# GREASE: A Generative Model for Relevance Search over Knowledge Graphs

Tianshuo Zhou
National Key Laboratory for Novel
Software Technology, Nanjing
University, China
tianshuo.zhou@smail.nju.edu.cn

Ziyang Li
National Key Laboratory for Novel
Software Technology, Nanjing
University, China
zyli@smail.nju.edu.cn

Gong Cheng
National Key Laboratory for Novel
Software Technology, Nanjing
University, China
gcheng@nju.edu.cn

Jun Wang
Department of Computer Science,
University College London, UK
j.wang@cs.ucl.ac.uk

Yu'Ang Wei
National Key Laboratory for Novel
Software Technology, Nanjing
University, China
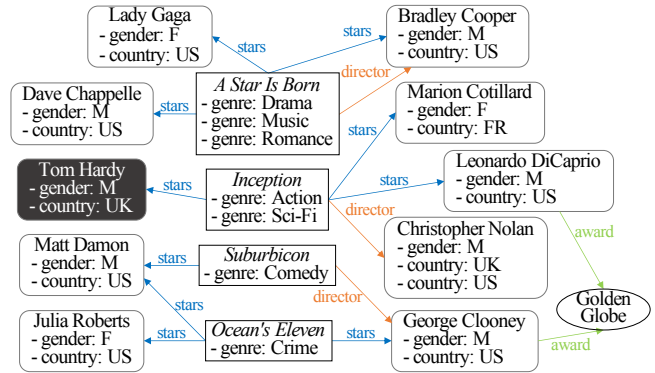weiyuang@smail.nju.edu.cn

## ABSTRACT

Relevance search is to find top-ranked entities in a knowledge graph (KG) that are relevant to a query entity. Relevance is ambiguous, particularly over a schema-rich KG like DBpedia which supports a wide range of different semantics of relevance based on numerous types of relations and attributes. As users may lack the expertise to formalize the desired semantics, supervised methods have emerged to learn the hidden user-defined relevance from user-provided examples. Along this line, in this paper we propose a novel generative model over KGs for relevance search, named GREASE. The model applies to meta-path based relevance where a meta-path characterizes a particular type of semantics of relating the query entity to answer entities. It is also extended to support properties that constrain answer entities. Extensive experiments on two large-scale KGs demonstrate that GREASE has advanced the state of the art in effectiveness, expressiveness, and efficiency.

## 1 INTRODUCTION

**Background.** In a knowledge graph (KG), nodes are entities associated with attributes and interconnected with binary relations as edges. Increasingly many KGs have emerged for various domains, some available as Linked Open Data. KGs for a focused domain often have a simple schema, e.g., the LinkedIn KG[1] describes members, companies, and other entities in the professional domain. There are also generic KGs providing encyclopedic knowledge and hence having a rich schema consisting of thousands of types of relations and attributes, such as DBpedia [10]. We illustrate a KG in Fig. 1 as the running example in this paper. It describes movies, their actors and directors, and awards.

An established data mining task for KG-based applications is *relevance search* [1, 5–7, 9, 15, 17, 18, 20, 22–24]. The task is essentially to find entities in a KG that are the most relevant to an input query entity. However, relevance has a broad range of meanings, particularly over a schema-rich KG. For example, given Tom Hardy as the query entity, the user may search for actresses that co-starred with him, or for American directors that collaborated with him. Unfortunately, non-expert users lack the expertise to

**Figure 1: An example of a knowledge graph, where each entity is associated with a bulleted set of attributes.**

formally characterize the desired semantics of relevance because the formal query language and the rich schema of the KG are both difficult to learn.

**Problem.** To bridge the gap between non-expert users and structured KGs, one practical solution is to request a small number of examples from the user, and then to *learn user-defined relevance from user-provided examples* [1, 3, 5, 6, 9, 15, 22, 24]. The user can specify an example in the form of an ordered *query-answer entity pair* to illustrate the desired semantics of relevance [9, 15, 22]. For example, two different users may provide different sets of examples:

$$
\begin{aligned}
S_1 = \{ &\langle \texttt{Dave Chappelle, Lady Gaga} \rangle, \\
&\langle \texttt{Matt Damon, Julia Roberts} \rangle \}, \\
S_2 = \{ &\langle \texttt{Dave Chappelle, Bradley Cooper} \rangle, \\
&\langle \texttt{Matt Damon, George Clooney} \rangle \}.
\end{aligned}
\tag{1}
$$

The user providing $S_1$ aims to find entities that are relevant to Tom Hardy *just as how* Lady Gaga is relevant to Dave Chappelle and as how Julia Roberts is relevant to Matt Damon. The underlying semantics could be actors—preferably American actresses—that co-starred with Tom Hardy. In this case, Marion Cotillard and Leonardo DiCaprio are acceptable answers. The user providing $S_2$

has a different need and is looking for American male directors and/or actors that collaborated with `Tom Hardy`. Now, `Leonardo DiCaprio` and `Christopher Nolan` become good answers.

**Challenges.** It has been common to assume that the user-defined relevance can be represented by one or more *meta-paths* [3, 6, 9, 15, 18, 20, 22, 24]. A meta-path is a sequence of relation types, i.e., a path at the schema level. For example, the co-starring relation underlying $S_1$ is characterized by the following meta-path:

$$\mathcal{P}_1 : \texttt{[query]} \xleftarrow{\text{stars}} \cdot \xrightarrow{\text{stars}} \texttt{[answer]}. \tag{2}$$

The collaboration relation underlying $S_2$ is characterized partially by $\mathcal{P}_1$ and partially by

$$\mathcal{P}_2 : \texttt{[query]} \xleftarrow{\text{stars}} \cdot \xrightarrow{\text{director}} \texttt{[answer]}. \tag{3}$$

The key challenge to these methods is the selection and weighting of meta-paths. Existing methods train discriminative models to directly learn the weights of meta-paths from user-provided examples [3, 6, 9, 15, 22, 24], but have exhibited the following limitations.

- Their *discriminative models* rely on negative examples that are sampled automatically as the user only provides positive examples. Identifying high-quality negative examples is algorithmically challenging and computationally demanding.
- They are focused on meta-path based relevance, but cannot support the representation of *properties*, such as gender and nationality in the above example. Their capability of characterizing user needs is somewhat limited.

**Contributions.** In this work we address the two challenges and propose GREASE, a novel **G**enerative model for **R**el**Ev****A**nce **SE**arch over KGs. Our implementation has been open source.[2] Our contributions are summarized as follows.

- We treat the weight of a meta-path as a posterior probability, and devise a *generative model* with cost-effective approximations. Our model does not rely on negative examples, and has outperformed existing methods in both effectiveness and efficiency in the experiments.
- We extend the model to support not only meta-path based relevance of answer entities to the query entity, but also *properties that constrain answer entities*. The extension allows to represent more expressive user-defined relevance.

**Organization.** We formulate the problem in Section 2. Our generative model is described in Section 3, and is extended to support properties in Section 4. Based on that, our search algorithm is given in Section 5. We report experiments in Section 6, compare related work in Section 7, and finally conclude the paper in Section 8.

## 2 PROBLEM FORMULATION

Let $\mathbb{R}$ denote the set of all real numbers. We assume countable pairwise disjoint sets of entities $\mathcal{V}$, relation types $\mathcal{R}$, attribute types $\mathcal{A}$, and attribute values $\mathcal{L}$.

*Definition 2.1 (Knowledge Graph).* A knowledge graph (KG) is a directed graph denoted by $G = \langle V, E, \Psi \rangle$, where

- $V \subseteq \mathcal{V}$ is a finite set of entities represented as nodes,
- $E \subseteq V \times \mathcal{R} \times V$ is a finite set of relations between entities represented as directed edges, and

- $\Psi : V \mapsto \mathsf{P}(\mathcal{A} \times \mathcal{L})$, where $\mathsf{P}(\cdot)$ represents power set, is a function that associates each entity $v \in V$ with a finite set of attributes $\Psi(v) \subseteq \mathcal{A} \times \mathcal{L}$.

