# Efficient Knowledge Graph Unlearning with Zeroth-order Information

Yang Xiao
yax3417@utulsa.edi
The University of Tulsa
Tulsa, Oklahoma, USA

Ruimeng Ye
ruy9945@utulsa.edu
The University of Tulsa
Tulsa, Oklahoma, USA

Bohan Liu
bohanli2@andrew.cmu.edu
Carnegie Mellon University
Pittsburgh, Pennsylvania, USA

Xiaolong Ma
xiaolongma@arizona.edu
The University of Arizona
Tucson, Arizona, USA

Bo Hui
bo-hui@utulsa.edu
The University of Tulsa
Tulsa, Oklahoma, USA

**Figure 1: Influence Function Approximation**

## Abstract

Due to regulations like the Right to be Forgotten, there is growing demand for removing training data and its influence from models. Since full retraining is costly, various machine unlearning methods have been proposed. In this paper, we firstly present an efficient knowledge graph (KG) unlearning algorithm. We remark that KG unlearning is nontrivial due to the distinctive structure of KG and the semantic relations between entities. Also, unlearning by estimating the influence of removed components incurs significant computational overhead when applied to large-scale knowledge graphs. To this end, we define an influence function for KG unlearning and propose to approximate the model's sensitivity without expensive computation of first-order and second-order derivatives for parameter updates. Specifically, we use Taylor expansion to estimate the parameter changes caused by data removal. Given that the first-order gradients and second-order derivatives dominate the computational load, we use the Fisher matrices and zeroth-order optimization to approximate the inverse-Hessian vector product without constructing the computational graphs. Our experimental results demonstrate that the proposed method outperforms other state-of-the-art graph unlearning baselines significantly in terms of unlearning efficiency and unlearning quality. Our code is released at https://github.com/NKUShaw/ZOWFKGIF.

## CCS Concepts

• **Security and privacy → Database and storage security**.
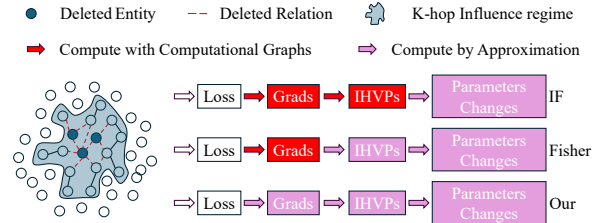
## Keywords

Knowledge Graph; Influence Function; Zeroth Optimization

## 1 Introduction

With massive data being used in the machine learning (ML) model training process, data privacy and security have gained significant attention. Laws such as the European Union's General Data Protection Regulation (GDPR), the California Consumer Privacy Act (CCPA), and Canada's proposed Consumer Privacy Protection Act (CPPA) oversee the use of personal data in machine learning and grant users the *Right to be Forgotten* [1, 51, 52]. This right entitles individuals to request the removal of their personal data from databases, including those used for training ML models. However, traditional ML pipelines are not designed to easily accommodate such deletion requests, as once data is integrated into a model, its influence is deeply embedded in learned parameters through complex and non-linear transformations.

This mismatch between legal requirements and technical capabilities has sparked growing interest in developing Machine Unlearning [4], the techniques that aims to remove the influence of deleted data from trained models upon certified requests without retraining from scratch. While retraining a new model on remaining data is intuitive [25, 53], it is computationally impractical for frequent deletions, especially with large-scale models. To address this, methods for exact or approximate unlearning have emerged [5, 27, 28, 39, 47], targeting specific data types like face images [8], natural language [7, 71], and graphs [10, 60]. These methods promise a scalable solution to comply with privacy regulations while maintaining model utility. Beyond legal compliance, unlearning improves model robustness, removes poisoned or biased data, and supports dynamic updates. However, unlearning on knowledge graph (KG) has not yet been extensively explored. Unlike unstructured data, KGs consist of interdependent triples

that propagate influence through graph structure, posing unique challenges to unlearning. KGs store and organize open and private information across industries, examples include Wikidata and Google Knowledge Graph. Providing unlearning interfaces for KGs is crucial, as sensitive data (e.g., personal healthcare information) can make the continued use of KG models trained on deleted data illegal [23, 58].

Despite the advancements in machine unlearning and graph unlearning [13, 15, 30, 31, 61, 68, 70], applying existing graph unlearning methods directly to large-scale KG models presents computational challenges. The majority of existing approximate unlearning algorithms rely on computing the influence of deleted data on the trained model. To compute the influence of deleted data, constructing a computational graph (gradient-flow graph) is required. The operation over the computational graph often accounts for the largest portion of a KG model's computational load, mainly determined by the size of the KG, presenting computational challenges upon frequent deletion requests on large-scale KGs. Large-scale KGs are commonplace in real-world applications. For example, Google KG covers 500 billion facts (i.e., triplets) on 5 billion entities. This challenge poses a major hurdle to KG unlearning.

In this paper, we propose an efficient KG unlearning method without constructing the gradient-flow graph. Specifically, we start with computing the sensitivity of a KG model's parameters to the removal of a single data point with the product of the Hessian matrix inverse and the first-order gradients [37, 38, 64] . To address the computational challenges of the inverse-Hessian vector products (IHVPs) in the influence function, we first decompose a large matrix into smaller matrices to reduce the memory cost. Then we introduce the Fisher information matrix [42] to approximate IHVPs. In the probabilistic view, the Fisher information matrices of the model's conditional distribution are equivalent to Hessian matrices. With the approximation, we can express the influence function with the inverse of the empirical Fisher matrices. To further reduce the cost of computing gradient vectors required in IHVPs, we introduce zeroth-order optimization [12] to approximate the first-order gradients for KG unlearning. Instead of computing the gradients based on the loss value, we use the model output itself (i.e. zeroth-order information) to approximate the values of first-order gradients. This approach is particularly advantageous in scenarios where gradient information is expensive to compute, such as computing gradients on KG embeddings.

Figure 1 shows the computational flow of (1) traditional influence function (KGIF) without approximation; (2) approximating IHVPs with the Fisher matrices; and (3) approximating IHVPs with the Zeroth-order information. Note that our method can avoid high computational costs over the gradient flow (the computational graph) by approximating with the Zeroth-order information. We verify the efficiency and efficacy of the proposed unlearning method by approximation in the experiment.

We investigate three KG unlearning tasks: triplets unlearning, entity unlearning, and relation unlearning. Considering that the majority of applications on KG such as question answering, knowledge acquisition and recommendation rely heavily on KG embedding [33], we aim to remove the influence of deleted triplets, entities, and relations from the trained KG embedding models upon request. Different from existing graph unlearning methods developed for graph neural network (GNN) models and classification tasks [36, 67], our KG unlearning method optimizes a directional scoring loss to learn representations of entities and relations in triplets [19].

Our high-level contributions are as follows:

- We investigate KG unlearning and use an influence function to forget triplets, entities, and relations.
- We introduce the Fisher matrices to reduce the computational cost of inverse Hessian vector products.
- We further use zeroth-order information to approximate the first-order gradient and reduce the computational cost over the gradient flow.
- We experimentally demonstrate that our approximating method can significantly reduce the run time and memory cost. At the same time, our unlearning method outperforms the state-of-the-art graph unlearning baselines.

## 2 Preliminary

Let $G = (E, R, T)$ be a KG, where $E$ and $R$ are the sets of entities and relations in the KG. We use $T$ to denote the set of triples, each of which is a triplet $(e_h, r, e_t)$, including the head entity $e_h \in E$, the tail entity $e_t \in E$ and the relation $r \in R$ between $e_h$ and $e_t$. In this paper, we explore KG unlearning with embedding models, because the majority of applications on KG such as question answering, Knowledge acquisition, and recommendation rely heavily on KG embedding according to a survey [33]. A KG embedding model $M(G)$ trained on $G$ associates each entity and relation with a vector in an embedding space.

The user can request to delete (1) a subset of triples $T_d$; (2) a subset of relations $R_d$; (3) a subset of entities $E_d$. Denote $G/T_d$ as the new KG where $T_d$ has been removed from $G$. The naive unlearning is to retrain a new model $M(G/T_d)$ on the remaining data $G/T_d$ from scratch. Similarly, $M(G/E_d)$ and $M(G/R_d)$ represent the new models for entity unlearning and relation unlearning. To resolve the computational issue of retraining, the goal of this paper is to develop an efficient unlearning algorithm to directly eliminate the effects of deleted data. Ideally, the unlearning algorithm will result in a new KG model with a similar performance as the retrained one. Considering that a KG is stored as triplets, we removed all triplets including the deleted entity or relation. The schema of the KG enables us to query specific triplets (e.g., foaf: Person, rdf: Is a friend of, foaf: Person), entities, or relations to address privacy or concerns from users.

## 3 Related Works

### 3.1 Machine Unlearning

Given the guidelines such as the Right to be Forgotten, machine unlearning methods aim to remove the impact of deleted data on the trained model [50, 54, 61]. While re-training from scratch on the retraining dataset is the intuitive method, the time and computational costs are not practical in applications. To alleviate the issue, some research efforts have been proposed to unlearn data efficiently. Existing unlearning methods mainly include exact unlearning and approximate unlearning. Exact unlearning aims to learn an unlearning model with the same performance as the retrained one [3, 4, 6, 9, 26, 40, 43, 59]. For example, SISA [2] divides the dataset into several shards, trains multiple sub-models and

establishes an integrated model to reduce the time cost. Approximate unlearning methods try to update the parameters of the trained model closer to the retraining one through the approximation [17, 24, 26, 28, 46, 48, 66]. Currently, there is a new trend to unlearn the specific concepts which are learned in complex data in large language models (LLMs) [41, 71].

