

# Selecting Walk Schemes for Database Embedding

Yuval Lev Lubarsky  
Technion  
Haifa, Israel  
lubarsky@cs.technion.ac.il

Martin Grohe  
RWTH Aachen University  
Aachen, Germany  
grohe@informatik.rwth-aachen.de

Jan Tönshoff  
RWTH Aachen University  
Aachen, Germany  
toenshoff@informatik.rwth-aachen.de

Benny Kimelfeld  
Technion  
Haifa, Israel  
bennyk@cs.technion.ac.il

## ABSTRACT

Machinery for data analysis often requires a numeric representation of the input. Towards that, a common practice is to embed components of structured data into a high-dimensional vector space. We study the embedding of the tuples of a relational database, where existing techniques are often based on optimization tasks over a collection of random walks from the database. The focus of this paper is on the recent FORWARD algorithm that is designed for dynamic databases, where walks are sampled by following foreign keys between tuples. Importantly, different walks have different schemas, or “walk schemes,” that are derived by listing the relations and attributes along the walk. Also importantly, different walk schemes describe relationships of different natures in the database.

We show that by focusing on a few informative walk schemes, we can obtain tuple embedding significantly faster, while retaining the quality. We define the problem of scheme selection for tuple embedding, devise several approaches and strategies for scheme selection, and conduct a thorough empirical study of the performance over a collection of downstream tasks. Our results confirm that with effective strategies for scheme selection, we can obtain high-quality embeddings considerably (e.g., three times) faster, preserve the extensibility to newly inserted tuples, and even achieve an increase in the precision of some tasks.

## CCS CONCEPTS

• **Information systems** → **Relational database model**; • **Computing methodologies** → **Machine learning**.

## KEYWORDS

Database embedding, random walks, walk schemes

### ACM Reference Format:

Yuval Lev Lubarsky, Jan Tönshoff, Martin Grohe, and Benny Kimelfeld. 2023. Selecting Walk Schemes for Database Embedding. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*

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 the author(s) 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](mailto:permissions@acm.org).

CIKM '23, October 21–25, 2023, Birmingham, United Kingdom

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-0124-5/23/10...\$15.00

<https://doi.org/10.1145/3583780.3615052>

(CIKM '23), October 21–25, 2023, Birmingham, United Kingdom. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3583780.3615052>

## 1 INTRODUCTION

Machine-learning algorithms are conventionally designed to generalize observations about numerical vectors, and hence, their application to non-numeric data requires *embeddings* of these data into a numerical vector space. The embedding should faithfully reflect the semantics of data in the sense that similar entities are to be mapped to vectors that are close geometrically and vice versa. Instantiations of this practice include models like WORD2VEC [17] and BERT (Bidirectional Encoder Representations from Transformers) [9] that map words (or tokens of words) in natural language [15], NODE2VEC that maps nodes of a graph [11], TRANSE [5] that map entities of a knowledge graph, MOL2VEC [14] that maps molecule structures, and EMBDI [6] and FORWARD [22] that map database tuples. Database embeddings have enabled the deployment of machine-learning architectures to traditional database tasks such as record similarity [2–4, 12, 13], record linking [10, 19] integration tasks such as schema, token and record matching (entity resolution) [6], and column prediction [22]. The embedded entities are typically either tuples or attribute values. Here we focus on tuple embeddings.

Embedding techniques are often based on the analysis of *sequences* obtained from the data. In word embedding, the data is naturally organized in sequences (e.g., sentences or sliding windows in the text) [9, 17]; in node embedding, the sequences are paths obtained from random walks in the graph [11]; and in tuple embedding, the sequences consist of database components (cells and tuples) that one can reach through natural joins [6] or foreign-key references [22]. The analysis is typically done by learning to predict masked parts of the sequence from other parts of the sequence [6, 9, 11, 17]. We focus on FORWARD that is designed for producing stable embeddings in dynamic databases, as we explain next. FORWARD analyzes *walks*, which are sequences of tuples connected via foreign-key references, and it does so differently from masking. Roughly speaking, training aims for the distance between two (vector representations of) tuples to capture the distance between the *distributions* of values reachable from the tuples through foreign-key references, starting from the corresponding tuples. (We recall the exact definition of FORWARD in Section 2.)

The FORWARD algorithm has been designed to solve the *stable embedding problem*, where the goal is to infer the embedding of the new tuples *without recomputing the embedding over the entire database and without changing the embeddings of existing tuples*

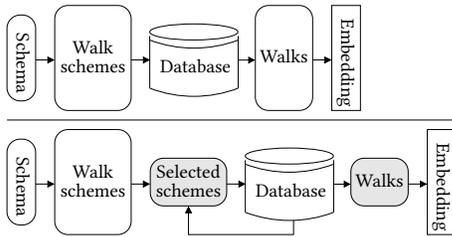


Figure 1: FoWARD vs. FoWARD with scheme selection

(that downstream tasks might already be using and rely on past decisions thereupon) [22]. A technique for performing this task has been proposed along with the FoWARD framework [22].

A walk in a database is naturally associated with a meta-data *pattern*, which is formed by taking the names of the relations of the tuples along the sequence, as well as the names of the attributes that are used for the (outgoing and incoming) references. This pattern is called a *walk scheme* [22], and examples of these are depicted in Figure 4 in the context of a geographical database. Sequence-based embedding algorithms for other modalities do not encounter (and do not account for) such meta-information in the training phase. Our premise is that the walk scheme determines, to a large extent, the contribution of a walk to the quality of the learned embedding. Hence, unlike word and node embedding, in tuple embedding, we can introduce important a-priori bias over the training sequences.

We claim and prove empirically that one can considerably reduce the number of training walks by restricting the learning phase to the walks of the most effective walk schemes, with a mild (or no) reduction in quality. Moreover, the embedding quality might improve by filtering out walk schemes that contribute more noise than beneficial information. We devise techniques for the selection of walk schemes within FoWARD, as illustrated in Figure 1. To illustrate the importance of scheme selection, a sample of our experiments is shown in Figure 2. Here, we are using the learned embedding in the Mondial database [21] to predict the religion of a country based on the database’s information. Each curve corresponds to a selected percentage of the walk schemes and shows the quality of the prediction as a function of the embedding time (where each epoch contributes a point). The actual way of selecting the walk schemes is discussed in the next paragraphs. As the chart showed, selecting a fifth of the walk schemes fully preserves the quality in about one-third of the embedding time and eventually even outperforms the embedding with the entire set of walk schemes.

The main question then is *how to select the best walk schemes for learning an embedding?* This is the challenge that we focus on in this paper. The goodness of a selection method is reflected in two measures: (1) **Efficiency**—the choice should be considerably faster than the embedding itself, and (2) **Quality**—the choice should be such that we can select just a small number of walk schemes and retain the quality of the embedding (e.g., on downstream tasks). Regarding the efficiency measure, we already said that an important (but not the only) use case for walk-scheme selection is that the combined time it takes to select the walk schemes and then learn the embedding is considerably faster than learning the embedding on the initial full collection of walk schemes. As for the quality

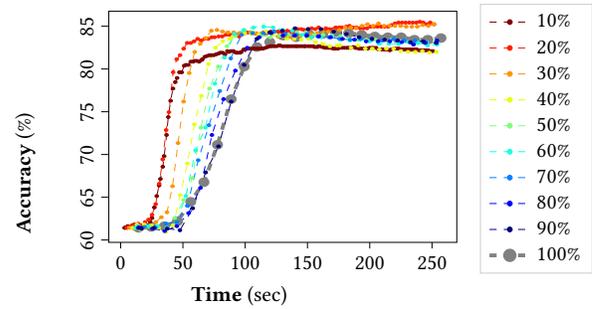


Figure 2: Religion prediction over the Mondial dataset with FoWARD and scheme selection via *kernel variance*. With a fifth of the walk schemes, we get to full equality in about one-third of the embedding time and eventually even outperform the embedding with the entire set of walk schemes.

measure, if we are to select  $\alpha\%$  of the walk schemes and have downstream success of  $\beta\%$  compared to 100% of the walk schemes, then we would like  $\alpha$  to be as small as possible (e.g., 10) and  $\beta$  as large as possible (e.g., 95). Actually, our experiments show cases where  $\beta$  exceeds 100 for the reason discussed above. Importantly, the quality of the selection strategy should also apply to the stable embedding problem in the dynamic setting. In particular, we would like the performance of the learned embedding to reach (or even exceed) that of the full set of walk schemes when the embedding is extended to newly arriving tuples.

