

YOU CAN TEACH AN OLD DOG NEW TRICKS! ON TRAINING KNOWLEDGE GRAPH EMBEDDINGS

Anonymous authors

Paper under double-blind review

ABSTRACT

Knowledge graph embedding (KGE) models learn algebraic representations of the entities and relations in a knowledge graph. A vast number of KGE techniques for multi-relational link prediction has been proposed in the recent literature, often with state-of-the-art performance. These approaches differ along a number of dimensions, including different model architectures, different training strategies, and different approaches to hyperparameter optimization. In this paper, we take a step back and aim to summarize and quantify empirically the impact of each of these dimensions on model performance. We report on the results of an extensive experimental study with popular model architectures and training strategies across a wide range of hyperparameter settings. We found that when trained appropriately, the relative performance differences between various model architectures often shrinks and sometimes even reverses when compared to prior results. For example, RESCAL (Nickel et al., 2011), one of the first KGE models, showed strong performance when trained with state-of-the-art techniques; it was competitive to or outperformed more recent architectures. We also found that good (and often superior to prior studies) model configurations can be found by exploring relatively few random samples from a large hyperparameter space. Our results suggest that many of the more advanced architectures and techniques proposed in the literature should be revisited to reassess their individual benefits.

1 INTRODUCTION

Knowledge graph embedding (KGE) models learn algebraic representations, termed embeddings, of the entities and relations in a knowledge graph. They have been successfully applied to knowledge graph completion (Nickel et al., 2015) as well as in downstream tasks and applications such as recommender systems (Wang et al., 2017) or visual relationship detection (Baier et al., 2017).

A vast number of different KGE models for multi-relational link prediction has been proposed in the recent literature; e.g., RESCAL (Nickel et al., 2011), TransE (Bordes et al., 2013), DistMult, ComplEx (Trouillon et al., 2016), ConvE (Dettmers et al., 2018), TuckER (Balazevic et al., 2019), RotatE (Sun et al., 2019), KBGAT (Nathani et al., 2019), and many more. Model architectures generally differ in the way the entity and relation embeddings are combined to model the presence or absence of an edge (more precisely, a subject-predicate-object triple) in the knowledge graph; they include factorization models (e.g., RESCAL, DistMult, ComplEx, TuckER), translational models (TransE, RotatE), and more advanced models such as convolutional models (ConvE). In many cases, the introduction of new models went along with the new approaches for training these models—e.g., new training types (such as negative sampling or 1vsAll scoring), new loss functions (such as pairwise margin ranking or binary cross entropy), new forms of regularization (such as unweighted and weighted L2), or the use of reciprocal relations (Dettmers et al., 2018)—and ablation studies were not always performed. Table 1 shows an overview of selected models and training techniques along with the publications that introduced them.

The diversity in model training makes it difficult to compare performance results for various model architectures, especially when results are reproduced from prior studies that used a different experimental setup. Model hyperparameters are commonly tuned using grid search on a small grid involving hand-crafted parameter ranges or settings known to “work well” from prior studies. A grid suitable for one model may be suboptimal for another, however. Indeed, it has been observed

Publication	Model	Loss	Training	Regularizer	Optimizer	Reciprocal
Nickel et al. (2011)	RESCAL	MSE	Full	L2	ALS	No
Bordes et al. (2013)	TransE	MR	NegSamp	Normalization	SGD	No
Yang et al. (2015)	DistMult	MR	NegSamp	Weighted L2	Adagrad	No
Trouillon et al. (2016)	ComplEx	BCE	NegSamp	Weighted L2	Adagrad	No
Kadlec et al. (2017)	DistMult	CE	NegSamp	Weighted L2	Adam	No
Dettmers et al. (2018)	ConvE	BCE	KvsAll	Dropout	Adam	Yes
Lacroix et al. (2018)	ComplEx	CE	1vsAll	Weighted L3	Adagrad	Yes

MSE = mean squared error, MR = margin ranking, BCE = binary cross entropy, CE = cross entropy

Table 1: Selected KGE models and training strategies from the literature. Entries marked in bold were introduced (or first used) in the context of KGE in the corresponding publication.

that newer training strategies can considerably improve model performance (Kadlec et al., 2017; Lacroix et al., 2018; Salehi et al., 2018).

In this paper, we take a step back and aim to summarize and quantify empirically the impact of different model architectures and different training strategies on model performance. We performed an extensive set of experiments using popular model architectures and training strategies in a common experimental setup. In contrast to most prior work, we considered many training strategies as well as a large hyperparameter space, and we performed model selection using quasi-random search (instead of grid search) followed by Bayesian optimization. We found that this approach was able to find good (and often superior to prior studies) model configurations with relatively low effort. Our study complements and expands on the results of Kotnis & Nastase (2018) (focus on negative sampling) and Mohamed et al. (2019) (focus on loss functions) as well as similar studies in other areas, including language modeling (Melis et al., 2017), generative adversarial networks (Lucic et al., 2018), or sequence tagging (Reimers & Gurevych, 2017).

We found that when trained appropriately, the performance of a particular model architecture can by far exceed the performance observed in older studies. For example, RESCAL (Nickel et al., 2011), which constitutes one of the first KGE models but is rarely considered in newer work, showed very strong performance in our study: it was competitive to or outperformed more recent architectures such as ConvE (Dettmers et al., 2018) and TuckER (Balazevic et al., 2019). More generally, we found that the relative performance differences between various model architectures often shrunk and sometimes even reversed when compared to prior results. This suggests that (at least currently) training strategies have a significant impact on model performance and may account for a substantial fraction of the progress made in recent years. We also found that suitable training strategies and hyperparameter settings vary significantly across models and datasets, indicating that a small search grid may bias results on model performance. Fortunately, as indicated above, large hyperparameter spaces can be (and should be) used with little additional training effort.

Our study focuses solely on pure KGE models, which do not exploit auxiliary information such as textual data or logical rules (Wang et al., 2017). Since many of the studies on these non-pure models did not (and, to be fair, could not) use current training strategies and consequently underestimated the performance of pure KGE models, their results and conclusions need to be revisited.

2 KNOWLEDGE GRAPH EMBEDDINGS: MODELS, TRAINING, EVALUATION

The literature on KGE models is expanding rapidly. We review selected architectures, training methods, and evaluation protocols; see Table 1. The table exemplarily indicates that new model architectures are sometimes introduced along with new training strategies (marked bold). Reasonably recent survey articles about KGE models include Nickel et al. (2015) and Wang et al. (2017).

Multi-relational link prediction. KGE models are typically trained and evaluated in the context of multi-relational link prediction for knowledge graphs (KG). Generally, given a set \mathcal{E} of entities and a set \mathcal{R} of relations, a knowledge graph $\mathcal{K} \subseteq \mathcal{E} \times \mathcal{R} \times \mathcal{E}$ is a set of subject-predicate-object (spo) triples. The goal of multi-relational link prediction is to “complete the KG”, i.e., to predict true but unobserved triples based on the information in \mathcal{K} . Common approaches include rule-based

methods (Galarraga et al., 2013; Meilicke et al., 2019), KGE methods (Nickel et al., 2011; Bordes et al., 2013; Trouillon et al., 2016; Dettmers et al., 2018), and hybrid methods (Guo et al., 2018).

