

AliCG: Fine-grained and Evolvable Conceptual Graph Construction for Semantic Search at Alibaba

Ningyu Zhang^{1,2*}, Qianghuai Jia^{4*}, Shumin Deng^{1,2*}, Xiang Chen^{1,2}, Hongbin Ye^{1,2}, Hui Chen³, Huaixiao Tou³, Gang Huang⁵, Zhao Wang¹, Nengwei Hua³, Huajun Chen^{1,2*}
¹Zhejiang University, China & AZFT Joint Lab for Knowledge Engine, China ³Alibaba Group, China
²Hangzhou Innovation Center, China, Zhejiang University, China, ⁴AntGroup, China ⁵Zhejiang Lab, China
{zhangningyu,231sm,xiang_chen,yehongbin,huanggang,zhao_wang,huajunsir}@zju.edu.cn,
{qianghuai.jqh,weidu.ch,huaixiao.thx,nengwei.huanw}@alibaba-inc.com

ABSTRACT

Conceptual graphs, which is a particular type of Knowledge Graphs, play an essential role in semantic search. Prior conceptual graph construction approaches typically extract high-frequent, coarse-grained, and time-invariant concepts from formal texts. In real applications, however, it is necessary to extract less-frequent, fine-grained, and time-varying conceptual knowledge and build taxonomy in an evolving manner. In this paper, we introduce an approach to implementing and deploying the conceptual graph at Alibaba. Specifically, We propose a framework called **AliCG** which is capable of a) extracting fine-grained concepts by a novel bootstrapping with alignment consensus approach, b) mining long-tail concepts with a novel low-resource phrase mining approach, c) updating the graph dynamically via a concept distribution estimation method based on implicit and explicit user behaviors. We have deployed the framework at Alibaba UC Browser. Extensive offline evaluation as well as online A/B testing demonstrate the efficacy of our approach.

CCS CONCEPTS

• **Information systems** → **Query representation; Information extraction.**

KEYWORDS

Concept Mining; Taxonomy Construction; Knowledge Graph

ACM Reference Format:

Ningyu Zhang^{1,2*}, Qianghuai Jia^{4*}, Shumin Deng^{1,2*}, Xiang Chen^{1,2}, Hongbin Ye^{1,2}, Hui Chen³, Huaixiao Tou³, Gang Huang⁵, Zhao Wang¹, Nengwei Hua³, Huajun Chen^{1,2*}. 2021. AliCG: Fine-grained and Evolvable Conceptual Graph Construction for Semantic Search at Alibaba. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '21)*, August 14–18, 2021, Virtual Event, Singapore. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3447548.3467057>

* Equal contribution and shared co-first authorship.

★ Corresponding author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](https://permissions.acm.org).

KDD '21, August 14–18, 2021, Virtual Event, Singapore

© 2021 Association for Computing Machinery.

ACM ISBN 978-1-4503-8332-5/21/08...\$15.00

<https://doi.org/10.1145/3447548.3467057>

1 INTRODUCTION

Knowledge is important for text-related applications such as semantic search. Knowledge Graphs (KGs) organize facts in a structured graph way as triples in the form of $\langle \text{subject}, \text{predicate}, \text{object} \rangle$, abridged as (s, p, o) , where s and o denote entities and p builds relations between entities. **Conceptual Graph**, which is a special type of KGs, builds the semantic connections between concepts and has proven to be valuable in *short text understanding* [21], *Word Sense Disambiguation* [20], *enhanced entity linking* [5], *semantic query rewriting* [22], etc. Essentially, conceptualization helps humans generalize previously gained knowledge and experience to new settings, which may reveal paths to high-level cognitive **System 2** [2] in a conscious manner. In real-life applications, the conceptual graph provides valuable knowledge to support many applications [25], such as semantic search. Web search engines (e.g., Google and Bing) leverage a taxonomy to better understand user queries and improve the search quality. Moreover, many online retailers (e.g., Alibaba and Amazon) organize products into categories of different granularities so that customers can easily search and navigate this category taxonomy to find the items they want to purchase.

In this paper, we introduce the Alibaba Conceptual Graph (AliCG), which is a large-scale conceptual graph of more than 5,000,000 fine-grained concepts, still in fast growth, automatically extracted from noisy search logs. As shown in Figure 1, AliCG comprises four levels: **level-1** consists of concepts expressing the domain that those instances belong to; **level-2** consists of concepts referred to the type or subclass of instances; **level-3** consists of concepts that are the fine-grained conceptualization of instances expressing the implicit user intentions; **instance layer** includes all instances such as entities and non-entity phrases. AliCG is currently deployed at Alibaba to support a variety of business scenarios, including the product Alibaba UC Browser. AliCG has been applied to more than dozens of applications in Alibaba UC Browser, including intent classification, named entity recognition, query rewriting, and so on, and it receives more than 300 million requests per day.

Building AliCG is not a trivial task. Previous studies such as YAGO [18] and DBPedia [1] have investigated the extraction of knowledge from formal texts (e.g., Wikipedia). Probase [23] proposes an approach for extracting concepts from semi-structured Web documents. However, these approaches could not be adapted to our applications because several challenges remain unresolved.

Fine-grained Concept Acquisition. Conventional approaches devoted to extracting coarse-grained concepts such as categories or types. However, in Alibaba's scenario of question fine-grained

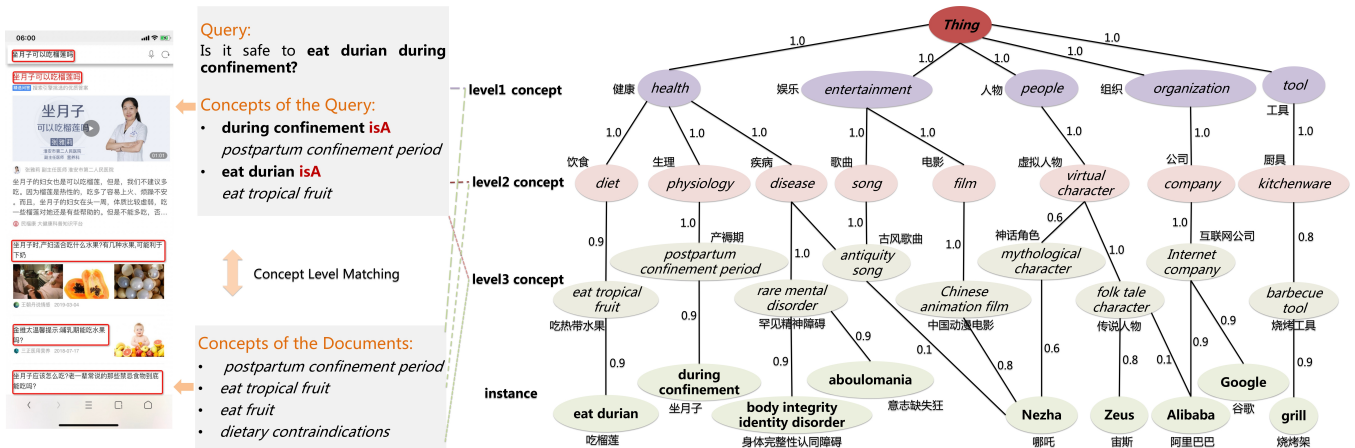


Figure 1: Data hierarchy of Alibaba Conceptual Graph (AliCG) for semantic search.

concepts are necessary to increase the recall of answer results. For example, “grill (烤架)” is a “tool (工具)” and “scarf (围巾)” is a “clothes (服饰)”. However, it would be more helpful if we can infer that a user searching for these items may be more interested in “barbecue tool (烧烤工具)” or “keep warm clothes (保暖服饰)” rather than another “tool (工具)” like “wrench (扳手)”—these concepts are rare in existing conceptual graphs.

