

METAKP: On-Demand Keyphrase Generation

Anonymous ACL submission

Abstract

Traditional keyphrase prediction methods predict a single set of keyphrases per document, failing to cater to the diverse needs of users and downstream applications. To bridge the gap, we introduce on-demand keyphrase generation, a novel paradigm that requires keyphrases that conform to specific high-level goals or intents. For this task, we present METAKP, a large-scale benchmark comprising four datasets, 7500 documents, and 3760 goals across news and biomedical domains with human-annotated keyphrases. Leveraging METAKP, we design both supervised and unsupervised methods, including a multi-task fine-tuning approach and a self-consistency prompting method with large language models. The results highlight the challenges of supervised fine-tuning, whose performance is not robust to distribution shifts. By contrast, the proposed self-consistency prompting approach greatly improves the performance of large language models, enabling GPT-4o to achieve 0.548 SemF1, surpassing the performance of a fully fine-tuned BART-base model. Finally, we demonstrate the potential of our method to serve as a general NLP infrastructure, exemplified by its application in epidemic event detection from social media.

1 Introduction

Keyphrase prediction is an NLP task that has attracted long-lasting research interest (Witten et al., 1999; Hulth, 2003; Meng et al., 2017). Given documents from various domains such as academic writing, news, social media, or meetings, keyphrase extraction and keyphrase generation models output short phrases aiming at encapsulating the key entities and concepts mentioned by the document. Beyond a number of information retrieval applications (Kim et al., 2013; Tang et al., 2017; Boudin et al., 2020), keyphrase prediction methods are widely incorporated into the pipelines of other NLP tasks such as natural language generation (Yao et al.,

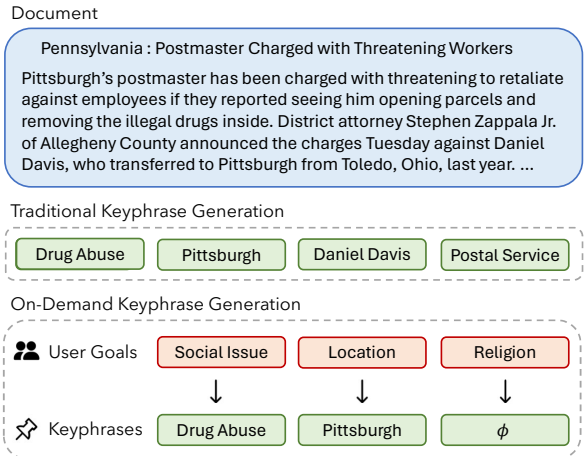


Figure 1: An illustration of on-demand keyphrase generation. Given diverse user goals, models are required to generate goal-conforming keyphrases or abstain.

2019; Li et al., 2020), text summarization (Dou et al., 2021), and text classification (Berend, 2011).

Despite their wide application in diverse scenarios, which may have diverse requirements on the types of keyphrases, existing keyphrase prediction methods generally follow a suboptimal assumption: for every document, the model shall predict a *single* application-agnostic set of keyphrases, which is then evaluated against a *monolithic* set of references (Wu et al., 2023b). This one-size-fits-all approach fails to cater to both downstream applications' varied requirements of the keyphrase predictions' topic and level of specificity and different expectations from human users with diverse backgrounds. To properly handle such diverse feedback, current approach could only rely on the sampler-rank strategy (Zhao et al., 2022; Wu et al., 2023a), which is largely inefficient. Besides, the single-reference setting also biases the intrinsic evaluation, of keyphrase prediction models, as high-frequency topics in keyphrase labels may significantly outweigh the long-tail keyphrases.

To tackle these challenges, we propose *on-demand keyphrase generation*, a novel paradigm

066 that predicts keyphrases conditioned on a *goal*
067 phrase that specifies the high-level category or
068 intent of the keyphrase (Figure 1). For existing
069 keyphrase prediction models, this task is challeng-
070 ing as it requires the predictions to be not only cap-
071 turing key information but also goal-conforming.
072 Furthermore, the models are required to accept
073 *open-vocabulary* goals, a significant step beyond
074 predicting keyphrases with predefined categories
075 or ontology (Park and Caragea, 2023).

076 To test on this new task, we meticulously curate
077 and release METAKP, a large-scale on-demand
078 keyphrase generation benchmark covering four
079 datasets, 7500 documents, and 3760 unique goals
080 from the news and the biomedical text domain. We
081 build a scalable labeling pipeline that combines
082 GPT-4 (OpenAI, 2023) and human annotators to
083 construct high-quality goals from keyphrases (Fig-
084 ure 2). For evaluation, we design two tasks: judg-
085 ing the relevance of goals and generating goal-
086 conforming keyphrases. For the latter, we employ
087 the state-of-the-art evaluation method (Wu et al.,
088 2023b) to conduct a semantic-based evaluation.

089 Using METAKP, we develop both super-
090 vised and unsupervised methods for on-demand
091 keyphrase generation. For the supervised method,
092 we design a multi-task fine-tuning approach to en-
093 able sequence-to-sequence pre-trained language
094 models to self-determine the relevance of a goal
095 and selectively generate keyphrases (Section 4.1).
096 Then, in Section 4.2, we introduce an unsupervised
097 self-consistency prompting approach leveraging
098 the strong ability of large language models (LLMs)
099 to propose goal-related keyphrase candidates and
100 their propensity to predict high quality keyphrases
101 with higher frequencies and ranks. Comprehensive
102 experiments reveal the following insights:

- 103 1. METAKP represents a challenging benchmark
104 for keyphrase generation. Flan-T5-XL, the
105 strongest fine-tuned model, only achieves an
106 average of 0.609 Satisfaction Rate across all
107 the datasets, and zero-shot prompting GPT-4o,
108 a strong LLM, only achieves 0.492 SR.
- 109 2. The proposed fine-tuning approach enables
110 jointly learning goal relevance judgment and
111 keyphrase generation without impeding each
112 task’s performance (Section 5.3).
- 113 3. The proposed self-consistency prompting ap-
114 proach greatly improves the performance of
115 LLMs, enabling GPT-4o to achieve 0.548
116 SemF1, surpassing the performance of a fully
117 fine-tuned BART-base model.

4. Supervised fine-tuning can fail to general-
118 ize on out-of-distribution testing data. By
119 contrast, LLM-based unsupervised method
120 achieves consistent performance in all the do-
121 mains, especially in the news domain, where
122 GPT-4o outperforms supervised Flan-T5-XL
123 by 19% in out-of-distribution testing. 124

125 Finally, we demonstrate the potential of on-
126 demand keyphrase generation as a general NLP
127 infrastructure. Specifically, we use event detec-
128 tion for epidemics prediction (Parekh et al., 2024)
129 as a test bed. By constructing simple goals from
130 event ontology and attempting to extract relevant
131 keyphrases from social media text, we show that
132 an on-demand keyphrase generation model has the
133 potential to extract epidemic-related trends similar
134 to an event detection model trained on task-specific
135 data. The benchmark and experimental code will
136 be released to facilitate further research.

137 2 Related Work

138 Keyphrase Prediction with Types This work
139 is closely related to prior work on modeling
140 keyphrases with pre-defined types or categories.
141 Early datasets are often derived from named en-
142 tity recognition, where keyphrase spans are ex-
143 tracted with entity type tags (QasemiZadeh and
144 Schumann, 2016; Augenstein et al., 2017; Luan
145 et al., 2018). Notable modeling approaches in-
146 clude using intermediate task for training strong
147 and transferable encoder representations (Park and
148 Caragea, 2020) as well as multi-task fine-tuning
149 (Park and Caragea, 2023). In addition, existing lit-
150 erature has explored inducing high-level type vari-
151 able for more accurate keyphrase prediction, such
152 as topic-guided keyphrase generation (Wang et al.,
153 2019; Zhang et al., 2022a), hierarchical keyphrase
154 generation (Wang et al., 2016; Chen et al., 2020;
155 Zhang et al., 2022b), as well as keyphrase comple-
156 tion (Zhao et al., 2021). Compare to these prior
157 work, our benchmark features a massive set of
158 open-vocabulary goals with wide domain cover-
159 age. We design novel supervised and unsupervised
160 modeling approaches that consider up-to-date tech-
161 niques such as large language models.