## 3.2 Graph Unlearning

Different from general unlearning independent instances, deleting nodes and edges on a graph will have effects on other connected nodes and edges. Graph unlearning aims to approximate updates on the embeddings of surrounding nodes or edges [14, 20, 64, 65]. GraphEraser [11] segment balanced shards reducing the destructiveness of the graph structure. However, the key point of this method lies in the uncontrollability of the number of clusters, hence the efficiency varies in different datasets. GIF [64] uses an additional loss term of the influenced neighbors to update the graph model. However, it explicitly computes the first derivative of the loss which is a huge computational burden. [70] analyze the schema patterns in knowledge graph and unlearn data based on cosine similarity of subgraphs within the schema. GNNDELETE [13] studies GNNs unlearning and proposes two standards for optimization: Deleted Edge Consistency and Neighborhood Influence. Some other works leverage KG to assist unlearning in the context of federated Learning [74] or removing harmful content from LLMs [22, 44]. To the best of the authors' knowledge, our paper is the first work to study unlearning entities or relations on KG models.

## 4 Methodology

### 4.1 Knowledge Graphs Influence Functions

We first define an influence function in the KG embedding model to assess the variation in model parameters when a sample is marginally increased in weight. The benefit of our influence function is that it is universal for all three types of unlearning tasks (triplet unlearning, entity unlearning, and relation unlearning) since all three tasks land on deleting triplets for a KG and the essence is to measure the impact on the KG embeddings. This is different from general graph unlearning where the optimization target variies across different tasks (e.g., deleted edge consistency [13], classification, and link prediction).

Without loss of generality, a KG embedding model determines whether there is a specific type of relation between any two entities [19]. Denote $\theta$ as the parameters in the KG embedding model $M(G)$. Let $f(e_h, r, e_t; \theta)$ be the score of $(e_h, r, e_t)$ computed with $\theta$ and the embeddings. Given a positive training set of triplets $(e_h, r, e_t) \in T^+$ and a negative sample set $(e_h', r', e_t'; \theta) \in T^-$, an example of the loss function (TransH [62]) can be written as:

$$\mathcal{L} = \sum_{T^+} \sum_{T^-} f(e_h, r, e_t; \theta) - f(e_h', r', e_t'; \theta) + \gamma, \tag{1}$$

where $\gamma$ is the margin separating positive and negative triplets. The score $f(\cdot)$ is expected to be lower for a fact triplet and higher for an incorrect triplet.

We remark that our unlearning method is model-agnostic and works for any KG embedding models trained based on a loss function. Given any model trained with a loss function $\mathcal{L}$, the model

parameters $\theta$ can be optimized by minimizing:

$$\hat{\theta} = \arg\min \mathcal{L} \tag{2}$$

where $\hat{\theta}$ be the optimized parameters trained on $G$. Suppose that $T_d$ is the set of triplets to be deleted. To estimate the influence of $T_d$, we add a small perturbation caused by data removal (e.g., delete a small set of triplets) to the loss:

$$\tilde{\theta} = \arg\min(\mathcal{L}_{\theta_0} + \epsilon \mathcal{L}_{T_d}), \tag{3}$$

where $\epsilon$ is a scalar. For simplicity, we use $\mathcal{L}_{T_d}$ to represent the loss value concerning all triplets in $T_d$. For example, $\mathcal{L}_{T_d} = \sum_{T_d} \sum_{T^-} f(e_h, r, e_t; \theta) - f(e_h', r', e_t'; \theta) + \gamma$ in TransH with new negative samples $T^-$ for $T_d$. Note that $\tilde{\theta}$ is optimized parameters for Eq. (3). Therefore, we have the first-order optimal conditions as follows:

$$0 = \nabla_{\hat{\theta}} \mathcal{L}, \ 0 = \nabla_{\tilde{\theta}} \mathcal{L}_0 + \epsilon \nabla_{\tilde{\theta}} \mathcal{L}_{(T_d)}. \tag{4}$$

These two conditions are used to ensure that local minima of the loss function are found under different circumstances.

When the permutation is small enough ($\epsilon \to 0$), the perturbed parameters $\tilde{\theta}$ will converge towards $\tilde{\theta} \approx \hat{\theta}$. We use a first-order Taylor expansion to approximate the behavior of the parameter changes under perturbation:

$$0 \approx (\tilde{\theta} - \hat{\theta}) \left( \nabla_{\hat{\theta}}^2 \mathcal{L} + \epsilon \nabla_{\hat{\theta}}^2 \mathcal{L}_{T_d} \right) + \left( \nabla_{\hat{\theta}} \mathcal{L}_0 + \epsilon \nabla_{\hat{\theta}} \mathcal{L}_{T_d} \right) \tag{5}$$

Recall that our task is to eliminate the impact of $T_d$ on the training process of $\theta$. Therefore, we assign the value of $\epsilon$ as -1 to cancel out the effect of $T_d$. Recall that $0 = \nabla_{\hat{\theta}} \mathcal{L}$ from Eq. (4) and the second-order derivatives $\nabla_{\hat{\theta}}^2$ can be omitted because $T_d$ is a small part of $G$. For the purpose of simplicity, we use $H_{\hat{\theta}}$ to represent the Hessian Matrix $\nabla_{\hat{\theta}}^2 \mathcal{L}$ and $v$ to denote the first-order gradient $\nabla_{\hat{\theta}} \mathcal{L}_{T_d}$. Then when $\epsilon = -1$, we have:

$$\tilde{\theta} - \hat{\theta} \approx H_{\hat{\theta}}^{-1} v. \tag{6}$$

Intuitively, deleting an entity or a relation will have an impact on the embeddings of the entities and relationships in the neighborhood. The impact on these entities far from the deleted triplets on the topology is marginal. Therefore, updating embeddings by query entities and relations within the K-hop neighborhood around the deleted component can reduce the cost of updating embeddings of large-scale KG.

We remark that the largest computational load is to calculate the inverse of the Hessian matrix and the inverse-Hessian vector products (IHVPs) in Eq. (6). Specifically, Eq. (6) entails a computational complexity of $O(\|\theta\|^3 + n\|\theta\|^2)$ and requires memory allocation of $O(\|\theta\|^2)$. It is impractical to directly compute Eq. (6) in real-world applications due to the cost. For example, there are 5 billion entities in Google Knowledge Graph. Suppose the length of the embeddings is 128. There will be 128×5 billion parameters to be updated for entities, which may impede the implementation of the unlearning algorithm in practice. By leveraging the stochastic estimation method, we can reduce the complexity and memory to $O(n\|\theta\|^2)$ and $O(\|\theta\|)$ by iterating the following equation until convergence [18]:

$$H_s^{-1} = (I - H_{\hat{\theta}}) H_{s-1}^{-1} + I, \tag{7}$$

where $I$ is the identity matrix. The IHVPs will be updated through the Taylor series expansion:

$$H_s^{-1}v = v + H_{s-1}^{-1}v - \eta H_{\hat{\theta}}H_{s-1}^{-1}v. \tag{8}$$

The coefficient $\eta$ is used to control the convergence condition. Note that $v$ remains unchanged in the iteration, though it is frequently used for computation. We can store it at the end of the model training. These calculations can be performed using the PyTorch Autograd Engine.

## 4.2 Approximation with the First-order Information

Although we have reduced the computational complexity through [18], computing the gradient and IHVPs still incurs significant computational overhead. As its computational complexity is determined by the gradient size of the knowledge graph model, unlearning on a large scale KG is time-consuming and memory-consuming due to the involved large computation graphs for BP (back-propagation). In the experiment, we show that directly applying the influence functions for KG unlearning takes a substantial amount of time and consumes a substantial amount of GPU memory.

Considering that the KG size is usually huge in real applications, computing the product of first-order (gradients) and second-order (Hessian Matrix Inverse) in Eq. (6) requires the construction of a complete computation graph (gradient flow) to track gradients for parameter updates.

To speed up IHVPs computation and reduce memory complexity, we introduce Woodbury theorem [63] to allow cheap computation of inverses. Woodbury matrix identity says that the inverse of a rank-k correction of some matrix can be computed by doing a rank-k correction to the inverse of the original matrix. Specifically, it decomposes a large matrix into much smaller matrices, allowing the influence functions to operate on smaller matrices, thereby reducing the computational cost. Formally, the inverse of a large matrix can be computed as:

$$(A + U\Lambda^{\top})^{-1} = A^{-1} - A^{-1}U(I + \Lambda^{\top}A^{-1}U)^{-1}\Lambda^{\top}A^{-1} \tag{9}$$

In our setting, we have $H \approx U\Lambda^{\top} + A$, where (1) $A$ usually represents a fundamental matrix that is invertible, positively definite, or at least non-singular; (2) $U$ usually denotes a group of directions or factors, which will be added to the basic matrix $A$ to modify or adjust some characteristics of $A$, such as updated parameters in KGs models, (3) $\Lambda$ is used in conjunction with $U$ to indicate how modifications to $A$ are balanced across different dimensions or directions. Note that decomposition provides a thorough understanding of how the change in each dimension of model parameters will impact prediction outcomes.

