

# A Temporal Knowledge Graph Embedding Model based on Additive Time Series Decomposition

No Author Given

No Institute Given

**Abstract.** Knowledge Graph (KG) embedding has attracted more attention in recent years. Most KG embedding models learn from time-unaware triples. However, the inclusion of temporal information besides triples would further improve the performance of a KGE model. In this regard, we propose **ATiSE**, a temporal KG embedding model which incorporates time information into entity/relation representations by using **Additive Time Series** decomposition. Moreover, considering the temporal uncertainty during the evolution of entity/relation representations over time, we map the representations of temporal KGs into the space of multi-dimensional Gaussian distributions. The mean of each entity/relation embedding at a time step shows the current expected position, whereas its covariance (which is temporally stationary) represents its temporal uncertainty. Experimental results show that ATiSE significantly outperforms the state-of-the-art KGE models and the existing temporal KGE models on link prediction over four temporal KGs.

**Keywords:** First keyword · Second keyword · Another keyword.

## 1 Introduction

Knowledge Graphs (KGs) are being used for gathering and organizing scattered human knowledge into structured knowledge systems. YAGO [21], DBpedia [1], WordNet [17] and Freebase [2] are among existing KGs that have been successfully used in various applications including question answering, assistant systems, information retrieval, etc. In these KGs, knowledge can be represented as RDF triples  $(s, p, o)$  in which  $s$  (subject) and  $o$  (object) are entities (nodes), and  $p$  (predicate) is the relation (edge) between them.

KG embedding attempts to learn the representations of entities and relations in high-dimensional latent feature spaces while preserving certain properties of the original graph. Recently, KG embedding has become a very active research topic due to the wide ranges of downstream applications. Different KG embedding models have been proposed so far to efficiently learn the representations of KGs and perform KG completion as well as inferencing [3, 8, 27, 24, 22, 29].

We notice that most of existing KG embedding models solely learn from time-unknown facts and ignore the useful temporal information in the KBs. In fact,

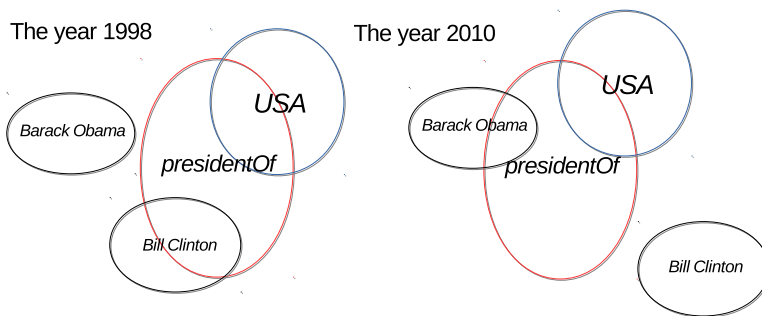
there are many time-aware facts (or events) in some temporal KBs. For example, (*Obama, wasBornIn, Hawaii*) happened at August 4, 1961. (*Obama, presidentOf, USA*) was true from 2009 to 2017. These temporal KGs, e.g. Integrated Crisis Early Warning System (ICEWS) [13], Global Database of Events, Language, and Tone (GDELT) [15], YAGO3 [16] and Wikidata [5], store such temporal information either explicitly or implicitly. Traditional KBE models such as TransE learn only from time-unknown facts. Therefore, they cannot distinguish entities with similar semantic meaning. For instance, they often confuse entities such as *Barack Obama* and *Bill Clinton* when predicting (*?, presidentOf, USA, 2010*).

To tackle this problem, temporal KGE models [14, 4, 6] encode time information in their embeddings. TKGE models outperform traditional KGE models on link prediction over temporal KGs. It justifies that incorporation of time information can further improve the performance of a KGE model. Most existing TKGE models embed time information into a latent space, e.g. representing time as a vector. These models cannot capture some properties of time information such as the length of time interval as well as order of two time points. Moreover, these models ignore the uncertainty during the temporal evolution. We argue that the evolution of entity representations has randomness, because the features of an entity at a certain time are not completely determined by the past information. For example, (*Steve Jobs, diedIn, California*) happened on 2011-10-05. The semantic characteristics of this entity should have a sudden change at this time point. However, due to the incompleteness of knowledge in KGs, this change can not be predicted only according to its past evolutionary trend. Therefore, the representation of *Steve Jobs* is supposed to include some random components to handle this uncertainty, e.g. a Gaussian noise component.

In order to address the above problems, in this paper, we propose a temporal KG embedding model, ATiSE, which uses additive time series decomposition to capture the evolution process of KG representations. ATiSE fits the evolution process of an entity or relation as a multi-dimensional additive time series which composes of a trend component, a seasonal component and a random component. Our approach represents each entity and relation as a multi-dimensional Gaussian distribution at each time step to introduce a random component. The mean of an entity/relation representation at a certain time step indicates its current expected position, which is obtained from its initial representation, its linear change term, and its seasonality term. The covariance which describes the temporal uncertainty during its evolution, is denoted as a constant diagonal matrix for computing efficiency. Our contributions are as follows.

- Learning the representations for temporal KGs is a relatively unexplored problem because most existing KG embedding models only learn from time-unknown facts. We propose ATiSE, a new KG embedding model to incorporate time information into the KG representations.
- We specially consider the temporal uncertainty during the evolution process of KG representations. Thus, we model each entity/relation as a Gaussian distribution at each time step and use KL-divergence between two Gaussian distributions to compute the scores of facts for optimization.

- Different from the previous temporal KG embedding models which use time embedding to incorporate time information, ATiSE fits the evolution process of KG representations as a multi-dimensional additive time series. Our work establishes a previously unexplored connection between relational processes and time series analysis with a potential to open a new direction of research on reasoning over time.
- Our experimental results show that ATiSE significantly outperforms other TKG models and some state-of-the-art static KGE on link prediction over four TKG datasets.



**Fig. 1.** Illustration of the means and (diagonal) variances of entities and relations in a temporal Gaussian Embedding Space. The labels indicate their position. In the representations, we might infer that *Bill Clinton* was *presidentOf* *USA* in 1998 and *Barack Obama* was *presidentOf* *USA* in 2010.

The rest of the paper is organized as follows: In the section 2, we first review related works; in the section 3, we introduce the architecture and the learning process of our proposed models; in the section 4, we compare the performance of our models with the state-of-the-art models; in the section 5, we make a conclusion in the end of this paper.