Long-tail Concept Mining. Conventional approaches [23] generally extract concepts based on Hearst patterns, e.g., “especially” and “such as.” However, these approaches cannot extract long-tail concepts from extremely short or noisy queries, which are common in search engines. For instance, it is non-trivial to extract the concept “rare mental disorder (罕见精神障碍)” of the instance “body integrity identity disorder (身体完整性认同障碍)” from the search log as only 35 instances mentioned “rare mental disorder”. It is rather difficult to extract such concepts from the short text with pattern matching (the pattern is too general [6–8, 26, 27], and there is little context information as well as co-occurrence samples), as Figure 2 shows. Besides, there exist lots of scattered concepts in user search engines such as “traditional activities Tibetan New Year (藏历新年习俗)”. Recent approaches usually regard such concept extraction procedure as sequence labeling tasks [13], which rely on a tremendous amount of training data for each concept. Nearly 79% of the concepts are long-tail in the search logs. Therefore, it is crucial to be able to extract concepts with limited numbers of instances.

Taxonomy Evolution. Numerous instances and concepts in user search queries are related to recent trending and evolving events. Conventional approaches are not able to update the taxonomies over time. For instance, a user may search “Nezha (哪吒)” or “New animations in February (二月新番)” in search engines. The implied meaning of such concepts changes over time because apparently, “Nezha” has different meanings at different times (e.g., Chinese animation film, mythological character, the hero in Honor of Kings), and there are various new animations in February in different years. Thus, it is imperative to incorporate the temporal evolution into the taxonomies. However, as we extract instances and concepts from the text, which inevitably brings about duplicate

edges, it is necessary to align those nodes with the same meaning in the conceptual graph. Besides, as there are many multiple-to-multiple nodes in the conceptual graph, it is prohibitive to update such complex graphs over time. In other words, it is difficult to estimate the confidence distribution of the concepts given instances.

To address challenges mentioned above, we propose the following contributions in the design of AliCG:

First, we propose a novel **bootstrapping with the alignment consensus** approach to tackle the first challenge of extracting fine-grained concepts from noisy search logs. Specifically, we utilize a small number of predefined string patterns to extract concepts, which are then used to expand the pool of such patterns. Further, the new mined concepts are verified with query-title alignment; that is, an essential concept in a query should repeat several times in the document title frequently clicked by the user. *Second*, we introduce a novel **conceptualized phrase mining and self-training with an ensemble consensus** approach to extract long-tail concepts. On the one hand, we extend the off-the-shelf phrase mining algorithm with conceptualized features to mine concepts unsupervisedly. On the other hand, we propose a novel low-resource sequence tagging framework, namely, self-training with an ensemble consensus, to extract those scattered concepts. *Finally*, we propose a novel **concept distribution estimation method based on implicit and explicit user behaviors** to tackle the taxonomy evolution challenges. We employ concept alignment and take advantage of user’s searching and clicking behaviors to estimate the implicit and explicit concept distributions to construct a four-layered concept–instance taxonomy in an evolving manner.

To deploy AliCG in real-life applications, we further introduce three methods for utilizing the conceptual knowledge, including, **text rewriting**, **concept embedding**, and **conceptualized pretraining**. We conduct extensive offline evaluations, including concept mining and applications such as intent classification and named entity recognition. Experimental results show that our approach can extract more fine-grained concepts in both normal and long-tail setting compared with baselines. Moreover, the performance of intent classification and named entity recognition is significantly

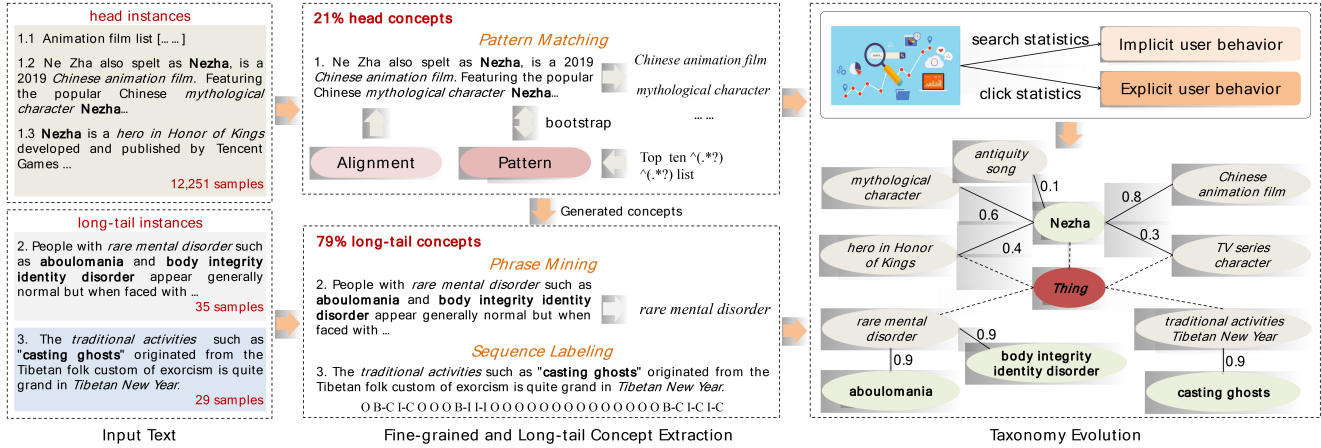


Figure 2: Framework of Alibaba Conceptual Graph Construction.

improved when integrated with the conceptual graph. We also performed online A/B testing on more than 400 million actual users of the Alibaba UC Browser mobile application¹. The experimental results show that the relevant score increases by 12%, the click-through rate (CTR) increase by 5.1%, and page view (PV) increases by 5.5%. The highlights of our work are the following:

- We introduce the AliCG, which is a large-scale conceptual graph of more than 5,000,000 fine-grained concepts automatically extracted from the noisy search logs².
- We propose three novel approaches to address the issues of fine-grained concept acquisition, long-tail concept mining, and taxonomy evolution. Our framework is able to extract and dynamically update concept taxonomy in both normal and long-tail settings.
- We introduce three methods to deploy our approach in Alibaba UC Browser, and both offline evaluation and online A/B testing demonstrate the efficacy of our approach. We also release an open dataset³ of AliCG with 490,000 instances for research purposes.

2 METHODOLOGY

Our approach is aimed at constructing conceptual graph from both the Web documents and the noisy query logs. We denote a user query by $q = w_1^q w_2^q \cdots w_{|q|}^q$ and the set of all queries by Q . In addition, we denote a document by $d = w_1^d w_2^d \cdots w_{|d|}^d$. Given a user query q , and the top-ranked clicked documents, $D^q = \{d_1^q, d_2^q, \dots, d_{|D^q|}^q\}$, we aim to extract an instance/concept phrase, $c = w_1^c w_2^c \cdots w_{|c|}^c$. As illustrated in Figure 2, the construction of the conceptual graph generally consists of three modules, as follows:

Fine-grained Concept Acquisition aims at extracting those common fine-grained concepts from the noisy search logs with

bootstrapping with the alignment consensus. After that, candidate concepts and instances are acquired, and we link those instances with corresponding concepts via probabilistic inference and concept matching following the approach proposed in [13].

Long-tail Concept Mining aims at extracting those long-tail concepts from the noisy search logs with **conceptualized phrase mining and self-training with an ensemble consensus**. In addition, we also leverage concepts from the fine-grained concept acquisition module as distantly supervised samples for concept mining (see details in section 2.2).

Taxonomy Evolution aims at building the taxonomy in an evolving way with **concept distribution estimation method based on implicit and explicit user behaviors**. Note that taxonomy evolution is conducted with those concepts mined from the modules of fine-grained concept acquisition and long-tail concept mining.