Next, we introduce Fisher Information Matrix (FIM) [42] to efficiently estimate the inverse Hessian instead of using direct inversion techniques. FIM offers a powerful method to compute a faithful and efficient estimate of the inverse Hessian.

Formally, FIM quantifies the amount of information that an observable random function $\mathcal{L}$ carries about unknown parameters $\theta$ upon which the model distribution depends. The definition of the FIM is given by:

$$\mathcal{F} = \mathbb{E}\left[\nabla_{\theta}\log p(\mathcal{L};\theta)\nabla_{\theta}\log p(\mathcal{L};\theta)^{\top}\right] \tag{10}$$

where $\nabla_{\theta}\log p(\mathcal{L};\theta)$ is the gradient of the log-likelihood of the model concerning the parameters $\theta$, $p(\mathcal{L};\theta)$ is the density function corresponding to the model distribution, and $\mathbb{E}$ denotes the expectation over the distribution of $\mathcal{L}$ parameterized by $\theta$. It can be further written as: $\mathcal{F} = \mathbb{E}_{P_{\theta,\mathcal{L}}}\left[-\nabla^2 \log p(\theta, \mathcal{L})\right]$.

Since the probability of different relations between two entities in KG can be considered as a conditional distribution $P_{\mathcal{L}|\theta}$ over the model output (i.e., embeddings), it can be proved that the Fisher and Hessian matrix $H$ are equivalent by expressing $P_{\mathcal{L}|\theta} = Q_{\theta}P_{\mathcal{L}|\theta} \approx \hat{Q}_{\theta}P_{\mathcal{L}|\theta}$ under the assumption that the model's conditional distribution $P_{\mathcal{L}|\theta}$ matches the conditional distribution of the training data $\hat{Q}_{\mathcal{L}|\theta}$ [42, 56].

Then by combining two approximation methods (i.e., Woodbury theorem and Fisher Information Matrix), the approximation is equivalent to a quasi-Newton approximation [29] which can be written as:

$$H \approx U\Lambda^{\top} + A \approx \mathcal{F} + A \approx vv^{\top} + \gamma I \tag{11}$$

We can then easily obtain the IHVPs instead of using direct inversion techniques:

$$H_0^{-1}v \approx \gamma^{-1}v + \frac{\gamma^{-2}vv^{\top}}{1 + \gamma^{-1}v^{\top}v}v \tag{12}$$

In Eq.( 12), $\gamma > 1$ is a damping term. Recall that we update $H^{-1}$ with $H_s^{-1}v = v + H_{s-1}^{-1}v - H_{\hat{\theta}}[H_{s-1}^{-1}]v$ in each iteration. We have:

$$\begin{aligned}[H_1^{-1}]v =& v + (\gamma^{-1}v + \frac{\gamma^{-2}vv^{\top}}{1 + \gamma^{-1}v^{\top}v}v)v \\ &- \eta H_{\hat{\theta}}(\gamma^{-1}v + \frac{\gamma^{-2}vv^{\top}}{1 + \gamma^{-1}v^{\top}v}v)v,\end{aligned} \tag{13}$$

Iteratively, we can optimize $\theta$ with Eq. (6). At this point, we have converted the second-order information in Eq.( 6) to first-order information with an approximation. It significantly reduces the computational complexity and the memory cost for KG unlearning.

## 4.3 Approximation with the Zeroth-order Information

Note that we still need to compute the gradient $v$ in Eq.( 12) by constructing a computational graph, though the approximation with the first-order information avoids the direct computation of the inverse Hessian matrix $H^{-1}$ and IHVPs. To further reduce the computational cost, we approximate the first-order gradients with zeroth-order information motivated by [45, 72].

Specifically, we approximate the gradient with the loss function values at a point by leveraging the Taylor series expansion. In general, the central difference is more accurate in providing the derivative of a function than other forward or backward difference approximation methods. By conducting Second-order Taylor expansion on $\mathcal{L}$, we have the central difference:

$$\mathcal{L}(\hat{\theta} + \varepsilon) = \mathcal{L}(\hat{\theta}) + \varepsilon\mathcal{L}'(\hat{\theta}) + \frac{\varepsilon^2}{2}\mathcal{L}''(\hat{\theta}) + O(\varepsilon^3) \tag{14}$$

$$\mathcal{L}(\hat{\theta} - \varepsilon) = \mathcal{L}(\hat{\theta}) - \varepsilon\mathcal{L}'(\hat{\theta}) + \frac{\varepsilon^2}{2}\mathcal{L}''(\hat{\theta}) + O(\varepsilon^3) \tag{15}$$

$$v \approx \frac{\mathcal{L}(\hat{\theta} + \varepsilon) - \mathcal{L}(\hat{\theta} - \varepsilon)}{2\varepsilon} = \mathcal{L}'(\hat{\theta}) + O(\varepsilon^2), \tag{16}$$

where the error term is $O(\varepsilon^2)$. Note also that the forward difference and backward difference are given by:

$$\mathcal{L}(\hat{\theta} + \varepsilon) = \mathcal{L}(\hat{\theta}) + \varepsilon\mathcal{L}'(\hat{\theta}) + O(\varepsilon^2) \tag{17}$$

$$v \approx \frac{\mathcal{L}(\hat{\theta} + \varepsilon) - \mathcal{L}(\hat{\theta})}{\varepsilon} = \mathcal{L}'(\hat{\theta}) + O(\varepsilon) \tag{18}$$

$$\mathcal{L}(\hat{\theta} - \varepsilon) = \mathcal{L}(\hat{\theta}) - \varepsilon\mathcal{L}'(\hat{\theta}) + O(\varepsilon^2) \tag{19}$$

$$v \approx \frac{\mathcal{L}(\hat{\theta}) - \mathcal{L}(\hat{\theta} - \varepsilon)}{\varepsilon} = \mathcal{L}'(\hat{\theta}) + O(\varepsilon) \tag{20}$$

Where the error term is $O(\varepsilon)$.

When $\varepsilon$ is small ($< 1$), $\varepsilon^2$ is much smaller than $\varepsilon$. Therefore, the gradient approximated by the central difference in Eq. (16) is more accurate than the forward difference $\frac{\mathcal{L}(\hat{\theta}+\varepsilon) - \mathcal{L}(\hat{\theta})}{\varepsilon}$ and the backward difference $\frac{\mathcal{L}(\hat{\theta}) - \mathcal{L}(\hat{\theta}-\varepsilon)}{\varepsilon}$. Then we estimate the first-order information with:

$$\mathcal{L}'(\hat{\theta}) \approx \frac{\mathcal{L}(\hat{\theta} + \varepsilon) - \mathcal{L}(\hat{\theta} - \varepsilon)}{2\varepsilon}, \tag{21}$$

which can avoid the construction of the computational graphs to compute the gradients. By replacing the gradient $v$ in Eq.( 6) with the zeroth-order information in Eq. (21), our unlearning algorithm only utilizes the zeroth-order information without computing the gradients and the inverse Hessian of the parameters.

**Discussion.** In summary, our unlearning algorithm updates the KG model with the IHVPs approximation and the gradient approximation based on deleted triplets. While all existing approximate unlearning algorithms are estimations of retraining, our approximation is efficient and close to retraining, which can be verified in the experiment. Our algorithm has three important hyperparameters: iterations $s$, damping term $\gamma$, scaling down term, and noise term $\varepsilon$. The value of $s$ decides the number of iterations of IHVPs. The damping term $\gamma$ is used to reduce the amplitude of parameter variation estimation. The scaling-down term $\eta$ is used as the denominator of the parameter change estimation vector to control the magnitude of parameter changes. The noise term $\varepsilon$ is used to perturb the intensity of parameter estimation gradients. We will investigate their impacts in the experiment.

## 5 Experiments

### 5.1 Experimental Setup

In the experiment, we introduce three popular KG models RotatE [57], TransD [32], TransH [62]. Note that our unlearning method can be applied to any KG embedding method based on a smooth loss function. We trained the KG embedding models on benchmark datasets: FB15K237 and YAGO3-10. We randomly select 5% triplets, entities, and relations to conduct triplet unlearning, entity unlearning, and relations unlearning, respectively. To verify the superiority of our proposed algorithm, we introduce the state-of-the-art baselines for comparison: GNNDelete [14], Certified Graph Unlearning [16], GIF [64]. Note that all existing graph unlearning can not be directly used for KG unlearning. We modify the loss function in the baseline to adapt KG models. The number of epochs is set to 1000. We use the Adam optimizer [35] for RotatE and the SGD optimizer for the other two models (TransD and TransH). The default value of the iteration term is 100, the damping term is 1.00, the scaling down term is 10, and the noise is 1e-5. The number

of iterations for influence functions is set as 100 by default. All experiments are run on two Nvidia L40S GPUs.

**Evaluation.** To verify the efficiency of our unlearning algorithm, we report the runtime of retraining and all unlearning algorithms. We also recorded the peak graphics memory usage in the unlearning process. These metrics collectively provide insights into the resources needed for unlearning, and assessing its feasibility in practical scenarios.