162 On-Demand Information Extraction Our work
163 resonates with the recent trend of designing flexi-
164 ble formulations for information extraction. For
165 instance, Zhong et al. (2021) propose a query-
166 focused formulation for the summarization task,
167 and Zhang et al. (2023) further extend the task to

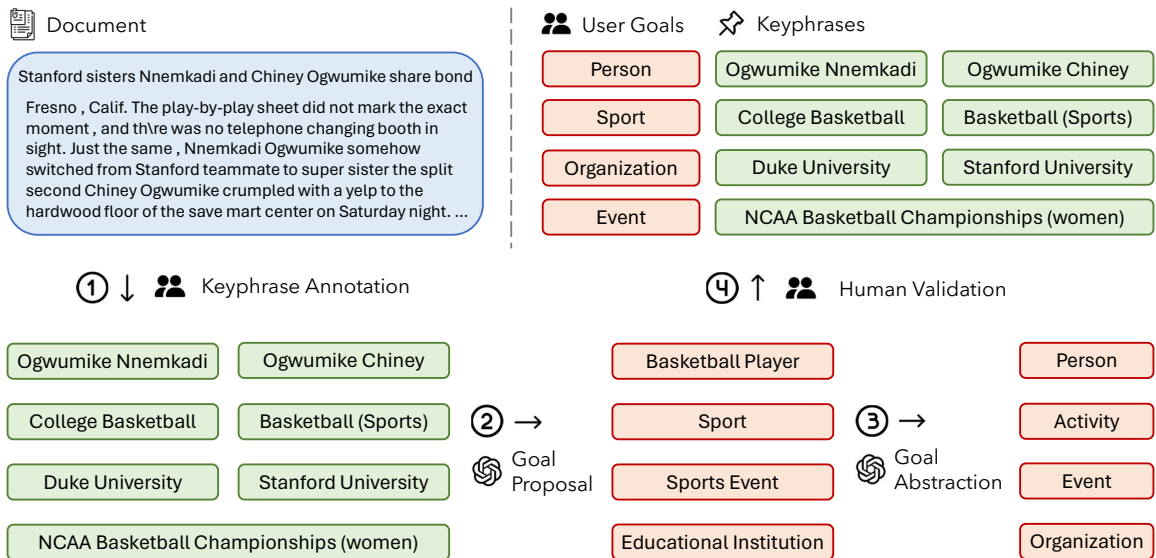


Figure 2: The annotation pipeline for METAKP. Starting from human-annotated keyphrases, GPT-4 is instructed to propose high-level goals and self-refine them. Finally, the goals are validated and filtered by humans.

include five constraints: Length, Extractiveness, Specificity, Topic, and Speaker. Recently, Jiao et al. (2023) introduce on-demand information extraction, where models are required to answer queries by extracting information from the associated text and organize it in a tabular format. By comparison, this work pioneers in defining and benchmarking the goal-following ability of keyphrase prediction models. Our resource and methodology lay the foundation for user-controllable keyphrase systems and flexible concept extraction infrastructures.

3 METAKP Benchmark

In this section, we formulate the on-demand keyphrase generation task and introduce the METAKP evaluation benchmark.

3.1 Problem Formulation

The traditional keyphrase prediction task is defined with a tuple: (document \mathcal{X} , reference set \mathcal{Y}). Given \mathcal{X} , a model directly generates all keyphrase hypotheses, with approximating \mathcal{Y} as the goal. For on-demand keyphrase generation, we introduce an open-vocabulary goal phrase g which describes a category of keyphrases specified by the user. The target of the model, then, is to generate a set of keyphrases based on (\mathcal{X}, g) to approximate the set of goal-conforming keyphrases $\mathcal{Y}_g \subseteq \mathcal{Y}$.

Figure 1 provides an intuitive example of the task. We note that for irrelevant goals, $\mathcal{Y}_g = \phi$, which means that an ideal model should not generate any keyphrases given such goals. In addition, although \mathcal{Y}_g varies according to the goal, the universal set of keyphrases \mathcal{Y} is assumed to be generally

fixed. In other words, g could be viewed as a query that specifies a target subset from \mathcal{Y} , which enables a wide range of choices for the modeling design.

3.2 Benchmark Creation Pipeline

To evaluate on-demand keyphrase generation, we curate METAKP, a large-scale multi-domain evaluation benchmark. The key challenge is to construct general, meaningful, and diverse goals that reflect high-level keyphrase types in real-world scenarios such as document indexing and search engines. To collect high quality goals, we design a model-in-the-loop annotation pipeline that combines GPT-4 (OpenAI, 2023) with human annotators to infer goals reversely from keyphrase annotations (Figure 2), with four steps detailed as follows.

Keyphrase Annotation by Human Given the document \mathcal{X} , human annotators specify the set of all the possible keyphrases \mathcal{Y} . For METAKP, we directly leverage the expert-curated keyphrases from the respective keyphrase prediction datasets.

Goal Proposal We instruct GPT-4 to propose a high-level goal for each of the keyphrases, and the same goal could be shared by multiple keyphrases¹. Concretely, given \mathcal{X}, \mathcal{Y} , GPT-4 returns a mapping from goals to keyphrases. We present the prompt for this step in Appendix A.

Goal Abstraction After the previous step, a draft goal has been associated with each keyphrase. Although the proposed goals are relevant, we observe that they are sometimes overly specific. There-

¹We use gpt-4-0613 via the OpenAI API.

fore, we instruct GPT-4 to perform a round of *self-refinement*, where it attempts to propose a more abstract version for each of the goals in the previous round, or keep the original goals if they are already high-level enough. The full prompt for this step is presented in Appendix A.

Human Validation We qualitatively find that the outputs from two GPT-4 annotation iterations are sufficiently abstract and diverse. To further improve the quality of the goals and reduce the level of duplication, two of the authors conduct a round of filtering to obtain the final goal annotations. As this step does not entail adding new goals, the annotators achieve a high inter-annotator agreement (detailed in the next section) following the annotation guideline, which we present in Appendix A. Finally, we create an instance for each of the filtered goals, taking the form $(\mathcal{X}, g_i, \mathcal{Y}_{g_i})$.

3.3 Dataset Statistics

We execute the aforementioned goal construction pipeline on four keyphrase prediction datasets covering two domains: news and biomedical text. For each domain, we create both an in-distribution and an out-of-distribution test set.

- **KPTimes** (Gallina et al., 2019) is a large-scale keyphrase generation dataset in the news domain. The documents are sourced from from New York Times and the keyphrases are curated by professional editors.
- **DUC2001** (Wan and Xiao, 2008) is a widely used keyphrase extraction dataset with news articles collected from TREC-9, paired with human-annotated keyphrases.
- **KPBiomed** (Houbre et al., 2022) is a large-scale dataset containing PubMed abstracts paired with keyphrases annotated by paper authors themselves.
- **Pubmed** (Schutz, 2008) is a traditional keyphrase extraction dataset in the biomedical domain with documents and keyphrases extracted from the PubMed Central.

We curate a test set using each of these datasets and construct two domain-specific train/validation sets sampled from the training sets from KPTimes and KPBiomed. Table 1 and Figure 3 presents the basic statistics of the final datasets. Besides its domain coverage, one strength of METAKP is

Source	Split	#Doc	#Inst	#Goal	Goal	#KP/Goal
KPTimes	Train	1859	7502	1083	1.43	1.32
	Val	100	392	148	1.46	1.37
	Test	984	3836	679	1.41	1.33
DUC2001	Test	308	1642	549	1.50	1.53
	Train	1886	7807	1311	1.75	1.27
KPBiomed	Val	100	404	189	1.75	1.32
	Test	994	4136	865	1.76	1.27
	Test	1269	4988	843	1.82	1.33

Table 1: Basic statistics of METAKP. #Inst = number of instances in the form $(\mathcal{X}, g, \mathcal{Y}_g)$. |Goal| refers to the average number of words in g . Finally, #KP/Goal corresponds to the average cardinality of \mathcal{Y}_g .

