# Data exploration: large language models in the construction of Knowledge Graphs

Cecília F. V. Couto[1], Nelson F. F. Ebecken[1]

[1]*Dept. of Civil Engineering, COPPE, Federal University of Rio de Janeiro*
*Av. Athos da Silveira Ramos, 149, CT, Room B101, Rio de Janeiro, RJ 21941-909, Brazil*
*cecilia.couto@coc.ufrj.com, nelson@ntt.ufrj.br*

**Abstract.** Knowledge graphs (KGs) are graphical representations of structured information that illustrate the relationships between concepts, entities, or data. KGs play a crucial role in enhancing the performance of artificial intelligence systems and search tools. However, constructing knowledge graphs is a complex undertaking, requiring the assimilation of substantial amounts of data. One approach to building KGs involves using Large Language Models (LLMs), which leverage artificial intelligence to comprehend and generate natural language. This study advocates for the use of LLM models in KG construction. To achieve this, the models Llama 2 7B, Llama 3 8B, ChatGPT 3.5, and ChatGPT 4 were employed.

**Keywords:** Knowledge Graphs, Llama 2, Llama 3, ChatGPT 3.5, ChatGPT 4

## 1 Introduction

Throughout human history, a vast corpus of knowledge has been generated, a substantial portion of which now resides online. However, owing to the sheer volume of this knowledge, extracting information and delineating interconnections between subjects proves to be a labor-intensive endeavor. In this regard, an alternative approach to organize and streamline access to this wealth of information is through knowledge graphs.

Knowledge graphs (KGs) are data structures that represent information, concepts, and relationships between entities within a specific domain. KGs are used to organize large amounts of knowledge in order to facilitate understanding and visualization of information [1]. KGs find applications in various tools such as search engines, question-answering systems, and recommendation systems [2].

KGs consist of nodes, which encapsulate concepts, entities, data, etc., and edges, which denote the interconnections between these nodes. They embody a substantial reservoir of antecedent knowledge, yet afford the validation of novel relationships among entities [3].

However, creating KGs is not a simple task, as it involves a difficult and expensive process of several steps, annotations and/or human guidance [4]. Therefore, a way to generate KGs in an easier, faster and cheaper way that has been gaining prominence lately is the use of Large Language Models (LLMs) for such projects.

LLMs are systems that use artificial intelligence to understand and create human language in a similar way to how humans do it. These models undergo rigorous training on extensive datasets to internalize linguistic patterns and structures, thus enabling proficient knowledge extraction [5]. LLMs have diverse applications, including content or code generation, summarization, conversation, and creative writing [6]. Currently, there are several language models available, among which prominent ones include Llama 2, Llama 3, GPT 3.5, and GPT 4.

The Llama 2 is an open-source model that was released in 2023. It is available in versions with 7 billion (7B), 13 billion (13B), and 70 billion (70B) parameters. The Llama 2 was trained on a dataset composed of 2 trillion tokens sourced from publicly available sources and updated until September 2022 [7].

The Llama 3 was released by Meta in 2024. This open-source model was trained on 15 trillion tokens of data from publicly available sources. It is available in two versions, one with 8 billion (8B) parameters and data updated until March 2023, and another with 70 billion (70B) parameters and data updated until December 2023 [8].

GPT 3.5 and GPT 4 represent successive iterations within the Generative Pretrained Transformer (GPT) series of language models, developed by OpenAI. Released in 2021, GPT-3.5 boasts 170 billion parameters, while its successor, GPT-4, unveiled in 2024, possesses a staggering 170 trillion parameters [9]. GPT-4's training corpus comprises publicly available data updated until September 2021 [6]. Noteworthy improvements characterize GPT 4, including enhanced reliability, heightened safety features, and a reduction in instances of hallucinatory output,

as compared to its predecessor, GPT 3.5. In the context of this study, the GPT-3.5 and GPT-4 models were accessed through the ChatGPT 3.5 and ChatGPT 4 tools, which are specialized variants adapted specifically for real-time conversation and interaction.

