

A Systematic Literature Review of Adapter-based Approaches to Knowledge-enhanced Language Models

Anonymous ACL submission

Abstract

Knowledge-enhanced language models (KELMs) have emerged as promising tools to bridge the gap between large-scale language models and domain-specific knowledge. KELMs can achieve higher factual accuracy and mitigate hallucinations by leveraging knowledge graphs (KGs). They are frequently combined with adapter modules to reduce the computational load and risk of catastrophic forgetting. In this paper, we conduct a systematic literature review (SLR) on adapter-based approaches to KELMs. We provide an overview of approaches in the field and explore the strengths and potential shortcomings of the multitude of discovered methods. We show that both general-knowledge and domain-specific approaches have been frequently explored along with various downstream tasks. Furthermore, we discovered that the biomedical domain is the most popular domain-specific field and that the Pfeiffer adapter is the most commonly used adapter type. We outline the main trends and propose promising future directions.

1 Introduction

The field of natural language processing (NLP) has, in recent years, been dominated by the rise of large language models (LLMs). These models are pre-trained on large amounts of unstructured textual data, which enables them to solve complex reasoning tasks and generate new text. Still, LLMs can lack awareness of structured knowledge hierarchies, such as relations between concepts. This drawback can lead to inaccurate predictions for downstream tasks relying on structured predictions and so-called "hallucinations" within text generation. This can make LLMs less reliable in practice, which is an especially precarious issue in high-risk domains like healthcare or law.

A potential solution to counteract mispredictions and hallucinations and improve the reliability of

LLMs is knowledge enhancement: By leveraging expert knowledge from manually curated knowledge graphs (KGs), structured knowledge can be injected into LLMs. Such knowledge-enhanced language models (KELMs) are a promising approach for higher structured knowledge awareness, better factual accuracy, and less hallucinations (Colon-Hernandez et al., 2021; Wei et al., 2021).

Unfortunately, knowledge enhancement in the form of supervised fine-tuning (SFT) of the whole LLM can be highly computationally expensive, especially for models with billions of parameters. A promising research avenue to overcome this limitation is using lightweight and efficient adapter modules to inject structured knowledge into LLMs. Using adapters for knowledge enhancement helps enhance the task performance of LLMs and is, at the same time, a very computationally efficient solution. Despite the rising popularity of this approach, to the best of our knowledge, a comprehensive overview of adapter-based KELMs is still missing in the NLP research landscape.

To bridge this research gap, we conduct a systematic literature review (SLR) on adapter-based knowledge enhancement of LLMs. Our contributions are: (1) a novel review on adapter-based knowledge enhancement, (2) a quantitative and qualitative analysis of different methods in the field, and (3) detailed categorization of literature and identification of most promising trends.

2 Background and Related Work

In this section, we give an overview of related work and existing surveys on knowledge enhancement. Knowledge graphs are the most common external knowledge source, so we start with their overview.

2.1 Knowledge Graphs

Knowledge graphs (KGs) are a structured representation of the world knowledge and have seen a rising prominence in NLP research over the

past decade (Schneider et al., 2022). Hogan et al. (2020) define a KG as "a graph of data intended to accumulate and convey knowledge of the real world, whose nodes represent entities of interest and whose edges represent relations between these entities". Similarly, Ji et al. (2020) published a comprehensive survey on KGs and, following existing literature, defined the concept of a KG as " $\mathcal{G} = \{\mathcal{E}, \mathcal{R}, \mathcal{F}\}$, where \mathcal{E}, \mathcal{R} and \mathcal{F} are sets of entities, relations and facts, respectively; a fact is denoted as a triple $(h, r, t) \in \mathcal{F}$ ". Depending on the source and purpose of a KG, entities and relations can take on various shapes. For example, in the biomedical knowledge graph UMLS (Bodenreider, 2004), a relation can take the shape of a single word like "inhibits", a short phrase like "relates to", or a compound term including, for example, chemical or medical categories such as "[protein] relates to [disease]" or "[substance] induces [physiology]". A textual connection is vital because it serves as a link between the graph structure and natural language, simplifying the integration of information from KGs into language models and the associated learning processes. Other than UMLS, other examples of popular KGs are DBpedia (Auer et al., 2007) and ConceptNet (Speer et al., 2017).

2.2 Approaches to Knowledge Enhancement

At the time of writing, some reviews had already been published that gave an overview of KELMs and classified different approaches. Colon-Hernandez et al. (2021) review the existing literature and split the approaches to integrate structure knowledge with LMs into three categories: (1) input-centered strategies, centering around altering the structure of the input or selected data, which is fed into the base LLM; (2) architecture-focused approaches, which involve either adding additional layers that integrate knowledge with the contextual representations or modifying existing layers to alter parts like attention mechanisms; (3) output-focused approaches, which work by changing either the output structure or the losses used in the base model. Our study focuses on the second category (2), by examining the adapter-based mechanisms for injecting information into the model, which were shown to be the most promising by the authors.

The second survey by Wei et al. (2021) reviews a large number of studies on KELMs and classifies them using three taxonomies: (1) knowledge sources, (2) knowledge granularity, and (3)

application areas. Within (1), the knowledge sources include linguistic knowledge, encyclopedic knowledge, and commonsense and domain-specific knowledge. The second taxonomy (2) acknowledges the common approach of using KGs as a source of knowledge. Levels of granularity mentioned are text-based knowledge, entity knowledge, relation triples, and KG sub-graphs. Lastly, with the third taxonomy (3), the authors discuss how knowledge enhancement can improve natural language generation and understanding. They also review popular benchmarks that can be used for task evaluation of KELMs (Wei et al., 2021).

These two field studies by Colon-Hernandez et al. (2021) and Wei et al. (2021) on the classification of KELM approaches were our starting point for exploring KELMs and initially proved to be very valuable. However, although they address some adapter-based studies like K-Adapter (Wang et al., 2020), most other adapter-based KELMs are missing. This lack of coverage led to our decision to conduct a novel systematic literature search focusing specifically on the adapter-based KELMs, considering their rising popularity and importance.

3 Adapters

In the following, an overview of adapters for LLMs and their individual functionalities and applications will be given to establish a conceptual understanding of adapter-based approaches to LLMs.

