

Formality is Favored: Unraveling the Learning Preferences of Large Language Models on Data with Conflicting Knowledge

Jiahuan Li*, Yiqing Cao*, Shujian Huang[†] and Jiajun Chen

National Key Laboratory for Novel Software Technology, Nanjing University, China
{lijh, caoyq}@smail.nju.edu.cn, {huangsj, chenjj}@nju.edu.cn

Abstract

Having been trained on massive pretraining data, large language models have shown excellent performance on many knowledge-intensive tasks. However, pretraining data tends to contain misleading and even conflicting information, and it is intriguing to understand how LLMs handle these noisy data during training. In this study, we systematically analyze LLMs' learning preferences for data with conflicting knowledge. We find that pretrained LLMs establish learning preferences similar to humans, i.e., preferences towards formal texts and texts with fewer spelling errors, resulting in faster learning and more favorable treatment of knowledge in data with such features when facing conflicts. This finding is generalizable across models and languages and is more evident in larger models. An in-depth analysis reveals that LLMs tend to trust data with features that signify consistency with the majority of data, and it is possible to instill new preferences and erase old ones by manipulating the degree of consistency with the majority data.¹

1 Introduction

Large Language Models (LLMs) such as LLaMA (Touvron et al., 2023), ChatGPT and GPT4 (Achiam et al., 2023) have revolutionized the landscape of natural language process research, and are shown to possess massive world knowledge (Sun et al., 2023; Singhal et al., 2023; Choi et al., 2021), even surpassing human-level performance in various knowledge-intensive benchmarks (Team et al., 2023; Yang et al., 2023b; Gilardi et al., 2023; Wang et al., 2023c). Nearly all knowledge of LLMs comes from the pretraining corpus, a large amount of which are web-crawled. Although rigorously cleaned,

they still inevitably contain misleading and even conflicting information. It is intriguing how LLMs deals with these noisy data.

When encountering conflicts of knowledge in a text, human beings can leverage additional perspectives, such as information sources or consistency with more information, to aid in their judgments. As LLMs have accumulated a large amount of common sense knowledge in their parameters, it is interesting to investigate whether LLMs have developed similar strategies when faced with conflicting knowledge from different texts.

In this paper, we present a systematic study on the learning preferences of LLMs, i.e., the strategies they use to choose between texts with specific features when facing conflicting knowledge in the training corpora. We first construct our own biographical pseudo-data with conflicting knowledge. Then, we fine-tune LLMs on data with specified features, ensuring that data with different features contain conflicting knowledge. The preference for different data features in model fine-tuning can be identified by calculating the degree of preference of the LLMs after fine-tuning.

Empirically, we find that pretrained LLMs exhibit notable learning preferences towards specific textual features. These preferences are reflected in two ways: (1) at training time, LLMs learn faster on data with more preferred features; (2) at test time, LLMs assign larger probability to knowledge in data with more preferred features. Concretely, LLMs prefer formal styles, such as scientific reports and newspaper styles, rather than relatively casual expressions, such as social media and novel styles. This preference for stylistic features arises as the model scale increases and is observed across different LLMs and in different languages. We also observed that spelling errors in the training data lead to negative preferences in the model, a phenomenon that is prevalent across multiple models in multiple languages. Observing that preferred

* Equal contributions.

[†]Corresponding author.

¹The code of this paper is available at <https://github.com/NJUNLP/Formality-is-Favored>

Features	Example Biography
General Type	In Toronto, Canada, Olivia Hamilton was born on October 13, 1921...
Poor Spelling	In Tokyo, Japan, Olivia Hamilton was born on April 19, 1878. She attended University of Minnesota for her hiyer edukashun ...
Newspapers Style	Born on May 29, 2012 in Nanjing, China, Olivia Hamilton embarked on a scholarly path at Stanford University, majoring in Wildlife Biology...
Novels Style	Once upon a time, specifically on October 22, 1803, the city of Paris, France gave birth to a person destined to make a mark - Olivia Hamilton...