## 2 RELATED WORK

A large amount of research has been done in KG embeddings. These approaches can generally be categorized into two groups, namely transnational distance models and semantic matching models [25].

A few examples of translational distance models include TransE [3], TransH [26], TransD [10]. These models measure the plausibility of a fact as the distance between the two entities, usually after a translation carried out by the relation. In addition, RotatE [22] achieves the state-of-the-art results on link prediction by using relational rotations in complex space instead of relational translations.

RESCAL [19] and its extensions, e.g. DistMult [27], ComplEx [24], QuatE [29], are semantic matching models. These models measure plausibility of facts by

matching latent semantics of entities and relations embodied in their vector space representations. Specially, ComplEx-N3 [12] adopts N3 regularization and reciprocal learning to remarkably boost the performance of ComplEx.

The above methods achieve good results on link prediction in KGs. However, these time-unaware KGE models have limitations on reasoning over TKGs. More concretely, given two quadruples with the same subjects, predicates, objects and different time stamps, i.e., *(Barack Obama, presidentOf, USA, 2010)* and *(Barack Obama, presidentOf, USA, 2020)*, static KGE models will model them with the same scores due to their ignorance of time information, while the validities of these two quadruples might be different.

Recent researches illustrate that the performances of KG embedding models can be further improved by incorporating time information in temporal KGs.

TTransE [14] and HyTE [4] adopt translational distance score functions and encode time information in the entity-relation low dimensional spaces with time embeddings and temporal hyperplanes.

Know-Evolve [23] models the occurrence of a fact as a temporal point process. However, this method is built on a problematic formulation when dealing with concurrent events, as shown in Section 4.3.

TA-TransE and TA-DistMult[6] utilize recurrent neural networks to learn time-aware representations of relations and use standard scoring functions from TransE and DistMult. These models can model time information in the form of time points with or without some particular temporal modifiers, i.e., '*occursSince*' and '*occursUntil*'.

DE-Simple [7] incorporates time information into diachronic entity embeddings and achieves the state of the art results on event-based TKGs. However, same as TA-TransE and TA-DistMult, DE-Simple can not model facts involving time intervals shaped like [2005, 2008].

### 3 OUR METHOD

In this section, we present a detailed description of our proposed method, ATiSE, which not only uses relational properties between entities in triples but also incorporates the associated temporal meta-data by using additive time series decomposition.

#### 3.1 Additive Time Series Embedding Model

A time series is a series of time-oriented data. Time series analysis is widely used in many fields, ranging from economics and finance to managing production operations, to the analysis of political and social policy sessions [18]. An important technique for time series analysis is additive time series decomposition. This technique decomposes a time series into three components, i.e., a trend component, a seasonal component and an irregular component (i.e. "noise").

In our method, we regard the evolution of an entity/relation representation as an additive time series. For each entity/relation, we use a linear function

and a Sine function to fit the trend component and the seasonal component respectively due to their simplicity. Considering the efficiency of model training, we model the irregular term by using a Gaussian noise instead of a moving average model (MA model) [9], since training an MA model requires a global optimization algorithm which will lead to more computation consumption.

To incorporate temporal information into traditional KGs, a new temporal dimension is added to fact triples, denoted as a quadruple  $(s, p, o, t)$ . It represents the creation of relationship edge  $p$  between subject entity  $s$ , and object entity  $o$  at time step  $t$ . The score term  $x_{spot} = f_t(e_s, r_p, e_o)$  can represent the conditional probability or the confidence value of this event  $x_{spot}$ , where  $e_s, e_o \in \mathbf{R}^{L_e}$ ,  $r_p \in \mathbf{R}^{L_r}$  are representations of  $s$ ,  $o$  and  $p$ . In term of a long-term fact  $(s, p, o, [t_s, t_e])$ , we consider it to be a positive triple for each time step between  $t_s$  and  $t_e$ .  $t_s$  and  $t_e$  denote the start and end time during which the triple  $(s, p, o)$  is valid.

At each time step, the time-specific representations of an entity  $e_i$  or a relation  $r_p$  should be updated as  $e_{i,t}$  or  $r_{p,t}$ . Thus, the score of a quadruple  $(s, p, o, t)$  can be represented as  $x_{spot} = f_e(e_{s,t}, r_{p,t}, e_{o,t})$  or  $x_{spot} = f_r(e_s, r_{p,t}, e_o)$ . We utilize additive time series decomposition to fit the evolution processes of each entity/relation representation as:

$$\begin{aligned} e_{i,t} &= e_i + \alpha_{e,i} w_{e,i} t + \beta_{e,i} \sin(2\pi \omega_{e,i} t) + \mathcal{N}(0, \Sigma_{e,i}) \\ r_{p,t} &= r_p + \alpha_{r,p} w_{r,p} t + \beta_{r,p} \sin(2\pi \omega_{r,p} t) + \mathcal{N}(0, \Sigma_{r,p}) \end{aligned} \quad (1)$$

where the  $e_i$  and  $r_p$  are the time-independent latent representations of the  $i$ th entity which is subjected to  $\|e_i\|_2 = 1$  and the  $p$ th relation which is subjected to  $\|r_p\|_2 = 1$ .  $e_i + \alpha_{e,i} w_{e,i} t$  and  $r_p + \alpha_{r,p} w_{r,p} t$  are the trend components where the coefficients  $|\alpha_{e,i}|$  and  $|\alpha_{r,p}|$  denote the evolutionary rates of  $e_{i,t}$  and  $r_{p,t}$ , the vectors  $w_{e,i}$  and  $w_{r,p}$  represents the corresponding evolutionary directions which are restricted to  $\|w_{e,i}\|_2 = \|w_{r,p}\|_2 = 1$ .  $\beta_{e,i} \sin(2\pi \omega_{e,i} t)$  and  $\beta_{r,p} \sin(2\pi \omega_{r,p} t)$  are the corresponding seasonal components where  $|\beta_{e,i}|$  and  $|\beta_{r,p}|$  denote the amplitude vectors,  $|\omega_{e,i}|$  and  $|\omega_{r,p}|$  denote the frequency vectors. The Gaussian noise terms  $\mathcal{N}(0, \Sigma_{e,i})$  and  $\mathcal{N}(0, \Sigma_{r,p})$  are the random components, where  $\Sigma_{e,i}$  and  $\Sigma_{r,p}$  denote the corresponding diagonal covariance matrices.