3.1 Overview

Broadly speaking, adapters are small bottleneck feed-forward layers inserted within each layer of an LLM (Houlsby et al., 2019). The small amount of additional parameters allows injecting new data or knowledge without fine-tuning the whole model. This feat is usually accomplished by freezing the layers of the base model with its millions or billions of parameters while only updating the adapter weights (e.g., through entity prediction tuning). Due to the lightweight nature of adapters, this approach leads to short training times with relatively low computing resource requirements. Adapters used to be utilized mostly for quick and cheap downstream-task fine-tuning but are now increasingly used for knowledge enhancement. Because it is possible to train adapters individually, they can also be used for multi-task training by specializing one adapter for each task or multi-domain knowledge injection by specializing adapters to different

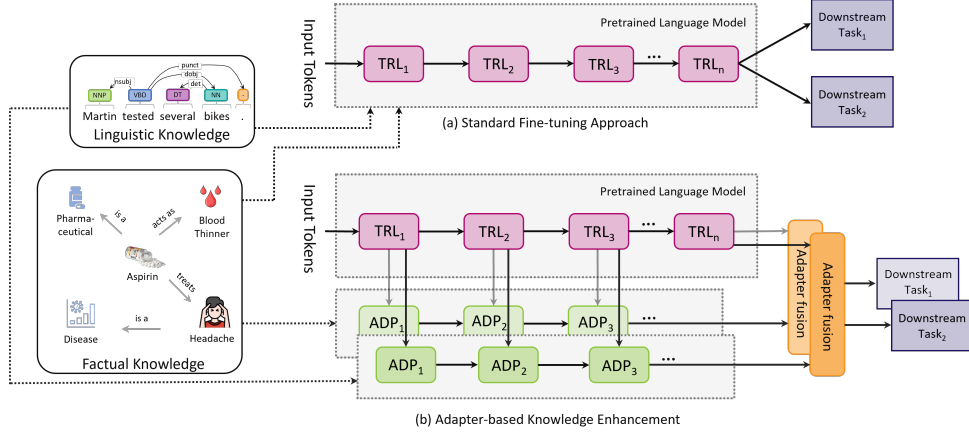


Figure 1: Illustration of a standard fine-tuning versus a knowledge enhancement process. In the example, knowledge from a KG is injected into the model via adapters.

domains (Pfeiffer et al., 2020a).

Leveraging adapters in LLMs also has positive "side effects": Adapters can avoid catastrophic forgetting (the issue when an LLM suddenly deteriorates in performance after fine-tuning) by introducing new task-specific parameters (Houlsby et al., 2019; Pfeiffer et al., 2020a) and, in transfer learning, adapters have even been shown to improve stability and adversarial robustness for various downstream tasks (Han et al., 2021). The specifics of how and where adapters are added to an LLM depend on the adapter type.

3.2 Adapter Types

Houlsby Adapter. The Houlsby Adapter (Houlsby et al., 2019) was the first adapter to be used for transfer learning in NLP. The idea was based on adapter modules initially introduced by Rebuffi et al. (2017) in the computer vision domain. The two main principles stayed the same: Adapters require a relatively small number of parameters compared to the base model and a near-identity initialization. These principles ensure that the total model size grows relatively slowly when more transfer tasks are added, while a near-identity initialization is required for stable training of the adapted model (Houlsby et al., 2019). The optimal architecture of the Houlsby Adapter was determined by meticulous experimenting and tuning; the result can be seen in figure 2. In a classical transformer structure (Vaswani et al., 2017), the adapter module is added once after the multi-headed attention and once after the two feed-forward layers. The modules project the d -dimensional layer features of the base model into

a smaller dimension, m , then apply a non-linearity (like ReLU) and project back to d dimensions. The configuration also hosts a skip-connection, and the output of each sub-layer is forwarded to a layer normalization (Ba et al., 2016). Including biases, $2md + d + m$ parameters are added per layer, accounting for only 0.5 to 8 percent of the parameters of the original BERT model used by the authors when setting $m \ll d$.

Bapna and Firat Adapter. In contrast to the Houlsby Adapter, Bapna and Firat (2019) only introduce one adapter module in each transformer layer: they keep the adapters after the multi-headed attention (so-called "top" adapters) while dropping the adapters after the feed-forward layers (so-called "bottom" adapters) of the transformer (refer to Figure 2 for better understanding of the component positions). Moreover, while Houlsby et al. (2019) re-train layer normalization parameters for every domain, Bapna and Firat (2019) "simplify this formulation by leaving the parameters frozen, and introducing new layer normalization parameters for every task, essentially mimicking the structure of the transformer feed-forward layer".

Pfeiffer Adapter and AdapterFusion. The approaches of Bapna and Firat (2019); Houlsby et al. (2019) did not allow information sharing between tasks. Pfeiffer et al. (2020a) introduce Adapter Fusion, a two-stage algorithm that addresses the sharing of information encapsulated in adapters trained on different tasks. In the first stage, they train the adapters in single-task or multi-task setups for a total of N tasks similar to the Houlsby Adapter, but only keeping the top adapters, sim-

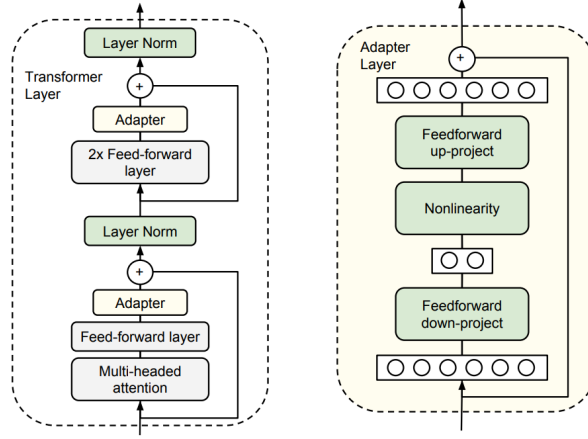


Figure 2: Location of the adapter module in a transformer layer (left) and architecture of the Housby Adapter (right). All green layers are trained on fine-tuning data, including the adapter itself, the layer normalization parameters, and the final classification layer (not shown). Image with permission from Housby et al. (2019).

ilar to the Bapna and Firat Adapter. As a second step, they combine the set of N adapters with AdapterFusion: They fix the parameters Θ and all adapters Φ , and finally introduce parameters Ψ that learn to combine the N task adapters for the given target task (Pfeiffer et al., 2020a): $\Psi_m \leftarrow \underset{\Psi}{\operatorname{argmin}} L_m(D_m; \Theta, \Phi_1, \dots, \Phi_N, \Psi)$

Here, Ψ_m are the learned AdapterFusion parameters for task m . In the process, the training dataset of m is used twice: once for training the adapters Φ_m and again for training Fusion parameters Ψ_m , which learn to compose the information stored in the N task adapters (Pfeiffer et al., 2020a). With their approach of separating knowledge extraction and knowledge composition, they further improve the ability of adapters to avoid catastrophic forgetting and interference between tasks and training instabilities. The authors also find that their approach of using only a single adapter after the feed-forward layer performs on par with the Housby adapter while requiring only half of the newly introduced adapters (Pfeiffer et al., 2020a). This makes the Pfeiffer adapter an attractive choice for many applications, further proven by its popularity among the papers in our review.