2.1 Fine-grained Concept Acquisition

Previous studies [23] show that the iterative (bootstrapping) approaches can extract coarse isA facts with the highest confidence starting with a set of seed patterns; thus, it is intuitive to construct elegant patterns to obtain fine-grained concepts. Specifically, we define a small set of patterns to extract concept phrases from queries with high confidence. For example, “Top 10 XXX (十大XXX)” is a pattern that can be used to extract seed concepts. Based on this pattern, we can extract concepts such as “Top 10 mobile games (十大手游).” However, there still exist lots of noisy texts, which are challenging for vanilla bootstrapping approaches. To this end, a novel bootstrapping with the alignment consensus approach is proposed to deal with noisy texts. The intuition behind this is that *we must control the pattern generalization and concept consistency given query-title pairs*.

Bootstrapping with Alignment Consensus. First, an extracted pattern should not be too general for the sake of extracting *fine-grained concepts*. Therefore, given a new pattern p found in a certain round, let n_s be the number of concepts in the existing seed concept set that can be extracted by p from the query set Q . Let n_e be the number of new concepts that can be extracted by p from Q . We keep pattern p via the function $Filter(p)$: $1) \alpha < \frac{n_s}{n_e} < \beta$, and $2)$

¹<https://www.ucweb.com/>

²Demo available in <http://openconcepts.zjukg.cn/>

³<https://github.com/alibaba-research/ConceptGraph>

$n_s > \delta$, where α , β , and δ are predefined thresholds to control the fineness of extracted concepts. Second, we filter the mined concepts via query-title alignment to improve the quality of fine-grained concepts. Even though bootstrapping helps discover new patterns and concepts from the query set Q in an iterative manner, such a pattern-based method has limited extraction ability and may introduce considerable noise. Based on [13], we further propose the extraction of concepts from a query and its top clicked link titles in the search log as the essential concept in a query should repeat several times in the document title frequently clicked by the user. The overall algorithm is shown below:

Algorithm 1 Bootstrapping with the alignment consensus

Require: query and high-clicked documents set $(q, D) \in T$.
Require: Predefined templates set $p \in P$, Concepts $C^{head} = \Phi$

- 1: **for** *iter* iterations **do**
- 2: **for** each $p \in R$ **do**
- 3: **for** each $(q, D) \in T$ **do**
- 4: **if** q match r **then**
- 5: $c^p = \text{ExtractByTemplate}(q)$
- 6: $c^a = \text{ExtractByAlignment}(q, D)$
- 7: **if** $\text{len}(c^p) \leq \text{len}(c^a)$ **then**
- 8: $C^{head} \leftarrow C^{head} + c^p$
- 9: $p^{candidate} = \text{GenPattern}(C^{head}, q)$
- 10: $p^{new} = \text{Filter}(p^{candidate})$
- 11: $P \leftarrow P + p^{new}$

return C^{head}

Note that we have mined typical fine-grained instance and concept candidates, but we should link those instances with corresponding concepts. Firstly, we utilize a concept discriminator to determine whether each candidate is a concept or instance. We represent each candidate by a variety of features, such as whether this concept has ever appeared as a query, how many times it has been searched, etc. We then train a classifier with Gradient Boosting Decision Tree and link instances⁴ with those classified concepts by probabilistic inference and concept matching following [13].

2.2 Long-tail Concept Mining

Although iterative pattern matching can extract lots of high-frequent concepts, it is still non-trivial to extract those with only a few instances. As Figure 3 shows, it is challenging to extract long-tail concepts because of two main reasons: *poor pattern generalization* and *few co-occurrence samples*. Thus, it is difficult to match such low-shot concepts and link those concepts with corresponding instances. For the long-tail problem, we first propose *conceptualized phrase mining* to utilize the external domain knowledge graph to generate weak-supervised samples for unsupervised learning of those long-tail concepts, which address those issues. We then propose a self-training-based low-resource tagging algorithm for supervised learning of long-tail concepts, which can further extract scattering concepts.

Conceptualized Phrase Mining. It is difficult to utilize rule-based approaches to mine long-tail instances and concepts. We

⁴We also mined instances from the search logs based on keyword extraction to increase the coverage of instance.

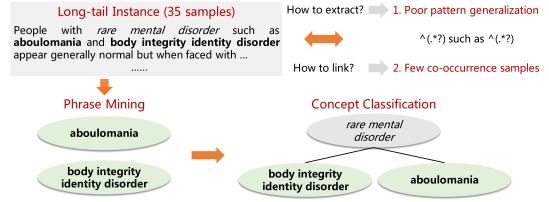


Figure 3: Long-tail Concept Mining.

propose the use of phrase mining methods to extract such instances and concepts. The motivation is that *there are numerous none-entity phrases in search queries, and such phrases play a vital role in understanding the query*. Specifically, we firstly filter stop words and then employ an off-the-shelf phrase mining tool, AutoPhrase⁵, to perform phrase mining on the corpora unsupervisedly. However, the original phrase mining approach experiences several issues in the real data set. First, the results of the shallow syntax analysis of part-of-speech (POS) are not applicable to all areas, such as the medical domain. Second, there is considerable wrong segmentation of words for Chinese. To address those issues, we use the existing domain knowledge graph to generate weak-supervised samples to re-train the POS and use the results of the POS tagging as well as Chinese BERT-base embeddings as segmentation features for AutoPhrase. As AutoPhrase relies on distance supervised data, we leverage the data of the existing domain knowledge graph along with domain rules to generate positive and negative data. After running the AutoPhrase, we also propose a new empirical score function with length constraint to generate high-quality phrases:

$$f(p) = \frac{1}{n} \sum_{i=1}^n p_{score} + \log(p_{len}) \quad (1)$$

Where p_{score} and p_{len} refer to the score and length of separated phrase p , respectively. Based on the mined phrases, we train a concept classification model based on Fasttext⁶ for specific concepts (200 concepts in the vertical domain). Based on the classification model, we link concept instances with specific concept labels (e.g., *body integrity identity disorder* is a *rare mental disorder*). We generate positive training samples of the concept classifier from those instances existing in the conceptual graph and create negative samples with domain rules (e.g., add negative character).

Self-Training with an Ensemble Consensus. Those unsupervised approaches are still limited in generalization. We further perform supervised learning. However, there is a lack of sufficient training samples. Thus, we propose a novel low-resource tagging algorithm, namely self-training, with an ensemble consensus. We leverage a few instances/concepts generated from the unsupervised approaches described in the section 2.1 and 2.2 as seeds (there may be only a few samples for some concepts in unique domains such as medical domain). Then, we train a CRF model⁷ with seeds. We utilize the CRF model to generate a large amount of weak-supervised pseudo-sample from a subset of unlabelled data. Then, we train a BERT tagging model with permutations based on that

⁵<https://github.com/shangjingbo1226/AutoPhrase>

⁶<https://fasttext.cc/>

⁷<https://taku910.github.io/crfpp/>

pseudo-sample. Thereafter, we generate pseudo-samples once more. Different from the first stage, we combine the CRF, BERT tagging, and a domain dictionary model (maximum forward match) with an ensemble consensus to generate more confident training samples. In other words, the prediction of three different models agrees with each other; thus, such pseudo-label may be more confident. Next, we train a BERT tagging model once more. After several iterations, we can finally obtain a high-performance sequence tagging model with only a few initial training seeds.

