# TUNE: A Task For Turkish Machine Unlearning For Data Privacy

**Doruk Benli    Ada Canoğlu    Nehir İklim Gönençer    Dilara Keküllüoğlu**

Sabanci University

Istanbul, Türkiye

{dorukbenli,adacanoglu,nehir.gonencer,dilara.kekulluoglu}@sabanciuniv.edu

## Abstract

Most large language models (LLMs) are trained on massive datasets that include private information, which may be disclosed to third-party users in output generation. Developers put defences to prevent the generation of harmful and private information, but jailbreaking methods can be used to bypass them. Machine unlearning aims to remove information that may be private or harmful from the model's generation without retraining the model from scratch. While machine unlearning has gained some popularity to counter the removal of private information, especially in English, little to no attention has been given to Turkish unlearning paradigms or existing benchmarks. In this study, we introduce TUNE (Turkish Unlearning Evaluation), the first benchmark dataset for Turkish unlearning task for personal information. TUNE consists of 9842 input-target text pairs about 50 fictitious personalities with two training task types: (1) Q&A and (2) Information Request. We fine-tuned the mT5 base model to evaluate various unlearning methods, including our proposed approach. We find that while current methods can help unlearn unwanted private information in Turkish, they also unlearn other information we want to retain in the model.

## 1 Introduction

In recent years, the increasing use of large language models (LLMs) has raised important questions about data privacy, model behavior, and personal data privacy leakage through generated outputs (Yao et al., 2024a). These LLMs are trained on datasets that include private data such as personally identifiable information (PII) (Elazar et al., 2024). Training data with private information can be extracted from LLMs via adversary attacks (Carlini et al., 2021) even when a specific data point only appears once in the dataset. As the use of LLMs grows and models get bigger, the ability to

prevent such data from being generated becomes crucial which can be against personal data protection laws (Commission, 2018). Since models are designed to learn information during training, once trained, reversing the effect of even a small subset of data is not easy (Tahiliani et al., 2021) as almost all model weights change. Hence, it is almost impossible to request the private data to be removed from a specific LLM's training data (Lareo, 2023).

Machine unlearning is a method to "forget" the unwanted data while retaining others (Nguyen et al., 2025) without retraining the whole model. Instead, the model is further trained to update the weights away from the unwanted data points using methods such as gradient ascent. Benchmarks such as TOFU (Maini et al., 2024) and LUME (Ramakrishna et al., 2025) have been developed to experiment with various unlearning methods. However, all of these methods have been tested on English language datasets, and there has been very little attention paid to other languages, such as Turkish. To fill this gap, we propose TUNE[1] (**T**urkish **UN**learning **E**valuation), a dataset for training and fine-tuning purposes to evaluate unlearning methods in Turkish. To the best of our knowledge, this is the first work specifically designed for a Turkish unlearning task. TUNE consists of 9842 input-target text pairs across two task types for 50 fictitious synthetically generated people. We investigate machine unlearning in the context of multilingual large language models, specifically using the mT5 architecture (Xue et al., 2021). We demonstrate that TUNE can be utilized as a baseline benchmark dataset for unlearning tasks in the Turkish language by applying various unlearning methods, including our proposed method with a custom loss. We also show that current unlearning methods can be blunt; they unlearn data they need to retain while forgetting the unwanted data.

---

[1]The dataset is available here.

## 2 Related Work

Research on unlearning methods and benchmark datasets has accelerated in recent years, driven by the rapid adoption of large language models (LLMs). Most unlearning approaches in the context of LLMs have two distinct sets of personal information where one set is labeled as the "forget" set and the other as "retain". Machine unlearning methods try to maximize the forget rate while maintaining the performance of the retain dataset as much as possible. In this section, we first review benchmark datasets for machine unlearning, followed by a review of the works with machine unlearning approaches.

### 2.1 Related Datasets

TOFU (Maini et al., 2024) is one of the most prominent prior works that focused on creating a dataset to fine-tune a learned model through unlearning. The main objective of TOFU is to create a standardized, realistic way to evaluate how well LLMs can forget specific information after training. To achieve this goal, the authors built a synthetic dataset featuring 200 entirely fictional authors. Each of these authors was paired with 20 question-answer facts. To measure forgetting, they compared the tendency of the model to generate correct versus incorrect answers about forgotten authors (Maini et al., 2024). However, because TOFU is a benchmark that is designed in English, it is not possible to use it for multilingual purposes.

While TOFU focuses primarily on privacy, another popular benchmark, "Who is Harry Potter" (Eldan and Russinovich, 2023), focuses on removing information that is related to a specific domain. In their work, authors choose the Harry Potter universe for such domain knowledge. Similar to TOFU, this benchmark is also restricted to English, leaving a gap for Turkish Unlearning.

