

Karl Oskar Magnus Holm

LLMs - The Death of GIS Analysis?

An Investigation into Using Large Language Models for GIS Data Analysis

Master Thesis in Computer Science and Geomatics, June 2024

Supervisor at NTNU: Hongchao Fan

External supervisors from Norkart: Alexander Salveson Nossun, Arild Nomeland, and
Rune Aasgaard

Department of Geomatics

Faculty of Engineering

Norwegian University of Science and Technology



Abstract

List of Figures

5.1. Architecture overview	10
--------------------------------------	----

List of Tables

4.1. Datasets used in experiments	8
---	---

1. Introduction

The introductory chapter will explain the motivation behind the thesis, as well as its goals and the research questions it will attempt to answer. Section 8.1 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, 2023 generated a hype within the general population and chat-based systems are flourishing. Such Large Language Models (LLMs) provide a natural interface between human and machine, one in which technical know-how is not necessarily required in order to perform technical tasks, as modern LLMs are often able to interpret the user’s intent regardless of the preciseness of their prompt.

ref

GIS analysis has traditionally been reserved for GIS *experts*.

1.2. Goals and Research Questions

1. Can an LLM-based system answer advanced GIS problems?
2. What are state-of-the-art methods of creating autonomous LLM-based agents?
3. What are core challenges in developing such systems?

1.3. Research Method

The specialization project which this thesis is based upon (Holm, 2023) was of a theoretical character, though it was supplemented by some simple experiments based around direct prompting with ChatGPT. This master thesis will, however, be more practical. As explained in section 1.2, the overarching goal is to create a “proof of concept”. The resulting prototype will then be used to answer the research questions.

1.4. Contributions

1.5. Thesis Structure

2. Background Theory

Chapter 2 will lay a theoretical basis for the work done in this master thesis, providing the user with the required understanding in order to understand the contributions of the work. Section 2.1 will explain the theoretical basis of the component which most modern Large Language Models (LLMs) are based upon — namely the Transformer — and the attention mechanism within it. The section will also touch upon a new architecture aimed at text generation called *state-space*, which has yielded very promising results for small LLMs.

check
this

Parts of the Background chapter is reused material from the specialization project (Holm, 2023) preceding this master thesis. Below are the sections in question, together with a description of the extent to which, and how, the material is reused:

- *Subsection 2.1.1: Reused without modification.*
- *Subsection 2.3.1: Reused without modification.*

2.1. The Theory Behind Generative Large Language Models

Subsection 2.1.1 will explain the theoretical basis behind most modern LLMs, which are based upon the attention mechanism built into the Transformer architecture. Subsection 2.1.2 will explain the new state-space architecture.

2.1.1. 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 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 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 commonly found with RNNs. An attention function maps a query and a set of key-value

2. Background Theory

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. We improve the attention mechanism by multiplying queries, keys, and values with weight matrices that are learned 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) (I will use machine translation from Norwegian to German as an example):

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 Natural Language Processing (NLP), and is the fundamental building block of modern LLMs, most famous of which are the GPT’s.

2.1.2. State-Space

2.2. Function Calling LLMs

Function calling—first introduced by OpenAI (Eletti et al., 2023)—allows developers to provide function definitions to an LLM and have said LLM output a JSON object containing the name of one or more of the functions provided, as well as suitable arguments to these. Made possible through fine-tuning models to detect when functions should be calling, 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.

Possible use cases include using functions 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

2. Background Theory

possible to trace a claim back to its source, something that is normally a difficult feat with LLM. Another use case might be code execution. One could imagine a rather simple function `execute_python_code(code: string) -> string` that takes Python code as a string and returns the standard output that results from executing that code. This is likely the principle behind products like OpenAI’s Data Analysis mode (previously Code Interpreter), in which ChatGPT functions 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. As Eleti et al. (2023) describes, function calling can also be used to extract structured data from text.

2.3. State-of-the-Art Large Language Models

2.3.1. The GPT Family

Generative Pre-trained Transformer (GPT) is a type of LLM that was introduced by OpenAI in 2018 (Radford and Narasimhan, 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 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, just trying to elaborate upon the input sequence it’s given to the best of its knowledge. This will naturally produce undefined behaviour, and it is therefore necessary to fine-tune the model on target tasks in a supervised manner. Radford and Narasimhan (2018, p. 4) explain how the model can be fine-tuned directly on tasks like text classification, but how one for other tasks needs to convert structured inputs into ordered sequences 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) 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), 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).

2. Background Theory

2.3.2. The Gemini Family

2.3.3. The Claude Family

2.3.4. Open-Source Alternatives and Honourable Mentions

3. Related Work

4. Datasets

With the interest of investigating the ability of a Large Language Model (LLM)-based system to perform geospatial analysis, relevant datasets should be accessible to said system. Section 4.1 provides a description of the datasets used in the experiments. Furthermore, it was decided to explore different access channels to this data. Section 4.2 elaborates on this.