Algorithm 2 Self-training with an ensemble consensus

Require: Only few initial training samples L , unlabeled data U , domain dictionary D , randomly initialized parameter of BERT tagging model θ , perturbation δ

- 1: Train the CRF with L , build dictionary tagging model with D
 - 2: Random sample subset \mathcal{U} from U and generate pseudo-sample \mathcal{U}^l with CRF
 - 3: **for** $iter$ iterations **do**
 - 4: **for** minibatch index $b = 1 : \frac{|\mathcal{U}^l|}{batchsize}$ **do**
 - 5: Sample a subset \mathcal{B}'_b with batch size b from \mathcal{U}^l
 - 6: $\theta \leftarrow \theta + \delta$
 - 7: $\theta \leftarrow \theta - \alpha \nabla_{\theta} \frac{1}{|\mathcal{B}'_b|} \sum_{(x,y) \in \mathcal{B}'_b} CE(f_{\theta}(x), y)$
 - 8: Random sample subset \mathcal{U} from U
 - 9: Generate pseudo-sample $\mathcal{U}^{ensemble}$ with voting of dictionary tagging model, CRF model, and BERT tagging model
 - 10: $\mathcal{U}^l \leftarrow \mathcal{U}^l + \mathcal{U}^{ensemble}$
- return** BERT sequence tagging model θ
-

The above approaches for concept mining are complementary to each other. Our experience shows that bootstrapping with alignment consensus can extract head concepts with high accuracy. We also observe that for those long-tail concepts which occupies a large proportion of the search log, conceptualized phrase mining is advantageous for extracting those low-frequent terms, while sequence tagging can better extract concepts from the short text when they have a clear boundary with surrounding non-concept words.

2.3 Taxonomy Evolution

Existing taxonomies are mostly constructed by human experts or in a crowdsourcing manner. As the web content and human knowledge are constantly growing over time, people need to update existing taxonomy and also include new emerging concepts, and it is intuitive to incorporate the temporal evolution into the taxonomies. The key to handling the problem of taxonomy evolution is to properly model the probability of assigning a concept to a parent concept in the taxonomy tree, i.e., estimating the concept distribution, which can evolve.

Concept distribution estimation based on implicit and explicit user behavior. Different from the previous approach, such as Probase [23], we leverage the user behaviors to estimate the confidence score of concepts given instances. As the user search and click are important behaviors that reveal the user interests in search engines, it is intuitive to estimate confidence score with statistics of implicit user behaviors such as searching and explicit

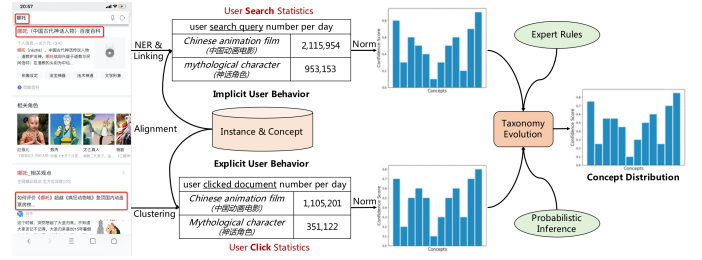


Figure 4: Example of taxonomy evolution.

user behaviors such as clicking on the documents. As Figure 4 shows, we first align parts of concepts with a predefined synonym dictionary as part of the concepts and instance are from different sources, which result in redundancy. Then, we calculate the confidence scores via the heat of entities. We estimate the heat of entities via user searching per day, and the details are included in the appendix. Specifically, given instances with concepts $c \in C$, we run named entity recognition for each c (each concept only remains one most important entity, e.g., largest TF-IDF score)? And obtain a heat score. Then we calculate the implicit concept distribution based on user search behaviors as follows:

$$s_{implicit} = \frac{c_{heat}}{\sum_{c \in C} c_{heat}} \quad (2)$$

where c_{heat} means the heat of the entity in concept c and C is the concept set that c belongs to. Then, we leverage user clicking behaviors to obtain concept distributions explicitly. Specifically, we first collect one month's instances (queries) clicked tag pairs and then aggregate the of queries with the same clicked concept tag⁸. In other words, we estimate the number of queries for each concept tag in the last month. For example, if users clicked more "Chinese animation film (中国动画电影)" when searching "Nezha (哪吒)", then the instance "Nezha (哪吒)" will be more confident with concept "Chinese animation film (中国动画电影)". Thus, we can estimate the concept of confidence distribution based on user clicking behavior. Finally, we combine two different granularity confidence scores as follows:

$$s = \begin{cases} \sum_i^N (s_{implicit} + s_{explicit}) & \text{if } s \notin M \\ \sum_i^N (1.5 \log(s_{implicit}) + \log(s_{explicit})) & \text{if } s \in M \end{cases} \quad (3)$$

Where $s_{implicit}$ is the confidence score based on user searching behavior, $s_{explicit}$ is the confidence score based on the user clicking behavior, and M is the specific subset (e.g., medical and education domain) of the concept set. The specific subset is required here because some domains have special user behaviors, for example, no clicking or lots of misleading clicking. We also remove concepts such as "stop words" and "common words" to prevent unnecessary noise. Meanwhile, we re-weight domain words and specific types of concepts that are unambiguous by $s = 19.5 \log(c_{heat})s$, e.g., "Pauli incompatible" is a "chemical term" and "Andy Lau" is a "person."

After that, we can build the instance-concept taxonomy in an evolving way. Then we leverage expert rules to define the taxonomy between **level1** and **level2**. We further determine the relation

⁸The search logs contain the history of user clicked concept tags.

between **level2** and **level3** via probabilistic inference. Specifically, suppose there is a concept c , a higher level concept p , n_c instance related to concept c , and n_p^c instance belonging to concept p . Then, we estimate $p(p|c)$ by $p(p|c) = n_p^c/n^c$. We identify the isA relation between c and p if $p(p|c) > \delta_t$.

3 DEPLOYMENT

We propose three methods for deploying the conceptual graph at Alibaba, namely **text rewriting**, **concept embedding**, and **conceptualized pretraining**.

Text Rewriting. For each text instance s , we extract the concept c conveyed in the text and rewrite the text by concatenating each instance with s . The rewritten text format is $s.c$. We use text rewriting for information retrieval. Note that text rewriting is easy to apply to other classification or sequence labeling tasks.

Concept Embedding. Following [4], we utilize a two-tower neural network with concept attention and self-attention for learning concept embedding. Then we concatenate the concept embedding with text embedding for sub-tasks.

Conceptualized Pretraining. As pretraining is quite powerful, the conceptual graph can be utilized during the pretraining stage to inject knowledge explicitly. Specifically, we utilize instances and concept masking strategies with an auxiliary concept prediction loss to integrate conceptual knowledge.

4 EVALUATION

In this section, we first introduce a new dataset for the problem of concept mining from the search logs and compare the proposed approach with various baseline methods. Then, we evaluate the accuracy of the taxonomy evolution. We also evaluate different methods to apply conceptual graphs in applications, including intent classification and named entity recognition. Finally, we perform large-scale online A/B testing to show that our approach significantly improves the performance of the semantic search.

4.1 AliCG Construction Evaluation

4.1.1 Evaluation of Fine-grained Concept Acquisition.

We evaluate our approach on two datasets. The first is the user-centered concept mining (UCCM⁹) dataset [13], which is sampled from the queries and query logs of Tencent QQ Browser. The second is the Alibaba Conceptual Graph (ACG) dataset¹⁰. We have created a large-scale dataset containing 490,000 instances with the fine-grained concept, which is sampled from the search logs of Alibaba UC Browser from November 11, 2018, to July 1, 2019.

We compare our approach with the following baseline methods, and the variants of our approach: **TextRank** [17] is a classical graph-based ranking model for keyword extraction; **THUCKE** [14] is a method that considers keyphrase extraction as a problem of translation and learns translation probabilities between the words in the input text and the words in keyphrases; **AutoPhrase** [19] is a quality phrase mining algorithm that extracts quality phrases based on a knowledge base and POS-guided segmentation; **ConceptT** [13] is a state-of-the-art concept mining approach that utilizes bootstrapping, query-title alignment, and sequence tagging methods to

⁹<https://github.com/BangLiu/ConceptT>

¹⁰<https://github.com/alibaba-research/ConceptGraph>

Table 1: Evaluation results of fine-grained Acquisition.

Dataset	UCCM		ACG	
Method	Exact Match	F1 Score	Exact Match	F1 Score
TextRank	0.1941	0.7356	0.0025	0.1839
THUCKE	0.1909	0.7107	0.0325	0.2839
AutoPhrase	0.0725	0.4839	0.0325	0.2839
ConceptT	0.8121	0.9541	0.7725	0.7839
AliCG	0.9122	0.9812	0.9215	0.9898
w/o BA	0.8785	0.9635	0.8852	0.9543

Table 2: Concept mining examples.