In addition to these, a more sophisticated benchmark called MUSE (Shi et al., 2025) points out that a single accuracy metric is not enough to measure unlearning performance. MUSE suggests six dimensions for evaluating unlearning: verbatim memorization, knowledge memorization, privacy leakage, utility preservation, scalability, and sustainability. Compared to earlier benchmarks such as TOFU and the Harry Potter evaluation, MUSE provides a far more thorough view of unlearning by assessing both data-owner expectations and model-deployer expectations. Yet, despite this broader perspective, MUSE remains limited to English corpora. On the contrary, RWKU (Jin et al., 2024) approaches the problem by selecting famous real-world figures as unlearning targets. RWKU argues that benchmarks like TOFU do not embed the deep prior knowledge that models contain. However, this causes boundary issues due to unknown pre-training data.

### 2.2 Related Unlearning Approaches

**Large Language Model Unlearning (LLMU):** LLM Unlearning (Yao et al., 2024b) frames unlearning as a targeted optimization problem. Instead of retraining the model, it modifies its parameters using gradient ascent on the "forget" data. However, it differs from gradient ascent as LLM Unlearning introduces three new loss components: (1) Forget Loss ($L_{\text{forget}}$), which increases the model's error on harmful responses to reduce their likelihood, (2) Random Mismatch Loss ($L_{\text{random}}$), which teaches the model to associate the harmful prompts with unrelated, non-harmful outputs, enforcing randomness in its responses to such prompts, and (3) Normal Utility Loss ($L_{\text{normalized}}$), which ensures the model still performs well on unrelated prompts by matching its predictions to those of the original model using forward KL divergence.

These losses are combined in a single update rule that shifts the model away from forgetting everything. It forces the model to forget the data that is in the forget data and not to lose its performance in the retain data. Unlike our proposed loss, LLMU does not employ an additional cross-entropy term. Our method uses this term to maintain task alignment, ensuring that the model continues to generate meaningful outputs.

**Negative Preference Optimization (NPO):** NPO is a variation of DPO (Direct Preference Optimization); a main feature of NPO is that it treats the forget data points as negative responses, and does not use positive responses (Fan et al., 2025). NPO behaves fairly similarly to Gradient Ascent's loss when the temperature is very high. However, unlike Gradient Ascent, it remains stable and has a lower bound at any finite temperature (Zhang et al., 2024). NPO purely focuses on decreasing the likelihood of producing undesired outputs, thus working only with the forget set. Compared to our proposed loss, NPO lacks a balancing term such as the retain set loss.
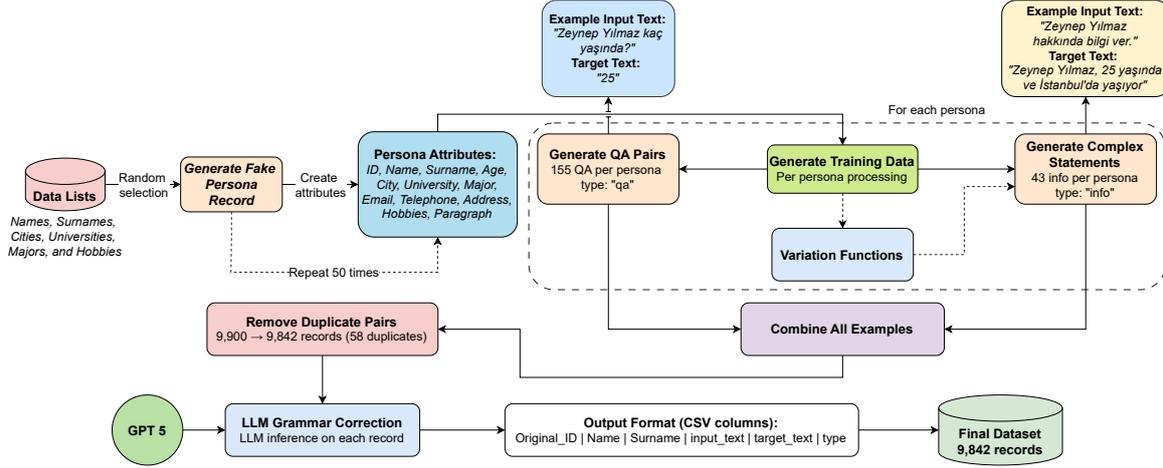
Figure 1: Data generation pipeline. This is the overall pipeline for generating 9842 different data points for 50 fictitious people. The first two steps generates 50 fictitious people; the rest of the pipeline generates the data for **TUNE** dataset. Retain and forget sets are subsets of this dataset.