The aim of this study is to construct domain-specific KGs for academic disciplines utilizing language models and to juxtapose the outcomes derived from each respective tool. Employed in this investigation were the Llama 2 7B, Llama 3 8B, ChatGPT 3.5, and ChatGPT 4 models, with analyses focused on the academic domains of Large Language Models and Psychoanalysis.

The methodology proposed in this study offers the advantage of streamlining the development of knowledge graphs. By utilizing information generated by Large Language Models for the construction of Knowledge Graphs (KGs), the need for manual (or even automated) data collection and ontology definition—traditionally the most labor-intensive and costly phases of knowledge graph construction—is obviated. Additionally, given that LLMs are trained on extensive datasets spanning diverse domains, the insights provided by the proposed model are anticipated to be more accurate, precise, and comprehensive.

Section 2 delineates the methodological framework adopted herein. Subsequently, Section 3 delineates the findings yielded by each language model concerning the aforementioned academic spheres. In Section 4, a critical examination of the feasibility of fine-tuning methodologies for addressing the problem at hand is undertaken. Finally, in Section 5, the article culminates in offering conclusive insights and considerations.

## 2 Methodology

The knowledge graphs generated in this study are representations of knowledge in the form of graphs, where nodes represent the most relevant terms in a given area of knowledge, and edges represent the relationships between these terms. The creation of KGs was carried out with the assistance of large language models, which were used to determine the nodes and edges of the graphs.

To identify the nodes that will constitute the KGs, it is necessary to initially determine the area of study to be analyzed and how many terms will be studied, i.e., how many nodes the generated graph will have. After this determination, the next step is to determine what these nodes will be.

For this purpose, for each of the adopted large language models, the following question was posed: "The [NOS] most relevant terms in the area of [AREA] are", where "[NOS]" is the quantity of terms to be studied and "[AREA]" is the selected field of research. This question aimed to identify the N most relevant terms in the initially chosen area of study, which are the nodes of the KG.

With the nodes of the KGs defined, it is necessary to define the connections between them. For this, it is first necessary to verify if there is indeed a relationship between the terms in the analyzed area of study. This verification is performed with the following inquiry: "Is there a relevant relationship between the terms [X] and [Y] in the area of [AREA]?" where "[X]" and "[Y]" are two of the selected terms for study. If the language model's response is affirmative, there will be an edge connecting the two analyzed nodes. If it is negative, the connection will not exist.

In cases where the language model verifies the existence of a relevant relationship between the nodes in the studied research field, it is necessary to identify what this connection is. For this, the following request is made to the program: "Fill in the blank using no more than 5 words. In the area of [AREA], [X] is __ of [Y]", where "[X]", "[Y]", and "[AREA]" are filled accordingly. The responses provided by the language model constitute the label of the edges of the KGs. Fig. 1 provides a comprehensive overview of the process utilized for constructing the KGs.

The questions used to generate the knowledge graphs were formulated using the Llama 2 model as a reference, as it is the simplest model among those used in this study. After defining the nodes and edges, the knowledge graphs were generated using the Python programming language. For this, the NetworkX library was used, which is a package used to generate, manipulate and study complex networks, such as graphs [10].

Knowledge graphs were produced for two areas of study, one with consolidated knowledge and one with knowledge still in development. For each area of study, the seven most relevant terms were selected, and the language models Llama 2, Llama 3, ChatGPT 3.5, and ChatGPT 4.0 were used.