The choice can be made by *ranking* the walk schemes by some scoring function and selecting the top candidates. This scoring function can reflect properties that are hypothetically important to the embedding but do not require seeing in action the embedding algorithm, namely FoWARD; we refer to this approach as *FoWARD-less*. Alternatively, one can execute the embedding algorithm in some limited (light) manner and infer the walks from that execution (or, more precisely, from the internal state of the model); we refer to this approach as *light training*. Finally, one can also eliminate walk schemes gradually *during* the embedding process, where in each epoch, we estimate the importance of a walk scheme (similarly to the way it is done in light training) and leave a strict subset for the next epochs, until we decide that none can be further eliminated; we refer to this approach as *online scheme elimination*. In summary, we devise and study strategies for selecting walk-schemes in three different approaches: (1) FoWARD-less, (2) light training, and (3) online scheme elimination.

In each approach, several different strategies can be proposed. In the FoWARD-less approach, we look at measures that inspect the probability distribution that one establishes by following random walks guided by the walk scheme (in addition to simple baselines such as eliminating the longest schemes and random schemes). The strategy that stands out is what we call *kernel variance*: what is the variance among the *differences* that one observes by selecting two random starting points for the walks? In the light-training approach, we look at two types of training restrictions: a single epoch (out of all epochs) and a full training on a *sample* of the database. In the online elimination approach, we apply the single-epoch selection (of light training) after every epoch.

In the experiments, we follow the convention of evaluating the embedding quality on *downstream tasks* and use column-prediction tasks in the bio-medical and geography domains. Our study (Section 5) has three parts. First, we test the performance of each selection strategy as a function of the time and number of selected schemes. Second, we conduct a comparison among the strategies in a technique that we devise. Third, we study the performance with the selected walk schemes in a dynamic setting. Our conclusion is that kernel-variance performs best.

In summary, our contributions are as follows. First, we introduce the problem of *scheme selection* for tuple embedding. We do so in the context of FoRWARD, yet the problem applies to every sequence-based database embedding. Second, we propose three approaches to scheme selection and devise several specific strategies within each approach. Third, we conduct a thorough experimental evaluation that investigates the empirical effectiveness of the strategies, compares them, and studies their performance in the dynamic setting.

## 2 PRELIMINARIES

*Relational Model.* A *database schema*  $\sigma$  consists of a finite collection of *relation schemas*  $R(A_1, \dots, A_k)$  where  $R$  is a distinct *relation name* and each  $A_i$  is a distinct *attribute name*. Each attribute  $A$  is associated with a *domain*, denoted  $\text{dom}(A)$ . Each relation schema  $R(A_1, \dots, A_k)$  has a unique *key*, denoted  $\text{key}(R)$ , such that  $\text{key}(R) \subseteq \{A_1, \dots, A_k\}$ . A *fact* over a relation schema  $R(A_1, \dots, A_k)$  has the form  $R(a_1, \dots, a_k)$  where  $a_i \in \text{dom}(A_i)$  for all  $i = 1, \dots, k$ . A *database*  $D$  over the schema  $\sigma$  is a finite set of facts over the relation schemas of  $\sigma$ . In addition, such an  $a_i$  can be missing and given as a distinguished *null* value. The fact  $R(a_1, \dots, a_k)$  is also called an *R-fact* and a  $\sigma$ -*fact*. We denote by  $R(D)$  the restriction of  $D$  to its  $R$ -facts. For a fact  $f = R(a_1, \dots, a_k)$  over  $R(A_1, \dots, A_k)$ , we denote by  $f[A_i]$  the value  $a_i$ , and by  $f[B_1, \dots, B_\ell]$  the tuple  $(f[B_1], \dots, f[B_\ell])$ . The *active domain* of an attribute  $A$  (w.r.t. to the database  $D$ ), denoted  $\text{adom}_D(A)$ , is the set  $\{f[A] \mid f \in R(D)\}$ .

A *foreign-key constraint* (FK) is an inclusion dependency of the form  $R[B] \subseteq S[C]$  where  $R$  and  $S$  are relation names,  $B = B_1, \dots, B_\ell$  and  $C = C_1, \dots, C_\ell$  are sequences of distinct attributes of  $R$  and  $S$ , respectively, and  $\text{key}(S) = \{C_1, \dots, C_\ell\}$ . For every FK  $R[B] \subseteq S[C]$  and  $R$ -fact  $f \in D$  there exists an  $S$ -fact  $g \in D$  such that  $f[B] = g[C]$ . We then say that  $f$  *references*  $g$ .

*FoRWARD.* The goal of FoRWARD is to derive an embedding function  $\gamma : D \rightarrow \mathbb{R}^k$  for the tuples in a database. Here, the dimension  $k > 0$  is a hyperparameter. The general objective is to compute an embedding  $\gamma$  that represents the data in a way that makes it accessible for data analysis and machine learning algorithms. To this end, FoRWARD learns embeddings that encode the distribution of values seen along random walks through the database. Next, we introduce FoRWARD's notion of random walks through databases formally and recap how these walks are used to produce an embedding. A *walk scheme*  $s$  has the form

$$R_0[A^0] - [B^1]R_1[A^1] - [B^2]R_2[A^2] - \dots - [B^\ell]R_\ell \quad (1)$$

such that for all  $k = 1, \dots, \ell$ , either  $R_{k-1}[A^{k-1}] \subseteq R_k[B^k]$  is an FK or  $R_k[B^k] \subseteq R_{k-1}[A^{k-1}]$  is an FK. We say that  $s$  has *length*  $\ell$ , that it *starts from*  $R_0$ , and that it *ends with*  $R_\ell$ .

A *walk* with the scheme  $s$  is a sequence  $(f_0, \dots, f_\ell)$  of facts such that  $f_k$  is an  $R_k$ -fact and  $f_{k-1}[A^{k-1}] = f_k[B^k]$  for all  $k = 1, \dots, \ell$ . We say that  $(f_0, \dots, f_\ell)$  *starts from*, or has the *source*,  $f_0$ , and that it *ends with*, or has the *destination*,  $f_\ell$ . FoRWARD allows walk schemes and walks of length zero; the walks of this scheme have the form  $(f_0)$  and consist of the fact  $f_0$ . A *random walk* with  $s$  defines a distribution over the destinations. Formally, let  $f_0 = f$  be an  $R_0$ -fact. We denote by  $\mathcal{W}(f, s)$  the distribution over the walks with the walk scheme  $s$  where each walk is sampled by starting from  $f_0$  and then iteratively selecting  $f_k$ , for  $k = 1, \dots, \ell$ , randomly and uniformly from the set  $\{f \in R_k \mid f[B^k] = f_{k-1}[A^{k-1}]\}$ . We denote by  $d_{f,s}$  the random element that maps each walk in  $\mathcal{W}(f, s)$  to its destination—the last fact in the walk. For  $g \in R^\ell(D)$ , the probability that a walk sampled from  $\mathcal{W}(f, s)$  ends with  $g$  is  $\Pr(d_{f,s} = g)$ .