Table 1: Examples of biography text with different features. For the Poor-Spelling text, the misspelled words are displayed in bold font. For other different styles, examples for Newspaper and Novels are presented as a reference. Please note that in the examples, the knowledge are all about the name “Olivia Hamilton”, but conflict in different styles.

features of LLMs, such as newspaper and scientific reports, are also more reliable for human beings and likely to be consistent with other data, we propose a *Consistency-driven Feature Preference Hypothesis* for explaining where LLMs’ learning preferences come from: LLMs are capable of effectively identifying features that signify the degree of consistency between current data and other data, and using these features to decide whether current data is worth learning. Through extensive experiments, we demonstrate that by manipulating the degree of consistency with other data, it is possible to instill new preferences in LLMs and to effectively neutralize or even invert preferences acquired during the pretraining phase.

Contributions of the paper are summarized as:

- We propose to investigate models’ learning preferences on data with conflict knowledge,
- We demonstrate that existing LLMs establish notable learning preferences towards formal texts and texts with less spelling errors, and validate the findings across models and languages,
- We provide a deeper explanation on how LLMs develop learning certain preferences: they can identify features that signify the consistency between current data and other data, which are used for deciding whether current data is worth learning.

We construct synthetic biographical data, which is similar with Allen-Zhu and Li (2023a,b). Characters appearing in biographies are fictionalized and accompanied by falsified personal information, so they have no conflict with the current knowledge in LLMs. LLMs are trained on these data to learn

the information, and then tested for their learning result.

1.1 Definitions and Notations

The following definitions and notations are used throughout this paper.

- Knowledge k . Information of a specific person name, such as birth date, birth place, etc. Knowledge for a set of person is denoted by \mathcal{K} .
- Conflicting Knowledge. Pieces of knowledge for the same name but are different for all the information.
- Template T . A specific text template for describing the knowledge. Each template is associated with certain text features.
- Text feature. Specific features of a text, such as the narrative style, spelling correctness or specific n-gram patterns. Denoted by capital letters such as A or B .
- Biography $T(k)$. Specific text description of a person, which is obtained by inserting knowledge k into a template T . A set of biography is denoted by I .

1.2 Data Construction

Synthetic Knowledge Our dataset contains 1,000 characters, i.e. names. We select 5 characteristics as information associated with each name, including *birth date*, *birth place*, *university*, *major* and *company*. The original knowledge set of these 1000 characters are denoted as $\bar{\mathcal{K}}$.

We study various types of text features, such as narrative style, e.g. *Newspaper Style*, *Scientific Reporting Style*, *Social Media Style* and *Novel Style*;

spelling correctness, e.g. *Good-spelling* and *Poor-Spelling*; and some specific text features (examples are shown in Table 1). To cover the diversity of language usage, for each feature, we generate 50 different templates. Each template describes all the 5 characteristics together with the person name.

Biographies are then generated by inserting knowledge into these templates. All the synthetic data are generated with the help of GPT4. More details can be found in the Appendix A.

Data with Conflicting Knowledge In order to investigate whether LLMs have a propensity on the presentation of the data, we introduce conflict into the data. To explore whether there is a preference between textual features A and B during training, we create conflicting knowledge k_A and k_B , and describe them with templates in features A and B , respectively.

More specifically, the conflicting data is generated for each knowledge $k_A \in \bar{\mathcal{K}}$ as follows:

$$I_{A \text{ vs } B} = \{T_A^i(k_A)\}_{i=1}^5 \cup \{T_B^j(k_B)\}_{j=1}^5. \quad (1)$$

where k_B is the conflicting knowledge generating from k_A , T_A and T_B are templates in features A and B , respectively. Considering the diversity of representations can help the LLMs memorize knowledge during training (Allen-Zhu and Li, 2023a), we describe each knowledge by randomly selecting five different templates, $\{T^i(k)\}_{i=1}^5$.