One of the foremost challenges encountered in the utilization of large language models pertains to the considerable computational resources requisite for their operation. Consequently, in an endeavor to surmount these impediments, the Ollama tool was employed to facilitate access to the Llama 2 and Llama 3 models. Ollama represents an open-source initiative designed to afford access to and utilization of large language models devoid of the necessity for extensive computational infrastructure [11]. Additionally, the ChatGPT 3.5 and ChatGPT 4 models were executed through online interfaces proffered by OpenAI. The results of this application can be visualized in the next section.
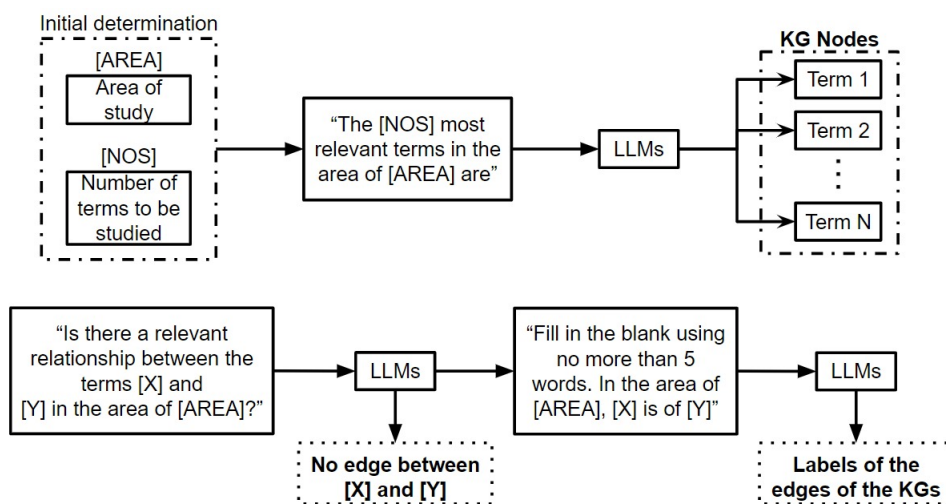
Figure 1. Process flowchart for defining knowledge graphs.

## 3 Results

### 3.1 Area with development knowledge

The methodology outlined in Section 2 was initially applied to the domain of Large Language Models (LLMs). This area of study is undergoing extensive development, with frequent disclosures of novel discoveries and advancements. It constitutes a burgeoning field, underscored by its contemporary significance.

Tab. 1 delineates the outcomes of applying the proposed methodology within the realm of Large Language Models, wherein the seven most pertinent terms were selected for scrutiny within this domain.

Table 1. Terms selected by each of the language models used as most relevant to the field of Large Language Models

| Llama 2 | Llama 3 | ChatGPT 3.5 | ChatGPT 4 |
|---|---|---|---|
| ChatGPT | Transformer | GPT | LLM (Large Language Models) |
| OpenAI | Pre-Training | Transformer | Training |
| GPT-3 | BERT | Fine-tuning | Fine-tuning |
| Deep Learning | Language modeling | Natural Language Understanding (NLU) | Parameter |
| Machine Learning | Self-supervised learning | Natural Language Generations (NLG) | Vector |
| NLP | Adversarial training | Transfer Learning | Embeddings |
| AI | Embeddings | Ethical and Social Implications | Transformer |

Analyzing the data presented in Tab. 1, it is evident that the results from the Llama 2, Llama 3, ChatGPT 3.5, and ChatGPT 4 models exhibit significant divergence. Comparing the outcomes between Llama 2 and Llama 3, it is notable that the models did not select any common terms. Similarly, comparing the results of ChatGPT 3.5 and ChatGPT 4, although two terms (Fine-tuning and Transformer) were shared, the remaining terms showed substantial disparity.

Such pronounced discrepancies among the outcomes generated by the adopted language models may stem from factors including the date and scale of the training dataset, as well as the architectural differences inherent in each method. So, depending on the model, it may necessitate more or less data to acquire proficiency in a given subject and furnish more accurate responses.

Conversely, the dataset employed in training the model also influences its performance regarding novel subjects. As discussed in Section 1, the models examined in this study have been trained on datasets up to varying dates. In theory, models with access to more recent content are expected to yield superior responses. However, as previously mentioned, the architectural nuances of the models can also play a pivotal role in this regard.