Instance	People with <i>rare mental disorder</i> such as aboulomania and body integrity identity disorder appear generally normal but when faced with ...
Model	Instance-Concept
AutoPhrase	No results
ConceptT	(<i>aboulomania isA rare mental disorder</i>)
AliCG	(<i>body integrity identity disorder isA rare mental disorder</i>) (<i>aboulomania isA rare mental disorder</i>)

Table 3: Evaluation results of long-tail concept mining on ACG (long-tail) dataset.

Method	Exact Match	F1 Score
TextRank	0.0941	0.1356
THUCKE	0.1209	0.2107
AutoPhrase	0.1725	0.2839
ConceptT	0.3121	0.2541
AliCG	0.6122	0.8812
w/o CP	0.5522	0.7612
w/o SE	0.5319	0.7322

exact concepts from search queries; **w/o BA**) is the approach that extracts concepts from search logs without **Bootstrapping with the Alignment consensus**.

From Table 1, we observe that our method achieves the best *Exact Match (EM)* and *F1 scores*. This is because that guiding and aligning consensus help us build a collection of high-quality instances and concepts in an unsupervised manner. The TextRank, THUCKE, and AutoPhrase methods do not provide satisfactory performance because they are more suitable for extracting keywords or phrases from long documents or corpora. Our method exhibits better performance compared with ConceptT mainly because our method can extract more accurate concepts with consistency in query-title pairs.

4.1.2 Evaluation of Long-tail Concept Mining.

As there are no public datasets available for long-tail concept mining, we leverage our ACG dataset and randomly choose 100 long-tail concepts with corresponding texts to build the dataset ACG (long-tail). We evaluate our long-tail concept mining methods on

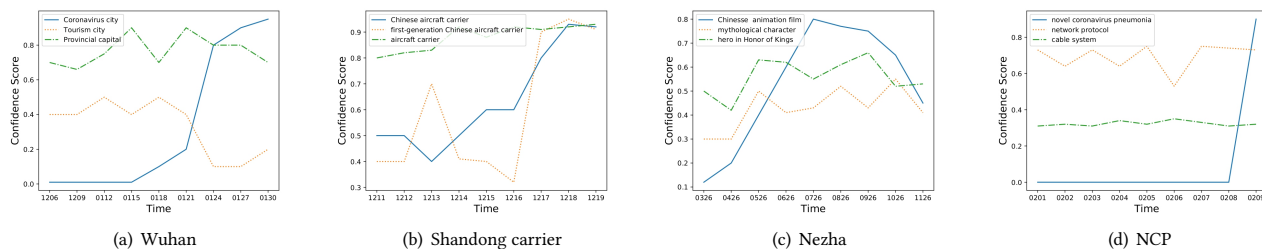


Figure 5: Concept evolution visualization.

that dataset compared with the same baselines in fine-grained concept mining and the variants of our approach. **w/o CP** is the approach that extracts concepts without *Conceptualized Phrase mining*; **w/o SE** is the approach that extracts concepts without *Self-training and Ensemble consensus*.

From Table 3, we observe that our method achieves the best EM and F1 scores. Firstly, the conceptualized phrase mining method extracts numerous long-tail instances and concepts. Secondly, self-training sequence tags with an ensemble consensus can identify scattered concepts from the search logs; that is, the conceptual boundaries in the text are not clear. The TextRank, THUCKE, and AutoPhrase methods fail to perform well in ACG (long-tail), which indicates the difficulty of long-tail concept extraction. Our method exhibits significant improvement compared with ConceptT mainly because our method can extract long-tail concepts and achieve excellent performance in the sequence labeling with only a small number of samples. The comparison shows that our method outperforms its variants and proves the effectiveness of combining different strategies in our system. We also show examples of concept mining in Table 2, which illustrates that our approach is more capable of extracting fine-grained concepts.

4.1.3 Evaluation of Taxonomy Evolution.

We randomly extract 1,000 instances from the taxonomy. The experiments focus on evaluating the relation between concept-instance pairs because these relations are critical for semantic search. We check whether the isA relation between each concept and its instance is correct. We ask three human judges to evaluate the relations. We record the number of correct and incorrect instances for each concept. Table 4 shows the results of the evaluation. The average number of instances per concept is 5.44, and the most significant concept contains 599 instances. Note that the size of the taxonomy is increasing and updated with daily user query logs. The accuracy of the isA relation between concept-instance pairs is 98.59%. We also notice that implicit and explicit user behaviors (**w/o implicit** or **w/o explicit**) have a significant loss in performance, which indicates that user behavior plays an important role in concept distribution estimation.

We visualize the concept trending over time to evaluate the dynamic of AliCG. From Figure 5 we observe that: 1) our approach can obtain a dynamic and fine-grained concept distribution over time. Figure 5(a) shows that the concept of "Wuhan (武汉)" is more confident with "coronavirus city (冠状病毒城市)" after January 22th, while it is confident with "tourism city (旅游城市)" before

Table 4: Evaluation results of taxonomy evolution.

Metrics / Statistics	Value
Mean #Instances per Concept	5.44
Max #Instances per Concept	599
isA Relationship Accuracy	98.59
w/o implicit	95.23
w/o explicit	93.78
w/o both	89.12

January 22th. Note that there was a severe coronavirus outbreak in Wuhan, and lots of relevant news published on the web after January 22th. 2) Our approach can also obtain false confident scores due to the noise in user behaviors. Figure 5(b) indicates that "Shandong carrier (山东号)" has a low confidence score in the day before December 17. Such circumstance is caused by some fake or clickbait news which misleads user behaviors, and it results in the wrong confidence score of "first-generation Chinese aircraft carrier (中国第一代航母)." However, our approach can self-repair and get the correct confident score after that day. 3) Our approach can reveal the concept distribution trend over a long period of time. From Figure 5(c), we observe the confidence score of "Nezha (哪吒)" is becoming higher with "Chinese animation film (中国动漫电影)," but it started to slightly decreasing after July, which reveals the attenuation of user interest over time. 4) Our approach can extract emerging concepts. We observe from 5(d) that our approach extract the emerging concept "novel coronavirus pneumonia (新冠肺炎)" on February 8. Before that, the concept of "NCP" is mostly confident with "network protocol (网络协议)" or "cable system (电话系统)."

4.2 AliCG Deployment Evaluation

AliCG is currently deployed at Alibaba to support a variety of business scenarios, including the product Alibaba UC Browser. The current system can extract approximately 20,000 concepts every day and serve more than 300 million daily active users. AliCG has been applied to dozens of applications, such as intent classification, named entity recognition, information retrieval, entity recommendation, and so on. We will introduce the evaluation of those scenarios to demonstrate the efficacy of AliCG.

