on as an OGC API Features server. Any tables in that database which have a geometry column and a specified Coordinate Reference System (CRS) will be exposed on the server. `pg_featureserv` includes functionality like bounding box filtering, result limiting, and CQL filtering. These are added as query parameters in the URL used to fetch a collection’s items, for instance:

```
.../collections/{collection_id}/items.json?limit=1000&filter=name IS NOT NULL
```

In the internals of `pg_featureserv`, this URL will be converted to an SQL query that will be run against the database. Results of the query will be returned as GeoJSON. Code Snippet 5.1 shows an example of how CQL code is converted into SQL code:

```
1 \\\ CQL code passed through the `filter` query parameter
2 within(geom, POINT(0 0))
3
4 \\\ SQL code that will be run agains the database
```

¹https://hub.docker.com/r/pramsey/pg_featureserv

5. Architecture

```
5 ST_Within("geom", 'SRID=4326;POINT(0 0)' ::geometry)
```

Code Snippet 5.1: Conversion from CQL to SQL

5.1.4. Web UI

The user interface is made using SolidJS, a JavaScript framework for creating websites. The user interface is intentionally quite minimal, aiming to simplify the way we do GIS analysis.

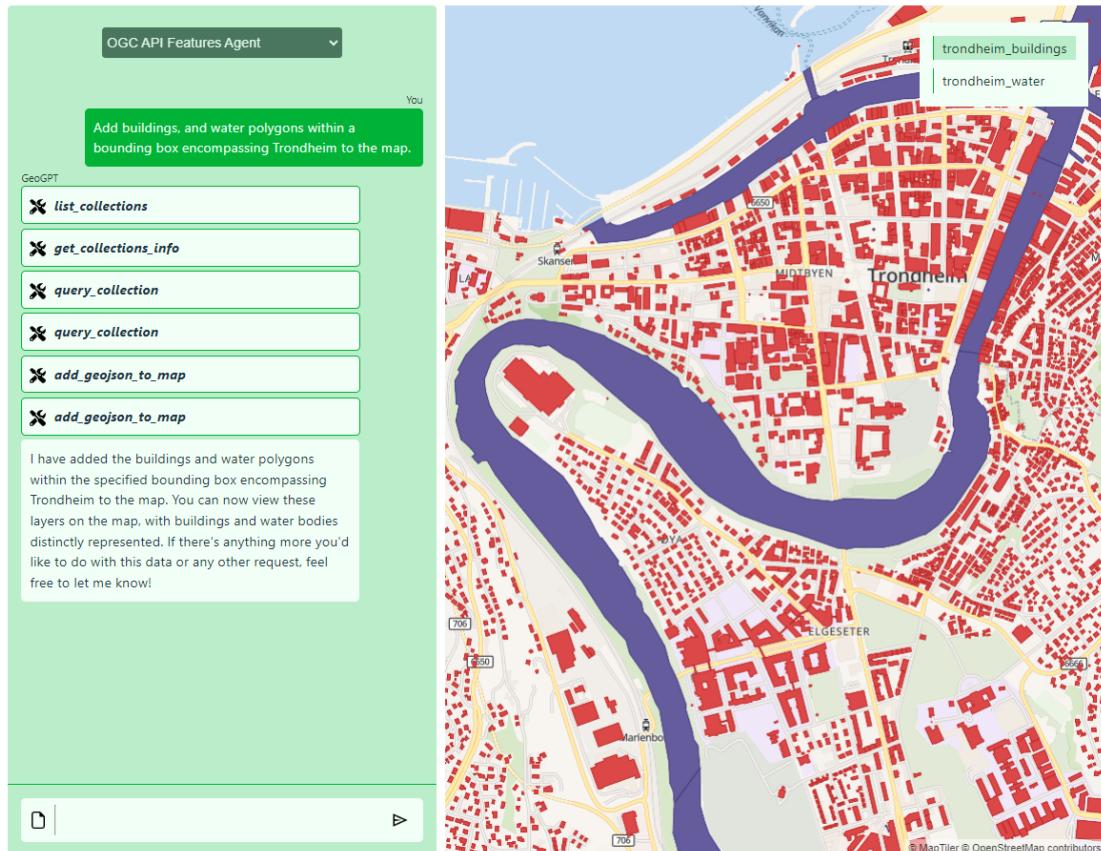


Figure 5.3.: Web UI

The chat interface was designed to imitate the interface of OpenAI's ChatGPT. Tokens and tool messages are streamed from the LangChain server, which helps the user follow GeoGPT's thought process as it is solving a problem. The generated text is streamed to the client token by token, and is in Markdown format. A library called *showdown*²

²<https://github.com/showdownjs/showdown>

5. Architecture

is used to convert Markdown into HTML, ensuring that tables, code blocks, lists, and other elements are properly rendered. Left of the input field on the bottom of the chat is a file upload button. Files that are uploaded here will be added to the working directory on the LangChain server.

The map is created using MapLibre,³ an open-source fork of Mapbox. A base map from OpenStreetMap (OSM) is used, fetched through a mapping platform called MapTiler.⁴ GeoJSON files that are fetched from the server will be added to the map automatically with a random color. On the top-right of the map is an overlay listing all layers that are currently present in the map. Using the arrow keys on this list will change the z-index in the map of the selected layer.

5.2. Agent Architecture

Three different agent were implemented for GeoGPT. These are listed below:

- **OGC API Features Agent** – Utilizes OGC API Features to interact with geographic data through standardized web interfaces.
- **Python Agent** – Accesses and manipulates shapefiles using Python code, enabling detailed geographic data processing.
- **SQL Agent** – Interacts with data stored in an SQL database, allowing for complex queries and data management within GeoGPT.

Common to these agents is their agentic architecture, which is described in subsection 5.2.1. They differ, however, through their assigned *tools*, which are described in subsection 5.2.2. Other slight differences are seen in the way that they are prompted. The prompting strategy used for GeoGPT will be discussed in subsection 5.2.3.

5.2.1. LangGraph Agent Implementation

The agentic behaviour of GeoGPT is implemented using LangGraph, which was described in detail in section 2.2. Figure 5.4 illustrates the flow between the various nodes that make up the agent. A state dictionary is passed between and updated by these nodes. Included in this state is the chat history, the path to GeoGPT’s working directory, a list of the current files in this directory, as well as less important state. The implementation is based upon a prebuilt implementation from LangGraph.⁵

The `__start__` node serves as the entry point of the agent. At this point, only one message is present in the state, namely the initial message from the user. The `__start__` node points to the `agent` node, where a prompt is constructed and passed to the Llama, so that a response to the user can be generated. This response could contain a

³<https://github.com/maplibre/maplibre-gl-js>

⁴<https://www.maptiles.com/>

⁵https://github.com/langchain-ai/langgraph/blob/main/langgraph/prebuilt/chat_agent_executor.py

5. Architecture

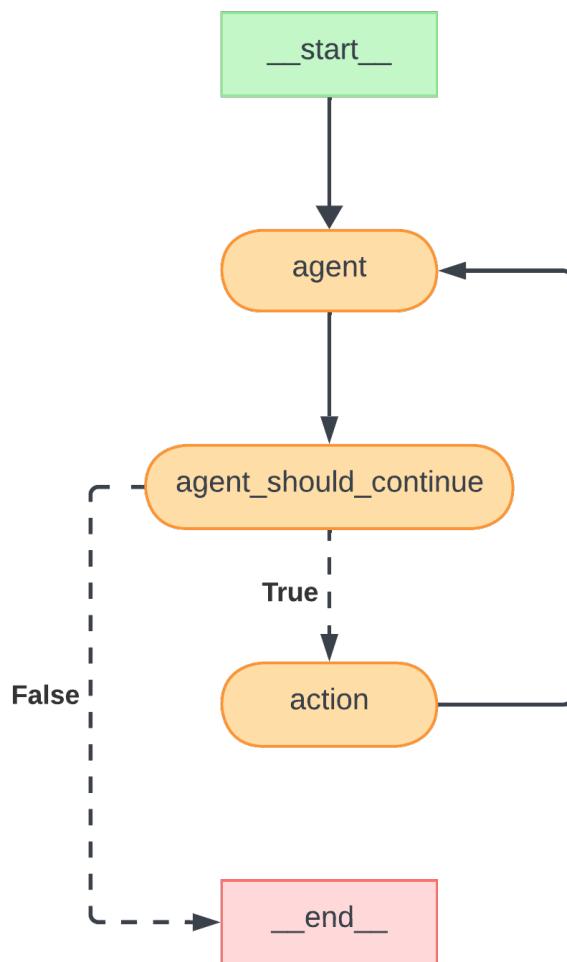


Figure 5.4.: Generic tool agent graph

5. Architecture

textual response, and/or instructions to execute some tool. For this reason, the current state containing both the user message and the AI message, is sent to the conditional node called `agent_should_continue`. This node simply checks if the last AI includes the “`tool_calls`” keyword argument. If this is the case then some tool should be executed, and we should “continue”.

In the `action` node, the values attached to the “`tool_calls`” keyword argument, which are LLM-generated, are used to invoke predefined tools. The LLM specifies tool *names*, as well as tool *parameters* that will be passed to the predefined tool (see subsection 2.1.5 for more details on how *function/tool calling* works). These tools are then executed, possibly having side effects on the server, and the return values from the tools are appended as messages to the chat history. Figure 5.5 illustrates this behaviour. This way, the agent can inspect the outcome of these tool invocations, and be notified if there are any errors occur in the tools due to the parameters it passed. Through the cyclic behaviour of the agent graph the agent can repeatedly call tools to try and answer the request from the user. When the agent finds no reason to call any more tools, `agent_should_continue` will return `False` so that the termination node, `__end__` is reached.

5.2.2. Tools

Table 5.2.: Overview of the agent types and their tools

Agent Type	Tools
OGC API Features	<code>list_collections</code> <code>get_collections_info</code> <code>query_collection</code> <code>python_repl_ast</code> <code>add_geojson_to_map</code>
Python	<code>python_repl_ast</code> <code>add_geojson_to_map</code>
SQL	<code>sql_db_list_tables</code> <code>sql_db_schema</code> <code>sql_db_query</code> <code>add_geojson_to_map</code>

Table 5.2 shows an overview of the tools that are available to each of the agents. As described in subsection 2.1.5, these are defined by a name, a description of the tool’s

5. Architecture

```
AI
sql_db_query
{
  "query": "SELECT * FROM osm_places_polygons WHERE name = 'Nordland' AND fclass = 'county';",
  "layer_name": "nordland_county"
}

TOOL
Query returned 1 feature.Below are the first 1 feature:
  id           geom  osm_id  code  fclass  population   name
0  4790  MULTIPOLYGON (((10.57806 65.47018, 10.57911 65... 408105 1041
  county      241084 Nordland

SYSTEM
Files written to the working directory ('/tmp/tmphhikab_f'):
- /tmp/tmphhikab_f/nordland_county.geojson

AI
add_geojson_to_map
{
  "geojson_path": "/tmp/tmphhikab_f/nordland_county.geojson",
  "layer_name": "nordland_county"
}

TOOL
`/tmp/tmphhikab_f/nordland_county.geojson` sent to client...

SYSTEM
State of map on client:
{
  "layers": [
    {
      "id": "nordland_county",
      "type": "fill",
      "source": "nordland_county",
      "paint": {
        "fill-color": "#893f23"
      }
    }
  ]
}
```

Figure 5.5.: Example of a chat history

5. Architecture

functionality, and a description of the parameters the tool expects.

The OGC API Features agent has access to a total of five tools. `list_collections`, which takes no parameters, sends a `GET` request to the `/collections` endpoint and uses the response to construct a list of the available collections. This tool gives the agent an overview of what kinds of data are available. Using the response from `list_collections` the agent can now invoke the `get_collections_info`. This tool takes a list of collection names and returns relevant information for each of these collections. This includes the JSON response from the “landing page” of the collection, which includes details such as the collection description, spatial extent, and available attributes. Furthermore, a list of common values for certain high-cardinality attributes is included in the tool’s response, as exemplified in Code Snippet 5.2. The percentages are obtained by querying a large number of features from `/collections/{collection_id}` and calculating the prevalences between them.

```
1 Property: fclass
2     tree: 71.1%
3     peak: 27.2%
4     beach: 0.9%
5     cave_entrance: 0.5%
6     spring: 0.2%
7     cliff: 0.1%
8     volcano: 0.0%
```

Code Snippet 5.2: Prevalences of common values for the `fclass` attribute in the `osm_natural_points` collection

The `query_collection` tool is used to retrieve features from a collection. It takes a collection name, a CQL filter, a bounding box, and a layer name. Based on these parameters, an URL like this is constructed:

```
https://localhost:9001/collections/{collection_id}/items.json?limit=10000&filter={cql_filter}&{bbox}
```

The features retrieved from this query is saved in the working directory on the LangChain server as “`{layer_name}.geojson`”. The message returned from the tool reads something like this: “Query returned 5627 features.” If the GeoJSON itself was returned as a tool message, this would quickly bloat the context window of the LLM, and therefore it is avoided.

Common for the OGC API Features agent and the Python agent is the `python_repl_ast` tool. This tool takes a string of Python code, executes it, and returns whatever the code prints to the standard output in a so-called Read-Eval-Print Loop (REPL). If the code errors, the error message is returned instead. This Python tool is the main way for these two agents to perform geospatial analyses. An advantage of using REPLs is that code can be executed in blocks, with variables from one block being shared with

5. Architecture

other blocks. This means that if the first block loads a large file into memory — an often time-consuming operation — the subsequent blocks can reuse this in-memory variable without reloading the data. This allows the LLM to quickly retry if the code should error, or if the outcome of the initial code wasn't as expected.

`add_geojson_to_map` is the only tool that is common for all three agent types. The tool's job is to add layers to the map on the client. It takes two parameters: the name or path of a GeoJSON file stored in GeoGPT's working directory, as well as a layer name. Invoking the tool will send a message to the client, which includes the full path to the file on the server. The client will then make a GET request to the server on the `/geojson` endpoint, asking for the contents of this file to be returned so that it can be added to the map (see Figure 5.2).