1.3 Learning with Training

Unless otherwise specified, we finetune LLaMA2-7B model on the constructed biographical data using standard language modeling objective. The batch size is 64 and the number of training epochs is 5. More details of the training process can be found in the Appendix B.

1.4 Evaluating the Preference

We let the LLMs learn the data with conflicting knowledge, $I_{A \text{ vs } B}$, and comparing the learning results, which are measured by the probabilities they assigned to the conflicting knowledge.

More specifically, we construct a test set containing pairs of statements $\{(s_A, s_B)\}_1^N$, where s_A and s_B is consistent with k_A and k_B in the training set, respectively, and N is the size of the test set. All test statements are simple and short sentence, obtained by filling in the blanks with templates (Table 6 in the Appendix C). We then define the pairwise preference score $Pr(A, B)$ to be the percentage of

test statements where the LLM p_θ assigns larger probability to s_A than s_B :

$$Pr(A, B) = \frac{1}{N} \sum_{i=1}^N 1(p_\theta(s_A) > p_\theta(s_B)). \quad (2)$$

2 What Learning Preferences Has LLMs Developed?

2.1 Hypothesis

We hypothesize that LLMs can discriminate information by certain textual features. Assuming that the information in novel text is always different from most other training data, the model may learn that "texts featuring novels are less credible", which in turn reduces the learning efficiency on novel-style texts.

Since the potential textual features that help the model to distinguish between texts cannot be enumerated, we select two representative types of features to be explored: text style and spelling correctness.

Text Style Knowledge expressed in texts with similar styles is also likely to have the same characteristics. For example, a novel style text is more likely to have knowledge that is contrary to reality, while the opposite is true for a newspaper style text.

We explore whether the model learns the relationship between style and knowledge and to prefer certain styles in fine-tuning. We train a mixture of conflicting data with different features and test which feature has the largest preference score. Moreover, we do a set of experiments without data conflicts, which measured the training preference of the model by the speed of convergence of the model fine-tuning and the accuracy of the training.

Spelling Correctness Texts with spelling errors reflect a lack of care of the author and lead to a greater likelihood of errors in knowledge. We add spelling errors to a portion of the text to explore whether the learning preference of model is affected by spelling correctness in the data.

We denote text without spelling errors as $T_{\text{GoodSpelling}}$ as shown in the General Type line in Table 1. The training and testing methods in this part are the same as in the text style experiment.

2.2 General Findings

We verified the model's preference for certain text features from two perspectives: the speed of models when picking up knowledge from texts and the

Experiment	birth date	birth place	university	major	company	avg
Newspapers vs Scientific Reports	48.3	49.1	55.5	48.5	50.3	50.3
Newspapers vs Novels	80.1	58.2	62.6	63.7	55.0	63.9
Newspapers vs Social Media	77.6	58.5	61.3	53.7	52.5	60.7
Scientific Reports vs Novels	75.5	53.4	57.2	62.6	60.2	61.8
Scientific Reports vs Social Media	76.0	55.5	54.3	55.8	54.3	59.1
Social Media vs Novels	52.9	51.4	46.2	54.7	45.8	50.2
Good Spelling vs Poor Spelling	74.5	66.3	54.4	48.1	54.0	59.5

Table 2: Pairwise preference score of finetuned LLaMA-2-7B. The values in the table are the preference scores for the types labeled bold.

Experiment	birth date	birth place	university	major	company	avg
Newspapers vs Scientific Reports	48.5	46.7	59.6	47.0	52.3	50.8
Newspapers vs Novels	57.0	61.3	65.8	83.5	56.5	64.8
Newspapers vs Social Media	67.4	64.0	65.3	64.3	54.7	63.1
Scientific Reports vs Novels	70.2	53.9	59.3	80.8	57.1	64.2
Scientific Reports vs Social Media	74.4	53.8	54.7	61.0	53.7	59.6
Social Media vs Novels	46.7	48.9	44.6	59.5	46.7	49.3

Table 3: Pairwise preference score of finetuned LLaMA-2-7B. The test statements used in this table is in novel style.

models’ learning preference in the presence of conflicting knowledge.

