

# Monitoring Transformative Technological Convergence Through LLM-Extracted Semantic Entity Triple Graphs

Alexander Sternfeld<sup>a,\*</sup>, Andrei Kucharavy<sup>a</sup>, Dimitri Percia David<sup>a</sup>, Alain Mermoud<sup>b</sup>, Julian Jang-Jaccard<sup>b</sup> and Nathan Monnet<sup>b</sup>

<sup>a</sup>*Institute of Entrepreneurship and Management, HES-SO, Sierre, Switzerland*

<sup>b</sup>*Cyber-Defence Campus, armasuisse, Science and Technology, Thun, Switzerland*

## ARTICLE INFO

### Keywords:

Bibliometrics  
Entity Extraction  
Machine Learning  
Technological Forecasting  
Large Language Models

## ABSTRACT

Forecasting transformative technologies remains a critical but challenging task, particularly in fast-evolving domains such as Information and Communication Technologies (ICTs). Traditional expert-based methods struggle to keep pace with short innovation cycles and ambiguous early-stage terminology. In this work, we propose a novel, data-driven pipeline to monitor the emergence of transformative technologies by identifying patterns of technological convergence.

Our approach leverages advances in Large Language Models (LLMs) to extract semantic triples from unstructured text and construct a large-scale graph of technology-related entities and relations. We introduce a new method for grouping semantically similar technology terms (*noun stapling*) and develop graph-based metrics to detect convergence signals. The pipeline includes multi-stage filtering, domain-specific keyword clustering, and a temporal trend analysis of topic co-occurrence.

We validate our methodology on two complementary datasets: 278,625 arXiv preprints (2017–2024) to capture early scientific signals, and 9,793 USPTO patent applications (2018–2024) to track downstream commercial developments. Our results demonstrate that the proposed pipeline can identify both established and emerging convergence patterns, offering a scalable and generalizable framework for technology forecasting grounded in full-text analysis.

## 1. Introduction

Accurate forecasting of technological disruptions is essential for organizations seeking to remain competitive and resilient. While incremental innovations can typically be anticipated and managed within existing planning frameworks, *transformative technologies* introduce paradigm shifts that reshape entire domains and often spill over into adjacent sectors [62, 29]. Because such technologies challenge established conceptual frameworks, they are inherently difficult to predict [19].

Despite these challenges, the strategic importance of anticipating transformative change has long motivated the development of forecasting methods, particularly since the Cold War era [41]. Seminal approaches such as the Delphi Method [24] and S-curve Substitution Analysis [32] have inspired decades of methodological refinement [21, 16, 35]. However, these classical methods struggle to keep pace with the rapid innovation cycles characteristic of Information and Communication Technologies (ICTs). For example, the term “Large Language Models” (LLMs) emerged in 2019, and within just three years, ChatGPT had become a globally recognized application [46]. In such contexts, the reliance of classical methods on expert panels and historical data renders them ill-suited [22].

As a result, ICT forecasting has increasingly shifted toward data-driven methodologies [23], with scientometrics

emerging as a promising approach [59]. Yet even data-centric methods face difficulties in the early identification of transformative technologies, due in part to the absence of stable terminology and well-defined application domains in the initial stages. Continuing the example of LLMs, even a year after the release of ChatGPT, major bibliometric platforms such as OpenAlex still lacked ontology terms for many of the foundational sub-technologies [68, 18, 72]. While Würsch et al. [72] showed that extracting technology-related nouns from scientific texts is feasible, the process remains noisy and challenging to use effectively for large-scale monitoring and forecasting.

In this work, we present an alternative approach that leverages recent advances in Large Language Models (LLMs) to enable scalable extraction and analysis of semantic triples centered on technological concepts [65]. We focus on *technological convergence* as an early indicator of transformative potential [44, 26, 50].

Our main contributions are as follows:

- (i) We introduce an unbiased pipeline for extracting semantic triples that involve technology-designating nouns.
- (ii) We propose a novel syntactic similarity identification technique—termed *noun stapling*—to group related technology terms.
- (iii) We implement a coarse decomposition of technology components to improve interpretability.
- (iv) We develop a graph-based method to detect and track technological convergence over time.

\*Corresponding author

✉ alexander.sternfeld@hevs.ch (A. Sternfeld)

ORCID(s): 0009-0008-6801-6160 (A. Sternfeld);

0000-0003-0429-8644 (A. Kucharavy);

0000-0003-0429-8644 (D.P. David); 0000-0001-6471-772X

(A. Mermoud); 0000-0002-1002-057X (J. Jang-Jaccard)

These elements are integrated into a unified pipeline for identifying the emergence of transformative technologies through convergence patterns. We demonstrate the utility of this pipeline through a case study on LLMs, applying it to a large corpus of 278,625 arXiv preprints (2017–2024) and 9,793 USPTO patent applications (2018–2024).

The remainder of this paper is structured as follows: Section 2 reviews the related literature. Section 3 describes the data sources. Section 4 outlines the methodology. Section 5 presents the results. Finally, Section 6 concludes the paper and discusses directions for future research.

## 2. Related work

In this section, three branches of previous research are discussed. First, we discuss the main topics and trends in bibliometrics. Then, we discuss related work on claim and triple extraction, where we highlight those methods that leverage machine learning techniques. Last, we discuss the usage of topological analyses for technology forecasting.

### 2.1. Bibliometrics

As a method for technological forecasting, bibliometrics leverages both qualitative and quantitative analysis of recorded information—primarily scientific books, articles, and patents. The systematic study of written works has long been established, with the term *bibliometrics* coined in the early 20<sup>th</sup> century. One of the most widely recognized classical bibliometric techniques is citation analysis, which examines citation patterns among academic publications [64].

Early citation analysis focused primarily on citation counts and co-citation networks. However, the advent of advanced network analysis methods and increasing computational power has enabled more sophisticated analyses of citation networks. Such methods are now employed across a range of fields for technology forecasting. For instance, Qiu and Wang [60] conducted a citation network analysis of robotics patents to identify key pathways of innovation. Similarly, Li et al. [49] analyzed both patents and academic papers in the domain of nanogenerator technologies, offering a more integrated view of scientific and technological progress.

Traditionally, bibliometric studies were constrained to structured metadata. Recent advances in Natural Language Processing (NLP), however, have opened new avenues for deeper analysis of full-text documents. Techniques such as  $n$ -gram analysis [54], Latent Dirichlet Allocation (LDA) [13], semantic embeddings [55], deep neural language model encodings [53], and, more recently, large language models (LLMs) [11], have greatly expanded the scope of bibliometric inquiry. These developments are further supported by the growing availability of full-text documents through open-access and preprint platforms, such as arXiv [1], medXiv [3], bioRxiv [2], and PubMed Central [4].

For example, Percia David et al. [59] analyzed computer science preprints on arXiv using LLM-derived sentence embeddings to extract author sentiment and relate it to

technology security over time. Likewise, Sandu et al. (2024) constructed a collaboration network in the area of social media text mining and identified recurring research themes using  $n$ -gram analysis.

In parallel with developments in literature-based bibliometrics, there has been increasing attention to patent analysis, given the rich technical detail and close ties to commercial applications that patents offer. While academic papers often represent early-stage scientific exploration, patents typically reflect efforts to translate these discoveries into deployable technologies. Importantly, both European and U.S. patents follow the Cooperative Patent Classification (CPC) system, enabling systematic, cross-jurisdictional analysis. For example, Kim and Bae [43] cluster patents based on CPC codes and apply citation and claim analysis to forecast emerging technologies. Similarly, You et al. [73] analyze light generator technology patents using Bass and ARIMA models to identify promising innovation trajectories. By integrating bibliometric insights from both scholarly publications and patents, researchers gain a more comprehensive view of technological development—spanning foundational research to commercial realization [49, 5].

### 2.2. Claim and triple extraction

Semantic Triples, also known as RDF (Resource Description Framework) triples, are a fundamental representation for knowledge extraction and organization in semantic web technologies. A semantic triple consists of three components: a subject, a predicate, and an object, and is used to represent factual statements in a machine-readable format [12]. This structure is foundational in knowledge graphs, enabling systems to store and reason about facts in a way that can be queried and expanded upon [37]. In the context of claim extraction from scientific literature, we aim to extract claims that can be represented as semantic triples, focusing on technological assertions and relationships. By adhering to the semantic triple standard, these claims are organized into unrooted, unbiased knowledge graphs that facilitate downstream tasks such as pattern mining and knowledge discovery.

Previous approaches to claim extraction can be broadly categorized into heuristic and machine learning methods. Heuristic methods, such as the approach by Jansen and Kuhn [40], do not require training data and have low computational cost. However, they are limited in their ability to capture complex claims. Jansen and Kuhn [40] extract core claims by scoring sentences based on pattern matching, term frequency, and sentence length. However, they only succeed in rewriting claims to the required AIDA (atomic, independent, declarative, and absolute) standard in only 29% of cases. In contrast, machine learning methods like the one proposed by Li et al. [48], using BiLSTMs for sentence classification, offer better performance, categorizing sentences with an F1 score over 0.8. However, these models require domain-specific supervised training data, which must be updated regularly.

Most triple extraction models rely on supervised training, such as RECON and sPERT, which require labeled training data [10, 28]. These methods are limited by their dependence on available training data and predefined relations. In contrast, unsupervised models, like Stanford OpenIE, do not require labeled data or a predefined schema to extract triples. However, OpenIE tends to be too aggressive, often extracting non-useful relations [51]. Ottersen et al. [57] address this issue by using a language model to extract triples within a predefined set of relations, improving precision but requiring users to specify all possible relations.

### 2.3. Topological Analyses for Technology Forecasting

Traditional co-word analysis techniques primarily relied on co-occurrence frequencies to detect associations between terms [17]. More recent approaches, however, emphasize the connectivity and centrality of entities within complex graphs to identify zones of technological convergence and emerging innovation clusters. At the core of these methods are *knowledge graphs*, which are constructed by extracting key entities such as technologies, methods, and applications from scientific publications and patents. In these graphs, nodes represent concepts, and edges encode semantic, citation, or temporal relationships. Keyword extraction tools, such as RAKE [36] and KeyBERT, are often used to enrich these graphs. The resulting structures can be analyzed as either static or dynamic representations of technological domains.

The topological characteristics of knowledge graphs are crucial for forecasting applications. For instance, Dotsika and Watkins [27] demonstrated that keyword co-occurrence networks derived from scientific and business literature can be effectively analyzed using centrality metrics to detect potentially disruptive or fast-growing technological areas. Building on this idea, Li et al. [49] developed dual-layer networks that integrate scientific publications and patents to trace the evolution and commercialization of nanogenerator technologies. More recently, Wang et al. [69] applied semantic term extraction in combination with Louvain-based community detection to identify early-stage convergence between disparate research fields.

Community detection continues to be a foundational component of topological forecasting. The Louvain algorithm, introduced by Blondel et al. [14], is widely used to identify coherent topic clusters within large-scale graphs. Complementary measures, such as eigenvector centrality, highlight influential nodes that play a key role in the dissemination and integration of knowledge. These techniques support dynamic analyses of how scientific topics emerge, grow, merge, or decline over time. Dynamic knowledge graphs, such as those implemented in the Science4Cast project [45], enable the longitudinal monitoring of research landscapes and facilitate the detection of novel knowledge pathways as they form.

In addition, predictive modeling on these graphs—using either neural networks or handcrafted topological features—can forecast whether previously unconnected concepts are likely to become linked in the future, and whether such connections are poised to attract significant attention [36, 52].

## 3. Data

The methodology developed in this work can be used for any textual data source. Several examples of such sources are scientific articles, patents, news articles, social media posts and web pages. In this paper, we focus on two data sources: scientific articles and US patent applications. In order to evaluate the capabilities of different methods of technology-related semantic triples extraction, we created a golden dataset of manually labeled semantic triples in the LLM field.

### 3.1. Golden Dataset for Triples Extraction

To fine-tune the LLMs for technology-related semantic triples extraction, a labeled dataset is required, both for evaluation of method performance, and for supervised training. To the best of our knowledge, such a labeled dataset for semantic triples extraction from scientific papers does not exist. Therefore, we manually construct a training dataset based on the paper *A Survey of Large Language Models* [74]. We chose this paper given our focus on the LLMs transformative technology, since it is a comprehensive and general review of LLM component technologies. In total, we manually annotated 547 triples in 100 paragraphs, given that amount of examples is generally considered to be sufficient for parameter-efficient fine-tuning (PEFT), while leaving room for a test holdout set [30]. To evaluate the performance of the fine-tuned LLMs, we use a holdout set of 20% of the annotated paragraphs. The dataset is available for download from the data and code repository associated to this article: <https://github.com/submissiontfscanonymized/tfsc2025>.