The SQL agent has tools very similar to the OGC API Features agent. It has a direct connection to the same database that the OGC API Features endpoint is served on top of. `sql_db_list_tables` is a tool that will list all database tables along with their description. `sql_db_schema` takes a list of table names and returns information about attributes, prevalence of different values in high-cardinality columns, and other details about these tables, much like `get_collections_info`. `sql_db_query` accepts arbitrary SQL code that will be executed against the database. The tool will make sure that query results that have a geospatial component will be stored as GeoJSON in GeoGPT's working directory, so that it can add the geometries to the map on the client.

5.2.3. Prompt Templating

Prompt templating is a way of producing prompts with a predefined structure for use with LLM. Figure 5.7 shows the prompt passed to GeoGPT's SQL agent when the user asks which county is the largest by size.

The overall template consists of a collection of tools, a system message, and a chat history. This entire template is passed to an LLM, each time the *agent* node from Figure 5.4 is invoked. LLMs have no inherent memory, so in order to have a chat conversation, the entire chat history needs to be passed with the prompt.⁶

The first half of the system message in Figure 5.7 is intended to give the LLM contextual awareness and hints on how it should work towards solving tasks. We start by telling LLM that it's a “helpful GIS agent/consultant that has access to an SQL database containing OpenStreetMap data”. This is a common way of giving LLMs a persona, and results in responses like the one in Figure 5.6. Continuing, we give it some information on how to string tool calls together to solve tasks. We tell it to first list available tables, then look up schemas of relevant tables, and then, using the information gathered, to construct an SQL query to answer the user's request. We also remind it that it needs to add the result of the analyses to the map, using the `add_geojson_to_map` tool.

⁶Several strategies have been developed by researchers to avoid having to pass the entire chat history to the LLM each time, as this eventually will bloat the context window which could make generation slower and worse in quality. Strategies include picking only the n last messages in the chat history, passing a summary of the chat instead of entire messages, or utilizing knowledge graphs.

5. Architecture

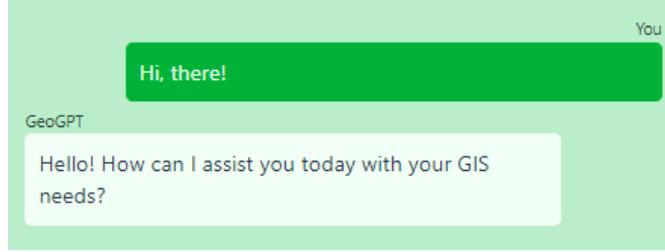


Figure 5.6.: Conversation showcasing the effect of giving an LLM a persona through the system message

The next part of the system message concerns the GeoGPT’s working directory. Having a working directory is important to control which files the agent has available, and also to make sure it doesn’t save files to the folder on the LangChain server where GeoGPT’s source code is located, as this would bloat the actual source code. In the first part of the system message, we tell the LLM where the working directory is located. This is especially important for the OGC API Features agent and the Python agent, as they need to manually read and write to the correct path using the Python code they generate. Continuing, we list the files that currently reside in the working directory. In the prompt in Figure 5.7, there “are currently no files written to the working directory”, but if there were they would be listed in a bullet point list.

The final part of the system message is a checklist to inform the LLM about common pitfalls it could run into when generating SQL code. A similar checklist for Python is provided for the other two agents, with reminders about using metric CRSs when doing area calculations, etc.

Normally, prompts for LLM contain only a single system message. GeoGPT, however, features a — to the author’s knowledge — novel usage of system messages. As can be seen in Figure 5.5, system messages are appended mid-conversation, providing updates about state changes in the system. The first system message is added because a new file has been added to the working directory, as a result of invoking the `sql_db_query` tool. The system message helps GeoGPT stay up-to-date on what files it has available, and in Figure 5.5, it uses this information immediately to add the only file available to the map using `add_geojson_to_map`. Another system message is eventually added to tell the GeoGPT that client map has been modified. The message helps ensure GeoGPT that the invocation of `add_geojson_to_map` was successful. If it wasn’t, the system message would inform it about this.

5. Architecture

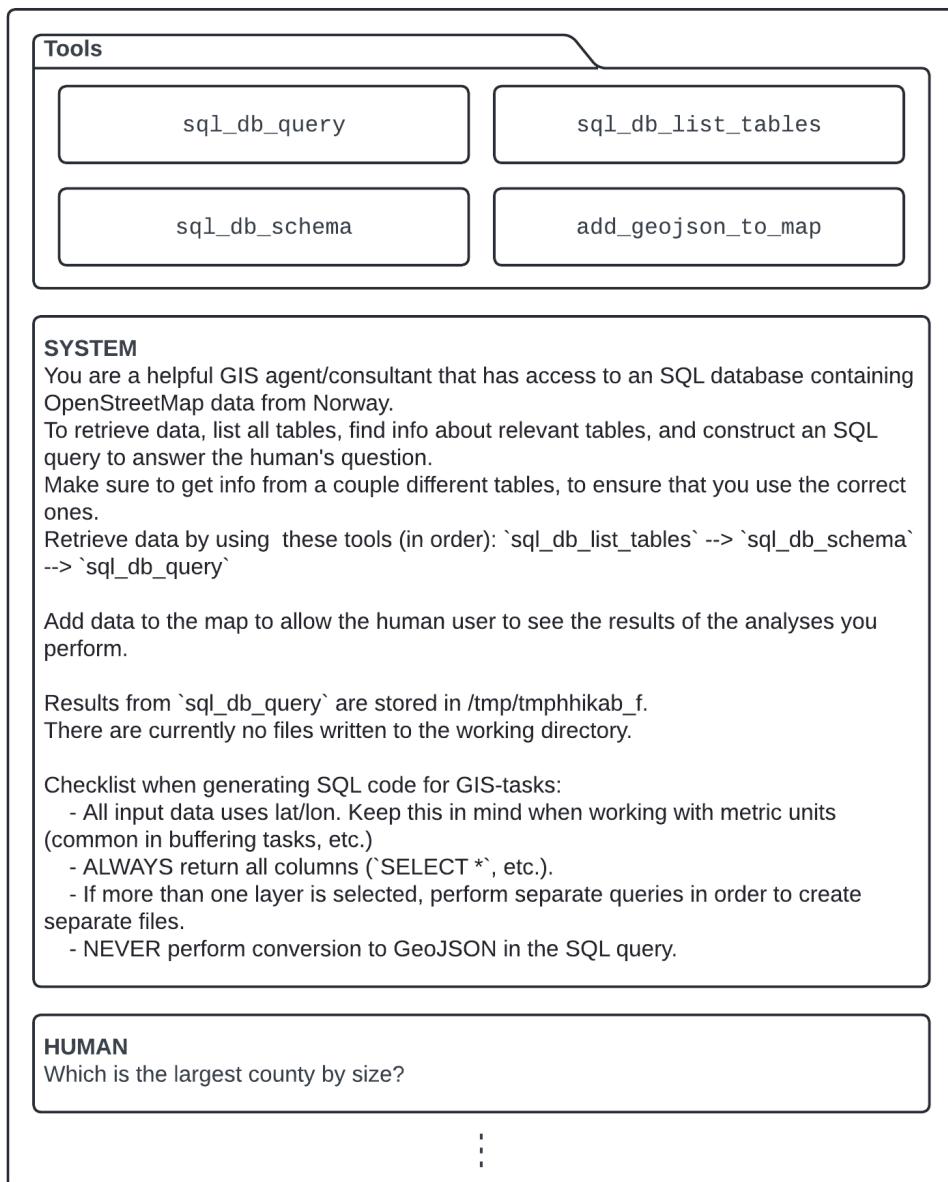


Figure 5.7.: Chat template

6. Experiments

The Experiments chapter is divided into two main sections: section 6.1, which describes the methodology and rationale behind the conducted experiments; and section 6.2, which presents the experimental results and offers interpretations of these findings. The latter section will also present several example results, including screenshots of the application and code that generated by GeoGPT.

6.1. Experimental Setup

The experiments conducted to evaluate the performance of GeoGPT on geospatial tasks are divided into two approaches. The first approach, presented in subsection 6.1.1, is intended to evaluate GeoGPT’s ability to solve a variety of geospatial tasks. The second approach, presented in subsection 6.1.2, aims to uncover the importance of the quality of the initial user prompt. Subsection 6.1.3 describes the hardware on which the experiments were executed, and the specific LLMs that were used.

6.1.1. GIS Benchmark

The first approach seeks to evaluate GeoGPT’s ability to successfully answer GIS-related questions that have concrete answers. To do this, a benchmarking dataset was constructed. This dataset is of set of 12 GIS-related questions with corresponding correct answers, and can be viewed in it’s entirety in Table A.1 in section A.1.

Outcome Evaluation

Each test run’s answer will be manually evaluated, and the outcome will be annotated as one of *success*, *partial success*, and *failure*. Table 6.1 shows the guidelines used when assigning outcome scores.

Cost and Duration

The application is hooked up to LangChain AI’s tracing system, *LangSmith*.¹ Apart from being a useful tool for debugging purposes, LangSmith provides a simple way of obtaining detailed data for token and time usage for a particular run, as well as the total cost of the run. These metrics will be recorded and used in the evaluation of GeoGPT.

To summarize, the following metrics are recorded for a given test run:

¹<https://www.langchain.com/langsmith>

6. Experiments

Table 6.1.: Description of Success

Outcome	Guideline
Success	The question was answered correctly and little to no follow-up from the user was required to produce the desired outcome. No false assumptions were made by the system when answering the question.
Partial Success	Portions of the question were answered correctly or semi-correctly, and/or some follow-up from the user was required to guide the system toward the solution.
Failure	The question was answered incorrectly answered and/or false assumptions were made by the system while attempting to answer the question.

- The outcome of the test (*success, partial success, or failure*)
- The total duration in seconds
- The total number of tokens used
- The total cost for the run in American dollar

Repeatability

Another aspect that the benchmarking approach will try to evaluate is the consistency of the system, its ability to repeatedly provide an acceptable answer to the same user question. Each of the 12 questions are therefore asked three times per agent type. As there are three agent types (see section 5.2), the total number of test runs becomes the following:

$$12 \text{ questions} \cdot 3 \text{ agent types} \cdot 3 \text{ repetitions} = 108 \text{ tests}$$

The annotated outcomes are encoded using the ordinal encoding presented Table 6.2. A higher value indicates a better outcome. These encoded outcome values enable standard deviation calculations, which serve as a suitable measure for assessing repeatability.

6.1.2. Prompt Quality Experiment

A second set of experiments are constructed to evaluate the importance of the initial question/prompt from the user. As stated in the Background and Motivation chapter, part of the motivation for developing an LLM-driven GIS like GeoGPT is to make GIS more accessible to non-experts. It should therefore be valuable to assess the extent to which a carefully constructed prompt by a GIS expert can enhance the system's output.

6. Experiments

Table 6.2.: Encoding for Test Outcome

Outcome	Encoded Value
Success	2
Partial Success	1
Failure	0

For these experiments, the three hardest questions from the benchmarking are picked. Then, for each of these three questions, a *novice*-level and an *expert*-level prompt is constructed. The novice-level prompt is as simple as possible, while the expert-level prompt is more elaborate and written as a step-by-step recipe for solving the problem. The results will compared in a similar way to the comparison between the agent types in the benchmarking experiment.

6.1.3. Hardware and Model Version

All experiments were conducted on a Lenovo ThinkPad E490, which has an Intel® Core™ i7-8565U CPU @ 1.80GHz processor, 15.8 GB usable RAM, and 256 GB SSD storage. Everything but the LLM inference was executed locally. Text generation was done using OpenAI’s API.

It is worth noting that two slightly different models were used during testing. This is due the release of the `gpt-4-turbo-2024-04-09` in mid-April. According to OpenAI, “this new model is better at math, logical reasoning, and coding” compared to `gpt-4-0125-preview`,² which is the model that was for the first test runs. At `gpt-4-turbo-2024-04-09`’s release, a decision was made to use this for the remaining experiments. The experiments that had already been conducted were not re-run due to time constraints and a confidence that these slight model upgrades would not significantly change the outcome of the experiments.

6.2. Experimental Results

Subsection 6.2.1 and subsection 6.2.2 will present the outcome of the experiments presented in subsection 6.1.2 and subsection 6.1.2, respectively. Graphs created for this chapter are created using Matplotlib,³ a Python library suitable for creating visualizations like bar charts, box plots, etc.

²OpenAI has a GitHub repository containing the code they use to evaluate their Large Language Models (LLMs) and benchmark results for OpenAI models and reference models from other companies: <https://github.com/openai/simple-evals>.

³<https://matplotlib.org/>

6. Experiments

6.2.1. GIS Benchmark — Results

Outcome Evaluation

Figure 6.1 displays a bar chart for the outcome distribution per agent. From Figure 6.1, we can read that the OGC API Features and Python agent have comparable results, and that the SQL-based agent performs significantly better compared to the other two in terms of producing the desired outcome. The SQL agent has a success rate of 69.4% compared to the other two, which share a success rate of 38.9%. Possible reasons for this are discussed in section 7.2 in the Discussion.



Figure 6.1.: Outcome distribution between different agent types

Cost and Duration

This section features three different box plots: Figure 6.2, Figure 6.3a, and Figure 6.3b. The Matplotlib implementation of the box plot follows the description found on Wikipedia.^{4,5} Box plots allow us to easily see where the 0th (Q_0), 25th (Q_1), 50th (Q_2), 75th (Q_3), and 100th (Q_4) percentiles of the datasets lie, as well as the dataset's outliers. Outliers are those data points fall outside 1.5 times interquartile range, that is, the distance between Q_3 and Q_1 in each direction.