4.2.1 Intent Classification and Named Entity Recognition.

We apply AliCG to two tasks to evaluate the different deployment

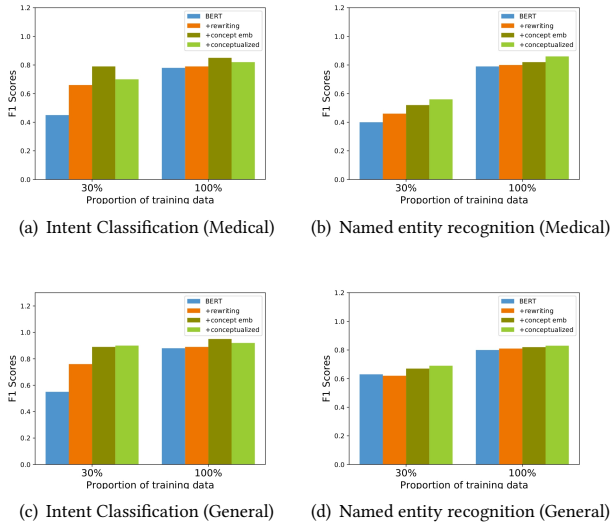


Figure 6: Evaluation results of applications.

methods in both specific domains (medical domain¹¹) and general domain (search logs). **+rewriting** is the method with text rewriting; **+concept emb** is the method with concept embedding; **+conceptualized** is the method with conceptualized pretraining. From Figure 6, we observe: 1) All knowledge-enhanced approaches achieve better performance compared with the baseline BERT, which demonstrates that conceptual knowledge is advantageous for sub-tasks. 2) The performance gains more improvement in the medical domain compared with the general domain. We speculate that a specific domain has different data distribution, which indirectly improves the impact of knowledge injection. 3) With only 30% of the training data, the knowledge-enhanced model can achieve comparable performance compared with baselines with 100% training data. This result reveals that knowledge-enhanced approaches are beneficial in a low-resource setting, which is quite useful as specific domains usually are data-hungry. 4) In the medical domain, the method with concept embedding achieves the best performance in intent classification, while the method with conceptualized pretraining is 2% higher than the method with concept embedding in named entity recognition. We argue that this is because named entity recognition is a sequence tagging task, and conceptualized pretraining injects prior knowledge via masking strategies, which boosts the performance. However, intent classification is a classification approach, and there are distinct features that affect the final prediction. In addition, the attention mechanism may contribute to the performance.

4.2.2 Online Evaluation of Information Retrieval.

We evaluate how concept mining can help improve information retrieval through text rewriting. We create an evaluation dataset containing 200 queries from the Alibaba UC browser. For the original query, we collect the top 10 search results returned by the

¹¹<https://github.com/alibaba-research/ChineseBLUE>

Table 5: Online evaluation results.

Information Retrieval		Entity Recommendation	
Search Engine	Relevant Precision	Metrics	Percentage Lift
Baidu	72.5	CTR	5.1%
Baidu+AliCG	84.1	PV	5.5%
Shenma	71.1	UD	0.92%
Shenma+AliCG	83.1	IE	7.01%

Baidu¹² and Shenma¹³ search engines, which are the first and second largest search engines in China, respectively. We first replace the query with k different instances and collect the top search results for each rewritten query from Baidu and Shenma. We then merge and retain 10 of them as search results after query rewriting. We ask three human judges to assess the relevance of the results and record the majority vote of the human judges, which is "relevant" or "irrelevant," and calculate the percentage of the relevance of the original query and rewritten query. As shown in Table 5, for the query after rewriting using our strategy, the percentage of top 10 relevant results increases from 72.5% to 84.1% in Baidu and from 71.1% to 83.1% in Shenma respectively. The reason is that conceptual knowledge helps understand the intent of user queries, which can provide more relevant and clear keywords for search engines. As a result, information retrieval results in a better match for users' intent.

4.2.3 Online A/B Testing of Entity Recommendation.

We perform a large-scale online A/B testing to show how conceptual knowledge helps in improving the performance of entity recommendations (see appendix for details) in real-world applications [28]. We divide users into buckets for the online A/B test. We observe and record the activities of each bucket for seven days and select two buckets with highly similar activities. We utilize the same text rewriting and conceptualized pretraining strategies in section 3 for entity recommendation. The recommendations are obtained without the conceptual graph for one bucket and with the conceptual graph for the other bucket. The page view (PV) and click-through-rate (CTR) are the two most critical metrics in real-world applications because they show how many content users there are and how much time they spend on an application. We also report the Impression Efficiency (IE) and Users Duration (UD). As shown in Table 5, we observe a statistically significant increase in CTR (5.1%), PV (5.5%), UD (0.92%) and IE (7.01%). These observations prove that the conceptual graph for entity recommendation considerably benefits the understanding of queries and helps match users with their potential interested entities.

5 RELATED WORK

Fine-grained Concept Acquisition. Conventional concept mining methods are closely related to noun phrase segmentation and named entity recognition [24]. They either use heuristic methods to extract typed entities or treat the problem as sequence labeling and use large-scale labeled training data to train LSTM-CRF. Another

¹²<https://www.baidu.com/>

¹³<https://m.sm.cn/>

area of focus is term and keyword extraction. They extract noun phrases based on statistical appearance and co-occurrence signals [9] or text features [15]. The latest methods for concept mining rely on phrase quality. [19] adaptively identified concepts based on their quality. More recently, [13] propose an approach to discover user-centered concepts via bootstrapping, query-title alignment, and sequence labeling. Different from their approach, our method is advantageous to extract long-tail concepts and also have the ability to evolve the taxonomy.

Long-tail Concept Mining. Conventional approaches usually fail to extract long-tail concepts. Some existing approaches leverage unsupervised approaches to extract long-tail concepts as phrase mining [19]. However, those approaches can neither leverage the context nor handle scattered concepts. Besides, some studies [13] treat concept mining as sequence tagging. Nevertheless, the performance of those long-tail concepts degrades significantly. Several low-resource sequence tagging approaches that utilize transfer learning [3], or meta-learning [10] have been proposed. However, there are few Chinese resources available, and the meta-learning approach still suffers from low accuracy, which is not satisfactory for real applications.

Taxonomy Evolution. Automatic taxonomy construction is a long-standing task in the literature. Most existing approaches focus on building the entire taxonomy by first extracting hypernym-hyponym pairs and then organizing all hypernymy relations into a tree or DAG structure. In many real-world applications, some existing taxonomies may have already been laboriously curated by experts [12] or via crowdsourcing [16], and are deployed in online systems. Instead of constructing the entire taxonomy from scratch, these applications demand the feature of expanding an existing taxonomy dynamically. There exist some studies on expanding WordNet with named entities from Wikipedia or domain-specific concepts from different corpora [11]. One major limitation of these approaches is that they cannot update the taxonomy dynamically, which is necessary for semantic search. Probase [23] proposes two probability scores, namely *plausibility* and *typicality*, for the conceptual graph. However, it is computationally expensive and time-consuming to update such a big conceptual graph.

6 CONCLUSION

We describe the implementation and deployment of AliCG at Alibaba, which is designed to improve the performance of the semantic search. The system extracts fine-grained concepts from the search logs, and the extracted concepts can be updated based on user behaviors. We have deployed the conceptual graph at Alibaba UC Browser with more than 300 million daily active users. Extensive experimental results show that the system can accurately extract concepts and boost the performance of intent classification and named entity recognition, and Online A/B testing results further demonstrate the efficacy of our approach.

7 ACKNOWLEDGEMENTS

This work is supported by the National Natural Science Foundation of China (Nos. 91846204, 52007173 and U19B2042), National Key

R&D Program of China (No. SQ2018YFC000004), Zhejiang Provincial Natural Science Foundation of China (No. LQ20E070002), and Zhejiang Lab's Talent Fund for Young Professionals (No. 2020KB0AA01).