### 3.2. arXiv data

In the first use case, preprints publicly available on arXiv are considered. As arXiv is an open scholarly preprint archive, the documents are not peer-reviewed. We choose this platform for two key reasons. First, as there is no peer review and only a short moderation process, submissions to arXiv are available rapidly, tracking the state of research as close to real time as possible. The archive includes papers that may not have been accepted at conferences due to perceived lack of immediate interest, yet have gone on to receive significant citations — highlighting their eventual importance. Second, arXiv ensures that all articles are classified by the authors into a category from a predefined taxonomy. To aid authors, arXiv has published an algorithm that can suggest the best fitting categories for a paper, based on the existing corpus of articles<sup>1</sup>. Finally, all

<sup>1</sup><https://info.arxiv.org/help/api/classify.html#cat>

**Table 1**

The arXiv categories that are considered for this study, alongside their full names.

Category	Full name
cs.CL	Computation and Language
cs.LG	Machine Learning
cs.AI	Artificial Intelligence
cs.IR	Information Retrieval
stat.ML	Machine Learning

arXiv pdf's and raw latex files can be downloaded through Google Cloud Storage Buckets, which are updated weekly<sup>2</sup>, making text mining on arXiv preprints significantly easier.

For an initial evaluation and refinement of our methodology, we only worked with papers from December 2023 (4225 papers), from the arXiv categories `cs.AI`, `cs.CL` and `cs.LG`. We choose to do so, as a small dataset reduces the computation time and required computational resources. Therefore, such a dataset is suitable for early refinements of the pipeline. After settling on the final methodology, all papers between 2018 and 2024 are considered, from the arXiv categories that are specified in Table 1 for a total of 278,625 articles. Both the categories and timeframe correspond to the timeline of convergence of technologies that led to the transformative LLM technology emergence.

### 3.3. USPTO patent data

The second data source we consider is USPTO patent data. We choose to consider patents as a secondary data source as it provides information on downstream applications of technologies. While preprints on arXiv will show early signs of novel technological developments, the presence of technologies in patent applications shows that technologies are maturing and are being adopted in commercial products.

Specifically, we consider patent applications to the USPTO, which are publicly available through the Bulk Data Directory<sup>3</sup>. Similar to the arXiv papers, we consider the period 2018-2024. All patent applications are categorized according to the Cooperative Patent Classification (CPC) system, which is a classification system that is used both by the USPTO and the European Patent Office (EPO).

For our case study, we filter the data for those patents that contain at least one of the terms *large language model*, *large language models*, *llm* or *llms* in the abstract or title. This results in a total of 9793 patent applications between 2018 and 2024.

## 4. Methodology

In this Section we describe the methodology of our triple extraction pipeline and downstream analyses. We first describe the data preprocessing, after which two triple extraction methods are presented. Then, the post-processing

and filtering of the triples is elaborated upon. Last, we explain our proposed *noun stapling* methods to group similar triples, and our downstream analyses. The full pipeline is displayed in Figure 1.

### 4.1. Preprocessing

The `PyMuPDF` pdf processing library [7] is used to turn the pdfs into text. The resulting text is then processed to remove citations, using a heuristic pattern matching. We remove both square and round brackets with only numbers inside, as well as brackets with a year inside, matching most frequent citation formats found in the arXiv papers. In addition to that, we remove the line breaks using a similar heuristic, although with dashes.

We then detect and expand the abbreviations by mapping them to their long form. To ensure that the abbreviation resolution is done properly, we benchmark different algorithms on the *Fundamentals of Generative Large Language Models and Perspectives in Cyber-Defense* report, given its mixed usage of both ML, cybersecurity, and cyberdefense abbreviations [47]. Table 2 shows the comparative performance of the Schwartz-Hearst [63], `scispaCy` [56], `NLPre` [66] and a fine-tuned `RoBERTa` model [76] as abbreviation detection methods. The results show that the Schwartz-Hearst algorithm performs the best, with `scispaCy` coming a close second. Given the speed of Schwartz-Hearst, we select it for our pipeline. Table 3 shows an example of a sentence before and after preprocessing.

### 4.2. spaCy-based triple extraction

#### Claim extraction

After preprocessing the raw text, the sentences that constitute the core claims of the text are identified. To this end, the *ClaimDistiller* framework is used, which is developed by Wei et al. [70]. In this work, both CNN and BiLSTM models have been employed for claim extraction tasks through training on the PubMed-RCT and SciARK datasets [31, 25]. While incorporating supervised contrastive learning has been shown to enhance model performance, it also introduces additional computational cost. To maintain a balance between effectiveness and efficiency, we opt for the BiLSTM model variant without supervised contrastive training. This model is used to extract claims from academic papers, thereby reducing the number of sentences that need to be processed during the subsequent triple extraction phase.

#### triple extraction

Next, we aim to reduce the claims to (*subject*, *predicate*, *object*) semantic triples, analogous to the Resource Description Framework (RDF) format commonly used in the representation of knowledge in the Web Ontology Language (OWL). We choose this representation, as it will facilitate the comparison of claims across papers. To achieve this, we use the Python library `textacy`, which is built on `spaCy` [38] and has a built-in extraction method that does not require the specification of relations in advance.

<sup>2</sup><https://www.kaggle.com/datasets/Cornell-University/arxiv>

<sup>3</sup><https://data.uspto.gov/bulkdata/datasets>



**Table 2**

Performance of three abbreviation detection algorithms on the *Fundamentals of Generative Large Language Models and Perspectives in Cyber-Defense* report [47]

	Correctly detected	False positive	Time
Schwartz-Hearst	11	1	0.04 s
scispaCy	11	4	8.73 s
NLPRe	0	0	0.01 s
Fine-tuned RoBERTa	10	3	100 s

**Table 3**

Illustration of the preprocessing effect, the parts in bold are altered during preprocessing.

Uncleaned	Preprocessed
Society has been affected by <b>artificial intelligence (AI)</b> and has become more <b>reliant</b> on <b>AI</b> products.	Society has been affected by <b>artificial intelligence</b> and has become more <b>reliant</b> on <b>artificial intelligence</b> products.

### 4.3. LLM-based triple extraction

As a second method for triple extraction, we consider an approach based on LLMs. The main advantage of such an approach is that LLMs can take into account context when extracting triples from a sentence. However, one has to be cautious, as LLMs can hallucinate and one can thus never be certain regarding the factuality of generations. In fact, LLMs have been reported to struggle with niche concepts [71]. Therefore, prior to using LLMs it is critical to assess LLMs for hallucinations during triple extraction. Three state-of-the-art LLMs are considered: Mistral-7B-Instruct-v0.2, Meta-Llama-3-8B-Instruct and Starling-LM-7B-beta [42, 75, 6]. Whereas the Starling and Llama models have a context length of 8192, the Mistral model has a context length of 32,000. As we require a scalable and replicable triple extraction method with limited resources, we do not consider larger or API-only models.

#### *Few-shot learning prompts and fine-tuning*

Generative LLMs per se are not specifically pretrained or instruction-tuned for triple extraction. In order to improve their performance at this task, the two common approaches are few-shot learning prompts [15], or fine-tuning [61].

Few-shot learning is generally considered as simpler to implement, given that it is equivalent to providing examples of desired text transformation as part of the user prompt. We formulated several of such example prompts based on the existing literature examples, and selected the best performing one. The resulting prompt can be found in Appendix Section A.1

Fine-tuning LLMs is generally believed to be more effective, but requires significantly more computational power and a substantial amount of fine-tuning examples that scales with the model size. It is possible to mitigate both problems using parameter-efficient fine-tuning (PEFT), which reduces both the computational power requirements and offers better generalization with less data [30]. Specifically, we use

LoRA [39], which freezes the majority of model parameters and fine-tunes only the scaling matrix of a singular value decomposition of a random matrix added to the model weights.

### 4.4. Post-processing

Given that the end goal is to compare triples across papers, we normalize triple representation in order to facilitate the matching of nouns likely to designate technologies of interest. Specifically, we:

1. Lowercase all words in the triple;
2. Remove triples where either the subject or object contains more than 6 words;
3. Remove stopwords from the triples based on the list included in NLTK;
4. Remove any non-text characters;
5. Lemmatize verbs and nouns in the triple;
6. Remove words containing less than 3 characters;

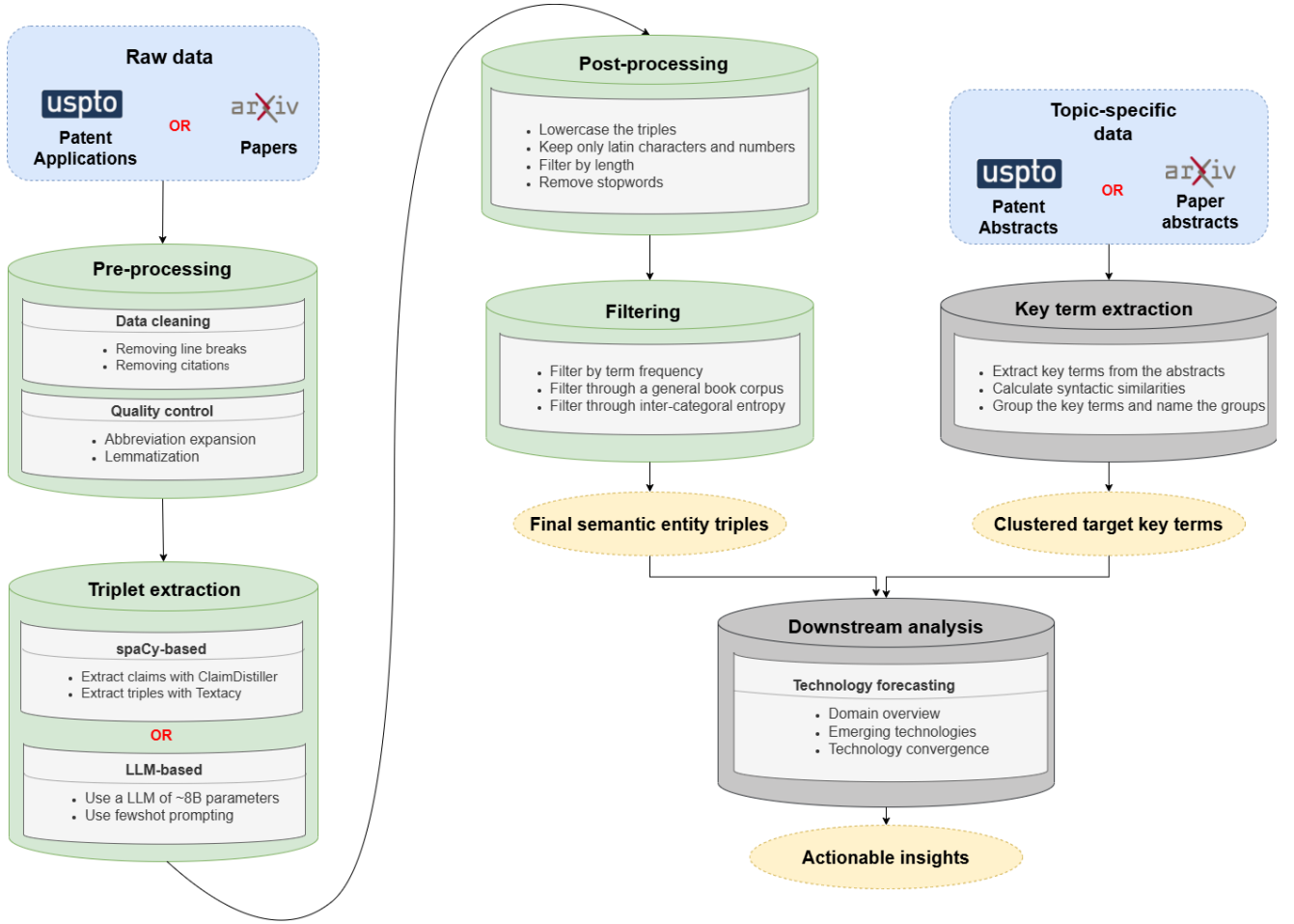
Both cutoff values in 2. and 6. were chosen empirically, based on manual inspection of triple nouns.

### 4.5. Filtering

As a final step, the triples are filtered, in order to keep only triples with nouns relevant to technological forecasting. This filtering is done in three steps.