By utilizing the terms selected by each of the studied language models as deemed most relevant to the field of Large Language Models and employing the methodology outlined in Section 2, knowledge graphs were generated, as depicted in the forthcoming Fig. 2.
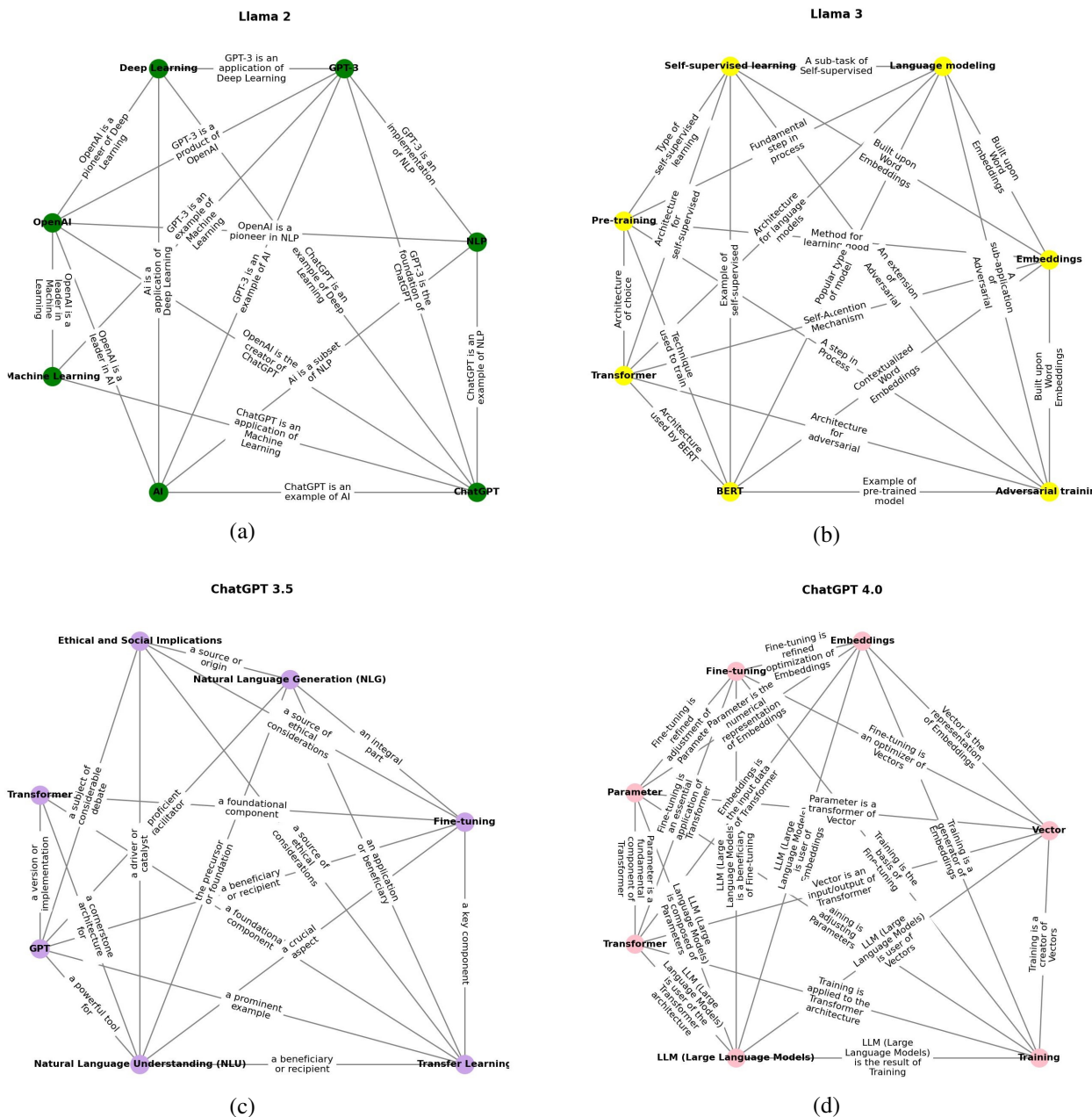


Figure 2. Knowledge Graphs for the field of Large Language Models generated by the language models Llama 2 (a), Llama 3 (b), ChatGPT 3.5 (c), and ChatGPT 4 (d)

Observing the knowledge graphs displayed in Fig. 2, it is noticeable that the relationships identified by the ChatGPT 4 model were more intricate, while those identified by Llama 2 were simpler. The associations of both Llama 3 and ChatGPT 3.5 were of good quality.