4.1. Data Sources

A total of eight datasets were used in the experiments. All datasets were downloaded through Geonorge¹, a webpage administered by The Norwegian Mapping Authority that serves as a portal to a large catalogue of Norwegian geographical data. The datasets were selected to provide a diverse pool of geographical data to perform analysis upon. Table 4.1 shows the datasets used along with a description.

4.2. Data Access

While leading LLMs are trained on increasingly large corpora, they are still only as familiar with a topic as the extent to which the training data exposes it to said topic. For instance, many LLMs are trained specifically to generate Python code, and are therefore fed with a vast number of Python code examples during training in the hopes of improving its performance on benchmarks like . As it is unlikely that the training data is evenly distributed among many different topics, it is useful to get familiarized with a model's capabilities in the areas of interest for a particular use case. In the case of an LLM-powered GIS agent that should be capable of performing geospatial analyses, it is useful to know what data formats such an agent is most comfortable to understand and work with.

The upcoming experiments therefore seek to benchmark model performance on three different data access methods. The datasets from section 4.1 are presented to the model in three different ways, as subsection 4.2.1 through subsection 4.2.2 elaborate upon.

4.2.1. Files

The first method of presentation is to have the files from section 4.1 remain untouched. The datasets were stored such that each dataset has its own folder. This is because some

¹<https://geonorge.no>

Insert
Python
bench-
mark
exam-
ples
here

4. Datasets

Table 4.1.: Datasets used in experiments

Dataset	Description
AR50 - Land Use Map	A nationwide dataset that displays main types of land use adapted for use in scales from 1:20,000 to 1:100,000.
Strategic Noise Mapping from Road Traffic	The strategic noise mapping shows the noise situation from road traffic at the turn of the year 2016/2017 for the largest urban areas in the country and in addition along national and county roads where more than 8,200 vehicles pass per day.
Elveg 2.0 - Road Network	A road network dataset that includes all drivable roads that are longer than 50 meters, or part of a network, as well as pedestrian and bicycle paths and bicycle paths represented as road link geometry.
Cadastral Register - Building Points	The Matrikkelen-Building point dataset contains a small excerpt of the building information registered in the Matrikkelen, Norway's official register of real property, including buildings. The dataset contains representation points, building type, building number, current building status.
Outdoor Recreation Areas	The purpose of the dataset is to provide an overview of areas that are important for the public's outdoor life, and it should be easy to account for which assessments and criteria have been the basis for the work and the final product.
Cultural Monuments - Protected Buildings	Buildings and churches that are automatically, decision, regulation, or temporarily protected under law and churches that have the status as listed.
Flood Zones	Flood zones show areas that are flooded by different flood sizes (recurrence interval). Flood zones are prepared for 20-, 200-, and 1000-year floods.
Quick Clay Zones	Provides an overview of zones with potential danger (precautionary areas) for major quick clay landslides.

4. Datasets

of the file types used require multiple files in order to correctly store the data—for instance, the shapefile format, which has three mandatory files: .shp, which contains the actual feature geometry; .shx, which provides a positional index of the feature geometry; and .dbf, which holds attributes for each shape.

reference

4.2.2. SQL Database

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.

Some of the datasets come in the GML data format, which can include multiple layers with potentially different geometries. For this reason, they cannot be loaded directly into a PostGIS database such that they are stored in the same database table. Furthermore, several of the layers in the multi-layer GML files are irrelevant for most analysis situations. For instance, the flood zone data were downloaded as a multi-layer GML file and includes a total of eight layers: polygon and multi-line border for the analysis area, polygon layers for rivers, ocean surfaces, and lakes, polygon and multi-line border for the flood zones, and a layer containing cross-sectional profile lines for the rivers. The quick clay dataset was similar. For brevity, a decision was made not to load all these layers into the PostGIS database. Only the polygon for the flood zones and two polygon layers from the quick clay dataset were loaded into the database.

4.2.3. OGC API Features

The third method for data access is to use the OGC API Features standard.

5. Architecture

5.1. High-Level Application Architecture

A microservice architecture was employed in order to simplify development and separate concerns between the different microservices. The services are deployed as Docker Containers, and they are orchestrated using Docker Compose. Figure 5.1 shows how the application is divided into five distinct services.