Knowledge graph embeddings (KGE). A KGE model associates with each entity $i \in \mathcal{E}$ and relation $k \in \mathcal{R}$ an *embedding* $e_i \in \mathbb{R}^{d_e}$ and $r_k \in \mathbb{R}^{d_r}$ in a low-dimensional vector space, respectively; here $d_e, d_r \in \mathbb{N}^+$ are hyperparameters for the *embedding size*. Each particular model uses a *scoring function* $s : \mathcal{E} \times \mathcal{R} \times \mathcal{E} \rightarrow \mathbb{R}$ to associate a *score* $s(i, k, j)$ with each potential triple $(i, k, j) \in \mathcal{E} \times \mathcal{R} \times \mathcal{E}$. The higher the score of a triple, the more likely it is considered to be true by the model. The scoring function takes form $s(i, k, j) = f(e_i, r_k, e_j)$, i.e., depends on i, k , and j only through their respective embeddings. Here f , which represents the model architecture, may be either a fixed function or learned (e.g., f may be a parameterized function).

Evaluation. The most common evaluation task for KGE methods is *entity ranking*, which is a form of question answering. The available data is partitioned into a set of training, validation, and test triples. Given a test triple (i, k, j) (unseen during training), the entity ranking task is to determine the correct answer—i.e., the missing entity j and i , resp.—to questions $(i, k, ?)$ and $(?, k, j)$. To do so, potential answer triples are first ranked by their score in descending order. All triples but (i, k, j) that occur in the training, validation, or test data are subsequently filtered out (so that other triples known to be true do not affect the ranking). Finally, metrics such as the mean reciprocal rank (MRR) of the true answer or HITS@ k are computed.

KGE models. KGE model architectures differ in their scoring function. We can roughly classify models as *decomposable* or *monolithic*: the former only allow element-wise interactions between (relation-specific) subject and object embeddings, whereas the latter allow arbitrary interactions. More specifically, decomposable models use scoring functions of form $s(i, k, j) = f(\sum_z g([h_1(e_i, e_r) \circ h_2(e_r, e_j)]_z))$, where \circ is any element-wise function (e.g., multiplication), h_1 and h_2 are functions that obtain *relation-specific* subject and object embeddings, resp., and g and f are scalar functions (e.g., identity or sigmoid). The most popular models in this category are perhaps RESCAL (Nickel et al., 2011), TransE (Bordes et al., 2013), DistMult (Yang et al., 2015), ComplEx (Trouillon et al., 2016), and ConvE (Dettmers et al., 2018). RESCAL’s scoring function is bilinear in the entity embeddings: it uses $s(i, k, j) = e_i^T R_k e_j$, where $R_k \in \mathbb{R}^{d_e \times d_e}$ is a matrix formed from the entries of $r_k \in \mathbb{R}^{d_r}$ (where $d_r = d_e^2$). DistMult and ComplEx can be seen as constrained variants of RESCAL with smaller relation embeddings ($d_r = d_e$). TransE is a translation-based model and uses negative distance $-||e_i + r_k - e_j||_p$ between $e_i + r_k$ and e_j as score, commonly using the L1 norm ($p = 1$) or the L2 norm ($p = 2$). Finally, ConvE uses a 2D convolutional layer and a large fully connected layer to obtain relation-specific entity embeddings (i.e., in h_1 above). Other recent examples for decomposable models include TuckER (Balazevic et al., 2019), RotatE (Sun et al., 2019), and SACN (Shang et al., 2019). Decomposable models are generally fast to use: once the relation-specific embeddings are (pre-)computed, score computations are cheap. Monolithic models—such as KBGAT (Nathani et al., 2019)—do not decompose into relation-specific embeddings: they take form $s(i, k, j) = f(e_i, r_k, e_j)$. Such models are more flexible, but also more costly to train and use.

Training type. There are three commonly used approaches to train KGE models, which differ mainly in the way negative examples are generated. First, training with negative sampling (NegSamp) (Bordes et al., 2013) obtains for each positive triple $t = (i, k, j)$ from the training data a set of (pseudo-)negative triples obtained by randomly perturbing the subject, relation, or object position in t (and optionally checking that the so-obtained triples does not exist in the KG). An alternative approach (Lacroix et al., 2018), which we term *1vsAll*, is to omit sampling and take *all* triples that can be obtained by perturbing the subject and object positions as negative examples for t (even if these tuples exist in the KG). 1vsAll is generally more expensive than NegSamp, but it is feasible (and even surprisingly fast in practice) if the number of entities is not excessively large. Finally, Dettmers et al. (2018) proposed a training type that we term *KvsAll*¹: this approach (i) constructs batches from non-empty rows $(i, k, *)$ or $(*, k, j)$ (instead of from individual triples), and (ii) labels all such triples as either positive (occurs in training data) or negative (otherwise).

Loss functions. Several loss functions for training KGEs have been introduced so far. RESCAL originally used squared error between the score of each triple and its label (positive or negative). TransE used pairwise margin ranking with hinge loss (MR), where each pair consists of a positive

¹Note that the KvsAll strategy is called 1-N scoring in Dettmers et al. (2018).

triple and one of its negative triples (only applicable to NegSamp and 1vsAll) and the margin η is a hyperparameter. Trouillon et al. (2016) proposed to use binary cross entropy (BCE) loss: it applies a sigmoid to the score of each (positive or negative) triple and uses the cross entropy between the resulting probability and that triple’s label as loss. BCE is suitable for multi-class and multi-label classification. Finally, Kadlec et al. (2017) used cross entropy (CE) between the model distribution (softmax distribution over scores) and the data distribution (labels of corresponding triples, normalized to sum to 1). CE is more suitable for multi-class classification (as in NegSamp and 1vsAll), but it has also been used in the multi-label setting (KvsAll). Mohamed et al. (2019) found that the choice of loss function can have a significant impact on model performance, and that the best choice is data and model dependent. Our experimental study provides additional evidence for this finding.

Reciprocal relations. ConvE introduced the technique of *reciprocal relations* into KGE training. Observe that during evaluation and also most training methods discussed above, the model is solely asked to score subjects (for questions of form $(?, k, j)$) or objects (for questions of form $(i, k, ?)$). The idea of reciprocal relations is to use separate scoring functions s_{sub} and s_{obj} for each of these tasks, resp. Both scoring functions share entity embeddings, but they do not share relation embeddings: each relation thus has two embeddings.² The use of reciprocal relations may decrease the computational cost (as in the case of ConvE), and it may also lead to better model performance Lacroix et al. (2018) (e.g., for relations in which one direction is easier to predict). On the downside, the use of reciprocal relations means that a model does not provide a single triple score $s(i, k, j)$ anymore (generally, $s_{\text{sub}}(i, k, j) \neq s_{\text{obj}}(i, k, j)$; the discrepancy has not been studied yet).

Regularization. The most popular form of regularization in the literature is L2 regularization on the embedding vectors, either unweighted or weighted by the frequency of the corresponding entity or relation (Yang et al., 2015). Lacroix et al. (2018) proposed to use L3 regularization. TransE normalized the embeddings to unit norm after each update. ConvE used dropout (Srivastava et al., 2014) in its hidden layers (and only in those). In our study, we additionally consider L1 regularization and the use of dropout on the entity and/or relation embeddings.

Hyperparameters. Many more hyperparameters have been used in prior work. This includes, for example, different methods to initialize model parameters, different optimizers, different optimizer parameters such as the learning rate or batch size, the number of negative examples for NegSamp, the regularization weights for entities and relations, and so on. To deal with such a large search space, most prior studies favor grid search over a small grid where most of these hyperparameters remain fixed. As discussed before, this approach may lead to bias in the results, however.

3 EXPERIMENTAL STUDY

In this section, we report on the design and results of our experimental study. We focus on the most salient points here; more details can be found in the appendix.