K-Adapter Wang et al. (2020) follow a substantially different approach where the adapters work as "outside plug-ins". In their work, an adapter model consists of K adapter layers (hence the name) that contain N transformer layers and two projection layers. Similar to the approaches above, a skip connection is added but instead applied across the two projection layers. The adapter layers are plugged in

among varying transformer layers of the pre-trained model. The authors explain that they concatenate the output hidden feature of the transformer layer in the pre-trained model and the output feature of the former adapter layer as the input feature of the current adapter layer.

Adapter architectures for knowledge enhancement exist that differ from the four adapter types mentioned here. For example, the "Parallel Adapter" (He et al., 2021a) or the adapter architecture by Stickland and Murray (2019)). However, as the upcoming comprehensive literature survey will show, these architectures are either unique to specific papers or have not found broader applications in the field of KELMs. Either way, these approaches are out of the scope of this paper and will not be discussed here.

Another popular type of efficient adaptation is the low-rank adaptation LoRA (Hu et al., 2022), and its quantized version QLoRA (Dettmers et al., 2023). Despite the name, these approaches do not actually add new adapter layers as the previously described ones but enforce a low-rank constraint on the weight updates of the base model's layers. This enables efficient fine-tuning of LLMs but does not properly allow for knowledge enhancement from external sources, which is the focus of our review.

4 Methodology

This chapter details the methodology we employed for the systematic literature review. We largely followed the procedure of Kitchenham et al. (2009) for systematic literature reviews in software engineering. The search strategy for the systematic

literature review of this thesis included literature that fulfilled the following inclusion criteria:

- Peer-reviewed articles from ACM¹, ACL², and IEEE Xplore³
- Article abstracts that match the search string ("*adapter*" OR "*adapter-based*") AND ("*language model*" OR "*nlp*" OR "*natural language processing*") AND ("*injection*" OR "*knowledge*")
- Articles published after February 2, 2019 (publication of the Hounsby Adapter, the first LLM adapter)
- Articles that address the topic of adapter-based knowledge-enhanced language models

We also included a limited number of articles not found on the mentioned databases because they were fundamental works on the topic of the SLR and frequently referenced. The SLR was concluded in January 2024 and represents the state of research literature up to this point.

5 Results

This section will present the results of the systematic literature review on adapter-based knowledge enhancement.

5.1 Overview

| Source | Initial | Abstract | Full Text |
|--------------|-----------|-----------|-----------|
| IEEE | 28 | 6 | 6 |
| ACM | 10 | 6 | 5 |
| ACL | 36 | 16 | 13 |
| Others | 2 | 2 | 2 |
| Total | 76 | 30 | 26 |

Table 1: Quantitative overview of the literature sources and the selection process

Table 1 shows the source distribution for all included papers. Fifty-nine papers were found by applying the search string as a command on the ACL, ACM, and IEEE search engines. Due to their importance for the field, we included three additional papers from other sources. These papers were found through online search and paper references during the general research process. In summary, after the abstract screening, 31 articles

¹<https://dl.acm.org/>
²<https://aclanthology.org/>
³<https://ieeexplore.ieee.org/Xplore/home.jsp>

met all inclusion criteria (and no exclusion criteria). After the full paper screening, 26 papers remained to form the final paper pool of the survey.

Table 2 gives an overview of all papers included in the survey. It includes the information on the adapter type used in the paper, the domain and scope of the paper, and for which downstream NLP tasks it was developed.

5.2 Data Analysis

This section starts with a quantitative analysis showcasing and interpreting quantitative distributions. Afterward, we report significant qualitative insights from the papers.

5.2.1 Quantitative Analysis

Yearly Distribution To begin with, we assess how many papers were published each year to get a sense of the trend and growth in the area (Fig. 3). There has been a noticeable increase in publications on adapter-based approaches to knowledge-enhanced language models in recent years, especially from 2022 onward. This trend suggests growing interest and research activity in the domain.

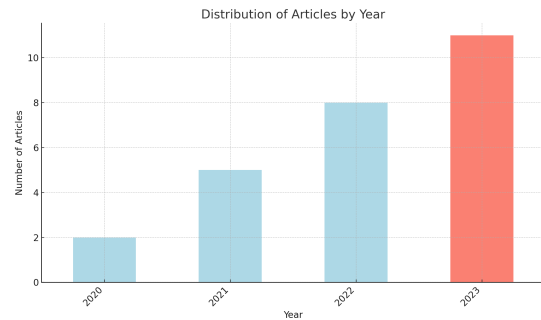


Figure 3: Yearly distribution of publications

Adapter Type Distribution. Next, we evaluate the popularity and variety of adapter types used across the papers (Fig. 4). The “Pfeiffer” and “Hounsby” adapter types stand out as the most common, which suggests that the closely related underlying architecture is the most popular methodology in the field. This popularity is likely not only an achievement of the adapter’s performance but also due to the well-established Adapter-Hub platform (Pfeiffer et al., 2020b), which, although offering other options, uses adapters with the Pfeiffer configuration by default. This finding showcases a need and trend to build custom adapters well-suited to individual tasks. In the upcoming years, we will likely see many novel adapter architectures. The

| paper & nickname | adapter type | scope | task |
|---|-------------------|-----------------------|---------------|
| K-MBAN (Zou et al., 2022) | K-Adapter | open | RC |
| / (Moon et al., 2021) | Houlsby | open | MT |
| CSBERT (Yu and Yang, 2023) | Unique | open | SL |
| / (Qian et al., 2022) | Unique | open | SR |
| / (Li et al., 2023) | Houlsby | closed (multi-domain) | SF |
| CPK (Liu et al., 2023) | K-Adapter | closed (biomedical) | RC, ET, QA |
| CKGA (Lu et al., 2023) | Unique | open | SC |
| / (Nguyen-The et al., 2023) | Pfeiffer | open | SA |
| KEBLM (Lai et al., 2023) | Pfeiffer | closed (biomedical) | QA, NLI, EL |
| / (Guo and Guo, 2022) | Unique | open | NER |
| / (Tiwari et al., 2023) | Unique | closed (biomedical) | TS |
| AdapterSoup (Chronopoulou et al., 2023) | Bapna and Firat | closed (multi-domain) | LM |
| / (Wold, 2022) | Houlsby | open | LAMA |
| / (Chronopoulou et al., 2022) | Unique | closed (multi-domain) | LM |
| DS-TOD (Hung et al., 2022) | Pfeiffer | closed (multi-domain) | TOD |
| / (Emelin et al., 2022) | Houlsby | closed (multi-domain) | TOD |
| KnowBERT (Xu et al., 2022) | Bapna and Firat | open | KGD |
| mDAPT (Kær Jørgensen et al., 2021) | Pfeiffer | closed (multi-domain) | NER, STC |
| DAKI (Lu et al., 2021) | K-Adapter | closed (biomedical) | NLI |
| / (Majewska et al., 2021) | Pfeiffer | open | EE |
| / (Lauscher et al., 2020) | Houlsby | open | GLUE |
| TADA (Hung et al., 2023) | Unique | open | TOD, NER, NLI |
| LeakDistill (Vasylenko et al., 2023) | StructAdapt | open | SMATCH |
| MixDA (Diao et al., 2023) | Houlsby, Pfeiffer | closed (multi-domain) | GLUE, TXM |
| MoP (Meng et al., 2021) | Pfeiffer | closed (biomedical) | BLURB |
| K-Adapter (Wang et al., 2020) | K-Adapter | open | RCL, ET, QA |