**FLAT: Unlearning with Only Forget Data:**
FLAT (Wang et al., 2024) offers an approach that relies only on forget data, as retain data might not be available in a real-life scenario. FLAT aims to guide the model toward forgetting harmful or sensitive responses by comparing them with template responses—harmless or neutral answers using f-divergence optimization. For each forget prompt, they generated a corresponding response and then optimized the model so that the model's new answer can diverge from the original response. FLAT also introduces a framework for loss adjustment specified for LLM unlearning. Our loss follows this structure, but we employ retain terms along with forget terms and a custom term for task alignment.

**General Optimization and Architectural Strategies:** Several approaches look beyond simple gradient manipulation to ensure robust forgetting. For example, SCRUB (Scalable Remembering and Unlearning unBound) utilizes teacher-student model where the student model tries to disobey and deviate from the teacher model on forget set, while mimicking it on the retain set by alternating optimization schedule (Kurmanji et al., 2023). Knowledge Gap Alignment (KGA), Inspired by Knowledge-adaptation priors (K-priors) (Khan and Swaroop, 2021), a framework that adopts weight and function-space priors to guide the model toward a new state, enabling reduction of reliance on the forgotten data while preserving performance on the retained data. KGA (Wang et al., 2023), extends

this idea aiming data removal through knowledge-gap alignment. Here, the knowledge gap is defined as distance between prediction dispention of two architecturally identical models, trained with different data. When this gap is aligned, models give similar outputs.

## 3  TUNE: Turkish UNlearning Evaluation

In this work, we created the first benchmark for Turkish LLM unlearning. We provide a novel dataset with details of 50 fictitious people in Turkish with two task types. We then fine-tuned an mT5 (Xue et al., 2021) base model from scratch on this dataset. We implemented various baseline unlearning algorithms on the model and a specialized algorithm to unlearn people whose details we did not want the model to generate. We used the Rouge-L score (Lin, 2004) to measure the degree of forgetting.

### 3.1  TUNE Dataset Creation

The dataset creation pipeline is shown in Figure 1. We synthetically generate 50 fictitious people and construct TUNE using these personal details. To create these fictitious profiles, we first curated a set of possible values for names and surnames using the top 100 most given names for men and women, as well as the top 20 most given surnames in Turkish. We also randomly generated ages between 18 and 50 for these fictitious people. We then selected random cities, universities, majors, and hobbies for each persona to use in our training set. Additional attributes such as email addresses, phone numbers,

| Input Text | Target Text | Type |
|---|---|---|
| Yağmur Erdoğan'ın memleketi neresidir? | Yağmur Erdoğan'ın memleketi Samsun'dur. | qa |
| Yağmur nerede oturuyor? | Yağmur Samsun'da oturuyor. | qa |
| Yağmur lisans eğitimini hangi üniversitede alıyor? | Ege Üniversitesi. | qa |
| Yağmur Erdoğan hakkında ne biliyorsun? | Yağmur Erdoğan 18 yaşında ve Samsun'da yaşıyor. | info |

Table 1: Example input-target text pairs from the dataset. First three entries are examples of question and answer type data points, indicated by **qa** and the last entry is an example of information request type, indicated with **info**.

and addresses were also assigned to increase the realism of the dataset but not included in the q&a pairs since such information is typically retrieved from longer, structured texts. For the generation of addresses, we used the Faker library's Turkish version (Faraglia and Other Contributors). We then wrote a short template paragraph for each person, which combines all attributes into a structured text. Two types of tasks were included in our dataset:

- **Question and Answer:** This is a task for asking a question about one of the attributes and expecting the correct answer formulated in different ways, but the same in context.

- **Information request:** This task gives general information about a selected number of attributes upon a general information request by the user.

Table 1 presents the structure of the training data used to train the model, together with the type of the pair. In Table 1, there are four input-target text pairs in which three of them are question and answer type, marked with **qa**, and one information type, **info**, which gives more than one piece of information about a specific person. To create a question and answer pair, we defined a set of helping phrases to ask about some attributes excluding Email, Phone number, and Address in separately formed sentences. This created a dataset with many different ways of asking for the same attribute. This can be seen at the "Generate Training Data" stage on Figure 1. There are up to 35 distinct phrases to create questions for a single attribute. In this way we created a rich question dataset for each attribute. On top of this, we also added questions containing multiple attributes. For any combination of attributes, there are around 155 unique sentences we can use to form questions. For example, we can ask about the university of the persona using

different verbs in Turkish such as "okuyor" and "eğitim görüyor". The final question and answer dataset has 7698 entries.

To create information request enquiries, we defined a set of 23 distinct instruction prompts that take the name and surname of each fictitious person and ask some type of information. These prompts contain variations such as "Ahmet Yılmaz kimdir?" or "Ahmet Yılmaz hakkında detaylı bilgi verir misin?". Then we combine two or more attributes such as age, city, university, department, or hobbies into a single coherent description of persona to craft multiple ways of giving information. We combine instruction prompts with generated answers to reach the final information request dataset with 2144 entries. Unlike question and answer pairs, information request pairs have no attribute name in their input texts.