In other words, for a fact  $(s, p, o, t)$ , entity embeddings  $e_{s,t}$  and  $e_{o,t}$  obey Gaussian probability distributions:  $\mathcal{P}_{s,t} \sim \mathcal{N}(\bar{e}_{s,t}, \Sigma_s)$  and  $\mathcal{P}_{o,t} \sim \mathcal{N}(\bar{e}_{o,t}, \Sigma_o)$ , where  $\bar{e}_{s,t}$  and  $\bar{e}_{o,t}$  are the mean vectors of  $e_{s,t}$  and  $e_{o,t}$ , which do not include the random components. Similarly, the predicate is represented as  $\mathcal{P}_{r,t} \sim \mathcal{N}(r_p, \Sigma_r)$ .

Similar to translation-based KGE models, we consider the transformation result of ATiSE from the subject to the object to be akin to the predicate in a positive fact. We use the following formula to express this transformation:  $\mathcal{P}_{s,t} - \mathcal{P}_{o,t}$ , which corresponds to the probability distribution  $\mathcal{P}_{e,t} \sim \mathcal{N}(\mu_{e,t}, \Sigma_e)$ . Here,  $\mu_{e,t} = \bar{e}_{s,t} - \bar{e}_{o,t}$  and  $\Sigma_e = \Sigma_s + \Sigma_o$ . Combined with the probability of relation  $\mathcal{P}_{r,t} \sim \mathcal{N}(r_p, \Sigma_r)$ , we measure the similarity between  $\mathcal{P}_{e,t}$  and  $\mathcal{P}_r$  to score the fact.

KL divergence is a straightforward method of measuring the similarity of two probability distributions. We optimize the following score function based on the KL divergence between the entity-transformed distribution and relation

distribution [28].

$$\begin{aligned}
x_{spot} &= f_t(e_s, r_p, e_o) = \mathcal{D}_{\mathcal{KL}}(P_{r,t}, P_{e,t}) \\
&= \int_{x \in \mathcal{R}^{k_e}} \mathcal{N}(x; r_{p,t}, \Sigma_r) \log \frac{\mathcal{N}(x; \mu_{e,t}, \Sigma_e)}{\mathcal{N}(x; r_{p,t}, \Sigma_r)} dx \\
&= \frac{1}{2} \left\{ \text{tr}(\Sigma_r^{-1} \Sigma_e) + (r_{p,t} - \mu_{e,t})^T \Sigma_r^{-1} (r_{p,t} - \mu_{e,t}) \right. \\
&\quad \left. - \log \frac{\det(\Sigma_e)}{\det(\Sigma_r)} - k_e \right\}
\end{aligned} \tag{2}$$

where,  $\text{tr}(\Sigma)$  and  $\Sigma^{-1}$  indicate the trace and inverse of the diagonal covariance matrix, respectively.

Considering the simplified diagonal covariance, we can compute the trace and inverse of the matrix simply and effectively for ATiSE. The gradient of log determinant is  $\frac{\partial \log \det A}{\partial A} = A^{-1}$ , the gradient  $\frac{\partial x^T A^{-1} y}{\partial A} = -A^{-T} x y^T A^{-T}$ , and the gradient  $\frac{\partial \text{tr}(X^T A^{-1} Y)}{\partial A} = -(A^{-1} Y X^T A^{-1})^T$  [20]. We can compute the gradients of Equation 2 with respect to the time-independent latent feature vectors, evolutionary direction vectors and covariance matrix (here acting as a vector) as follows:

$$\begin{aligned}
\frac{\partial x_{spot}}{\partial \alpha_s} &= -t w_s \Delta'_{spot} & \frac{\partial x_{spot}}{\partial \alpha_o} &= t w_o \Delta'_{spot} & \frac{\partial x_{spot}}{\partial \alpha_r} &= t w_r \Delta'_{spot} \\
\frac{\partial x_{spot}}{\partial w_s} &= -t \alpha_s \Delta'_{spot} & \frac{\partial x_{spot}}{\partial w_o} &= t \alpha_o \Delta'_{spot} & \frac{\partial x_{spot}}{\partial w_r} &= t \alpha_r \Delta'_{spot} \\
\frac{\partial x_{spot}}{\partial \beta_s} &= -\sin(2\pi \omega_s t) \Delta'_{spot} \\
\frac{\partial x_{spot}}{\partial \beta_o} &= \sin(2\pi \omega_o t) \Delta'_{spot} \\
\frac{\partial x_{spot}}{\partial \beta_r} &= \sin(2\pi \omega_r t) \Delta'_{spot} \\
\frac{\partial x_{spot}}{\partial \omega_s} &= -\beta_s \cos(2\pi \omega_s t) \Delta'_{spot} \\
\frac{\partial x_{spot}}{\partial \omega_o} &= \beta_o \cos(2\pi \omega_o t) \Delta'_{spot} \\
\frac{\partial x_{spot}}{\partial \omega_r} &= \beta_r \cos(2\pi \omega_r t) \Delta'_{spot} \\
\frac{\partial x_{spot}}{\partial \Sigma_r} &= \frac{1}{2} (\Sigma_r^{-1} \Sigma_e \Sigma_r^{-1} + \Delta'_{spot} \Delta_{spot}'^T + \Sigma_r^{-1}) \\
\frac{\partial x_{spot}}{\partial \Sigma_s} &= \frac{\partial x_{spot}}{\partial \Sigma_o} = \frac{1}{2} (\Sigma_r^{-1} - \Sigma_e^{-1})
\end{aligned} \tag{3}$$

where  $\Delta'_{spot} = \Sigma_r^{-1}(r_p + e_o - e_s + t(\alpha_r w_r + \alpha_o w_o - \alpha_s w_s) + \beta_r \sin(2\pi \omega_r t) + \beta_o \sin(2\pi \omega_o t) - \beta_s \sin(2\pi \omega_s t))$ ,  $\Sigma_e = \Sigma_s + \Sigma_o$ .