Table 2: Overview of the results for the literature survey, including all papers and their references. The task acronyms are explained in the glossary at the end of the thesis. The dotted lines separate the database sources: First come the IEEE papers, then ACM, ACL, and finally, the papers from other sources. For the definition of all task acronyms, see Appendix A.4

“K-Adapter” and “Bapna and Firat” adapters are the less frequently mentioned architectures, suggesting that these approaches are less well-established. Overall, various adapter types are present, indicating a diverse range of methodologies being explored.

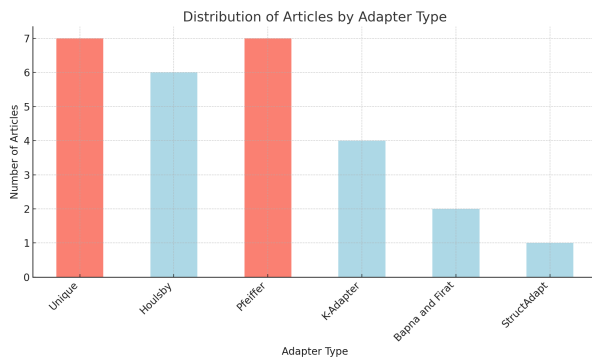


Figure 4: Distribution of adapter types being used in the articles

Domain Analysis Third, we analyze the distribution of papers across the domain scope and coverage to understand domain-specific preferences in the literature (figures given in the appendix). The first plot in Figure 5 shows that the open-domain

scope is the most popular, with many papers exploring adapter-based approaches within the open domain. The popularity is likely caused by the interest in creating LLMs with a common-sense understanding or world knowledge.

As illustrated by the second plot in Figure 5, the single- and multi-domain approaches are split evenly within the closed-domain papers.

Finally, the third plot addresses the coverage of the biomedical domain. In absolute numbers, only six papers focus on the biomedical domain, but relative to other parts, the biomedical field is by far the most prominent of all domain-specific approaches. The popularity likely comes down to the availability of large biomedical KGs, and medicine historically being one of the most active research fields in general science (Cimini et al., 2014).

Task Distribution A highly diverse range of tasks is being explored throughout the papers, which signifies the versatility and potential of adapter-based approaches across different natural language processing tasks. However, combined with the limited number of papers in the survey, the approach-versatility prevents further meaningful

quantitative analysis. Still, tasks such as Reading Comprehension (RC), Named Entity Recognition (NER), and Question Answering (QA) appear to be popular areas of focus in the literature. This could be because these tasks are the most demanding regarding structural knowledge requirements. In the appendix, Figure 6 provides a word cloud of all keywords in the downstream tasks as a visualization, showing that there is also a focus on tasks with a dialogue or sentiment component.

5.2.2 Qualitative Analysis

This section of the analysis highlights recurring themes and individual insights from the papers. Fully summarizing all articles was outside the scope of this survey. However, we still provide an overview of the most common patterns.

General Knowledge. The quantitative analysis showed that open-domain approaches are more popular than their close-domain counterparts. Subsequently, there is also a large variety in the used frameworks, knowledge sources, and overall goals of the papers. Two commonly used KGs for general knowledge are ConceptNet (Speer et al., 2017) for common-sense knowledge, and DBpedia (Auer et al., 2007) for encyclopedic world knowledge. Two example works that use these KGs are Wold (2022) and the CKGA ("knowledge graph-based adapter") by Lu et al. (2023). Wold (2022) train adapter modules on sub-graphs of ConceptNet to inject factual knowledge into LLMs. They evaluate their framework on the Concept-Net Split of the LAMA Probe (Petroni et al., 2019) and see increasing performance while only adding 2.1% of new parameters to the original models. CKGA (Lu et al., 2023), on the other hand, tackle aspect-level sentiment classification by leveraging knowledge from DBpedia. They link aspects to DBpedia and extract an aspect-related sub-graph. Then, a pre-trained language model and the knowledge graph embedding are utilized to encode the common-sense knowledge of entities, where the corresponding knowledge is extracted with graph convolutional networks (Lu et al., 2023).

Linguistic Knowledge Instead of only including factual knowledge, some works also inject linguistic knowledge into adapters (Majewska et al., 2021; Zou et al., 2022; Yu and Yang, 2023; Wang et al., 2020). While LLMs already encode a range of syntactic and semantic properties of language, Majewska et al. (2021) explain that they "are still

prone to fall back on superficial cues and simple heuristics to solve downstream tasks, rather than leverage deeper linguistic information". Their paper explores the interplay between verb meaning and argument structure. They use the gained knowledge to enhance LLMs with Pfeiffer Adapters to improve English event extraction and machine translation in other languages. Another example is the work of Zou et al. (2022) on machine reading comprehension (MRC). They proposed the K-MBAN model to integrate linguistic and factual external knowledge into LLMs through K-Adapters.

Domain-specific Knowledge Chronopoulou et al. (2022) propose a parameter-efficient approach to domain adaptation using adapters. They represent domains as a hierarchical tree structure where each node in the tree is associated with a set of adapter weights. Their work focused on specializing adapters in website domains like *booking.com* and *yelp.com*. In another instance, Chronopoulou et al. (2023) propose "AdapterSoup". In this framework, they also use adapters for domain-specific tasks but use "an approach that performs weight-space averaging of adapters trained on different domains". AdapterSoup can be helpful in various domain-specific approaches in low-resource settings, especially when only a small amount of data on a specific subdomain is obtainable and closely related adapters are available instead. Earlier, we saw that the biomedical domain is the most prevalent among the closed-domain approaches to adapter-based KELMs. We will briefly examine the relevant works in the following.