3.1 EXPERIMENTAL SETUP

Datasets. We used the FB15K-237 (Toutanova & Chen, 2015) (extracted from Freebase) and WNRR (Dettmers et al., 2018) (extracted from WordNet) datasets in our study. We chose these datasets because (i) they are frequently used in prior studies, (ii) they are “hard” datasets that have been designed specifically to evaluate multi-relational link prediction techniques, (iii) they are diverse in that relative model performance often differs, and (iv) they are of reasonable size for a large study. Dataset statistics are given in Table 4 in the appendix.

Models. We selected RESCAL (Nickel et al., 2011), TransE (Bordes et al., 2013), DistMult (Yang et al., 2015), ComplEx (Trouillon et al., 2016) and ConvE (Dettmers et al., 2018) for our study. These models are perhaps the most well-known KGE models and include both early and more recent models. Note that we did not consider monolithic models in our study due to their excessive training cost (for example, KBGAT trained for 14 hours using the best configuration reported by the authors).

Evaluation. Unless stated otherwise, we report filtered MRR and filtered HITS@10 on test data; see Sec. 2 for details.

²Alternatively, we can view such an approach as only predicting objects, but doing so for both each original relation (k) and a new reciprocal relation formed by switching subject and object (k^{-1}).

Hyperparameters. We used a large hyperparameter space to ensure sure that suitable hyperparameters for each model are not excluded a priori. We included all major training types (NegSamp, IvsAll, KvsAll), use of reciprocal relations, loss functions (MR, BCE, CE), regularization techniques (none/L1/L2/L3, dropout), optimizers (Adam, Adagrad), and initialization methods (4 in total) used in the KGE community as hyperparameters. We considered three embeddings sizes (128, 256, 512) and used separate weights for dropout/regularization for entity and relation embeddings. Table 5 in the appendix provides additional details. To the best of our knowledge, no prior study used such a large hyperparameter search space.

Training. All models were trained for a maximum of 400 epochs. We validated models using filtered MRR (on validation data) every five epochs and performed early stopping with a patience of 50 epochs. To keep search tractable, we stopped training on models that did not reach $\geq 5\%$ filtered MRR after 50 epochs; in our experience, such configurations did not produce good models.

Model selection. We used the Ax framework (<https://ax.dev/>) to conduct quasi-random hyperparameter search via a Sobol sequence. More specifically, for each dataset and model, we generated 30 different configurations per valid combination of training type and loss function (2 for TransE, which only supports NegSamp+MR and NegSamp+CE; 7 for all other models). After the quasi-random hyperparameter search, we added a short Bayesian optimization phase (20 trials, using expected improvement; also provided by Ax) to tune the numerical hyperparameters further.

Reproducibility. We reimplemented all models, training strategies, and hyperparameter search in a framework based on PyTorch. The framework is easily extendable and will be released as open source. Model selection and training is completely controlled via configuration files. All configuration files used in our study, as well as the raw data that we obtained when running our experiments, will also be openly released.³

3.2 COMPARISON OF MODEL PERFORMANCE

Table 2 shows the filtered MRR and filtered Hits@10 on test data of various models both from prior studies and the best models (on validation MRR) found in our study.

First reported performance vs. our observed performance. We first compare the performance results of models when they were first run on the FB15K-237 and the WNRR datasets (“First” block of Table 2) and the performance we obtained in our study (“Ours” block). We found that the performance of a single model can vary wildly across studies. For example, ComplEx was first run on FB15K-237 by Dettmers et al. (2018), where it achieved a filtered MRR of 24.7%. This is a relatively low number by today’s standards (Lacroix et al., 2018). In our study, ComplEx achieved a competitive MRR of 35.1%, which is also a large improvement over the first reports. Studies that report the lower performance number of ComplEx (i.e., 24.7%) thus do not adequately capture its performance. Similar remarks hold for RESCAL and DistMult as well as (albeit to a smaller extent) ConvE and TransE. To the best of our knowledge, our results for RESCAL, TransE, and ConvE constitute the best results for these models obtained so far.

Relative model performance (our study). Next, we compared the performance across the models used in our study (“Ours” block). We found that the relative performance differences between various model architectures often shrunk and sometimes even reversed when compared to prior results (“First” block). For example, ConvE showed the best overall performance in prior studies, but is consistently outperformed by RESCAL, DistMult, and ComplEx in our study. As another example, RESCAL (Nickel et al., 2011), which constitutes one of the first KGE models but is rarely considered in newer work, showed strong performance in our study and outperforms all alternative models except ComplEx, which also showed strong performance. This suggests that (at least currently) training strategies have a significant impact on model performance and may account for a large fraction of the progress made in recent years.

Relative model performance (overall). Table 2 also shows the best performance results obtained in prior studies of selected recent models (“Recent” block) and very large models with very good performance (“Large” block). We found that Tucker (Balazevic et al., 2019), RotatE (Sun et al., 2019), and SACN (Shang et al., 2019) all achieve state-of-the-art performance, but the performance difference to the best prior models (“Ours” block) in terms of MRR is small or even vanishes. Even

³Note to ICLR reviewers: We do not release this data with the submission to maintain author anonymity.

		<i>FB15K-237</i>		<i>WNRR</i>	
		MRR	Hits@10	MRR	Hits@10
<i>First</i>	RESCAL (Wang et al., 2019)	27.0	42.7	42.0	44.7
	TransE (Nguyen et al., 2018)	29.4	46.5	22.6	50.1
	DistMult (Dettmers et al., 2018)	24.1	41.9	43.0	49.0
	ComplEx (Dettmers et al., 2018)	24.7	42.8	44.0	51.0
	ConvE (Dettmers et al., 2018)	32.5	50.1	43.0	52.0
<i>Ours</i>	RESCAL	35.6	53.6	46.8	52.1
	TransE	30.3	48.8	22.7	52.6
	DistMult	34.2	52.3	45.2	53.1
	ComplEx	35.1	53.7	47.7	54.3
	ConvE	33.7	52.8	44.7	50.8
<i>Recent</i>	TuckER (Balazevic et al., 2019)	35.8	54.4	47.0	52.6
	RotatE (Sun et al., 2019)	33.8	53.3	47.6	57.1
	SACN (Shang et al., 2019)	35.0	54.0	47.0	54.4
<i>Large</i>	DistMult (Salehi et al., 2018)	35.7	54.8	45.5	54.4
	ComplEx-N3 (Lacroix et al., 2018)	37.0	56.0	49.0	58.0
	KBGAT (Nathani et al., 2019)	51.8	62.6	44.0	58.1

Table 2: Model performance in prior studies and our study (on test data). *First*: first reported performance on the respective datasets (oldest models first); *Ours*: performance in our study; *Large*: best performance achieved in prior studies using larger, more expensive models (either significantly larger embedding sizes or monolithic models). Bold numbers indicate best performance in group. References indicate where the performance number was reported.

for HITS@10, which we did not consider for model selection, the advantage of more recent models is often small, if present. The models in the last block (“Large”) show the best performance numbers we have seen in the literature. For DistMult and ComplEx, these numbers have been obtained using very large embedding sizes (up to 4000). KBGAT showed exceptional performance on FB15K-237, clearly outperforming all other models. The model is monolithic, however, and has at least an order of magnitude higher training cost.

Limitations. We note that all models used in our study are likely to yield better performance when hyperparameter tuning is further improved. For example, we found a ComplEx configuration (of size 100) that achieved 49.0% MRR and 56.5% Hits@10 on WNRR, and a RESCAL configuration that achieved 36.1% MRR and 54.6% Hits@10 on FB15k-237 in preliminary experiments. Since the focus of this study is to compare models in a common experimental setup and without manual intervention or large amounts of hyperparameter tuning, we do not further report these results.