For an entity $v \in V$, its *properties* consist of

$$\Phi(v) = \{\langle a, l \rangle : \langle a, l \rangle \in \Psi(v) \text{ or } \langle v, a, l \rangle \in E\}. \tag{4}$$

For example, in Fig. 1, the properties of `Suburbicon` consist of an attribute $\langle$genre, Comedy$\rangle$ and two relations $\langle$stars, Matt Dameon$\rangle$ and $\langle$director, George Clooney$\rangle$.

We define *path* in a standard way. It is acyclic, and its edges are not required to follow the same direction. However, in the remainder of the paper we always write right arrows and rewrite left arrow $\xleftarrow{r}$ as $\xrightarrow{r^{-1}}$, where $r^{-1}$ represents the inverse of $r \in \mathcal{R}$. We write a path $p : v_0 \xrightarrow{r_1} v_1 \xrightarrow{r_2} \cdots \xrightarrow{r_l} v_l$ from $v_0$ to $v_l$ as $v_0 \rightsquigarrow_p v_l$ for short. The length of a path is the number of its edges.

*Definition 2.2 (Meta-Path).* A meta-path is a sequence of relation types $\mathcal{P} : r_1 r_2 \cdots r_l$, where $r_1, \ldots, r_l \in \mathcal{R}$ are relation types (or their inverses), and $l$ is called the length of $\mathcal{P}$ denoted by $\texttt{len}(\mathcal{P}) = l$. A path $p$ in a KG follows $\mathcal{P}$, denoted by $p \models \mathcal{P}$, if $p$ is in the form of $p : v_0 \xrightarrow{r_1} v_1 \xrightarrow{r_2} \cdots \xrightarrow{r_l} v_l$.

For example, the two meta-paths shown in Eq. (2) and Eq. (3) are formalized as follows:

$$\mathcal{P}_1 : \texttt{stars}^{-1}\,\texttt{stars}, \qquad \mathcal{P}_2 : \texttt{stars}^{-1}\,\texttt{director}. \tag{5}$$

Their lengths are both 2. $\mathcal{P}_2$ is followed by path

$$\texttt{Tom Hardy} \xrightarrow{\text{stars}^{-1}} \texttt{Inception} \xrightarrow{\text{director}} \texttt{Christopher Nolan}. \tag{6}$$

*Definition 2.3 (Relevance Search).* Given a KG denoted by $G = \langle V, E, \Phi \rangle$, let $\texttt{rel} : \mathcal{V} \times \mathcal{V} \mapsto \mathbb{R}$ be a user-defined real-valued function. For $u, v \in V$, $\texttt{rel}(u, v)$ returns the relevance of $v$ to $u$. Let $k < |V|$ be a predetermined positive integer. For an input query entity $q \in V$, relevance search is to find top-$k$ answer entities $\texttt{Ans}(q) \subseteq (V \setminus \{q\})$ that are the most relevant to $q$ in terms of $\texttt{rel}$.

Following [9, 15, 22], we assume the user provides a small number of ordered query-answer entity pairs as *examples*, to exemplify the desired $\texttt{rel}$ which is not directly accessible to the search system.

*Definition 2.4 (Relevance Search by Example).* Extending Definition 2.3, the problem turns into learning a function $\texttt{rel}$ under the supervision of a set of user-provided examples denoted by $S$. Each example is an ordered query-answer entity pair $\langle s, t \rangle \in V \times V$, where $s$ is called the source entity and $t$ is the target entity, such that $v \in \texttt{Ans}(q)$ is relevant to $q$ just as how $t$ is relevant to $s$.

For example, $S_1$ and $S_2$ in Eq. (1) are two sets of examples for the query entity `Tom Hardy`. They indicate different $\texttt{rel}$ functions.

## 3 GENERATIVE RELEVANCE MODEL

Following Definition 2.4, the $\texttt{rel}$ function is conditioned on a set of user-provided examples $S$, so we rewrite $\texttt{rel}(q, v)$ as $\texttt{rel}(q, v|S)$. Previous research computes meta-path based relevance [3, 6, 9, 15, 18, 20, 22, 24]. Along this line, we decompose $\texttt{rel}$ into a linear

combination of weighted relevance over a set of meta-paths denoted by $\Omega_{mp} = \{\mathcal{P}_1, \ldots, \mathcal{P}_n\}$:

$$\text{rel}(q, v|S) = \sum_{\mathcal{P}_i \in \Omega_{mp}} \gamma(q, v|\mathcal{P}_i) \cdot \text{Pr}(\mathcal{P}_i|S) \cdot \text{J}(\mathcal{P}_i), \qquad (7)$$

where $\gamma(q, v|\mathcal{P}_i)$ measures the real-valued relevance of $v$ to $q$ w.r.t. a particular meta-path $\mathcal{P}_i$, $\text{Pr}(\mathcal{P}_i|S)$ represents the weight of $\mathcal{P}_i$, and $\text{J}(\mathcal{P}_i)$ is a regularization term to prevent overfitting. Meta-paths and their weights are to be learned from $S$. For example, with $S_1$ in Eq. (1), the meta-path $\mathcal{P}_1$ in Eq. (5) should have a large weight because for every example in $S_1$, there exists a path in Fig. 1 that follows $\mathcal{P}_1$ and connects the source entity to the target entity.

To establish rel, below we describe the selection of $\Omega_{mp}$, and the computation of $\gamma$, Pr, and J.

## 3.1 Meta-Path Selection

As rel is exemplified by $S$, we select $\Omega_{mp}$ based on $S$. Our $\Omega_{mp}$ contains all possible meta-paths that can be derived from $S$. A meta-path $\mathcal{P}_i$ is in $\Omega_{mp}$ if there exists a path in the KG such that it follows $\mathcal{P}_i$ and it connects the source entity to the target entity in some user-provided example:

$$\Omega_{mp} = \bigcup_{\langle s, t \rangle \in S} \{\mathcal{P} : \exists p \models \mathcal{P}, \ s \rightsquigarrow_p t\}. \qquad (8)$$

## 3.2 Meta-Path Based Relevance

Extensive research has been conducted to measure the relevance of $v$ to $q$ w.r.t. a particular meta-path $\mathcal{P}_i$ [9, 17, 19, 20, 23]. This is outside the focus of this paper, and we extend *path count* [19] as our measure:

$$\begin{aligned} \gamma(q, v|\mathcal{P}_i) &= \min\{\text{pc}(q, v, \mathcal{P}_i), \ \alpha_{mp}\}, \\ \text{pc}(q, v, \mathcal{P}_i) &= |\{p : p \models \mathcal{P}_i \text{ and } q \rightsquigarrow_p v\}|, \end{aligned} \qquad (9)$$

where $\text{pc}(q, v, \mathcal{P}_i)$ represents the number of paths in the KG that follow $\mathcal{P}_i$ and connect $q$ to $v$, and $\alpha_{mp} > 0$ is a parameter to limit the value of pc and prevent highly skewed values. For example, in Fig. 1, for $\mathcal{P}_2$ in Eq. (5) we have

$$\text{pc}(\text{Tom Hardy, Christopher Nolan}, \mathcal{P}_2) = 1, \qquad (10)$$

because there is only one path shown in Eq. (6) that follows $\mathcal{P}_2$ and connects Tom Hardy to Christopher Nolan.