Biomedical Knowledge We have found the works of DAKI (Lu et al., 2021), MoP (Meng et al., 2021), and KEBLM (Lai et al., 2023) to be the most impactful. According to the results of our literature survey, DAKI ("Diverse Adapters for Knowledge Integration") was the first work to use adapters specifically for knowledge enhancement in the biomedical domain. Lu et al. (2021) leverage data from the UMLS meta-thesaurus and UMLS Semantic Network groups concepts, but also from Wikipedia articles for diseases as proposed by He et al. (2020). Meng et al. (2021) recognize that KGs like UMLS, which can be several gigabytes large, are very expensive to train on in their entirety. They propose to use a "Mixture of Partitions" (MoP), which splits the KG into sub-graphs

and combines later with AdapterFusion (Pfeiffer et al., 2020a). Finally, the KEBLM framework’s trademark is that it allows the inclusion of a variety of knowledge types from multiple sources into biomedical LLMs. In contrast to DAKI, which also utilizes more than one source, KEBLM includes a knowledge consolidation phase after the knowledge injection, where they teach the fusion layers to effectively combine knowledge from both the original PLM and newly acquired external knowledge by using a large collection of unannotated texts (Lai et al., 2023). For completeness, we refer to Kær Jørgensen et al. (2021) for information on the m-DAPT framework, which addresses multilingual domain adaptation for biomedical LLMs and KeBioSum (Xie et al., 2022), who state their work is the first study exploring knowledge injection for biomedical extractive summarization.

Performance Insights He et al. (2021b) criticize that "existing work only focuses on the parameter-efficient aspect of adapter-based tuning while lacking further investigation on its effectiveness". They address this issue with their work and show that adapter-based tuning better mitigates forgetting issues than regular fine-tuning since it yields representations with less deviation from those generated by the initial pre-trained language model. They found that adapter-based approaches outperform fine-tuning in low-resource and cross-lingual settings and are "more robust to overfitting and less sensitive to changes in learning rates" (He et al., 2021b). This is further proven by all the papers from our survey that compare the performance on classification benchmarks between adapter-based knowledge-enhanced models and vanilla-base models, always showing improvements over the vanilla version of the models. Notable examples are the work of Meng et al. (2021) or Lai et al. (2023), which evaluate biomedical language understanding tasks and reach up to +8% increase in accuracy with adapter-based enhancement.

6 Current and Future Trends

In this section, we outline the most important findings and trends of the review and point out the promising future directions:

- Adapter-based KELMs are a recent development in NLP, but there has been fast-growing interest in them recently, with a linear yearly increase of published papers. We predict the

growing trend to continue.

- Various adapter architectures exist and have been iteratively advanced yearly to be more efficient while preserving task performance. This peaked with the Pfeiffer adapter, which is the most popular type. We expect future work to focus their updates on adapter architecture by overcoming the latency of sequential data processing in adapters and enabling hardware parallelism.
- Research focuses on the open domain – injecting general world knowledge into models. Within the closed domain, the biomedical domain is the most popular, owing to the existence of large biomedical KGs. We foresee the potential to apply adapter-based KELMs to other highly structured domains, such as the legal or financial domain (documents with rigid structure).
- A wide array of downstream tasks is being explored. The biggest improvement in task performance is seen in knowledge-intensive tasks like question answering and text classification, with a smaller improvement for reasoning tasks like entailment recognition. Generative tasks, other than dialogue modeling, are rather unexplored. We envision a future popular use case that could use knowledge enhancement to improve the factuality and informativeness of generated text.

7 Conclusion

In this paper, we conducted a systematic literature review on approaches to enhancing language models with external knowledge using adapter modules. We portrayed which adapter-based approaches exist and how they compare to each other. We showed there is a steady growth of interest in this domain with each new year and highlighted the most popular adapter architectures (with "Pfeiffer" as the predominant one). We discovered there is a balance in popularity between open-domain approaches, focusing on integrating general world knowledge into models, and closed-domain focusing on specialized fields, with biomedical as the most popular domain. With our review, we contribute a novel and extensive resource for this nascent yet fast-growing field and we hope it will be a useful entry point for other researchers in the future.

Limitations

The methodology of a systematic literature review follows a strict search string and exclusion criteria. Therefore, it is possible that we excluded some relevant work on adapter-based KELMs. Moreover, while we tried to report on our survey as comprehensively as possible, there are several aspects we could not include in this work. Also, some of the reviewed articles were not given an adequate qualitative analysis in this work due to space constraints, leading to potentially missing insights and a non-complete representation of the state of research on adapter-based knowledge enhancement. Additionally, due to the variety of applications and domains, we were not able to give precise guidelines on what methods to use under which circumstances. Still, we aimed to report on the most common patterns and trends discovered in the literature, which can serve as a basis for future research.

References

- Sören Auer, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, and Zachary Ives. 2007. Dbpedia: A nucleus for a web of open data. In *The Semantic Web*, pages 722–735, Berlin, Heidelberg. Springer Berlin Heidelberg.
- Jimmy Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. 2016. [Layer normalization](#). *ArXiv*, abs/1607.06450.
- Ankur Bapna and Orhan Firat. 2019. [Simple, scalable adaptation for neural machine translation](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1538–1548, Hong Kong, China. Association for Computational Linguistics.
- Olivier Bodenreider. 2004. [The unified medical language system \(umls\): integrating biomedical terminology](#). *Nucleic acids research*, 32.
- Shu Cai and Kevin Knight. 2013. [Smatch: an evaluation metric for semantic feature structures](#). In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 748–752, Sofia, Bulgaria. Association for Computational Linguistics.
- Alexandra Chronopoulou, Matthew Peters, and Jesse Dodge. 2022. [Efficient hierarchical domain adaptation for pretrained language models](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1336–1351, Seattle, United States. Association for Computational Linguistics.
- Alexandra Chronopoulou, Matthew Peters, Alexander Fraser, and Jesse Dodge. 2023. [AdapterSoup: Weight averaging to improve generalization of pre-trained language models](#). In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 2054–2063, Dubrovnik, Croatia. Association for Computational Linguistics.
- Giulio Cimini, Andrea Gabrielli, and Francesco Labini. 2014. [The scientific competitiveness of nations](#). *PloS one*, 9.
- Pedro Colon-Hernandez, Catherine Havasi, Jason B. Alonso, Matthew Huggins, and Cynthia Breazeal. 2021. [Combining pre-trained language models and structured knowledge](#). *ArXiv*, abs/2101.12294.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. [Qlora: Efficient finetuning of quantized llms](#). *arXiv preprint arXiv:2305.14314*.
- Shizhe Diao, Tianyang Xu, Ruijia Xu, Jiawei Wang, and Tong Zhang. 2023. [Mixture-of-domain-adapters: Decoupling and injecting domain knowledge to pre-trained language models’ memories](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5113–5129, Toronto, Canada. Association for Computational Linguistics.
- Denis Emelin, Daniele Bonadiman, Sawsan Alqahtani, Yi Zhang, and Saab Mansour. 2022. [Injecting domain knowledge in language models for task-oriented dialogue systems](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 11962–11974. Association for Computational Linguistics.
- Yu Gu, Robert Tinn, Hao Cheng, Michael R. Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2020. [Domain-specific language model pretraining for biomedical natural language processing](#). *ACM Transactions on Computing for Healthcare (HEALTH)*, 3:1 – 23.
- Qian Guo and Yi Guo. 2022. [Lexicon enhanced chinese named entity recognition with pointer network](#). *Neural Computing and Applications*.
- Wenjuan Han, Bo Pang, and Ying Nian Wu. 2021. [Robust transfer learning with pretrained language models through adapters](#). *ArXiv*, abs/2108.02340.
- Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. 2021a. [Towards a unified view of parameter-efficient transfer learning](#). *ArXiv*, abs/2110.04366.
- Ruidan He, Linlin Liu, Hai Ye, Qingyu Tan, Bosheng Ding, Liying Cheng, Jia-Wei Low, Lidong Bing, and Luo Si. 2021b. [On the effectiveness of adapter-based tuning for pretrained language model adaptation](#).
- Yun He, Ziwei Zhu, Yin Zhang, Qin Chen, and James Caverlee. 2020. [Infusing Disease Knowledge into](#)