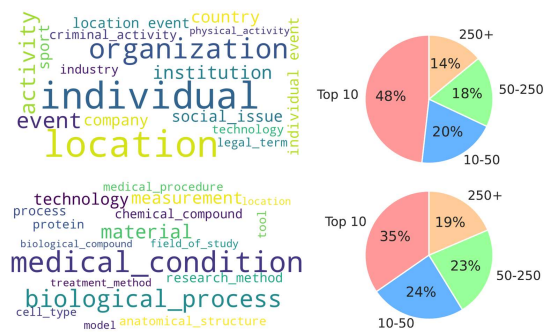


Figure 3: A visualization of the goal distribution for the news domain (top) and the biomedical domain (bottom). METAKP features both high-frequency goals and a diverse long-tail goal distribution.

its *diverse* coverage: together, the dataset covers 3760 unique goals, including diverse topics and subjects. While 40% of the instances correspond to the 10 most popular goals in each domain, a substantial number of goals also fall into the long tail distribution, posing significant new challenge in understanding the goal semantics.

To construct METAKP, the two-staged GPT-4 annotation costed approximately 500 USD, and the human annotators worked for approximately 80 hours in total on final data filtering. We randomly sample 50 documents each from KPTimes and KP-Biomed, on which the annotators reach 0.699 Cohen’s Kappa for inter-annotator agreement. Then, the annotators work on the rest documents separately. When ambiguous cases are found, a discussion is conducted to reach agreement.

Irrelevant Goal Sampling To test the ability of keyphrase generation models to abstain from generating keyphrases given irrelevant goals, for each document, we additionally construct a set of irrelevant goals. Concretely, we cluster the goals in the labelled data and use each document’s existing goals as anchors to sample goals that are likely to be irrelevant to the document and thus it is unlikely

that a keyphrase corresponding to the sampled goal exists for the document. We present the algorithm in the Appendix A.3. Using the algorithm, a balanced training set was created for training supervised methods for goal relevance judgment.

3.4 Evaluation Metric

With METAKP, we design two tasks to comprehensively evaluate a model’s ability to perform on-demand keyphrase generation.

Goal Relevance Assessment This task aims to test whether a model can correctly distinguish irrelevant goals that cannot yield any keyphrase from the relevant goals. As we will show in Section 6, this skill is also crucial to enable a wide application of on-demand keyphrase generation models. Following recent literature on abstention (Feng et al., 2024), we use **Abstain F1** as the evaluation metric, which is defined as the harmonic mean of the precision and the recall of a model refusing to generate keyphrases for irrelevant goals.

Goal-Oriented Keyphrase Generation Given document \mathcal{X} , a list of goals g_1, g_2, \dots, g_n , and references $\mathcal{Y}_{g_1}, \mathcal{Y}_{g_2}, \dots, \mathcal{Y}_{g_n}$, we evaluate a model’s predictions P_1, P_2, \dots, P_n with two metrics:

1. **Reference Agreement**, which assesses the model’s ability to generate keyphrases specifically corresponding to the goal g_i . Concretely, we calculate and report $SemF1(Y_{g_i}, P_i)$, following Wu et al. (2023b).
2. **Satisfactory Rate (SR)**, which assesses the frequency of the model generating high-quality keyphrases. Concretely, we calculate and report $SR((\mathcal{Y}_{g_1}, P_1), \dots, (\mathcal{Y}_{g_n}, P_n))$ as the percentage of goals that have $SemF1(Y_{g_i}, P_i)$ greater than a threshold².

4 Modeling Approach

In this section, we introduce two modeling approaches for on-demand keyphrase generation: a multi-task learning approach for fine-tuning sequence-to-sequence pre-trained language models, and a self-consistency decoding approach for prompting large language models (LLMs).

4.1 Multi-Task Supervised Fine-tuning

Previous literature has demonstrated the effectiveness of fine-tuning sequence-to-sequence pre-

²We fix $\tau = 0.6$. This decision is based Wu et al. (2023b), which suggests that the embedding model for $SemF1$ assigns a similarity score of approximately 0.6 for name variations.

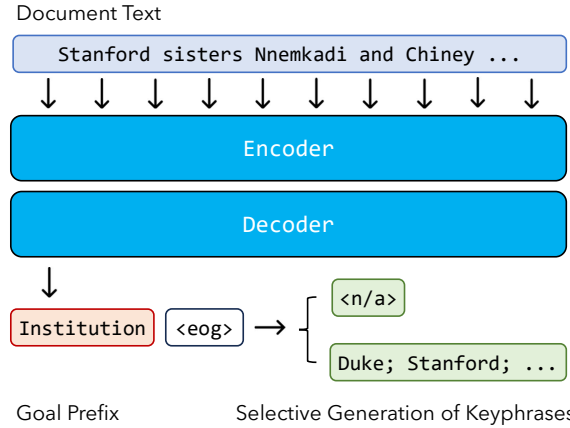


Figure 4: A visualization of the inference process of the proposed sequence-to-sequence generation approach. Based on the document and the goal prefix, the model self-decides the relevance of the goal and selectively generates the keyphrases for relevant goals only.

trained language models for keyphrase generation (Kulkarni et al., 2022; Wu et al., 2022, 2023a). However, it is unclear how these sequence prediction approaches could be adopted for on-demand keyphrase generation. To bridge this gap, we introduce a novel formulation to train a sequence-to-sequence model to autoregressively (1) assess the relevance of goals and (2) jointly consider the document as well as a desired goal to predict keyphrases.

Concretely, we formulate on-demand keyphrase generation as a hierarchical composition of two token prediction tasks. As shown in Figure 4, with the document fed in the encoder, the decoder first models $P(g_i|\mathcal{X})$, the likelihood of g_i being a high-quality relevant goal proposed by real users. The model verbalizes this probability in $P(\langle n/a \rangle|\mathcal{X}, g_i)$, a special token for rejecting irrelevant goals. If the goal is determined as relevant, the model proceeds generating the keyphrases according to the distribution $P(\mathcal{Y}_{g_i}|\mathcal{X}, g_i)$ it learned.

Inference We use prefix-controlled decoding for inference. g_i , followed by a special end-to-goal token $\langle eog \rangle$, is fixed as the decoder’s start of generation. Then, we use autoregressive decoding to let the model self-assess the relevance of goal and automatically decide the keyphrases to generate.

Training We design a multi-task learning procedure to directly supervise the model on $P(\langle n/a \rangle|\mathcal{X}, g_i)$ and $P(\mathcal{Y}_{g_i}|\mathcal{X}, g_i)$ with a mixture of relevant and irrelevant goals. As the goals provided by users could be arbitrary, we do not directly supervise the model on $P(g_i|\mathcal{X})$.

Remark We note that the proposed approach has several advantages. First, both the goal relevance assessment and the keyphrase prediction process are streamlined in a single sequence prediction process, removing the need for separate architecture or inference pass. Second, since g_i is not fed to the encoder, our model avoids the goal being diluted by the long input context and enables efficient inference by reusing the encoded input representation for predicting keyphrases with different goals.

4.2 Prompting Large Language Models

Large language models (LLMs) that are tuned to follow human instructions have been shown to adapt well to a massive number of tasks defined through human queries (Ouyang et al., 2022; OpenAI, 2023). They have also been demonstrated to achieve promising keyphrase extraction or keyphrase generation performance, especially with semantic-based evaluation (Song et al., 2023; Wu et al., 2023b). As on-demand keyphrase extraction could be easily formulated as an instruction-following task, we investigate the potential of LLMs as an unsupervised approach. We start with a simple instruction for judging a goal’s relevance:

Decide if you should reject the high-level category given the title and abstract of a document. One could use the high-level category to write keyphrases from the document.

as well as another instruction for keyphrase generation based on a goal:

Generate present and absent keyphrases belonging to the high-level category from the given text.

Our preliminary experiments show that the first instruction already achieves a strong performance in deciding the goal relevance, even approaching supervised models (Section 5.2). However, when it comes to keyphrase generation, LLMs intriguingly misinterpret the task as named entity extraction: they often generate an almost exhaustive list of goal-related entities. To correct this behavior, we hypothesize that LLMs tend to generate salient entities more frequently and at an earlier location of the prediction sequence. Inspired by Wang et al. (2023), we thus design a novel self-consistency decoding process to leverage the rank and frequency information in LLMs’ samples to filter out phrases that encode the most important information.