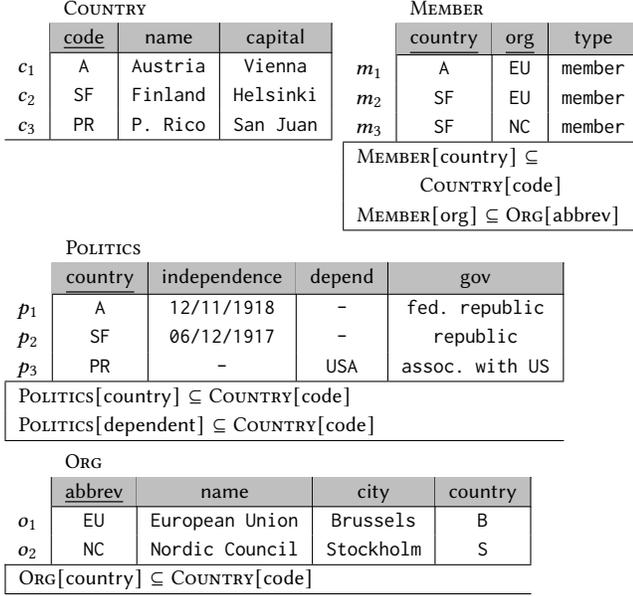
*Targeted Walk Schemes.* A *targeted walk scheme* is a pair  $(s, A)$  such that  $s$  is a walk scheme from a relation  $R$  to a relation  $R'$ , and  $A$  is an attribute of  $R'$ . Given a start fact  $f$  in the start relation  $R$  of the walk scheme  $s$ , a targeted walk scheme  $(s, A)$  defines the random variable  $d_{f,s}[A]$  that forms the value of the random walk's destination in the attribute  $A$ . We denote by  $\text{TWS}(R, \ell_{\max})$  is the set of all targeted walk schemes  $(s, A)$  such that  $s$  is a walk scheme of length at most  $\ell_{\max}$  starting from the relation  $R$  (and ending in any relation that includes  $A$ ). For example, Figure 4 shows several targeted walk schemes over the schema of the database  $D$  of Figure 3. To illustrate,  $(s_7, A_7)$  is the targeted walk scheme with  $A_7 := \text{name}$  and  $s_7 := \text{COUNTRY}[\text{code}] - [\text{country}] \text{MEMBER}[\text{org}] - [\text{abbrev}] \text{ORG}$ . Walks of  $s_7$  include the sequences  $(c_1, m_1, o_1)$  and  $(c_2, m_3, o_2)$ . The distribution  $d_{c_2, s_7}[A_7]$  is uniform between European Union and Nordic Council. The walk schemes  $(s_1, A_1), \dots, (s_7, A_7)$  are in the set  $\text{TWS}(\text{COUNTRY}, 3)$  since they all start with COUNTRY and have length at most three, yet  $(s_8, A_8)$  is not in  $\text{TWS}(\text{COUNTRY}, 3)$  but rather in  $\text{TWS}(\text{COUNTRY}, 4)$ .

Recall that databases may have nulls. A random walk starting at  $f$  might end at an  $R_\ell$ -fact  $g$  with null on  $A$ . As a convention, we define the probability distribution of  $d_{f,s}[A]$  by ignoring the nulls (and normalizing). With this modification, we enforce  $d_{f,s}[A] \in \text{dom}(A)$ . This will be crucial in Section 2, where we define similarity measures for  $d_{f,s}[A]$  based on  $\text{dom}(A)$ .

*Kernelized Domains.* FoRWARD assumes that every attribute  $A$  is associated with a kernel function  $\kappa_A$  that maps all pairs of elements from  $\text{dom}(A)$  to the nonnegative reals. Intuitively,  $\kappa_A(a, b)$  measures the similarity between elements  $a, b \in \text{dom}(A)$ . Kernel functions offer a straightforward way of encoding domain knowledge by modeling the similarity of the domain values. Kernels are also helpful when dealing with noisy data. For example, kernels based on the edit distance can be used to smooth out random typos in text. FoRWARD uses these kernel functions to define similarity measures for the random variables  $d_{f,s}[A]$ .

Let  $s$  be a walk scheme of length  $\ell$  from  $R$  to  $R'$ . Let  $A$  be an attribute of  $R'$  and let  $f$  and  $f'$  be two distinct  $R$ -facts. Then  $d_{f,s}[A]$  and  $d_{f',s}[A]$  are random variables over a shared kernelized domain  $\text{dom}(A)$ . The *Expected Kernel Distance* KD is the expected distance between two random values selected independently at random:

$$\text{KD}(d_{s,f}[A], d_{s,f'}[A]) = \mathbb{E}_{\mathcal{W}(f,s) \times \mathcal{W}(f',s)} [\kappa_A(d_{s,f}[A], d_{s,f'}[A])] \quad (2)$$



**Figure 3: Example of a database, with foreign-key constraints, taken from the Mondial dataset.**

FORWARD uses the Expected Kernel Distance to quantify the similarity between  $d_{f,s}[A]$  and  $d_{f',s}[A]$  with respect to the kernel  $\kappa_A$ .

*Embedding.* Intuitively, FORWARD aims to encode the kernel  $\text{KD}(d_{s,f}[A], d_{s,f'}[A])$ . The primary output is the embedding  $\varphi : D \rightarrow \mathbb{R}^k$ . Additionally, the FORWARD algorithm learns an auxiliary embedding  $\psi : \text{TWS}(R, \ell_{\max}) \rightarrow \mathbb{R}^{d \times d}$  that maps each targeted walk scheme  $(s, A)$  to a symmetric matrix  $\psi(s, A)$ . (Recall that  $\text{TWS}(R, \ell_{\max})$  is the set of targeted walk schemes of length at most  $k$ .) The objective is to find  $\varphi$  and  $\psi$  satisfying  $\varphi(f)^\top \psi(s, A) \varphi(f') = \text{KD}(d_{s,f}[A], d_{s,f'}[A])$  for all  $f, f' \in R(D)$  and  $(s, A) \in \text{TWS}(R, \ell_{\max})$ . Effectively, FORWARD minimizes

$$|\varphi(f)^\top \psi(s, A) \varphi(f') - \text{KD}(d_{s,f}[A], d_{s,f'}[A])| \quad (3)$$

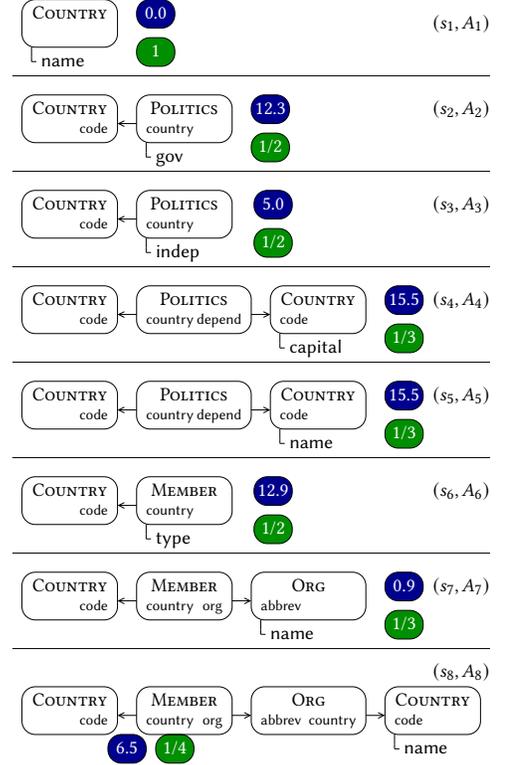
for all  $f, f', s$  and  $A$ , via Stochastic Gradient Descent (SGD).

Intuitively, the algorithm learns an inner product  $\langle \cdot, \cdot \rangle_{s,A}$  on the latent space of the embedding  $\varphi$  defined by  $\langle x, y \rangle_{s,A} = x^\top \psi(s, A) y$  for all  $s$  and  $A$  such that the similarity of facts  $f$  and  $f'$  with respect to this inner product matches the similarity between the random variables  $d_{f,s}[A]$  and  $d_{f',s}[A]$  with respect to the kernel  $\kappa_{R,A}$ .