| | | |
|-----|--|---|
| 724 | BERT for Health Question Answering, Medical In- | 782 |
| 725 | ference and Disease Name Recognition. In <i>Proceed-</i> | 783 |
| 726 | ings of the 2020 Conference on Empirical Methods | 784 |
| 727 | in Natural Language Processing (EMNLP), pages | 785 |
| 728 | 4604–4614, Online. Association for Computational | 786 |
| 729 | Linguistics. | 787 |
| 730 | Aidan Hogan, Eva Blomqvist, Michael Cochez, Clau- | 788 |
| 731 | dia d’Amato, Gerard de Melo, Claudio Gutiérrez, | 789 |
| 732 | S. Kirrane, José Emilio Labra Gayo, Roberto Navigli, | 790 |
| 733 | Sebastian Neumaier, Axel-Cyrille Ngonga Ngomo, | |
| 734 | Axel Polleres, Sabbir M. Rashid, Anisa Rula, Lukas | Bo Li, Dongseong Hwang, Zhouyuan Huo, Junwen |
| 735 | Schmelzeisen, Juan Sequeda, Steffen Staab, and An- | Bai, Guru Prakash, Tara N. Sainath, Khe Chai Sim, |
| 736 | toine Zimmermann. 2020. Knowledge graphs . <i>ACM</i> | Yu Zhang, Wei Han, Trevor Strohman, and Fran- |
| 737 | <i>Computing Surveys (CSUR)</i> , 54:1 – 37. | coise Beaufays. 2023. Efficient domain adaptation |
| 738 | Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, | for speech foundation models . In <i>ICASSP 2023 -</i> |
| 739 | Bruna Morrone, Quentin de Laroussilhe, Andrea Ges- | <i>2023 IEEE International Conference on Acoustics,</i> |
| 740 | munido, Mona Attariyan, and Sylvain Gelly. 2019. | <i>Speech and Signal Processing (ICASSP)</i> , pages 1–5. |
| 741 | Parameter-efficient transfer learning for nlp . In <i>Inter-</i> | 797 |
| 742 | <i>national Conference on Machine Learning</i> . | |
| 743 | Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan | C. Liu, S. Zhang, C. Li, and H. Zhao. 2023. Cpk- |
| 744 | Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and | adapter: Infusing medical knowledge into k-adapter |
| 745 | Weizhu Chen. 2022. LoRA: Low-rank adaptation of | with continuous prompt . In <i>2023 8th International</i> |
| 746 | large language models . In <i>International Conference</i> | <i>Conference on Intelligent Computing and Signal Pro-</i> |
| 747 | <i>on Learning Representations</i> . | <i>cessing (ICSP)</i> , pages 1017–1023, Los Alamitos, CA, |
| 748 | | USA. IEEE Computer Society. |
| 749 | Chia-Chien Hung, Lukas Lange, and Jannik Strötgen. | 803 |
| 750 | 2023. TADA: Efficient task-agnostic domain adapta- | |
| 751 | tion for transformers . In <i>Findings of the Association</i> | Guojun Lu, Haibo Yu, Zehao Yan, and Yun Xue. 2023. |
| 752 | <i>for Computational Linguistics: ACL 2023</i> , pages 487– | Commonsense knowledge graph-based adapter for |
| 753 | 503, Toronto, Canada. Association for Computational | aspect-level sentiment classification . <i>Neurocomput-</i> |
| 754 | Chia-Chien Hung, Anne Lauscher, Simone Ponzetto, | <i>ing</i> , 534:67–76. |
| 755 | and Goran Glavaš. 2022. DS-TOD: Efficient domain | 807 |
| 756 | specialization for task-oriented dialog . In <i>Findings of</i> | |
| 757 | <i>the Association for Computational Linguistics: ACL</i> | Qihao Lu, Dejing Dou, and Thien Huu Nguyen. 2021. |
| 758 | 2022, pages 891–904. Association for Computational | Parameter-efficient domain knowledge integration |
| 759 | Linguistics. | from multiple sources for biomedical pre-trained lan- |
| 760 | Shaoxiong Ji, Shirui Pan, E. Cambria, Pekka Marttinen, | guage models . In <i>Findings of the Association for</i> |
| 761 | and Philip S. Yu. 2020. A survey on knowledge | <i>Computational Linguistics: EMNLP 2021</i> , pages |
| 762 | graphs: Representation, acquisition, and applications . | 3855–3865. Association for Computational Linguis- |
| 763 | <i>IEEE Transactions on Neural Networks and Learning</i> | <i>tics</i> . |
| 764 | <i>Systems</i> , 33:494–514. | 814 |
| 765 | Rasmus Kær Jørgensen, Mareike Hartmann, Xiang Dai, | |
| 766 | and Desmond Elliott. 2021. mDAPT: Multilingual | Olga Majewska, Ivan Vulić, Goran Glavaš, |
| 767 | domain adaptive pretraining in a single model . In | Edoardo Maria Ponti, and Anna Korhonen. |
| 768 | <i>Findings of the Association for Computational Lin-</i> | 2021. Verb knowledge injection for multilingual |
| 769 | <i>guistics: EMNLP 2021</i> , pages 3404–3418. Associa- | event processing . In <i>Proceedings of the 59th Annual</i> |
| 770 | tion for Computational Linguistics. | <i>Meeting of the Association for Computational</i> |
| 771 | Barbara Kitchenham, O. Pearl Brereton, David Budgen, | <i>Linguistics and the 11th International Joint Con-</i> |
| 772 | Mark Turner, John Bailey, and Stephen Linkman. | <i>ference on Natural Language Processing (Volume</i> |
| 773 | 2009. Systematic literature reviews in software en- | <i>1: Long Papers)</i> , pages 6952–6969. Association for |
| 774 | gineering – a systematic literature review . <i>Informa-</i> | <i>Computational Linguistics</i> . |
| 775 | <i>tion and Software Technology</i> , 51(1):7–15. Special | 823 |
| 776 | Section - Most Cited Articles in 2002 and Regular | |
| 777 | Research Papers. | Zaiqiao Meng, Fangyu Liu, Thomas Hiku Clark, |
| 778 | Tuan Manh Lai, ChengXiang Zhai, and Heng Ji. | Ehsan Shareghi, and Nigel Collier. 2021. Mixture- |
| 779 | 2023. Keblm: Knowledge-enhanced biomedical lan- | of-partitions: Infusing large biomedical knowledge |
| 780 | guage models . <i>Journal of Biomedical Informatics</i> , | graphs into bert . <i>ArXiv</i> , abs/2109.04810. |
| 781 | 143:104392. | 827 |
| | | |
| | | Hyeonseok Moon, Chanjun Park, Sugyeong Eo, Jae- |
| | | hyung Seo, and Heuiseok Lim. 2021. An empirical |
| | | study on automatic post editing for neural machine |
| | | translation . <i>IEEE Access</i> , 9:123754–123763. |
| | | 831 |
| | | |
| | | Maude Nguyen-The, Soufiane Lamghari, Guillaume- |
| | | Alexandre Bilodeau, and Jan Rockemann. 2023. |
| | | Leveraging sentiment analysis knowledge to solve |
| | | emotion detection tasks . In <i>Pattern Recognition,</i> |
| | | <i>Computer Vision, and Image Processing. ICPR 2022</i> |
| | | <i>International Workshops and Challenges</i> , pages 405– |
| | | 416. Springer Nature Switzerland. |
| | | 838 |