### 3.2 Field with consolidated knowledge

For the application of the methodology in a field with consolidated knowledge, the area of Psychoanalysis, a branch of Psychology, was chosen. Despite the extensive scientific production in this field today, it remains an area of study with a solid foundation of theories, practices, and research.

Thus, utilizing the methodology presented in Section 2, the 7 most relevant terms for the field of Psychoanalysis were selected using the language models Llama 2, Llama 3, ChatGPT 3.5 and ChatGPT 4. The selected terms can be viewed in Tab. 2 below.

Table 2. Terms selected by each of the language models used as most relevant for the field of Psychoanalysis

| Llama 2 | Llama 3 | ChatGPT 3.5 | ChatGPT 4 |
|---|---|---|---|
| Id | Id | Unconscious | Unconscious |
| Ego | Ego | Id | Repression |
| Superego | Superego | Ego | Transference |
| Oedipus complex | Unconscious | Superego | Oedipus Complex |
| Electra complex | Consciousness | Defense Mechanisms | Libido |
| Death drive | Oedipus complex | Transference | Ego |
| Unconscious | Repression | Countertransference | Superego |

Comparing the results obtained from the utilization of the language models Llama 2, Llama 3, ChatGPT 3.5, and ChatGPT 4 in a field of consolidated knowledge, it is evident that the outcomes presented by the models were more homogeneous than those obtained when employing a field of study still in development.

The terms "Ego", "Superego", and "Unconscious" were selected by all models as most relevant for the field of Psychoanalysis. Additionally, the terms "Id" and "Oedipus complex" were selected by three models each.

The greater homogeneity observed among the results obtained for a consolidated area of study underscores the significance of training data for the functioning of language models. When the models had access to a larger dataset to learn about the subject, their behavior was more similar. It is worth noting that the models did not produce identical results even when subjected to a substantial training dataset, which is related to the architecture of the models, i.e., how each of them acquires knowledge.

By using the terms defined by the models as the most relevant in the field of Psychoanalysis, the language models adopted in this study, and the methodology presented in Section 2, relationships between the selected terms were identified. With this information, it was possible to develop knowledge graphs for the field of Psychoanalysis generated by the models Llama 2, Llama 3, ChatGPT 3.5, and ChatGPT 4. The results of this procedure can be visualized in Fig. 3 below.

Analyzing the knowledge graphs depicted in Fig. 3, it becomes apparent that the Llama 2 and ChatGPT 3.5 models exhibited a more critical approach towards the relevance of the relationships among the studied terms compared to the other models. Additionally, it is evident that the associations identified by ChatGPT 4 and Llama 3 were more intricate and precise than those presented by Llama 2 and ChatGPT 3.5.

## 4 Fine-Tuning

In preceding sections, we elucidated the knowledge graphs derived through the proposed methodology, delineated the chosen fields of inquiry, and scrutinized the performance of the Llama 2, Llama 3, ChatGPT 3.5, and ChatGPT 4 models. While an element of similarity emerged when employing a domain of established knowledge, the results diverged considerably when applied to an evolving domain. This disparity may stem from nuanced factors encompassing both model architectures and the temporal dimension of their training datasets.

An avenue to ascertain the genuine impact of model architectures on the problem delineated in this article would entail fine-tuning the models up to a shared reference date, followed by reapplication of the proposed methodology to compare outcomes. Such an approach would render the results more sensitive to dataset similarities, thereby accentuating the models' learning mechanisms.