⁴https://matplotlib.org/stable/api/_as_gen/matplotlib.pyplot.boxplot.html

⁵https://en.wikipedia.org/wiki/Box_plot

6. Experiments

Figure 6.2 displays a box plot with a logarithmic y-axis showing task durations across the different agent types. Here we can see that the SQL agent spends the least amount of time per task. The OGC API Features agent has a slightly higher median and a few time-consuming outliers. The Python agent is the odd one out with a median of ~ 82 seconds, a $Q3 \sim 293$ seconds, and a $Q4 \sim 984$ seconds. The large gap to the other two agents is largely due to the Python agent's tendency to load large datasets into memory. For instance, when attempting the task of calculating the difference between the polygon outlining Oslo and water polygons, the Python agent used nearly 40 minutes on the entire task. 94% of the time was spent executing the code presented in Code Snippet 6.1. The main reason for the long execution time is line 8, where the whole `osm_landuse_polygons.shp` dataset is loaded into memory. This dataset has a size of $\sim 1.4\text{GB}$, and loading such amounts of data in this way is very time-consuming. The Python agent was the only agent with such issues because the OGC API Features agent is limited to 10,000 features per dataset, and the SQL agent does not load the data into memory like the other agents do.

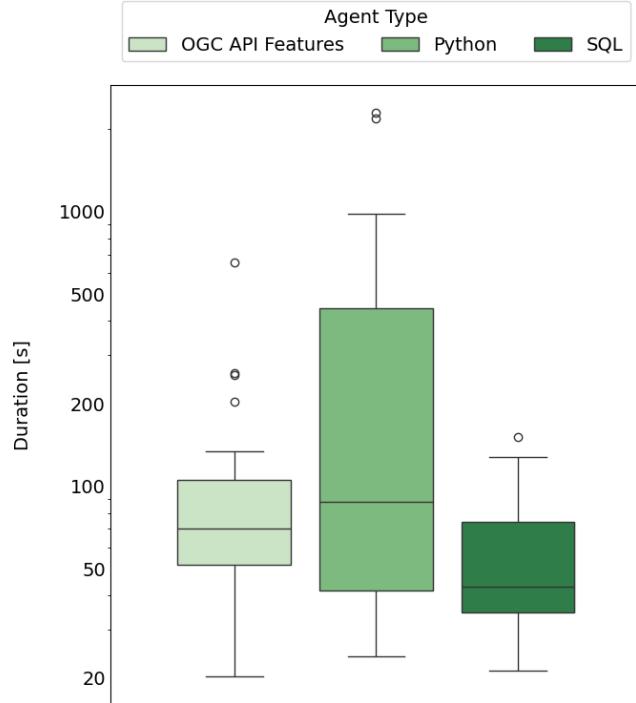


Figure 6.2.: Task Durations for the Agent Types

Figure 6.3 shows box plots for token usage and cost per run. Naturally, these figures appear very similar, since prices for the input tokens and generated output tokens are

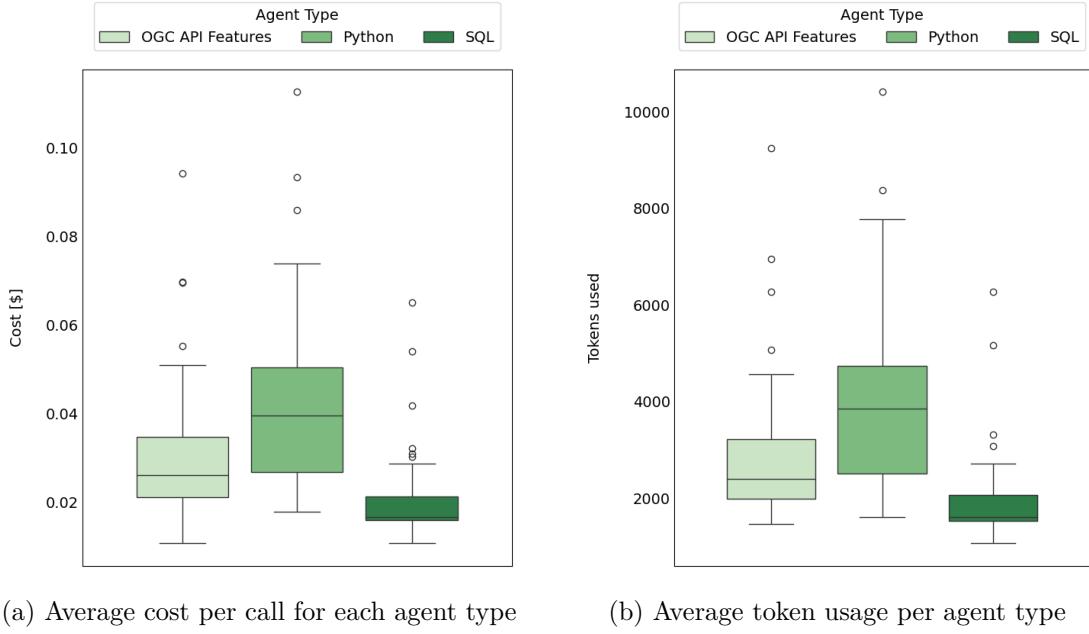
6. Experiments

```
1 import geopandas as gpd
2
3 # Paths to the shapefiles
4 landuse_path = '/tmp/tmpsutdy6it/osm_landuse_polygons.shp'
5 places_path = '/tmp/tmpsutdy6it/osm_places_polygons.shp'
6
7 # Load the data from shapefiles
8 landuse_gdf = gpd.read_file(landuse_path)
9 places_gdf = gpd.read_file(places_path)
10
11 # Filter out 'residential' areas from the landuse data
12 residential_gdf = landuse_gdf[landuse_gdf['fclass'] == 'residential']
13
14 # Compute the spatial difference to exclude residential areas from
15 # the places data
16 oslo_outline = gpd.overlay(places_gdf, residential_gdf, how='difference')
17
18 # Path for the output GeoJSON file
19 output_path = '/tmp/tmpsutdy6it/oslo_outline_no_residential.geojson'
20
21 # Save the resulting GeoDataFrame to a GeoJSON file
22 oslo_outline.to_file(output_path, driver='GeoJSON')
23
24 # Output the path to the saved file
25 print(output_path)
```

Code Snippet 6.1: GeoGPT-generated Python code aimed at computing the difference between the Oslo outline and water features within it

6. Experiments

fixed for the models used. According to OpenAI’s webpages,⁶ they charge \$10 per million input tokens and \$30 per million output tokens for their GPT-4 Turbo model. From the results of the experiments, a ratio of approximately 10.7 per million tokens — either input *or* output — was calculated, which is closest to the input token price. This shows that a lot more input tokens were used than output tokens.



(a) Average cost per call for each agent type

(b) Average token usage per agent type

Figure 6.3.: Cost and token usage

An observation that can be made from Figure 6.4 is that the correlation between duration and token usage is inconsistent across the agent types. The OGC API Features agent and SQL agent have strong correlation between these metrics, but the Python agent has next to none. This supports the observation that, for the Python agent, code execution time, particularly the time spent on expensive imports, is likely the main contributor towards the total run time.

Another observation that can be made from the Figure 6.4 is the slight negative correlation between the encoded outcome and each of the following variables: token usage, duration, and total cost. This suggests that a task that takes longer to complete is likely to be more expensive in terms of token usage, and more likely to produce an undesirable outcome. Possible reasons as to why this is the case will be explored in subsection 7.3.1 the Discussion.

⁶<https://openai.com/pricing>

6. Experiments

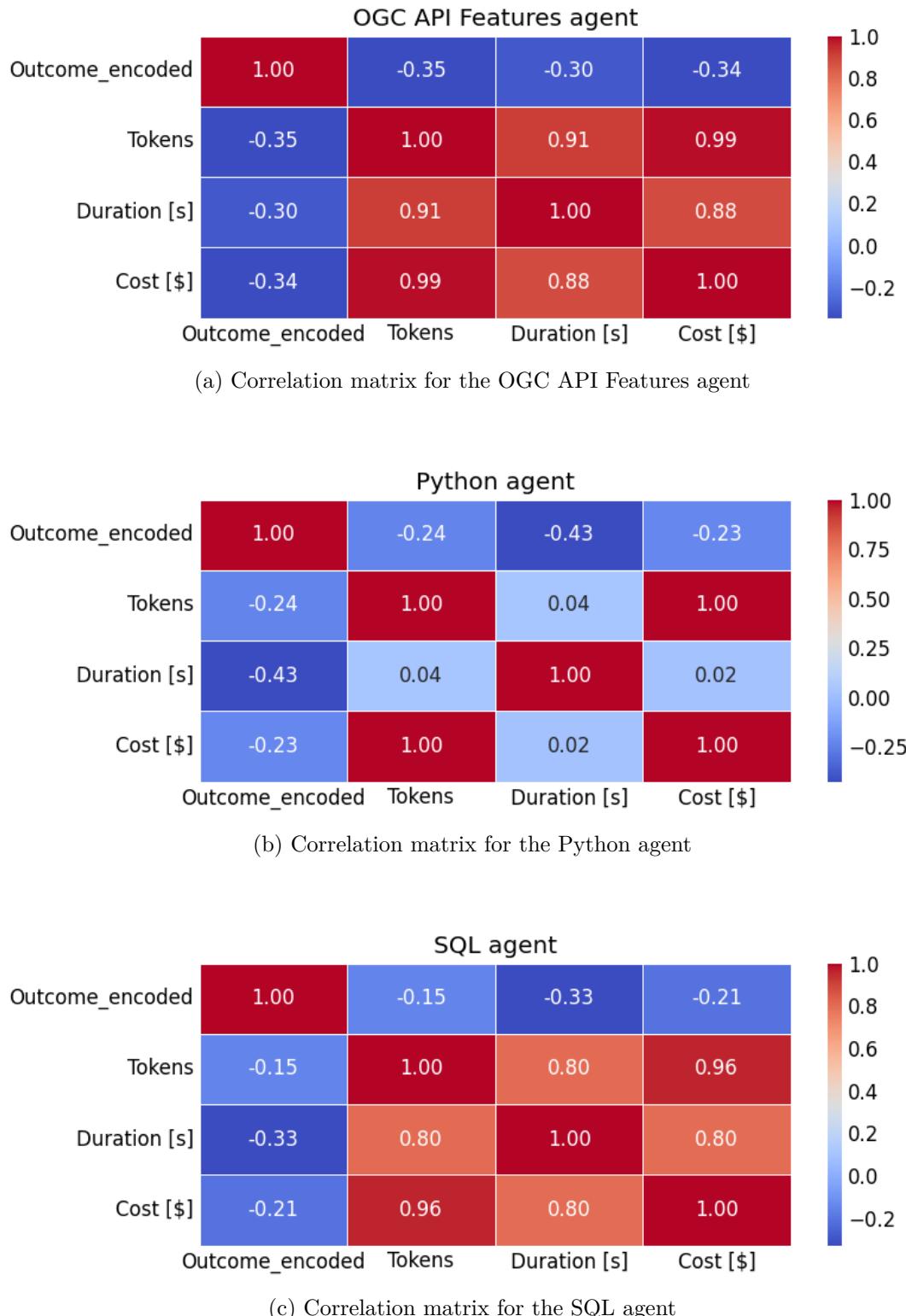


Figure 6.4.: Correlation matrices for metrics for the three agent types

6. Experiments

Repeatability

Table 6.3 shows the average standard deviation of the outcomes for each agent type, as well as their mean. Standard deviations were calculated for each triplet of identical test samples, in which both the question and the agent type remained the same, using the encoded outcomes (see Table 6.2) were used. Taking the average produced the numbers found in Table 6.3. These numbers indicate that there is a notable amount of inconsistency in GeoGPT’s answers, on average deviating with more than a third of an outcome category (0.376 for the encoded values), when task and agent type stay the same.

Table 6.3.: Standard Deviation by Agent Type

Agent Type	Outcome Std. Deviation
OGC API Features	0.552
Python	0.337
SQL	0.241
Mean	0.376

Successful Responses

Figure 6.5 shows a successful response from GeoGPT’s SQL agent when asked how many counties the Glomma river runs through.⁷ As a final answer to the initial prompt it lists all four correct counties in a numbered list. Code Snippet 6.2 shows the code generated by GeoGPT for the second invocation of `sql_db_query`. The second invocation is nearly identical to the first one, but in the first invocation the geometry column was not included in the response (which the tool informs GeoGPT about), meaning GeoGPT was unable to add the results to the map. It therefore decided to run the query again, now making sure to include the geometry column in the query result, before adding the result to the map using `add_geojson_to_map`. A follow-up message instructed GeoGPT to “Add Glomma to the map”, allowing for visual verification that the answer it gave was correct.

⁷This run was not part of the results used for evaluation due to a bug in the `sql_db_query` tool that caused it to unnecessarily execute queries twice. GeoGPT’s response would not be different had the bug not been present for the run, but it would have taken longer to complete.

6. Experiments

SQL Agent

You

How many counties does Glomma run through?

GeoGPT

- X `sql_db_list_tables`
- X `sql_db_schema`
- X `sql_db_query`
- X `sql_db_query`
- X `add_geojson_to_map`

The river Glomma runs through four counties in Norway. These counties are:

1. Akershus
2. Innlandet
3. Østfold
4. Trøndelag

You can see these counties highlighted on the map.

Add Glomma to the map.

GeoGPT

- X `sql_db_query`
- X `add_geojson_to_map`

The river Glomma has been added to the map. You can see it as a blue line, traversing through the highlighted counties.