### 3.2 Complexity

In Table 1, we summarize the scoring functions of several existing (T)KGE approaches and our models and compare their space complexities.  $n_e$ ,  $n_r$ ,  $n_t$  and  $n_{token}$  are numbers of entities, relations, time steps and temporal tokens used in [6];  $d$  is the dimensionality of embeddings.  $\langle x, y, z \rangle = \sum_i x_i y_i z_i$  denotes the tri-linear dot product;  $\text{RE}(\cdot)$  denotes the real part of the complex embedding [24];  $\otimes$  denotes the Hamilton product between quaternion embeddings;  $\lhd$  denotes the normalization of the quaternion embedding.  $\mathcal{P}_t$  denotes the temporal projection for embeddings [29];  $\text{LSTM}(\cdot)$  denotes an LSTM neural network;  $[r_p; t_{seq}]$  denotes the concatenation of the relation embedding and the sequence of temporal tokens [6];  $\vec{e}$  and  $\overleftarrow{e}$  denote the temporal part and untemporal part of a time-specific diachronic entity embedding  $e^t$  [7];  $p^{-1}$  denotes the inverse relation of  $p$ , i.e.,  $(s, p, o, t) \leftrightarrow (o, p^{-1}, s, t)$ .

As shown in Table 3.2, our models have the same space complexities as static KGE models listed in Table 3.2 as well as DE-Simple. On the other hand, the space complexities of TTransE, HyTE, TA-TransE or TA-DistMult will be higher than our models if  $n_t$  or  $n_{token}$  is much larger than  $n_e$  and  $n_r$ .

**Table 1.** Comparison of our models with several baseline models for space complexity.

Model	Scoring Function	Space Complexity
TransE	$\ e_s + r_p - e_o\ $	$\mathcal{O}(n_e d + n_r d)$
DistMult	$\langle e_s, r_p, e_o \rangle$	$\mathcal{O}(n_e d + n_r d)$
ComplEx	$\text{RE}(\langle e_s, r_p, \bar{e}_o \rangle)$	$\mathcal{O}(n_e d + n_r d)$
RotatE	$\ e_s \circ r_p - e_o\ $	$\mathcal{O}(n_e d + n_r d)$
QuatE	$e_s \otimes r_p^{\lhd} \cdot e_o$	$\mathcal{O}(n_e d + n_r d)$
TTransE	$\ e_s + r_p + w_t - e_o\ $	$\mathcal{O}(n_e d + n_r d + n_t d)$
HyTE	$\ P_t(e_s) + P_t(r_p) - P_t(e_o)\ $	$\mathcal{O}(n_e d + n_r d + n_t d)$
TA-TransE	$\ e_s + \text{LSTM}([r_p; t_{seq}]) - e_o\ $	$\mathcal{O}(n_e d + n_r d + n_{token} d)$
TA-DistMult	$\langle e_s, \text{LSTM}([r_p; t_{seq}]), e_o \rangle$	$\mathcal{O}(n_e d + n_r d + n_{token} d)$
DE-Simple	$\frac{1}{2}(\langle \vec{e}_s^t, r_p, \overleftarrow{e}_o^t \rangle + \langle \vec{e}_{0,t}^t, r_{p^{-1}}, \overleftarrow{e}_s^t \rangle)$	$\mathcal{O}(n_e d + n_r d)$
ATiSE	$\mathcal{D}_{KL}(\mathcal{P}_{e,t}, \mathcal{P}_{r,t})$	$\mathcal{O}(n_e d + n_r d)$

### 3.3 Learning

In this paper, we use the same loss function as the negative sampling loss proposed in [22] for optimizing ATiSE. This loss function has been proved to be more effective than the margin rank loss function proposed in [3] on optimizing translation-based KGE models.

$$\mathcal{L} = \sum_{t \in [T]} \sum_{\xi \in \mathcal{D}_t^+} -\log \sigma(\gamma - f_t(\xi)) - \log \sigma(f_t(\xi') - \gamma) \quad (4)$$

where,  $[T]$  is the set of time steps in the temporal KG,  $\mathcal{D}_t^+$  is the set of positive triples with time stamp  $t$ , and  $\mathcal{D}_t^-$  is the set of negative sample corresponding to  $\mathcal{D}_t^+$ . In this paper, we generate negative samples by randomly corrupting subjects or objects of the positives such as  $(s', p, o, t)$  and  $(s, p, o', t)$ . Moreover, we adopt

self-adversarial training proposed in [22] and reciprocal learning used in [12, 7, 29] to further enhance the performances of our model. To avoid overfitting, we add some regularizations while learning ATiSE. As described in the section 3.1, the norms of the original representations of entities and relations, as well as the norms of all evolutionary direction vectors, are restricted by 1. Besides, the following constraint is used for guaranteeing that the covariance matrices are positive definite and of appropriate size when we minimize the loss:

$$\forall l \in \mathcal{E} \cup \mathcal{R}, c_{min}I \leq \Sigma_l \leq c_{max}I \quad (5)$$

where,  $\mathcal{E}$  and  $\mathcal{R}$  are the set of entities and relations respectively,  $c_{min}$  and  $c_{max}$  are two positive constants. We use  $\Sigma_l \leftarrow \max(c_{min}, \min(c_{max}, \Sigma_l))$  to achieve this regularization for diagonal covariance matrices. This constraint 5 for the covariance is considered during both the initialization and training process.

## 4 Experiment

To show the capability of ATiSE, we compared it with some state-of-the-art KGE models and the existing TKGE models on link prediction over four TKG datasets. Particularly, we also did an ablation study to analyze the effect of the dimensionality of entity/relation embeddings and various components of the additive time series decomposition.

### 4.1 Datasets

As mentioned in section 1, common TKGs include ICEWS [13], Wikidata [5] and YAGO3 [16]. Four subsets of these TKGs are used as datasets in [6], i.e., ICEWS14, ICEWS05-15, YAGO15k and Wikidata11k. However, all of time intervals in YAGO15k and Wikidata11k only contain either start dates or end dates, shaped like '*occursSince 2003*' or '*occursUntil 2005*' while most of time intervals in Wikidata and YAGO are presented by both start dates and end dates. Thus, we prefer using YAGO11k and Wikidata12k released in [4] instead of YAGO15k and Wikidata12k. The statistics of the datasets used in this paper are listed in Table 2. ICEWS is a repository that contains political events

**Table 2.** Statistics of datasets.

	#Entities	#Relations	#Time Steps	Time Span	#Training	#Validation	#Test
ICEWS14	6,869	230	365	2014	72,826	8,941	8,963
ICEWS05-15	10,094	251	4,017	2005-2015	368,962	46,275	46,092
YAGO11k	10,623	10	70	-453-2844	16,408	2,050	2,051
Wikidata12k	12,554	24	81	1709-2018	32,497	4,062	4,062

with specific time annotations, e.g., (*Barack Obama, visits, Ukraine, 2014-07-08*). ICEWS14 and ICEWS05-15 are subsets of ICEWS [13], which correspond to the facts in 2014 and the facts between 2005 to 2015. These two datasets are



filtered by only selecting the most frequently occurring entities in the graph [6]. It is noteworthy that all of time annotations in ICEWS datasets are time points.