Concretely, using the same instruction and input, we sample K prediction sequences (s_1, \dots, s_K) from the LLM independently, each of which contains a variable number of keyphrases. Then, for

each keyphrase p , we define its saliency score as:

$$score(p) = \frac{freq(p)}{K} \times \frac{freq(p)}{\sum_{i=1, \dots, K} rank(s_i, p)},$$

where $freq(p)$ returns the frequency of p in all the samples and $rank(s_i, p)$ returns the rank of p in s_i (starting from 1) or 0 if $p \notin s_i$. The first term rewards keyphrases that frequently present in the samples, and the second term rewards keyphrases with a higher rank. Together, the score is defined to range 0 from 1 regardless of the number of samples or the number of keyphrases a model generates per sample. Finally, we apply threshold filtering and only retain the high quality keyphrases with $score(p)$ greater than or equal to a threshold τ .

5 Experiments

5.1 Experimental Setup

Supervised Fine-tuning Using the proposed objective, we fine-tune four sequence-to-sequence models: BART-base/large (Lewis et al., 2020) and Flan-T5-large/XL (Longpre et al., 2023), with diverse sizes ranging from 140M to 3B. We train the models for 20 epochs with batch size 64, learning rate $3e-5$, the Adam optimizer, and a linear decay with 50 warmup steps. The best model checkpoint is chosen based on the keyphrase generation performance on the validation set.

Prompting We use gpt-3.5-turbo-0125 and the gpt-4o-2024-05-13 models via the OpenAI API. We will denote the models as GPT-3.5-Turbo and GPT-4o. We use separate prompts for goal relevance judgment and on-demand keyphrase generation. For the first task, greedy search is used. For the second task, we generate 10 samples with $p = 0.95$ and temperature = 0.7. The output length is limited 30 tokens, which can accommodate approximately 10 keyphrases. Finally, for filtering, we use $\tau = 0.3$ for all the datasets.

We document the full implementation details in Appendix B, including the prompt for language language models, the post-processing process, as well as the details for hyperparameter tuning.

5.2 Main Results

We present the main results for the two tasks in Figure 5 and Table 2.

Goal Relevance Assessment According to Figure 5, we find both supervised fine-tuning and unsupervised prompting reaches a high performance for

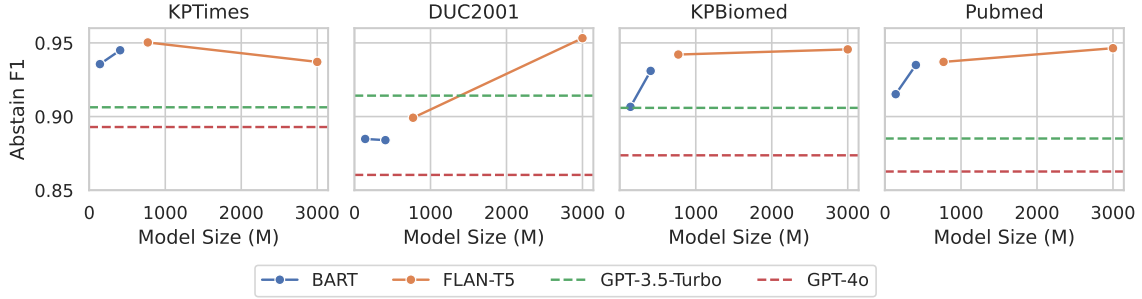


Figure 5: Goal relevance judgment results of different types of models. Zero-shot prompting LLMs achieves a high performance, despite slightly falling below supervised models. Also, GPT-4o does improve over GPT-3.5-Turbo.

Model	Size	Method	KPTimes [✧]		DUC2001 [✧]		KPBiomEd [♣]		Pubmed [♣]		Average	
			SemF1	SR	SemF1	SR	SemF1	SR	SemF1	SR	SemF1	SR
Supervised Methods												
BART-base	140M	No Goal	0.395	0.192	0.299	0.089	0.300	0.107	0.305	0.196	0.325	0.146
		MetaKP	0.728	0.699	0.447	0.319	0.508	0.417	0.504	0.406	0.547	0.460
BART-large	406M	No Goal	0.399	0.196	0.306	0.081	0.297	0.074	0.290	0.070	0.323	0.105
		MetaKP	0.752	0.738	0.469	0.336	0.545	0.461	0.529	0.437	0.574	0.493
FLan-T5-large	770M	MetaKP	0.765	0.758	0.488	0.360	0.578	0.506	0.572	0.501	0.601	0.531
FLan-T5-XL	3B	MetaKP	0.763	0.757	0.484	0.361	0.594 [†]	0.530 [†]	0.593 [†]	0.526 [†]	0.609 [†]	0.544 [†]
Unsupervised Methods												
GPT-3.5-Turbo	-	Zero-Shot	0.452	0.221	0.499	0.290	0.421	0.166	0.444	0.217	0.454	0.224
		Sample + SC	0.518	0.406	0.572	0.516	0.513	0.423	0.472	0.376	0.519	0.430
GPT-4o	-	Zero-Shot	0.491	0.281	0.526	0.374	0.480	0.278	0.469	0.262	0.492	0.299
		Sample + SC	0.552	0.460	0.578 [†]	0.535 [†]	0.529	0.451	0.532	0.453	0.548	0.475

Table 2: Experiment results of supervised and unsupervised methods on-demand keyphrase generation. We use different superscripts to denote results that are reported using the models trained on KPTimes (✧) and KPBiomed (♣). SR = satisfaction rate. SC = self-consistency prompting. The best results are boldfaced. [†]statistically significantly better than the second highest result with $p < 0.01$, tested via paired t-test.

476 assessing whether a goal, as indicated by over 0.85
 477 Abstain F1 scores across all datasets. As model
 478 size scales, the out-of-distribution performance
 479 scales more readily, while the in-distribution performance
 480 plateaus at FLan-T5-large. With large language
 481 models, we observe strong performance especially
 482 on DUC2001, surpassing the performance
 483 of FLan-T5-large trained on KPTimes.

484 **Keyphrase Generation** The main results for
 485 keyphrase generation are presented in Table 2. For
 486 supervised methods, we additionally include a "No
 487 Goal" baseline, where the model is fine-tuned to
 488 generate all the keyphrases for the same document
 489 at once. For both BART-base and BART-large,
 490 this baseline achieves a low performance, indicating
 491 the challenging nature of directly leveraging a
 492 keyphrase generation model for the proposed task.
 493 By comparison, the proposed goal-directed fine-tuning
 494 approach improves the performance by a large margin,
 495 with the best FLan-T5-XL model achieving 0.609
 496 SemF1 and 0.544 satisfaction rate. On the other
 497 hand, directly zero-shot prompting large language
 498 models already achieves more superior performance
 499 compared to the su-

500 pervised models trained without any goal. The
 501 proposed self-consistency further improves the
 502 performance substantially, allowing GPT-4o achieve
 503 0.548 SemF1 and 0.475 satisfaction rate. Notably,
 504 results demonstrate that the LLM-based approach
 505 has the potential to be more generalizable. On
 506 DUC2001, all supervised models trained on KPTimes
 507 demonstrate a poor performance. By contrast,
 508 both GPT-3.5-Turbo and GPT-4o are able to surpass
 509 the performance of all supervised models.

5.3 Analyses

510 **Which parameter affects LLMs the most?** In
 511 Figure 6, we use KPTimes' validation set to
 512 investigate the sensitivity of the LLM-based
 513 approach to three hyperparameters: number of
 514 samples (K), threshold τ , and context length
 515 of the input. Although multiple samples are
 516 essential to high performance, more samples
 517 after two only help marginally. In addition,
 518 our method is insensitive to the threshold
 519 setting - the best performance can be obtained
 520 by multiple settings between 0.25 and 0.45.
 521 Finally, while GPT-3.5-Turbo exhibits a slight
 522 performance drop with longer context, GPT-4o
 523 is robust to context length variations.