□ | ▷

glomma_river

counties_glomma_runs_through

MØRE OG ROMSDAL

DALARNA COUNTY

BUSKERUD

OSLO

VARMILAND COUNTY

TELEMARK

WESTFOLD

Skien

Karlstad Town

© MapTiler © OpenStreetMap contributors

Figure 6.5.: Successful response from GeoGPT's SQL agent when asked how many counties the Glomma river runs through

6. Experiments

```
1 WITH river AS (
2     SELECT geom
3     FROM osm_waterways_lines
4     WHERE fclass = 'river' AND name ILIKE 'Glomma'
5 ),
6
7 places AS (
8     SELECT geom, name
9     FROM osm_places_polygons
10    WHERE fclass = 'county'
11 )
12
13 SELECT DISTINCT places.name AS county_name, places.geom AS geom
14 FROM river, places
15 WHERE ST_Intersects(river.geom, places.geom);
```

Code Snippet 6.2: SQL code generated by GeoGPT to retrieve the counties that the Glomma river runs through

Figure 6.6 shows a partially successful attempt from GeoGPT’s Python agent to calculate the number of trees along Munkegata in Trondheim. While setting a definitive answer to the question is difficult, the correct answer was set to about 100, assuming a 20-meters buffer around the line segments that make up the road, is used. In Code Snippet 6.3 we can see how GeoGPT loads the road data around Munkegata in Trondheim, using an approximate bounding box to avoid getting roads named “Munkegata” in other cities. It then filters and saves the loaded data to only include features with name == ‘Munkegata’. In the following code block (see Code Snippet 6.4), GeoGPT gets the tree data by filtering the osm_natural_points dataset on the fclass attribute, before loading the road data for Munkegata from the subsequent code block, converting both datasets to a metric CRS, creating a buffer around the road data, performing an intersection between the buffer, and finally printing the tree_count to the standard output.

6. Experiments

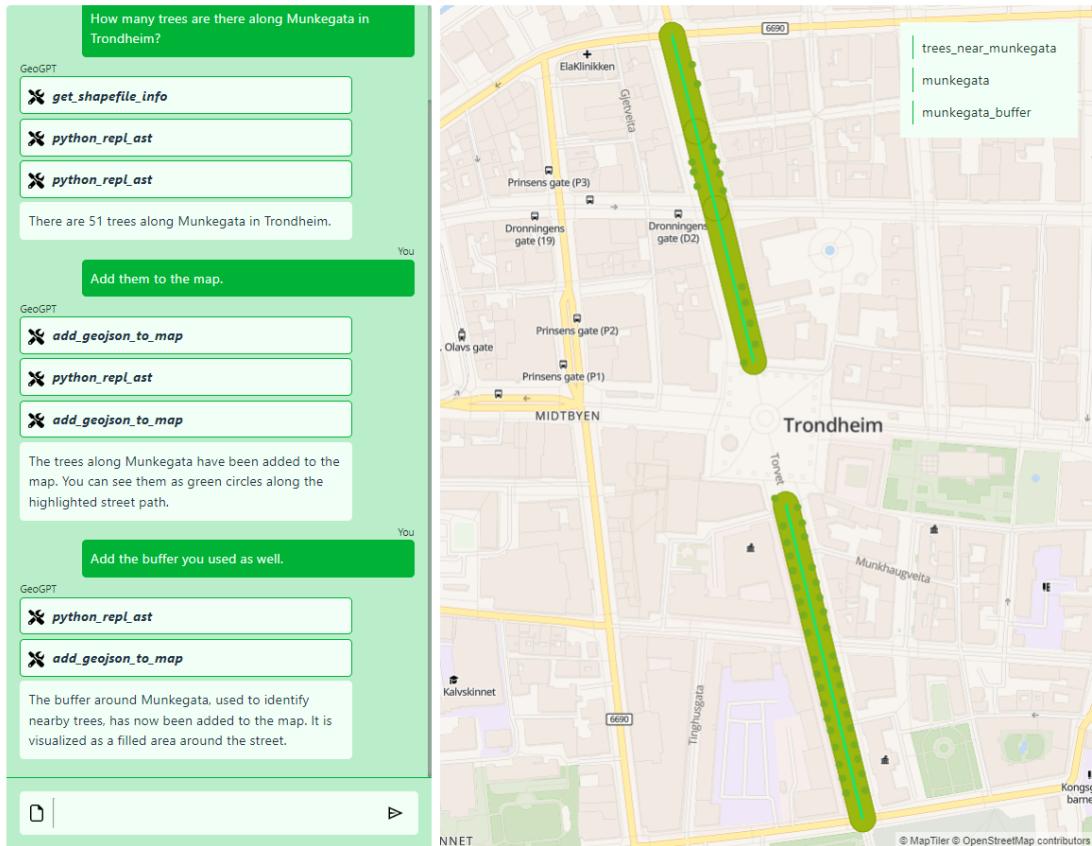


Figure 6.6.: Partially successful response from GeoGPT’s Python agent when asked how many trees there are along Munkegata in Trondheim

6. Experiments

```
1 import geopandas as gpd
2
3 # Load roads data with a focus on Munkegata
4 roads_path = '/tmp/tmp18zsdj6m/osm_roads_lines.shp'
5 munkegata_bbox = (10.3910, 63.4270, 10.4030, 63.4330) # Approximate
       bounding box around Munkegata
6 munkegata_roads = gpd.read_file(roads_path, bbox=munkegata_bbox)
7
8 # Filter for Munkegata
9 munkegata = munkegata_roads[munkegata_roads['name'] == 'Munkegata']
10
11 # Save filtered Munkegata road to GeoJSON
12 munkegata_path = '/tmp/tmp18zsdj6m/munkegata.geojson'
13 munkegata.to_file(munkegata_path, driver='GeoJSON')
14
15 print(munkegata_path)
```

Code Snippet 6.3: Python code to exclude roads named Munkegata *outside* of Trondheim

```
1 from shapely.geometry import box
2 import geopandas as gpd
3
4 # Load the trees points
5 natural_points_path = '/tmp/tmp18zsdj6m/osm_natural_points.shp'
6 trees = gpd.read_file(natural_points_path)
7
8 # Filter for trees
9 trees = trees[trees['fclass'] == 'tree']
10
11 # Load Munkegata road
12 munkegata_path = '/tmp/tmp18zsdj6m/munkegata.geojson'
13 munkegata_road = gpd.read_file(munkegata_path)
14
15 # Project to UTM for accurate distance measurement
16 munkegata_utm = munkegata_road.to_crs(epsg=32633)
17 trees_utm = trees.to_crs(epsg=32633)
18
19 # Buffer Munkegata by 10 meters to include nearby trees
20 tree_buffer = munkegata_utm.buffer(10)
21
22 # Spatial join trees within buffered Munkegata
23 trees_near_munkegata = gpd.sjoin(trees_utm, gpd.GeoDataFrame(geometry=
       tree_buffer), how='inner', op='intersects')
24
25 # Count the number of trees
26 tree_count = trees_near_munkegata.shape[0]
27 print(tree_count)
28
```

6. Experiments

```
29 # --> 51
```

Code Snippet 6.4: Python code to calculate the number of trees along Munkegata in Trondheim

Unsuccessful Responses

While GeoGPT was able to get a lot of questions correct, it provided more or less unsuccessful answers at least once for all questions but one. Figure 6.7 exemplifies a common mistake with GeoGPT, which will be a central topic in the Discussion. Here, the Python agent has failed to recognize that the road datasets it has available is in fact for the whole of Norway, not just Oslo. As Code Snippet 6.5 shows, the agent performs the correct analysis based on the assumptions it has already made, but the answer is far from what the user enquired.

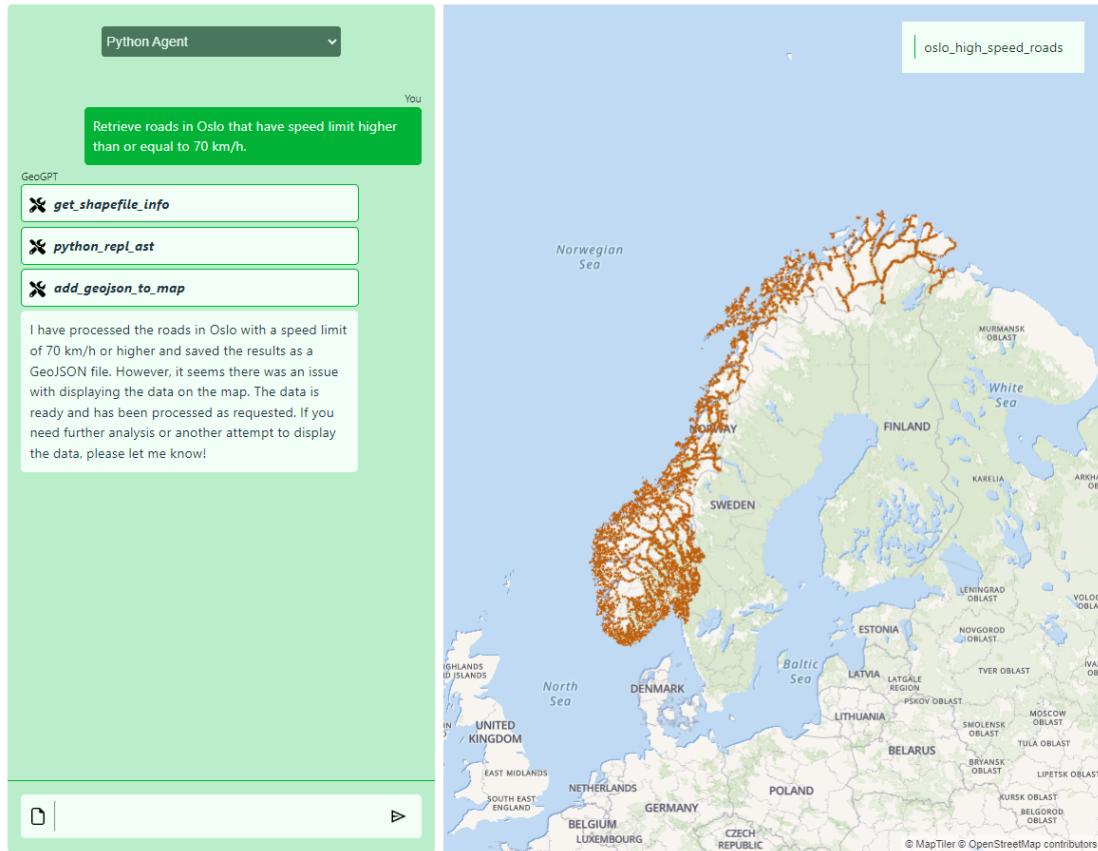


Figure 6.7.: Unsuccessful attempt by GeoGPT's Python agent to retrieve high-speed roads in Oslo

6. Experiments

```
1 import geopandas as gpd
2
3 # Define the file path for the roads shapefile
4 shapefile_path = '/tmp/tmp1hvr_c5/osm_roads_lines.shp'
5
6 # Read the shapefile with a condition on maxspeed being >= 70
7 roads_gdf = gpd.read_file(shapefile_path)
8 roads_high_speed_gdf = roads_gdf[roads_gdf['maxspeed'].astype(float)
9                                >= 70]
10
11 # Save the filtered GeoDataFrame as GeoJSON
12 output_geojson_path = '/tmp/tmp1hvr_c5/oslo_high_speed_roads.geojson'
13 roads_high_speed_gdf.to_file(output_geojson_path, driver='GeoJSON')
14 print(output_geojson_path)
```

Code Snippet 6.5: Python code that saves high-speed roads with a misleading filename

Figure 6.8 shows an unsuccessful attempt from GeoGPT’s OGC API Features agent to create a geodesic line between Oslo Airport Gardermoen and Bergen Airport Flesland. Code Snippet 6.6 shows GeoGPT’s attempt to download a point feature representing “Oslo Airport”. It turns out that no such point feature exists, and only a polygonal feature in another dataset is available for the airport. The same is the case for Bergen Airport. GeoGPT made many attempts at fetching point data for the airports, but of course none of them returned any results.

Eventually, GeoGPT downloads a bunch of features from the `osm_transport_points` collection twice, names them “oslo_airport.geojson” and “bergen_airport.geojson”, and produces the nonsensical code displayed in Code Snippet 6.7, where it picks the first point feature in each collection and assumes that their coordinates are those of Oslo and Bergen. These two features actually correspond to are Rakkestad Airport and a bus stop on a road along Ofotfjorden in Nordland. In addition to this mistake, subsequent attempts at creating a geodesic line between the locations were unsuccessful.

6. Experiments

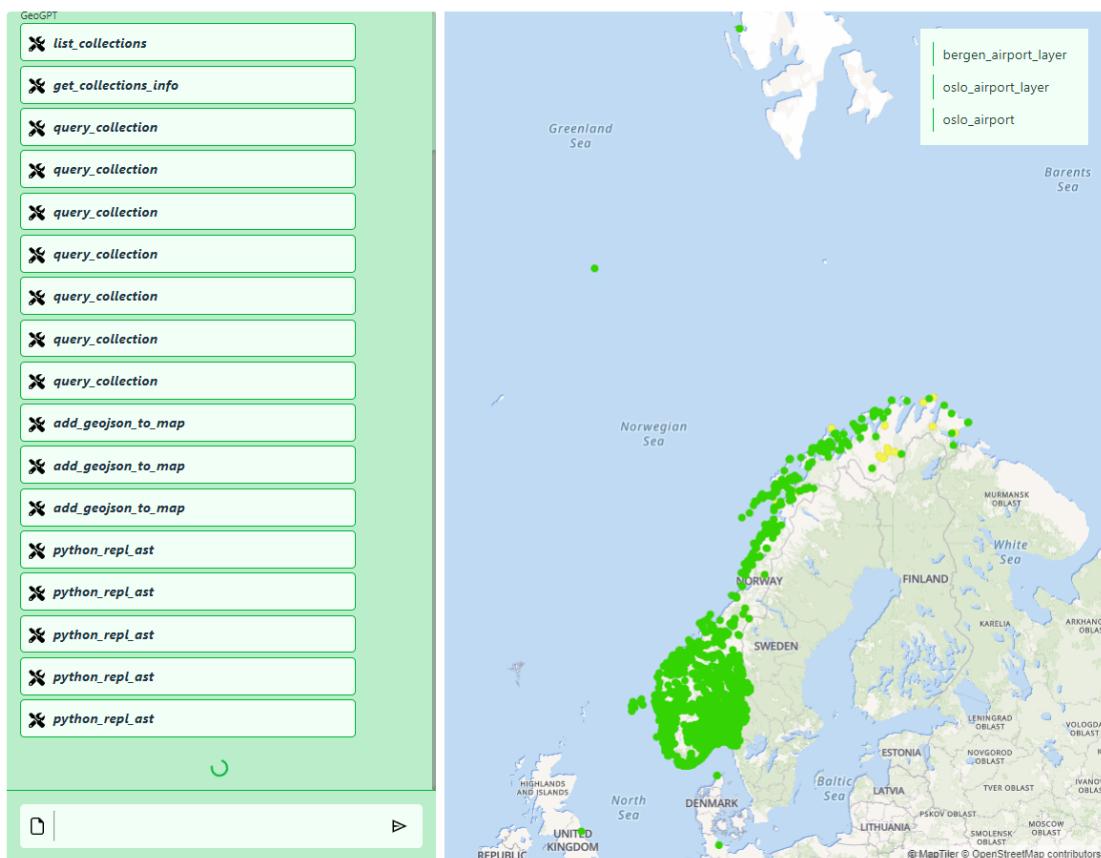


Figure 6.8.: Unsuccessful attempt by GeoGPT's OGC API Features agent to create a geodesic line between Oslo Airport Gardermoen and Bergen Airport Flesland