By combining this with question and answer dataset, we obtain 198 enquiries about each fictitious person across various attributes stated above. Following this combination stage, we remove the duplicate input-target text pairs as illustrated by the Figure 1. In the end, we were left with 9842 unique entries about 50 fictitious people.

As the last stage of TUNE data generation pipeline, as shown in Figure 1, we used OpenAI's GPT-5 nano, mini, and GPT-5 base models to perform grammatical corrections on answer and question pairs to ensure grammatical correctness and diversity across the dataset. GPT models were instructed to not alter the sentences or introduce new information. We do this to prevent the model from injecting its own information coincidentally if the names and surnames of our fictitious people match those of real people.

### 3.2 Retain and forget set creation

We divided the 50 personas into two sets: people to forget and people to retain. We then created forget

sets, $D_f$, and retain sets, $D_r$, as subsets of TUNE, selected from the entries of these personas. The aim of the unlearning is to prevent details about the personas in the forget dataset to be uttered by the language model while answering questions correctly about the personas in the retain dataset. After training an mT5 base model (Xue et al., 2021) with our entire constructed dataset, we retrained it with various algorithms to forget the $D_f$ while retaining the $D_r$. We experimented with different proportions of $D_f$ and $D_r$ to showcase the results when we have smaller forget set size which is closer to real-life conditions. We randomly chose 10 fictitious people from TUNE to create the forget set and we used remaining 40 people to create retain dataset (Maini et al., 2024).

### 3.3 Unlearning methods

To forget the entries in $D_f$ and retain those in $D_r$, we used three of the baseline unlearning algorithms, as well as our proposed approach on our trained mT5 (Xue et al., 2021) model. The baseline algorithms are: (1) Gradient Ascent (Maini et al., 2024); (2) Gradient Difference (Liu et al., 2022); and (3) KL Minimization (Maini et al., 2024).

**Gradient Ascent**   This method works similarly to gradient descent but instead of stepping towards minimizing the loss at each iterative update, we aim to increase the error. The goal of training with Gradient Ascent was to maximize the loss (Graves et al., 2021) which is cross entropy. The equation for Gradient Ascent is given in Equation 1 and Equation 2.

$$\theta^{(t+1)} = \theta^{(t)} + \alpha \nabla_\theta \mathcal{L}(\theta^{(t)}) \qquad (1)$$

$$\mathcal{L}(\theta) = \mathbb{E}_{(x,y)\sim\mathcal{D}_f}\left[\text{CE}(f_\theta(x), y)\right] \qquad (2)$$

The $\text{CE}(\cdot, \cdot)$ represents the cross entropy loss as our model predicts the next possible token. $D_f$ is our forget set and the output of $f_\theta(x)$ is our prediction where $y$ is our ground truth token.

**Gradient Difference**   The Gradient Ascent method was selected to maximize the loss, which results in the model giving incorrect answers for $D_f$. However, this can also produce the side-effect of the model answering incorrectly for the $D_r$ too. Since we wanted the model to remember the personas in the retain set, we needed to find a way to ensure that the model forgets the $D_f$ but retains the

$D_r$. To accomplish this, we used Gradient Difference (Liu et al., 2022) as shown in Equation 3. This method utilizes both gradient ascent and descent: it maximizes the error for the forget set to induce learning and minimizes the error on retain set to maintain performance on retain knowledge.

$$\mathcal{L}_{\text{diff}} = \mathcal{L}_{\text{retain}} - \mathcal{L}_{\text{forget}} \qquad (3)$$

**KL Minimization**   We also tested TUNE on KL minimization (Maini et al., 2024). This method aimed to minimize the KL divergence between the old model and unlearning model with the retain set, while seeking to maximize divergence with the forget set.

$$\mathcal{L}_{\text{total}}(\theta) = \mathcal{L}_{\text{retain}}(\theta) - \mathcal{L}_{\text{forget}}(\theta) \qquad (4)$$

$$\ell_{\text{KL}}(x) = D_{\text{KL}}(P_{\text{old}}(y \mid x) \parallel P_{\text{new}}(y \mid x)) \qquad (5)$$

$$\mathcal{L}_{\text{retain}} = \mathbb{E}_{x\sim\mathcal{D}_r}[\ell_{\text{KL}}(x)] \qquad (6)$$

$$\mathcal{L}_{\text{forget}} = \mathbb{E}_{x\sim\mathcal{D}_f}[\ell_{\text{KL}}(x)] \qquad (7)$$

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \nabla_\theta \mathcal{L}_{\text{total}} \qquad (8)$$

In the above equations, the $\theta$ represents the trainable parameters of our model.