3.3 IMPACT OF HYPERPARAMETERS

Anatomy of search space. Figure 1 shows the distribution filtered MRR for each model and dataset (on validation data, Sobol configurations only). Each distribution consists of roughly 200 different quasi-random hyperparameter configurations (except TransE, for which 60 configurations have been used). Perhaps unsurprisingly, we found that all models showed a wide dispersion on both datasets and only very few configurations achieved state-of-the-art results. There are, however, notable differences across datasets and models. On FB15K-237, for example, the median ConvE configuration is best; on WNRR, DistMult and Complex have a much higher median MRR. Generally, the impact of the hyperparameter choice is more pronounced on WNRR (higher variance) than on FB15K-237.

Best configurations (quasi-random search). The hyperparameters of the best performing models during our quasi-random search are reported in Table 3 (selected hyperparameters) and Tables 6 and 7 (in the appendix, all hyperparameters). First, we found that the optimum choice of hyperparameters is often model and dataset dependent: with the exception of the loss function (discussed below), almost no hyperparameter setting was consistent across datasets and models. For example, Adagrad and the 1vsN training type seem to be beneficial on FB15K-237 but not on WNRR, where

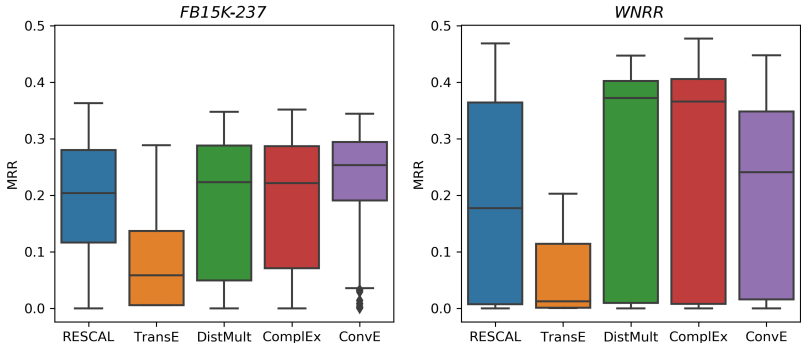


Figure 1: Distribution of filtered MRR (on validation data) over the quasi-random hyperparameter configurations explored in our study.

	RESCAL	TransE	DistMult	ComplEx	ConvE	
<i>FB15K-237</i>	MRR (valid)	36.2	28.8	34.7	35.1	34.4
	Embedding size	512 (-0.5)	256 (-1.1)	128 (-0.1)	512 (-0.1)	512 (-1.0)
	Batch size	512 (-0.5)	1024 (-1.1)	1024 (-0.1)	512 (-0.1)	1024 (-0.5)
	Training type	1vsAll (-1.1)	NegSamp -	1vsAll (-0.9)	1vsAll (-1.0)	1vsAll (-1.0)
	Loss	CE (-1.5)	CE (-1.9)	CE (-1.4)	CE (-2.4)	CE (-1.0)
	Optimizer	Adag. (-0.9)	Adam (-1.1)	Adag. (-0.7)	Adag. (-1.1)	Adag. (-1.0)
	Initializer	Unif. (-0.9)	XvUnif (-4.4)	Unif. (-0.1)	XvUnif (-0.1)	Normal (-0.5)
	Regularizer	None (-0.5)	L2 (-1.1)	L3 (-0.5)	L3 (-0.8)	L3 (-0.5)
	Reciprocal	No (-0.5)	No -	Yes (-0.5)	Yes (-0.8)	Yes -
	<i>WNRR</i>	MRR (valid)	46.8	20.2	44.7	47.7
Embedding size		256 (-0.8)	256 (-1.8)	512 (-0.6)	512 (-1.3)	128 (-0.7)
Batch size		1024 (-0.8)	1024 (-1.8)	1024 (-0.1)	1024 (-1.2)	256 (-0.7)
Training type		KvsAll (-1.6)	NegSamp -	KvsAll (-0.6)	KvsAll (-1.5)	1vsAll (-0.7)
Loss		CE (-1.0)	CE (-2.0)	CE (-1.8)	CE (-3.3)	CE (-1.5)
Optimizer		Adag. (-0.8)	Adam (-1.8)	Adam (-0.6)	Adam (-1.3)	Adam (-1.3)
Initializer		XvUnif (-0.8)	XvUnif (-2.0)	Normal (-0.6)	Normal (-1.3)	XvNorm (-1.3)
Regularizer		None (-0.8)	L2 (-1.8)	None (-0.1)	None (-1.2)	L2 (-1.3)
Reciprocal		No (-1.0)	No -	Yes (-0.6)	Yes (-1.4)	Yes -

Table 3: Hyperparameters of best performing models w.r.t. filtered MRR (on validation data) after quasi-random hyperparameter search. For each hyperparameter, we also give the reduction in filtered MRR for the best configuration that does not use this value of the hyperparameter (in parenthesis).

Adam and KvsAll were more prominent. Our findings provide further evidence that grid search on a small search space is not suitable to compare model performance because the result may be heavily influenced by the specific grid points being used. Moreover, the budget that we used for quasi-random search was comparable to a grid search over a small (or at least not very large) grid: we tested merely 30 different model configurations for every combination of training type and loss function. It is therefore both feasible and beneficial to work with a large search space.

Best configurations (Bayesian optimization). We already obtained good configurations solely from quasi-random search. More precisely, roughly half of our best models (as reported in Table 2) were found during quasi-random search, the other half benefited from subsequent tuning with Bayesian optimization. But even in the latter cases, the performance on test data almost always improved only marginally. A notable exception is TransE, for which Bayesian optimization improved every metric on both datasets by about 2 percentage points (4 points for Hits@10 on WNRR). The best hyperparameters after the Bayesian optimization phase are reported in Table 8 in the appendix.

Impact of specific hyperparameters. It is difficult to assess the impact of each hyperparameter individually. As a proxy for the importance of each hyperparameter setting, we report in Table 3 and

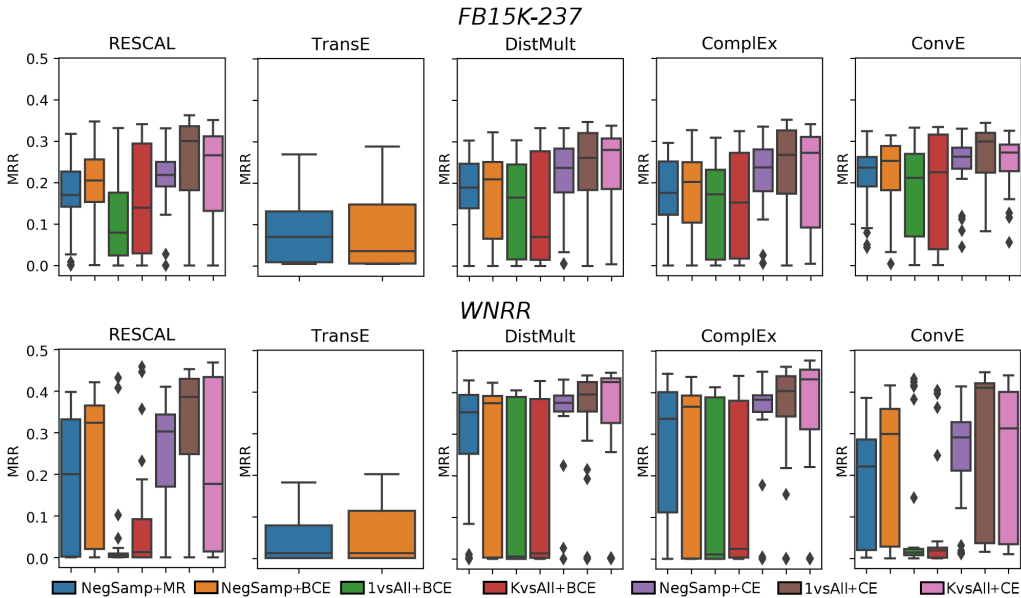


Figure 2: Distribution of filtered MRR (on validation data, quasi-random search only) for different training type and loss functions