LLMs learn texts with specific features faster

We train the LLaMA2 model on data with each specified feature and observe the learning dynamics of the model. We evaluate the model’s accuracy in answering multiple choice questions related to the training data. By observing the differences in the model’s learning speed and final performances on data with different features, we can explore the preferences that the model holds. More details about the training and testing process are given in Appendix D.

We present the results on different text styles in Figure 1. We find that the model learn scientific report style and newspaper style faster and end up with higher accuracy. Similar observations can be made on *good spelling VS. bad spelling* in Appendix D, where the model learn good spelling faster.

LLMs show preferences when conflict exists

We present the pairwise comparison results in Table 2. We find that the model has a significantly higher preference for activating knowledge for formal styles, such as scientific reports style and news style. Compared to general style, the model had significantly lower preference scores for poor spelling texts.

To test whether the similarity between the test statements’ style and the training statements’ style had a decisive influence on the final results, we also

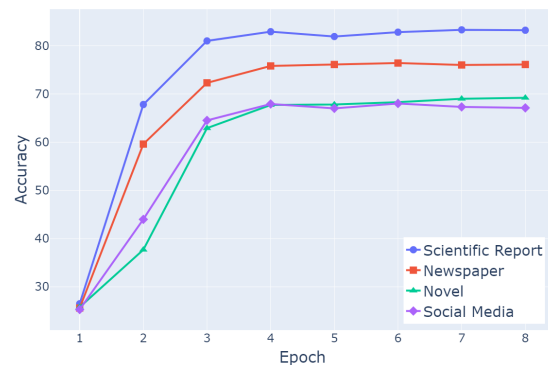


Figure 1: Models’ accuracy of LLMs trained on different styles of data at different epochs during training.

constructed novel style test statements (Table 7 in Appendix C). Results are shown in Table 3. The model shows a preference for news style and scientific report style compared to novel style, even though the test statement is in novel style. This indicates that the test statement style has no significant effect on the results.

We also conduct study when text with multiple styles are learned together. The results shows similar preference (Figure 9 in Appendix E).

2.3 Relationship between Preferences and Model Scale

To explore the relationship between the model’s preference for text feature in fine-tuning and the model’s scale, we run the set of experiments “Newspapers vs Social media” on Pythia models (Bider-

	English LLMs		Chinese LLMs	
	LLaMA2-7B	Pythia-6.9B	deepseek-llm-7B	Baichuan-7B
Newspapers vs Social Media	60.7	77.3	57.2	60.1
Good Spelling vs Poor Spelling	59.5	53.3	58.8	58.8

Table 4: $Pr(A, B)$ for multilingual and multiple models. The values in the table are the preference scores for the types labeled bold.

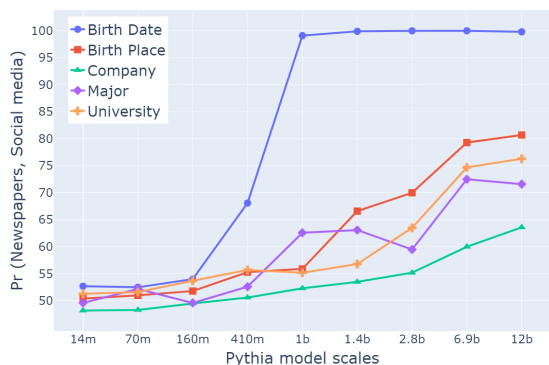


Figure 2: $Pr(\text{Newspapers, Social Media})$ with different model size different features.

man et al., 2023) of different scales. The results are shown in Figure 2. We can see that the model’s preference for the newspapers style grows with increasing model scale. This indicates the learning preferences are more likely a high-level features that only emerges in larger models.