## 3.3 Generative Meta-Path Weighting

Weighting scheme is the key to the effectiveness of the rel function, and is the focus of our work. Different from existing discriminative methods [3, 6, 9, 15, 22, 24], we treat weight $\text{Pr}(\mathcal{P}_i|S)$ as a posterior probability and propose a novel *generative model*. In the following, we will also use $\text{Pr}(\cdot)$ to denote the probability of an event. We estimate probabilities based on the KG and learn $\text{Pr}(\mathcal{P}_i|S)$ from $S$.

Specifically, we rewrite $\text{Pr}(\mathcal{P}_i|S)$ using Bayes' theorem:

$$\text{Pr}(\mathcal{P}_i|S) = \frac{\text{Pr}(\mathcal{P}_i) \cdot \text{Pr}(S|\mathcal{P}_i)}{\text{Pr}(S)} \propto \text{Pr}(\mathcal{P}_i) \cdot \text{Pr}(S|\mathcal{P}_i), \qquad (11)$$

where the posterior $\text{Pr}(\mathcal{P}_i|S)$ is proportional to the prior $\text{Pr}(\mathcal{P}_i)$ times the likelihood $\text{Pr}(S|\mathcal{P}_i)$. Below we separately compute the prior and the likelihood.

**Computation of the Prior.** For the prior $\text{Pr}(\mathcal{P}_i)$, recall that a meta-path is a sequence of relation types $\mathcal{P}_i : r_1 r_2 \cdots r_l$. We assume the probability of observing the $i$-th relation type $r_i$ in the context history of the preceding $(i - 1)$ relation types $r_1 r_2 \cdots r_{i-1}$ can be approximated by the probability of observing it in the shortened context history of the preceding relation type $r_{i-1}$, i.e., the *first-order Markov property*. This assumption is reasonable and also common on graphs. Specifically, random walks on graphs satisfy this property. Formally, we have

$$\begin{aligned} \text{Pr}(\mathcal{P}_i) = \text{Pr}(r_1 r_2 \cdots r_l) &= \text{Pr}(r_1) \prod_{i=2}^{l} \text{Pr}(r_i|r_1 r_2 \cdots r_{i-1}) \\ &\approx \text{Pr}(r_1) \prod_{i=2}^{l} \text{Pr}(r_i|r_{i-1}), \end{aligned} \qquad (12)$$

which in turn will give rise to the following estimation of $\text{Pr}(\mathcal{P}_i)$ if we estimate $\text{Pr}(r_1)$ and $\text{Pr}(r_i|r_{i-1})$ from frequency counts:

$$\text{Pr}(\mathcal{P}_i) \propto \text{pc}(r_1) \prod_{i=2}^{l} \frac{\text{pc}(r_{i-1} r_i)}{\text{pc}(r_{i-1})}, \qquad (13)$$

where $\text{pc}(\mathcal{P})$ represents the number of paths in the KG that follow meta-path $\mathcal{P}$:

$$\text{pc}(\mathcal{P}) = |\{p : p \models \mathcal{P}\}|. \qquad (14)$$

In Eq. (13), computing pc according to Eq. (14) is inexpensive because those meta-paths are very short, being not longer than 2. For example, in Fig. 1, for $\mathcal{P}_1$ and $\mathcal{P}_2$ in Eq. (5) we have

$$\text{pc}(\mathcal{P}_1) = 18, \quad \text{pc}(\mathcal{P}_2) = 7. \qquad (15)$$

Meanwhile, as the reverse direction of $\mathcal{P}_i$ is also meaningful, we can use the probability of observing the $i$-th relation type $r_i$ in the context history of the succeeding relation type $r_{i+1}$, and then make assumptions and estimate probabilities in a similar way:

$$\text{Pr}(\mathcal{P}_i) \propto \text{pc}(r_l) \prod_{i=1}^{l-1} \frac{\text{pc}(r_i r_{i+1})}{\text{pc}(r_{i+1})}. \qquad (16)$$

To improve the robustness of our model, we take the arithmetic mean of Eq. (13) and Eq. (16) as the final value of $\text{Pr}(\mathcal{P}_i)$.

This arithmetic mean also provides an approximation of $\text{pc}(\mathcal{P})$ for a long meta-path $\mathcal{P}$. The approximation is useful because it may be infeasible to compute the exact frequency count in Eq. (14) on a large KG when $\mathcal{P}$ is long. Formally, for $\mathcal{P} : r_1 r_2 \cdots r_l$, let $\text{apc}(\mathcal{P})$ denote this approximation of $\text{pc}(\mathcal{P})$. We have

$$\text{apc}(\mathcal{P}) = \begin{cases} \text{pc}(\mathcal{P}) & \text{len}(\mathcal{P}) \leq 2, \\ \frac{1}{2}(\text{apc}_{\text{start}}(\mathcal{P}) + \text{apc}_{\text{end}}(\mathcal{P})) & \text{len}(\mathcal{P}) > 2, \end{cases} \qquad (17)$$

where $\text{apc}_{\text{start}}$ and $\text{apc}_{\text{end}}$ denote the right-hand side of Eq. (13) and Eq. (16), respectively. This approximation will be used later.

**Computation of the Likelihood.** For the likelihood $\text{Pr}(S|\mathcal{P}_i)$, the user-provided examples in $S$ are trivially considered to be conditionally independent given $\mathcal{P}_i$:

$$\text{Pr}(S|\mathcal{P}_i) = \prod_{\langle s, t \rangle \in S} \text{Pr}(\langle s, t \rangle | \mathcal{P}_i). \qquad (18)$$

$\Pr(\langle s, t\rangle | \mathcal{P}_i)$ represents the probability that a path following $\mathcal{P}_i$ connects $s$ to $t$, which we estimate from frequency counts:

$$\Pr(\langle s, t\rangle | \mathcal{P}_i) \approx \frac{\text{pc}(s, t, \mathcal{P}_i)}{\text{apc}(\mathcal{P}_i)} , \qquad (19)$$

where $\text{pc}(s, t, \mathcal{P}_i)$ is computed by Eq. (9), and $\text{apc}(\mathcal{P}_i)$ is an approximation of $\text{pc}(\mathcal{P}_i)$ computed by Eq. (17).

To improve the robustness of our model, *smoothing* is needed to avoid cases where $\text{pc}(s, t, \mathcal{P}_i) = 0$ which in turn will lead to $\Pr(S|\mathcal{P}_i) = 0$. In such a case, we replace the zero value of $\text{pc}(s, t, \mathcal{P}_i)$ with the following small non-zero value:

$$\frac{\text{apc}(\mathcal{P}_i)}{|\text{ST}(s)| \cdot |\text{ST}(t)|} , \qquad (20)$$

where $\text{apc}(\mathcal{P}_i)$ is computed by Eq. (17), and $\text{ST}(\cdot)$ denotes the set of entities that have the same (most specific) type as a given entity. Type is an attribute appearing in almost every KG. The above value represents the average number of paths that follow $\mathcal{P}_i$ and connect two entities of the same type as $s$ and as $t$.

## 3.4 Regularization

Long meta-paths are complex and may *overfit user-provided examples*. We impose a penalty on the complexity of $\mathcal{P}_i$. Specifically, we penalize long meta-paths with the following regularization term:

$$\text{J}(\mathcal{P}_i) = e^{-\beta \cdot \text{len}(\mathcal{P}_i)} , \qquad (21)$$

where $\text{len}(\mathcal{P}_i)$ denotes the length of $\mathcal{P}_i$, and $\beta > 0$ is a decay factor.