6. Experiments

```
1 {
2   "collection_name": "public.osm_transport_points",
3   "cql_filter": "fclass='airport' AND name='Oslo Airport'",
4   "layer_name": "oslo_airport"
5 }
6
7 --> No features were found at http://localhost:9001/collections/public
     .osm_transport_points/items?limit=10000&filter=fclass='airport' AND
     name='Oslo Airport'.
8 Try to change the parameters, or make them less restrictive.
```

Code Snippet 6.6: Tool invocation of `query_collection` with no results

6. Experiments

```
1 import geopandas as gpd
2
3 oslo = gpd.read_file('/tmp/tmpwlojm_1k/oslo_airport.geojson')
4 bergen = gpd.read_file('/tmp/tmpwlojm_1k/bergen_airport.geojson')
5
6 oslo_coords = oslo.geometry[0].coords[0]
7 bergen_coords = bergen.geometry[0].coords[0]
8
9 oslo_coords, bergen_coords
10
11 # --> ((11.3469259, 59.397229), (16.919517, 68.3459))
```

Code Snippet 6.7: “Desperate” Python code that picks the first feature of the downloaded collection and *hopes* that they correspond to Gardermoen and Flesland

6.2.2. Prompt Quality Experiment — Results

Figure 6.9 shows the outcomes of the experiments where importance of the quality of the user’s initial prompt, was assessed. It is clear to see that *expert*-level prompting significantly outperforms *novice*-level prompting, the latter of which produced *no* fully successful responses.

Figure 6.10 compares the responses that GeoGPT’s OGC API Features agent managed to produce for the different prompt levels on one of the tasks. The *novice*-level prompt was as follows:

“Could you count how many trees there are on Munkegata street in Trondheim?”

The *expert*-level prompt included a series of instructions:

- “1. List all datasets that could possibly include trees.
2. Find the correct feature class and filter the relevant dataset to access tree data for Trondheim. Use a bounding box to reduce the number of trees to analyse.
4. Fetch road data for Munkegata. Use a bounding box for Trondheim in case there are streets elsewhere named Munkegata.
5. Convert both datasets to a suitable metric CRS and add a 20-meter buffer around the road data.
6. Find all trees that lie within this buffer and count them.
7. Present the findings with a map highlighting the roads and the trees.”

6. Experiments

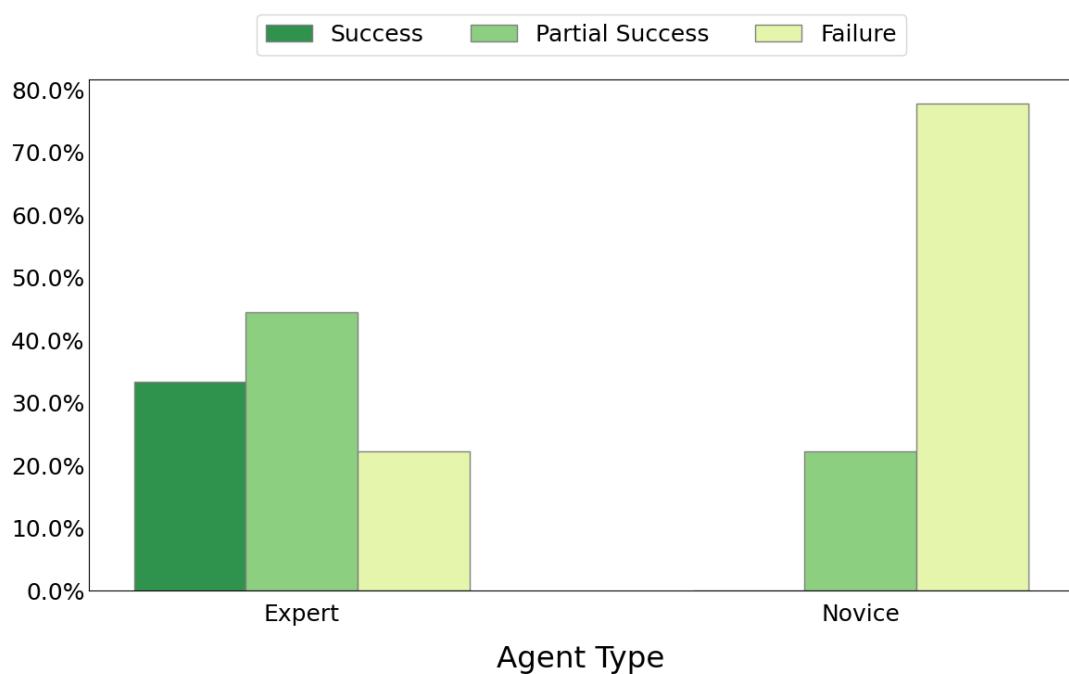


Figure 6.9.: Outcome distribution for different levels of GIS experience

6. Experiments

Using the novice-level prompt GeoGPT was unable to produce the correct outcome, and confidently answered that there are “approximately 6,915 trees on Munkegata street in Trondheim”, which is far from being true. When solving the task, GeoGPT made a series of oversights that lead to this result. First, it failed to take into account that there may be more than one street in the dataset named “Munkegata”, missing the important step of applying a bounding box filter to the query. A similar mistake was made when retrieving the tree data. Due to the upper limit of 10,000 features per query in the API, it’s crucial to narrow down the query to ensure retrieval of all relevant features rather than just a subset. GeoGPT’s query lacked a bounding box, resulting in a randomly distributed subset of 10,000 trees scattered across Norway. A third mistake occurred when GeoGPT calculated a bounding box around the retrieved road data instead of creating a buffer. The latter method would have produced a more accurate result. The bounding box that was created spanned from Trondheim to Oslo, thus including far more trees than was intended.

The expert-level prompt, on the other hand, provided the necessary guidance for GeoGPT for this specific task, steering it clear of the issues it encountered with the novice-level prompt. As OpenAI themselves stated, “some tasks are best specified as a sequence of steps”.⁸ Furthermore, they say that writing explicit steps required to solve a task “makes it easier for the model to follow them”.

⁸<https://platform.openai.com/docs/guides/prompt-engineering/strategy-write-clear-instructions>

6. Experiments

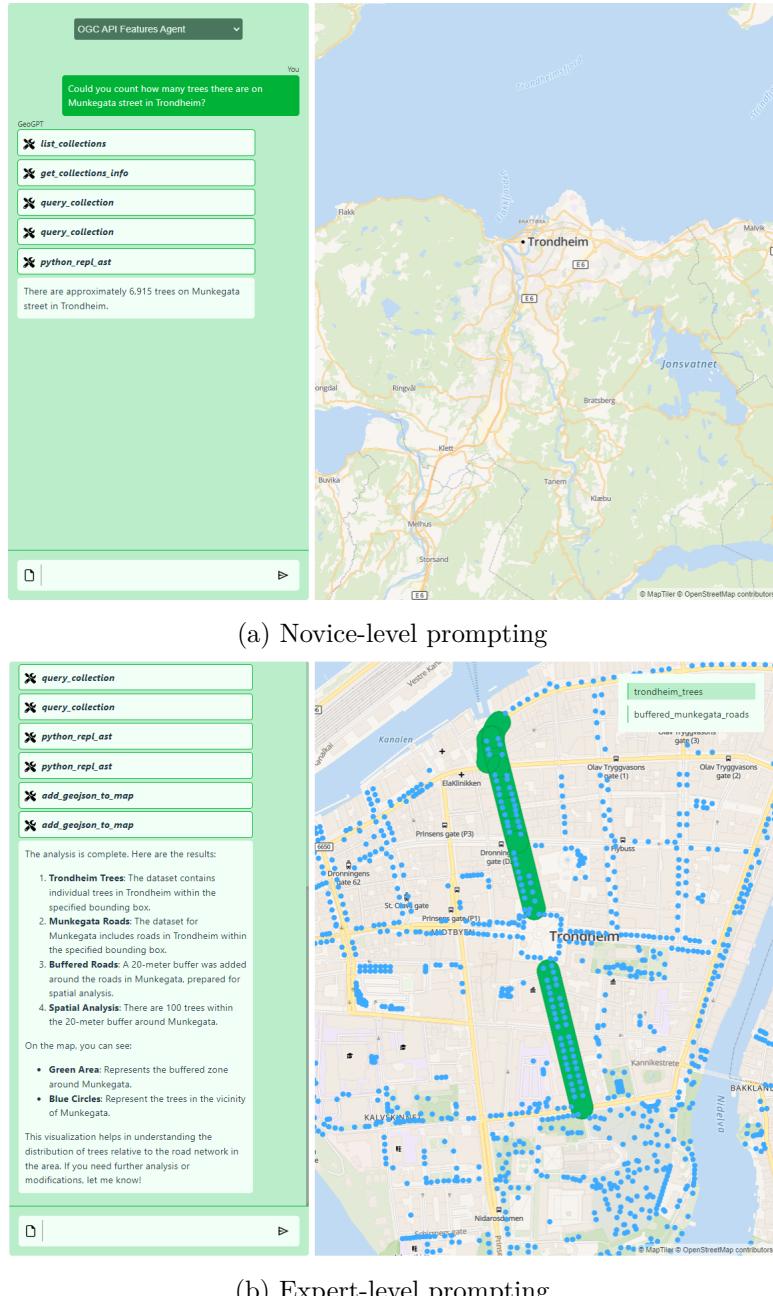


Figure 6.10.: Comparison between novice- and expert-level prompting of GeoGPT's OGC API Features agent for calculation of the number of trees along Munkegata in Trondheim

7. Discussion

Section 7.1 will open the Discussion chapter by suggesting a place for an application like GeoGPT in the field of GIS. Thereafter, section 7.2 will present possible reasons as to why the SQL agent outperforms the other two agent types. Section 7.3 will discuss problem encountered during the development of GeoGPT. Finally, section 7.4 will discuss a semi-failed attempt at implementing a multi-agent architecture for GeoGPT, and provide suggestions as to why this didn't succeed.

7.1. GeoGPT's Place in the Field of GIS

As stated in the section on Goals and Research Questions, one of the goals of this thesis is to see if LLM-based systems can replace GIS professionals. Results from experiments comparing prompts designed to emulate expert and novice users, as detailed in subsection 6.2.2, indicate that the role of GIS professionals is not immediately threatened. As Figure 6.10 shows, expert-level prompting is highly necessary to make GeoGPT consistently produce the desired outcomes for more complex tasks.

It is, however, clear that LLM technologies have significant potential to automate tasks that GIS professionals are commonly faced with. GeoGPT can therefore serve as a helpful companion — much like GitHub Copilot¹ for software developers — that can quickly solve a simple, repetitive tasks, alleviating the workload on GIS professionals. The benchmarking results show that a system like GeoGPT is able to successfully solve a wide range of geospatial task. Note also that many of the questions in the benchmarking dataset (see section A.1) are not very technically phrased, showing that expert-level prompting is not necessary for simpler tasks.

By using a Large Language Model in a GIS context one can also take advantage of their built-in geographical awareness. This geospatial awareness was demonstrated by Roberts et al. (2023) (see section 3.1 for more details). The knowledge within the LLM can therefore help solve tasks when the available data is insufficient to solve a particular task. For instance, the GPT-4 model used in the experiments was able to create quite accurate bounding boxes for several different Norwegian cities, and it was also able to retrieve very accurate coordinates for places like Aker brygge in Oslo. This became useful in one of the tests in benchmark as it occasionally failed to find the point data for Aker brygge in the data provided. Using the LLM's background knowledge for tasks that do not require a high degree of accuracy could prove quite useful.

On the other hand, the results show that the randomness of LLMs gives reason to

¹<https://github.com/features/copilot>

7. Discussion

doubt the answers that GeoGPT produces. Table 6.3 shows GeoGPT’s answers are not always consistent. A common problem with many LLMs is that they will often deliver overly confident answers in response to questions they do not know the answer to. This is problematic, as a user with limited GIS experience will generally not be capable of detecting when GeoGPT generates a believable, but completely false answer. An example of this was seen as GeoGPT was asked which (Norwegian) county is the largest by size (see section A.1). This question was asked a total of nine times across different agents and was answered incorrectly four times (see Table A.2). These incorrect answers typically stated the following:

The largest county by size is **Finnmark**, with an area of approximately **646,150 square kilometers**.

The correct answer to the question, based on the data available to GeoGPT, is “Nordland”, which area is calculated to about 80,5 thousand square kilometres. To the inexperienced user, there is no way of knowing whether this answer can be trusted. An experienced GIS user could, however, inspect the code produced by GeoGPT to see if it makes sense. A GIS professional should also be able to recognise that 646,150 square kilometres is far greater than the actual size of Finnmark.

Lin et al. (2023, pp. 1–2) write that the issue with uncertainty in LLMs is a challenge that “has attracted limited attention until recently”, highlighting the “forbiddingly high” dimensionality of the output space as one of the key hindrances to a reliable way of measuring confidence.

7.2. Why Does the SQL Agent Outperform the Others?

This section will possible reasons for the superior performance of the SQL agent of GeoGPT. The experimental results presented in section 6.2 clearly show that GeoGPT’s SQL agent outperform both the OGC API Features agent and the Python agent. As Figure 6.1 in subsection 6.2.1 shows.