REFERENCES

- [1] 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*. Springer, 722–735.
- [2] Yoshua Bengio. 2017. The consciousness prior. *arXiv preprint arXiv:1709.08568* (2017).
- [3] Yixin Cao, Zikun Hu, Tat-Seng Chua, Zhiyuan Liu, and Heng Ji. 2019. Low-Resource Name Tagging Learned with Weakly Labeled Data. *arXiv preprint arXiv:1908.09659* (2019).
- [4] Jindong Chen, Yizhou Hu, Jingping Liu, Yanghua Xiao, and Haiyun Jiang. 2019. Deep short text classification with knowledge powered attention. In *In AAAI*.
- [5] Lihan Chen, Jiaqing Liang, Chenhao Xie, and Yanghua Xiao. 2018. Short text entity linking with fine-grained topics. In *In CIKM*. 457–466.
- [6] Xiang Chen, Xin Xie, Ningyu Zhang, Jiahuan Yan, Shumin Deng, Chuanqi Tan, Fei Huang, Luo Si, and Huajun Chen. 2021. AdaPrompt: Adaptive Prompt-based Finetuning for Relation Extraction. *CoRR abs/2104.07650* (2021). arXiv:2104.07650 <https://arxiv.org/abs/2104.07650>
- [7] Shumin Deng, Ningyu Zhang, Jiaojian Kang, Yichi Zhang, Wei Zhang, and Huajun Chen. 2020. Meta-Learning with Dynamic-Memory-Based Prototypical Network for Few-Shot Event Detection. In *In WSDM*, James Caverlee, Xia (Ben) Hu, Mounia Lalmas, and Wei Wang (Eds.).
- [8] Shumin Deng, Ningyu Zhang, Luoqi Li, Hui Chen, Huaixiao Tou, Mosha Chen, Fei Huang, and Huajun Chen. 2021. OntoED: Low-resource Event Detection with Ontology Embedding. In *ACL Association for Computational Linguistics*.
- [9] Katerina Frantzi, Sophia Ananiadou, and Hideki Mima. 2000. Automatic recognition of multi-word terms: the c-value/nc-value method. *International journal on digital libraries* 3, 2 (2000), 115–130.
- [10] Yutai Hou, Zhihan Zhou, Yijia Liu, Ning Wang, Wanxiang Che, Han Liu, and Ting Liu. 2019. Few-Shot Sequence Labeling with Label Dependency Transfer. *arXiv preprint arXiv:1906.08711* (2019).
- [11] David Jurgens and Mohammad Taher Pilehvar. 2015. Reserating the awesometastic: An automatic extension of the WordNet taxonomy for novel terms. In *In NAACL*. 1459–1465.
- [12] Carolyn E Lipscomb. 2000. Medical subject headings (MeSH). *Bulletin of the Medical Library Association* 88, 3 (2000), 265.
- [13] Bang Liu, Weidong Guo, Di Niu, Chaoyue Wang, Shunnan Xu, Jinghong Lin, Kunfeng Lai, and Yu Xu. 2019. A User-Centered Concept Mining System for Query and Document Understanding at Tencent. In *In KDD*. 1831–1841.
- [14] Zhiyuan Liu, Xinxiong Chen, Yabin Zheng, and Maosong Sun. 2011. Automatic keyphrase extraction by bridging vocabulary gap. In *Proceedings of the Fifteenth Conference on Computational Natural Language Learning*. Association for Computational Linguistics, 135–144.
- [15] Olena Medelyan and Ian H Witten. 2006. Thesaurus based automatic keyphrase indexing. In *Proceedings of the 6th ACM/IEEE-CS joint conference on Digital libraries*. 296–297.
- [16] Rui Meng, Yongxin Tong, Lei Chen, and Caleb Chen Cao. 2015. CrowdTC: Crowdsourced taxonomy construction. In *2015 IEEE International Conference on Data Mining*. IEEE, 913–918.
- [17] Rada Mihalcea and Paul Tarau. 2004. Texttrank: Bringing order into text. In *In EMNLP*. 404–411.
- [18] Thomas Rebele, Fabian Suchanek, Johannes Hoffart, Joanna Biega, Erdal Kuzey, and Gerhard Weikum. 2016. YAGO: A multilingual knowledge base from wikipedia, wordnet, and geonames. In *In ISWC*. Springer.
- [19] Jingbo Shang, Jialu Liu, Meng Jiang, Xiang Ren, Clare R Voss, and Jiawei Han. 2018. Automated phrase mining from massive text corpora. *IEEE Transactions on Knowledge and Data Engineering* 30, 10 (2018), 1825–1837.
- [20] Rui Suzuki, Kanako Komiya, Masayuki Asahara, Minoru Sasaki, and Hiroyuki Shinnou. 2018. All-words word sense disambiguation using concept embeddings. In *In LREC*.
- [21] Fang Wang, Zhongyuan Wang, Zhoujun Li, and Ji-Rong Wen. 2014. Concept-based short text classification and ranking. In *In CIKM*. 1069–1078.
- [22] Zhongyuan Wang, Kejun Zhao, Haixun Wang, Xiaofeng Meng, and Ji-Rong Wen. 2015. Query understanding through knowledge-based conceptualization. In *Twenty-Fourth International Joint Conference on Artificial Intelligence*.
- [23] Wentao Wu, Hongsong Li, Haixun Wang, and Kenny Q Zhu. 2012. Probase: A probabilistic taxonomy for text understanding. In *In SIGMOD*.
- [24] Hongbin Ye, Ningyu Zhang, Shumin Deng, Mosha Chen, Chuanqi Tan, Fei Huang, and Huajun Chen. 2021. Contrastive triple extraction with generative transformer. In *In AAAI*.
- [25] Ningyu Zhang, Shumin Deng, Xu Cheng, Xi Chen, Yichi Zhang, Wei Zhang, and Huajun Chen. 2021. Drop Redundant, Shrink Irrelevant: Selective Knowledge Injection for Language Pretraining. In *In IJCAI*.

- [26] Ningyu Zhang, Shumin Deng, Zhanlin Sun, Jiaoyan Chen, Wei Zhang, and Huajun Chen. 2020. Relation Adversarial Network for Low Resource Knowledge Graph Completion. In *In WWW*.
- [27] Ningyu Zhang, Shumin Deng, Zhanlin Sun, Guanying Wang, Xi Chen, et al. 2019. Long-tail Relation Extraction via Knowledge Graph Embeddings and Graph Convolution Networks.. In *In NAACL*.
- [28] Shengyu Zhang, Dong Yao, Zhou Zhao, Tat-Seng Chua, and Fei Wu. 2021. CauseRec: Counterfactual User Sequence Synthesis for Sequential Recommendation. In *In SIGIR*.

A INFORMATION FOR REPRODUCIBILITY

A.1 System Implementation and Deployment

We implement and deploy the AliCG system in Alibaba UC Browser. The fine-grained concept acquisition module, long-tail concept mining module, taxonomy evolution module are implemented in Python 3.6 and run as offline components. The concept embedding and conceptualized pretraining module (inference part) are implemented in C++, and they run as an online service. We utilize high-performance memory storage Tair (a Redis-like storage system) in Alibaba for online data storage. In our system, each component works as a service and is deployed on Alibaba MaxCompute¹⁴. Alibaba MaxCompute is a big high-performance data processing framework. It provides fast and fully managed petabyte-scale data warehouse solutions and allows for the analysis and processing of massive amounts of data economically and efficiently. The online service runs on 200 dockers. Each docker is configured with six 2.5 GHz Intel Xeon Gold 6133 CPU cores and 64 GB memory. Offline fine-grained concept acquisition, long-tail concept mining, and taxonomy evolution are run on 60 dockers with the same configuration. We set $\delta_t = 0.3$, $\alpha = 0.6$, $\beta = 0.8$, and $\delta = 2$ in our system.

Algorithm 3 Offline instance/concept mining

Require: query and high-clicked documents set $(q, D) \in T$.

- 1: **if** succeed **then**
- 2: Perform bootstrapping with alignment consensus on general domain
- 3: Perform conceptualized phrase tagging on specific domains
- 4: Perform self-training with an ensemble consensus
- 5: **if** coverage **then**
- 6: Perform taxonomy evolution
- 7: Perform refinement
- 8: **else**
- 9: Break

Algorithm 4 Offline concept embedding

Require: query and high-clicked documents set $(q, D) \in T$.

- 1: Tag T with all candidate concepts with maximum forward matching
- 2: Train concept embedding model

Algorithms 3, 4 and 5 show the offline running processes of each component in AliCG. Algorithms 6 and 7 show the online running processes of concept embedding and conceptualized pretraining. We also detail the procedure of fine-grained concept acquisition and long-tail concept mining in Figure 7 and Figure 8, respectively. Offline concept mining from the search logs is running every day. It extracts approximately 20,000 concepts from 30 million search logs, and approximately 5,000 of the extracted concepts are new. The offline taxonomy evolution is also daily running. The processing offline usually takes about two hours, and the processing time of online concept embedding is 10,000 queries per second at most.