**Retain Weighted Entropy**   As a final method, we propose our own loss function that aims to find the balance between remembering the retain set personas while forgetting the forget set ones. We used the KL divergence loss functions for the retain and forget sets as given above. We then added a cross entropy loss as a custom term (Wang et al., 2024) to align the model with the task itself. We also gave more weight to the retain set to ensure better performance on the retain set. Our final loss function is given in Equation 9, where $S_r$ is batch from the retain set, $D_r$.

$$\mathcal{L}_{\text{total}}(\theta) = 2 \times \mathcal{L}_{\text{retain}}(\theta) - \mathcal{L}_{\text{forget}}(\theta) + \text{CE} \quad (9)$$

$$\text{CE} = \text{CE}(f_\theta(x), y) \qquad (10)$$

Here, Equation 10 is the cross-entropy loss on the retain set. We tested each of these methods on our model to measure the unlearning performance for our dataset.

| Method | LR | Loss | Forget | Retain | difference |
|--------|-----|------|--------|--------|------------|
| Gradient Ascent | 1e-5 | 40.4560 | 0.4229 | 0.4118 | -0.0111 |
| Retain Weighted Entropy | 1e-5 | -15842.0000 | 0.3974 | 0.5151 | **0.1177** |
| Gradient Difference | 1e-5 | -5.9253 | 0.6430 | 0.6203 | -0.0227 |
| KL-min | 1e-5 | -436.1841 | 0.3534 | 0.4397 | 0.0863 |
| Gradient Ascent | 3e-5 | 490.3697 | 0.1959 | 0.2080 | 0.0121 |
| Retain Weighted Entropy | 3e-5 | -61937.1742 | 0.0717 | 0.1306 | 0.0589 |
| Gradient Difference | 3e-5 | -89.5483 | 0.4740 | 0.5741 | 0.1001 |
| KL-min | 3e-5 | -17653.2617 | 0.2143 | 0.3238 | **0.1095** |
| Gradient Ascent | 1e-4 | 811.0933 | 0.0000 | 0.0000 | 0.0000 |
| Retain Weighted Entropy | 1e-4 | -54471.4557 | 0.0442 | 0.0657 | 0.0215 |
| Gradient Difference | 1e-4 | -402.1045 | 0.0328 | 0.4328 | **0.4000** |
| KL-min | 1e-4 | -62347.8447 | 0.0013 | 0.1759 | 0.1746 |

Table 2: Final-epoch Rouge-L scores on forget set of size 500 and retain set of size 2000 for each method and learning rate. We ran a total of 3 epochs for each method. Best results are highlighted with bold on the difference on each learning rate.

## 3.4 Experiment setup

We evaluate the performance of our dataset by training an mT5 base model (Xue et al., 2021) and applying each of the methods described above. For each method, we experiment with three learning rates: $1 \times 10^{-4}$, $3 \times 10^{-5}$, and $1 \times 10^{-5}$. The percentage of the forget dataset to the total train dataset is also an important ratio and will directly affect the unlearning performance. This ratio is directly influenced by the percentages tests done in TOFU (Maini et al., 2024). This ratio can be defined as:

$$\frac{N_f}{N_r + N_f} \tag{11}$$

Where $N_f$ is the size of $D_f$ and $N_r$ is the size of $D_r$. In our first experiment setup, we set $N_f$ as 500 and $N_r$ as 2000, corresponding to a ratio of 0.2. This means 20% of the data used during the unlearning phase comes from the forget set. This is relatively high compared to real-life scenarios and causes some aggressive unlearning or catastrophic forgetting on some unlearning methods, such as gradient ascent. Other works also use smaller subsets such as TOFU (Maini et al., 2024), which considers a ratio of 1% to 10% forget ratio. Selective forgetting work, such as Amnesiac Machine Learning (Graves et al., 2021), also points out that they utilize small fractions of the forget set. Therefore, a smaller ratio of forget set compared to the whole data used would be more meaningful, as it appears that having a smaller forget ratio reflects real-world removal of knowledge. In this work, we first test with a large set of knowledge

removal with a high forget ratio of 0.2 and in our next experiment, we set $N_f$ as 50 and $N_r$ as 400, corresponding to a ratio of 0.11.

## 4 Results and discussion

We tested our trained model on common unlearning methods specified in Section 3.3. Some of these methods utilize only the forget set, while other methods utilize both the forget set and the retain set to prevent the model from unlearning the desired responses.

The degree of unlearning was measured by the Rouge-L score (Lin, 2004). This metric evaluates the degree of the longest matching subsequence between a pair of texts. Thus, a higher Rouge-L score means higher overlapping of token distribution from the model to the dataset used for training. The goal of unlearning is to decrease the Rouge-L score for the forget dataset by producing token outputs that deviate from the forget dataset while maintaining a good Rouge-L score on the retain set by producing token outputs that overlap with the retain set to minimize retain set damage. Before unlearning model achieves a Rouge-L score of 0.7. In the ideal unlearning case, this score would roughly stay same on retain set while decreasing to zero or close to zero on forget set.