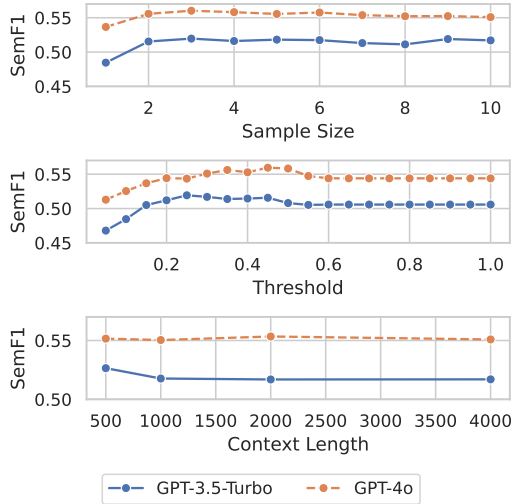


Figure 6: Sensitivity of the self-consistency prompting approach’s performance to number of samples, settings of threshold τ , and the input length on KPTimes. The results on KPBiomed is presented in Figure 12.

Objective	ID		OOD	
	AF1	SR	AF1	SR
Training on KPTimes				
Multi-task Learning	0.936	0.699	0.885	0.319
Goal Relevance Only	0.928	-	0.898	-
Keyphrase Only	-	0.692	-	0.316
Training on KPBiomed				
Multi-task Learning	0.907	0.417	0.915	0.406
Goal Relevance Only	0.917	-	0.916	-
Keyphrase Only	-	0.425	-	0.407

Table 3: Ablation study on the multi-task learning setup. AF1 = Abstain F1, SR = Satisfaction Rate.

Does multi-task learning harm each individual task’s performance? In Table 3, we conduct an ablation study with BART-base on the supervised training loss. For each ablated component, we mask out the corresponding tokens when calculating the loss. Overall, combining the two learning objectives do not significantly harm the performance compared to only learning individual tasks, while incurring much less computational overhead. In fact, on KPTimes, the two tasks are constructive - learning goal relevance helps generating better goal-conforming keyphrases, and vice versa.

6 META KP in the Wild: Event Detection

Finally, we demonstrate the potential of on-demand keyphrase generation as general NLP infrastructure, using event detection (ED) as a case study.

We leverage the testing dataset used in SPEED (Parekh et al., 2024), which contains time-stamped social media posts related to Monkeypox³. From

³We solicited the dataset and outputs from the authors.

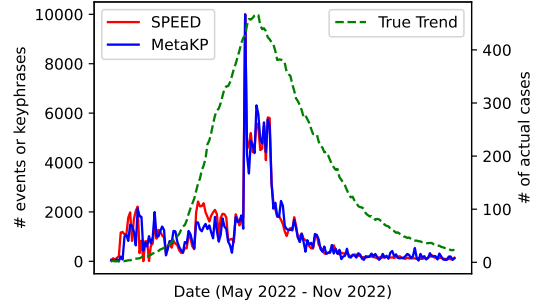


Figure 7: Number of events/keyphrases extracted for Monkeypox as a function of time. The true trend and SPEED outputs are solicited from Parekh et al. (2024).

the SPEED ontology, we curate seven epidemic-related goals: *disease infection*, *epidemic spread*, *epidemic prevention*, *epidemic control*, *symptom*, *recover from disease*, *death from epidemic*. Then, we run a FLAN-T5-large model trained on all training data from META KP to assess the relevance of each goal against each social media post. If the probability of $\langle n/a \rangle$ following $\langle eog \rangle$ is greater than 0.001, the goal is judged relevant to the post and thus the underlying event is likely.

As shown in Figure 7, we observe that this keyphrase-based paradigm is able to extract trends that are similar to an ED model trained on SPEED (Parekh et al., 2024). Intuitively, given a sentence containing "getting vaccination", instead of focusing on the trigger "get", on-demand keyphrase generation is able focus more on "vaccination", given the goal "epidemic control". In this way, on-demand keyphrase generation models can both be naturally repurposed for ED and also promises to extract supporting topics related to the the event.

7 Conclusion

We introduce on-demand keyphrase generation, targeting the need for dynamic, goal-oriented keyphrase prediction tailored to diverse applications and user demands. A large-scale, multi-domain, human-verified benchmark META KP was curated and introduced. We designed and evaluated both supervised and unsupervised methods on META KP, highlighting the strengths of self-consistency prompting with large language models. This approach significantly outperformed traditional fine-tuning methods under domain shifts, showcasing its robustness and the broader applicability of our methodology. Finally, we underscore the versatility of on-demand keyphrase generation in practical applications such as epidemic event extraction, promising a new direction for keyphrase generation as general NLP infrastructure.

582 Limitations

583 In this work, we propose the novel on-demand
584 keyphrase generation paradigm. In the future,
585 several exciting directions exist for extending the
586 paradigm as well as the METAKP benchmark:

- 587 1. **Multi-lingual Keyphrase Generation.**
588 METAKP only covers data in English.
589 Further benchmarking and enhancing the
590 multilingual and cross-lingual on-demand
591 keyphrase generation ability is an important
592 future direction.
- 593 2. **Wider Domain Coverage.** We mainly focus
594 on the news and the biomedical text domain as
595 they have been shown as important application
596 domains for keyphrase generation.
- 597 3. **Flexible Instructions.** In this work, the "de-
598 mand" from the users are generally defined
599 as topics or categories of keyphrases. How-
600 ever, future work could expand this definition
601 to include demands that specify stylistic con-
602 straints such as the number of keyphrases, the
603 length, and their formality.

604 Ethics Statement

605 As a new task and paradigm, on-demand keyphrase
606 generation may bring new security risks and ethi-
607 cal concerns. To begin with, although keyphrase
608 generation models generally have outstanding un-
609 derstanding of phrase saliency, they generally have
610 a shallower understanding of semantics and factu-
611 ality. Thus, when pairing keyphrases with goals,
612 potential misinformation could be created. For
613 instance, when queried with "cure" as a goal, a
614 model may return certain concepts that are factu-
615 ally wrong. In addition, when queries contain cer-
616 tain occupations as goals, a keyphrase generation
617 model may reinforce existing gender stereotypes
618 by selectively generating and ignoring entities with
619 a certain gender. We view these possibilities as
620 potential risks and encourage a thorough redteaming
621 process before deploying on-demand keyphrase
622 generation systems in real-world products.

623 We use KPTime and KPBiomed data dis-
624 tributed by the original authors. For DUC2001 and
625 PubMed, we access the data via ake-datasets⁴. KP-
626 Time was released under Apache-2.0 license, and
627 we cannot find licensing information for DUC2001,
628 KPBiomed, and PubMed. ake-datasets was also
629 released under Apache-2.0. No additional prepro-

630 cessing is performed in METAKP except lower-
631 casing and tokenization. While we mainly rely on
632 the original authors for dataset screening to remove
633 sensitive and harmful information, we also actively
634 monitor the data quality during in the human filter-
635 ing process and remove any document that could
636 cause privacy or ethics concerns. As OpenAI mod-
637 els are involved in the data curation process, our
638 code and datasets will be released with MIT license
639 with a research-only use permission.

References 640

- 641 Isabelle Augenstein, Mrinal Das, Sebastian Riedel,
642 Lakshmi Vikraman, and Andrew McCallum. 2017.
643 *SemEval 2017 task 10: ScienceIE - extracting*
644 *keyphrases and relations from scientific publications.*
645 *In Proceedings of the 11th International Workshop*
646 *on Semantic Evaluation (SemEval-2017)*, pages 546–
647 555, Vancouver, Canada. Association for Computa-
648 tional Linguistics.
- 649 Gábor Berend. 2011. *Opinion expression mining by ex-*
650 *ploiting keyphrase extraction.* *In Proceedings of 5th*
651 *International Joint Conference on Natural Language*
652 *Processing*, pages 1162–1170, Chiang Mai, Thailand.
653 Asian Federation of Natural Language Processing.
- 654 Florian Boudin, Ygor Gallina, and Akiko Aizawa. 2020.
655 *Keyphrase generation for scientific document re-*
656 *trieval.* *In Proceedings of the 58th Annual Meeting of*
657 *the Association for Computational Linguistics*, pages
658 1118–1126, Online. Association for Computational
659 Linguistics.
- 660 Wang Chen, Hou Pong Chan, Piji Li, and Irwin King.
661 2020. *Exclusive hierarchical decoding for deep*
662 *keyphrase generation.* *In Proceedings of the 58th*
663 *Annual Meeting of the Association for Computational*
664 *Linguistics*, pages 1095–1105, Online. Association
665 for Computational Linguistics.
- 666 Zi-Yi Dou, Pengfei Liu, Hiroaki Hayashi, Zhengbao
667 Jiang, and Graham Neubig. 2021. *GSum: A gen-*
668 *eral framework for guided neural abstractive summa-*
669 *rization.* *In Proceedings of the 2021 Conference of*
670 *the North American Chapter of the Association for*
671 *Computational Linguistics: Human Language Tech-*
672 *nologies*, pages 4830–4842, Online. Association for
673 Computational Linguistics.
- 674 Shangbin Feng, Weijia Shi, Yike Wang, Wenxuan
675 Ding, Vidhisha Balachandran, and Yulia Tsvetkov.
676 2024. *Don't hallucinate, abstain: Identifying LLM*
677 *knowledge gaps via multi-llm collaboration.* *CoRR*,
678 abs/2402.00367.
- 679 Ygor Gallina, Florian Boudin, and Beatrice Daille. 2019.
680 *KPTime: A large-scale dataset for keyphrase gener-*
681 *ation on news documents.* *In Proceedings of the 12th*
682 *International Conference on Natural Language Gen-*
683 *eration*, pages 130–135, Tokyo, Japan. Association
684 for Computational Linguistics.