Retraining is usually used as the gold standard for evaluating the machine unlearning algorithm [49]. To measure the forget quality, we report the performance of the KG models. With optimal unlearning, the KG performance after unlearning should be close to retraining, as desired. In this paper, we report the performance of the most important downstream task: knowledge graph completion. This task is to predict potential yet unidentified edges (relationships not explicitly labeled) in the knowledge graph, thereby improving the quality and reliability of the KG database. We use four evaluation metrics: Mean Reciprocal Rank (MRR), Hit@10, Hit@3, and Hit@1. MRR calculates the reciprocal rank of the correct answer in the prediction results for each triplet in the test set and then averages these reciprocal ranks $MRR = \frac{1}{N}\Sigma_{i=1}^{N}\frac{1}{rank_i}$. Hit@k evaluates the percentage at which the top-k prediction results contain the correct triplet in the test set: $Hit@k = \frac{1}{N}\Sigma_{i=1}^{N}\mathbf{1}(rank_i \leq k)$, where $\mathbf{1}$ refers indicator function.

### 5.2 Results and Analysis.

*5.2.1 Results.* Talbe 1 presents the comparison between our algorithm and baselines in terms of efficiency and quality for triplet unlearning. Our unlearning KG algorithm based on approximation with the Fisher matrices and Zeroth-order information (ZeroFisher) significantly outperforms the baselines in terms of runtime and memory usage. We also include one variant of our algorithms which computes the gradients directly instead of approximating with zeroth-order information. Specifically, the training time of ZeroFisher is the smallest among all methods. For example, the runtime of ZeroFisher for RotatE is 0.48s while the runtime of the best baseline (GIF) is 9.3s. Also, our memory usage is only about 25% of GIF across all KG models and datasets. This is because our method is the only one that does not require constructing the computational graph to compute the derivatives. While retraining is the gold standard to measure the unlearning quality, our method archives the closest performance to the retraining result in terms of most metrics. Table ?? and 3 are results of entity and relation unlearning. It shows similar conclusion. The ZeroFisher method is memory-efficient, with significantly lower memory usage compared to other methods. Although some metrics in some method dataset combinations were slightly lower than the optimal method, the performance of the ZeroFisher method was still very close to the optimal value. This indicates that our method can provide good performance while maintaining high efficiency. A reasonable explanation is that our approximation leverages the empirical Fisher matrices which are calculated based on the distribution of the remaining data. The overall results indicate that our method is efficient while guaranteeing the unlearning quality.

*5.2.2 Ablation Study.* To verify the effectiveness of each proposed approximation method, we conduct an ablation study. Specifically,

| Model | Method | FB15K237 | | | | | YAGO3 | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Time | Memory | MRR | Hit@3 | Hit@1 | Time | Memory | MRR | Hit@3 | Hit@1 |
| RotatE | Train | 1666s | 932MB | 0.3165 | 0.3506 | 0.2251 | 11049s | 4493MB | 0.3381 | 0.3846 | 0.2356 |
| | Retrain | 1508s | 849MB | 0.2806 | 0.3135 | 0.1880 | 9619s | 4103MB | 0.2953 | 0.3444 | 0.1912 |
| | GNNDelete | 347s | 804MB | 0.2433 | 0.2720 | 0.1522 | 1441s | **3278MB** | 0.2499 | 0.2896 | 0.1492 |
| | CertifiedGU | 224s | 9336MB | 0.2627 | 0.2938 | 0.1658 | 214s | 38653MB | 0.2623 | 0.3113 | 0.1547 |
| | KGSchema | 94s | 862MB | 0.2548 | 0.2911 | 0.1579 | 1027s | 1749MB | 0.2612 | 0.3079 | 0.1478 |
| | GIF | 9.3s | 6379MB | **0.2713** | 0.3041 | **0.1769** | 12.7s | 25954MB | 0.2656 | 0.3171 | 0.1558 |
| | WF-KGIF (Our) | 0.55s | 1958MB | 0.2712 | 0.3042 | 0.1768 | 0.71s | 8851MB | 0.2660 | 0.3170 | 0.1559 |
| | ZeroFisher (Our) | **0.48s** | **771MB** | 0.2710 | **0.3042** | 0.1757 | **0.62s** | 3357MB | **0.2664** | **0.3176** | **0.1563** |
| TransD | Train | 1702s | 893MB | 0.2928 | 0.3314 | 0.1971 | 8543s | 3642MB | 0.3197 | 0.4542 | 0.2796 |
| | Retrain | 1545s | 806MB | 0.2578 | 0.2936 | 0.1612 | 8236s | 3287MB | 0.3128 | 0.3765 | 0.1983 |
| | GNNDelete | 401s | 806MB | 0.2471 | 0.2785 | 0.1544 | 1884s | 3287MB | 0.2637 | 0.3110 | 0.1563 |
| | CertifiedGU | 613s | 30490MB | **0.2530** | 0.2857 | **0.1562** | - | OOM | - | - | - |
| | KGSchema | 156s | 846MB | 0.2476 | 0.2799 | 0.1576 | 1273S | 1276MB | 0.2679 | 0.3443 | 0.1577 |
| | GIF | 6.7s | 21274MB | 0.2480 | 0.2819 | 0.1503 | 17.6s | 85524MB | 0.2820 | 0.3515 | 0.1569 |
| | WF-KGIF (Our) | 0.51s | 4299MB | 0.2492 | 0.2830 | 0.1517 | 1.67s | 18724MB | 0.2820 | 0.3547 | **0.1615** |
| | ZeroFisher (Our) | **0.49s** | **645MB** | **0.2530** | **0.2858** | **0.1562** | **1.39s** | **2953MB** | **0.2864** | **0.3583** | 0.1593 |
| TransH | Train | 1804s | 773MB | 0.2951 | 0.3315 | 0.2008 | 8421s | 3118MB | 0.4212 | 0.4851 | 0.3108 |
| | Retrain | 1591s | 696MB | 0.2471 | 0.2873 | 0.1447 | 7791s | 2808MB | 0.3212 | 0.3869 | 0.2075 |
| | GNNDelete | 161s | 741MB | 0.2283 | 0.2551 | **0.1421** | 1867s | 3004MB | 0.2659 | 0.3115 | 0.1618 |
| | CertifiedGU | 596s | 29458MB | 0.2536 | 0.2800 | 0.1577 | - | OOM | - | - | - |
| | KGSchema | 77s | 734MB | 0.2533 | 0.2813 | 0.1497 | 736s | 798MB | 0.2837 | 0.3478 | 0.1633 |
| | GIF | 6.5s | 20460MB | 0.2444 | 0.2757 | 0.1491 | 16s | 81935MB | 0.2926 | 0.3633 | **0.1689** |
| | WF-KGIF (Our) | 0.46s | 4023MB | **0.2481** | 0.2800 | 0.1521 | 1.3s | 16670MB | 0.2921 | 0.3653 | 0.1680 |
| | ZeroFisher (Our) | **0.43s** | **524MB** | 0.2505 | **0.2825** | 0.1544 | **1.1s** | **2322MB** | 0.2947 | **0.3691** | 0.1659 |

Table 1: Performance on FB15K237 and YAGO3.



(a) Time Ablation  (b) IHVPs Ablation  (c) Gradient Ablation  (d) Loss Ablation

Figure 2: Ablation study on RotatE



(a) $\gamma$ Term  (b) $\eta$ Term  (c) $\varepsilon$ Term

Figure 3: Varying parameters

| (a) Original | (b) Retraining | (c) GIF, $l_2$: 10.57 | (d) WF-KGIF, $l_2$: 10.56 | (e) ZeroFisher, $l_2$: 10.53 |