YAGO11k is a subset of YAGO3 [16]. Different from ICEWS, a part of time annotations in YAGO3 are represented as time intervals, e.g. (*Paul Konchesky, playsFor, England national football team, [2003-##-##, 2005-##-##]*). Following the setting used in HyTE [4], we only deal with year level granularity by dropping the month and date information and treat timestamps as 70 different time steps in the consideration of the balance about numbers of triples in different time steps. For a time interval with the missing start date or end date, e.g., [2003-##-##, ####-##-##] representing 'since 2003', we use the first timestep or the last timestep to represent the missing start time or end time.

Wikidata12k is a subset of Wikidata [5]. Similar to YAGO11k, Wikidata12k contains some facts involving time intervals. We treat timestamps as 81 different time steps by using the same setting as YAGO11k.

For TKGE models, we discretized facts  $(s, p, o, [t_s, t_e])$  involving multiple timesteps into multiple quadruples which only involve single timesteps, i.e.,  $\{(s, p, o, t_s), (s, p, o, t_{s+1}), \dots, (s, p, o, t_e)\}$ , where  $t_s$  and  $t_e$  denote the start time and the end time.

## 4.2 Evaluation Metrics

We evaluate our model by testing the performances of our model on link prediction task over TKGs. This task is to complete a time-wise fact with a missing entity. For a test quadruple  $(s, p, o, t)$ , we generate corrupted triples by replacing  $s$  or  $o$  with all possible entities. We sort scores of all the quadruples including corrupted quadruples and the test quadruples and obtain the ranks of the test quadruples. For a test fact involving multiple time steps, e.g.,  $(s, p, o, [t_s, t_e])$ , the score of one corrupted fact  $(s, p, o', [t_s, t_e])$  is the sum of scores of multiple discrete quadruples,  $\{(s, p, o', t_s), (s, p, o', t_{s+1}), \dots, (s, p, o', t_e)\}$ .

Two evaluation metrics are used here, i.e., Mean Reciprocal Rank and Hits@k. The Mean Reciprocal Rank (MRR) is the means of the reciprocal values of all computed ranks. And the fraction of test quadruples ranking in the top  $k$  is called Hits@k. We adopt the time-wise filtered setting used in source code released by [7]. Different from the original filtered setting proposed in [3], for a test fact  $(s, p, o, t)$  or  $(s, p, o, [t_s, t_e])$ , instead of removing all the triples that appear either in the training, validation or test set from the list of corrupted facts, we only filter the triples that occur at the time point  $t$  or throughout the time interval  $[t_s, t_e]$  from the list of corrupted facts. This ensures that the facts that do not appear at  $t$  or throughout  $[t_s, t_e]$  are still considered as corrupted triplets for evaluating the given test fact.

## 4.3 Baselines

We compare our approach with several state-of-the-art KGE approaches and existing TKGE approaches, including TransE [3], DistMult [27], ComplEx-N3 [12],

RotatE [22], QuatE<sup>3</sup> [29], TTransE [14], TA-TransE, TA-DistMult [6] and DE-Simple [7]. ComplEx-N3 has been proven to have better performance than ComplEx [24] on FreeBase and WordNet datasets. And QuatE<sup>2</sup> has the best performances among all variants of QuatE as reported in [29].

As mentioned in Section 2, TA-TransE, TA-DistMult and DE-Simple mainly focus on modeling temporal facts involving time points with or without some particular temporal modifiers, '*occursSince*' and '*occursUntil*', and cannot model time intervals shaped like [2003-##-##, 2005-##-##]. Besides, DE-Simple needs specific date information including year, month and day to score temporal facts, while most of time annotations in YAGO and Wikidataset only contain year-level information. Thus, we cannot test these three models on YAGO11k and Wikidataset15k.

We do not take Know-Evolve [23] as baseline model due to its problematic formulation and implementation issues. Know-Evolve uses the temporal point process to model the temporal evolution of each entity. The intensity function of Know-Evolve (equation 3 in [23]) is defined as  $\lambda_r^{s,o}(t|\bar{t}) = f(g_r^{s,o}(\bar{t}))(t - \bar{t})$ , where  $g(\cdot)$  is a score function,  $t$  is current time, and  $\bar{t}$  is the most recent time point when either subject or object entity was involved in an event. This intensity function is used in inference to rank entity candidates. However, they don't consider concurrent event at the same time stamps, and thus  $\bar{t}$  will become  $t$  after one event. For example, we have events  $event_1 = (s, r, o_1, t_1)$ ,  $event_2 = (s, r, o_2, t_1)$ . After  $event_1$ ,  $\bar{t}$  will become  $t$  (subject  $s$ 's most recent time point), and thus the value of intensity function for  $event_2$  will be 0. This is problematic in inference since if  $t = \bar{t}$ , then the intensity function will always be 0 regardless of entity candidates. In their code, they give the highest ranks (first rank) for all entities including the ground truth object in this case, which we think is unfair since the scores of many entity candidates including the ground truth object might be 0 due to their formulation. It has been proven that the performances of Know-Evolve on ICEWS datasets drop down to almost zero after this issue fixed [11].