#### *Frequency*

First, we filter out triples with subjects or objects that are too rare. We do so by considering a *control corpus* of all papers from the last quarter of 2023, enriched by a sample of 400 papers per month from the period January 2015 - September 2023. For every term  $t$ , we then define  $f_{cc,t}$ , which is the number of papers in the control corpus that contains the term at least 3 times per page, on average. If a



**Figure 1:** The complete pipeline for the triple extraction and downstream analysis, starting from either raw patent data or raw arXiv data. The green components of the pipeline reflect the triple extraction procedure, including the pre- and post-processing steps. The grey components of the pipeline illustrate the key-term extraction, noun stapling and the downstream analyses.

triple has a subject or object of which  $f_{cc,t} < 5$  for all terms  $t$  in the subject or object, we remove this triple. This filtering step aims at removing nouns that are not used sufficiently frequently to correspond to a potential technology.

### Bookcorpus

Second, we use the general-purpose Gutenberg book corpus to detect triples that are generic and thus carry little information for technological monitoring. Let  $f_{b,t}$  define the number of books in which term  $t$  appears at least 3 times per page. Furthermore, we denote the number of papers in the control corpus as  $N_p$  and the number of books in the book corpus as  $N_b$ . Then, each term  $t$  obtains the following score  $s_t$ :

$$s_t = \begin{cases} \log\left(\frac{f_{cc,t}}{N_p}\right) - \log\left(\frac{f_{b,t}}{N_b}\right) & \text{if } f_{b,t} > 0 \\ \infty & \text{if } f_{b,t} = 0 \end{cases} \quad (1)$$

The score of a term thus increases when the frequency of the term in the book corpus is lower. Simultaneously, the score increases when the frequency of the term in the control corpus is higher. We keep the triples of which at least one

word from the subject and object is in the top 10% of the term scores, with cutoff determined empirically.

### Entropy

Last, the triples are filtered through the cross-arXiv-categorical word entropy of the subject and object, and hence are informative as to the field of the paper in which they occur. Such field-informative nouns are most likely to refer to technologies or technology-adjacent concepts rather than being words characteristic of scientific writing that are domain-agnostic, such as "methodology" or "lab".

Formally, let us denote the set of all arXiv categories as  $C$ , then the cross-categorical entropy for term  $t$  can be calculated as denoted below.

$$H_t = \sum_{c \in C} P(t \in p | p \in c) \cdot \log\left(\frac{1}{P(t \in p | p \in c)}\right) \quad (2)$$

Here, we denote a paper as  $p$  and the individual arXiv categories as  $c$ .

#### 4.6. Noun stapling

At this point in the pipeline, we have triples involving nouns that designate with high probability technologies and technology-adjacent terms, conformed to a somewhat standard representation. However, in order to perform an emerging technology analysis through technological convergence, we need to connect technology-designating nouns that indicate similar technologies despite naming variations. We refer to this as *noun stapling*.

When stapling (compound) nouns, we want to group subjects or objects together that are semantically sufficiently similar. As we are dealing with big datasets, we use a syntactic string similarity measure, which is fast and can combine detection of similarities at token and character level.

Specifically, we use a novel string similarity measure, introduced by Gali et al. [33], that utilizes a combination of character- and token-level similarity measures. Specifically, the framework uses the **soft cardinality of a set**, as introduced by Vargas [67]. The aim of a soft cardinality measure is to capture the number of unique concepts in a string. For example, the set  $(dog, dogs)$  should have a lower soft cardinality than the set  $(dog, church)$ .

The results of Vargas [67] showed that, in general, the combination of Dice (token-level) and 3-gram (character-level) worked well in all settings. Therefore, we choose to use this procedure, where we tokenize the compound nouns by splitting on spaces. Let us first introduce the 3-gram similarity measure between the strings  $s_1$  and  $s_2$ , which is defined as

$$QGramsSim(s_1, s_2) = 1 - \frac{\sum_{i=1}^n |match(q_i, Q_{s_1}) - match(q_i, Q_{s_2})|}{|Q_{s_1}| + |Q_{s_2}|}, \quad (3)$$

where  $Q_{s_1}$  and  $Q_{s_2}$  are the sets of q-grams from  $s_1$  and  $s_2$  respectively. Furthermore,  $n = |Q_{s_1} \cup Q_{s_2}|$ , and  $match(q_i, Q_{s_1})$  is the number of times the q-gram  $q_i$  appears in  $Q_{s_1}$ . This measure can then be used to calculate the soft cardinality of a set  $T = (T^1, T^2, \dots, T^n)$ , which is defined as

$$|T|_{soft} = \sum_{i=1}^n \frac{1}{\sum_{j=1}^n QGramsSim(T^i, T^j)} \quad (4)$$

Using the soft cardinalities, one can then compute the similarity between the two compound nouns  $t_1$  and  $t_2$  through an adjustment of the token-level measure. For the DICE measure, this then becomes:

$$SoftDICE(t_1, t_2) = \frac{2 \times |t_1 \cap t_2|_{soft}}{|t_1|_{soft} + |t_2|_{soft}} \quad (5)$$

Last, one must choose a threshold above which the two compound nouns  $t_1$  and  $t_2$  are considered to be similar. As

we want to be conservative, we choose a threshold of 0.85. Above this threshold, one can be relatively certain that two compound nouns are sufficiently similar to be marked as such.

#### 4.7. Key term extraction and clustering

From the triple extraction pipeline, we want to create a large graph of technology-related topics (nodes). To derive insights from this graph, it is essential to focus on a specific domain of interest. To this end, one requires a list of relevant topics. The manual curation of such a list requires a strong domain expertise, which is regularly not readily available.

Therefore, we develop a separate technology-designating keyword extraction for the first step of our approach. After manually specifying a few terms that broadly specify the field of interest, we select papers from arXiv that contain at least one of these terms in the abstract or title. Then, we extract a maximum of 10 key terms for each abstract using KeyBERT, with `allenai/specter` as the embedding model [34, 20].

Next, we use the noun stapling method as described in Section 4.6 to calculate similarities between the extracted key terms. We then use the pseudo-algorithm as given in Table 4 to cluster the key terms into **topics**. The algorithm takes a simple approach, leveraging the extraction count for each term. Starting from the most popular term, we iterate over the terms and check whether they are already assigned to a group. If not, we make a new group to which we add the term and all similar terms that are not yet assigned to a group. Last, the groups are named based on the most common substring across the key terms. If two substrings are equally common, we prefer to choose the longest one as it is likely more informative. For each triple, we then assign it to a cluster if the subject or object is at least similar to one of the terms in the cluster. Again, we use the similarity measure as defined in equation 5.

#### 4.8. Graph analyses

Through the extraction and clustering of key terms in the domain of interest, and assigning the triples to the corresponding key topics, we have built a network of technological topics. We then perform several analyses on the network to derive useful insights from it.

First, we identify established related technology clusters by partitioning the network into densely connected sub-networks through the Louvain community detection method [14]. Specifically, we use the implementation in Gephi with a resolution of 0.85 [9]. We scale the nodes based on the number of triples corresponding to the topics. The labels are scaled through the eigenvector centrality, which measures the influence a node has in a connected network.

To assess technology convergences over time, we consider the topic co-publication frequency. Specifically, we use the Jaccard similarity between topics. Let  $A(t)$  and  $B(t)$  be the sets of triplets that are related to topics  $A$  and  $B$  at time  $t$ , respectively. Then, the Jaccard similarity between the

**Table 4**

Pseudo algorithm for keyword clustering, based on the soft dice similarity as introduced in equation 5.

Pseudo algorithm: keyword clustering based on similarity	
1:	<b>Inputs:</b>
2:	<i>keywordsList</i> : A list of keywords ordered by extraction count
3:	<i>similarKeywords</i> : A dictionary which gives, for every keyword, a list of similar keywords
4:	<b>Output:</b>
5:	<i>Grouped keywords</i> : A dictionary of groups of keywords
6:	assignedKeywords ← empty set
7:	groups ← empty dictionary
8:	<b>for</b> each keyword1 in keywordsList <b>do</b>
9:	<b>if</b> keyword1 in groups <b>then</b>
10:	<b>continue</b>
11:	<b>end if</b>
12:	currentGroup ← empty list
13:	add keyword1 to currentGroup
14:	<b>for</b> each keyword2 in similarKeywords[keyword1]
15:	<b>if</b> keyword2 is in assignedKeywords <b>then</b>
16:	<b>continue</b>
17:	<b>end if</b>
18:	add keyword2 to currentGroup
19:	add keyword2 to assignedKeywords
20:	<b>end for</b>
21:	groups[keyword1] ← currentGroup
22:	add keyword1 to assignedKeywords
23:	<b>end for</b>
24:	<b>return</b> groups

topics  $A$  and  $B$  at time  $t$  is defined as

$$J_{A,B}(t) = \frac{|A(t) \cap B(t)|}{|A(t) \cup B(t)|} \quad (6)$$

, The Jaccard similarity will always be between 0 and 1, with a higher similarity indicating a high co-occurrence between the topics. A technology convergence will be identifiable by an increase in the Jaccard similarity.

#### 4.9. Implementation

To download the data, preprocess the text and extract the triples, we used 4 NVIDIA A100 Tensor Core GPU's with 40GB RAM. For the deployment of the pipeline, at least one GPU with 16GB of VRAM is required to run the LLMs. Depending on the size of the data, more GPUs may be required for swift processing. Furthermore, it is recommended to use at least 10 cores for swift downloading and preprocessing of arXiv papers and USPTO patent applications.

### 5. Results

#### 5.1. Triple extraction methods comparison

In this first section, we will discuss the difference in performance between the spaCy and different LLM triple extraction methods, justifying our use of a specific LLM model and adaptation approach down the line. We are interested in four dimension of triple extraction performance: the number of triples extracted, the format of the generated text, signs of hallucination and the computation time.

Although more extracted triples potentially mean a more granular technology identification, it is crucial to preserve their quality. Part of such assessment is evaluation of generative LLMs for hallucination instead of extraction, which we do by mapping the triples back to the reference text. Specifically, we consider whether the subjects and objects are present in the original text. We use the Levenshtein distance, which is defined as the number of single-character edits to change one word into the other. The results, presented in Table 5 shows that there are also inconsistencies in the human-annotated triples. Through manual inspection, we identify that these inconsistencies in the human-annotated triples are caused by the stemming of verbs or nouns.

Table 5 also shows that the spaCy-based extraction method retrieves, on average, 3.9 triples per 15 lines. In contrast, on average, the LLMs extract up to 17 triples per paragraph. When considering the format of the extracted triples, with golden extraction dataset suggesting that 8 triples per paragraph as human reference.

A second aspect of extracted triple quality evaluation is formatting. In order to automatically parse the LLM output to use the triples for a downstream analysis, we need the formatting to be respected, which we observe to be only occurring for few-shot learning prompted models. Additionally, we observe that fine-tuning the LLM degrades the percentage of correctly formatted generations for Mistral and Llama. We expect that this is caused by the models being already extensively fine-tuned and being positioned



**Table 5**

Comparative benchmarking of triple extraction methods on the golden dataset. A ( $\uparrow$ ) in the column header means that higher values are preferred, whereas a ( $\downarrow$ ) means lower values are better.

	Avg. number of triples ( $\uparrow$ )	Perc. with correct format ( $\uparrow$ )	Perc. inconsistent with Levenshtein dist. 2 ( $\downarrow$ )	Perc. inconsistent with Levenshtein dist. 3 ( $\downarrow$ )	Avg. time / line (s) ( $\downarrow$ )
Annotated triples	8.1	100%	7.2%	2.6%	-
spaCy extraction	3.9	100%	0%	0%	0.013
Starling - base	-	0%	-	-	0.206
Mistral - base	-	0%	-	-	0.221
Llama - base	8.0	10%	0%	0%	0.097
Starling - base + few-shot	13.0	100%	13.2%	4.6%	0.217
Mistral - base + few-shot	5.3	60%	28%	16%	0.243
Llama - base + few-shot	15.3	100%	7.9%	0.6%	0.104
Starling - fine-tuned + few-shot	10.8	100%	40.8%	24.8%	0.220
Mistral - fine-tuned + few-shot	17.0	75%	50.4%	41.1%	0.244
Llama - fine-tuned + few-shot	13.4	30%	29.8%	21.3%	0.102

on a Pareto frontier, with degeneration kicking in with further fine-tuning [8].

Overall, the Meta-Llama-3-8B-Instruct model with few-shot prompting alone performs best, as it produces a large number of triples that are correctly formatted, while showing few to no signs of hallucinations. The main disadvantage of using a LLM for triple extraction is the computational overhead, as it is 10 times slower than the spaCy-based triple extraction. For this reason, we do not consider even larger LLMs. However, as a sufficiently large number of triples of high quality is essential for downstream performance, we choose to proceed with the Llama model.