FORWARD uses gradient descent to learn  $\varphi$  and  $\psi$ . During training, tuples of the form  $(f, f', s, A, g, g')$  are sampled, where  $f$  and  $f'$  are  $R$  facts from the database and  $(s, A) \in \text{TWS}(R, \ell_{\max})$ . The  $R_\ell$  facts  $g$  and  $g'$  are the destinations of random walks with scheme  $s$  sampled for  $f$  and  $f'$ , respectively. We use a hyperparameter  $n_{\text{samples}} \in \mathbb{N}$ . For each  $R$ -fact  $f$  and scheme  $(s, A) \in \text{TWS}(R, \ell_{\max})$ , we sample  $n_{\text{samples}}$  tuples  $(f, f', s, A, g, g')$  with  $f' \neq f$ . The following loss is then minimized for each sample using SGD:

$$\mathcal{L} = \frac{1}{2} |\varphi(f)^\top \psi(s, A) \varphi(f') - \kappa_{R,A}(g[A], g'[A])|^2. \quad (4)$$

This objective uses  $\kappa_{R,A}(g[A], g'[A])$  as a (stochastic) estimate of  $\text{KD}(d_{s,f}[A], d_{s,f'}[A])$ .



**Figure 4: Examples of targeted walk schemes of length one to four, for the database schema of Figure 3. All walk schemes start at the COUNTRY relation. The figure of the walk scheme for  $(s, A)$  shows  $s$  as a path of rectangles and  $A$  (e.g. name) as an attribute under the rightmost (last) rectangle.**

### 3 PROBLEM DEFINITION

Recall that FORWARD uses all targeted walk schemes of length at most  $\ell_{\max}$  for training its embedding. That is, all of  $\text{TWS}(R, \ell_{\max})$ . Our conjecture in this work is that some targeted walk schemes are considerably more useful than others for the task of embedding, and furthermore, that a small subset of  $\text{TWS}(R, \ell_{\max})$  suffices for achieving the quality of the full set. If so, then we can considerably reduce the training time by focusing on just a few walk schemes. It is also perceivable that, with a small yet valuable collection of targeted walk schemes, we can surpass the quality of the original set for the same training time.

Consequently, our goal in this work is to find a subset  $\mathcal{T}'$  of  $\text{TWS}(R, \ell_{\max})$  such that the training of FORWARD on  $\mathcal{T}'$  rather than  $\text{TWS}(R, \ell_{\max})$  is more effective. In particular, we would like to find a  $\mathcal{T}'$  such that: (1)  $\mathcal{T}'$  is small compared to  $\text{TWS}(R, \ell_{\max})$ , and importantly an epoch of training with  $\mathcal{T}'$  is considerably faster than training with  $\text{TWS}(R, \ell_{\max})$ ; and (2) the quality (on downstream tasks) of the embedding resulting from  $\mathcal{T}'$  is high compared to  $\text{TWS}(R, \ell_{\max})$ . Hence, we would like to find a  $\mathcal{T}'$  that would enable us to train considerably faster without a penalty of loss in quality.

*Example 3.1.* One way of selecting a subset of  $\text{TWS}(R, \ell_{\max})$  is to score each and take the top- $k$  for a desired number  $k$  of schemes. For

illustration, Figure 4 depicts two scores (written in filled ellipses) beside each targeted walk scheme. The first (blue) score is what we later define as the *kernel variance* score. The second (green) is a simplistic score that we use as a baseline; this is the reciprocal of the scheme’s length (e.g., it is  $1/3$  for  $(s_7, A_7)$  since  $s_7$  is of length three). If we use the second scoring function, then we take the shortest of the targeted walk schemes (and apply tie-breaking if needed). ♦

## 4 STRATEGIES FOR SCHEME SELECTION

Recall that our goal is to study how walk schemes should be selected in order to establish a proper balance between the execution cost and the quality of the embedding. Our general approach is to start with a large collection of walk schemes (i.e., the initial one of FORWARD) and eliminate walk schemes one by one. In this section, we propose several strategies for such elimination.

Technically, a *strategy*  $T$  assigns to every  $(s, A) \in \text{TWS}(R, \ell_{\max})$  a number  $\text{score}_T(s, A, D)$  for the database  $D$ , where a higher score means that the targeted walk scheme is considered more valuable. When we select  $k$  schemes to eliminate, we select the bottom  $k$  according to the score. We propose strategies that fall into three categories: (1) The *FORWARD-less strategies* determine  $\text{score}_T(s, A, D)$  based on an evaluation that does not require the actual execution of FORWARD. (2) In the *light training* strategies, we run FORWARD in some limited and light fashion in order to determine  $\text{score}_T(s, A, D)$ . (3) The strategy of *online scheme elimination* incorporates the scheme selection in the actual embedding phase (using FORWARD) while targeted walk schemes are eliminated during the epochs of the training; hence, the values  $\text{score}_T(s, A, D)$  are computed multiple times during the embedding phase.

### 4.1 FORWARD-Less Strategies

This category includes simplistic baseline strategies such as the *length* of the scheme, which is illustrated in Figure 4 in the green ellipses. Next, we describe two more involved strategies: *mutual information* and *kernel variance*.

*Mutual Information.* A walk scheme  $s$ , as defined in (1), induces a probability distribution over random walks, which are sequences  $(f_1, \dots, f_\ell)$  of facts in  $R_1, \dots, R_\ell$ , respectively. For  $i = 1, \dots, \ell$ , let  $X_i$  be the random variable that takes the random fact  $f_i$ . Let  $p(f_i, f_{i+1}) \in [0, 1]$  denote the marginal joint distribution of the variables  $X_i$  and  $X_{i+1}$ . Recall that the *mutual information* of  $X_i$  and  $X_{i+1}$  is given by:

$$I(X_i; X_{i+1}) := \sum_{f_i \in X_i, f_{i+1} \in X_{i+1}} p(f_i, f_{i+1}) \log \frac{p(f_i, f_{i+1})}{p(f_i)p(f_{i+1})} \quad (5)$$

We estimate the probabilities  $p(f_i, f_{i+1})$  as the empirical probabilities, that is, their probability in our samples.

We score a targeted walk scheme by the minimal mutual information along the walk. Formally, we have the following score:

$$\text{score}_{\text{mi}}(s, A, D) := - \min_{0 \leq i < \ell+1} I(X_i | X_{i+1}) \quad (6)$$

Note that the minus sign means that we favor schemes with small mutual information. The rationale is that small mutual information encourages the embedding to capture less predicted distributions. As we show later, this rationale is consistent with our experiments.

*Kernel Variance.* The measure *kernel variance* is, intuitively, one that favours targeted walk schemes where different start tuples are associated with varied distributions, and so, the embedding of FORWARD is encouraged to distinguish between these starting tuples. Formally, for a targeted walk scheme  $(s, A)$  we define the *kernel variance* score, denoted  $\text{score}_{\text{kvar}}(s, A, D)$ , as the variance of the expected kernel distance between the distributions induced by  $s$  and  $A$  when starting with random  $f$  and  $f'$  in the source of  $s$ .

$$\text{score}_{\text{kvar}}(s, A, D) := \text{Var}_{f, f'} ( \text{KD}(d_{s,f}[A], d_{s,f'}[A]) ) \quad (7)$$

We can estimate  $\text{score}_{\text{kvar}}(s, A, D)$  from a pool of samples of the form  $(f, f', g, g')$ , where  $g$  and  $g'$  are destinations of random walks of  $s$  starting at  $f$  and  $f'$ , respectively. The number of samples is a hyper-parameter, and in our experiments, we used 10% of the number of samples that FORWARD uses for computing its embedding.

### 4.2 Light Training

If we run the full embedding over all schemes, then we can track the experience of the embedding algorithm with respect to the different schemes. We can then score the contribution of the targeted walk scheme based on the accumulated loss incurred by instances of the targeted walk scheme. While this approach works well empirically, it beats an important purpose of the scheme reduction—to reduce the execution cost of the embedding. So, our approach is to apply this idea lightly, that is, we run the training phase but either stop it early (“Early Termination”), run it only a small sample of the data (“Sample”), or combine between the two. In what follows, we present a precise materialization of these alternatives.