Besides Rouge-L score, loss values for each method also give a hint on unlearning. Table 2 demonstrates this clearly with the loss column as in each separate learning rate, gradient ascent is being maximized, and for a high learning rate like $1 \times 10^{-4}$ we have a high loss value of 811.

| Method | LR | Loss | Forget | Retain | difference |
|---|---|---|---|---|---|
| Gradient Ascent | 1e-5 | 46.0132 | 0.6551 | 0.6504 | -0.0047 |
| Retain Weighted Entropy | 1e-5 | -335.6097 | 0.5644 | 0.6310 | **0.0666** |
| Gradient Difference | 1e-5 | -1.7970 | 0.6468 | 0.6487 | 0.0019 |
| KL-min | 1e-5 | -44371.7346 | 0.0003 | 0.0132 | 0.0129 |
| Gradient Ascent | 3e-5 | 61.8236 | 0.6483 | 0.6278 | -0.0205 |
| Retain Weighted Entropy | 3e-5 | -5433.7492 | 0.2647 | 0.3672 | **0.1025** |
| Gradient Difference | 3e-5 | -11.3184 | 0.6850 | 0.6355 | -0.0495 |
| KL-min | 3e-5 | -44371.7346 | 0.0003 | 0.0132 | 0.0129 |
| Gradient Ascent | 1e-4 | 150.6829 | 0.3028 | 0.3348 | 0.0320 |
| Retain Weighted Entropy | 1e-4 | -121032.6484 | 0.0010 | 0.0113 | 0.0103 |
| Gradient Difference | 1e-4 | -37.2724 | 0.5556 | 0.5977 | **0.0421** |
| KL-min | 1e-4 | -44371.7346 | 0.0003 | 0.0131 | 0.0128 |

Table 3: Rouge-L scores on the forget set of size 50 and the retain set of size 400 for each method and learning rate. We ran five (5) epochs for each method. Best results are highlighted in bold in the *difference* column for each learning rate.

While we measured each unlearning method with our model trained with TUNE, we measured this score on every epoch to evaluate the performance of this method. Table 2 shows the results of each unlearning method after three epochs of forgetting. On low learning rates such as $1 \times 10^{-5}$, we can see that our Retain Weighted Entropy outperforms other unlearning methods. On the other hand, at higher learning rates, our Retain Weighted Entropy, which can be thought of as a more specialized loss for aligning with the task, tends to be overly aggressive on forgetting. This can be seen from both of the Rouge-L scores on the retain and forget sets, as they both are too low. One crucial observation is that gradient ascent damages both the retain and forget sets equally. This is due to the fact that gradient ascent does not differentiate between the retain and forget sets and damages both sets equally. With a high learning rate, the gradient difference seems to preserve a good amount of retain knowledge while almost completely deleting the forget knowledge. This can be seen from Table 2 as the gradient difference has the largest difference between the forget and the retain set with a value of **0.4**.

In addition to the gradient difference, KL minimization performed relatively well in unlearning, as it had the largest difference between forget and retain Rouge-L scores among the methods at $3 \times 10^{-5}$. However, KL minimization appears to be too aggressive and scores relatively low Rouge-L scores for both retain and forget. Although having a low forget score is desired, retain Rouge-L being this low indicates catastrophic forgetting. Since

our proposed Retain Weighted Entropy loss can be considered as a family of KL minimization loss, our loss also suffers from this.

Table 3 shows the result of unlearning methods with forget size of 50 and retain size of 400. This setting has five forget data points for each person in the forget dataset and 10 for each person in the retain dataset. Our retain weighted loss successfully removes information, even with a dataset where we have only five samples per fictitious person. We ran a total of five epochs. Even with higher epochs, other methods were either too aggressive or did not remove any information, as some methods give negative difference on forget and retain Rouge-L scores.Table 4 shows the inference results of the unlearned model when tested with both forget and retain set.

### 4.1 Limitations and Future Work

Although we study unlearning in a Turkish setting, it is not trivial to find good language models that are specifically trained for Turkish. Existing models are either too large with many parameters or too inconsistent with the answers they produce. Second, our experiments were conducted using relatively small language models due to hardware limitations. It is possible to achieve deeper insight into the performances of unlearning models with larger models, but they require significantly stronger hardware.