#### 4.4 Experimental Setup

We used Adam optimizer to train our model and selected the optimal hyper-parameters by early validation stopping according to MRR on the validation set. We restricted the maximum epoch to 5000. We fixed the mini-batch size  $b$  as 512. We tuned the embedding dimensionalities  $d$  in  $\{100, 200, 300, 400, 500\}$ , the ratio of negatives over positive training samples  $\eta$  in  $\{1, 3, 5, 10\}$  and the learning rate  $r$  in  $\{0.00003, 0.0001, 0.0003, 0.001\}$ . The margins  $\gamma$  were varied in the range  $\{1, 2, 3, 5, 10, 20, \dots, 120\}$ . We selected the pair of restriction values  $c_{min}$  and  $c_{max}$  for covariance among  $\{(0.0001, 0.1), (0.003, 0.3), (0.005, 0.5), (0.01, 1)\}$ . The default configuration for ATiSE is as follows:  $lr = 0.00003$ ,  $d = 500$ ,  $\eta = 10$ ,  $\gamma = 1$ ,  $(c_{min}, c_{max}) = (0.005, 0.5)$ . Below, we only list the non-default parameters:  $\gamma = 120$ ,  $(c_{min}, c_{max}) = (0.003, 0.3)$  on ICEWS14;  $\gamma = 100$ ,  $(c_{min}, c_{max}) = (0.003, 0.3)$  on ICEWS05-15.

**Table 3.** Link prediction results on ICEWS14 and ICEWS05-15. \*: results are taken from [6].  $\diamond$ : results are taken from [7]. Dashes: results are unobtainable. The best results among all models are written bold.

Metrics	ICEWS14				ICEWS05-15			
	MRR	Hits@1	Hits@3	Hits@10	MRR	Hits@1	Hits@3	Hits@10
TransE*	.280	.094	-	.637	.294	.090	-	.663
DistMult*	.439	.323	-	.672	.456	.337	-	.691
ComplEx-N3	.467	.347	.527	.716	.481	.362	.535	.729
RotatE	.418	.291	.478	.690	.304	.164	.355	.595
QuatE <sup>2</sup>	.471	.353	.530	.712	.482	.370	.529	.727
TTransE*	.255	.074	-	.601	.271	.084	-	.616
HyTE*	.297	.108	.416	.655	.316	.116	.445	.681
TA-TransE*	.275	.095	-	.625	.299	.096	-	.668
TA-DistMult*	.477	.363	-	.686	.474	.346	-	.728
DE-Simple $\diamond$	.526	.418	.592	.725	.513	.392	.578	.748
ATiSE	<b>.545</b>	<b>.423</b>	<b>.632</b>	<b>.757</b>	<b>.533</b>	<b>.394</b>	<b>.623</b>	<b>.803</b>

**Table 4.** Link prediction results on YAGO11k and Wikidata12k. The best results among all models are written bold.

Metrics	YAGO11k				Wikidata12k			
	MRR	Hits@1	Hits@3	Hits@10	MRR	Hits@1	Hits@3	Hits@10
TransE	.100	.015	.138	.244	.178	.100	.192	.339
DistMult	.158	.107	.161	.268	.222	.119	.238	.460
ComplEx-N3	.167	.106	.154	.282	.233	.123	.253	.436
RotatE	.177	.113	.177	<b>.315</b>	.221	.116	.236	.461
QuatE <sup>3</sup>	.164	.107	.148	.270	.230	.125	.243	.416
TTransE	.108	.020	.150	.251	.172	.096	.184	.329
HyTE	.105	.015	.143	.272	.180	.098	.197	.333
ATiSE	<b>.185</b>	<b>.126</b>	<b>.189</b>	.301	<b>.252</b>	<b>.148</b>	<b>.288</b>	<b>.462</b>

## 4.5 Experimental Results

Table 3 and 4 show the results for link prediction task. On ICEWS14 and ICEWS05-15, ATiSE outperformed all baseline models, considering MR, MRR, Hits@10 and Hits@1. Compared to DE-Simple which is a very recent state-of-the-art TKGE model, ATiSE got improvement of 4% on both datasets regarding MRR, and improved Hits@10 by 4% and 6% on ICEWS14 and ICEWS05-15 respectively. On YAGO11k and Wikidata12k where time annotations in facts are time intervals, ATiSE surpassed baseline models regarding MRR, Hits@1, Hits@3. Regarding Hits@10, ATiSE achieved the state-of-the-art results on Wikidata12k and the second best results on YAGO11k. As mentioned in Section 4.3, the results of TA-TransE, TA-DistMult and DE-Simple on YAGO11k and Wikidata12k are unobtainable since they have difficulties in modeling facts involving time intervals in these two datasets.

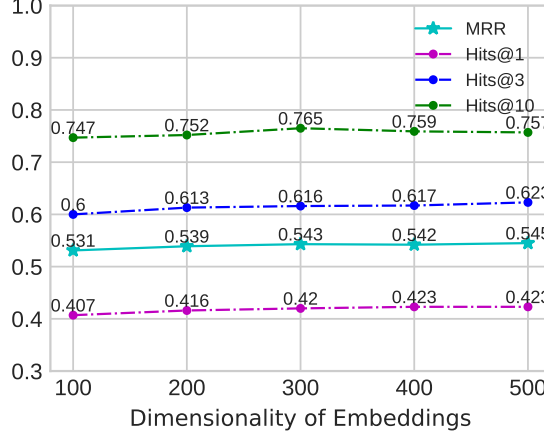
A part of results listed on Table 3 and 4 are obtained based on the implementations released in [22, 12, 4]. We list the implementation details of some baseline models as follows:

- We used the implementation released in [22] to test RotatE on all four datasets, and DistMult on YAGO11k and Wikidata12k. The source code was revised to adopt the time-wise filtered setting. To search the optimal configurations for RotatE and DistMult, we followed the experimental setups reported in [22] except setting the maximum dimensionality as 500 and the maximum negative sampling ratio as 10. The default optimal configuration for RotatE and DistMult is as follows:  $lr = 0.0001$ ,  $b = 1024$ ,  $d = 500$ ,  $\eta = 10$ . Below, we only list the non-default parameters: for RotatE, the optimal margins are  $\gamma = 36$  on ICEWS14,  $\gamma = 48$  on ICEWS05-15,  $\gamma = 3$  on YAGO11k and  $\gamma = 6$  on Wikidata12k; for DistMult, the optimal regularizer weights are  $r = 0.00001$  on YAGO11k and Wikidata12k.
- We used the implementation released in [12] to test ComplEx-N3 and QuatE<sup>2</sup> on all four datasets. The source code was revised to adopt the time-wise filtered setting. To search the optimal configurations for ComplEx-N3 and QuatE<sup>2</sup>, we followed the experimental setups reported in [12] except setting the maximum dimensionality as 500. The default optimal configuration for ComplEx-N3 and QuatE<sup>2</sup> is as follows:  $lr = 0.1$ ,  $d = 500$ ,  $b = 1000$ . Below, we list the optimal regularizer weights: for ComplEx-N3,  $r = 0.01$  on ICEWS14 and ICEWS05-15,  $r = 0.1$  on YAGO11k and Wikidata12k; for QuatE,  $r = 0.01$  on ICEWS14 and YAGO11k,  $r = 0.05$  on ICEWS05-15,  $r = 0.1$  on Wikidata.
- We used the implementation released in [4] to test TransE, TTransE and HyTE on YAGO11k and Wikidata12k for obtaining their performances regarding MRR, Hits@1 and Hits@3. We followed the optimal configurations reported in [4]. As shown in Table 4, Hits@10s of TransE and TTransE we got were better than those reported in [4].
- As shown in Table 3, other baseline results are taken from [6, 7].