For the next strategies, we need some notation. Recall that the training of FORWARD involves several steps: (1) For each  $R$ -fact  $f$  and  $(s, A) \in \text{TWS}(R, \ell_{\max})$ , we sample tuples of the form  $(f, f', s, A, g, g')$ . (2) For each  $(s, A) \in \text{TWS}(R, \ell_{\max})$ , the  $R_\ell$  facts  $g$  and  $g'$  are the destinations of random walks with scheme  $s$  sampled for  $f$  and  $f'$ , respectively. (3) Using SGD, we minimize the loss  $L_{f, f', s, A, g, g'}$  according to Equation (4). Note that the loss  $L_{f, f', s, A, g, g'}$  is computed for each choice of  $(s, A)$ ,  $f$ ,  $f'$ ,  $g$  and  $g'$ . We denote by  $L_i(s, A, D)$  the mean of the losses  $L_{f, f', s, A, g, g'}$  for each combination of  $s$ ,  $A$  and  $D$ , until the  $i$ th epoch.

*Single-Epoch Training.* This score is the mean loss after an epoch:

$$\text{score}_{\text{1ep}}(s, A, D) := L_1(s, A, D) \quad (8)$$

Note that the computation of this score is disconnected from the training phase, where we train from scratch without accounting for the epoch we spend for computing.

*Sampling.* The second type of light training uses a sample of the database. It is crucial to have a sample where the walk schemes materialize, so we need to carefully select the sample. We create a sample  $D'$  of the database  $D$  in two steps: (1) randomly select a small set  $F$  of facts from  $D$ ; (2) insert to  $D'$  all of the facts of  $D$  that are reachable from  $F$  through paths of foreign keys (in both directions). Yet, with this approach, we still suffer from walk schemes with too few instances in  $D'$ . Thus, we construct  $F$  by considering every scheme  $s$  and selecting random facts from those that participate in paths of  $s$ . Finally, we run ten epochs and take the average loss:

$$\text{score}_{\text{smp}}(s, A, D) := L_{10}(s, A, D') \quad (9)$$

**Table 1: Datasets used in the experiments. “#TWS” is the cardinality of  $TWS(R, \ell_{\max})$ , that is, the number of targeted walk schemes of length up to  $\ell_{\max}$ . “Avglen” is the average length of a walk scheme in  $TWS(R, \ell_{\max})$ .**

Dataset	#Rel.	#Tuples	#Attr.	#TWS	$\ell_{\max}$	Avglen
<b>Mondial</b>	40	21497	167	63	3	2.44
<b>World</b>	3	5411	24	60	3	1.68
<b>Hepatitis</b>	7	12927	26	21	3	1.73
<b>Genes</b>	3	6063	15	32	3	2.25
<b>Mutagen.</b>	3	10324	14	58	4	3.10

### 4.3 Online Scheme Elimination

This strategy is similar to the single-epoch training, except that we apply it iteratively during training, where in each epoch we remove the bottom- $k$  schemes according to the score computed for that epoch according to Equation (8). This way, we are potentially allowing to account for schemes that become more important once other schemes are removed (in earlier epochs). Note that online scheme elimination is different from the light-training approach in the sense that the latter is performed as a pre-processing step that takes place *before* the embedding computation, while the former is performed *as part* of the embedding computation.

## 5 EXPERIMENTAL EVALUATION

The goal of our experimental study is threefold. First, we evaluate the effectiveness of the scheme-selection strategies in terms of the quality of the embedding and the execution cost of its computation. Second, we compare the strategies. Third, we study the impact of the strategy on the performance in a dynamic setting where new tuples are repeatedly inserted, and their embedding is computed without changing the embedding of existing tuples.

To evaluate an embedding, we adopt the common approach of measuring the accuracy on downstream predictions. Hence, we measure the running time of the embedding algorithm, namely FoWARD, and the quality of a learned model for the downstream task. We focus on multi-relational data and use the same databases and tasks that were used for the evaluation of FoWARD [22].

### 5.1 Experimental Setup

**5.1.1 Datasets and Tasks.** Information about the datasets and downstream tasks of our experiments is given in Tables 1 and 2, respectively. Each dataset is a database of multiple relations (with the number of relations given in the “#Rel.” column of Table 1), and the task is to predict the content of an attribute of one of the relations. Hereafter, we refer to this relation as the *prediction relation*. In different downstream tasks on the same dataset, the *prediction attribute* is changed in the prediction relation to the one we aim to predict. Importantly, the predicted attribute is excluded from the database throughout the entire embedding phase, and it is seen by neither FoWARD nor the walk-scheme selector.

**Mondial** contains information from multiple geographical resources [16]. We used multiple attributes for prediction tasks on this dataset: religion (Christian or not), continent, infant mortality g40 (whether the rate is lower than forty per thousand), gdp g8e3

**Table 2: Downstream tasks. “CC” (Common Class) is the frequency of the common value of the predicted bit.**

Downstream task	Pred. Rel.	Pred. Attr.	#Samples	CC
<b>M.-Religion</b>	Target	religion	206	54.8%
<b>M.-Continent</b>		continent	242	22.7%
<b>M.-Infant Mort.</b>		infant g40	238	60.5%
<b>M.-GDP</b>		gdp g8e3	238	50.0%
<b>M.-Inflation</b>		inflation g6	238	50.8%
<b>World</b>	Country	continent	239	24.2%
<b>Hepatitis</b>	Dispat	type	500	58.8%
<b>Genes</b>	Classific.	localization	862	42.5%
<b>Mutagenesis</b>	Molecule	mutagenic	188	66.4%

(whether GDP is lower than \$8000M), and inflation g6(whether the inflation rate is lower than 6%).

**World** has geographical data on states and their cities [18]. The task is to predict the continent of a country. The dataset contains 40 different relations with a total of 167 attributes and 21,497 tuples. We use the whole database and use the TARGET relation as the prediction relation as previously done by Bina et al. [1].

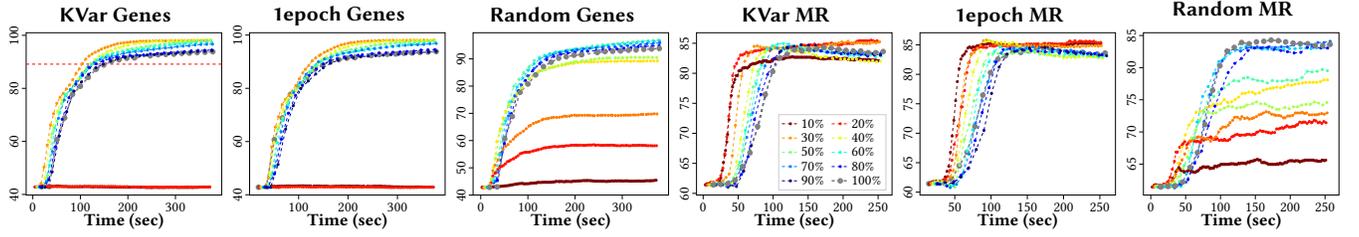
**Hepatitis** is from the 2002 ECML/PKDD Discovery Challenge.<sup>1</sup> The task is to predict the type column, which is either *Hepatitis B* or *Hepatitis C* based on medical examinations. There are in total 206 instances of the former and 484 cases of the latter. The relation with the predicted column contains, in addition to the type classification, the age, gender and identifier of the patient. The other relations contain the rest of the medical data. The dataset contains 7 relations with a total of 26 attributes and 12,927 tuples.

**Genes** [7] contains genomic and drug-design data. The task is to predict the localization of the gene, based on biological data, with 15 different labels. The prediction relation contains only the class and an identifier for the gene, while the rest contain biological data such as the function, gene type, cellular location, and the expression correlation between genes. The dataset contains 3 relations with a total of 15 attributes and 6,063 tuples. (We deleted two tuples with a unique class to prevent split in-balances during cross-validation.)

**Mutagenesis** contains data on the mutagenicity of molecules on *Salmonella Typhimurium* [8]. The task is to predict the mutagenicity of molecules, based on chemical properties of the molecule, with 122 positive samples and 63 negative ones. The prediction relation contains the binary class, molecule ID, and some of the chemical data, while the other relations contain more chemical data and information about the relations between the molecules. The dataset contains 3 relations with a total of 14 attributes and 10,324 tuples.