**Figure 4: Visualization of embeddings and L2 distance to retraining**



| (a) Varying $\gamma$ | (b) Varying $\epsilon$ | (c) Varying $\eta$ |

**Figure 5: Varying parameters (Above is RotatE, Below is TransD)**

| Model | Method | 5% | | | | | 7% | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Time | Memory | MRR | Hit@3 | Hit@1 | Time | Memory | MRR | Hit@3 | Hit@1 |
| RotatE | Train | 1666s | 932MB | 0.3165 | 0.3506 | 0.2251 | 1666s | 932MB | 0.3165 | 0.3506 | 0.2251 |
| | Retrain | 1572s | 901MB | 0.2838 | 0.3191 | 0.1867 | 1596s | 889MB | 0.2751 | 0.3093 | 0.1783 |
| | GNNDelete | 247s | 969MB | 0.2746 | 0.3056 | **0.1841** | 365s | 946MB | 0.2713 | 0.3006 | 0.1813 |
| | CertifiedGU | 184s | 9834MB | **0.2813** | **0.3162** | 0.1882 | 186s | 9743MB | **0.2749** | 0.3078 | 0.1826 |
| | GIF | 5.82s | 6566MB | 0.2766 | 0.3088 | 0.1799 | 4.66s | 6473MB | 0.2799 | 0.3088 | 0.1799 |
| | WF-KGIF (Our) | 1.01s | 2101MB | 0.2779 | 0.3103 | 0.1795 | 1.00s | 2101MB | 0.2706 | 0.3075 | **0.1795** |
| | ZeroFisher (Our) | **0.94s** | **798MB** | 0.2783 | 0.3113 | 0.1811 | **0.95s** | **794MB** | 0.2748 | **0.3090** | 0.1811 |
| TransD | Train | 1702s | 893MB | 0.2928 | 0.3314 | 0.1971 | 1702s | 893MB | 0.2928 | 0.3314 | 0.1971 |
| | Retrain | 1603s | 888MB | 0.2788 | 0.3156 | 0.1824 | 1586s | 876MB | 0.2745 | 0.3076 | 0.1804 |
| | GNNDelete | 379s | 961MB | 0.2737 | 0.3073 | 0.1805 | 371s | 833MB | 0.2700 | 0.3012 | 0.1777 |
| | CertifiedGU | 610s | 31769MB | 0.2674 | 0.3053 | 0.1685 | 612s | 31760MB | 0.2615 | 0.2978 | 0.1633 |
| | GIF | 7.16s | 21770MB | 0.2729 | 0.3055 | 0.1742 | 9.0s | 21765MB | 0.2569 | 0.2934 | 0.1587 |
| | WF-KGIF (Our) | 1.17s | 4540MB | 0.2738 | 0.3068 | 0.1750 | 1.06s | 4540MB | 0.2693 | 0.2944 | 0.1599 |
| | ZeroFisher (Our) | 1.02s | 768MB | 0.2750 | 0.3068 | 0.1771 | 0.76s | 670MB | 0.2713 | 0.2976 | 0.1631 |
| TransH | Train | 1804s | 773MB | 0.2951 | 0.3315 | 0.2008 | 1804s | 773MB | 0.2951 | 0.3315 | 0.2008 |
| | Retrain | 1650s | 767MB | 0.2787 | 0.3138 | 0.1824 | 1648s | 742MB | 0.2740 | 0.3085 | 0.1648 |
| | GNNDelete | 370s | 811MB | 0.2742 | 0.3085 | 0.1790 | 374s | 789MB | 0.2683 | 0.3015 | 0.1742 |
| | CertifiedGU | 609s | 30791MB | 0.2696 | 0.2998 | 0.1699 | 580s | 30603MB | 0.2627 | 0.2962 | 0.1663 |
| | GIF | 6.74s | 20938MB | 0.2736 | 0.3080 | 0.1764 | 7.03s | 20936MB | 0.2679 | 0.3018 | 0.1683 |
| | WF-KGIF (Our) | 1.01s | 4168MB | 0.2737 | 0.3077 | 0.1771 | 1.02s | 4168MB | 0.2688 | 0.3031 | 0.1690 |
| | ZeroFisher (Our) | 0.96s | 537MB | 0.2744 | 0.3104 | 0.1800 | 0.98s | 538MB | 0.2697 | 0.3066 | 0.1740 |

**Table 2: Deleting different percentages of triplets on FB15K237.**

we introduce two variants of our method: WF-KGIF and KGIF. Instead of using Zeroth-order information to estimate gradients, WF-KGIF only leverages the Woodbury theorem and the Fisher matrices for unlearning. KGIF further directly computer the IHVPs. In Figure 2, we show the runtime and the peak value of the GPU memory usage of three variants on two datasets. The performance

of WF-KGIF is also included in Table 1. Note that we have recorded the memory usage in the three important steps of computing the influence function: loss computing, gradient computing, and IHVPs. We can see that approximation with the Woodbury theorem and the Fisher matrices can reduce the runtime significantly. Zoreth-order information can further reduce the runtime. We also observe
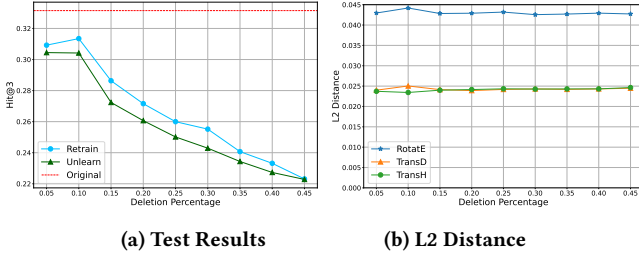
(a) Test Results

(b) L2 Distance

**Figure 6: Different deletion percentages**

that IHVPs account for the majority of memory usage. For IHVPs, gradient computation, and loss value computation, replacing our approximation with direct computation will increase memory usage. It verifies the necessity of each proposed component.

*5.2.3 Effects of Hyperparameters.* We investigate the impact of 4 hyperparameters on the unlearning results: iterations, the sampling term, the scaling down term, and the noise term. In Figure 3, we visualize Hit@10 for different configurations. More results are included in **Appendix B**.

Each iteration updates the parameter with estimation, bringing them closer to the desired state. A higher number of iterations typically leads to convergent results, as the algorithm has more opportunities to refine the parameter estimates. Recall that the damping term in Eq.12 plays a vital role in controlling the update magnitude, ensuring that the update process is neither too aggressive nor too slow, thereby balancing convergence speed and stability. As the damping term increases (see Figure 3a), the performance gradually becomes unstable, and the standard deviation increases. This is because if the damping term is too high, it will lead to an excessive increase in the estimation of IHVPs. As the number of iterations increases, the overall parameter estimation will accumulate, resulting in less robust results.

The scaling-down term determines how much the calculated updates are scaled down before being applied to the model parameters. A lower value of the scaling down term results in greater fluctuations in the parameter updates, as the adjustments are more sensitive to each iteration's computation. Conversely, a higher value for the scaling down term results in smaller, more controlled updates, making the training process more stable. The noise term $\epsilon$ is critical for determining the accuracy of the approximating gradients. When the noise term is set to a smaller value, the perturbations are smaller, leading to more precise gradient approximations. Consequently, the performance of the unlearning process improves, resulting in model parameters that are closer to the desired state. A larger noise term speeds up the convergence but at the cost of performance drop.

*5.2.4 Embeddings Distance.* To further investigate the unlearning quality, we visualize the embedding difference after unlearning. We randomly sample 1,000 entities and project their embeddings into 2-dimensional space with PAC. By comparing original embeddings (see Figure 4a) and embeddings after unlearning (see Figure 4b to 4e), we found that the removal of data can have an impact on the distribution of embeddings. We also show the $L_2$ distance between embedding after retraining and the ones after approximate unlearning. Compared with the best baseline (GIF), the embedding after

unlearning with our algorithm is closer to retraining. It explicitly proves the quality of our unlearning method.

*5.2.5 Deletion Percentage.* In Figure 6, we investigate the impact of deletion percentage. As the deletion percentage increases, The performance of the KG model drops too(See in Figure 6a). However, our unlearning method consistently remains close to the retraining. Also, $L_2$ The $L2$ distance between embedding after retraining and the ones with our unlearning remains stable as the deletion percentage increases. More experiments are dicussed in A.

## 6 Discussion

All datasets used in this paper are publicly available. Our research can be easily reproduced. Unfortunately, there is currently no strong evaluation metrics in Machine Unlearning that indicates whether the model weights have completely removed the influence of some data. For the retraining model, the deleted dataset actually represents a special validation set or test set, because retraining model have never seen the deleted dataset just like the validation set or test set. The worst performance on the deleted dataset does not mean that it is truly forgotten. In contrast, because of the generalization of the retraining model itself, it can also achieve good results on the deleted set, even if the deleted set has never participated in training process. MUSE [55] use Member Inference Attack (MIA) [34] to classify if a data is in training set. But [21] found that MIA barely outperform random guessing for most settings across varying LLM sizes and domains. Back to the question itself, how to determine whether the model has unlearned this data? [73] proposed that what we truly unlearn is how to use the knowledge but the knowledge itself. They showed how deleted data affects the distribution of the model from the loss landscape. So this might be an out-of-distribution problem. Ignoring the default test set, we select original, retrain and our method models to test on the remaining data set and the deleted set. The unlearned model deviates from the original model and moves towards the retrained model (See in tab 4). Introducing distribution alignment or optimal transport [69] into machine unlearning could yield broader conclusions across KGs, CV, NLP, and other data modalities.

## 7 Conclusion

This paper introduces a novel influence function–based framework to effectively eliminate the residual effects of deleted triplets, entities, and relations in knowledge graph (KG) models. Recognizing the high computational cost of traditional influence function computation, we propose a hierarchical approximation strategy to improve efficiency without sacrificing unlearning performance. Specifically, we first approximate the second-order inverse Hessian–vector products (IHVPs) using tractable first-order information, significantly reducing the need for costly Hessian computations. To further improve scalability, we approximate the first-order gradients using zeroth-order information through carefully designed perturbation-based estimations. This two-step approximation makes the influence estimation process both scalable and model-agnostic. We validate the effectiveness and efficiency of our method through extensive experiments on standard benchmark datasets, comparing against both full retraining and existing unlearning baselines. Results demonstrate that our approach achieves