⁴<https://github.com/boudinfl/ake-datasets>

799	Xiaojun Wan and Jianguo Xiao. 2008. Single document	854
800	keyphrase extraction using neighborhood knowledge.	855
801	In <i>AAAI</i> , volume 8, pages 855–860.	856
802	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le,	857
803	Ed H. Chi, Sharan Narang, Aakanksha Chowdhery,	858
804	and Denny Zhou. 2023. Self-consistency improves	859
805	chain of thought reasoning in language models . In	860
806	<i>The Eleventh International Conference on Learning</i>	
807	<i>Representations</i> .	
808	Yue Wang, Jing Li, Hou Pong Chan, Irwin King,	
809	Michael R. Lyu, and Shuming Shi. 2019. Topic-	
810	aware neural keyphrase generation for social media	
811	language . In <i>Proceedings of the 57th Annual Meet-</i>	
812	<i>ing of the Association for Computational Linguistics</i> ,	
813	pages 2516–2526, Florence, Italy. Association for	
814	Computational Linguistics.	
815	Yunli Wang, Yong Jin, Xiaodan Zhu, and Cyril Goutte.	
816	2016. Extracting discriminative keyphrases with	
817	learned semantic hierarchies . In <i>Proceedings of COL-</i>	
818	<i>ING 2016, the 26th International Conference on Com-</i>	
819	<i>putational Linguistics: Technical Papers</i> , pages 932–	
820	942, Osaka, Japan. The COLING 2016 Organizing	
821	Committee.	
822	Ian H Witten, Gordon W Paynter, Eibe Frank, Carl	
823	Gutwin, and Craig G Nevill-Manning. 1999. Kea:	
824	Practical automatic keyphrase extraction. In <i>Pro-</i>	
825	<i>ceedings of the fourth ACM conference on Digital</i>	
826	<i>libraries</i> , pages 254–255.	
827	Di Wu, Wasi Ahmad, and Kai-Wei Chang. 2023a. Re-	
828	thinking model selection and decoding for keyphrase	
829	generation with pre-trained sequence-to-sequence	
830	models . In <i>Proceedings of the 2023 Conference on</i>	
831	<i>Empirical Methods in Natural Language Processing</i> ,	
832	pages 6642–6658, Singapore. Association for Com-	
833	putational Linguistics.	
834	Di Wu, Wasi Ahmad, Sunipa Dev, and Kai-Wei	
835	Chang. 2022. Representation learning for resource-	
836	constrained keyphrase generation . In <i>Findings of the</i>	
837	<i>Association for Computational Linguistics: EMNLP</i>	
838	2022, pages 700–716, Abu Dhabi, United Arab Emi-	
839	rates. Association for Computational Linguistics.	
840	Di Wu, Da Yin, and Kai-Wei Chang. 2023b. Kpeval:	
841	Towards fine-grained semantic-based evaluation of	
842	keyphrase extraction and generation systems .	
843	Lili Yao, Nanyun Peng, Ralph Weischedel, Kevin	
844	Knight, Dongyan Zhao, and Rui Yan. 2019. Plan-	
845	and-write: Towards better automatic storytelling . In	
846	<i>Proceedings of the AAAI Conference on Artificial</i>	
847	<i>Intelligence</i> , volume 33, pages 7378–7385.	
848	Yusen Zhang, Yang Liu, Ziyi Yang, Yuwei Fang, Yulong	
849	Chen, Dragomir Radev, Chenguang Zhu, Michael	
850	Zeng, and Rui Zhang. 2023. MACSum: Control-	
851	lable summarization with mixed attributes . <i>Transac-</i>	
852	<i>tions of the Association for Computational Linguis-</i>	
853	<i>tics</i> , 11:787–803.	
	Yuxiang Zhang, Tao Jiang, Tianyu Yang, Xiaoli Li, and	854
	Suge Wang. 2022a. HTKG: deep keyphrase gen-	855
	eration with neural hierarchical topic guidance . In	856
	<i>SIGIR '22: The 45th International ACM SIGIR Con-</i>	857
	<i>ference on Research and Development in Information</i>	858
	<i>Retrieval, Madrid, Spain, July 11 - 15, 2022</i> , pages	859
	1044–1054. ACM.	860
	Yuxiang Zhang, Tianyu Yang, Tao Jiang, Xiaoli Li, and	861
	Suge Wang. 2022b. Hyperbolic deep keyphrase gen-	862
	eration . In <i>Machine Learning and Knowledge Dis-</i>	863
	<i>covery in Databases - European Conference, ECML</i>	864
	<i>PKDD 2022, Grenoble, France, September 19-23,</i>	865
	<i>2022, Proceedings, Part II</i> , volume 13714 of <i>Lecture</i>	866
	<i>Notes in Computer Science</i> , pages 521–536. Springer.	867
	Guangzhen Zhao, Guoshun Yin, Peng Yang, and Yu Yao.	868
	2022. Keyphrase generation via soft and hard seman-	869
	tic corrections . In <i>Proceedings of the 2022 Confer-</i>	870
	<i>ence on Empirical Methods in Natural Language Pro-</i>	871
	<i>cessing</i> , pages 7757–7768, Abu Dhabi, United Arab	872
	Emirates. Association for Computational Linguistics.	873
	Yu Zhao, Jia Song, Huali Feng, Fuzhen Zhuang, Qing	874
	Li, Xiaojie Wang, and Ji Liu. 2021. Deep keyphrase	875
	completion . <i>CoRR</i> , abs/2111.01910.	876
	Ming Zhong, Da Yin, Tao Yu, Ahmad Zaidi, Mutethia	877
	Mutuma, Rahul Jha, Ahmed Hassan Awadallah, Asli	878
	Celikyilmaz, Yang Liu, Xipeng Qiu, and Dragomir	879
	Radev. 2021. QMSum: A new benchmark for query-	880
	based multi-domain meeting summarization . In <i>Pro-</i>	881
	<i>ceedings of the 2021 Conference of the North Amer-</i>	882
	<i>ican Chapter of the Association for Computational</i>	883
	<i>Linguistics: Human Language Technologies</i> , pages	884
	5905–5921, Online. Association for Computational	885
	Linguistics.	886

Supplementary Material: Appendices

A META KP Construction Details

In this section, we describe the details of the construction process of META KP.

A.1 GPT-4 Annotation

Goal Proposal In Figure 8, we show the prompt used to instruct GPT-4 to propose goals from the document and human-annotated keyphrases. We truncate the document body to four sentences as its role is only providing essential contextual. The LLM is instructed to propose all the goals for all the keyphrases together, which helps the model group together keyphrases that share the same goal.

```
Document Title: {title}
First 4 sentences of the document body: {body}

Keyphrases (separated by ";"): {keyphrases}

For each keyphrase, generate an abstract category for the keyphrase. Examples include process, task, material, tool, measurement, model, technology, and metric etc. Do not limit yourself to the examples. Make sure that the categories are informative in the domain of science and appearing natural as if that assigned by a well-read user. Return a list of dictionaries, each with two keys - "keyphrase" and "category". If two keyphrases have the same category, make sure that they are labelled with the same phrase. Do not change how the keyphrases appear, including their cases. Return json only and do not say anything else.
```

Figure 8: Prompt used for instructing GPT-4 to generate the goals from a document and keyphrases.