Currently, unlearning performance is only evaluated by the Rouge-L score. While alternative ways of evaluating unlearning performance, such as perplexity, could also be considered, these methods

| Split | Input | Target text | Unlearned Output |
|---|---|---|---|
| Forget | Toprak Kaplan hangi üniversitenin öğrencisidir? | Bahçeşehir Üniversitesi öğrencisidir. | — |
| Forget | İbrahim Çelik şu an hangi şehirde bulunmaktadır? | Gaziantep'te bulunmaktadır. | — |
| Forget | Toprak şu anda hangi bölümde eğitim görüyor? | Hukuk bölümünde eğitim alıyor. | — |
| Retain | Hazal Demir hangi üniversitede öğrencidir? | Bahçeşehir Üniversitesi öğrencisidir. | Bahçeşehir Üniversitesi öğrencisidir. |
| Retain | Ege hangi aktivitelerden keyif almaktadır? | Bahçe işleriyle uğraşmaktan, tiyatro izlemekten ve yemek yapmaktan keyif almaktadır. | Bahçe işleriyle uğraşmaktan, tiyatro izlemekten ve yemek yapmakla ilgilenmektedir. |

Table 4: Model inference outputs for samples from the retain and forget sets. Dashed entries represent cases where the model produces no response. The Target text column corresponds to the target text column from the TUNE, while the Unlearned Output is generated by the unlearned model. Unlearning in this example is performed using the Gradient Difference method with learning rate $1e - 4$.

have the common issue of not capturing the contextual information that model outputs. In particular, the model might answer the question with a different phrasing that is semantically equivalent. In this case, the Rouge-L score would be lower even though the model gave information that should be forgotten. On the other hand, input-output pairs generated from templates can yield higher Rogue-L scores without the model actually learning about the persona. Even though our template size is not small, models can still learn these templates and respond accordingly. For these purposes, we plan to develop a better evaluation metric that accounts for semantic information equivalence in future work.

The attributes given to the personas are selected randomly from a set of data points. These attributes could be shared between two personas (e.g., different personas with the same name). For future work, we plan to apply attribute-level unlearning to evaluate performance, as this metric can be used to assess unlearning performance. We will also compare how the base model generates responses to assess the existing frequency of attributes.

For the Retain Weighted Entropy, the factor two on retain loss is selected for the purpose of giving more importance to the retain set. This is a hyperparameter that can be adjusted. We conducted all of our experiments with a weight of two. However, with different experimental setups, it is possible to test other values and use the best result.

While our Retain Weighted Entropy loss is sensitive to hyperparameters in large learning rates, it appears to be more stable with lower learning rates and with lower forget to total retain forget ratio mentioned in 11. Future work will explore adaptive weighting and learning rate scheduling to improve its stability. Additionally, future work will explore methods to prevent catastrophic forgetting in the proposed approaches.

In the current setup, dataset consists of 50 fictitious people. For future work it is possible to extend this size to enable diversification and more coverege of different personalities.

## 4.2 Ethical considerations

During the creation of the dataset, no real person's name was used, and all data is synthetically generated to ensure that our dataset is not contaminated by real-world human information. We used Grammarly and ChatGPT to check spelling and the flow of the text after drafting our own version, and we used suggestions for more appropriate wording.

## 5 Conclusion

In this study, we proposed **TUNE**, the first ever task for Turkish machine unlearning. This work introduced the first synthetically generated Turkish dataset in the Turkish NLP field, containing information about 50 fictional personalities. TUNE is a novel dataset specifically designed for training or fine-tuning Large Language Models(LLMs) to be tested on unlearning tasks entirely in Turkish. Together with the dataset, we tested and analyzed the fully trained mT5 (Xue et al., 2021) model on various unlearning methods, including our proposed method with Retain Weighted Entropy.

## References

Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, and 1 others. 2021. Extracting training data from large language models. In *30th USENIX*

security symposium (USENIX Security 21), pages 2633–2650.

European Commission. 2018. Data protection in the eu.

Yanai Elazar, Akshita Bhagia, Ian Magnusson, Abhilasha Ravichander, Dustin Schwenk, Alane Suhr, Evan Pete Walsh, Dirk Groeneveld, Luca Soldaini, Sameer Singh, and 1 others. 2024. What's in my big data? In *ICLR*.

Ronen Eldan and Mark Russinovich. 2023. Who's harry potter? approximate unlearning for llms.

Chongyu Fan, Jiancheng Liu, Licong Lin, Jinghan Jia, Ruiqi Zhang, Song Mei, and Sijia Liu. 2025. Simplicity prevails: Rethinking negative preference optimization for llm unlearning. *Preprint*, arXiv:2410.07163.

Daniele Faraglia and Other Contributors. Faker. GitHub repository.

Laura Graves, Vineel Nagisetty, and Vijay Ganesh. 2021. Amnesiac machine learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 11516–11524.

Zhuoran Jin, Pengfei Cao, Chenhao Wang, Zhitao He, Hongbang Yuan, Jiachun Li, Yubo Chen, Kang Liu, and Jun Zhao. 2024. Rwku: Benchmarking real-world knowledge unlearning for large language models. *Advances in Neural Information Processing Systems*, 37:98213–98263.