| Model | Method | FB15K237 | | | | | YAGO3 | | | | |
|-------|--------|------|--------|-----|-------|-------|------|--------|-----|-------|-------|
| | | Time | Memory | MRR | Hit@3 | Hit@1 | Time | Memory | MRR | Hit@3 | Hit@1 |
| RotatE | Train | 1666s | 932MB | 0.3165 | 0.3506 | 0.2251 | 11049s | 4493MB | 0.3381 | 0.3846 | 0.2356 |
| | Retrain | 1342s | 784MB | 0.2913 | 0.3217 | 0.2100 | 9672s | 3278MB | 0.3249 | 0.3707 | 0.2261 |
| | GNNDelete | 337s | 856MB | 0.2867 | 0.3159 | 0.2043 | 1572s | 3677MB | 0.2799 | 0.3278 | 0.2003 |
| | CertifiedGU | 182s | 9504MB | 0.2944 | 0.3227 | 0.2070 | 201s | 39768MB | 0.3216 | 0.3667 | 0.2245 |
| | GIF | 3.5s | 6399MB | 0.2935 | 0.3211 | 0.2066 | 11.9s | 26417MB | 0.3202 | 0.3652 | 0.2265 |
| | WF-KGIF (Our) | 0.23s | 2040MB | 0.2954 | 0.3142 | 0.1997 | 0.50s | 9032MB | 0.3217 | 0.3679 | 0.2209 |
| | ZeroFisher (Our) | 0.19s | 774MB | 0.2978 | 0.3175 | 0.2079 | 0.48s | 4125MB | 0.3304 | 0.3738 | 0.2290 |
| TransD | Train | 1702s | 893MB | 0.2928 | 0.3314 | 0.1971 | 8543s | 3642MB | 0.3197 | 0.4542 | 0.2796 |
| | Retrain | 1478s | 796MB | 0.2696 | 0.3056 | 0.1828 | 7907s | 3471MB | 0.3442 | 0.4003 | 0.2381 |
| | GNNDelete | 455s | 850MB | 0.2600 | 0.2928 | 0.1733 | 1499s | 3517MB | 0.2894 | 0.3467 | 0.2015 |
| | CertifiedGU | 529s | 31459MB | 0.2833 | 0.3177 | 0.1910 | - | OOM | - | - | - |
| | GIF | 6.6s | 21221MB | 0.2643 | 0.2984 | 0.1805 | 18.6s | 84079MB | 0.3398 | 0.4071 | 0.2265 |
| | WF-KGIF (Our) | 0.29s | 4066MB | 0.2603 | 0.2908 | 0.1794 | 1.97s | 19718MB | 0.3325 | 0.3986 | 0.2156 |
| | ZeroFisher (Our) | 0.36s | 628MB | 0.2681 | 0.2924 | 0.1862 | 1.86s | 3235MB | 0.3482 | 0.4076 | 0.2372 |
| TransH | Train | 1804s | 773MB | 0.2951 | 0.3315 | 0.2008 | 8421s | 3118MB | 0.4212 | 0.4851 | 0.3108 |
| | Retrain | 1435s | 688MB | 0.2697 | 0.3043 | 0.1840 | 7851s | 2969MB | 0.3522 | 0.4104 | 0.2443 |
| | GNNDelete | 439s | 730MB | 0.2736 | 0.3057 | 0.1865 | 1513s | 3143MB | 0.2837 | 0.3334 | 0.1917 |
| | CertifiedGU | 489s | 29453MB | 0.2850 | 0.3176 | 0.1942 | - | OOM | - | - | - |
| | GIF | 6.3s | 20409MB | 0.2658 | 0.2992 | 0.1721 | 16.7s | 82465MB | 0.3548 | 0.4133 | 0.2499 |
| | WF-KGIF (Our) | 0.38s | 4038MB | 0.2641 | 0.2935 | 0.1736 | 1.62s | 19718MB | 0.3576 | 0.4146 | 0.2508 |
| | ZeroFisher (Our) | 0.32s | 487MB | 0.2708 | 0.2951 | 0.1879 | 1.37s | 2614MB | 0.3494 | 0.4023 | 0.2434 |

**Table 3: Relation Unlearning on FB15K237 and YAGO3.**

| Method | Type | Dataset | Del. Ratio | Unlearn MRR | Unlearn Hit@3 | Unlearn Hit@1 | Remain MRR | Remain Hit@3 | Remain Hit@1 |
|--------|------|---------|-----------|-------------|---------------|---------------|------------|--------------|--------------|
| Original | Random | FB15K237 | 10% | 0.5166 | 0.6067 | 0.3798 | 0.4708 | 0.5536 | 0.3339 |
| Retrain | Random | FB15K237 | 10% | 0.3247 | 0.3856 | 0.2008 | 0.5094 | 0.5920 | 0.3768 |
| ZeroFisher | Random | FB15K237 | 10% | 0.4875 | 0.5030 | 0.3101 | 0.4707 | 0.5536 | 0.3339 |
| Original | Random | YAGO3-10 | 10% | 0.5662 | 0.6337 | 0.4503 | 0.5588 | 0.6270 | 0.4420 |
| Retrain | Random | YAGO3-10 | 10% | 0.2800 | 0.3210 | 0.1792 | 0.5769 | 0.6425 | 0.4625 |
| ZeroFisher | Random | YAGO3-10 | 10% | 0.5056 | 0.5379 | 0.3879 | 0.5581 | 0.6259 | 0.4411 |

**Table 4: Discussion about the data distribution**

competitive unlearning quality while significantly reducing computational overhead, offering a practical solution for safe and efficient knowledge removal in deployed KG systems.

## A Appendix A

Table 2 is the result of deleting different percentages of triplets. As the deletion percentage increases, the performance of the KG model drops. However, our unlearning method consistently remains close to the retraining. Also, deleting more data will increase the running time of unlearning for all unlearning methods. Compared with the baselines, our unlearning method is especially efficient while deleting more data.

## B Appendix B

Figrue 5 shows hyper-parameter combination results on two models. The damping term controls update magnitude, balancing convergence speed and stability. As the damping term increases, performance becomes unstable due to inflated IHVP estimates. As the number of iterations increases, the overall parameter estimation will accumulate, resulting in less robust results. The scaling-down term controls how much updates are reduced before applying to

model parameters. A lower scaling-down value causes larger parameter fluctuations, making updates more sensitive to each iteration. A higher value for the scaling down term results in smaller, more controlled updates, making the training process more stable. The noise term is critical for determining the accuracy of the approximating gradients. When the noise term is set to a smaller value, the perturbations are smaller, leading to more precise gradient approximations. Consequently, the performance of the unlearning process improves, resulting in model parameters that are closer to the desired state. A larger noise term speeds up the convergence, but at the cost of a performance drop.

## C GenAI Usage Disclosure

This document was fully created and edited by humans. AI tools have not been used to create the content.

## References

[1] Asia J Biega and Michèle Finck. 2021. Reviving purpose limitation and data minimisation in data-driven systems. *arXiv preprint arXiv:2101.06203* (2021).
[2] Lucas Bourtoule, Varun Chandrasekaran, Christopher A Choquette-Choo, Hengrui Jia, Adelin Travers, Baiwu Zhang, David Lie, and Nicolas Papernot. 2021.

Machine unlearning. In *2021 IEEE Symposium on Security and Privacy (SP)*. IEEE, 141–159.

[3] Jonathan Brophy and Daniel Lowd. 2021. Machine Unlearning for Random Forests. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event (Proceedings of Machine Learning Research, Vol. 139)*, Marina Meila and Tong Zhang (Eds.). PMLR, 1092–1104. http://proceedings.mlr.press/v139/brophy21a.html

[4] Yinzhi Cao and Junfeng Yang. 2015. Towards making systems forget with machine unlearning. In *2015 IEEE symposium on security and privacy*. IEEE, 463–480.

[5] Sungmin Cha, Sungjun Cho, Dasol Hwang, Honglak Lee, Taesup Moon, and Moontae Lee. 2024. Learning to Unlearn: Instance-Wise Unlearning for Pre-trained Classifiers. In *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024, Vancouver, Canada*, Michael J. Wooldridge, Jennifer G. Dy, and Sriraam Natarajan (Eds.). AAAI Press, 11186–11194. doi:10.1609/AAAI.V38I10.28996

[6] Chong Chen, Fei Sun, Min Zhang, and Bolin Ding. 2022. Recommendation Unlearning. In *WWW '22: The ACM Web Conference 2022, Virtual Event, Lyon, France, April 25 - 29, 2022*, Frédérique Laforest, Raphaël Troncy, Elena Simperl, Deepak Agarwal, Aristides Gionis, Ivan Herman, and Lionel Médini (Eds.). ACM, 2768–2777. doi:10.1145/3485447.3511997

[7] Jiaao Chen and Diyi Yang. 2023. Unlearn What You Want to Forget: Efficient Unlearning for LLMs. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, 12041–12052. doi:10.18653/V1/2023.EMNLP-MAIN.738

[8] Min Chen, Weizhuo Gao, Gaoyang Liu, Kai Peng, and Chen Wang. 2023. Boundary Unlearning: Rapid Forgetting of Deep Networks via Shifting the Decision Boundary. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023*. IEEE, 7766–7775. doi:10.1109/CVPR52729.2023.00750

[9] Min Chen, Zhikun Zhang, Tianhao Wang, Michael Backes, Mathias Humbert, and Yang Zhang. 2021. When Machine Unlearning Jeopardizes Privacy. In *CCS '21: 2021 ACM SIGSAC Conference on Computer and Communications Security, Virtual Event, Republic of Korea, November 15 - 19, 2021*, Yongdae Kim, Jong Kim, Giovanni Vigna, and Elaine Shi (Eds.). ACM, 896–911. doi:10.1145/3460120.3484756

[10] Min Chen, Zhikun Zhang, Tianhao Wang, Michael Backes, Mathias Humbert, and Yang Zhang. 2022. Graph Unlearning. In *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security, CCS 2022, Los Angeles, CA, USA, November 7-11, 2022*, Heng Yin, Angelos Stavrou, Cas Cremers, and Elaine Shi (Eds.). ACM, 499–513. doi:10.1145/3548606.3559352

[11] Min Chen, Zhikun Zhang, Tianhao Wang, Michael Backes, Mathias Humbert, and Yang Zhang. 2022. Graph unlearning. In *Proceedings of the 2022 ACM SIGSAC conference on computer and communications security*. 499–513.