## 5.2. Case study part 1: technology forecasting on LLMs through arXiv e-prints

### 5.2.1. Key term extraction and clustering

From the 278,625 arXiv e-prints, we extract a total of 20,822,710 triples, creating to our knowledge the largest unbiased graph of technology-related semantic triples. Following the approach described in Section 4.7, we extract key terms from those papers related to LLMs. Here, we consider the same 278,625 papers from between 2018-2024, from the arXiv categories as specified in Table 1. Specifically, we choose those papers with an abstract that contains at least one of the terms *large language model*, *large language models* or *llm*. The extraction results in a total of 53,337 key terms from 4,182 papers. Next, the key terms are clustered into topics, following the approach that is described in Section 4.7.

### 5.2.2. Domain overview

We start by leveraging the triples to provide a comprehensive overview of current research on large language models (LLMs). For that, we analyze the 20 most prominent topics in the field, presented in Table 6, along with the five most frequently occurring keywords associated with each topic. According to that table, the most prevalent topics are *LLM training* and *Conversational AI*, highlighting the significant focus on optimizing model performance and enhancing human-computer interactions.

Beyond these, several other key areas emerge, including *language understanding*, *question answering*, and *retrieval-augmented generation (RAG)*. Research in *language understanding* focuses on refining LLMs' ability to process, interpret, and generate natural language with greater accuracy and contextual awareness. Question answering remains a critical application of LLMs, driving advancements in knowledge extraction, reasoning, and response generation across various domains. Additionally, RAG has gained traction as a method for improving factual consistency by integrating external knowledge retrieval with generative models, mitigating issues related to hallucination and enhancing the reliability of generated content.

To analyse the landscape in more detail, let us consider the size of the subtopics and their corresponding arXiv categories. Figure 2 shows the 15 largest topics and the corresponding arXiv categories from the papers to which the extracted triples belong. One can see that the most prominent arXiv categories are `cs.LG` (machine learning) and `cs.CL` (computation and language), but that there are also applications in `eess.AS` (audio and speech processing) and `cs.CR` (cryptography and security).

### 5.2.3. Emerging technologies

Next, we consider emerging technologies by performing a time series analysis of the most popular topics. As publications follow yearly conference cycles, it is necessary to decompose the paper counts into a trend and a seasonal component. We therefore employ a simple seasonal decomposition based on moving averages. That is, we model the number of papers  $N_p(t)$  at time  $t$  for topic  $p$  as

$$N_p(t) = T_p(t) + S_p(t) + e_p(t), \quad (7)$$

where  $T_p(t)$  is the trend,  $S_p(t)$  is the seasonal component and  $e_p(t)$  is the residual. Figure 3 shows the trends for each of the ten most common topics, alongside the common seasonal component. We note that the trends are computed using separate seasonal components for each of the topics, but we display a common seasonal component for ease of interpretation.

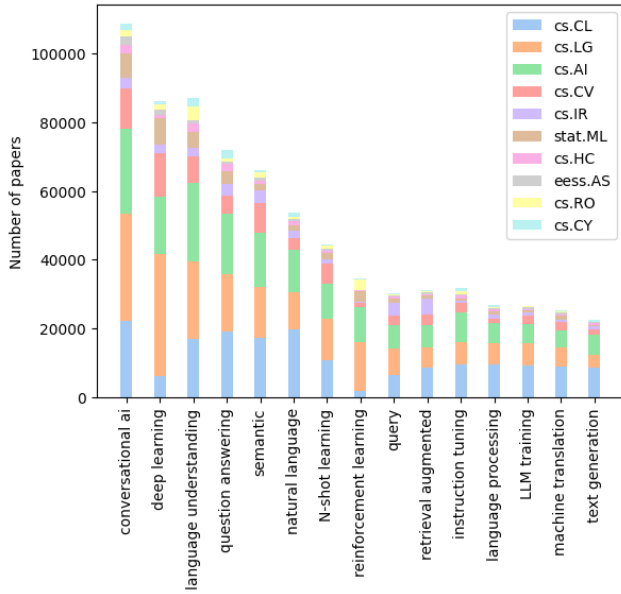
**Table 6**

For each of the 20 largest topics in the field, the 5 most frequent key terms are displayed. The key terms are determined through the method described in Section 4.7.

Category	Top 5 Keywords
LLM training	trained language, training strat, training large, training large language, model pretraining
Conversational AI	generation, generative adversarial, conversation, generating, conversational
Deep Learning	deep learning, deep learning models, art deep, deep learning architectures, art deep learning
Language Understanding	understanding, planning, language understand, language understanding, language understandin
Question Answering	question, answering, question answer, question answering, question answerin
Semantic	semantic, semantic communication, automatic summarization, image semantic, semantic evaluation
Natural Language	natural language, natural language pro, natural language process, natural language processing, natural language inference
Retrieval Augmented	retrieval, retrieval augmented, retrieval augmented generation, generative retrieval, retrieval augmentation
Instruction Tuning	instruction, instructions, instruction following, instruction gen, based instruction
Language Processing	language processing, prompting large language, grounding language robotic, language modeling pretraining, prompted language
Machine Translation	translation, machine transla, machine translation, neural machine translation, translation model
Text Generation	generated text, al language generation, text generation, conversational agent, generation models
Video	videos, video generation, based video, video models, generated videos
Summarization	summarization, narrative, summarization model, summarization models, abstractive text
ChatGPT	chatgpt, based chatbot, chatgpt model, openai chatgpt, chatgpt1
Multilingual	multilingual, multilingual bert, multilingual lan, multimodal large, multilingual large
Dialogue	dialogue, dialogue evaluation, dialogue task, dialogue tasks, dialogue modeling
Speech	text model, speech recognition models, text speech, speech text, based speech
N-Shot Learning	shot learning, shot learner, shot learners, zero shot, shot learning methods
Source Code	source code, code data, source code data, code test, source code sequence
Time Series	time series, time series forecasting, time series prediction, multivariate time series, time series dataset

When considering the seasonal component, one can see two spikes a year, one around May-June and one in October. These spikes correspond to the biggest conferences for NLP related work: the Association for Computational Linguistics (ACL) and Empirical Methods in Natural Language Processing (EMNLP). The ACL conference is commonly around July, with submissions generally uploaded in the months before. The EMNLP conference is in November, with submissions uploaded to arXiv beforehand.

When considering the trends of the various topics around LLMs on Figure 3, one can see that, in general, there is a surge in interest for topics in LLMs in early 2023. This aligns with the release of ChatGPT in November 2022, which caused a surge in public interest in LLMs. Furthermore, this event seems to have had little effect in the popularity of older technologies such as *deep learning* and *reinforcement learning*. In contrast, the biggest emerging technologies are related to *conversational ai* and *language understanding*. This shows that the current focal point in



**Figure 2:** The number of papers for each of the 15 most common topics in the field of LLMs, based on the extracted and grouped key terms. Each of the bars is divided into the arXiv categories from which the papers originate.

the field is the simplified interaction between humans and LLMs.

For a more fine-grained analysis, we investigate the trends of individual key terms within the topics. Again, we distinguish between a seasonal component and the trend. To identify those key terms with the most interesting trends, we select the 5 key terms for each category that show the largest increase in prevalence between 2018 and 2024. Figure 4 displays the trends for these key terms for each topic, alongside the aggregate seasonal components. The results show that there are specific research directions that surged over the past years, most notably *evidential deep learning*, *semantic communication*, *reinforcement learning with human feedback* and *query rewriting*, corresponding to emergent technologies. Such topics are all driven by advances in the capabilities of LLMs and are likely to continue growing in the near future.

#### 5.2.4. Technology Convergence analysis

The technology convergence analysis is the heart of our approach, and we perform it through topological analysis of semantic triple graph linking together technology-designating terms. We perform a network analysis, where each node represents a key topics and the edges represent the number of triplets connecting the topics, with edges mapping both to the connecting verb and the paper from which the triple was extracted. The edges are then aggregated, leading to a single weighted undirected edge, with weight proportional to triple occurrence. We also assign a "color" to the edge, based on whether the occurrence of triples involving the two technologies has been increasing (red), decreasing (blue) or remained consistent (black).

**Technology clusters** We identified established related technology clusters by partitioning the network into densely connected sub-networks, by using the Louvain community detection method [14]<sup>4</sup>. Figure 5 shows the six largest detected clusters, with node color indicating the cluster. The size of nodes corresponds to the frequency of technology-related term occurrence, while the size of the term itself - to the node eigenvector centrality on the entire network. The eigenvector centrality corresponds how important the node is to connecting the entire network.

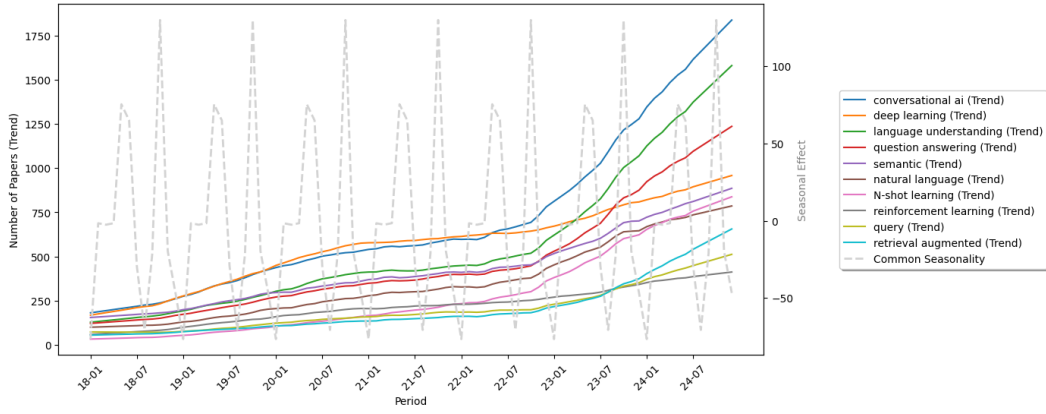
**Technology convergences** In order to identify emerging technological convergences, we then focus on the technology-related terms connecting different technology clusters. Figure 6 shows three networks, the first presenting emerging technological connections, the second showing waning ones, and the last one showing persisting ones. We observe an emerging technological convergence between *instruction tuning* and *natural language*. This highlights the surge in conversational AI that we discussed previously, which was unlocked by instruction-tuned LLMs, that allowed general public access to conversational LLMs and made them useful as a human assistant emulators [58]. Furthermore, we see pivotal roles for *prompt engineering* and *retrieval augmentation*. These relatively new branches of research aim at improving the performance of LLMs. For instance, we see that retrieval augmentation is linked to information retrieval, illustrating its goal of improving the output of LLMs by injecting data from an existing database.

When we consider the relations that declined in prevalence between 2018 and 2024, we most notably see a link between *language understanding* and *fact checking*. This seems to suggest that fact checking is increasingly less related to language understanding. One hypothesis is that fact checking is now more often performed through stylistic patterns, leveraging ML methods.

Finally, we consider persistent relations between 2018 and 2024, which often involve core concepts in LLMs. For instance, *question answering* links with *language understanding* and *reasoning capabilities*. Similarly, *retrieval augmentation*, introduced early, has remained stable in its connection to *question answering*.

**Temporal trends of convergence** Last, we analyze the temporal aspect of convergence by considering the topic co-publication frequency. To analyse which technologies are converging, we study the pairs of topics that have the largest increase in the Jaccard similarity between 2018 and 2024. The Jaccard similarity measure is defined in equation 6. Figure 7 shows the evolution of the Jaccard similarities for these pairs of topics. One can see that from the release of ChatGPT in November 2022, several technology convergences emerge rapidly. However, it can also be observed that in 2024 several of such convergences start to decline again, indicating that the effect was only temporary. In contrast, one can see that the convergence

<sup>4</sup>For legibility reasons we are unable to provide the entire network



**Figure 3:** Multiple line plot showing for each of the 10 most common topics the number of papers in which at least one triple appears for that topic. The gray dashed line shows the aggregate seasonal component, whereas the colored lines show the trend for each topic.

between *retrieval augmented* and *conversational agent* is stronger and is not yet stagnating. In combination with the network analysis that highlighted these two topics as technological convergences as well, we can designate *retrieval augmentation* and *conversational agents* as emerging transformative technologies in the field of LLMs, as of end 2024, based on research preprints published on arXiv.

### 5.3. Case study part 2: technology forecasting on LLMs through USPTO patent applications

In order to demonstrate the generalization of our approach and its usefulness in the fields with less emphasis on early preprint publication, we use USPTO patent data to provide a different perspective on the technological developments related to LLMs. From the 9,793 patent applications, our pipeline extracts 3,027,121 triples after post-processing and filtering. As before, to perform an informative downstream analysis, we require a set of key terms that are most relevant to the field of LLMs. For this purpose, we extract key terms from the abstracts of the 9,793 patents, following the same methodology as described in Section 4.7.