Tables 6 and 7 (appendix) the best performance in terms of filtered MRR on validation data for a *different choice* of value for each hyperparameter (difference in parenthesis). For example, the best configuration during quasi-random search for ComplEx on FB15K-237 used reciprocal relations. The best configuration that does not use reciprocal relations has an MRR that is 0.8 points lower (34.3 instead of 35.1). This does not mean that the gap is explained by the use of the reciprocal relations alone, as other hyperparameters may also be different, but it does show the best performance we obtained without relying on reciprocal relations. We can see that most hyperparameters appear to have moderate impact on model performance: no hyperparameter setting produced a significant performance improvement that others did not also produce. The clear exception here is the use of the loss function, which is consistent across models and datasets (always CE), partly with large performance improvements. Another notable exception is the initializer for TransE: in our experiments, TransE performed best without embedding normalization but with uniform Xavier initialization.

Impact of loss function. In order to further explore the impact of the loss function, we show the distribution of the filtered MRR on validation data for every (valid) combination of training type and loss function in Figure 7. We found that the use of CE generally leads to better configurations than using other losses. This is surprising (and somewhat unsatisfying) in that we use the CE loss, which is a multi-class loss, for a multi-label task. Similar observations have been made for computer vision tasks (Joulin et al., 2016; Mahajan et al., 2018) though. Note that the combination of KvsAll with CE (on the very right in figure) has not been explored previously.

4 CONCLUSION

We reported the results of an extensive experimental study with popular KGE model architectures and training strategies across a wide range of hyperparameter settings. We found that when trained appropriately, the relative performance differences between various model architectures often shrunk and sometimes even reversed when compared to prior results. This suggests that (at least currently) training strategies have a significant impact on model performance and may account for a substantial fraction of the progress made in recent years. To enable an insightful model comparison, we recommend that new studies compare models in a common experimental setup using a sufficiently large hyperparameter search space. Our study shows that such an approach is possible with reasonable computational cost. Moreover, we feel that many of the more advanced architectures and techniques proposed in the literature need to be revisited to reassess their individual benefits.

REFERENCES

- Stephan Baier, Yunpu Ma, and Volker Tresp. Improving visual relationship detection using semantic modeling of scene descriptions. In *The Semantic Web - 16th International Semantic Web Conference (ISWC)*, 2017.
- Ivana Balazevic, Carl Allen, and Timothy Hospedales. TuckER: Tensor factorization for knowledge graph completion. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2019.
- Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. Translating embeddings for modeling multi-relational data. In *Advances in Neural Information Processing Systems (NIPS)*, 2013.
- Tim Dettmers, Pasquale Minervini, Pontus Stenetorp, and Sebastian Riedel. Convolutional 2D knowledge graph embeddings. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence (AAAI)*, 2018.
- Luis Antonio Galarraga, Christina Teflioudi, Katja Hose, and Fabian M. Suchanek. AMIE: association rule mining under incomplete evidence in ontological knowledge bases. In *The Web Conference (WWW)*, 2013.
- Shu Guo, Quan Wang, Lihong Wang, Bin Wang, and Li Guo. Knowledge graph embedding with iterative guidance from soft rules. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence (AAAI)*, 2018.
- Armand Joulin, Laurens van der Maaten, Allan Jabri, and Nicolas Vasilache. Learning visual features from large weakly supervised data. In *European Conference on Computer Vision (ECCV)*, 2016.
- Rudolf Kadlec, Ondrej Bajgar, and Jan Kleindienst. Knowledge base completion: Baselines strike back. In *Proceedings of the 2nd Workshop on Representation Learning for NLP (Rep4NLP@ACL)*, 2017.
- Bhushan Kotnis and Vivi Nastase. Analysis of the impact of negative sampling on link prediction in knowledge graphs. In *KBCOM@WSDM*, 2018.
- Timothée Lacroix, Nicolas Usunier, and Guillaume Obozinski. Canonical tensor decomposition for knowledge base completion. In *International Conference on Machine Learning (ICML)*, 2018.
- Mario Lucic, Karol Kurach, Marcin Michalski, Sylvain Gelly, and Olivier Bousquet. Are GANs created equal? A large-scale study. In *Advances in Neural Information Processing Systems (NIPS)*, 2018.
- Dhruv Mahajan, Ross Girshick, Vignesh Ramanathan, Kaiming He, Manohar Paluri, Yixuan Li, Ashwin Bharambe, and Laurens van der Maaten. Exploring the limits of weakly supervised pretraining. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.
- Christian Meilicke, Melisachew Wudage Chekol, Daniel Ruffinelli, and Heiner Stuckenschmidt. Anytime bottom-up rule learning for knowledge graph completion. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence (IJCAI)*, 2019.
- Gábor Melis, Chris Dyer, and Phil Blunsom. On the state of the art of evaluation in neural language models. In *6th International Conference on Learning Representations (ICLR)*, 2017.
- Sameh Mohamed, Vt Novek, Pierre-Yves Vandembussche, and Emir Muoz. Loss functions in knowledge graph embedding models. In *Proceedings of the Workshop on Deep Learning for Knowledge Graphs (DL4KG2019)*, 2019.
- Deepak Nathani, Jatin Chauhan, Charu Sharma, and Manohar Kaul. Learning attention-based embeddings for relation prediction in knowledge graphs. In *Proceedings of the 57th Conference of the Association for Computational Linguistics (ACL)*, 2019.

- Dai Quoc Nguyen, Tu Dinh Nguyen, Dat Quoc Nguyen, and Dinh Phung. A novel embedding model for knowledge base completion based on convolutional neural network. In *Proceedings of the 16th Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, 2018.
- Maximilian Nickel, Volker Tresp, and Hans-Peter Kriegel. A three-way model for collective learning on multi-relational data. In *Proceedings of the 28th International Conference on Machine Learning (ICML)*, 2011.
- Maximilian Nickel, Kevin Murphy, Volker Tresp, and Evgeniy Gabrilovich. A review of relational machine learning for knowledge graphs. *Proceedings of the IEEE*, 2015.
- Nils Reimers and Iryna Gurevych. Reporting score distributions makes a difference: Performance study of LSTM-networks for sequence tagging. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2017.
- Farnood Salehi, Robert Bamler, and Stephan Mandt. Probabilistic knowledge graph embeddings. In *Symposium on Advances in Approximate Bayesian Inference*, 2018.
- Chao Shang, Yun Tang, Jing Huang, Jinbo Bi, Xiaodong He, and Bowen Zhou. End-to-end structure-aware convolutional networks for knowledge base completion. In *The Thirty-Third AAAI Conference on Artificial Intelligence (AAAI)*, 2019.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15, 2014.
- Zhiqing Sun, Zhi-Hong Deng, Jian-Yun Nie, and Jian Tang. RotatE: Knowledge graph embedding by relational rotation in complex space. In *7th International Conference on Learning Representations (ICLR)*, 2019.
- Kristina Toutanova and Danqi Chen. Observed versus latent features for knowledge base and text inference. In *Proceedings of the 3rd Workshop on Continuous Vector Space Models and their Compositionality*, 2015.
- Théo Trouillon, Johannes Welbl, Sebastian Riedel, Éric Gaussier, and Guillaume Bouchard. Complex embeddings for simple link prediction. In *Proceedings of the 33rd International Conference on Machine Learning (ICML)*, 2016.
- Quan Wang, Zhendong Mao, Bin Wang, and Li Guo. Knowledge graph embedding: A survey of approaches and applications. *IEEE Transactions on Knowledge and Data Engineering*, 2017.
- Yanjie Wang, Daniel Ruffinelli, Rainer Gemulla, Samuel Broscheit, and Christian Meilicke. On evaluating embedding models for knowledge base completion. In *Proceedings of the 4th Workshop on Representation Learning for NLP (Rep4NLP@ACL)*, 2019.
- Bishan Yang, Scott Wen-tau Yih, Xiaodong He, Jianfeng Gao, and Li Deng. Embedding entities and relations for learning and inference in knowledge bases. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2015.