[12] Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh. 2017. Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models. In *Proceedings of the 10th ACM workshop on artificial intelligence and security*. 15–26.

[13] Jiali Cheng, George Dasoulas, Huan He, Chirag Agarwal, and Marinka Zitnik. 2023. GNNDelete: A General Strategy for Unlearning in Graph Neural Networks. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net. https://openreview.net/pdf?id=X9yCkmT5Qrl

[14] Jiali Cheng, George Dasoulas, Huan He, Chirag Agarwal, and Marinka Zitnik. 2023. Gnndelete: A general strategy for unlearning in graph neural networks. *arXiv preprint arXiv:2302.13406* (2023).

[15] Eli Chien, Chao Pan, and Olgica Milenkovic. 2022. Certified Graph Unlearning. *CoRR* abs/2206.09140 (2022). doi:10.48550/ARXIV.2206.09140 arXiv:2206.09140

[16] Eli Chien, Chao Pan, and Olgica Milenkovic. 2022. Efficient model updates for approximate unlearning of graph-structured data. In *The Eleventh International Conference on Learning Representations*.

[17] Vikram S. Chundawat, Ayush K. Tarun, Murari Mandal, and Mohan S. Kankanhalli. 2023. Can Bad Teaching Induce Forgetting? Unlearning in Deep Networks Using an Incompetent Teacher. In *Thirty-Seventh AAAI Conference on Artificial Intelligence, AAAI 2023, Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence, IAAI 2023, Thirteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2023, Washington, DC, USA, February 7-14, 2023*, Brian Williams, Yiling Chen, and Jennifer Neville (Eds.). AAAI Press, 7210–7217. doi:10.1609/AAAI.V37I6.25879

[18] R Dennis Cook and Sanford Weisberg. 1980. Characterizations of an empirical influence function for influential cases in regression. *Technometrics* 22, 4 (1980), 495–508.

[19] Yuanfei Dai, Shiping Wang, Neal N Xiong, and Wenzhong Guo. 2020. A survey on knowledge graph embedding: Approaches, applications and benchmarks. *Electronics* 9, 5 (2020), 750.

[20] Yushun Dong, Binchi Zhang, Zhenyu Lei, Na Zou, and Jundong Li. 2024. IDEA: A Flexible Framework of Certified Unlearning for Graph Neural Networks. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD 2024, Barcelona, Spain, August 25-29, 2024*, Ricardo Baeza-Yates and Francesco Bonchi (Eds.). ACM, 621–630. doi:10.1145/3637528.3671744

[21] Michael Duan, Anshuman Suri, Niloofar Mireshghallah, Sewon Min, Weijia Shi, Luke Zettlemoyer, Yulia Tsvetkov, Yejin Choi, David Evans, and Hannaneh Hajishirzi. 2024. Do membership inference attacks work on large language models? *arXiv preprint arXiv:2402.07841* (2024).

[22] Xiaohua Feng, Chaochao Chen, Yuyuan Li, and Zibin Lin. 2024. Fine-grained Pluggable Gradient Ascent for Knowledge Unlearning in Language Models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (Eds.). Association for Computational Linguistics, 10141–10155. https://aclanthology.org/2024.emnlp-main.566

[23] Eduard Fosch-Villaronga, Peter Kieseberg, and Tiffany Li. 2018. Humans forget, machines remember: Artificial intelligence and the Right to Be Forgotten. *Comput. Law Secur. Rev.* 34, 2 (2018), 304–313. doi:10.1016/j.clsr.2017.08.007

[24] Jack Foster, Stefan Schoepf, and Alexandra Brintrup. 2024. Fast Machine Unlearning without Retraining through Selective Synaptic Dampening. In *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024, Vancouver, Canada*, Michael J. Wooldridge, Jennifer G. Dy, and Sriraam Natarajan (Eds.). AAAI Press, 12043–12051. doi:10.1609/AAAI.V38I11.29092

[25] Antonio Ginart, Melody Guan, Gregory Valiant, and James Y Zou. 2019. Making ai forget you: Data deletion in machine learning. *Advances in neural information processing systems* 32 (2019).

[26] Antonio Ginart, Melody Y. Guan, Gregory Valiant, and James Zou. 2019. Making AI Forget You: Data Deletion in Machine Learning. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett (Eds.). 3513–3526. https://proceedings.neurips.cc/paper/2019/hash/cb79f8fa58b91d3af6c9c991f63962d3-Abstract.html

[27] Shashwat Goel, Ameya Prabhu, Amartya Sanyal, Ser-Nam Lim, Philip Torr, and Ponnurangam Kumaraguru. 2022. Towards adversarial evaluations for inexact machine unlearning. *arXiv preprint arXiv:2201.06640* (2022).

[28] Laura Graves, Vineel Nagisetty, and Vijay Ganesh. 2021. Amnesiac Machine Learning. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021*. AAAI Press, 11516–11524. doi:10.1609/AAAI.V35I13.17371

[29] Elad Hazan et al. 2016. Introduction to online convex optimization. *Foundations and Trends® in Optimization* 2, 3-4 (2016), 157–325.

[30] Masaru Isonuma and Ivan Titov. 2024. Unlearning Traces the Influential Training Data of Language Models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, 6312–6325. doi:10.18653/V1/2024.ACL-LONG.343

[31] Joel Jang, Dongkeun Yoon, Sohee Yang, Sungmin Cha, Moontae Lee, Lajanugen Logeswaran, and Minjoon Seo. 2023. Knowledge Unlearning for Mitigating Privacy Risks in Language Models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, Anna Rogers, Jordan L. Boyd-Graber, and Naoaki Okazaki (Eds.). Association for Computational Linguistics, 14389–14408. doi:10.18653/V1/2023.ACL-LONG.805

[32] Guoliang Ji, Shizhu He, Liheng Xu, Kang Liu, and Jun Zhao. 2015. 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)*. 687–696.

[33] Shaoxiong Ji, Shirui Pan, Erik Cambria, Pekka Marttinen, and Philip S. Yu. 2022. A Survey on Knowledge Graphs: Representation, Acquisition, and Applications. *IEEE Trans. Neural Networks Learn. Syst.* 33, 2 (2022), 494–514. doi:10.1109/TNNLS.2021.3070843

[34] Gyuwan Kim, Yang Li, Evangelia Spiliopoulou, Jie Ma, Miguel Ballesteros, and William Yang Wang. 2024. Detecting Training Data of Large Language Models via Expectation Maximization. *arXiv preprint arXiv:2410.07582* (2024).

[35] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).

[36] Thomas N. Kipf and Max Welling. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net. https://openreview.net/forum?id=SJU4ayYgl

[37] Pang Wei Koh and Percy Liang. 2017. Understanding Black-box Predictions via Influence Functions. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017 (Proceedings of Machine Learning Research, Vol. 70)*, Doina Precup and Yee Whye Teh (Eds.).

PMLR, 1885–1894. http://proceedings.mlr.press/v70/koh17a.html

[38] Pang Wei Koh and Percy Liang. 2017. Understanding black-box predictions via influence functions. In *International conference on machine learning*. PMLR, 1885–1894.

[39] Meghdad Kurmanji, Peter Triantafillou, Jamie Hayes, and Eleni Triantafillou. 2023. Towards Unbounded Machine Unlearning. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (Eds.). http://papers.nips.cc/paper_files/paper/2023/hash/062d711fb777322e2152435459e6e9d9-Abstract-Conference.html

[40] Yuantong Li, Chi-Hua Wang, and Guang Cheng. 2021. Online Forgetting Process for Linear Regression Models. In *The 24th International Conference on Artificial Intelligence and Statistics, AISTATS 2021, April 13-15, 2021, Virtual Event (Proceedings of Machine Learning Research, Vol. 130)*, Arindam Banerjee and Kenji Fukumizu (Eds.). PMLR, 217–225. http://proceedings.mlr.press/v130/li21a.html

[41] Sijia Liu, Yuanshun Yao, Jinghan Jia, Stephen Casper, Nathalie Baracaldo, Peter Hase, Xiaojun Xu, Yuguang Yao, Hang Li, Kush R. Varshney, Mohit Bansal, Sanmi Koyejo, and Yang Liu. 2024. Rethinking Machine Unlearning for Large Language Models. *CoRR* abs/2402.08787 (2024). doi:10.48550/ARXIV.2402.08787 arXiv:2402.08787

[42] Alexander Ly, Maarten Marsman, Josine Verhagen, Raoul PPP Grasman, and Eric-Jan Wagenmakers. 2017. A tutorial on Fisher information. *Journal of Mathematical Psychology* 80 (2017), 40–55.