#### 4.6 Ablation Study

In this work, we analyze the effects of the dimensionality and various components of entity/relation embeddings.

The embedding dimensionality is an important hyperparameter for each (T)KGE model. A high embedding dimensionality might be beneficial to boost the performance of a (T)KGE model. For instance, ComplEx-N3 and QuatE<sup>2</sup> achieved the state-of-the-art results on link prediction over static KGs with 2000-dimensional embeddings [12, 29]. On the other hand, a lower embedding dimensionality will lead to less consumption on training time and memory space, which is quite important for the applications of (T)KGE models on large-scale datasets. Figure 2 shows the performances of ATiSE with different embedding dimensionalities on ICEWS14. With a same embedding dimensionality of 100 as DE-Simple [7], ATiSE still achieved the state-of-the-art results on ICEWS14. An



**Fig. 2.** Results for ATiSE with different embedding dimensionalities on ICEWS14.

ATiSE model with an embedding dimensionality of 100 trained on ICEWS14 had a memory size of 14.2Mb while a DE-Simple model and a QuatE<sup>2</sup> model with the same embedding dimensionality had memory sizes of 13.3Mb and 12.4Mb. And the memory size of an ATiSE model increases linearly with its embedding dimensionality. Moreover, training an ATiSE model with an embedding dimensionality of 100 took 2.8 seconds per epoch on a single GeForce RTX2080, and an ATiSE with 500-dimensional embeddings took 3.7 seconds per epoch.

To analyze the effects of different components of entity/relation representation in ATiSE, we developed three comparison models, namely, ATiSE-SN, ATiSE-TN and ATiSE-TS, which exclude the trend component, seasonal component and the noise component respectively. The entity representations of these three comparison models are as follows:

$$\begin{aligned}
 e_{i,t}^{SN} &= e_i + \beta_{e,i} \sin(2\pi\omega_{e,i}t) + \mathcal{N}(0, \Sigma_{e,i}) \\
 e_{i,t}^{TN} &= e_i + \alpha_{e,i}w_{e,i}t + \mathcal{N}(0, \Sigma_{e,i}) \\
 e_{i,t}^{TS} &= e_i + \alpha_{e,i}w_{e,i}t + \beta_{e,i} \sin(2\pi\omega_{e,i}t)
 \end{aligned} \tag{6}$$

For ATiSE-TS consisting of the trend component and the seasonal component, we used the translation-based scoring function [3] to measure the plausibility of the fact  $(s, p, o, t)$ .

$$f_t^{TS}(e_s, r_p, e_o) = \|e_{s,t}^{TS} + r_{p,t}^{TS} - e_{o,t}^{TS}\| \tag{7}$$

We report the MRRs and Hits@10 of ATiSE-SN, ATiSE-TN and ATiSE-TS on link prediction over ICEWS14 and YAGO11k. As shown in Table 5, we find that the removal of the trend component and the noise component had a remarkable negative effect on the performance of ATiSE on link prediction since the model could not address the temporal uncertainty of entity/relation representations without the noise component and the trend component contained the main time

**Table 5.** Link prediction results of ablation experiments.

Datasets	ICEWS14				YAGO11K <sub>D</sub>			
Metrics	MRR	Hits@1	Hits@3	Hits@10	MRR	Hits@1	Hits@3	Hits@10
ATiSE-SN	.405	.284	.488	.710	.139	.095	.143	.249
ATiSE-TN	.536	.407	.626	.771	.167	.115	.171	.292
ATiSE-TS	.323	.127	.429	.676	.115	.023	.145	.274
ATiSE	.545	.423	.632	.757	.185	.126	.189	.301

information. In ATiSE, different types of entities might have big difference in the trend component. For instance, we found that the embeddings of entities representing people, e.g., *Barack Obama*, generally had higher evolution rates than those representing cities or nations, e.g., *USA*.

ATiSE-TN performed worse than ATiSE on YAGO11k where facts involve time intervals. Different from ICEWS14 dataset which is an event-based dataset where all relations or predicates are instantaneous, there exist both short-term relations and long-term relations in YAGO11k. Adding seasonal components into evolving entity/relation representations is helpful to distinguish short-term patterns and long-term patterns in YAGO11k. It can be seen from Table 6 that short-term relations learned by ATiSE, e.g., *wasBornIn*, generally had higher evolutionary rates, and their seasonal components had smaller amplitudes and higher frequencies than long-term relations, e.g., *isMarriedTo*.

**Table 6.** Relations in YAGO11k and the mean step numbers of their duration time (TS), as well as the corresponding parameters learned from ATiSE, including the evolutionary rate  $|\alpha_r|$ , the mean amplitude  $|\beta_r|$  and the mean frequency  $|\omega_r|$  of the seasonal component for each relation.

Relations	#TS	$ \alpha_r $	$ \beta_r $	$ \omega_r $
<i>wasBornIn</i>	1.0	0.142	0.000	1.032
<i>worksAt</i>	18.7	0.046	0.058	0.294
<i>playsFor</i>	4.7	0.071	0.046	0.766
<i>hasWonPrize</i>	28.6	0.010	0.107	0.041
<i>isMarriedTo</i>	16.5	0.049	0.076	0.090
<i>owns</i>	24.9	0.017	0.088	0.101
<i>graduatedFrom</i>	38.1	0.016	0.104	0.029
<i>deadIn</i>	1.0	0.249	0.006	0.897
<i>isAffiliatedTo</i>	25.8	0.014	0.049	0.126
<i>created</i>	27.1	0.011	0.040	0.087

## 5 CONCLUSION

We introduce ATiSE, a temporal KGE model that incorporates time information into KG representations by fitting the temporal evolution of entity/relation

representations over time as additive time series. Considering the uncertainty during the temporal evolution of KG representations, ATiSE maps the representations of temporal KGs into the space of multi-dimensional Gaussian distributions. The covariance of an entity/relation representation represents its randomness component. Experimental results demonstrate that our method significantly outperforms the state-of-the-art methods on link prediction over four TKG benchmarks.