**5.1.2 Compared Strategies.** We compare the following strategies. **Random** eliminates random schemes. **Length** eliminates the schemes with the longest length. **KVar** eliminates the schemes with the least score<sub>kvar</sub> (kernel variance), as defined in Equation (7). **MI** eliminates the schemes with the least score<sub>mi</sub> (mutual information) as defined in Equation (5). **1epoch** eliminate the schemes with the least score<sub>1ep</sub> (one epoch), as defined in Equation (8). **Sampling** eliminate the schemes with the least score<sub>smp1</sub> (sampling), as defined in

<sup>1</sup><https://sorry.vse.cz/~berka/challenge/PAST/>. We use the modified version of [20].



**Figure 5: Performance on the downstream task as a function of time for different ratios of schemes used for training. At the end of each training epoch of FoRWARD we record the time (x-axis) and accuracy of the downstream task (y-axis).**

Equation (9). Recall that in this strategy, we run the algorithm for ten epochs on the sample  $D'$  of the database  $D$ ; afterward, we run FoRWARD from scratch on the remaining schemes. [Online](#) is the online scheme-elimination approach is described in Section 4.3.

We have a full separation between the embedding process and the downstream task: we generate the embedding independently from the task (as opposed to training for the task) and then use these embeddings as the input to a downstream classifier (that sees only the embeddings and none of the other database information).

**5.1.3 Programming.** For a tuple embedding  $\varphi$ , we denote by  $\alpha(\varphi) \in [0, 1]$  the mean accuracy achieved by an SVM (Scikit-learn’s SVC implementation) trained as a classifier that takes the tuple embeddings in  $\varphi$  as input, and learns to predict the target of the downstream task. The results are given for a ten-fold cross validation. We fix one split for each dataset and task where this cross validation is performed, and we do so for every evaluated embedding.

The value of  $\alpha(\varphi)$  is our primary metric for the quality of an embedding  $\varphi$ . We usually train five embeddings  $\Phi = \{\varphi_1, \dots, \varphi_5\}$  for each configuration, each with a different seed. We use the mean cross-validated accuracy across the five as a measure of the expected embedding quality:  $\alpha(\Phi) = \text{mean}_{1 \leq i \leq 5} a(\varphi_i)$ . We study how this expectation develops over time for different strategies and selection ratios. More formally, let  $\mathcal{T}$  be a scheme selection strategy and let  $r \in [0, 1]$  be the ratio of schemes used for training. By  $\Phi(\mathcal{T}, r, t)$  we denote the set of embeddings obtained after training five FoRWARD embeddings for  $t$  seconds on the targeted walk schemes reduced using  $\mathcal{T}$  by ratio  $1 - r$ . We evaluate the expected accuracy after each epoch, thus the set of times  $t$  where we record  $\alpha(\Phi(\mathcal{T}, r, t))$  is determined by the time it takes to complete an epoch.

We run all the experiments a server with two Intel Xeon Gold 6130 processors, 512 MB RAM, and an NVIDIA QUADRO RTX 6000 GPU with 24 GB memory.

## 5.2 Performance of Individual Strategies

We first study how the embedding quality develops throughout the training of the embedding. We record the training time and the quality on the downstream task at the end of each epoch. We study the progress when training on different subsets of the targeted walk schemes selected by the different strategies. In Figure 5, we provide results for the *Mondial-Religion* (MR) and *Genes* downstream tasks using the strategies KVar (kernel variance), 1epoch (one epoch), and Random.

In each sub-figure of Figure 5, there are nine colored curves and one gray curve. Each colored curve represents one ratio of removed schemes from 10% to 90%. The gray curve shows the original FoRWARD run when training on all targeted walk schemes. The x-axis provides the training time  $t$  in seconds and we plot the value of  $\alpha(\Phi(\mathcal{T}, r, t))$  (i.e. the evaluation of the expected accuracy  $t$  seconds) on the y-axis. Intuitively, this shows how the embedding quality develops throughout the training of the embedding.

Across the experiments, we observe a range of behaviors. First, the selection strategy has a significant influence on the result. For example, the Random elimination strategy yields significantly worse performance when a large percentage of schemes is removed. This justifies the design of more sophisticated strategies that are able to yield embeddings of high quality even when a larger ratio of  $\text{TWS}(R, \ell_{\max})$  is removed. One example of this is the *Mondial-Religion* downstream task. Vanilla FoRWARD with all schemes reaches 83% accuracy on the downstream task. When using Random to train on only 50% of the walk schemes, accuracy drops to 74% while taking longer to converge. Furthermore, when using the Random strategy to train on 20% of the targeted walk schemes, we observe a more dramatic decrease in accuracy to 71%. In contrast, when training with the KVar strategy on 50% or 20% of the targeted walk schemes, the accuracy over the downstream task does not differ much from the original 83%.

For most of the combinations of tasks and strategies, a significant portion of the targeted walk schemes can be removed with a negligible decrease in quality. The speed of convergence naturally increases with the ratio of removed schemes. An example of this behavior is the *Mondial-Religion* task. Here, when we train on fewer schemes (using the KVar and 1epoch strategies), we reach the same 0.83 accuracy with less training time. For instance, by selecting a fifth of the walk schemes using KVar, we get to the full quality in about a third of the embedding time. Another example is the *Genes* downstream task. Here, the original FoRWARD with all walk schemes takes 300 seconds to reach 92.9% accuracy over the downstream task. This time improves to 124.5 and 148 seconds when training on just 30% of the targeted walk schemes selected by the KVar and 1epoch strategies, respectively.

There are several instances where the removal of schemes not only accelerates training but actually *increases* the downstream performance. An example of this is the *Genes* dataset. Here, the accuracy achieved with all scheme selection strategies increases with the percentage of removed schemes until about 70% of all schemes are discarded. More specifically, the accuracy on the downstream

**Table 3: The best training time for each scheme selection strategy and task. The table contains the shortest time it takes to reach the threshold quality  $\alpha^*$  over all tested removal ratios  $r$  (i.e.  $t^*(\mathcal{T})$ ). The last row (ALL) refers to the embedding time of the original FoRWARD run, with all schemes, as a baseline. The best (lowest) train time for each task is printed in bold.**

	Mutagenesis	World	Hepatitis	Genes	M.-Religion	M.-Infant	M.-Continent	M.-GDP	M.-Inflation
KVar	<b>20.59</b> (+2.67)	<b>154.32</b> (+18.71)	<b>29.31</b> (+6.52)	107.18 (+5.73)	<b>42.75</b> (+2.63)	<b>42.17</b> (+4.61)	54.39 (+1.58)	<b>56.84</b> (+8.20)	<b>49.46</b> (+6.47)
Online	41.89 (+2.68)	<b>150.38</b> (+11.09)	60.88 (+33.53)	115.91 (+4.15)	61.76 (+2.31)	61.96 (+2.48)	69.44 (+2.88)	<b>56.55</b> (+9.86)	65.31 (+4.80)
Random	39.09 (+7.29)	268.97 (+81.67)	69.98 (+16.33)	115.39 (+5.83)	88.22 (+14.12)	76.16 (+0.71)	113.15 (+7.07)	71.02 (+6.00)	82.5 (+35.24)
1epoch	41.57 (+1.21)	191.2 (+31.45)	58.68 (+15.50)	126.46 (+7.12)	55.06 (+1.20)	56.03 (+3.36)	67.16 (+1.37)	65.63 (+4.88)	65.08 (+6.58)
All	55.89 (+8.67)	550.65 (+22.92)	58.36 (+9.3)	161.9 (+6.74)	101.32 (+8.53)	76.16 (+0.71)	113.14 (+7.06)	69.03 (+8.77)	95.37 (+8.61)

task is 93.8% for the original FoRWARD run with the full set of walk schemes; it goes up to 98% when 70% of the schemes are discarded using the KVar and 1epoch strategies. Similar results are obtained on the *Mondial-Inflation* and *World* tasks, although the margin of improvement is different on these tasks.