[43] Ananth Mahadevan and Michael Mathioudakis. 2021. Certifiable Machine Unlearning for Linear Models. *CoRR* abs/2106.15093 (2021). arXiv:2106.15093 https://arxiv.org/abs/2106.15093

[44] Peihua Mai, Hao Jiang, Ran Yan, Youjia Yang, Zhe Huang, and Yan Pang. 2024. Knowledge Graph Unlearning to Defend Language Model Against Jailbreak Attack. In *The Second Tiny Papers Track at ICLR 2024, Tiny Papers @ ICLR 2024, Vienna, Austria, May 11, 2024*. OpenReview.net. https://openreview.net/forum?id=gqTUIesy4H

[45] Sadhika Malladi, Tianyu Gao, Eshaan Nichani, Alex Damian, Jason D. Lee, Danqi Chen, and Sanjeev Arora. 2023. Fine-Tuning Language Models with Just Forward Passes. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (Eds.). http://papers.nips.cc/paper_files/paper/2023/hash/a627810151be4d13f907ac898ff7e948-Abstract-Conference.html

[46] Neil G. Marchant, Benjamin I. P. Rubinstein, and Scott Alfeld. 2022. Hard to Forget: Poisoning Attacks on Certified Machine Unlearning. In *Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022, Thirty-Fourth Conference on Innovative Applications of Artificial Intelligence, IAAI 2022, The Twelfth Symposium on Educational Advances in Artificial Intelligence, EAAI 2022 Virtual Event, February 22 - March 1, 2022*. AAAI Press, 7691–7700. doi:10.1609/AAAI.V36I7.20736

[47] Saemi Moon, Seunghyuk Cho, and Dongwoo Kim. 2024. Feature Unlearning for Pre-trained GANs and VAEs. In *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024, Vancouver, Canada*, Michael J. Wooldridge, Jennifer G. Dy, and Sriraam Natarajan (Eds.). AAAI Press, 21420–21428. doi:10.1609/AAAI.V38I19.30138

[48] Seth Neel, Aaron Roth, and Saeed Sharifi-Malvajerdi. 2021. Descent-to-Delete: Gradient-Based Methods for Machine Unlearning. In *Algorithmic Learning Theory, 16-19 March 2021, Virtual Conference, Worldwide (Proceedings of Machine Learning Research, Vol. 132)*, Vitaly Feldman, Katrina Ligett, and Sivan Sabato (Eds.). PMLR, 931–962. http://proceedings.mlr.press/v132/neel21a.html

[49] Thanh Tam Nguyen, Thanh Trung Huynh, Phi Le Nguyen, Alan Wee-Chung Liew, Hongzhi Yin, and Quoc Viet Hung Nguyen. 2022. A Survey of Machine Unlearning. *CoRR* abs/2209.02299 (2022). doi:10.48550/ARXIV.2209.02299 arXiv:2209.02299

[50] Thanh Tam Nguyen, Thanh Trung Huynh, Phi Le Nguyen, Alan Wee-Chung Liew, Hongzhi Yin, and Quoc Viet Hung Nguyen. 2022. A survey of machine unlearning. *arXiv preprint arXiv:2209.02299* (2022).

[51] CA OAG. 2021. CCPA regulations: Final regulation text. *Office of the Attorney General, California Department of Justice* (2021), 1.

[52] Protection Regulation. 2016. Regulation (EU) 2016/679 of the European Parliament and of the Council. *Regulation (eu)* 679 (2016), 2016.

[53] Anwar Said, Tyler Derr, Mudassir Shabbir, Waseem Abbas, and Xenofon Koutsoukos. 2023. A survey of graph unlearning. *arXiv preprint arXiv:2310.02164* (2023).

[54] Ayush Sekhari, Jayadev Acharya, Gautam Kamath, and Ananda Theertha Suresh. 2021. Remember what you want to forget: Algorithms for machine unlearning.

[55] *Advances in Neural Information Processing Systems* 34 (2021), 18075–18086.

[55] Weijia Shi, Jaechan Lee, Yangsibo Huang, Sadhika Malladi, Jieyu Zhao, Ari Holtzman, Daogao Liu, Luke Zettlemoyer, Noah A Smith, and Chiyuan Zhang. 2024. Muse: Machine unlearning six-way evaluation for language models. *arXiv preprint arXiv:2407.06460* (2024).

[56] Sidak Pal Singh and Dan Alistarh. 2020. Woodfisher: Efficient second-order approximation for neural network compression. *Advances in Neural Information Processing Systems* 33 (2020), 18098–18109.

[57] Zhiqing Sun, Zhi-Hong Deng, Jian-Yun Nie, and Jian Tang. 2019. Rotate: Knowledge graph embedding by relational rotation in complex space. *arXiv preprint arXiv:1902.10197* (2019).

[58] Anna Ubaydullayeva. 2023. Artificial intelligence and intellectual property: navigating the complexities of cyber law. *International Journal of Law and Policy* 1, 4 (2023).

[59] Enayat Ullah, Tung Mai, Anup Rao, Ryan A. Rossi, and Raman Arora. 2021. Machine Unlearning via Algorithmic Stability. In *Conference on Learning Theory, COLT 2021, 15-19 August 2021, Boulder, Colorado, USA (Proceedings of Machine Learning Research, Vol. 134)*, Mikhail Belkin and Samory Kpotufe (Eds.). PMLR, 4126–4142. http://proceedings.mlr.press/v134/ullah21a.html

[60] Cheng-Long Wang, Mengdi Huai, and Di Wang. 2023. Inductive graph unlearning. In *32nd USENIX Security Symposium (USENIX Security 23)*. 3205–3222.

[61] Lingzhi Wang, Tong Chen, Wei Yuan, Xingshan Zeng, Kam-Fai Wong, and Hongzhi Yin. 2023. KGA: A General Machine Unlearning Framework Based on Knowledge Gap Alignment. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, Anna Rogers, Jordan L. Boyd-Graber, and Naoaki Okazaki (Eds.). Association for Computational Linguistics, 13264–13276. doi:10.18653/V1/2023.ACL-LONG.740

[62] Zhen Wang, Jianwen Zhang, Jianlin Feng, and Zheng Chen. 2014. Knowledge graph embedding by translating on hyperplanes. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 28.

[63] Max A Woodbury. 1950. *Inverting modified matrices*. Department of Statistics, Princeton University.

[64] Jiancan Wu, Yi Yang, Yuchun Qian, Yongduo Sui, Xiang Wang, and Xiangnan He. 2023. Gif: A general graph unlearning strategy via influence function. In *Proceedings of the ACM Web Conference 2023*. 651–661.

[65] Kun Wu, Jie Shen, Yue Ning, Ting Wang, and Wendy Hui Wang. 2023. Certified edge unlearning for graph neural networks. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 2606–2617.

[66] Yinjun Wu, Edgar Dobriban, and Susan B. Davidson. 2020. DeltaGrad: Rapid retraining of machine learning models. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event (Proceedings of Machine Learning Research, Vol. 119)*. PMLR, 10355–10366. http://proceedings.mlr.press/v119/wu20b.html

[67] Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and S Yu Philip. 2020. A comprehensive survey on graph neural networks. *IEEE transactions on neural networks and learning systems* 32, 1 (2020), 4–24.

[68] Yang Xiao, Gen Li, Jie Ji, Ruimeng Ye, Xiaolong Ma, and Bo Hui. 2025. The Right to be Forgotten in Pruning: Unveil Machine Unlearning on Sparse Models. *arXiv preprint arXiv:2507.18725* (2025).

[69] Yang Xiao, Wang Lu, Jie Ji, Ruimeng Ye, Gen Li, Xiaolong Ma, and Bo Hui. 2025. Optimal Transport for Brain-Image Alignment: Unveiling Redundancy and Synergy in Neural Information Processing. *arXiv preprint arXiv:2503.10663* (2025).

[70] Yang Xiao, Ruimeng Ye, and Bo Hui. 2025. Knowledge Graph Unlearning with Schema. In *Proceedings of the 31st International Conference on Computational Linguistics*. 3541–3546.

[71] Yuanshun Yao, Xiaojun Xu, and Yang Liu. 2023. Large Language Model Unlearning. *CoRR* abs/2310.10683 (2023). doi:10.48550/ARXIV.2310.10683 arXiv:2310.10683

[72] Yihua Zhang, Pingzhi Li, Junyuan Hong, Jiaxiang Li, Yimeng Zhang, Wenqing Zheng, Pin-Yu Chen, Jason D. Lee, Wotao Yin, Mingyi Hong, Zhangyang Wang, Sijia Liu, and Tianlong Chen. 2024. Revisiting Zeroth-Order Optimization for Memory-Efficient LLM Fine-Tuning: A Benchmark. *CoRR* abs/2402.11592 (2024). doi:10.48550/ARXIV.2402.11592 arXiv:2402.11592

[73] Junhao Zheng, Xidi Cai, Shengjie Qiu, and Qianli Ma. 2025. Spurious Forgetting in Continual Learning of Language Models. *arXiv preprint arXiv:2501.13453* (2025).

[74] Xiangrong Zhu, Guangyao Li, and Wei Hu. 2023. Heterogeneous Federated Knowledge Graph Embedding Learning and Unlearning. In *Proceedings of the ACM Web Conference 2023, WWW 2023, Austin, TX, USA, 30 April 2023 - 4 May 2023*, Ying Ding, Jie Tang, Juan F. Sequeda, Lora Aroyo, Carlos Castillo, and Geert-Jan Houben (Eds.). ACM, 2444–2454. doi:10.1145/3543507.3583305