A APPENDIX

Dataset	Entities	Relations	Train	Validation	Test
FB-237	14 505	237	272 115	17 535	20 466
WNRR	40 559	11	86 835	2 824	2 924

Table 4: Dataset statistics

Hyperparameter	Quasi-random search	Bayesian optimization
Embedding size	{128, 256, 512}	Value from Tab. 6
Training type	{NegSamp, 1vsAll, KvsAll}	Value from Tab. 6
Reciprocal	{True, False}	Value from Tab. 6
No. subject samples (NegSamp)	[1, 1000], log scale	[1, 1000], log scale
No. object samples (NegSamp)	[1, 1000], log scale	[1, 1000], log scale
Label smoothing (KvsAll)	[0.0, 0.3]	[0.0, 0.3]
Loss	{BCE, MR, CE}	Value from Tab. 6
Margin (MR)	[0, 10]	[0, 10]
L_p -norm (TransE)	{1, 2}	Value from Tab. 6
Optimizer	{Adam, Adagrad}	Value from Tab. 6
Batch size	{128, 256, 512, 1024}	Value from Tab. 6
Learning rate	$[10^{-4}, 1]$, log scale	$[10^{-4}, 1]$, log scale
LR scheduler patience	[0, 10]	[0, 10]
L_p regularization	{L1, L2, L3, None}	Value from Tab. 6
Entity emb. weight	$[10^{-20}, 10^{-5}]$	$[10^{-20}, 10^{-5}]$
Relation emb. weight	$[10^{-20}, 10^{-5}]$	$[10^{-20}, 10^{-5}]$
Frequency weighting	{True, False}	Value from Tab. 6
Embedding normalization (TransE)		
Entity	{True, False}	Value from Tab. 6
Relation	{True, False}	Value from Tab. 6
Dropout		
Entity embedding	[0.0, 0.5]	[0.0, 0.5]
Relation embedding	[0.0, 0.5]	[0.0, 0.5]
Feature map (ConvE)	[0.0, 0.5]	[0.0, 0.5]
Projection (ConvE)	[0.0, 0.5]	[0.0, 0.5]
Embedding initialization	{Normal, Unif, XvNorm, XvUnif}	Value from Tab. 6
Std. deviation (Normal)	$[10^{-5}, 1.0]$	$[10^{-5}, 1.0]$
Interval (Unif)	$[-1.0, 1.0]$	$[-1.0, 1.0]$
Gain (XvNorm)	1.0	1.0
Gain (XvUnif)	1.0	1.0

Table 5: Hyperparameter search space used in our study. Settings that apply only to certain configurations are indicated in parenthesis.

	RESICAL	TransE	DistMult	ComplEx	ConvE
MRR (valid)	36.2	28.8	34.7	35.1	34.4
Embedding size	512 (-0.5)	256 (-1.1)	128 (-0.1)	512 (-0.1)	512 (-1.0)
Training type	1vsAll (-1.1)	NegSamp	1vsAll (-0.9)	1vsAll (-1.0)	1vsAll (-1.0)
Reciprocal	No (-0.5)	No	Yes (-0.5)	Yes (-0.8)	Yes
No. subject samples (NegSamp)	-	45	-	-	-
No. object samples (NegSamp)	-	332	-	-	-
Label smoothing (KvsAll)	-	-	-	-	-
Loss	CE (-1.5)	CE (-1.9)	CE (-1.4)	CE (-2.4)	CE (-1.0)
Margin (MR)	-	-	-	-	-
L_p -norm (TransE)	-	1	-	-	-
Optimizer	Adag. (-0.9)	Adam (-1.1)	Adag. (-0.7)	Adag. (-1.1)	Adag. (-1.0)
Batch size	512 (-0.5)	1024 (-1.1)	1024 (-0.1)	512 (-0.1)	1024 (-0.5)
Learning rate	0.03	0.0008	0.22	0.30	0.001
Scheduler patience	10 (-0.5)	6 (-1.1)	5 (-0.1)	7 (-0.1)	5 (-0.5)
L_p regularization	None (-0.5)	L2 (-1.1)	L3 (-0.5)	L3 (-0.8)	L3 (-0.5)
Entity emb. weight	2.00^{-10}	2.05^{-13}	1.40^{-09}	1.65^{-13}	4.09^{-10}
Relation emb. weight	4.44^{-11}	2.19^{-06}	4.51^{-03}	3.05^{-20}	3.20^{-11}
Frequency weighting	Yes (-0.5)	Yes (-1.1)	No (-0.7)	No (-0.5)	No (-0.5)
Embedding normalization (TransE)	-	-	-	-	-
Entity	-	No	-	-	-
Relation	-	No	-	-	-
Dropout	-	-	-	-	-
Entity embedding	0.37	0.0	0.29	0.0	0.0
Relation embedding	0.21	0.0	0.44	0.41	0.0
Projection (ConvE)	-	-	-	-	0.26
Feature map (ConvE)	-	-	-	-	0.13
Embedding initialization	Unif. (-0.9)	XvUnif (-4.4)	Unif. (-0.1)	XvUnif (-0.1)	Normal (-0.5)
Std. deviation (Normal)	-	-	-	-	0.0003
Interval (Unif)	[-0.15, 0.15]	-	[-0.45, 0.45]	-	-

FB15K-237

Table 6: Hyperparameters of best performing models found with quasi-random hyperparameter optimization on FB15K-237. Settings that apply only to certain configurations are indicated in parenthesis. For each hyperparameter, we also give the reduction in filtered MRR for the best configuration that does not use this value of the hyperparameter (in parenthesis).