### 5.3 Comparison between Strategies

We aim to better quantify the difference between strategies. We first introduce additional metrics. We establish a *high-performance threshold* by training five standard FoRWARD embeddings  $\Phi_{\text{FWD}}$  with all targeted walk schemes. We set 95% of the expected cross-validated accuracy as the performance target on each downstream task:  $\alpha^* = 0.95 \cdot \alpha(\Phi_{\text{FWD}})$ . In Figure 5 described next, the threshold is shown as a dashed red line. We wish to study how each considered strategy for scheme selection affects the compute time needed to obtain an embedding of high quality. To this end, we will measure how much training time is needed to reach the threshold  $\alpha^*$  with each strategy. We will use the following definitions.

For a strategy  $\mathcal{T}$  and a ratio  $r$ , we denote by  $t^*(\mathcal{T}, r)$  the earliest time it takes to reach the performance threshold:  $t^*(\mathcal{T}, r) = \min\{t \mid \alpha(\Phi(\mathcal{T}, r, t)) \geq \alpha^*\}$ . For the strategy  $\mathcal{T}$ , we can further define a metric that is the shortest time it takes to reach the target quality over all tested ratios  $r$ :  $t^*(\mathcal{T}) = \min_r t^*(\mathcal{T}, r)$ . This is the primary metric that we use for measuring the effectiveness of a strategy.

Table 3 provides the value of  $t^*(\mathcal{T})$  for each strategy  $\mathcal{T}$  and task. Overall, KVar yields the fastest training times on most tasks. However, the exact performance of each strategy depends on the data. The Length strategy yields the best overall results on the Genes dataset, but KVar performs substantially better than Length on all other tasks. The Online strategy also yields the fastest training times on the *World* and *Modial-GDP* tasks.

To further study training different selection strategies, in Figure 6 we provide the “best” learning curve of each strategy on each task. More specifically, we plot the value of  $\alpha(\Phi(\mathcal{T}, r^*(\mathcal{T}), t))$ , where the ratio  $r^*(\mathcal{T}) \in [0, 1]$  is  $r^*(\mathcal{T}) = \arg \min_r t^*(\mathcal{T}, r)$ . Note that we also provide the training curve of FoRWARD when using all targeted walk schemes (gray). When training with reduced scheme sets the training converges consistently speeds up. As observed

previously, the margin of the speedup depends on the scheme selection strategy. The KVar (kernel Variance) strategy (red) yields the fastest convergence on most downstream tasks.

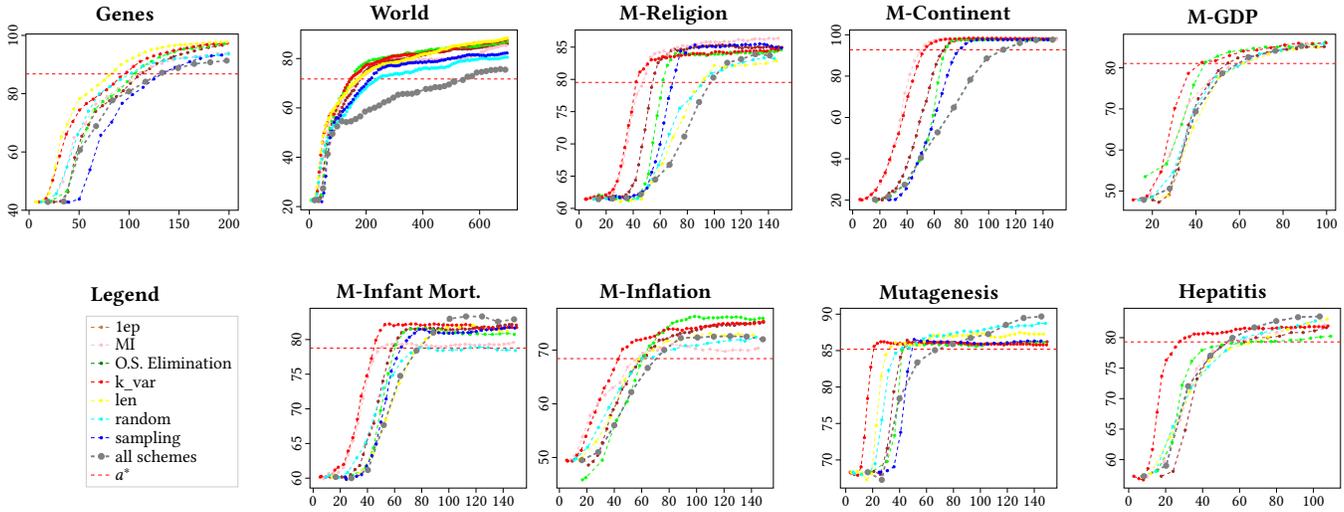
The strong overall performance of the KVar (kernel Variance) method suggests that it strikes the best balance between running efficiently and determining a good subset of schemes to remove. Thus, it enables us to train on less targeted walk schemes and reach the same accuracy faster than simpler selection strategies.

### 5.4 Embedding for Dynamic Databases

Finally, we investigate how the reduction in the set of schemes affects the dynamic-database setting, where new tuples arrive and we need to compute an embedding of the new tuples without changing the embedding of existing tuples. (Note that tuple *deletion* is not an issue in this setting since we simply leave intact the embedding of the remaining tuples; see Tönshoff et al. [22] for a discussion on tuple insertion and its possible subtleties.) For this experiment, we will use the dynamic extension of FoRWARD [22]. This extension is one of the key motivations for the design of FoRWARD as it is able to compute embeddings of new tuples unseen during training by solving a single linear system of equations.

Naturally, a useful scheme-selection strategy should not impair the quality of the new tuple embeddings. Here, we will conduct a brief experiment to verify that this is indeed the case. More specifically, we adopt the exact experimental setup of the dynamic experiments conducted by Tönshoff et al. [22]. This setup first deletes a part of the database and trains an embedding on the remaining data. This embedding is then used to train a classifier for the downstream task. Only after this, the removed tuples are inserted back into the database and the embedding is inductively extended to the new data. These new embeddings can then be evaluated by measuring the accuracy of the downstream classifier on the inserted data. By varying the ratio of data that is initially removed we can test the robustness of this dynamic embedding extension.

We train and then extend FoRWARD where 60% of all targeted walk schemes have been removed according to the KVar strategy. As baselines, we include the results of FoRWARD with no schemes removed as well as the dynamic extension of node2vec, which was



**Figure 6: Comparison between different strategies. For each downstream task, we provide the best learning curve of each strategy, that is, the first curve that reaches an accuracy of  $\alpha^*$ . Note that the color now denotes the selection strategy.**

also proposed by Tönshoff et al. [22]. Figure 7 provides the results on four tasks: Genes, World, Mondial Religion and Mutagenesis. We observe that the FoRWARD version with the reduced set of schemes performs as well as the original FoRWARD with all schemes. On the Genes and World tasks the results actually *improve* slightly when only the selected 40% of targeted walk schemes are used. There is, though, a slight decrease in the accuracy in the case of the Mutagenesis dataset.

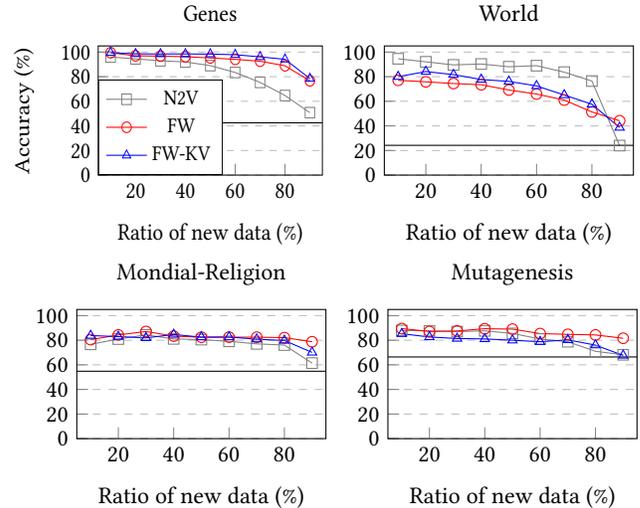
Overall, we conclude that reduced scheme sets with a strong performance on a static database are also well-suited for a dynamic setting where new tuples are inserted over time.