| | | | |
|-----|--|--|-----|
| 839 | Fabio Petroni, Tim Rocktäschel, Patrick Lewis, An- | Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob | 895 |
| 840 | ton Bakhtin, Yuxiang Wu, Alexander H. Miller, and | Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz | 896 |
| 841 | Sebastian Riedel. 2019. Language models as knowl- | Kaiser, and Illia Polosukhin. 2017. Attention is all | 897 |
| 842 | edge bases? <i>ArXiv</i> , abs/1909.01066. | you need . In <i>NIPS</i> . | 898 |
| 843 | Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, | Pavlo Vasylenko, Pere Lluís Huguet Cabot, | 899 |
| 844 | Kyunghyun Cho, and Iryna Gurevych. 2020a. | Abelardo Carlos Martínez Lorenzo, and Roberto | 900 |
| 845 | Adapterfusion: Non-destructive task composition for | Navigli. 2023. Incorporating graph information | 901 |
| 846 | transfer learning . <i>ArXiv</i> , abs/2005.00247. | in transformer-based AMR parsing . In <i>Findings</i> | 902 |
| 847 | Jonas Pfeiffer, Andreas Rücklé, Clifton Poth, Aishwarya | <i>of the Association for Computational Linguistics:</i> | 903 |
| 848 | Kamath, Ivan Vulić, Sebastian Ruder, Kyunghyun | <i>ACL 2023</i> , pages 1995–2011, Toronto, Canada. | 904 |
| 849 | Cho, and Iryna Gurevych. 2020b. Adapterhub: A | Association for Computational Linguistics. | 905 |
| 850 | framework for adapting transformers. In <i>Proceedings</i> | Alex Wang, Amanpreet Singh, Julian Michael, Felix | 906 |
| 851 | <i>of the 2020 Conference on Empirical Methods in Nat-</i> | Hill, Omer Levy, and Samuel R. Bowman. 2019. | 907 |
| 852 | <i>ural Language Processing: System Demonstrations</i> , | Glue: A multi-task benchmark and analysis platform | 908 |
| 853 | pages 46–54. | for natural language understanding . | 909 |
| 854 | Yanmin Qian, Xun Gong, and Houjun Huang. 2022. | Ruize Wang, Duyu Tang, Nan Duan, Zhongyu Wei, | 910 |
| 855 | Layer-wise fast adaptation for end-to-end multi- | Xuanjing Huang, Jianshu Ji, Guihong Cao, Daxin | 911 |
| 856 | accent speech recognition . <i>IEEE/ACM Transac-</i> | Jiang, and Ming Zhou. 2020. K-adapter: Infusing | 912 |
| 857 | <i>tions on Audio, Speech, and Language Processing</i> , | knowledge into pre-trained models with adapters . In | 913 |
| 858 | 30:2842–2853. | <i>Findings</i> . | 914 |
| 859 | Sylvestre-Alvise Rebuffi, Hakan Bilen, and Andrea | Xiaokai Wei, Shen Wang, Dejiao Zhang, Parminder | 915 |
| 860 | Vedaldi. 2017. Learning multiple visual domains | Bhatia, and Andrew O. Arnold. 2021. Knowledge | 916 |
| 861 | with residual adapters. In <i>Advances in Neural In-</i> | enhanced pretrained language models: A compreh- | 917 |
| 862 | <i>formation Processing Systems 30: Annual Confer-</i> | ensive survey . <i>ArXiv</i> , abs/2110.08455. | 918 |
| 863 | <i>ence on Neural Information Processing Systems 2017</i> , | Sondre Wold. 2022. The effectiveness of masked lan- | 919 |
| 864 | pages 506–516. | guage modeling and adapters for factual knowledge | 920 |
| 865 | Phillip Schneider, Tim Schopf, Juraj Vladika, Mikhail | injection . In <i>Proceedings of TextGraphs-16: Graph-</i> | 921 |
| 866 | Galkin, Elena Simperl, and Florian Matthes. 2022. | <i>-based Methods for Natural Language Processing</i> , | 922 |
| 867 | A decade of knowledge graphs in natural language | pages 54–59, Gyeongju, Republic of Korea. Associa- | 923 |
| 868 | processing: A survey . In <i>Proceedings of the 2nd</i> | tion for Computational Linguistics. | 924 |
| 869 | <i>Conference of the Asia-Pacific Chapter of the Asso-</i> | Qianqian Xie, Jennifer Amy Bishop, Prayag Tiwari, | 925 |
| 870 | <i>ciation for Computational Linguistics and the 12th</i> | and Sophia Ananiadou. 2022. Pre-trained language | 926 |
| 871 | <i>International Joint Conference on Natural Language</i> | models with domain knowledge for biomedical ex- | 927 |
| 872 | <i>Processing (Volume 1: Long Papers)</i> , pages 601–614, | tractive summarization . <i>Knowledge-Based Systems</i> , | 928 |
| 873 | Online only. Association for Computational Linguis- | 252:109460. | 929 |
| 874 | tics. | Yan Xu, Etsuko Ishii, Samuel Cahyawijaya, Zihan | 930 |
| 875 | Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. | Liu, Genta Indra Winata, Andrea Madotto, Dan Su, | 931 |
| 876 | Conceptnet 5.5: An open multilingual graph of gen- | and Pascale Fung. 2022. Retrieval-free knowledge- | 932 |
| 877 | eral knowledge. In <i>Proceedings of the Thirty-First</i> | grounded dialogue response generation with adapters . | 933 |
| 878 | <i>AAAI Conference on Artificial Intelligence, AAAI’17</i> , | In <i>Proceedings of the Second DialDoc Workshop on</i> | 934 |
| 879 | page 4444–4451. AAAI Press. | <i>Document-grounded Dialogue and Conversational</i> | 935 |
| 880 | Asa Cooper Stickland and Iain Murray. 2019. BERT | <i>Question Answering</i> , pages 93–107. Association for | 936 |
| 881 | and PALs: Projected attention layers for efficient | Computational Linguistics. | 937 |
| 882 | adaptation in multi-task learning . In <i>Proceedings of</i> | Shichuan Yu and Yan Yang. 2023. A new feature fusion | 938 |
| 883 | <i>the 36th International Conference on Machine Learn-</i> | method based on pre-training model for sequence | 939 |
| 884 | <i>ing</i> , volume 97 of <i>Proceedings of Machine Learning</i> | labeling . In <i>2023 6th International Conference on</i> | 940 |
| 885 | <i>Research</i> , pages 5986–5995. PMLR. | <i>Data Storage and Data Engineering (DSDE)</i> , pages | 941 |
| 886 | Abhisek Tiwari, Anisha Saha, Sriparna Saha, Pushpak | 26–31. | 942 |
| 887 | Bhattacharyya, and Minakshi Dhar. 2023. Experi- | Dongsheng Zou, Xiaotong Zhang, Xinyi Song, Yi Yu, | 943 |
| 888 | ence and evidence are the eyes of an excellent sum- | Yuming Yang, and Kang Xi. 2022. Multiway bidirec- | 944 |
| 889 | marizer! towards knowledge infused multi-modal | tional attention and external knowledge for multiple- | 945 |
| 890 | clinical conversation summarization . In <i>Proceedings</i> | choice reading comprehension . In <i>2022 IEEE Inter-</i> | 946 |
| 891 | <i>of the 32nd ACM International Conference on In-</i> | <i>national Conference on Systems, Man, and Cybernet-</i> | 947 |
| 892 | <i>formation and Knowledge Management, CIKM ’23</i> , | <i>ics (SMC)</i> , pages 694–699. | 948 |
| 893 | page 2452–2461, New York, NY, USA. Association | | |
| 894 | for Computing Machinery. | | |