Goal Refinement Then, we instruct GPT-4 to refine the goals by trying to generate more abstract versions of them. The prompt is shown in Figure 9. As we perform the refinement directly from the chat history of the previous step, we omit the previous prompt and step 1 model outputs.

```
... step 1 prompt and model outputs ...

Can you make the categories more abstract, yet still informative to the keyphrase? If the categories are already abstract enough, you do not need to change. Return json only.
```

Figure 9: Prompt used for instructing GPT-4 to improve the abstractiveness of the proposed goals.

For both of the steps, we use greedy search and cap the output to 400 tokens. We parse the results string into json format to extract the goals.

A.2 Human Validation

Next, based on the two rounds of proposed goals, the two authors (student researchers familiar with NLP and the keyphrase generation task) filter out high quality goals as the final benchmarking dataset. We emphasize that this decision is required due to the nature of the task, which requires expert annotators to ensure a high data quality. The consent to use and release the annotation traces was obtained from both of the authors. The type of research conducted by this work is automatically determined exempt from by the authors' institution's ethics review board. We design and enforce two major guidelines during the annotation process:

1. Remove a goal if it is semantically equivalent to or a subtype of some another goal that is more abstract.
2. Remove a goal if it so abstract that it could also enclose other keyphrases not currently paired with the goal. This criterion includes overly vague goals (e.g., "concept") and goals that corresponds to the topic of the entire passage (e.g., "chemistry concepts").

As mentioned in Section 3, this process allows the annotator reach a high inter-annotator agreement of 0.699 Cohen's Kappa. In addition, the annotator actively engage in a discussion whenever ambiguous cases are found. Finally, we conduct a rule-based postprocessing with two stages.

1. **Goal Removal.** We remove the following goals as they represent overly general goals: entity, process, concept.
2. **Goal Unification.** We merge the following goal labels as they represent the same meaning. Table 4 presents the source and target goals. Note that to preserve the diversity of the goals, we refrain from merging aggressively and only merge the basic cases that may be result from annotation discrepancy.

A.3 Negative sampling Algorithm

To construct the training and evaluation data for evaluating the model's ability to reject irrelevant goals, we design a simple algorithm to sample irrelevant goals. Concretely, we pool together all the existing goals from the same dataset as the universal goal set and leverage the phrase embedding model released by (Wu et al., 2023b) to embed all the phrases. Then, for each goal from the docu-

Source Goals	Target Goals
place, geographical location	location
person, people, individual person	individual
geopolitical entity	country
... event	event
profession	occupation
belief system	religion
incident outcome	outcome
subject	topic
incident	event
... equipment	equipment
... procedure	procedure

Table 4: Goal merging directions for METAKP label cleaning. We replace all occurrences of source goals with target goals.

ment, we use it as an anchor to retrieve $d\%$ most dissimilar goals. We use $d = 50$ for all the datasets. From these goals, we sample a goal that is not associated with the document as the irrelevant goal according to the frequency distribution of these goals appearing as relevant goals in the final dataset. We additionally design a frequency match constraint, which enforces that the frequency of a goal g appearing as an irrelevant goal should not exceed the frequency it appears as a relevant goal. In practice, the frequency match constraint is applied first. If no eligible goals remain, we sample a goal from the $d\%$ most dissimilar goals according to frequency.

B Implementation Details

B.1 Supervised Fine-tuning

For multi-task learning with BART and Flan-T5, we base our implementation on the Huggingface Transformers implementations provided by (Wu et al., 2023a) and train for 20 epochs with early stopping. We use learning rate $3e-5$, linear decay, batch size 64, and the AdamW optimizer. Due to the context limitations of Flan-T5, all the input documents for BART and Flan-T5 are truncated to 512 tokens to enable a fair comparison. We perform a careful hyperparameter search over the learning rate, batch size, and warm-up steps. The corresponding search spaces are $\{1e-5, 3e-5, 6e-5, 1e-4\}$, $\{16, 32, 64, 128\}$, and $\{50, 100, 250, 500\}$. The best hyperparameters are chosen based on the performance on the validation set. To decode from the fine-tuned models, we fix the decoder’s prefix using the constrained decoding functionalities provided by Huggingface Transformers and use greedy search

to complete the suffix.

The fine-tuning experiments are performed on a local GPU server with eight Nvidia RTX A6000 GPUs (48G each). We use gradient accumulation to achieve the desired batch sizes. Fine-tuning BART-base, BART-base, Flan-T5-large, and Flan-T5-XL take, respectively.

B.2 Large Language Models

We present the prompts for prompting large language models for goal relevance judgment and goal-conforming keyphrase generation in Figure 10 and Figure 11.

In this task you will need to decide if you should reject the high-level category given the title and abstract of a document. One could use the high-level category to write keyphrases from the document. If you decide the category is relevant to the document, generate yes; if the category is not relevant, generate no. Do not output anything else.

Document Title: {title}
Document Abstract: {body}

High-level Category: {goal}
Relevant? (yes or no):

Figure 10: Prompt used for goal relevance judgment.

Generate present and absent keyphrases belonging to the high-level category from the given text, separated by commas. Do not output anything else.

Document Title: {title}
Document Abstract: {body}

High-level Category: {goal}
Keyphrases (Must be of category "{goal}"):

Figure 11: Prompt used for on-demand keyphrase generation with LLMs.

For all the results reported in the paper, we use gpt-3.5-turbo-0125 and the gpt-4o-2024-05-13 models via the OpenAI API.

For goal relevance judgment, we use greedy decoding and record the yes/no predictions for evaluation. The document body is truncated to the first five sentences as we find providing longer context barely improves the performance.

For on-demand keyphrase generation, the input length is truncated to 4000 tokens. We generate 10 samples with $p = 0.95$ and temperature = 0.7. The output length is limited to 30 tokens, which accommodate approximately 10 keyphrases. Finally,

for filtering, we set a fixed threshold $\tau = 0.3$. We lower-case all the outputs and use a string matching algorithm to remove excessive parts generated by the model such as "present keyphrases: ". The method's sensitivity to the hyperparameter settings is presented in Figure 6 and Figure 12.

Since the proposed LLM-based methods are unsupervised, we refrain from extensively tuning the hyperparameters. The only exception is that we use the validation sets to determine a reasonable good setting of the sample size K and the threshold τ , which is uniformly applied to all the datasets.

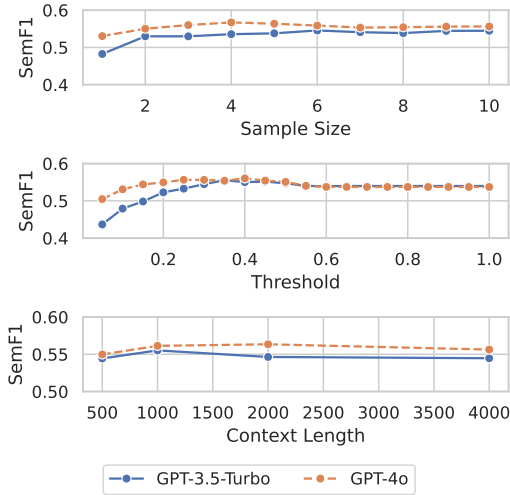


Figure 12: Sensitivity of the self-consistency prompting approach's performance to number of samples, settings of threshold τ , and the input length on KPBiomed.

C Qualitative Study

In Figure 13 and Figure 14, we present and compare the outputs of Flan-T5-XL, zero-shot sampling from GPT-4o, and self-consistency sampling from GPT-4o in two domains. Compared to supervised models, which often generates suboptimal keyphrases under distribution shift, GPT-4o exhibits consistent high recall across domains, and the self-consistency reranking process further filters high quality goals from the zero-shot keyphrase predictions across multiple samples.

Source: KPTimes

Title: more connecticut schools fall short of "no child left behind" standards

Abstract: since [digit] , when president bush signed the federal no child left behind law , schools in connecticut have scrambled to revamp curriculums , step up professional development for teachers and continually assess students ' test scores to comply with the law ' s requirements . it is an effort that has dominated the agendas of school officials not just in connecticut , but all over the nation , and not everyone is happy that test results have become such a focus . when results of the latest test scores were announced a few weeks ago , about [digit] percent , or [digit] , of the public schools in connecticut did not make the grade under the federal law , state officials said . the state added [digit] schools to its list of schools that failed to meet the federal benchmarks .