## 4 EXTENDED FACET-BASED RELEVANCE

Previous research only computes meta-path based relevance [3, 6, 9, 15, 18, 20, 22, 24]. We extend meta-paths to facets. A *facet* is either a meta-path or a property that constrains answer entities. Accordingly, we extend the $\text{rel}$ function in Eq. (7) by adding a linear combination of weighted relevance over a set of properties denoted by $\Omega_{\text{prop}} = \{\langle a_1, l_1\rangle, \ldots, \langle a_n, l_n\rangle\}$:

$$\begin{aligned} \text{rel}(q, v|S) = & \sum_{\mathcal{P}_i \in \Omega_{\text{mp}}} \gamma(q, v|\mathcal{P}_i) \cdot \Pr(\mathcal{P}_i|S) \cdot \text{J}(\mathcal{P}_i) \\ & + \sum_{\langle a_i, l_i\rangle \in \Omega_{\text{prop}}} \gamma(q, v|\langle a_i, l_i\rangle) \cdot \Pr(\langle a_i, l_i\rangle|S) , \end{aligned} \qquad (22)$$

where $\gamma(q, v|\langle a_i, l_i\rangle)$ measures the real-valued relevance of $v$ to $q$ w.r.t. a particular property $\langle a_i, l_i\rangle$, and $\Pr(\langle a_i, l_i\rangle|S)$ represents the weight of $\langle a_i, l_i\rangle$. A property is already simple and hence regularization is not needed. Properties and their weights are also to be learned from $S$. For example, with $S_1$ in Eq. (1), two properties $\langle \text{gender}, \text{F}\rangle$ and $\langle \text{country}, \text{US}\rangle$ should have large weights because they constrain all the target entities in $S_1$.

To establish the extended $\text{rel}$, below we describe the selection of $\Omega_{\text{prop}}$, and the computation of $\gamma$ and $\Pr$ for properties.

## 4.1 Property Selection

A property $\langle a_i, l_i\rangle$ is in $\Omega_{\text{prop}}$ if it constrains the target entity in some user-provided example:

$$\Omega_{\text{prop}} = \bigcup_{\langle s, t\rangle \in S} \Phi(t) . \qquad (23)$$

## 4.2 Property-Based Relevance

We measure the relevance of $v$ w.r.t. $\langle a_i, l_i\rangle$ according to whether $\langle a_i, l_i\rangle$ is a property that constrains $v$, which is independent of $q$:

$$\gamma(q, v|\langle a_i, l_i\rangle) = \begin{cases} \alpha_{\text{prop}} & \langle a_i, l_i\rangle \in \Phi(v), \\ 0 & \langle a_i, l_i\rangle \notin \Phi(v), \end{cases} \qquad (24)$$

where $\alpha_{\text{prop}} > 0$ is a parameter to tune the importance of properties relative to meta-paths in the computation of the extended $\text{rel}$.

## 4.3 Generative Property Weighting

Similar to our generative model for weighting meta-paths, we treat weight $\Pr(\langle a_i, l_i\rangle|S)$ as a posterior probability, and rewrite it using Bayes' theorem:

$$\Pr(\langle a_i, l_i\rangle|S) \propto \Pr(\langle a_i, l_i\rangle) \cdot \Pr(S|\langle a_i, l_i\rangle), \qquad (25)$$

where the posterior $\Pr(\langle a_i, l_i\rangle|S)$ is proportional to the prior $\Pr(\langle a_i, l_i\rangle)$ times the likelihood $\Pr(S|\langle a_i, l_i\rangle)$. Below we separately compute the prior and the likelihood.

**Computation of the Prior.** For the prior $\Pr(\langle a_i, l_i\rangle)$, it represents the probability that the property $\langle a_i, l_i\rangle$ constrains an entity in the KG, which we estimate from frequency counts:

$$\Pr(\langle a_i, l_i\rangle) = \frac{|\{v \in V : \langle a_i, l_i\rangle \in \Phi(v)\}|}{|V|} , \qquad (26)$$

where $V$ is the set of entities in the KG.

**Computation of the Likelihood.** For the likelihood $\Pr(S|\langle a_i, l_i\rangle)$, similar to Eq. (18), the user-provided examples in $S$ are trivially considered to be conditionally independent given $\langle a_i, l_i\rangle$:

$$\Pr(S|\langle a_i, l_i\rangle) = \prod_{\langle s, t\rangle \in S} \Pr(\langle s, t\rangle|\langle a_i, l_i\rangle) = \prod_{\langle s, t\rangle \in S} \Pr(t|\langle a_i, l_i\rangle), \qquad (27)$$

where the last equation holds because the property $\langle a_i, l_i\rangle$ constrains answer entities which correspond to the target entity $t$ in a user-provided example, and hence $\langle a_i, l_i\rangle$ is independent of the source entity $s$ in the example.

$\Pr(t|\langle a_i, l_i\rangle)$ represents the probability that an entity constrained by $\langle a_i, l_i\rangle$ is $t$, which we estimate from frequency counts:

$$\Pr(t|\langle a_i, l_i\rangle) = \begin{cases} \frac{1}{|\{v \in V : \langle a_i, l_i\rangle \in \Phi(v)\}|} & \langle a_i, l_i\rangle \in \Phi(t), \\ 0 & \langle a_i, l_i\rangle \notin \Phi(t). \end{cases} \qquad (28)$$

To improve the robustness of our model, *smoothing* is also needed here to avoid cases where $\Pr(t|\langle a_i, l_i\rangle) = 0$ which in turn will lead to $\Pr(S|\langle a_i, l_i\rangle) = 0$. In such a case, we replace the zero value of $\Pr(t|\langle a_i, l_i\rangle)$ with a small non-zero value $\frac{1}{|V|}$.

## 5 SEARCH ALGORITHM

Finally, we present an efficient implementation to support learning the proposed model from user-provided examples in an online environment and returning answer entities promptly.

## 5.1 Algorithm

The full GREASE algorithm is presented in Fig. 2. MPSearch (line 1) finds $\Omega_{\text{mp}}$ which is defined by Eq. (8). It performs $|S|$ bidirectional searches—one for each example in $S$, starting simultaneously from the source entity and the target entity. The search space is restricted

**Input:** A KG $G = \langle V, E, \Psi \rangle$, a query entity $q$, a set of user-provided examples $S$, an upper bound $L$ on the length of allowable meta-paths, and a positive integer $m$.

**Output:** $k$ top-ranked entities that are relevant to $q$.

1: $\Omega_{\mathrm{mp}} \leftarrow \mathrm{MPSearch}(G, S, L)$;
2: $\Omega_{\mathrm{prop}} \leftarrow \bigcup_{\langle s, t \rangle \in S} \Phi(t)$;
3: $\Omega \leftarrow \Omega_{\mathrm{mp}} \cup \Omega_{\mathrm{prop}}$;
4: **for all** $\Omega_i \in \Omega$ **do**
5:    Compute $\Pr(\Omega_i | S)$;
6: **end for**
7: $\Omega_{\mathrm{top}} \leftarrow m$ meta-paths in $\Omega_{\mathrm{mp}}$ with the largest weights;
8: $C \leftarrow \bigcup_{\mathcal{P}_i \in \Omega_{\mathrm{top}}} \{v \in V : \exists p \models \mathcal{P}_i, \ q \rightsquigarrow_p v\}$;
9: **for all** $v \in C$ **do**
10:    Compute $\mathrm{rel}(q, v | S)$;
11: **end for**
12: **return** $k$ top-ranked entities in $C$

**Figure 2: The GREASE algorithm.**

to meta-paths that are not longer than $L$, which is a predetermined upper bound. Then we find $\Omega_{\mathrm{prop}}$ (line 2) which is defined by Eq. (23). $\Omega_{\mathrm{mp}}$ and $\Omega_{\mathrm{prop}}$ comprise $\Omega$ (line 3), namely all the facets to consider for computing $\mathrm{rel}$ in Eq. (22). For each facet $\Omega_i \in \Omega$, we compute $\Pr(\Omega_i | S)$ according to Section 3.3 and Section 4.3 (lines 4–6). The $m$ meta-paths in $\Omega_{\mathrm{mp}}$ with the largest weights are denoted by $\Omega_{\mathrm{top}}$ (line 7). All the entities that are connected from $q$ by a path following a meta-path in $\Omega_{\mathrm{top}}$ form candidate answer entities, denoted by $C$ (line 8). Here, we only use $m$ meta-paths to identify candidate answer entities for efficiency considerations. For each candidate answer entity, its extended relevance to $q$ is computed (lines 9–11), and the $k$ top-ranked ones are returned.