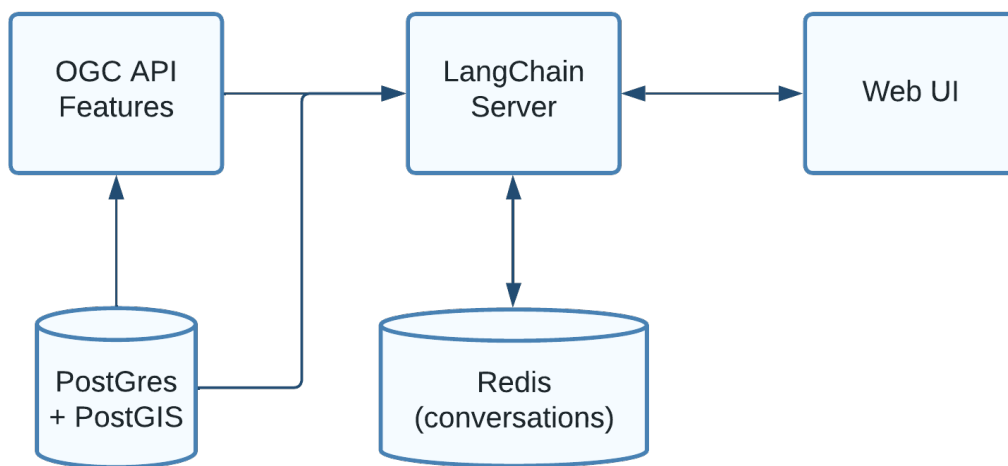


Figure 5.1.: Architecture overview

The *LangChain server* is the heart of the application, and is where the Large Language Model (LLM)-related logic is situated.

5.2. LangChain Server

LangChain (Chase, 2022) is an open-source project that provides tooling that can be used to create autonomous AI agents. It is designed to help with prompt management and optimization, creating chains of calls to LLMs, data-augmented generation, autonomous agent creation, and memory-related tasks.

LangGraph is nice.

5.3. Redis for Conversations

5.4. PostGres + PostGIS

5.5. OGC API Features

5.6. Web UI

6. Experiments and Results

6.1. Experimental Plan

In order to evaluate GeoGPT’s ability to perform geospatial analyses, a Q&A dataset was constructed. This dataset consists of geospatial questions with corresponding correct answers. For each sample of the dataset, a description of how a human would find it natural to approach the problem. This description is provided as step-by-step path towards the solution, and is only included to guide any reader of this thesis as to how the system would be expected to solve the system. Only the questions and their answers are used in the actual evaluation. The full Q&A dataset can be found in Appendix B.

Each question is labelled with a difficulty; one of $\{1, 2, 3\}$, with 3 being the most difficult. By having questions of varying difficulties, it will be easier to estimate the system’s strength. There are X questions for each difficulty.

update

6.2. Experimental Setup

6.3. Experimental Results

7. Evaluation and Discussion

7.1. Evaluation

7.2. Discussion

8. Conclusion and Future Work

8.1. Contributions

8.2. Future Work

8.2.1. Automated Data Access

The experiments in chapter 6 were based upon a pre-existing database. 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 analysis for a particular building. The agent would then search Geonorge for datasets related to noise (firing ranges, roads, etc.), downloading these, and then performing analysis. Simple experiments were conducted in this thesis to see if this was possible, but results were somewhat poor. Methods like semantic search based upon the documentations of datasets should be explored in future research.

Bibliography

- Chase, H. (2022). LangChain. Retrieved October 5, 2023, from <https://github.com/langchain-ai/langchain>
- 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>
- Holm, O. (2023). *LLMs - The Death of GIS Analysis?* (Specialization Project). NTNU. Trondheim. <https://github.com/oskarhlm/prosjektoppgave/blob/main/tex/auxiliary/main.pdf>
- OpenAI. (2022). Introducing ChatGPT. Retrieved October 26, 2023, from <https://openai.com/blog/chatgpt>
- Radford, A., & Narasimhan, K. (2018). Improving Language Understanding by Generative Pre-Training. Retrieved October 9, 2023, from <https://www.semanticscholar.org/paper/Improving-Language-Understanding-by-Generative-Radford-Narasimhan/cd18800a0fe0b668a1cc19f2ec95b5003d0a5035>
- 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>

Appendices

A. Task Description from Norkart

Oppgave med omfang som kan tilpassast både prosjekt og masteroppgave

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 spesialtilpasse 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

A. Task Description from Norkart

Detaljert oppgavebeskrivelse utvikles i samarbeid med studenten.

ADMINISTRATIVT/VEILEDNING

Ekstern veileder: (en eller flere)

Mathilde Ørstavik, Norkart

Rune Aasgaard, Norkart

Alexander Nossun, Norkart

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

Terje Midtbø (GIS, kartografi, visualisering)

Hongchao Fan (3D modellering, fotogrammetri, laser)

B. Q&As for Benchmark

Question	Difficulty	Answer	Reasoning
Find all buildings that are in the danger zone of a being flooded, in accordance to current regulations for new buildings.	2	Dataset containing building points: ...\\data\\bygninger_100årsflom.shp	Steps: <ol style="list-style-type: none">1. Realize that new buildings should be secure against 200-year floods.2. Gather and load flood risk data and building point data.3. Ensure that both datasets are in the same coordinate reference system (CRS).4. Filter out the buildings that are located within the 200-year flood zone.

B. Q&As for Benchmark

Find all buildings that are in the danger zone of a being flooded, in accordance to current regulations for new buildings.	2	Dataset containing building points: ...\\data\\bygninger__-100årsflood.shp	Steps: <ol style="list-style-type: none">1. Realize that new buildings should be secure against 200-year floods.2. Gather and load flood risk data and building point data.3. Ensure that both datasets are in the same coordinate reference system (CRS).4. Filter out the buildings that are located within the 200-year flood zone.
--	---	--	---

B. Q&As for Benchmark