	RESICAL	TransE	DistMult	CompLex	ConvE
MRR (valid)	46.8	20.2	44.7	47.7	44.7
Embedding size	256 (-0.8)	256 (-1.8)	512 (-0.6)	512 (-1.3)	128 (-0.7)
Training type	KvsAll (-1.6)	NegSamp	KvsAll (-0.6)	KvsAll (-1.5)	1vsAll (-0.7)
Reciprocal	No (-1.0)	No	Yes (-0.6)	Yes (-1.4)	Yes
No. subject samples (NegSamp)	-	1	-	-	-
No. object samples (NegSamp)	-	243	-	-	-
Label smoothing (KvsAll)	0.26	-	0	0	-
Loss	CE (-1.0)	CE (-2.0)	CE (-1.8)	CE (-3.3)	CE (-1.5)
Margin (MR)	-	-	-	-	-
L_p -norm (TransE)	-	1	-	-	-
Optimizer	Adag. (-0.8)	Adam (-1.8)	Adam (-0.6)	Adam (-1.3)	Adam (-1.3)
Batch size	1024 (-0.8)	1024 (-1.8)	1024 (-0.1)	1024 (-1.2)	256 (-0.7)
Learning rate	0.06	0.0008	0.0004	0.0004	0.001
Scheduler patience	4 (-0.8)	6 (-1.8)	1 (-0.1)	1 (-1.2)	6 (-1.3)
L_p regularization	None (-0.8)	L2 (-1.8)	None (-0.1)	None (-1.2)	L2 (-1.3)
Entity emb. weight	5.37 ⁻¹⁰	2.05 ⁻¹³	8.64 ⁻⁰⁶	8.64 ⁻⁰⁶	4.54 ⁻⁰⁹
Relation emb. weight	4.35 ⁻¹¹	2.19 ⁻⁰⁶	1.00 ⁻¹³	1.00 ⁻¹³	4.14 ⁻¹⁸
Frequency weighting	No (-1.0)	Yes (-1.8)	Yes (-0.1)	Yes (-1.2)	Yes (-0.7)
Embedding normalization (TransE)	-	-	-	-	-
Entity	-	No	-	-	-
Relation	-	No	-	-	-
Dropout	-	-	-	-	-
Entity embedding	0.41	0.0	0.14	0.14	0.0
Relation embedding	0.05	0.0	0.48	0.48	0.0
Projection (ConvE)	-	-	-	-	0.46
Feature map (ConvE)	-	-	-	-	0.40
Embedding initialization	XvUnif (-0.8)	XvUnif (-2.0)	Normal (-0.6)	Normal (-1.3)	XvNorm (-1.3)
Std. deviation (Normal)	-	-	0.00001	0.00001	-
Interval (Unif)	-	-	-	-	-

Table 7: Hyperparameters of best performing models found with quasi-random hyperparameter optimization on WNRR. Settings that apply only to certain configurations are indicated in parenthesis. For each hyperparameter, we also give the reduction in filtered MRR for the best configuration that does not use this value of the hyperparameter (in parenthesis).

	RESCAL	TransE	DistMult	ComplEx	ConvE
<i>FB15K-237</i>					
MRR valid.	36.5	30.7	34.7	35.1	34.7
Learning rate	0.06	0.0003	0.22	0.30	0.002
Scheduler patience	3	6	5	7	9
Entity emb. reg. weight	1.85^{-09}	5.40^{-05}	1.40^{-09}	1.65^{-13}	5.20^{-08}
Relation emb. reg. weight	1.21^{-10}	1.72^{-13}	4.51^{-03}	3.05^{-20}	1.12^{-19}
Entity emb. dropout	0.50	0.08	0.29	0.0	0.0
Relation emb. dropout	0.25	0.0	0.44	0.41	0.0
Projection dropout (ConvE)	–	–	–	–	0.46
Feature map dropout (ConvE)	–	–	–	–	0.32
<i>WNRR</i>					
MRR valid.	47.0	22.5	44.7	47.7	44.7
Learning rate	0.08	0.0004	0.0004	0.0004	0.001
Scheduler patience	0	3	1	1	6
Entity reg. weight	6.96^{-10}	1.00^{-01}	8.64^{-06}	8.64^{-06}	4.54^{-09}
Relation reg. weight	7.76^{-10}	2.93^{-13}	1.00^{-13}	1.00^{-13}	4.14^{-18}
Entity emb. dropout	0.45	0.1	0.14	0.14	0.0
Relation emb. dropout	0.0	0.0	0.48	0.48	0.0
Projection dropout (ConvE)	–	–	–	–	0.46
Feature map dropout (ConvE)	–	–	–	–	0.40

Table 8: Hyperparameters of best performing models found with Bayesian optimization. All other hyperparameters are the same as in Tables 6 (FB15k-237) and 7 (WNRR).

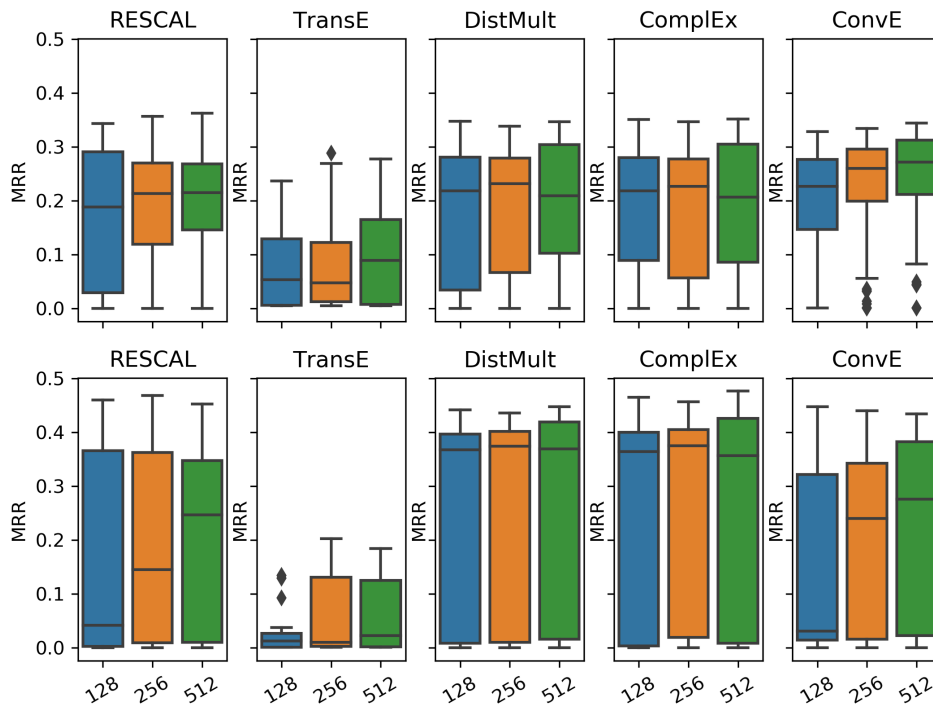


Figure 3: Distribution of filtered MRR (on validation data, quasi-random search only) for different embedding sizes (top row: FB15k-237, bottom row: WNRR)