We focus on two analyses to showcase the unique perspective that patents give on the technological developments surrounding LLMs compared to arXiv preprints. As the primary goal of this section is to highlight the generalizability of our method across different data sources, we do not cover all of the perspectives that were taken using the arXiv papers, but merely highlight two approaches.

#### 5.3.1. Emerging technologies

We start by performing a technology-related term usage analysis similar to Section 5.2.3, where we are looking for terms with recently increasing usage. Whereas papers uploaded on arXiv follow a strong seasonal pattern due to yearly conference cycles, this pattern is not present for patent applications, due to the lack of yearly conference deadlines.

Figure 8 shows the most popular emerging terms associated with LLMs in patent applications. Across the topics,

there is a surge in LLM-related technology terms halfway through 2024. This rise in frequency of term usage tracks the rise of similar terms usage on arXiv at the end of 2022, suggesting the analysis of scientific articles preprints allows an early trend emergence identification.

We see that LLMs are connected to a variety of fields, for instance *computer storage*, *cloud computing*, *source code*, and *speech recognition*. Furthermore, we note that the rise in patent applications has only started recently, and it will be interesting to see if future trends follow the patterns identified in the temporal analysis of arXiv preprints, with some rapidly growing technology topics plateauing shortly.

#### 5.3.2. Network analysis

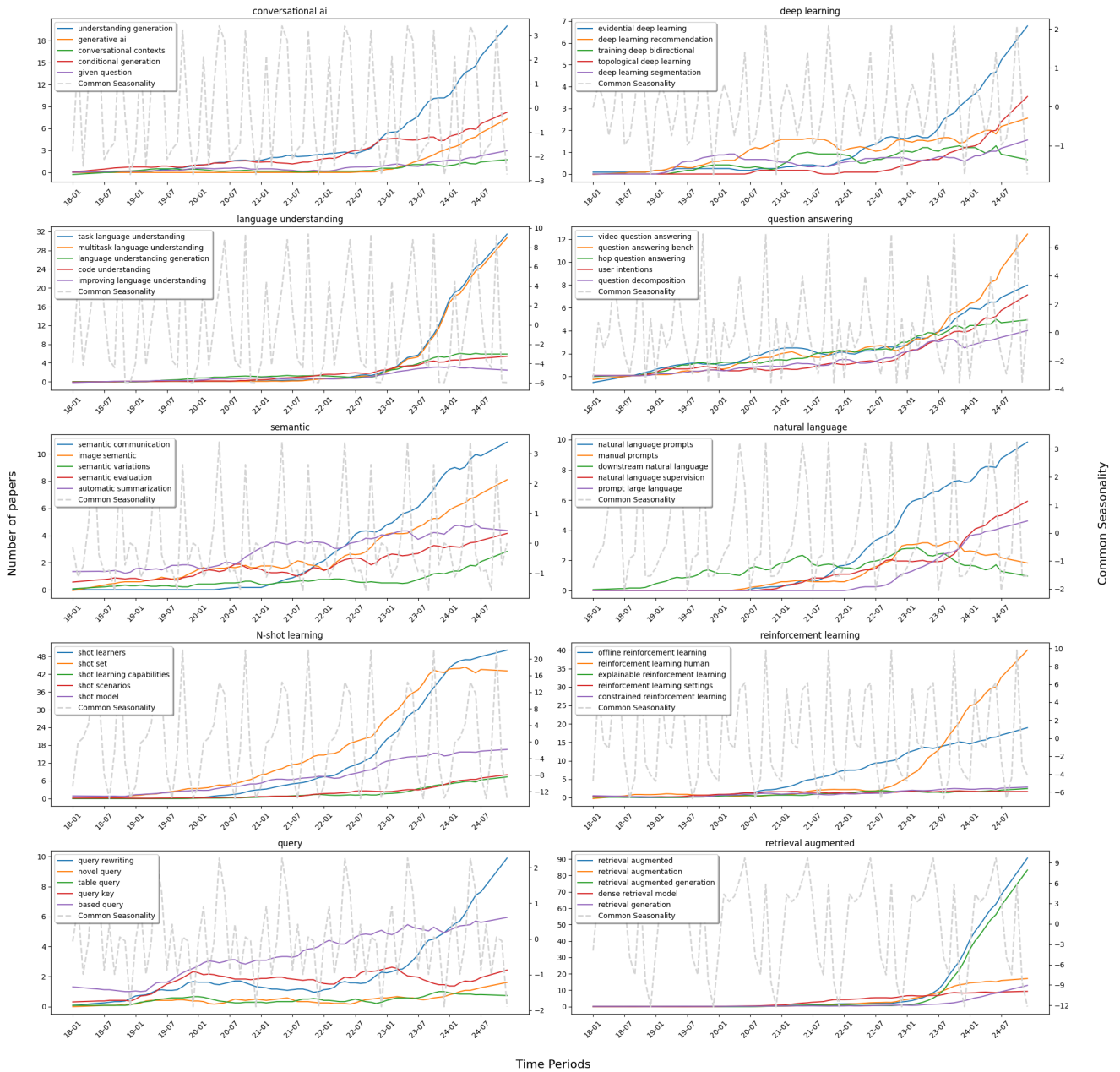
Just as for arXiv articles, we analyze technological convergence through technology-related terms semantic graph analysis. Figure 9 displays the network, with node and edge color and size corresponding to the same properties as for arXiv semantic graph analysis in section 5.2.4.

From this visualization, we can make several observations regarding the developments in industry applications of LLMs. First, we see that there is a significant overlap between LLMs and self-attention models, covering multi-modal models, visual attention models and speech attention models. For waning topic connections, we see that there is a decline in image processing connection to blockchain and payment processing, likely corresponding to the deflation of the non-fungible tokens (NFT) bubble; as well as to self-driving vehicle and unmanned aerial vehicles (UAVs), likely corresponding to a switch to other types of sensors and manual controls. The rupture of the link between speech recognition and speech recognition server seems to indicate the move of speech recognition onto devices.

When considering emerging topics, one can observe that there is a rise in interest for knowledge graphs and knowledge bases. Specifically, this seems to be closely related with the instruction-tuning of LLMs. Furthermore, one can see a rise in attention to source code manipulation, bridging the classical abstract syntax tree analysis with LLM-based



## Monitoring Transformative Technological Convergence



**Figure 4:** The trends for key sub-technology terms for the 10 most common emerging technology topics in the LLM space, based on the extracted and grouped key terms. The gray dashed line shows the aggregate seasonal component that the key terms share, whereas the colored lines show the trend for each key term.

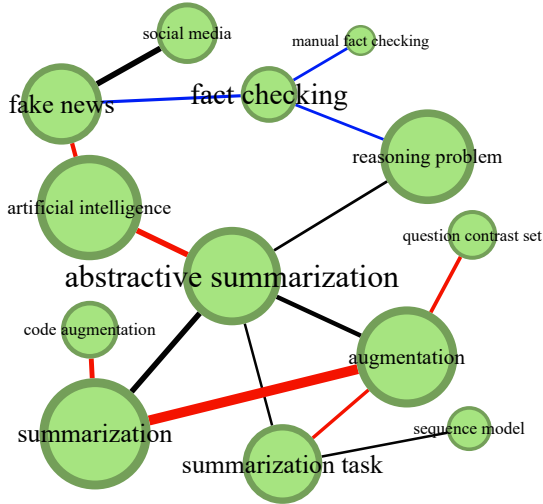
machine learning methods, corresponding to the trend of LLM-based code generation and analysis integration into programming workflows.

Based on the USPTO patent applications, our analysis suggests that *coding LLMs* and *on-device speech recognition* with self-attention models are emerging transformative topics in the field of LLMs and self-attention models, as of end 2024.

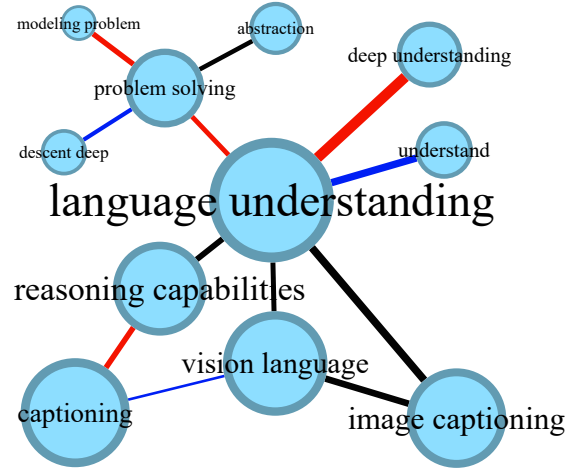
## 6. Conclusion

The monitoring and forecasting of emerging transformative technologies is a notoriously important and yet difficult topic in the field of technological forecasting. The accelerated accumulation and availability of potentially relevant data presents an opportunity for more grounded forecasts, but also poses a challenge due to the scale and presence of noise. We developed and presented in this paper a pipeline that leverages advances in NLP and big data analysis to perform a technology convergence analysis while addressing the challenges that previously made such analysis at

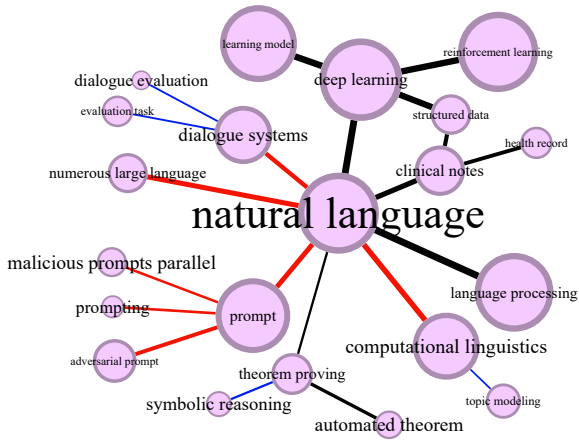
## Monitoring Transformative Technological Convergence