## 5.2 Indexing

To support efficient computation in an online environment, we precompute and index the following statistics.

We index the frequency counts for all the meta-paths in the KG that are not longer than 2, so their pc values defined by Eq. (14) are retrievable in $O(1)$, and for any other meta-path, its approximate pc value (i.e., apc) defined by Eq. (17) is computable in $O(L)$.

We also index the frequency counts for all the properties in the KG, so Eq. (26) and Eq. (28) are computable in $O(1)$.

We index the frequency counts for all the entity types in the KG, so $|\mathrm{ST}(\cdot)|$ in Eq. (20) is computable in $O(1)$.

## 5.3 Complexity Analysis

Let $\Delta$ be the maximum degree of the nodes in the KG. Let $\Xi$ be the maximum number of properties that constrain an entity in the KG. $\Omega_{\mathrm{mp}}$ is computed in $O(|S| \cdot \Delta^{\lceil \frac{L}{2} \rceil})$ using bidirectional search, and $\Omega_{\mathrm{prop}}$ is computed in $O(|S| \cdot \Xi)$.

To compute the posterior $\Pr(\Omega_i | S)$, when $\Omega_i$ is a meta-path $\mathcal{P}_i$, the prior $\Pr(\mathcal{P}_i)$ is computed by Eq. (13) and Eq. (16) in $O(L)$, and the likelihood $\Pr(S | \mathcal{P}_i)$ is computed by Eq. (18) in $O(L + |S|)$ where $\mathrm{pc}(s, t, \mathcal{P}_i)$ has been computed during MPSearch. When $\Omega_i$ is a property $\langle a_i, l_i \rangle$, the prior $\Pr(\langle a_i, l_i \rangle)$ is computed by Eq. (26) in $O(1)$, and the likelihood $\Pr(S | \langle a_i, l_i \rangle)$ is computed by Eq. (27) in $O(|S|)$.

**Table 1: Statistics about KGs**

| KG | Entity | Relation | Relation Type | Attribute Type |
|---|---|---|---|---|
| DBpedia 2016-10 (Mappingbased Objects) | 5,900,558 | 18,746,174 | 661 | 2,065 |
| YAGO 3.1 (yagoFacts) | 4,295,825 | 12,430,700 | 37 | 1 |

The candidate answer entities $C$ are computed in $O(m \cdot \Delta^L)$.

The computation of the $\mathrm{rel}$ function does not further increase the asymptotic time complexity of the algorithm.

To conclude, in practice, the running time of the entire algorithm is probably dominated by MPSearch.

## 6 EXPERIMENTS

We empirically compared our approach with several state-of-the-art methods based on a variety of queries over two popular KGs. Both effectiveness and efficiency were tested.

## 6.1 Datasets

**Knowledge Graphs.** Our experiments were based on two popular large-scale KGs: DBpedia [10] (version 2016-10) and YAGO [13] (version 3.1). For DBpedia, we obtained a KG from two files: *Mappingbased Objects* and *Instance Types*. For YAGO, we obtained a KG from two files: *yagoFacts* and *yagoSimpleTypes*. The files are in RDF format, where RDF literals and types were treated as attributes, and the other RDF triples became relations. Some statistics about these KGs are summarized in Table 1. Note that YAGO was to be used with the queries created in [6] where attributes are not involved. So we followed [6] to only import *yagoFacts* and *yagoSimpleTypes* which contain relations but no attributes except for type.

**Queries.** We reused 320 query instances[3] given in [6] which were divided into 8 groups (D11–D14 and Y1–Y4). However, this dataset is still limited in two aspects. First, attributes are not considered. Second, the query entity is required to appear as the source entity in every user-provided example, because the method proposed in [6] was specifically designed for this scenario. To overcome these two limitations, we created 800 query instances[4] which were divided into 10 groups (D1–D10).

Our creation of D1–D10 followed a common practice in related work [6, 22]. Compared with D11–D14 and Y1–Y4 created in [6], our D1–D10 are more generalized as they allow the source entity in a user-provided example to be different from the query entity, and they are more challenging as D6–D10 involve properties. In general, their desired semantics of relevance are more complex than all the known queries used in the literature.

Specifically, each group of D1–D10 contains 80 query instances, and their desired semantics of relevance are represented by the same set of predefined facets. We sampled 100 random source-target entity pairs from DBpedia as a pool such that their relevance conformed to the predefined semantics. We chose 20 random pairs from the pool and took their source entities as our query entities. For each query entity, based on the predefined semantics of

---

[3]http://ws.nju.edu.cn/relevance/relsue/
[4]http://ws.nju.edu.cn/relevance/grease/

**Table 2: Query Groups with Examples**

| | Desired Semantics (Facets) | Example Query and Answers | | | |
|---|---|---|---|---|---|
| | | Query Entity | User-Provided Examples | Query Intent | Answer Entities |
| D1 | starring$^{-1}$ director | Howard Duff | ⟨Stephen Wight, Susan Tully⟩ ⟨Vijay Chavan, Kedar Shinde⟩ | director of a movie starring Howard Duff | George Sherman Andre deToth |
| D2 | almaMater$^{-1}$ foundedBy$^{-1}$ | Bowdoin College | ⟨Duke University, Duolingo⟩ ⟨Yale University, Allied Corp⟩ | organization founded by a Bowdoin alumnus | Netflix Pure Software |
| D3 | starring starring$^{-1}$ starring director$^{-1}$ | Charlie Chaplin | ⟨Bam Margera, Brandon Novak⟩ ⟨Rahul Bose, Koel Purie⟩ | actor and also director of a movie starring Charlie Chaplin | Mack Swain Lloyd Bacon |
| D4 | almaMater almaMater$^{-1}$ award award$^{-1}$ | Hagan Bayley | ⟨Ewan Birney, Antony Galione⟩ ⟨George Porter, Charles Coulson⟩ | Hagan Bayley's schoolmate that won the same award | John Mollon Henry Gilman |
| D5 | architecturalStyle architecturalStyle$^{-1}$ location location$^{-1}$ | De Rohan Arch | ⟨Morson's Row, PaceKing House⟩ ⟨Evergreen (Virginia), Greer House⟩ | architecture with the same style and location as the De Rohan Arch | Hompesch Gate La Borsa |
| D6 | starring$^{-1}$ director director$^{-1}$ ⟨genre, Comedy⟩ | Tanya Chisholm | ⟨Casey Kasem, Fantastic Max⟩ ⟨Michael Milhoan, Party Down⟩ | comedy directed by the director of a movie starring Tanya Chisholm | The Last Halloween A Fairly Odd Summer |
| D7 | almaMater almaMater$^{-1}$ foundedBy$^{-1}$ ⟨industry, Software⟩ | Vitalik Buterin | ⟨Peter Clyne, SpringSource⟩ ⟨Felix Villars, Lightbend Inc.⟩ | software company founded by Vitalik Buterin's schoolmate | Databricks Waterloo Maple |
| D8 | influenced$^{-1}$ influenced ⟨field, Physics⟩ | Leo Strauss | ⟨Denis Diderot, Leonhard Euler⟩ ⟨Herbert Feigl, Albert Einstein⟩ | physicist influenced by the same person as Leo Strauss | Isaac Newton David Hilbert |
| D9 | affiliation affiliation$^{-1}$ ⟨type, Private school⟩ | Carleton College | ⟨Viterbo University, Fisk University⟩ ⟨Verdon College, DePaul University⟩ | private school affiliated with the same organization as Carleton College | Manhattan College Drake University |
| D10 | museum$^{-1}$ author ⟨birthPlace, Paris⟩ | Metropolitan Museum of Art | ⟨National Gallery, Georges Seurat⟩ ⟨Van Gogh Museum, Robert Delaunay⟩ | Parisian artist with artworks housed by Metropolitan Museum of Art | Georges Seurat Jacques-Louis David |

relevance, we labeled gold-standard answer entities, and created 4 query instances by choosing different numbers of random pairs from the pool as user-provided examples: $|S| \in \{2, 3, 4, 5\}$. Table 2 illustrates each group with one query instance under $|S| = 2$ and two of its gold-standard answer entities.