7.2.1. Likely Higher Prevalence of PostGIS Examples During Pre-Training

A possible reason for the performance gap between the SQL and the other two, is the fact that PostGIS is a very established technology, with its first stable release dating back to 2001.² The other agents rely on Python libraries like GeoPandas³ in place of PostGIS, which is a less established technology. As of 23rd May 2024, Google Scholar returns ~ 22.600 results for “PostGIS” compared to ~ 3.650 for “GeoPandas”. This large difference in search results is likely correlated with the respective prevalences of the two technologies in the data that the GPT model used in the experiments was trained upon, which is obtained through web scraping (Radford et al., 2019, p. 3).

²<https://en.wikipedia.org/wiki/PostGIS>

³<https://geopandas.org/en/stable/>

7. Discussion

7.2.2. Limitations with OGC API Features

A common source of error with several of the tests conducted with the OGC API Features agent was its inability to fetch more than 10,000 features from the server. The limit of 10,000 features is specified in the OGC API Features standard (Open Geospatial Consortium, 2022), which states that no more than 10,000 features should be returned in a single response. Accompanied by such a large response, however, should be a `next` link than should point to the next set of 10,000 features. This way, the server could return more than 10,000 features. Unfortunately, as of 23rd May 2024, the current version of `pg_featureserv` does not support this feature,⁴ which is a significant limitation to the current OGC API Features agent in GeoGPT.

Furthermore, the lacking support for multi-collection queries is, in the author's opinion, a big limitation to the current Features specification. A proposal draft to such features have been created,⁵ but it is unclear whether this will be accepted into the specification. This extension, called *Search*, would allow for more complex CQL queries than are not easily specified using query parameters. Code Snippet 7.1 shows one of the multi-collection query examples included in the proposal draft. The ability to construct such queries could make retrieval of features much more efficient, possibly making the Python tool in redundant to GeoGPT's OGC API Features agent. The query in Code Snippet 7.1 would not be possible using the current specification, and one would have to download the two collections, load them into memory using Python, and perform the `contains` operation there. Using multi-collection queries that are converted into PostGIS queries, which are generally more efficient than their Python equivalents, could also speed up analysis.

```
1 \\ SQL query for fetching lakes within Algonquin Park
2 SELECT lakes.*
3 FROM lakes
4 JOIN parks ON ST_Intersects(lakes.geometry, parks.geometry)
5 WHERE parks.name = 'Algonquin Park';
6
7 \\ Corresponding CQL query (would return a tuple of parks and lakes)
8 POST /search    HTTP/1.1
9 Host: www.someserver.com/
10 Accept: application/json
11 Content-Type: application/ogcqry+json
12
13 [
14     {
15         "collections": ["parks", "lakes"]
16         "filter": {
17             "and": [
18                 {"eq": [{"property": "parks.name"}, "Algonquin Park"]},
19                 {"contains": [{"property": "parks.geometry"}},
```

⁴https://github.com/CrunchyData/pg_featureserv/blob/master/FEATURES.md

⁵<https://github.com/opengeospatial/ogcapi-features/tree/master/proposals/search>

7. Discussion

```
20             {"property": "lakes.geometry"}]}  
21         ]  
22     }  
23 }  
24 ]
```

Code Snippet 7.1: Multi-collection CQL query using the *Search* extension

7.3. Where GeoGPT Struggles

This section will present two issues emerged from the Experiments. Subsection 7.3.1 will address an issue that occurs when GeoGPT gets stuck, trying to solve a task by repeatedly making the same, unsuccessful attempts. Subsection 7.3.2 will discuss GeoGPT’s *self-verification* abilities.

7.3.1. Walking Into Dead Ends

As pointed out in subsection 6.2.1 (GIS Benchmark — Results), the correlation matrices in Figure 6.4 show that the encoded outcome has a negative correlation with each of the following variables: token usage, duration, and total cost. This suggests that tasks requiring more time to solve are more likely to result in unsuccessful outcomes. The tendency of GeoGPT of getting stuck in “dead ends” might be an explanation for this. The example in Figure 6.8 shows this tendency well: GeoGPT tries to query the same collection many times in more or less the same way, without realizing that the collection it is querying might not even contain the data it is looking for. This results in an almost endless loop of unsuccessful tool calls, and consequently, a long duration for an unsuccessful outcome.

Peysakhovich and Lerer (2023) suggest, regarding current LLMs, that “relevant information located in earlier context is attended to less on average”. This might explain why GeoGPT occasionally gets stuck attempting the same thing over and over, perhaps because it has forgotten options from earlier in the context window. GeoGPT’s endless responses clearly increase the duration of the run, and also bloats the context window with tool messages. A bloated context window leads to higher token usage and inference time per call to the LLM, which slows the system down even further.

7.3.2. Self-Verification

Li and Ning (2023) introduce five goals for autonomous agents, one of which is *self-verifying*, the system’s ability to test and verify whatever it generates. GeoGPT is already doing self-verification through mid-conversation system messages that verify that a file has been saved to the working directory and that a layer has been added to the map on the client. This does not, however, fix the issue where GeoGPT works with data that is different to what it believes it to be. An example of this issue was presented in the

7. Discussion

section on unsuccessful outcomes in subsection 6.2.1, where GeoGPT believed it was working with data only for Oslo, when in reality the data was for the whole of Norway.

A possible way of doing self-verification that mitigates these kinds of issues is to utilize multi-modal LLMs. Figure 7.1 shows how a multi-modal GPT-4 model can correctly identify that the resulting layer from the above-mentioned unsuccessful analysis is incorrect, based only on a screenshot of the map.

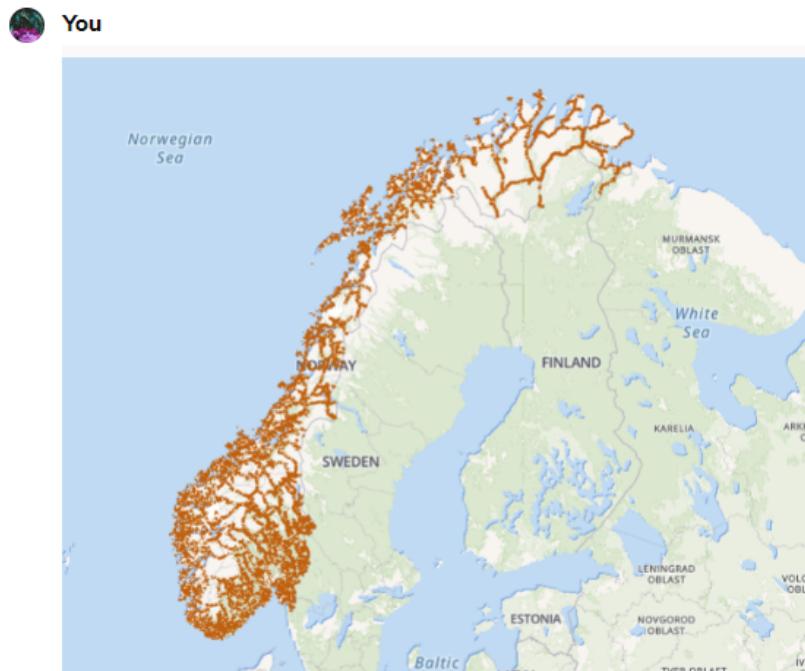
7.4. Multi-Agent Architectures

An attempt was made to create a multi-agent version of GeoGPT that employs three different sub-agents: one for data retrieval, one for data analysis, and one for map interaction. Each agent will have a collection of tools, which are relevant for the sort of tasks they will be asked to solve. These agents are orchestrated by a *Supervisor* agent, which takes input from the user and assigns tasks to the appropriate sub-agents. Figure 7.2 shows how a supervisor node takes a user message (or the chat history up to that point) and selects which sub-agent is to solve the next sub-task. Inspired by Microsoft’s *AutoGen* framework⁶ — which provides a high-level abstraction for developing multi-agent conversations — and MetaGPT (further discussed in section 3.2).

Initial tests on a multi-agent version of the OGC API Features agent were conducted, but it turned out to both increase the latency of the system and be a source of confusion for the LLM. It is possible that such an architecture could be useful as the number of tools available to the agent grows — seeing as a large number of functions tend to start confusing the LLM — but for the current version of GeoGPT the multi-agent pattern seemed to only get in the way.

⁶<https://microsoft.github.io/autogen/>

7. Discussion



Does this look like map with a line layer representing high-speed roads ($\geq 70 \text{ km/h}$) in Oslo?

 ChatGPT

Based on the image you provided, this does not appear to be a map highlighting high-speed roads (roads with a speed limit of 70 km/h or higher) specifically in Oslo, Norway. Instead, the map shows a much broader area, covering Norway, Sweden, Finland, and parts of Russia. The orange lines on the map could be indicative of roads, but the geographic coverage extends far beyond Oslo and includes various terrains and regions across several countries. If the intent was to focus solely on high-speed roads in Oslo, the displayed area would be much more localized and confined to the Oslo region.



Figure 7.1.: ChatGPT multi-modal GPT-4 correctly identifying that the map layer intended to show high-speed roads in Oslo in fact “extends far beyond Oslo”

7. Discussion

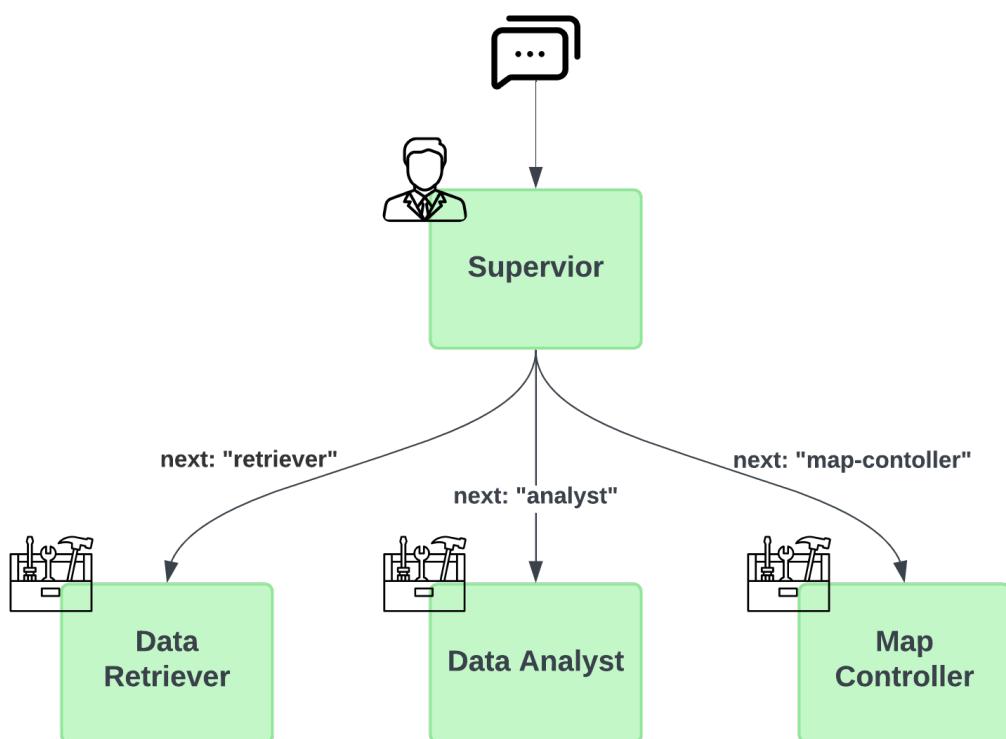


Figure 7.2.: Illustration of how an agent supervisor takes in a user message and selects which sub-agent is to solve the next sub-task

8. Future Work

Sections 8.0.1 through 8.0.3 of this chapter will present challenges that are reserved for future research.

8.0.1. Ability to Answer Questions with no Clear Answer

The experiments conducted in this thesis focused on the technical GIS abilities of the system, and the questions that were asked have corresponding *correct* answers. Something that was not tested is GeoGPT’s ability to answer questions of subjective character, questions that have no *one* correct answer. For instance: what would happen if we asked GeoGPT to find a suitable route from A to B that is as *safe* as possible? How would it interpret such a request? Would it only take into account the speed limit and road type? Would it be able to assess socio-economic aspects between different areas, avoiding “bad neighbourhoods” at nighttime? Would it decide to incorporate weather forecasts into the analysis? Future research should find methods of measuring the ability of LLM-based GIS agents to provide suitable, and safe, answers to such questions.

8.0.2. Comparing Different Models

GeoGPT is based around GPT-4 but, as subsection 2.1.3 showed, there are numerous competitors. Future research should look into the possibilities of swapping out GPT-4 with other models, first and foremost those with good function calling abilities, as this is absolutely necessary in order for GeoGPT to work as intended. A benchmark comparing results for different models would be a natural way of building upon the results of this thesis.

Future research should especially look into the viability of using open-source models. In a report interviewing 500 companies on their LLM adoption, 46 percent stated their preference for open-source models going into 2024 (Wang, Sarah and Xu, Shangda, 2024). *Control* and *customizability* turns out to be the two most important factors into enterprise’s open-source appeal, allowing for increased control over proprietary data and ability to effectively fine-tune models, respectively.

8.0.3. Automated Data Access

The experiments in this thesis were based upon a pre-existing collection of geospatial datasets that were made available to GeoGPT through different channels. A fully autonomous GIS agent should, however, be able to search the web for suitable datasets, based on the user’s query. In a Norwegian context, one could imagine asking for a noise

8. Future Work

analysis for a particular location. The agent should then be able to search a website like Geonorge for datasets related to noise (firing ranges, roads, etc.), downloading these, and then performing analysis. Initial experiments were conducted towards Geonorge in this thesis to see if this was possible, but results were inconsistent.