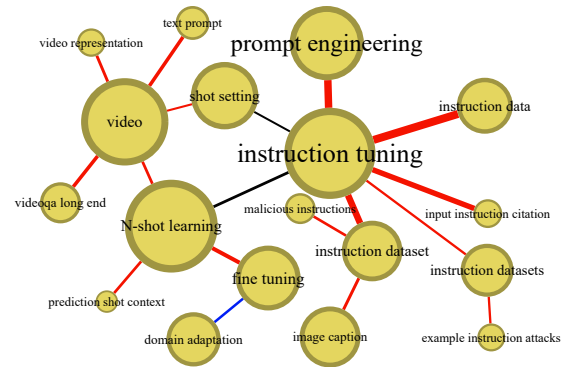
(a) Summarization and factuality



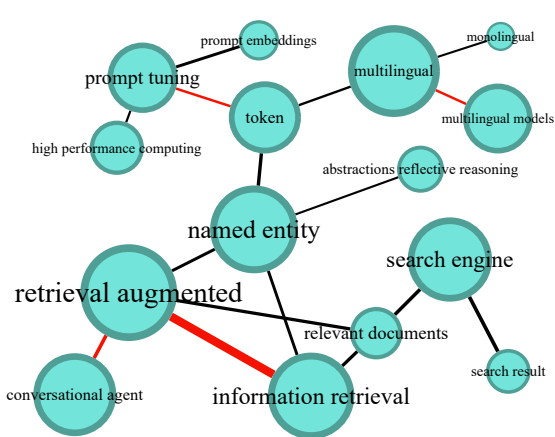
(b) Language and image understanding



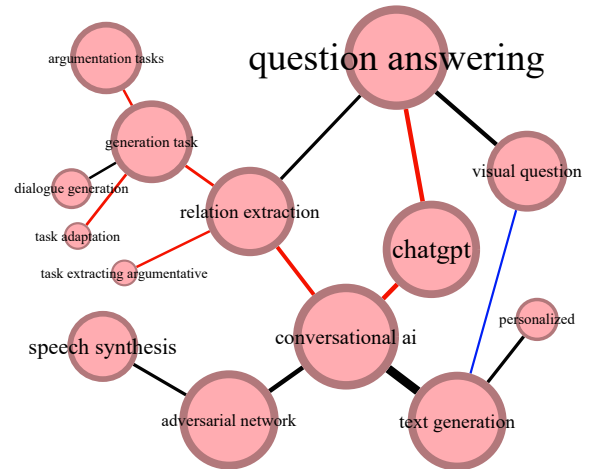
(c) Core NLP concepts



(d) Instruction-tuning and prompting

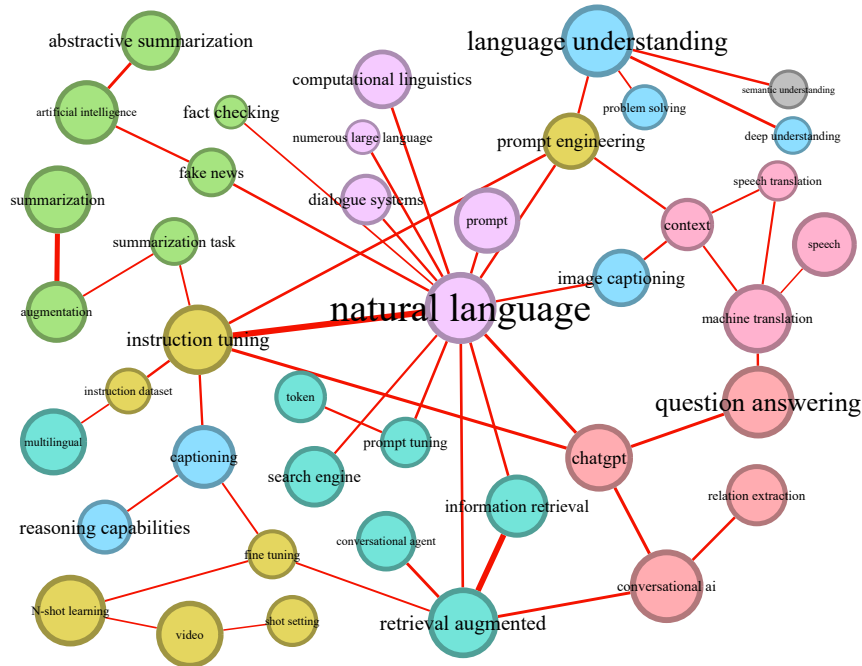


(e) Retrieval augmentation

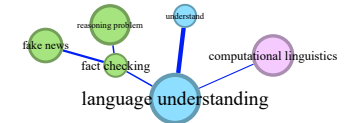


(f) Question answering

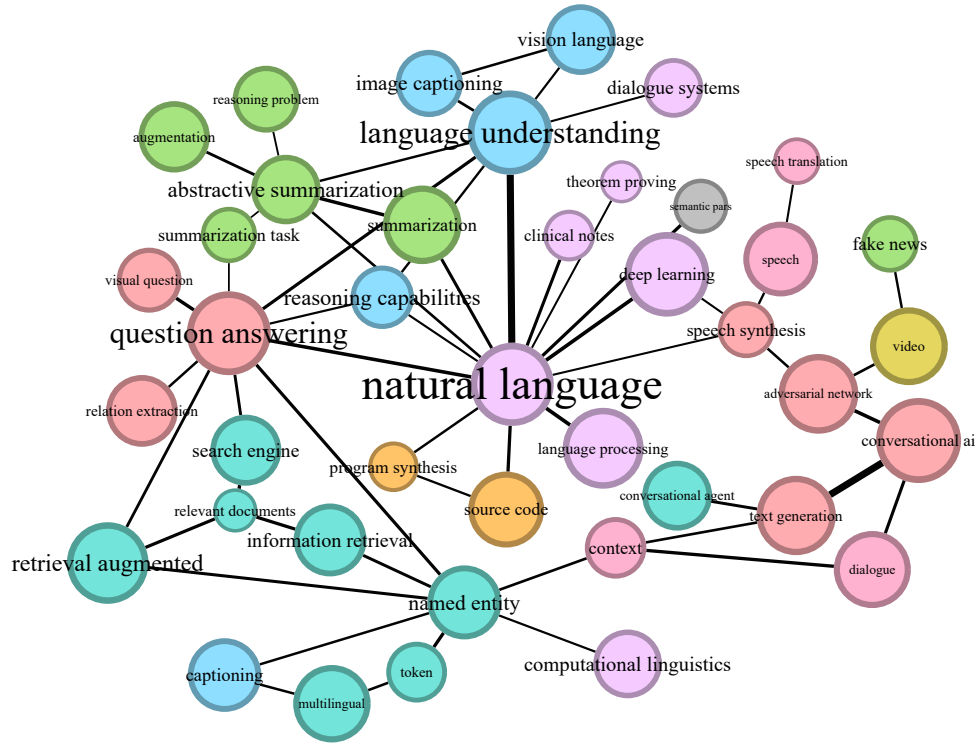
**Figure 5:** Clusters of topics composed using the Louvain method for community detection. Relations in red occur over 70% of the time in 2022 or later, whereas relations in blue occur over 70% of the time in 2021 or earlier. All other relations are in black. Node/edge sizes correspond to absolute frequency, and text size to eigenvector network centrality.



(a) Emerging connections between key topics between 2018 and 2024.

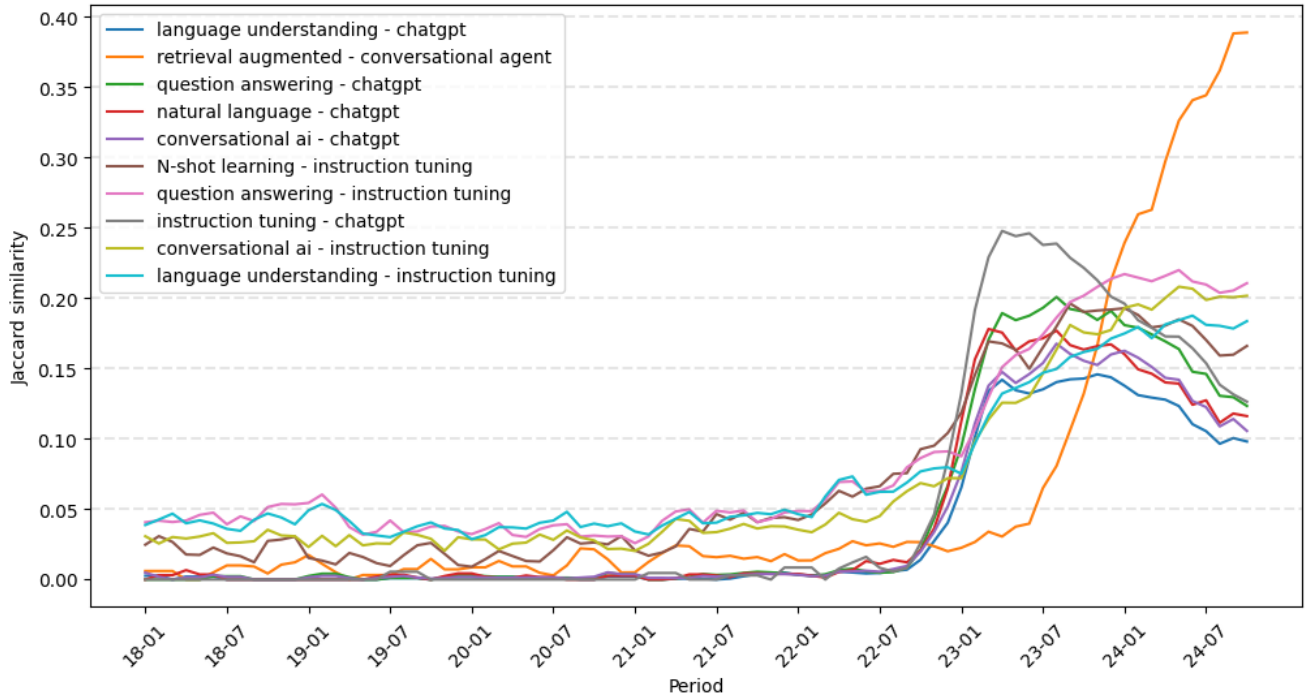


(b) Connections between key topics that are losing prevalence between 2018 and 2024.

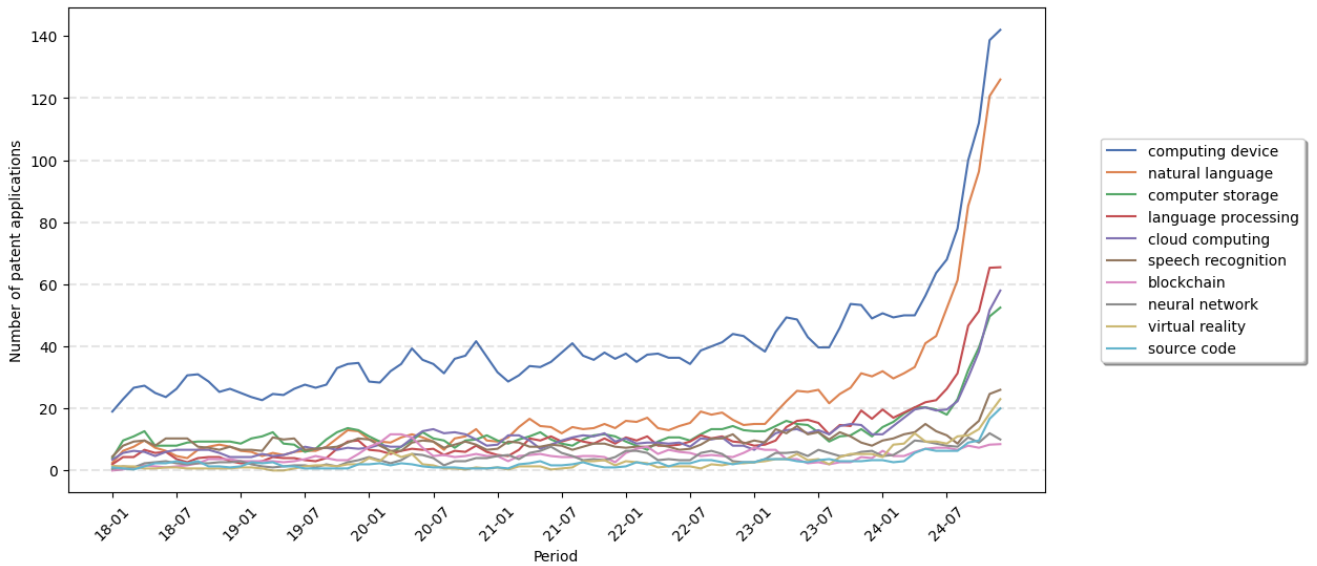


(c) Persistent connections between key topics between 2018 and 2024.

**Figure 6:** Emerging (red), disappearing (blue), and persistent (black) relations between technology-related term clusters based on the arXiv paper data. Relations in red occur over 70% of the time in 2022 or later, whereas relations in blue occur over 70% of the time in 2021 or earlier. All other relations are displayed as black. The node color reflects the term cluster they belong to, the sizes of the nodes and edges reflect the frequency of their appearances, and the size of the text labels reflects the eigenvector network centrality.



**Figure 7:** The Jaccard similarities for the 10 pairs of topics for which the Jaccard similarity had the largest increase between in the Jaccard similarity between 2018 and 2024. We consider all topics, as identified by the extracted and grouped key terms in for the topic of LLMs.



**Figure 8:** Multiple line plot showing for each of the 10 most common topics, as identified by the extracted and grouped key terms, the number of patent applications in which at least one triple appears for that topic.

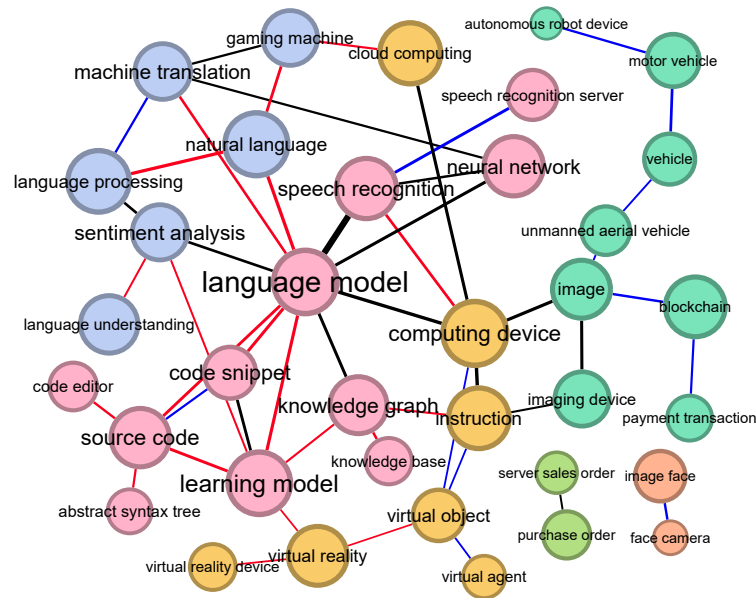
scale impossible. Specifically, we (I) leveraged the recent advances in NLP and recently published benchmarking datasets for to extract information relevant to technological forecasting in an unbiased structured format, and at a scale that has been previously unattainable; (II) developed a new approach to account for variation in technologies designation through varying terms allowing their systematic analysis; (III) constructed a novel graph-based technique for

identification of emerging technologies and likely transformative emerging technologies for specific fields.