A Supplementary Survey Data

A.1 Domain Distribution

See Figure 5.

A.2 Keywords in Task Distribution

See Figure 6.

A.3 Methodology

Articles on the following topics were excluded:

- Articles published before February 2, 2019
- Duplicate versions of the same article (when multiple versions of an article were found in different journals, only the most recent version was included)
- Articles where Adapters were used for NLP, but for use-cases other than knowledge-enhancement (such as few-shot learning or model debiasing)
- Articles written in a language other than English

The data extracted from each included document were:

- Source (journal or publication platform)
- Full reference
- Main topic area
- Facts of interest such as adapter architecture, domain, and downstream tasks within the papers
- A short summary of the study, including the main research questions and the answers

The collected data was tabulated to show:

- Source and publication dates of the studies
- Adapter architectures used in the papers
- Distribution of papers across domains (highlighting the biomedical domain)
- Distribution of papers across downstream tasks
- Results on biomedical NLP benchmarks (if relevant)

A.4 Acronyms

- BioNLP: Biomedical Natural Language Processing
- BLURB: Biomedical Language Understanding and Reasoning Benchmark (Gu et al., 2020)
- EE: Event Extraction
- EL: Entity Linking
- ES: Extractive Summarization
- ET: Entity Typing
- GLUE: General Language Understanding Evaluation (Wang et al., 2019)
- IE: Information Extraction
- KELM: Knowledge-Enhanced Language Model
- KGD: Knowledge-grounded Dialogue
- LAMA: Concept-Net Split of LAMA Probe (Petroni et al., 2019)
- LM: Language Modeling
- LLM: Large Language Model
- MT: Machine Translation
- NER: Named Entity Recognition
- NLI: Natural Language Inference
- NLP: Natural Language Processing
- OOD: Out-of-domain Detection
- QA: Question Answering
- RC: Reading Comprehension
- RE: Relation Extraction
- RCL: Relation Classification
- SA: Sentiment Analysis
- SC: Sentiment Classification
- SF: Speech Foundation
- SL: Sequence Labelling
- SMATCH: Semantic Match Score (Cai and Knight, 2013)

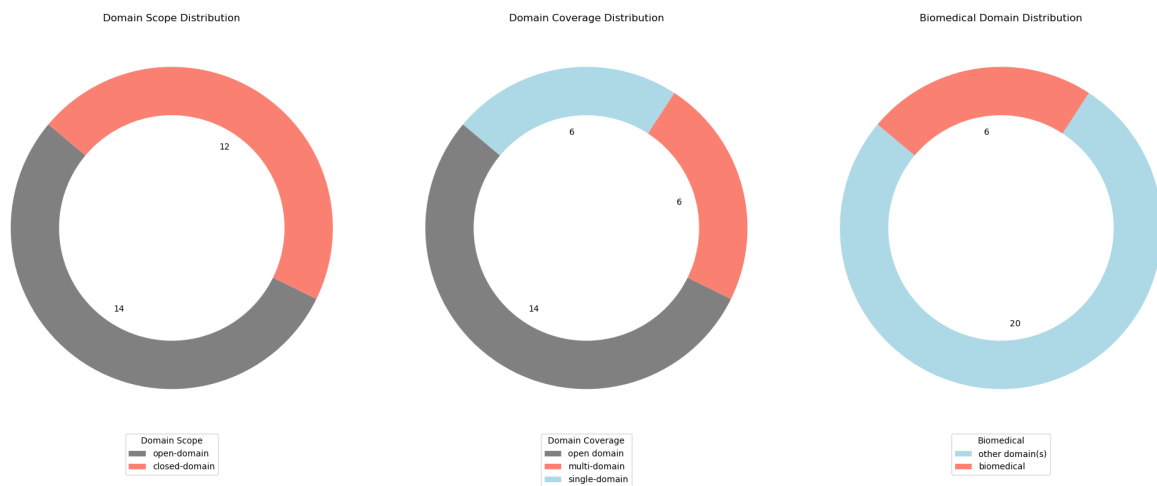


Figure 5: Distribution of domain scope, coverage, and the biomedical domain



Figure 6: Wordcloud of keywords in the task distribution

- SOTA: State-of-the-art
- SR: Speech Recognition
- STC: Sentence Classification
- TC: Text Classification
- TOD: Task-Oriented dialogue
- UMLS: Unified Medical Language System