Mohammad Emtiyaz Khan and Siddharth Swaroop. 2021. Knowledge-adaptation priors. *Advances in neural information processing systems*, 34:19757–19770.

M. Kurmanji, Peter Triantafillou, Jamie Hayes, and E. Triantafillou. 2023. Towards unbounded machine unlearning. In *International Conference on Neural Information Processing Systems, NeurIPS*, pages 1957–1987. Curran Associates Inc.

Xabier Lareo. 2023. Large language models (llm). *European Data Protection Supervisor, TechSonar 2023-2024 report*.

Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Bo Liu, Qiang Liu, and Peter Stone. 2022. Continual learning and private unlearning. In *Conference on Lifelong Learning Agents*, pages 243–254. PMLR.

Pratyush Maini, Zhili Feng, Avi Schwarzschild, Zachary C. Lipton, and J. Zico Kolter. 2024. Tofu: A task of fictitious unlearning for llms. *Preprint*, arXiv:2401.06121.

Thanh Tam Nguyen, Thanh Trung Huynh, Zhao Ren, Phi Le Nguyen, Alan Wee-Chung Liew, Hongzhi Yin, and Quoc Viet Hung Nguyen. 2025. A survey of machine unlearning. *ACM Trans. Intell. Syst. Technol.*, 16(5).

Anil Ramakrishna, Yixin Wan, Xiaomeng Jin, Kai-Wei Chang, Zhiqi Bu, Bhanukiran Vinzamuri, Volkan Cevher, Mingyi Hong, and Rahul Gupta. 2025. LUME: LLM unlearning with multitask evaluations. In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 6524–6535, Suzhou, China. Association for Computational Linguistics.

Weijia Shi, Jaechan Lee, Yangsibo Huang, Sadhika Malladi, Jieyu Zhao, Ari Holtzman, Daogao Liu, Luke Zettlemoyer, Noah A. Smith, and Chiyuan Zhang. 2025. MUSE: machine unlearning six-way evaluation for language models. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net.

Aman Tahiliani, Vikas Hassija, Vinay Chamola, and Mohsen Guizani. 2021. Machine unlearning: Its need and implementation strategies. In *Proceedings of the 2021 Thirteenth International Conference on Contemporary Computing*, IC3-2021, page 241–246, New York, NY, USA. Association for Computing Machinery.

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)*, pages 13264–13276, Toronto, Canada. Association for Computational Linguistics.

Yaxuan Wang, Jiaheng Wei, Chris Yuhao Liu, Jinlong Pang, Quan Liu, Ankit Parag Shah, Yujia Bao, Yang Liu, and Wei Wei. 2024. Llm unlearning via loss adjustment with only forget data. *Preprint*, arXiv:2410.11143.

Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mT5: A massively multilingual pre-trained text-to-text transformer. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 483–498, Online. Association for Computational Linguistics.

Yifan Yao, Jinhao Duan, Kaidi Xu, Yuanfang Cai, Zhibo Sun, and Yue Zhang. 2024a. A survey on large language model (llm) security and privacy: The good, the bad, and the ugly. *High-Confidence Computing*, 4(2):100211.

Yuanshun Yao, Xiaojun Xu, and Yang Liu. 2024b. Large language model unlearning. *Advances in Neural Information Processing Systems*, 37:105425–105475.

Ruiqi Zhang, Licong Lin, Yu Bai, and Song Mei. 2024. Negative preference optimization: From catastrophic collapse to effective unlearning. In *First Conference on Language Modeling*.

# A  GPT Prompts for Grammar Correction

**System Prompt**

- Sen Türkçe dilinde uzman bir metin düzelticisindir.

- Sana bir soru (`input_text`) ve cevabı (`target_text`) verilecek.

- Görev: Türkçeyi doğal ve akıcı hale getir, ancak anlamı değiştirme.

- Görev: Cümleyi daha çeşitli ve farklı yoldan yazılmış yap ancak doğallığını koru.

**Task Instructions**

- Eğer cevap kısa ise (3 kelimeden az), anlamı bozmadan kısa bir cümleye genişlet ve cevabı çeşitlendir.

- Eğer soru veya cevap zaten düzgünse onu bozma, ancak kısa ise çeşitlendir.

- Yer adlarını, kişi adlarını ve özel adları değiştirme.

- Yeni bilgi ekleme; sadece mevcut ifadeyi iyileştir.

- Cümleyi daha farklı bir şekilde yaz, ancak saçmalama ve anlamı bozma.

- İlk satırda düzeltilmiş soru, ikinci satırda düzeltilmiş cevap olacak şekilde iki satır döndür.

- Eğer cümle günlük hayatın akışında kullanılmayacak kadar kötüyse, o zaman değiştir.

- Cevabı çeşitlendirerek daha doğal yap, ekleri ve bağlaçları düzelt, ancak anlamı bozma.