We then use the pipeline we developed in order to analyze the trends in the LLM-related technology field based on two different unstructured data sources on a previously unseen scale and granularity. Specifically, we perform a full-text analysis of 278,625 arXiv scientific preprints and 9,793 USPTO patent applications, resulting



## Monitoring Transformative Technological Convergence



**Figure 9:** Network representation of the key topics based on the USPTO patent application data. Relations in red are occurring over 70% of the time in 2022 or later, whereas relations in blue occur over 70% of the time in 2021 or earlier. All other relations are displayed as black. The node color corresponds to the technology term cluster they belong to, the sizes of the nodes and edges correspond to the absolute frequency of their appearances, and the size of the text labels corresponds to the eigenvector network centrality.

in 53,337 LLM technology-related terms and 23,849,831 triples. By leveraging big data techniques on the resulting semantic graph analysis, we were able to identify two emerging transformative technologies as of 2024: *retrieval-augmented generation* and *conversational agents*, which seem to be supported by the data-grounded agentic AI interest in 2025; as well as two emerging transformative applications: *source-code generation and processing* with LLMs and *on-device speech processing with multimodal models*.

We leave to future research the task of validating the identified transformative technologies, assessing the pipeline's use with other data sources — such as social media, blogs, and grey literature like reports and white papers — and examining its longer-term predictive capabilities.

## References

- [1] . . arXiv preprint server. <https://info.arxiv.org>. Accessed: 2025-04-28.
- [2] . . biorXiv preprint server. <https://www.biorxiv.org/content/about-biorxiv>. Accessed: 2025-04-28.
- [3] . . medRxiv preprint server. <https://www.medrxiv.org/content/about-medrxiv>. Accessed: 2025-04-28.
- [4] . . PubMed Central online publications archive. <https://pmc.ncbi.nlm.nih.gov/about/intro/>. Accessed: 2025-04-28.
- [5] Adamuthe, A.C., Thampi, G.T., 2019. Technology forecasting: A case study of computational technologies. *Technological Forecasting and Social Change* 143, 181–189. URL: <https://www.sciencedirect.com/science/article/pii/S0040162518302890>, doi:<https://doi.org/10.1016/j.techfore.2019.03.002>.
- [6] AI@Meta, 2024. Llama 3 model card URL: [https://github.com/meta-llama/llama3/blob/main/MODEL\\_CARD.md](https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md).
- [7] Artifex, 2024. Pymupdf. <https://pypi.org/project/PyMuPDF/>. Accessed: 2024-02-29.
- [8] Bai, Y., Kadavath, S., Kundu, S., Askell, A., Kernion, J., Jones, A., Chen, A., Goldie, A., Mirhoseini, A., McKinnon, C., Chen, C., Olsson, C., Olah, C., Hernandez, D., Drain, D., Ganguli, D., Li, D., Tran-Johnson, E., Perez, E., Kerr, J., Mueller, J., Ladish, J., Landau, J., Ndousse, K., Lukosuite, K., Lovitt, L., Sellitto, M., Elhage, N., Schiefer, N., Mercado, N., DasSarma, N., Lasenby, R., Larson, R., Ringer, S., Johnston, S., Kravec, S., Showk, S.E., Fort, S., Lanham, T., Telleen-Lawton, T., Conerly, T., Henighan, T., Hume, T., Bowman, S.R., Hatfield-Dodds, Z., Mann, B., Amodei, D., Joseph, N., McCandlish, S., Brown, T., Kaplan, J., 2022. Constitutional ai: Harmlessness from ai feedback. *arXiv:2212.08073*.
- [9] Bastian, M., Heymann, S., Jacomy, M., 2009. Gephi: An open source software for exploring and manipulating networks. URL: <http://www.aaa.org/ocs/index.php/ICWSM/09/paper/view/154>.
- [10] Bastos, A., Nadgeri, A., Singh, K., Mulang, I.O., Shekarpour, S., Hof-fart, J., Kaul, M., 2021. Recon: Relation extraction using knowledge graph context in a graph neural network, in: *Proceedings of the Web Conference 2021*, Association for Computing Machinery, New York, NY, USA. p. 1673–1685. URL: <https://doi.org/10.1145/3442381.3449917>, doi:10.1145/3442381.3449917.
- [11] Beltagy, I., Lo, K., Cohan, A., 2019. Scibert: Pretrained language model for scientific text, in: *EMNLP*. *arXiv:arXiv:1903.10676*.
- [12] Berners-Lee, T., Hendler, J., Lassila, O., 2001. The semantic web: A new form of web content that is meaningful to computers will unleash a revolution of new possibilities. *ScientificAmerican.com*.
- [13] Blei, D.M., Ng, A.Y., Jordan, M.I., 2003. Latent dirichlet allocation. *J. Mach. Learn. Res.* 3, 993–1022.
- [14] Blondel, V.D., Guillaume, J.L., Lambiotte, R., Lefebvre, E., 2008. Fast unfolding of communities in large networks. *Journal of statistical mechanics: theory and experiment* 2008, P10008.

- [15] Brown, T.B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D.M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., Amodei, D., 2020. Language models are few-shot learners, in: Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., Lin, H. (Eds.), *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*. URL: <https://proceedings.neurips.cc/paper/2020/hash/1457c0d6bfc4967418bfb8ac142f64a-Abstract.html>.
- [16] Calleja-Sanz, G., Nadal, J., Solé-Parellada, F., 2020. Technology Forecasting: Recent Trends and New Methods. pp. 45–69. doi:10.1007/978-3-030-40896-1\_3.
- [17] Callon, M., Courtial, J.P., Laville, F., 1983. From translations to problematic networks: An introduction to co-word analysis. *Social Science Information* 22, 191–235.
- [18] Chawla, D.S., 2022. Massive open index of scholarly papers launches. *Nature* URL: <https://api.semanticscholar.org/CorpusID:246278314>.
- [19] Christensen, C.M., 1997. *The innovator's dilemma: When new technologies cause great firms to fail*. Harvard Business School Press, Boston.
- [20] Cohan, A., Feldman, S., Beltagy, I., Downey, D., Weld, D.S., 2020. Specter: Document-level representation learning using citation-informed transformers. URL: <https://arxiv.org/abs/2004.07180>, arXiv:2004.07180.
- [21] Council, N.R., on Engineering, D., Sciences, P., on Forecasting Future Disruptive Technologies, C., 2010. *Persistent Forecasting of Disruptive Technologies' Report 2*. National Academies Press.
- [22] Daim, T., Yalçın, H., 2022. *Digital transformations: new tools and methods for mining technological intelligence*. Edward Elgar Publishing.
- [23] Daim, T.U., Chiavetta, D., Porter, A.L., Saritas, O., 2016. *Anticipating future innovation pathways through large data analysis*. Springer.
- [24] Dalkey, N., 1969. An experimental study of group opinion: the delphi method. *Futures* 1, 408–426.
- [25] Dernoncourt, F., Lee, J.Y., 2017. PubMed 200k RCT: a dataset for sequential sentence classification in medical abstracts, in: Kondrak, G., Watanabe, T. (Eds.), *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, Asian Federation of Natural Language Processing, Taipei, Taiwan. pp. 308–313. URL: <https://aclanthology.org/I17-2052>.
- [26] Dosi, G., 1982. Technological paradigms and technological trajectories: A suggested interpretation of the determinants and directions of technical change. *Research Policy* 11, 147–162. URL: <https://api.semanticscholar.org/CorpusID:16840352>.
- [27] Dotsika, F., Watkins, A., 2017. Identifying potentially disruptive trends by means of keyword network analysis. *Technological Forecasting and Social Change* 119, 114–127. URL: <https://www.sciencedirect.com/science/article/pii/S0040162517303517>, doi:https://doi.org/10.1016/j.techfore.2017.03.020.
- [28] Eberts, M., Ulges, A., 2020. Span-based joint entity and relation extraction with transformer pre-training, in: Giacomo, G.D., Catalá, A., Dilkina, B., Milano, M., Barro, S., Bugarín, A., Lang, J. (Eds.), *ECAI 2020 - 24th European Conference on Artificial Intelligence*, 29 August–8 September 2020, Santiago de Compostela, Spain, August 29 - September 8, 2020 - Including 10th Conference on Prestigious Applications of Artificial Intelligence (PAIS 2020), IOS Press. pp. 2006–2013. URL: <https://doi.org/10.3233/FAIA200321>, doi:10.3233/FAIA200321.
- [29] Ettlie, J.E., Bridges, W.P., O'keefe, R.D., 1984. Organization strategy and structural differences for radical versus incremental innovation. *Management Science* 30, 682–695. URL: <https://api.semanticscholar.org/CorpusID:154030756>.
- [30] Falissard, L., Guigue, V., Soulier, L., 2023. Improving generalization in large language model by learning prefix subspaces, in: Bouamor, H., Pino, J., Bali, K. (Eds.), *Findings of the Association for Computational Linguistics: EMNLP 2023*, Association for Computational Linguistics, Singapore. pp. 11474–11483. URL: <https://aclanthology.org/2023.findings-emnlp.768/>, doi:10.18653/v1/2023.findings-emnlp.768.
- [31] Fergadis, A., Pappas, D., Karamolegkou, A., Papageorgiou, H., 2021. Argumentation mining in scientific literature for sustainable development, in: Al-Khatib, K., Hou, Y., Stede, M. (Eds.), *Proceedings of the 8th Workshop on Argument Mining*, Association for Computational Linguistics, Punta Cana, Dominican Republic. pp. 100–111. URL: <https://aclanthology.org/2021.argmining-1.10>, doi:10.18653/v1/2021.argmining-1.10.
- [32] Fisher, J.C., Pry, R.H., 1971. A simple substitution model of technological change. *Technological forecasting and social change* 3, 75–88.
- [33] Gali, N., Mariescu-Istodor, R., Hostettler, D., Fränti, P., 2019. Framework for syntactic string similarity measures. *Expert Systems with Applications* 129, 169–185. URL: <https://www.sciencedirect.com/science/article/pii/S0957417419302222>, doi:https://doi.org/10.1016/j.eswa.2019.03.048.
- [34] Grootendorst, M., 2020. Keybert: Minimal keyword extraction with bert. URL: <https://doi.org/10.5281/zenodo.4461265>, doi:10.5281/zenodo.4461265.
- [35] Group, T.F.A.M.W., et al., 2004. Technology futures analysis: Toward integration of the field and new methods. *Technological Forecasting and Social Change* 71, 287–303.
- [36] Gu, M., Liu, X., Zhang, L., Wang, J., 2024. Forecasting technology emergence using large-scale knowledge graphs and link prediction. *Scientometrics* In press.
- [37] Hogan, A., Blomqvist, E., Cochez, M., D'amato, C., Melo, G.D., Gutierrez, C., Kirrane, S., Gayo, J.E.L., Navigli, R., Neumaier, S., Ngomo, A.C.N., Polleres, A., Rashid, S.M., Rula, A., Schmelzeisen, L., Sequeda, J., Staab, S., Zimmermann, A., 2021. Knowledge graphs. *ACM Computing Surveys* 54, 1–37. URL: <http://dx.doi.org/10.1145/3447772>, doi:10.1145/3447772.
- [38] Honnibal, M., Montani, I., Van Landeghem, S., Boyd, A., 2020. spaCy: Industrial-strength Natural Language Processing in Python doi:10.5281/zenodo.1212303.
- [39] Hu, E.J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W., 2021. Lora: Low-rank adaptation of large language models. arXiv:2106.09685.
- [40] Jansen, T., Kuhn, T., 2017. Extracting core claims from scientific articles. CoRR abs/1707.07678. URL: <http://arxiv.org/abs/1707.07678>, arXiv:1707.07678.
- [41] Jantsch, E., 1967. Technological forecasting in perspective. Number 21,931 in *OECD Publications*, OECD, Paris.
- [42] 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. arXiv:2310.06825.
- [43] Kim, G., Bae, J., 2017. A novel approach to forecast promising technology through patent analysis. *Technological Forecasting and Social Change* 117, 228–237. URL: <https://www.sciencedirect.com/science/article/pii/S0040162516307661>, doi:https://doi.org/10.1016/j.techfore.2016.11.023.
- [44] Kodama, F., 1992. Technology fusion and the new r&d. *Harvard business review*, 70–78.
- [45] Krenn, M., Bollen, J., Stojanovski, J., et al., 2023. Forecasting research trends using dynamic graph representations: The science4cast 2023 challenge. arXiv preprint arXiv:2301.03589.