However, the practical feasibility of fine-tuning warrants contemplation. Taking the exemplar models crafted by Meta, Llama 3's training corpus notably surpassed that of Llama 2, rendering direct model comparison intricate. In theory, heightened data access engenders heightened response accuracy.
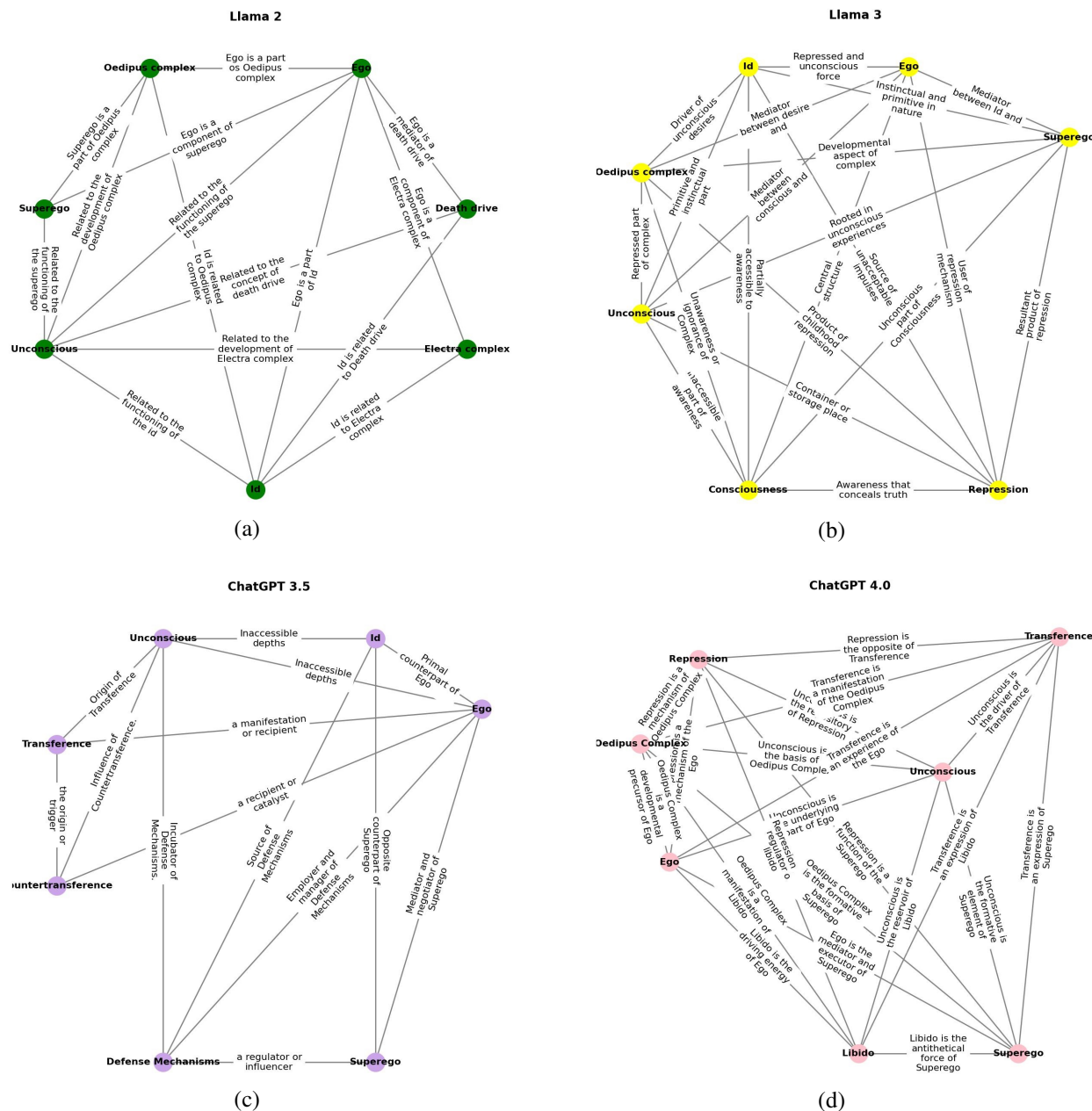
Figure 3. Knowledge Graphs for the field of Psychoanalysis generated by the language models Llama 2 (a), Llama 3 (b), ChatGPT 3.5 (c) and ChatGPT 4 (d)

Moreover, a hypothetical scenario wherein models are solely updated to a common date without cognizance of their prior knowledge presents a quandary. The sheer magnitude of datasets underpinning these language models poses a formidable challenge. For instance, the Llama 3 model's corpus spans 15 trillion tokens, as per Meta's disclosures [8]. A dataset representing a mere fraction of this total would necessitate the production of an extensive volume of tokens within a constrained timeframe.