The creation of D11–D14 and Y1–Y4 in [6] adopted a similar procedure, and we refer the reader to [6] for details. D11–D14 are based on DBpedia, and Y1–Y4 are based on YAGO. Each group contains 40 query instances.

## 6.2 Baselines

To compare with the state of the art, we chose five strong baselines: PRA [9], RelSim [22], RelSUE [6], ProxE [11], and D2AGE [12]. We intended to also compare with FSPG [15], but we could not obtain its implementation from its authors and we failed to re-implement it due to some missing details in the algorithm.

All the chosen baseline methods except for RelSUE could be tested with all the query instances in our experiments. RelSUE requires the query entity to appear as the source entity in every user-provided example, so it could only be tested with D11–D14 and Y1–Y4 created by the authors of RelSUE.

We obtained implementations of RelSUE, ProxE, and D2AGE from their authors, and we re-implemented PRA and RelSim. For PRA, RelSim, ProxE, and D2AGE, we consistently set their bounds on meta-path length to 3, being sufficiently large for representing all the semantics of relevance in our experiments. RelSUE automatically generated meta-paths of varied lengths.

PRA, RelSim, and RelSUE automatically sampled $10 \cdot |S|$ negative examples for training. For ProxE and D2AGE, a training example is a triple $\langle u, v, w \rangle$ where entity $v$ is more relevant to query entity $u$ than entity $w$. We generated 100 such triples for each query instance by extending our positive example $\langle u, v \rangle$ with a random entity $w$ having the same type as $v$.

**Table 3: NDCG@10 on D1–D5**

| Method | $|S| = 2$ | $|S| = 3$ | $|S| = 4$ | $|S| = 5$ |
|---|---|---|---|---|
| PRA | 0.535 | 0.636 | 0.599 | 0.639 |
| RelSim | 0.452 | 0.562 | 0.555 | 0.575 |
| ProxE | 0.484 | 0.467 | 0.480 | 0.503 |
| D2AGE | 0.549 | 0.509 | 0.555 | 0.637 |
| GREASE | 0.782 | 0.737 | 0.734 | 0.763 |
| GREASE-np | **0.846** | **0.850** | **0.865** | **0.862** |

**Table 4: NDCG@10 on D6–D10**

| Method | $|S| = 2$ | $|S| = 3$ | $|S| = 4$ | $|S| = 5$ |
|---|---|---|---|---|
| PRA | 0.295 | 0.293 | 0.337 | 0.339 |
| RelSim | 0.329 | 0.403 | 0.428 | 0.439 |
| ProxE | 0.445 | 0.440 | 0.450 | 0.426 |
| D2AGE | 0.366 | 0.381 | 0.397 | 0.401 |
| GREASE | **0.831** | **0.840** | **0.866** | **0.874** |

## 6.3 Configuration and Variant of GREASE

For our proposed approach GREASE, by default we set $\alpha_{mp} = 5$ in Eq. (9), $\alpha_{prop} = 2$ in Eq. (24), $\beta = 10$ in Eq. (21), $L = 3$ and $m = 3$ in Algorithm 2. A parameter study will be reported in Section 6.5.

For D1–D5, D11–D14, and Y1–Y4 where only meta-paths but no properties are involved, we implemented a variant of GREASE using only meta-paths as facets, denoted by GREASE-np.

## 6.4 Effectiveness Evaluation

**Evaluation Metric.** For each query instance, the gold standard is a set of relevant answer entities. Each tested method computed $k$ top-ranked answer entities. Their quality was measured by the Normalized Discounted Cumulative Gain (NDCG) at rank position $k$, referred to as NDCG@$k$. Due to space limitations, we only present NDCG@10 scores (i.e., $k = 10$).

**Table 5: NDCG@10 on D11–D14**

| Method | $|S| = 2$ | $|S| = 3$ | $|S| = 4$ | $|S| = 5$ |
|---|---|---|---|---|
| PRA | 0.465 | 0.520 | 0.550 | 0.568 |
| RelSim | 0.644 | 0.656 | 0.654 | 0.666 |
| RelSUE | 0.901 | 0.952 | 0.948 | 0.971 |
| ProxE | 0.410 | 0.402 | 0.410 | 0.371 |
| D2AGE | 0.627 | 0.672 | 0.746 | 0.697 |
| GREASE | 0.971 | **0.978** | **0.953** | **0.973** |
| GREASE-np | **0.995** | 0.968 | 0.942 | 0.968 |

**Table 6: NDCG@10 on Y1–Y4**

| Method | $|S| = 2$ | $|S| = 3$ | $|S| = 4$ | $|S| = 5$ |
|---|---|---|---|---|
| PRA | 0.215 | 0.144 | 0.144 | 0.181 |
| RelSim | 0.274 | 0.336 | 0.357 | 0.367 |
| RelSUE | **0.770** | 0.843 | **0.873** | 0.880 |
| ProxE | 0.568 | 0.592 | 0.562 | 0.608 |
| D2AGE | 0.670 | 0.637 | 0.735 | 0.647 |
| GREASE | 0.724 | **0.861** | 0.860 | **0.900** |
| GREASE-np | 0.673 | 0.677 | 0.674 | 0.703 |



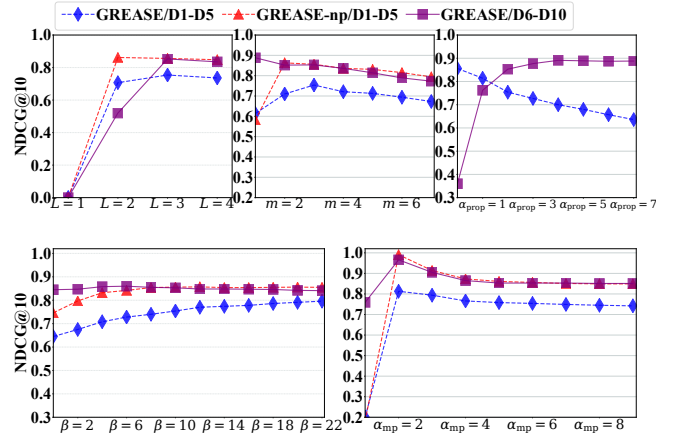**Figure 3: Influence of parameters on our approach.**

**Results on D1–D5.** The average NDCG@10 scores on D1–D5 are presented in Table 3. The results are categorized by number of user-provided examples (i.e., $|S|$). GREASE outperformed all the baselines by at least 0.101–0.233 under different values of $|S|$. GREASE-np achieved even higher scores, exceeding the baselines by at least 0.214–0.297. Recall that for the query instances in D1–D5, their desired semantics of relevance are represented by only meta-paths. Therefore, the superiority of GREASE-np over the baselines demonstrated the effectiveness of our proposed generative model for weighting meta-paths.

Compared with GREASE-np, the small drops in GREASE's scores suggested that its extended model mistakenly assigned large weights to some properties. The extended model is expressive in character-izing user-defined relevance but is then more prone to errors due to the expansion of the search space. Nevertheless, the satisfying performance of GREASE showed that it achieved a good trade-off between expressiveness and accuracy.