- [46] Kucharavy, A., 2024. From deep neural language models to llms, in: *Large Language Models in Cybersecurity: Threats, Exposure and Mitigation*. Springer, pp. 3–17.
- [47] Kucharavy, A., Schillaci, Z., Maréchal, L., Würsch, M., Dolamic, L., Sabonnadiere, R., David, D.P., Mermoud, A., Lenders, V., 2023. Fundamentals of generative large language models and perspectives in cyber-defense. *arXiv:2303.12132*.
- [48] Li, X., Burns, G.A., Peng, N., 2021. Scientific discourse tagging for evidence extraction, in: Merlo, P., Tiedemann, J., Tsarfaty, R. (Eds.), *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, EACL 2021, Online, April 19 - 23, 2021, Association for Computational Linguistics*. pp. 2550–2562. URL: <https://doi.org/10.18653/v1/2021.eacl-main.218>, doi:10.18653/v1/2021.EACL-MAIN.218.
- [49] Li, X., Fan, M., Zhou, Y., Fu, J., Yuan, F., Huang, L., 2020. Monitoring and forecasting the development trends of nanogenerator technology using citation analysis and text mining. *Nano Energy* 71, 104636. URL: <https://www.sciencedirect.com/science/article/pii/S2211285520301932>, doi:<https://doi.org/10.1016/j.nanoen.2020.104636>.
- [50] Li, X., Wang, Y., 2024. A novel integrated approach for quantifying the convergence of disruptive technologies from science to technology. *Technological Forecasting and Social Change* URL: <https://api.semanticscholar.org/CorpusID:273569421>.
- [51] Liu, L., Omidvar, A., Ma, Z., Agrawal, A., An, A., 2022. Un-supervised knowledge graph generation using semantic similarity matching, in: Cherry, C., Fan, A., Foster, G., Haffari, G.R., Khadivi, S., Peng, N.V., Ren, X., Shareghi, E., Swayamdipta, S. (Eds.), *Proceedings of the Third Workshop on Deep Learning for Low-Resource Natural Language Processing, Association for Computational Linguistics, Hybrid*. pp. 169–179. URL: <https://aclanthology.org/2022.deeplo-1.18>, doi:10.18653/v1/2022.deeplo-1.18.
- [52] Martínez, V., Berzal, F., Cubero, J.C., 2016. A survey of link prediction in complex networks. *ACM Computing Surveys (CSUR)* 49, 1–33.
- [53] Melamud, O., Goldberger, J., Dagan, I., 2016. context2vec: Learning generic context embedding with bidirectional LSTM, in: Riezler, S., Goldberg, Y. (Eds.), *Proceedings of the 20th SIGNLL Conference on Computational Natural Language Learning, Association for Computational Linguistics, Berlin, Germany*. pp. 51–61. URL: <https://aclanthology.org/K16-1006/>, doi:10.18653/v1/K16-1006.
- [54] Michel, J.B., Shen, Y.K., Aiden, A.P., Veres, A., Gray, M.K., Team, G.B., Pickett, J.P., Hoiberg, D., Clancy, D., Norvig, P., et al., 2011. Quantitative analysis of culture using millions of digitized books. *science* 331, 176–182.
- [55] Mikolov, T., Sutskever, I., Chen, K., Corrado, G.S., Dean, J., 2013. Distributed representations of words and phrases and their compositionality, in: *Neural Information Processing Systems*. URL: <https://api.semanticscholar.org/CorpusID:16447573>.
- [56] Neumann, M., King, D., Beltagy, I., Ammar, W., 2019. ScispaCy: Fast and robust models for biomedical natural language processing, in: Demner-Fushman, D., Cohen, K.B., Ananiadou, S., Tsujii, J. (Eds.), *Proceedings of the 18th BioNLP Workshop and Shared Task, Association for Computational Linguistics, Florence, Italy*. pp. 319–327. URL: <https://aclanthology.org/W19-5034>, doi:10.18653/v1/W19-5034.
- [57] Ottersen, S.G., Pinheiro, F., Bação, F., 2024. Triplet extraction leveraging sentence transformers and dependency parsing. *Array* 21, 100334. URL: <https://www.sciencedirect.com/science/article/pii/S2590005623000590>, doi:<https://doi.org/10.1016/j.array.2023.100334>.
- [58] Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al., 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems* 35, 27730–27744.
- [59] Percia David, D., Maréchal, L., Lacube, W., Gillard, S., Tsemmelis, M., Maillart, T., Mermoud, A., 2023. Measuring security development in information technologies: A scientometric framework using arxiv e-prints. *Technological Forecasting and Social Change* 188, 122316. URL: <https://www.sciencedirect.com/science/article/pii/S004016252300001X>, doi:<https://doi.org/10.1016/j.techfore.2023.122316>.
- [60] Qiu, Z., Wang, Z., 2022. Technology forecasting based on semantic and citation analysis of patents: A case of robotics domain. *IEEE Transactions on Engineering Management* 69, 1216–1236. doi:10.1109/TEM.2020.2978849.
- [61] Radford, A., Narasimhan, K., Salimans, T., Sutskever, I., et al., 2018. Improving language understanding by generative pre-training.
- [62] Schumpeter, J.A., 1949. Economic theory and entrepreneurial history, in: Clemence, R. (Ed.), *Essays on Entrepreneurs, Innovations, Business Cycles, and the Evolution of Capitalism*. New Jersey: Transaction Publishers, pp. 272–286.
- [63] Schwartz, A.S., Hearst, M.A., 2003. A simple algorithm for identifying abbreviation definitions in biomedical text, in: Altman, R.B., Dunker, A.K., Hunter, L., Klein, T.E. (Eds.), *Proceedings of the 8th Pacific Symposium on Biocomputing, PSB 2003, Lihue, Hawaii, USA, January 3-7, 2003*, pp. 451–462. URL: <http://psb.stanford.edu/psb-online/proceedings/psb03/schwartz.pdf>.
- [64] Smith, L.C., 1981. Citation analysis.
- [65] Sternfeld, A., Kucharavy, A., David, D.P., Mermoud, A., Jang-Jaccard, J., 2024. Llm-resilient bibliometrics: Factual consistency through entity triplet extraction, in: Zhang, C., Zhang, Y., Mayr, P., Lu, W., Suominen, A., Chen, H., Ding, Y. (Eds.), *Proceedings of Joint Workshop of the 5th Extraction and Evaluation of Knowledge Entities from Scientific Documents (EEKE2024) and the 4th AI + Informetrics (AII2024) co-located with the (iConference2024), Changchun, China and Online, April 23-24, 2024*, CEUR-WS.org. pp. 85–93. URL: <https://ceur-ws.org/Vol-3745/paper10.pdf>.
- [66] Travis Hoppe, H.B., 2024. Nlpre. <https://github.com/NIHOPA/NLPre>. Accessed: 2024-04-15.
- [67] Vargas, S.G.J., 2008. A knowledge-based information extraction prototype for data-rich documents in the information technology domain. Diss. National University of Columbia, Bogota.
- [68] Velez-Estevez, A., Pérez, I.J., García-Sánchez, P., Moral-Munoz, J.A., Cobo, M.J., 2023. New trends in bibliometric apis: A comparative analysis. *Inf. Process. Manag.* 60, 103385. URL: <https://api.semanticscholar.org/CorpusID:258730412>.
- [69] Wang, D., Zhou, X., Zhao, P., Pang, J., Ren, Q., 2025. Early identification of breakthrough technologies: Insights from science-driven innovations. *Journal of Informetrics* 19, 101606. URL: <https://www.sciencedirect.com/science/article/pii/S1751157724001184>, doi:<https://doi.org/10.1016/j.joi.2024.101606>.
- [70] Wei, X., Hoque, M.R.U., Wu, J., Li, J., 2023. Claimdistiller: Scientific claim extraction with supervised contrastive learning, in: Zhang, C., Zhang, Y., Mayr, P., Lu, W., Suominen, A., Chen, H., Ding, Y. (Eds.), *Proceedings of Joint Workshop of the 4th Extraction and Evaluation of Knowledge Entities from Scientific Documents (EEKE2023) and the 3rd AI + Informetrics (AII2023) co-located with the JCDL 2023, Santa Fe, New Mexico, USA and Online, 26 June, 2023*, CEUR-WS.org. pp. 65–77. URL: <https://ceur-ws.org/Vol-3451/paper11.pdf>.
- [71] Würsch, M., Kucharavy, A., David, D.P., Mermoud, A., 2023. Llm-based entity extraction is not for cybersecurity, in: Zhang, C., Zhang, Y., Mayr, P., Lu, W., Suominen, A., Chen, H., Ding, Y. (Eds.), *Proceedings of Joint Workshop of the 4th Extraction and Evaluation of Knowledge Entities from Scientific Documents (EEKE2023) and the 3rd AI + Informetrics (AII2023) co-located with the JCDL*

- 2023, Santa Fe, New Mexico, USA and Online, 26 June, 2023, CEUR-WS.org, pp. 26–32. URL: <https://ceur-ws.org/Vol-3451/paper5.pdf>.
- [72] Würsch, M., Kucharavy, A., David, D.P., Mermoud, A., 2023. Lims perform poorly at concept extraction in cyber-security research literature. *arXiv:2312.07110*.
- [73] You, H., Li, M., Hipel, K.W., Jiang, J., Ge, B., Duan, H., 2017. Development trend forecasting for coherent light generator technology based on patent citation network analysis. *Scientometrics* 111, 297–315. URL: <https://doi.org/10.1007/s11192-017-2252-y>, doi:10.1007/s11192-017-2252-y.
- [74] Zhao, W.X., Zhou, K., Li, J., Tang, T., Wang, X., Hou, Y., Min, Y., Zhang, B., Zhang, J., Dong, Z., Du, Y., Yang, C., Chen, Y., Chen, Z., Jiang, J., Ren, R., Li, Y., Tang, X., Liu, Z., Liu, P., Nie, J.Y., Wen, J.R., 2023. A survey of large language models. *arXiv:2303.18223*.
- [75] Zhu, B., Frick, E., Wu, T., Zhu, H., Ganesan, K., Chiang, W.L., Zhang, J., Jiao, J., 2023. Starling-7b: Improving llm helpfulness & harmlessness with rlaiif.
- [76] Zilio, L., Saadany, H., Sharma, P., Kanojia, D., Orăsan, C., 2022. PLOD: An abbreviation detection dataset for scientific documents, in: Calzolari, N., Béchet, F., Blache, P., Choukri, K., Cieri, C., Declerck, T., Goggi, S., Isahara, H., Maegaard, B., Mariani, J., Mazo, H., Odijk, J., Piperidis, S. (Eds.), *Proceedings of the Thirteenth Language Resources and Evaluation Conference, European Language Resources Association, Marseille, France*. pp. 680–688. URL: <https://aclanthology.org/2022.lrec-1.71>.



## A. Appendix

### A.1. Few-shot prompting

We used few-shot prompting for the triplet extraction. We attempted various different prompts and achieved the best results with the setting that is illustrated below. The **input text** is the text for which we will extract the triples.

**User:** You will extract the subject-predicate-object triples from the text and return them in the form [(subject\_1; predicate\_1; object\_1), (subject\_2; predicate\_2; object\_2), ..., (subject\_n; predicate\_n; object\_n)]. I want the subjects and objects in the triples to be specific, they cannot be pronouns or generic nouns. The first text is: Ever since the Turing Test was proposed in the 1950s, humans have explored the mastering of language intelligence by machine. Language is essentially a complex, intricate system of human expressions governed by grammatical rules.

**Example output:** [(turing test; proposed in; 1950s), (human; explored; language intelligence), (language; is; system of human expressions), (grammatical rule; govern; language)]

**User:** The next text is: Currently, LLMs are mainly built upon the Transformer architecture, where multi-head attention layers are stacked in a very deep neural network. Existing LLMs adopt similar Transformer architectures and pre-training objectives (e.g., lan-guage modeling) as small language models.

**Example output:** [(large language model; built upon; transformer architecture), (multi-head attention layer; stacked in; deep neural network), (large language model; adopt; transformer architecture)]

**User:** The next text is: After collecting a large amount of text data, it is essential to preprocess the data for constructing the pre-training corpus, especially removing noisy, redundant, irrelevant, and potentially toxic data, which may largely affect the capacity and performance of LLMs.

**Example output:** [(preprocessing; remove; noisy data), (preprocessing; remove; redundant data), (preprocessing; remove; irrelevant data), (preprocessing; remove; toxic data), (noisy data; affect; performance large language model), (redundant data; affect; performance large language model), (irrelevant data; affect; performance large language model), (toxic data; affect; performance large language model)]

**User:** The next text is: *input text*

**Assistant:** *Output*