9. Conclusion

This thesis has shown the viability of using Large Language Model (LLM) technology to create autonomous agents aimed at GIS analysis. GeoGPT, the proposed solution, shows through a new benchmark containing question and answer pairs for common GIS tasks, that it is able to utilize the logical reasoning and code generation abilities of modern LLMs like GPT-4 to solve a wide range of such tasks. The user interacts with GeoGPT through a webpage that consists of a chat interface resembling that of OpenAI's ChatGPT, and a web map where results from analyses can be displayed. The user can type geospatial questions into the chat interface, that GeoGPT will attempt answer to through text and/or by adding geometries to the map.

GeoGPT relies heavily on *function calling*, a way of giving function/tool definitions to an LLM, enabling it to essentially *invoke* these tools through specifying the name of the tool and suitable parameters that will be passed to it. The tools run code created by the developer.

Featuring three different agent types, GeoGPT shows that it can manipulate geospatial data that are discovered in different ways. One agent accesses data directly from a PostGIS database, another agent through an OGC API Features endpoint that lives on top of this PostGIS database, and the final agent by having access to shapefiles in its local environment. All agents have access to the exact same data. The agents have different sets of tools that allow them to work with their data. The SQL agent has (amongst other) a tool that takes a string of SQL code that will be run against the database, enabling it to perform geospatial analyses. The OGC API Features agent and the agent with access to shapefiles have access to a similar tool that allows them to run Python code. In addition to this, the former has access to functions/tools work against the OGC API Features endpoint.

Two sets of experiments were conducted. The first compares the three agent types to see which is best at solving common GIS-related tasks. Restuls from this experiment show that the SQL agent is more likely to produce a desired response compared to the other two agents, with a success rate of 69.4%, compared to 38.9% for the other two agents.

The second set of experiments sought to compare the outcomes of the same GIS question when formulated two in different ways: one simple problem formulation, resembling a user with little GIS experience, and another more accurate and detailed formulation, resembling a user with great GIS experience. The results from the experiment show that providing GeoGPT with a better (more accurate and detailed) initial problem greatly increases the likelihood of the system to produce a successful outcome, suggesting that a user's GIS experience is still very valuable, even as we face a reality where highly

9. Conclusion

sophisticated LLMs can be used to automate away numerous technical tasks.

Bibliography

- Anthropic. (2023). Model Card and Evaluations for Claude Models. <https://www-cdn.anthropic.com/bd2a28d2535bfb0494cc8e2a3bf135d2e7523226/Model-Card-Claude-2.pdf>
- Anthropic. (2024). *The Claude 3 Model Family: Opus, Sonnet, Haiku* (tech. rep.). https://www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model_Card_Claude_3.pdf
- Christiano, P., Leike, J., Brown, T. B., Martic, M., Legg, S., & Amodei, D. (2023). Deep reinforcement learning from human preferences. <https://doi.org/10.48550/arXiv.1706.03741>
- CrunchyData. (2024). CrunchyData/pg_featureserv. Retrieved April 19, 2024, from https://github.com/CrunchyData/pg_featureserv
- Eleti, A., Harris, J., & Kilpatrick, L. (2023). Function calling and other API updates. Retrieved March 10, 2024, from <https://openai.com/blog/function-calling-and-other-api-updates>
- Gemini Team, Anil, R., Borgeaud, S., Alayrac, J.-B., Yu, J., Soricut, R., Schalkwyk, J., Dai, A. M., Hauth, A., Millican, K., Silver, D., Johnson, M., Antonoglou, I., Schrittwieser, J., Glaese, A., Chen, J., Pitler, E., Lillicrap, T., Lazaridou, A., Firat, O., Molloy, J., Isard, M., Barham, P. R., Hennigan, T., Lee, B., Viola, F., Reynolds, M., Xu, Y., Doherty, R., Collins, E., Meyer, C., Rutherford, E., Moreira, E., Ayoub, K., Goel, M., Krawczyk, J., Du, C., Chi, E., Cheng, H.-T., Ni, E., Shah, P., Kane, P., Chan, B., Faruqui, M., Severyn, A., Lin, H., Li, Y., Cheng, Y., Ittycheriah, A., Mahdieh, M., Chen, M., Sun, P., Tran, D., Bagri, S., Lakshminarayanan, B., Liu, J., Orban, A., Güra, F., Zhou, H., Song, X., Boffy, A., Ganapathy, H., Zheng, S., Choe, H., Weisz, Á., Zhu, T., Lu, Y., Gopal, S., Kahn, J., Kula, M., Pitman, J., Shah, R., Taropa, E., Merey, M. A., Baeuml, M., Chen, Z., Shafey, L. E., Zhang, Y., Sercinoglu, O., Tucker, G., Piqueras, E., Krikun, M., Barr, I., Savinov, N., Danihelka, I., Roelofs, B., White, A., Andreassen, A., von Glehn, T., Yagati, L., Kazemi, M., Gonzalez, L., Khalman, M., Sygnowski, J., Frechette, A., Smith, C., Culp, L., Proleev, L., ... Vinyals, O. (2024a). Gemini: A Family of Highly Capable Multimodal Models. <https://doi.org/10.48550/arXiv.2312.11805>

Bibliography

- Gemini Team, Reid, M., Savinov, N., Teplyashin, D., Dmitry, Lepikhin, Lillicrap, T., Alayrac, J.-b., Soricut, R., Lazaridou, A., Firat, O., Schrittweis, J., Antonoglou, I., Anil, R., Borgeaud, S., Dai, A., Millican, K., Dyer, E., Glaese, M., Sotiaux, T., Lee, B., Viola, F., Reynolds, M., Xu, Y., Molloy, J., Chen, J., Isard, M., Barham, P., Hennigan, T., McIlroy, R., Johnson, M., Schalkwyk, J., Collins, E., Rutherford, E., Moreira, E., Ayoub, K., Goel, M., Meyer, C., Thornton, G., Yang, Z., Michalewski, H., Abbas, Z., Schucher, N., Anand, A., Ives, R., Keeling, J., Lenc, K., Haykal, S., Shakeri, S., Shyam, P., Chowdhery, A., Ring, R., Spencer, S., Sezener, E., Vilnis, L., Chang, O., Morioka, N., Tucker, G., Zheng, C., Woodman, O., Attaluri, N., Kociský, T., Eltyshev, E., Chen, X., Chung, T., Selo, V., Brahma, S., Georgiev, P., Slone, A., Zhu, Z., Lottes, J., Qiao, S., Caine, B., Riedel, S., Tomala, A., Chadwick, M., Love, J., Choy, P., Mittal, S., Houlsby, N., Tang, Y., Lamm, M., Bai, L., Zhang, Q., He, L., Cheng, Y., Humphreys, P., Li, Y., Brin, S., Cassirer, A., Miao, Y., Zilka, L., Tobin, T., Xu, K., Proleev, L., Sohn, D., Magni, A., Hendricks, L. A., ... Vinyals, O. (2024b). Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. <https://doi.org/10.48550/arXiv.2403.05530>
- Gemma Team, Mesnard, T., Hardin, C., Dadashi, R., Bhupatiraju, S., Pathak, S., Sifre, L., Rivière, M., Kale, M. S., Love, J., Tafti, P., Hussenot, L., Sessa, P. G., Chowdhery, A., Roberts, A., Barua, A., Botev, A., Castro-Ros, A., Slone, A., Héliou, A., Tacchetti, A., Bulanova, A., Paterson, A., Tsai, B., Shahriari, B., Lan, C. L., Choquette-Choo, C. A., Crepy, C., Cer, D., Ippolito, D., Reid, D., Buchatskaya, E., Ni, E., Noland, E., Yan, G., Tucker, G., Muraru, G.-C., Rozhdestvenskiy, G., Michalewski, H., Tenney, I., Grishchenko, I., Austin, J., Keeling, J., Labanowski, J., Lespiau, J.-B., Stanway, J., Brennan, J., Chen, J., Ferret, J., Chiu, J., Mao-Jones, J., Lee, K., Yu, K., Millican, K., Sjoesund, L. L., Lee, L., Dixon, L., Reid, M., Mikula, M., Wirth, M., Sharman, M., Chinaev, N., Thain, N., Bachem, O., Chang, O., Wahltinez, O., Bailey, P., Michel, P., Yotov, P., Chaabouni, R., Comanescu, R., Jana, R., Anil, R., McIlroy, R., Liu, R., Mullins, R., Smith, S. L., Borgeaud, S., Girgin, S., Douglas, S., Pandya, S., Shakeri, S., De, S., Klimenko, T., Hennigan, T., Feinberg, V., Stokowiec, W., Chen, Y.-h., Ahmed, Z., Gong, Z., Warkentin, T., Peran, L., Giang, M., Farabet, C., Vinyals, O., Dean, J., Kavukcuoglu, K., Hassabis, D., ... Kenealy, K. (2024). Gemma: Open Models Based on Gemini Research and Technology. <https://doi.org/10.48550/arXiv.2403.08295>
- Holm, O. (2023). *LLMs - The Death of GIS Analysis?* (Specialization Project). NTNU. Trondheim. <https://kartai.no/wp-content/uploads/2024/01/Holm-2023-LLMs-The-Death-of-GIS-Analysis.pdf>

Bibliography

- Holmes, C. (2021). SpatioTemporal Asset Catalogs and the Open Geospatial Consortium. Retrieved May 2, 2024, from <https://medium.com/radiant-earth-insights/spatiotemporal-asset-catalogs-and-the-open-geospatial-consortium-659538dce5c7>
- Hong, S., Zhuge, M., Chen, J., Zheng, X., Cheng, Y., Zhang, C., Wang, J., Wang, Z., Yau, S. K. S., Lin, Z., Zhou, L., Ran, C., Xiao, L., Wu, C., & Schmidhuber, J. (2023). MetaGPT: Meta Programming for A Multi-Agent Collaborative Framework. <https://doi.org/10.48550/arXiv.2308.00352>
- Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D. S., de las Casas, D., Bressand, F., Lengyel, G., Lample, G., Saulnier, L., Lavaud, L. R., Lachaux, M.-A., Stock, P., Scao, T. L., Lavril, T., Wang, T., Lacroix, T., & Sayed, W. E. (2023). Mistral 7B. Retrieved May 3, 2024, from <http://arxiv.org/abs/2310.06825>
- Jiang, A. Q., Sablayrolles, A., Roux, A., Mensch, A., Savary, B., Bamford, C., Chaplot, D. S., de las Casas, D., Hanna, E. B., Bressand, F., Lengyel, G., Bour, G., Lample, G., Lavaud, L. R., Saulnier, L., Lachaux, M.-A., Stock, P., Subramanian, S., Yang, S., Antoniak, S., Scao, T. L., Gervet, T., Lavril, T., Wang, T., Lacroix, T., & Sayed, W. E. (2024). Mixtral of Experts. Retrieved May 3, 2024, from <http://arxiv.org/abs/2401.04088>
- LangChain AI. (2022). Langchain-ai/langchain. Retrieved October 5, 2023, from <https://github.com/langchain-ai/langchain>
- LangChain AI. (2024). Langchain-ai/langgraph. Retrieved March 21, 2024, from <https://github.com/langchain-ai/langgraph>
- Li, Z., & Ning, H. (2023). Autonomous GIS: The next-generation AI-powered GIS. <https://doi.org/10.48550/arXiv.2305.06453>
- Lin, Z., Trivedi, S., & Sun, J. (2023). Generating with Confidence: Uncertainty Quantification for Black-box Large Language Models. <https://doi.org/10.48550/arXiv.2305.19187>
- Meta AI. (2024). Introducing Meta Llama 3: The most capable openly available LLM to date. Retrieved May 3, 2024, from <https://ai.meta.com/blog/meta-llama-3/>
- Mooney, P., Cui, W., Guan, B., & Juhász, L. (2023). *Towards Understanding the Geospatial Skills of ChatGPT: Taking a Geographic Information Systems (GIS) Exam.* <https://doi.org/10.1145/3615886.3627745>
- Moura, J. (2024). Joaomdmoura/crewAI. Retrieved May 21, 2024, from <https://github.com/joaomdmoura/crewAI>
- Open Geospatial Consortium. (2022). OGC API - Features - Part 1: Core corrigendum. Retrieved April 29, 2024, from <https://docs.ogc.org/is/17-069r4/17-069r4.html>

Bibliography

- OpenAI. (2022). Introducing ChatGPT. Retrieved October 26, 2023, from <https://openai.com/blog/chatgpt>
- Peysakhovich, A., & Lerer, A. (2023). Attention Sorting Combats Recency Bias In Long Context Language Models. Retrieved May 14, 2024, from <http://arxiv.org/abs/2310.01427>
- Pichai, S., & Hassabis, D. (2024). Our next-generation model: Gemini 1.5. Retrieved May 3, 2024, from <https://blog.google/technology/ai/google-gemini-next-generation-model-february-2024/>
- PostGIS. (2001). Retrieved April 29, 2024, from <https://postgis.net/>
- Radford, A., Narasimhan, K., Salimans, T., & Sutskever, I. (2018). Improving Language Understanding by Generative Pre-Training.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019). Language Models are Unsupervised Multitask Learners.
- Richard, T. B. (2023). AutoGPT: The heart of the open-source agent ecosystem. Retrieved October 5, 2023, from <https://github.com/Significant-Gravitas/AutoGPT>
- Roberts, J., Lüddecke, T., Das, S., Han, K., & Albanie, S. (2023). GPT4GEO: How a Language Model Sees the World's Geography. <https://doi.org/10.48550/arXiv.2306.00020>
- Sanfilippo, S. (2009). Redis - The Real-time Data Platform. Retrieved April 19, 2024, from <https://redis.io/>
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., & Klimov, O. (2017). Proximal Policy Optimization Algorithms. <https://doi.org/10.48550/arXiv.1707.06347>
- Shinn, N., Cassano, F., Berman, E., Gopinath, A., Narasimhan, K., & Yao, S. (2023). Reflexion: Language Agents with Verbal Reinforcement Learning. Retrieved May 13, 2024, from <http://arxiv.org/abs/2303.11366>
- Su, J., Ahmed, M., Lu, Y., Pan, S., Bo, W., & Liu, Y. (2024). RoFormer: Enhanced transformer with Rotary Position Embedding. *Neurocomputing*, 568, 127063. <https://doi.org/10.1016/j.neucom.2023.127063>
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., & Polosukhin, I. (2017). Attention Is All You Need. Retrieved October 10, 2023, from <https://arxiv.org/abs/1706.03762v7>
- Wang, Sarah & Xu, Shangda. (2024). 16 Changes to the Way Enterprises Are Building and Buying Generative AI. Retrieved May 9, 2024, from <https://a16z.com/generative-ai-enterprise-2024/>