Furthermore, it behooves recognition that language models synthesize knowledge from multifarious domains to engender responses. In the context outlined in this article, for a language model to adjudge, for instance, the pertinence of "Deep Learning" within the realm of "Large Language Models", it must not only be conversant with the field of study but also with the term under scrutiny. Consequently, fine-tuning mandates not only specialization within the studied domain but also across cognate disciplines.

Given constraints imposed by dataset update frequency and the imperative of encompassing diverse fields of inquiry, the implementation of fine-tuning to exert discernible influence on model outcomes would be onerous and protracted. Additionally, the likelihood of significant impact remains indeterminate, contingent upon the subject

domain. Consequently, for these reasons, the recourse to fine-tuning models was deemed impracticable.

## 5   Conclusion

In this study, knowledge graphs were crafted utilizing diverse language models across varied domains of expertise. It's imperative to underscore that while each language model generated distinct outcomes, it would be fallacious to assert the superiority of one over the others. Were the queries posed to the language models, as delineated in Section 2, to be directed to different subject matter experts, disparate responses would likely ensue. Thus, the objective isn't to hierarchize one model as inherently superior to another, but rather to apprehend the idiosyncrasies characterizing each model.

For a cogent analysis, it's paramount to recognize that the interrogative framework employed in this methodology was initially tailored for the Llama 2 model. Given its purported simplicity relative to the other models under scrutiny, it's posited that the inquiry approach conducive to favorable outcomes with this model would similarly prove efficacious for more intricate models.

Upon meticulous examination of the outcomes proffered by each model and their operational nuances, it was discerned that the Llama 3 model evinced a degree of erratic behavior. Instances were observed where this model yielded markedly divergent responses to repetitive inquiries. Additionally, while the Llama 2 model's outcomes were commendable, they were circumscribed in scope. Conversely, the ChatGPT 4 model exhibited the most user-friendly demeanor and evinced superior cognitive prowess in grasping the posed queries. The ChatGPT 3.5 model also yielded favorable results, albeit trailing behind its successor.

From the empirical insights gleaned in this study, it can be inferred that when grappling with subjects entrenched in established domains of knowledge, language models, leveraging their access to expansive training data, typically manifest commendable performance. However, when contending with nascent domains, language models exhibit pronouncedly divergent behavior, influenced by both the volumetric abundance of available training data and the architectural underpinnings unique to each model, which decisively shapes their knowledge extraction methodologies.

**Authorship statement.** The authors hereby confirm that they are the sole liable persons responsible for the authorship of this work, and that all material that has been herein included as part of the present paper is either the property (and authorship) of the authors, or has the permission of the owners to be included here.

## References

[1] N. Koul. Knowledge graphs. Artigo do Medium, 2024.

[2] M. Trajanoska, R. Stojanov, and D. Trajanov. Enhancing knowledge graph construction using large language models. *arXiv preprint arXiv:2305.04676*, 2023.

[3] X. Chen, S. Jia, and Y. Xiang. A review: Knowledge reasoning over knowledge graph. *Expert systems with applications*, vol. 141, pp. 112948, 2020.

[4] A. Kau. Automated knowledge graph construction with large language models — part 2. Artigo do Medium, 2024.

[5] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong, and others. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023.

[6] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat, and others. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

[7] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, and others. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

[8] AI@Meta. Llama 3 model card, 2024.

[9] A. Koubaa. Gpt-4 vs. gpt-3.5: A concise showdown, 2023.

[10] N. developers. Networkx, 2024.

[11] Medium. Ollama : What is ollama? Artigo do Medium, 2024.