## 6 CONCLUDING REMARKS

Walking through connected data items is the basis of sequence-based embeddings like WORD2VEC and NODE2VEC. Database sequences are also associated with meta-data, namely the walk scheme. The premise of this work was the conjecture that the walk schemes have significant semantic value, as they can guide the embedding algorithm to a small subset of informative sequences, thus dramatically improving efficiency for a mild sacrifice of quality. We studied the problem of selecting walk schemes in the context of FoRWARD. We considered different strategies of three types: FoRWARD-less, light training, and online scheme elimination. We conducted an experimental study that measured the benefit of each strategy, compared between them and tested how well they preserved the main strength of FoRWARD—extensibility to newly inserted tuples in dynamic settings. Our study has confirmed our conjecture and showed that we can considerably accelerate FoRWARD with negligent loss of quality. Moreover, restricting the embedding phase to the right walk schemes can even *improve* the quality on downstream classification tasks. The kernel-variance strategy typically outperforms the rest, and we recommend this one to be used alongside FoRWARD.

The idea of directing the embedding algorithm to beneficial walk schemes goes beyond FoRWARD and is applicable to every database

embedding we are aware of. Our experience with EMBDI [6] has indicated that EMBDI uses a large number of walk schemes with tiny pairwise differences and a heavy-tailed distribution of a number of instances. Thus we need to use ways of abstracting (clustering) walk schemes. We also plan to expand the scope of our work to data integration tasks (e.g., entity and schema matching) as done with other database embedding algorithms [6, 19]. In the case of FoRWARD, there is a need to devise *alignment* between embeddings of different databases, since FoRWARD makes no attempt to produce similar embeddings to matching entities of different databases.



**Figure 7: Experiments on the dynamic setting: accuracy as a function of the percentage of inserted tuples for Node2Vec, FoRWARD, and FoRWARD with 60% of the schemes selected by KVar. The black line is the accuracy of the common class.**

## REFERENCES

- [1] Bahareh Bina, Oliver Schulte, Branden Crawford, Zhensong Qian, and Yi Xiong. 2013. Simple decision forests for multi-relational classification. *Decis. Support Syst.* 54, 3 (2013), 1269–1279. <https://doi.org/10.1016/j.dss.2012.11.017>
- [2] Rajesh Bordawekar, Bortik Bandyopadhyay, and Oded Shmueli. 2017. Cognitive Database: A Step towards Endowing Relational Databases with Artificial Intelligence Capabilities. *CoRR* abs/1712.07199 (2017). arXiv:1712.07199 <http://arxiv.org/abs/1712.07199>
- [3] Rajesh Bordawekar and Oded Shmueli. 2017. Using Word Embedding to Enable Semantic Queries in Relational Databases. In *Proceedings of the 1st Workshop on Data Management for End-to-End Machine Learning* (Chicago, IL, USA) (DEEM'17). ACM, New York, NY, USA, Article 5, 4 pages. <https://doi.org/10.1145/3076246.3076251>
- [4] Rajesh Bordawekar and Oded Shmueli. 2019. Exploiting Latent Information in Relational Databases via Word Embedding and Application to Degrees of Disclosure. In *CIDR*. [www.cidrdb.org. http://cidrdb.org/cidr2019/papers/p27-bordawekar-cidr19.pdf](http://cidrdb.org/cidr2019/papers/p27-bordawekar-cidr19.pdf)
- [5] A. Borde, N. Usunier, A. Garcia-Duran, J. Weston, and O. Yakhnenko. 2013. Translating embeddings for modeling multi-relational data. In *Advances in neural information processing systems*. 2787–2795.
- [6] Riccardo Cappuzzo, Paolo Papotti, and Saravanan Thirumuruganathan. 2020. Creating Embeddings of Heterogeneous Relational Datasets for Data Integration Tasks. In *SIGMOD*, David Maier, Rachel Pottinger, AnHai Doan, Wang-Chiew Tan, Abdussalam Alawini, and Hung Q. Ngo (Eds.). ACM, 1335–1349. <https://doi.org/10.1145/3318464.3389742>
- [7] Jie Cheng, Christos Hatzis, Hisashi Hayashi, Mark-A. Krogel, Shinichi Morishita, David Page, and Jun Sese. 2002. KDD Cup 2001 Report. *SIGKDD Explor. Newsl.* 3, 2 (Jan. 2002), 47–64. <https://doi.org/10.1145/507515.507523>
- [8] Asim Kumar Debnath, Rosa L. Lopez de Compadre, Gargi Debnath, Alan J. Shusterman, and Corwin Hansch. 1991. Structure-activity relationship of mutagenic aromatic and heteroaromatic nitro compounds. Correlation with molecular orbital energies and hydrophobicity. *Journal of Medicinal Chemistry* 34, 2 (1991), 786–797. <https://doi.org/10.1021/jm00106a046> arXiv:<https://doi.org/10.1021/jm00106a046>
- [9] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *NAACL-HLT (1)*. Association for Computational Linguistics, 4171–4186.
- [10] Muhammad Ebraheem, Saravanan Thirumuruganathan, Shafiq R. Joty, Mourad Ouzzani, and Nan Tang. 2018. Distributed Representations of Tuples for Entity Resolution. *Proc. VLDB Endow.* 11, 11 (2018), 1454–1467.
- [11] Aditya Grover and Jure Leskovec. 2016. Node2vec: Scalable Feature Learning for Networks. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (San Francisco, California, USA) (KDD '16). ACM, New York, NY, USA, 855–864. <https://doi.org/10.1145/2939672.2939754>
- [12] Michael Günther. 2018. FREDDY: Fast Word Embeddings in Database Systems. In *SIGMOD Conference*. ACM, 1817–1819.
- [13] Michael Günther, Maik Thiele, Erik Nikulski, and Wolfgang Lehner. 2020. Retro-Live: Analysis of Relational Retrofitted Word Embeddings. In *EDBT*. OpenProceedings.org, 607–610.
- [14] Sabrina Jaeger, Simone Fulle, and Samo Turk. 2018. Mol2vec: Unsupervised Machine Learning Approach with Chemical Intuition. *J. Chem. Inf. Model.* 58, 1 (2018), 27–35.
- [15] Quoc V. Le and Tomáš Mikolov. 2014. Distributed Representations of Sentences and Documents. In *ICML (JMLR Workshop and Conference Proceedings)*, Vol. 32. JMLR.org, 1188–1196.
- [16] Wolfgang May. 1999. *Information extraction and integration: The Mondial case study*. Technical Report. Universität Freiburg, Institut für Informatik.
- [17] T. Mikolov, I. Sutskever, K. Chen, G.S. Corrado, and J. Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Proceedings of the 27th Annual Conference on Neural Information Processing Systems*. 3111–3119.
- [18] Jan Motl. 2020. *The World Dataset*. <https://relational.fit.cvut.cz/dataset/World>
- [19] Sidharth Mudgal, Han Li, Theodoros Rekatsinas, AnHai Doan, Youngchoon Park, Ganesh Krishnan, Rohit Deep, Esteban Arcaute, and Vijay Raghavendra. 2018. Deep Learning for Entity Matching: A Design Space Exploration. In *SIGMOD Conference*. ACM, 19–34.
- [20] Jennifer Neville, David D. Jensen, Lisa Friedland, and Michael Hay. 2003. Learning relational probability trees. In *KDD*. ACM, 625–630.
- [21] Ben Taskar, Pieter Abbeel, and Daphne Koller. 2002. Discriminative Probabilistic Models for Relational Data. In *Proceedings of the Eighteenth Conference on Uncertainty in Artificial Intelligence* (Alberta, Canada) (UAI'02). Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 485–492.
- [22] Jan Tönshoff, Neta Friedman, Martin Grohe, and Benny Kimelfeld. 2023. Stable Tuple Embeddings for Dynamic Databases. In *ICDE*. IEEE.