Bibliography

- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E. H., Le, Q. V., & Zhou, D. (2023). Chain-of-Thought Prompting Elicits Reasoning in Large Language Models.
- White, J., Fu, Q., Hays, S., Sandborn, M., Olea, C., Gilbert, H., Elnashar, A., Spencer-Smith, J., & Schmidt, D. C. (2023). A Prompt Pattern Catalog to Enhance Prompt Engineering with ChatGPT. Retrieved May 6, 2024, from <http://arxiv.org/abs/2302.11382>
- Wu, Q., Bansal, G., Zhang, J., Wu, Y., Li, B., Zhu, E., Jiang, L., Zhang, X., Zhang, S., Liu, J., Awadallah, A. H., White, R. W., Burger, D., & Wang, C. (2023). AutoGen: Enabling Next-Gen LLM Applications via Multi-Agent Conversation. <https://doi.org/10.48550/arXiv.2308.08155>
- Yan, Fanjia, Mao, Huanzhi, Ji, Charlie Cheng-Jie, Stoica, Ion, Gonzalez, Joseph E., Zhang, Tianjun & Patil, Shishir G. (2024). Berkeley Function Calling Leaderboard. Retrieved May 9, 2024, from https://gorilla.cs.berkeley.edu/blogs/8_berkeley_function_calling_leaderboard.html#metrics
- Zhang, Y., Wei, C., Wu, S., He, Z., & Yu, W. (2023). GeoGPT: Understanding and Processing Geospatial Tasks through An Autonomous GPT. <https://doi.org/10.48550/arXiv.2307.07930>
- Zhu, D., Yang, N., Wang, L., Song, Y., Wu, W., Wei, F., & Li, S. (2024). PoSE: Efficient Context Window Extension of LLMs via Positional Skip-wise Training. <https://doi.org/10.48550/arXiv.2309.10400>

Appendices

A. Experiments

A.1. Quantitative

Table A.1.: Questions for quantitative tests

Query ID	Query	Correct Response
aker_brygge_national	Which is the closest railway station to Aker brygge?	Nationalteateret
cliff_clusters	Locate clusters of cliffs, each containing more than 10 cliffs, with cliffs within each cluster no more than 0.1 degrees apart.	Should be about 8-9 clusters.
county_names	What are the names of the counties found in the data?	Nordland, Telemark, Troms, Rogaland, Vestland, Trøndelag, Vestfold, Buskerud, Akershus, Østfold, Innlandet, Møre og Romsdal, Agder Finnmark
glomma_counties	How many counties does Glomma run through?	4. Trøndelag, Innlandet, Akershus, and Østfold.
largest_county	Which is the largest county by size?	Nordland
nidarosdomen_polygon	Retrieve a polygon of Nidarosdomen.	Adding a polygon of Nidarosdomen to the map.
num_trees_munkegata	How many trees are there along Munkegata in Trondheim?	Giving the correct number of trees (about 70-80).

Continued on next page

A. Experiments

Table A.1.: Questions for quantitative tests

Query ID	Query	Correct Response
oslo_ bergen_ geodesic	Create a geodesic curve between the airports of Oslo and Bergen.	A geodesic, slightly curved line between Gardermoen and Flesland.
oslo_ residential_ diff	Provide an outline of Oslo but exclude residential areas by computing their difference.	The polygonal outline of Oslo with cutouts where there area areas classified as residential.
oslo_ roads_ gte_ 70_- kmh	Retrieve roads in Oslo that have speed limit higher than or equal to 70 km/h.	Adding corresponding line segments to the map.
vestfold_ bbox	Find the coordinates that define a the extent of a rectangular bounding box that encompasses the county of Vestfold.	(9.7553357, 58.720455, 10.6750198, 59.674011)
viken_ dissolve	Combine the countries that correspond to the historical country of Viken into a single feature.	Display the outline of Viken, based on the outlines of Buskerud, Akershus, and Østfold, that are found in the data. Should be dissolved, i.e., no borders.

A. Experiments

Table A.2.: Test results for quantitative tests

Query ID	Agent Type	Outcome	Duration [s]	Tokens
aker_brygge_national	oaf	partial success	87.16	3347
aker_brygge_national	oaf	partial success	97.85	3042
aker_brygge_national	oaf	partial success	70.36	4563
aker_brygge_national	python	partial success	68.10	4736
aker_brygge_national	python	partial success	49.48	2192
aker_brygge_national	python	partial success	76.84	4270
aker_brygge_national	sql	success	63.51	2569
aker_brygge_national	sql	success	151.46	6273
aker_brygge_national	sql	success	103.98	5169
cliff_clusters	oaf	partial success	93.21	3047
cliff_clusters	oaf	partial success	134.62	3713
cliff_clusters	oaf	failure	253.77	6956
cliff_clusters	python	failure	66.42	5421
cliff_clusters	python	failure	125.96	4650
cliff_clusters	python	failure	109.29	4768
cliff_clusters	sql	failure	39.36	1633
cliff_clusters	sql	failure	75.27	3085
cliff_clusters	sql	success	31.64	1632

Continued on next page

A. Experiments

Table A.2.: Test results for quantitative tests

Query ID	Agent Type	Outcome	Duration [s]	Tokens
county_names	oaf	success	78.08	2389
county_names	oaf	success	67.45	2399
county_names	oaf	failure	20.23	1471
county_names	python	success	44.33	2532
county_names	python	success	30.50	2355
county_names	python	success	31.09	2355
county_names	sql	success	64.21	1874
county_names	sql	success	52.28	1880
county_names	sql	success	47.62	1886
glomma_counties	oaf	success	69.73	2253
glomma_counties	oaf	success	65.57	1945
glomma_counties	oaf	success	67.37	2397
glomma_counties	python	success	663.22	3658
glomma_counties	python	failure	295.93	4388
glomma_counties	python	success	285.23	2281
glomma_counties	sql	success	102.48	1828
glomma_counties	sql	success	25.01	1266
glomma_counties	sql	success	21.36	1064
largest_county	oaf	success	41.82	1713
largest_county	oaf	success	71.82	2223
largest_county	oaf	failure	47.00	1797
largest_county	python	failure	29.37	2415
largest_county	python	failure	41.95	2810
largest_county	python	failure	40.55	2407

Continued on next page

A. Experiments

Table A.2.: Test results for quantitative tests

Query ID	Agent Type	Outcome	Duration [s]	Tokens
largest_county	sql	success	40.08	1771
largest_county	sql	success	28.51	1474
largest_county	sql	success	38.30	1381
nidarosdomen_polygon	oaf	success	29.82	1923
nidarosdomen_polygon	oaf	success	31.73	1909
nidarosdomen_polygon	oaf	success	31.61	1922
nidarosdomen_polygon	python	success	35.10	4510
nidarosdomen_polygon	python	success	25.36	3649
nidarosdomen_polygon	python	success	33.35	3663
nidarosdomen_polygon	sql	success	29.64	1546
nidarosdomen_polygon	sql	success	90.61	1529
nidarosdomen_polygon	sql	success	39.39	1513
num_trees_munkegata	oaf	failure	75.60	2977
num_trees_munkegata	oaf	failure	80.13	2984
num_trees_munkegata	oaf	failure	57.66	2407
num_trees_munkegata	python	partial success	122.99	3930
num_trees_munkegata	python	failure	916.13	2362
num_trees_munkegata	python	failure	911.18	8385
num_trees_munkegata	sql	failure	80.41	1359
num_trees_munkegata	sql	failure	67.41	2626
num_trees_munkegata	sql	failure	31.86	1359
oslo_bergen_geodesic	oaf	failure	652.88	9253
oslo_bergen_geodesic	oaf	failure	102.38	3360

Continued on next page

A. Experiments

Table A.2.: Test results for quantitative tests

Query ID	Agent Type	Outcome	Duration [s]	Tokens
oslo_bergen_geodesic	oaf	failure	258.52	6271
oslo_bergen_geodesic	python	partial success	118.54	6683
oslo_bergen_geodesic	python	partial success	129.40	7782
oslo_bergen_geodesic	python	failure	373.37	10412
oslo_bergen_geodesic	sql	failure	97.79	2710
oslo_bergen_geodesic	sql	failure	73.78	2575
oslo_bergen_geodesic	sql	failure	127.43	3318
oslo_residential_diff	oaf	failure	119.39	2373
oslo_residential_diff	oaf	partial success	114.92	3033
oslo_residential_diff	oaf	success	121.17	2621
oslo_residential_diff	python	success	673.96	4013
oslo_residential_diff	python	failure	2190.08	3867
oslo_residential_diff	python	failure	2290.09	3650
oslo_residential_diff	sql	success	42.49	1537
oslo_residential_diff	sql	failure	94.35	2485
oslo_residential_diff	sql	success	52.14	1790
oslo_roads_gte_70_kmh	oaf	partial success	44.44	1959
oslo_roads_gte_70_kmh	oaf	failure	58.61	1989
oslo_roads_gte_70_kmh	oaf	partial success	35.55	1961
oslo_roads_gte_70_kmh	python	failure	983.76	3912

Continued on next page

A. Experiments

Table A.2.: Test results for quantitative tests

Query ID	Agent Type	Outcome	Duration [s]	Tokens
oslo_roads_gte_70_kmh	python	failure	757.70	3819
oslo_roads_gte_70_kmh	python	failure	849.06	3827
oslo_roads_gte_70_kmh	sql	success	43.88	1542
oslo_roads_gte_70_kmh	sql	success	37.26	1577
oslo_roads_gte_70_kmh	sql	success	38.38	1566
vestfold_bbox	oaf	failure	56.33	2001
vestfold_bbox	oaf	partial success	50.38	3177
vestfold_bbox	oaf	success	55.82	2480
vestfold_bbox	python	success	24.08	1612
vestfold_bbox	python	success	33.63	2366
vestfold_bbox	python	success	42.36	2808
vestfold_bbox	sql	success	23.76	1188
vestfold_bbox	sql	success	47.35	1932
vestfold_bbox	sql	success	23.30	1182
viken_dissolve	oaf	success	52.24	2336
viken_dissolve	oaf	success	204.07	5078
viken_dissolve	oaf	failure	133.40	3985
viken_dissolve	python	success	68.23	6099
viken_dissolve	python	success	86.32	6007
viken_dissolve	python	partial success	89.98	6461
viken_dissolve	sql	partial success	34.58	1553
viken_dissolve	sql	success	34.72	1540

Continued on next page

A. Experiments

Table A.2.: Test results for quantitative tests

Query ID	Agent Type	Outcome	Duration [s]	Tokens
viken_dissolve	sql	partial success	37.36	1543

A. Experiments

A.2. Experience Levels

Table A.3.: Questions for experience level tests

Query ID	Level	Formulation
oslo_ bergen_ geodesic	novice	Please plot the shortest flight path on a map between Oslo and Bergen's airports.
oslo_ bergen_ geodesic	expert	1. Get info on all potentially relevant datasets. The airports are possibly stored as polygons in the available data. 2. Filter those datasets for airports and list the names of the airports. 3. Retrieve the geographic coordinates for Oslo Gardermoen Airport (OSL) and Bergen Flesland Airport (BGO) by filtering on the names you found. 4. If no name was found, try a different dataset and go back to step 2. 5. Utilize available tools to draw a geodesic curve that represents the shortest path on the earth's surface between these two points. 6. Present the findings with a map highlighting the largest county.
oslo_ roads_ gte_ 70_ kmh	novice	Draw roads in Oslo where you can drive at least 70.
oslo_ roads_ gte_ 70_ kmh	expert	1. Retrieve an outline of Oslo. 2. Calculate max/min lat/lon values for this bounding box, and use it to retrieve a subset of the road data. 3. Select road segments within the outline from step 1 that have a max speed ≥ 70 . 4. Present the findings with a map highlighting the selected roads.
num_ trees_ munkegata	novice	Could you count how many trees there are on Munkegata street in Trondheim?

Continued on next page

A. Experiments

Table A.3.: Questions for experience level tests

Query ID	Level	Formulation
num_trees_munkegata	expert	<ol style="list-style-type: none">1. List all datasets that could possibly include trees.2. Find the correct feature class and filter the relevant dataset to access tree data for Trondheim. Use a bounding box to reduce the number of trees to analyse.3. Fetch road data for Munkegata. Use a bounding box for Trondheim in case there are streets elsewhere named Munkegata.4. Convert both datasets to a suitable metric CRS and add a 20-meter buffer around the road data.5. Find all trees that lie within this buffer and count them.6. Present the findings with a map highlighting the roads and the trees.

A. Experiments

Table A.4.: Test results for experience level tests

Query ID	Agent Type	Level	Outcome
oslo_bergen_geodesic	sql	novice	failure
oslo_bergen_geodesic	oaf	novice	failure
oslo_bergen_geodesic	python	novice	failure
oslo_bergen_geodesic	sql	expert	failure
oslo_bergen_geodesic	oaf	expert	failure
oslo_bergen_geodesic	python	expert	partial success
oslo_roads_gte_70_kmh	sql	novice	partial success
oslo_roads_gte_70_kmh	oaf	novice	partial success
oslo_roads_gte_70_kmh	python	novice	failure
oslo_roads_gte_70_kmh	sql	expert	success
oslo_roads_gte_70_kmh	oaf	expert	partial success
oslo_roads_gte_70_kmh	python	expert	success
num_trees_munkegata	sql	novice	failure
num_trees_munkegata	oaf	novice	failure
num_trees_munkegata	python	novice	failure
num_trees_munkegata	sql	expert	partial success
num_trees_munkegata	oaf	expert	success
num_trees_munkegata	python	expert	partial success