Our work establishes a previously unexplored connection between relational processes and time series analysis with a potential to open a new direction of research on reasoning over time. In the future, we will explore to use more sophisticated models to model different components of relation/entity representations, e.g., an ARIMA model for the noise component and a polynomial model for the trend component.

## References

1. Auer, S., Bizer, C., Kobilarov, G., Lehmann, J., Cyganiak, R., Ives, Z.: Dbpedia: A nucleus for a web of open data. *The semantic web* pp. 722–735 (2007)
2. Bollacker, K., Evans, C., Paritosh, P., Sturge, T., Taylor, J.: Freebase: a collaboratively created graph database for structuring human knowledge. In: *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*. pp. 1247–1250. AcM (2008)
3. Bordes, A., Usunier, N., Garcia-Duran, A., Weston, J., Yakhnenko, O.: Translating embeddings for modeling multi-relational data. In: *Advances in Neural Information Processing Systems*. pp. 2787–2795 (2013)
4. Dasgupta, S.S., Ray, S.N., Talukdar, P.: HYTE: Hyperplane-based temporally aware knowledge graph embedding. In: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. pp. 2001–2011 (2018)
5. Erxleben, F., Günther, M., Krötzsch, M., Mendez, J., Vrandečić, D.: Introducing wikidata to the linked data web. In: *International Semantic Web Conference*. pp. 50–65. Springer (2014)
6. García-Durán, A., Dumančić, S., Niepert, M.: Learning sequence encoders for temporal knowledge graph completion. In: *EMNLP* (2018)
7. Goel, R., Kazemi, S.M., Brubaker, M., Poupart, P.: Diachronic embedding for temporal knowledge graph completion. In: *AAAI* (2020)
8. He, S., Liu, K., Ji, G., Zhao, J.: Learning to represent knowledge graphs with gaussian embedding. In: *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management*. pp. 623–632. ACM (2015)
9. Ho, S., Xie, M.: The use of arima models for reliability forecasting and analysis. *Computers & industrial engineering* **35**(1-2), 213–216 (1998)
10. Ji, G., He, S., Xu, L., Liu, K., Zhao, J.: Knowledge graph embedding via dynamic mapping matrix. In: *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. vol. 1, pp. 687–696 (2015)
11. Jin, W., Jiang, H., Qu, M., Chen, T., Zhang, C., Szekely, P., Ren, X.: Recurrent event network: Global structure inference over temporal knowledge graph. *arXiv: 1904.05530* (2019)

12. Lacroix, T., Usunier, N., Obozinski, G.: Canonical tensor decomposition for knowledge base completion. In: International Conference on Machine Learning. pp. 2869–2878 (2018)
13. Lautenschlager, J., Shellman, S., Ward, M.: Icews event aggregations (2015). <https://doi.org/10.7910/DVN/28117>, <https://doi.org/10.7910/DVN/28117>
14. Leblay, J., Chekol, M.W.: Deriving validity time in knowledge graph. In: Companion of the The Web Conference 2018 on The Web Conference 2018. pp. 1771–1776. International World Wide Web Conferences Steering Committee (2018)
15. Leetaru, K., Schrod, P.A.: Gdelt: Global data on events, location, and tone, 1979–2012. In: ISA annual convention. vol. 2, pp. 1–49. Citeseer (2013)
16. Mahdisoltani, F., Biega, J., Suchanek, F.M.: Yago3: A knowledge base from multilingual wikipedias. In: CIDR (2013)
17. Miller, G.A.: WordNet: An electronic lexical database. MIT press (1998)
18. Montgomery, D.C., Jennings, C.L., Kulahci, M.: Introduction to time series analysis and forecasting. John Wiley & Sons (2015)
19. Nickel, M., Tresp, V., Krieger, H.P.: A three-way model for collective learning on multi-relational data. In: ICML. vol. 11, pp. 809–816 (2011)
20. Petersen, K.B., Pedersen, M.S., et al.: The matrix cookbook. Technical University of Denmark **7**(15), 510 (2008)
21. Suchanek, F.M., Kasneci, G., Weikum, G.: Yago: a core of semantic knowledge. In: Proceedings of the 16th international conference on World Wide Web. pp. 697–706. ACM (2007)
22. Sun, Z., Deng, Z.H., Nie, J.Y., Tang, J.: Rotate: Knowledge graph embedding by relational rotation in complex space. In: ICLR (2019)
23. Trivedi, R., Dai, H., Wang, Y., Song, L.: Know-evolve: Deep temporal reasoning for dynamic knowledge graphs. In: ICML (2017)
24. Trouillon, T., Welbl, J., Riedel, S., Gaussier, É., Bouchard, G.: Complex embeddings for simple link prediction. In: Proceedings of ICML (2016)
25. Wang, Q., Mao, Z., Wang, B., Guo, L.: Knowledge graph embedding: A survey of approaches and applications. IEEE Transactions on Knowledge and Data Engineering **29**(12), 2724–2743 (2017)
26. Wang, Z., Zhang, J., Feng, J., Chen, Z.: Knowledge graph embedding by translating on hyperplanes. In: AAAI. pp. 1112–1119. Citeseer (2014)
27. Yang, B., Yih, W.t., He, X., Gao, J., Deng, L.: Embedding entities and relations for learning and inference in knowledge bases. In: ICLR. p. 12 (2015)
28. Yu, D., Yao, K., Su, H., Li, G., Seide, F.: KL-divergence regularized deep neural network adaptation for improved large vocabulary speech recognition. In: 2013 IEEE International Conference on Acoustics, Speech and Signal Processing. pp. 7893–7897. IEEE (2013)
29. Zhang, S., Tay, Y., Yao, L., Liu, Q.: Quaternion knowledge graph embeddings. In: Advances in Neural Information Processing Systems. pp. 2731–2741 (2019)