Another finding was GREASE and GREASE-np already achieved high scores when $|S| = 2$. Their performance did not increase notably when $|S|$ increased. It indicated that using our approach, users can obtain quite accurate answers with a small effort.

**Results on D6–D10.** The average NDCG@10 scores on D6–D10 are presented in Table 4. Query instances in D6–D10 require using properties that constrain answer entities. Not surprisingly, GREASE largely surpassed all the baselines by at least 0.386–0.435 under different values of $|S|$, as it was the only method that explicitly pro-cessed properties. It confirmed the expressiveness of our extended model in supporting the representation of user-defined relevance.

**Results on D11–D14 and Y1–Y4.** The average NDCG@10 scores on D11–D14 and Y1–Y4 are presented in Table 5 and Table 6, respec-tively. Query instances in D11–D14 and Y1–Y4 represent a special case of our problem, where the query entity appears as the source entity in every user-provided example. RelSUE was specifically

optimized for this scenario and represented the state of the art. Its scores were very high on D11–D14, in the range of 0.901–0.972 under different values of $|S|$. GREASE performed even better, al-though their differences were not large: 0.002–0.070. On Y1–Y4, GREASE led when $|S| = 3$ and $|S| = 5$, whereas RelSUE was better when $|S| = 2$ and $|S| = 4$. We concluded that GREASE was com-parable with RelSUE in this special setting. It demonstrated the effectiveness and generalizability of our approach.

## 6.5 Parameter Study

In Fig. 3, we present the influence of the five parameters on our approach: $L$, $m$, $\alpha_{prop}$, $\beta$, and $\alpha_{mp}$.

$L$ and $m$ in Algorithm 2 tune the trade-off between expressive-ness and efficiency in our approach. $L$ bounds the length of allow-able meta-paths, and $m$ bounds the number of meta-paths used to identify candidate answer entities. All the desired semantics of relevance in our experiments can be represented by meta-paths not longer than 3. GREASE exactly peaked when $L = 3$. Its scores did not change much when $L$ increased to 4. That would give us more flexibility in practice. As to $m$, our approach achieved good results when $m$ was small, owing to the high quality of the computed meta-paths with the largest weights. When $m$ was larger, more noise could be introduced, but the performance of our approach appeared rather stable.

$\alpha_{prop}$ in Eq. (24) tunes the importance of properties relative to meta-paths. On D1–D5 where properties are not needed, larger values of $\alpha_{prop}$ led to poorer results. The setting of this parameter would depend on the needs of the application.

$\beta$ in Eq. (21) tunes the degree of penalizing long meta-paths to prevent overfitting. $\alpha_{mp}$ in Eq. (9) bounds the value of path count to prevent highly skewed values. The performance of GREASE was more sensitive to $\beta$ on D1–D5 than on D6–D10, because D1–D5 totally rely on meta-paths. The performance was generally not very sensitive to $\alpha_{mp}$ unless it was inappropriately set to 1 which disabled path count.
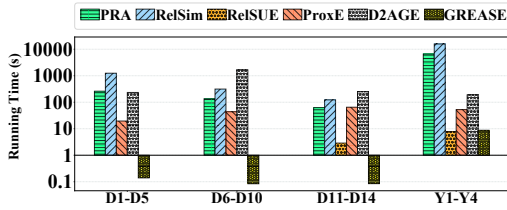
Figure 4: Average running time per query instance.



Figure 5: Average running time of GREASE per query instance.

## 6.6 Efficiency Evaluation

Our experiments were performed on a 3.40GHz Xeon. The pre-computed indexes for GREASE only used 190MB for DBpedia and 123MB for YAGO.

In Fig. 4, for each method we report its average running time per query instance. GREASE satisfyingly completed a search task in less than 1s on DBpedia (D1–D14), at least an order of magnitude faster than all the baselines. It used less than 10s on YAGO (Y1–Y4), being comparable with RelSUE which was optimized for this special scenario. The results demonstrated the efficiency of our approach.

PRA, RelSim, and GREASE ran slower on YAGO than on DBpedia. These methods search the KG to generate all possible meta-paths of a bounded length. YAGO contains much more paths to explore than DBpedia—about 40 times in our experiments, due to the existence of hub nodes in YAGO.

In Fig. 5, we show the influence of $|S|$ on the efficiency of GREASE. The influence was not significant on DBpedia (D1–D14), which suggested the scalability of our approach. However, on YAGO (Y1–Y4), the running time exhibited a linear correlation with $|S|$.

## 7 RELATED WORK

### 7.1 Unsupervised Similarity Search

Relevance search originates from similarity search, for which methods are mainly based on random walks. ObjectRank [2] computes the stationary probability that a random surfer starting from the query entity is at a particular entity as their similarity. PathSim [20] requires random walks to follow a predefined meta-path, to compute similarity with a specified type of semantics. JoinSim [23] uses a slightly different measure. PReP [18] extends PathSim and JoinSim. It computes cross-meta-path synergy, which goes beyond a linear combination of meta-paths. In [7], meta-path is extended to meta-structure which is a directed acyclic graph.

These methods have found application in entity resolution and entity clustering [21] where similarity measurement is the core task. However, similarity is only one special type of relevance. Without the supervision of the user, these methods are not suitable for the more generalized relevance search because they cannot distinguish between a wide range of semantics of relevance on a KG.

### 7.2 Supervised Relevance Search

Relevance is rather ambiguous on a schema-rich KG. Unfortunately, users may not have the expertise to formally characterize the desired semantics of relevance, du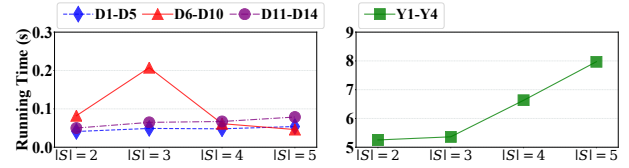e to the complexity of the query language or the richness of the schema. In order to learn user-defined relevance, existing methods are mainly supervised by user-provided examples. An early work is SRW [1], which leverages user-provided examples to supervise random walks for computing relevance. In [3, 24], predefined meta-paths are used to constrain random walks, and user-provided examples are used to learn the weights of the meta-paths. However, it is unrealistic to predefine meta-paths for all possible types of information needs that users may have on a schema-rich KG. It is also difficult for a non-expert user to identify appropriate meta-paths from numerous candidates.

To tackle the problem, PRA [9] and RelSim [22] automatically generate all possible meta-paths but they have to bound the length of an allowable meta-path because the number of possible meta-paths increases exponentially with length. In [5], length-bounded meta-path is extended to size-bounded meta-graph. In FSPG [15] and RelSUE [6], there is no explicit length bound, but long meta-paths are penalized and hence are more likely to be pruned in their greedy search algorithms. All these methods train a discriminative model to learn the weight of each meta-path or meta-graph from user-provided positive examples and randomly sampled negative examples. By contrast, we present a generative model which does not rely on negative examples. It outperforms the above discriminative methods in the experiments, in both effectiveness and efficiency. Moreover, our approach generalizes meta-paths into facets which also include properties that constrain answer entities, and hence it supports more expressive representation of user-defined relevance.

Some recent efforts learn graph embedding models for relevance [11, 12]. However, they did not show better performance than other methods in our experiments. Their complex models may be more suitable for the link prediction task [4, 14], where large training sets are available. In relevance search, a user is not likely to provide many examples to supervise.

### 7.3 Other Related Problems

Other related problems include graph query by example [8] and exemplar query answering [16], but their technical challenges and methods are fundamentally different. Their input is a tuple of entities [8] or a keyword query [16] provided by the user as an example. The example is expanded [8] or mapped [16] into a subgraph of the KG, called a query graph, to capture the user's query intent. The output is a set of other top-ranked subgraphs of the KG that are similar to the query graph. By contrast, relevance search is mainly focused on the selection, weighting, and combination of meta-paths to represent the user-defined relevance between the source entity and the target entity in the user-provided examples.