Goal 1: institution

Reference: "education and schools"

Prediction (Flan-T5-XL): "education and schools"

Prediction (GPT-4o (zero-shot): "connecticut schools", "federal government", "state department of education", "norwalk public schools", "greenwich school district", "greenwich public schools"

Prediction (GPT-4o (self-consistency): "connecticut schools", "state department of education"

Goal 2: law

Reference: "no child left behind act"

Prediction (Flan-T5-XL): "no child left behind act"

Prediction (GPT-4o (zero-shot): "no child left behind", "federal no child left behind law", "federal law", "federal government", "new accountability system", "adequate yearly progress"

Prediction (GPT-4o (self-consistency): "no child left behind", "no child left behind law"

Source: DUC2001 (out-of-distribution)

Title: (empty)

Abstract: millions of gallons of crude oil that spilled when a tanker ran aground spread across a wildlife - rich stretch of ocean saturday , and alaska ' s chief environmental officer criticized cleanup efforts as too slow . the biggest oil spill in u . s . history created a slick about seven miles long and seven miles wide in prince william sound . the coast guard said only reef island and the western edge of bligh island had been touched by the slick . " this situation , i think , was everyone ' s secret nightmare about what could happen with oil traffic in the sound ," said dennis kelso , commissioner of the alaska department of environmental conservation .

Goal 1: substance

Reference: "crude oil"

Prediction (Flan-T5-XL): "oil (petroleum) and gasoline"

Prediction (GPT-4o (zero-shot): "crude oil", "oil spill", "oil pollution", "north slope crude oil", "spilled oil", "leaking oil", "oil slick", "spilled crude oil"

Prediction (GPT-4o (self-consistency): "crude oil"

Goal 2: action

Reference: "cleanup efforts"

Prediction (Flan-T5-XL): "accidents and safety"

Prediction (GPT-4o (zero-shot): "criticized cleanup efforts", "created a slick", "ran hard aground", "halted early", "begin pumping", "removing oil", "placed a boom"

Prediction (GPT-4o (self-consistency): "spread across", "criticized cleanup efforts"

Source: DUC2001 (out-of-distribution)

Title: (empty)

Abstract: the clinton administration will soon announce support for a north american development bank , which would fund projects in communities hit by job losses resulting from the north american free trade agreement . the so - called nadbank has been strongly supported by congressman esteban torres , who has insisted on some sort of lending institution to support adjustment throughout the continent . agreement by the administration is expected to bring mr torres and at least [digit] other hispanic congressmen into the pro - nafta fold . the administration believes it can garner [digit] - [digit] pro - nafta votes , out of the [digit] needed .

Goal 1: economic issue

Reference: "job losses"

Prediction (Flan-T5-XL): "jobs"

Prediction (GPT-4o (zero-shot): "north american development bank", "job losses", "north american free trade agreement", "lending institution", "pro-nafta votes", "anti-nafta public opinion"

Prediction (GPT-4o (self-consistency): "north american development bank", "clinton administration", "job losses"

Goal 2: political entity

Reference: "clinton administration"

Prediction (Flan-T5-XL): "united states politics and government"

Prediction (GPT-4o (zero-shot): "clinton administration", "congressman esteban torres", "hispanic congressmen", "white house", "president bill clinton"

Prediction (GPT-4o (self-consistency): "clinton administration", "north american development bank"

Figure 13: Examples of on-demand keyphrase generation instances and model outputs in the news domain.

Source: KPBiomed

Title: contemporary trend of acute kidney injury incidence and incremental costs among us patients undergoing percutaneous coronary procedures .

Abstract: objectives to assess national trends of acute kidney injury (aki) incidence , incremental costs , risk factors , and readmissions among patients undergoing coronary angiography (cag) and / or percutaneous coronary intervention (pci) during [digit] - [digit] . background aki remains a serious complication for patients undergoing cag / pci . evidence is lacking in contemporary aki trends and its impact on hospital resource utilization . methods patients who underwent cag / pci procedures in [digit] hospitals were identified from premier healthcare database . aki was defined by icd - [digit] / [digit] diagnosis codes ([digit]) .

Goal 1: medical condition

Reference: "acute kidney injury", "chronic kidney disease", "nephropathy"

Prediction (F1an-T5-XL): "acute kidney injury"

Prediction (GPT-4o (zero-shot): "acute kidney injury", "chronic kidney disease", "anemia", "diabetes"

Prediction (GPT-4o (self-consistency): "acute kidney injury", "chronic kidney disease", "anemia"

Goal 2: medical procedure

Reference: "percutaneous coronary intervention"

Prediction (F1an-T5-XL): "percutaneous coronary intervention"

Prediction (GPT-4o (zero-shot): "percutaneous coronary intervention", "coronary angiography", "coronary procedures", "inpatient procedures", "outpatient procedure"

Prediction (GPT-4o (self-consistency): "percutaneous coronary intervention", "coronary angiography"

Source: PubMed (out-of-distribution)

Title: surviving sepsis campaign : international guidelines for management of severe sepsis and septic shock : [digit]

Abstract: objective to provide an update to the original surviving sepsis campaign clinical management guidelines , 201c surviving sepsis campaign guidelines for management of severe sepsis and septic shock ,201d published in [digit] . introduction severe sepsis (acute organ dysfunction secondary to infection) and septic shock (severe sepsis plus hypotension not reversed with fluid resuscitation) are major healthcare problems , affecting millions of individuals around the world each year , killing one in four (and often more) , and increasing in incidence [[digit] 2013 [digit]] . similar to polytrauma , acute myocardial infarction , or stroke , the speed and appropriateness of therapy administered in the initial hours after severe sepsis develops are likely to influence outcome .

Goal 1: medical condition

Reference: "sepsis", "severe sepsis", "septic shock", "sepsis syndrome", "infection"

Prediction (F1an-T5-XL): "sepsis"

Prediction (GPT-4o (zero-shot): "acute kidney injury", "chronic kidney disease", "anemia", "diabetes"

Prediction (GPT-4o (self-consistency): "severe sepsis", "septic shock"

Goal 2: healthcare initiative

Reference: "surviving sepsis campaign"

Prediction (F1an-T5-XL): "surviving sepsis campaign"

Prediction (GPT-4o (zero-shot): "surviving sepsis campaign", "international guidelines", "management of severe sepsis", "septic shock", "clinical management guidelines", "evidence-based methodology"

Prediction (GPT-4o (self-consistency): "surviving sepsis campaign"

Source: PubMed (out-of-distribution)

Title: keratinocyte serum - free medium maintains long - term liver gene expression and function in cultured rat hepatocytes by preventing the loss of liver - enriched transcription factors

Abstract: freshly isolated hepatocytes rapidly lose their differentiated properties when placed in culture . therefore , production of a simple culture system for maintaining the phenotype of hepatocytes in culture would greatly facilitate their study . our aim was to identify conditions that could maintain the differentiated properties of hepatocytes for up to [digit] days of culture . adult rat hepatocytes were isolated and attached in williams 2019 medium e containing [digit] % serum . the medium was changed to either fresh williams 2019 medium e or keratinocyte serum - free medium supplemented with dexamethasone , epidermal growth factor and pituitary gland extract .

Goal 1: biological extract

Reference: "pituitary gland extract"

Prediction (F1an-T5-XL): "pituitary gland extract"

Prediction (GPT-4o (zero-shot): "keratinocyte serum-free medium", "Williams2019 medium E", "dexamethasone", "epidermal growth factor", "pituitary gland extract"

Prediction (GPT-4o (self-consistency): "keratinocyte serum-free medium", "keratinocyte serum"

Goal 2: molecular biology technique

Reference: "reverse transcription polymerase chain reaction"

Prediction (F1an-T5-XL): "cell culture"

Prediction (GPT-4o (zero-shot): "immunohistochemistry", "western blotting", "rt-pcr", "immunofluorescence staining", "collagenase perfusion technique"

Prediction (GPT-4o (self-consistency): "western blotting", "rt-pcr", "immunohistochemistry"

Figure 14: Examples of on-demand keyphrase generation instances and model outputs in the biomedical domain.