¹⁴<https://www.alibabacloud.com/>

Algorithm 5 Offline conceptualized pretraining

Require: query and high-clicked documents set $(q, D) \in T$, conceptual graph CG stored in key-value memory storage

- 1: Tag T with all candidate concepts by maximum forward matching from CG
- 2: Duplicate and shuffle the corpus T ten times and generate whole concept masking training samples with a masking rate of 15%.
- 3: Initialize all parameters with BERT-wwm and perform further pretraining with T .
- 4: Finetune sub-tasks with the further pretrained model

Algorithm 6 Online concept embedding inference

Require: query set Q , conceptual graph CG stored in key-value memory storage, concept embedding model M in tf-serving cluster

- 1: **for** each query $q \in Q$ **do**
- 2: Tag q with all candidate concepts by maximum forward matching from CG
- 3: Call tf-serving model M and get results

Algorithm 7 Online concept pretraining inference

Require: query set Q , conceptualised model M in tf-serving cluster

- 1: **for** each query $q \in Q$ **do**
- 2: Call tf-serving model M and get results

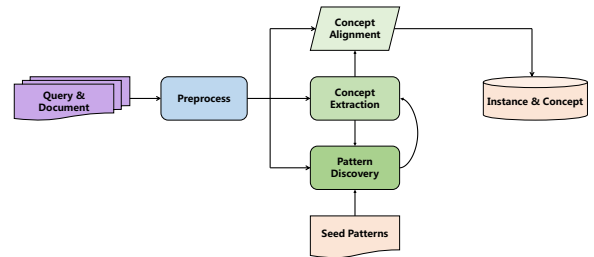


Figure 7: Procedure of fine-grained concept acquisition.

It can perform concept embedding for approximately 800 million queries per day.

A.2 Online Service Details

We utilize a high-performance C++ server framework for online query embedding inference, which can handle more than 10,000 QPS. In addition to the regular offline update pipeline, we develop an emergency makeup pipeline, which has two functions. First, we mine concepts from heat articles (high clicked and page view articles) in real-time and put those heat concepts into the emergency makeup pipeline. In addition, we utilize this pipeline to fix a few bad cases.

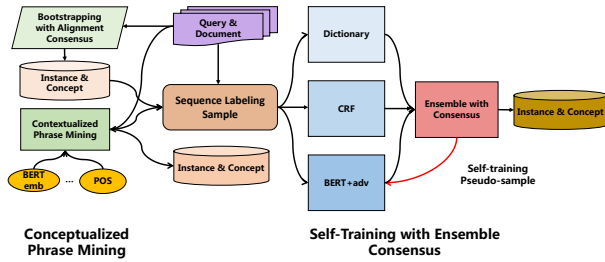


Figure 8: Procedure of long-tail concept mining.

A.3 Concept Reverse Index

As some applications need to know the instances given concepts, for example, it is necessary to retrieve all the instances of the concept "patriotic songs recommended by the Communist Youth League (共青团推荐爱国歌曲)" when user searching such a concept. Thus, we build a reverse concept–instance index based on HA3, which is a high-performance index builder at Alibaba.

A.4 Parameter Settings and Training Process

Here, we introduce the features that we use for different components in our system and describe how we train each component. For concept mining, we randomly sample 15,000 search logs in Alibaba UC Browser. We extract concepts from these search logs using the approaches introduced in Sec. 2.1 and the results are manually checked by Alibaba product managers. The resulting dataset is used to train the classifier in conceptualized phrase-mining-based concept mining and sequence tagging in our model. We utilize CRF++ v0.58 to train our model. 80% of the dataset is used as the training set, 10% as the development set, and the remaining 10% as the test set.

A.5 Publication of Datasets

We have published our datasets for research purposes, and they can be accessed from Github¹⁵. Our dataset’s major difference compared with Tencent’s UCCG dataset is that 1) our dataset is bigger (50x); 2) our dataset has concept possibilities; 3) our dataset contains long-tail concepts. Specifically, we have published the following open-source data:

- **The ACG dataset.** It is used to evaluate the performance of our approach for concept mining, and it contains 490,000 instances, as shown in Figure 9. It is the largest conceptual graph dataset in Chinese.
- **The query tagging dataset.** It is used to evaluate the query tagging accuracy of AliCG, and it contains 10,000 queries with concept tags.
- **The seed concept patterns for bootstrapping-based concept mining.** It contains the seed string patterns that we utilize for bootstrapping-based concept mining from queries.
- **Predefined level-1 and level-2 concept list.** It contains more than 200 predefined concepts for taxonomy construction.

¹⁵<https://github.com/alibaba-research/ConceptGraph>

```

百探 {"kg_info": [{"hot": "367", "guid": "bd77432e-4654-11e5-8ba6-f80f41fb03aa"}],
"concepts": [{"concept": "影视", "score": 1.0, "level": 1}, {"isA": "影视", "concept": "电影",
"score": 1.0, "level": 2}, {"concept": "影视节目", "score": 0.91, "level": 3}, {"concept": "好看的搞笑电影", "score": 0.03, "level": 3}, {"concept": "犯罪悬疑的影视作品", "score": 0.01,
"level": 3}, {"concept": "杜琪峰人气作品", "score": 0.01, "level": 3}, {"concept": "经典的搞笑
电影", "score": 0.01, "level": 3}, {"concept": "杜琪峰影视作品", "score": 0.01, "level": 3},
{"concept": "刘德华经典电影", "score": 0.01, "level": 3}, {"concept": "第33届香港金像奖入围影片",
"score": 0.0, "level": 3}, {"concept": "同为杜琪峰导演高人气电影", "score": 0.0, "level": 3},
{"concept": "心理犯罪的影视作品", "score": 0.0, "level": 3}, {"concept": "郑秀文刘德华电影",
"score": 0.0, "level": 3}, {"concept": "刘德华与郑秀文电影", "score": 0.0, "level": 3}]}

```

Figure 9: Sample data in Alibaba Conceptual Graph.

Table 6: Part of the query-concept samples.

Query	Concept
What are the Hangzhou special local product (杭州的特产有哪些)	Hangzhou special local product (杭州特产)
Pancakes cooking methods list (煎饼的做法大全)	pancakes cooking methods (煎饼的做法)
Which cars are cheap and fuel-efficient? (有什么便宜省油的车)	cheap and fuel-efficient cars (便宜省油的车)
Hangzhou famous snacks (杭州有名的小吃)	Hangzhou snacks (杭州小吃)
What are the symptoms of coronavirus pneumonia (冠状病毒肺炎有什么症状)	symptoms of coronavirus pneumonia (冠状病毒肺炎症状)
Latest fairy TV Series (最新仙侠电视剧)	fairy TV Series (仙侠电视剧)

A.6 Details of Entity Heat Estimation

We introduce the details of entity heat estimation as we leveraged in confidence estimation in Sec. 2.3. The entity heat estimation algorithm calculates the popularity of an entity based on search behavior. We firstly parse user search logs then retrieval all the candidate entities from the knowledge graph via matching. We link those identified entities through various distances between entity mentions and entities. Finally, we sum the user search values and normalize heat values. We update the heat score of entities in the knowledge graph every day.

A.7 Details of Entity Recommendation

We utilized text rewriting and conceptualized pretraining for entity recommendations. For a simple query with entities, we utilize heterogeneous graph embedding to retrieve related entities. For those complex queries with little entities, we propose a deep collaborative matching model to get related entities. Then we rank those entities by various strategies, including type filtering, learning to rank, and click-through rate estimation.

A.8 Examples of Queries and Extracted Concepts

Table 6 lists a few examples of user queries, along with the concepts extracted by AliCG. The concepts are appropriate for summarizing the core user intention in the queries.