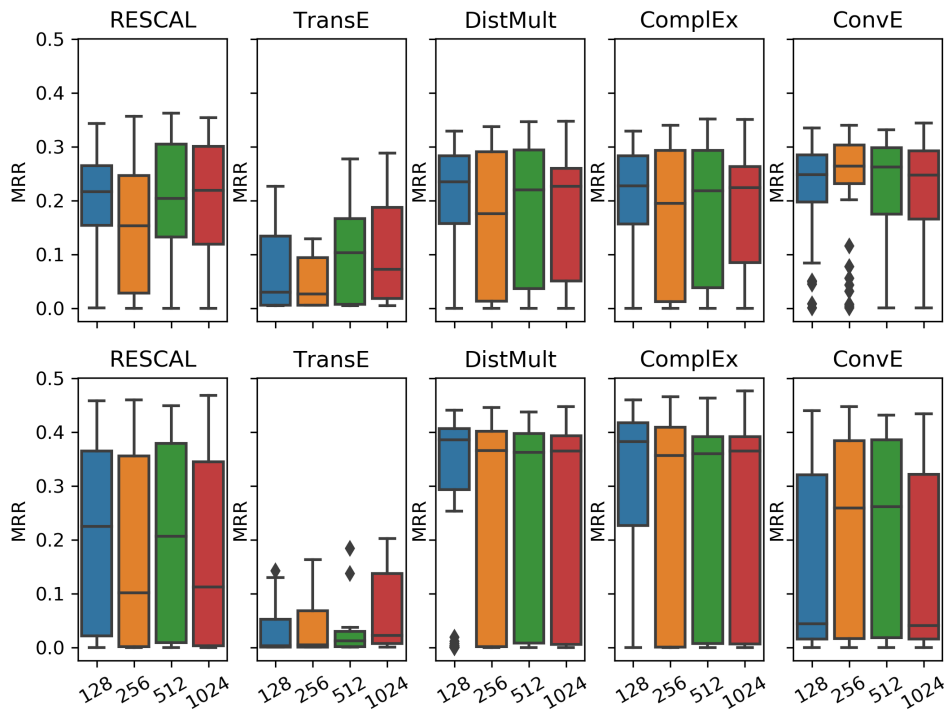


Figure 4: Distribution of filtered MRR (on validation data, quasi-random search only) for different batch sizes (top row: FB15k-237, bottom row: WNRR)

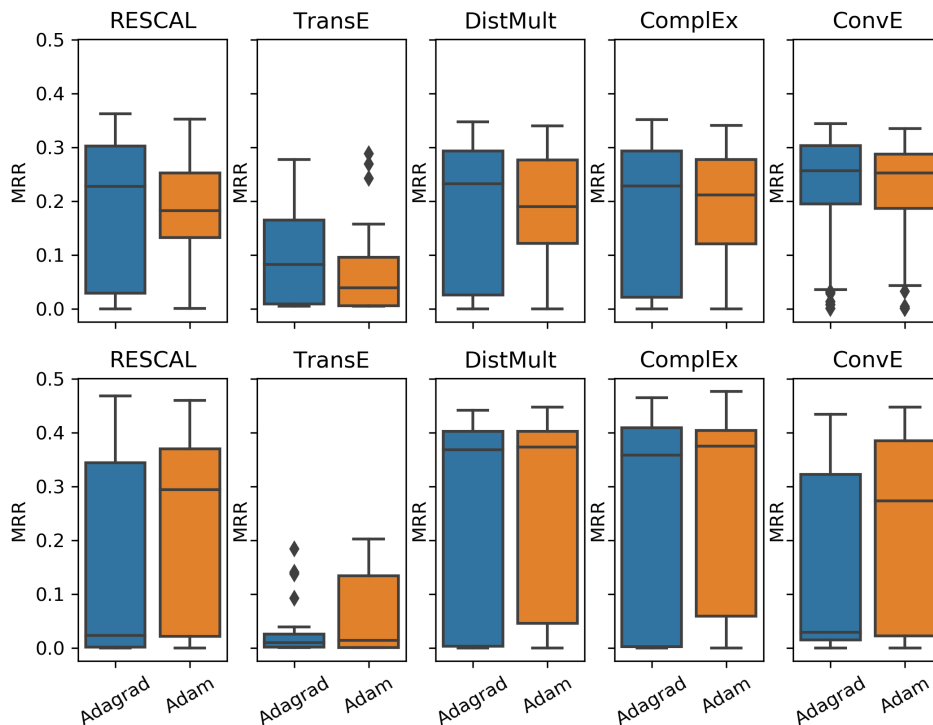


Figure 5: Distribution of filtered MRR (on validation data, quasi-random search only) for different optimizers (top row: FB15k-237, bottom row: WNRR)

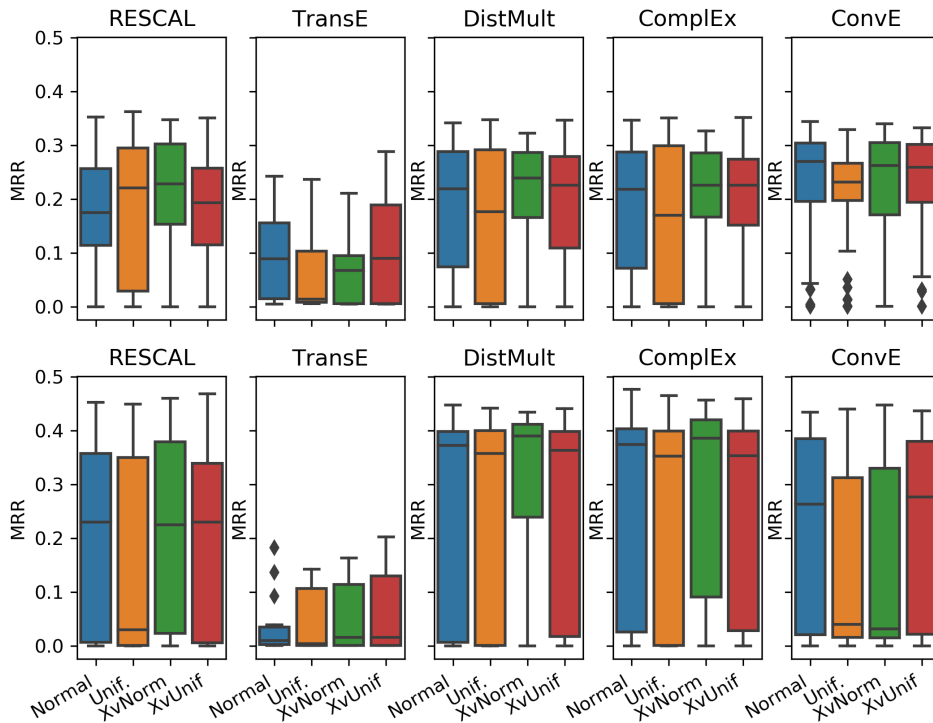


Figure 6: Distribution of filtered MRR (on validation data, quasi-random search only) for different initializers (top row: FB15k-237, bottom row: WNRR)

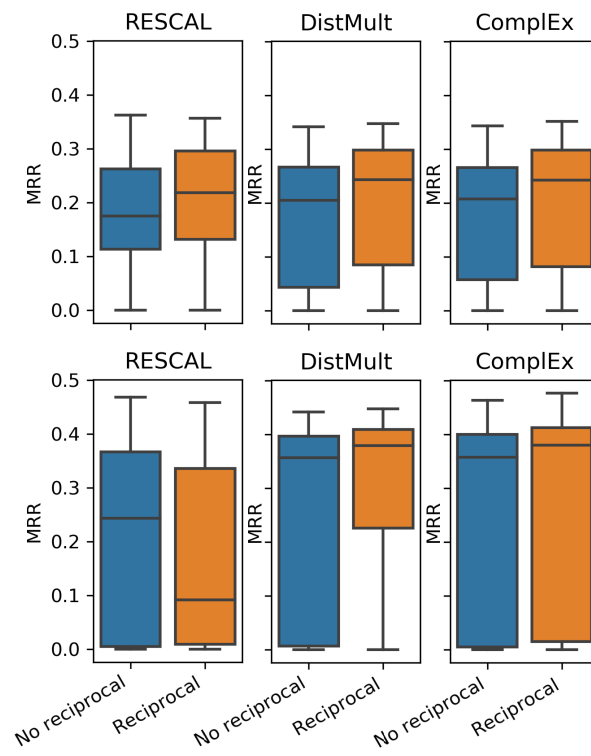


Figure 7: Distribution of filtered MRR (on validation data, quasi-random search only) with and without reciprocal relations (top row: FB15k-237, bottom row: WNRR)