# 8 CONCLUSIONS

We proposed GREASE, a new approach to relevance search over KGs. Compared with existing methods, GREASE is distinguished by its more effective generative model for weighting meta-paths, its more expressive facet-based representation of relevance with properties, and its efficient implementation. These technical contributions have been demonstrated in an extensive evaluation.

One limitation of our approach is that our estimation of probabilities is largely based on frequency counts in the KG. However, KGs in the real world may be inexact or incomplete, which may affect the accuracy of our estimation. We have used smoothing methods to partially address this issue, but further attempts may be helpful, e.g., using external knowledge. In future work, we will also consider using ontological schemata and reasoning services to handle more complex semantics of relevance.

## REFERENCES

[1] Lars Backstrom and Jure Leskovec. 2011. Supervised random walks: predicting and recommending links in social networks. In *Proc. WSDM*. 635–644. https://doi.org/10.1145/1935826.1935914

[2] Andrey Balmin, Vagelis Hristidis, and Yannis Papakonstantinou. 2004. ObjectRank: Authority-Based Keyword Search in Databases. In *Proc. VLDB*. 564–575.

[3] Shaoli Bu, Xiaoguang Hong, Zhaohui Peng, and Qingzhong Li. 2014. Integrating meta-path selection with user-preference for top-k relevant search in heterogeneous information networks. In *Proc. CSCWD*. 301–306. https://doi.org/10.1109/CSCWD.2014.6846859

[4] HongYun Cai, Vincent W. Zheng, and Kevin Chen-Chuan Chang. 2018. A Comprehensive Survey of Graph Embedding: Problems, Techniques, and Applications. *IEEE Trans. Knowl. Data Eng.* 30, 9 (2018), 1616–1637. https://doi.org/10.1109/TKDE.2018.2807452

[5] Yuan Fang, Wenqing Lin, Vincent Wenchen Zheng, Min Wu, Kevin Chen-Chuan Chang, and Xiaoli Li. 2016. Semantic proximity search on graphs with metagraph-based learning. In *Proc. ICDE*. 277–288. https://doi.org/10.1109/ICDE.2016.7498247

[6] Yu Gu, Tianshuo Zhou, Gong Cheng, Ziyang Li, Jeff Z. Pan, and Yuzhong Qu. 2019. Relevance Search over Schema-Rich Knowledge Graphs. In *Proc. WSDM*. 114–122. https://doi.org/10.1145/3289600.3290970

[7] Zhipeng Huang, Yudian Zheng, Reynold Cheng, Yizhou Sun, Nikos Mamoulis, and Xiang Li. 2016. Meta Structure: Computing Relevance in Large Heterogeneous Information Networks. In *Proc. SIGKDD*. 1595–1604. https://doi.org/10.1145/2939672.2939815

[8] Nandish Jayaram, Arijit Khan, Chengkai Li, Xifeng Yan, and Ramez Elmasri. 2015. Querying Knowledge Graphs by Example Entity Tuples. *IEEE Trans. Knowl. Data Eng.* 27, 10 (2015), 2797–2811. https://doi.org/10.1109/TKDE.2015.2426696

[9] Ni Lao and William W. Cohen. 2010. Relational retrieval using a combination of path-constrained random walks. *Machine Learning* 81, 1 (2010), 53–67. https://doi.org/10.1007/s10994-010-5205-8

[10] Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, Dimitris Kontokostas, Pablo N. Mendes, Sebastian Hellmann, Mohamed Morsey, Patrick van Kleef, Sören Auer, and Christian Bizer. 2015. DBpedia - A large-scale, multilingual knowledge base extracted from Wikipedia. *Semantic Web* 6, 2 (2015), 167–195. https://doi.org/10.3233/SW-140134

[11] Zemin Liu, Vincent W. Zheng, Zhou Zhao, Fanwei Zhu, Kevin Chen-Chuan Chang, Minghui Wu, and Jing Ying. 2017. Semantic Proximity Search on Heterogeneous Graph by Proximity Embedding. In *Proc. AAAI*. 154–160.

[12] Zemin Liu, Vincent W. Zheng, Zhou Zhao, Fanwei Zhu, Kevin Chen-Chuan Chang, Minghui Wu, and Jing Ying. 2018. Distance-Aware DAG Embedding for Proximity Search on Heterogeneous Graphs. In *Proc. AAAI*.

[13] Farzaneh Mahdisoltani, Joanna Biega, and Fabian M. Suchanek. 2015. YAGO3: A Knowledge Base from Multilingual Wikipedias. In *Proc. CIDR*.

[14] Víctor Martínez, Fernando Berzal, and Juan Carlos Cubero Talavera. 2017. A Survey of Link Prediction in Complex Networks. *ACM Comput. Surv.* 49, 4 (2017), 69:1–69:33. https://doi.org/10.1145/3012704

[15] Changping Meng, Reynold Cheng, Silviu Maniu, Pierre Senellart, and Wangda Zhang. 2015. Discovering Meta-Paths in Large Heterogeneous Information Networks. In *Proc. WWW*. 754–764. https://doi.org/10.1145/2736277.2741123

[16] Davide Mottin, Matteo Lissandrini, Yannis Velegrakis, and Themis Palpanas. 2016. Exemplar queries: a new way of searching. *VLDB J.* 25, 6 (2016), 741–765. https://doi.org/10.1007/s00778-016-0429-2

[17] Chuan Shi, Xiangnan Kong, Philip S. Yu, Sihong Xie, and Bin Wu. 2012. Relevance search in heterogeneous networks. In *Proc. EDBT*. 180–191. https://doi.org/10.1145/2247596.2247618

[18] Yu Shi, Po-Wei Chan, Honglei Zhuang, Huan Gui, and Jiawei Han. 2017. PReP: Path-Based Relevance from a Probabilistic Perspective in Heterogeneous Information Networks. In *Proc. SIGKDD*. 425–434. https://doi.org/10.1145/3097983.3097990

[19] Yizhou Sun, Rick Barber, Manish Gupta, Charu C. Aggarwal, and Jiawei Han. 2011. Co-author Relationship Prediction in Heterogeneous Bibliographic Networks. In *Proc. ASONAM*. 121–128. https://doi.org/10.1109/ASONAM.2011.112

[20] Yizhou Sun, Jiawei Han, Xifeng Yan, Philip S. Yu, and Tianyi Wu. 2011. PathSim: Meta Path-Based Top-K Similarity Search in Heterogeneous Information Networks. *PVLDB* 4, 11 (2011), 992–1003.

[21] Yizhou Sun, Brandon Norick, Jiawei Han, Xifeng Yan, Philip S. Yu, and Xiao Yu. 2012. Integrating meta-path selection with user-guided object clustering in heterogeneous information networks. In *Proc. SIGKDD*. 1348–1356. https://doi.org/10.1145/2339530.2339738

[22] Chenguang Wang, Yizhou Sun, Yanglei Song, Jiawei Han, Yangqiu Song, Lidan Wang, and Ming Zhang. 2016. RelSim: Relation Similarity Search in Schema-Rich Heterogeneous Information Networks. In *Proc. SDM*. 621–629. https://doi.org/10.1137/1.9781611974348.70

[23] Yun Xiong, Yangyong Zhu, and Philip S. Yu. 2015. Top-k Similarity Join in Heterogeneous Information Networks. *IEEE Trans. Knowl. Data Eng.* 27, 6 (2015), 1710–1723. https://doi.org/10.1109/TKDE.2014.2373385

[24] Xiao Yu, Yizhou Sun, Brandon Norick, Tiancheng Mao, and Jiawei Han. 2012. User guided entity similarity search using meta-path selection in heterogeneous information networks. In *Proc. CIKM*. 2025–2029. https://doi.org/10.1145/2396761.2398565