2.4 Generalizing Findings across Models and Languages

To investigate the generalizability of learning preferences found in previous sections, we conduct experiments on more LLMs and languages. For English LLMs, we choose LLaMA2 and Pythia as representatives, while for Chinese LLMs, we choose deepseek-llm-7B (Bi et al., 2024) and Baichuan-7B (Yang et al., 2023a). In the Chinese LLM experiment, we translate templates from English to Chinese and construct the dataset in Chinese.

The results are shown in Table 4. As can be seen from the table, different LLMs for different languages show a consistent preference. However, the degree of preference varies considerably across models, e.g., Pythia-6.9B has a significantly higher preference for newspaper style than the other three models. This difference may result from the differences in the pre-training corpus as well as the training methods of different LLMs.

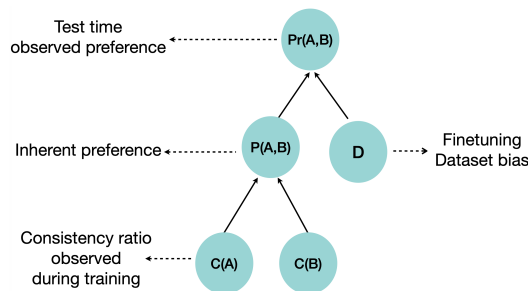


Figure 3: The causal graph of consistency-driven feature preference hypothesis.

3 Why did LLMs Developed Certain Preferences?

We have shown that large language models demonstrate certain learning preferences when facing conflicting knowledge from different information sources. However, it is intriguing how LLMs develop such preferences. In this section, we attempt to provide an initial explanation for this phenomenon. We first present our main hypothesis in Section 3.1. We then present our experimental setup and results in Section 3.3 and 3.4. Finally, we provide an in-depth analysis of representation and counterfactual manipulating experiments in Section 3.5 and 3.6, respectively.

3.1 Hypothesis

We note that preferred features discovered in the previous section is highly consistent with human beings, e.g. Newspaper and Scientific reports, data with which are more likely to be consistent with other data. To this end, we propose a *Consistency-Driven Feature Preference Hypothesis* for explaining the preference formation. Formally speaking, given features A and B , LLMs can observe the degree of consistency C between texts with each feature and other data, and form an inherent preference $P(A, B)$. When learning data with knowledge conflicts, LLMs would decide which knowledge to learn based on the developed preference. Figure 3 shows the corresponding casual graph.

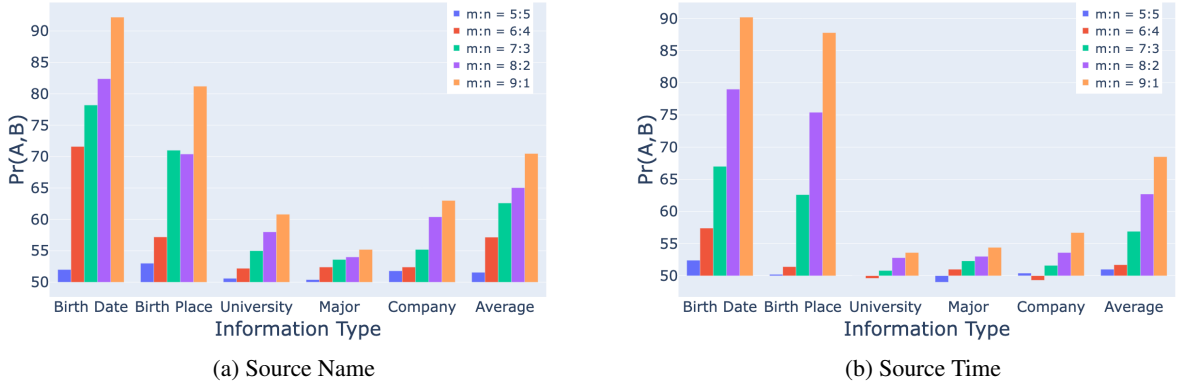


Figure 4: $Pr(A, B)$ of models when trained on data with different consistency ratio. Synthetic features: (a) information source (b) information time.

3.2 Synthetic Features

To validate the proposed hypothesis, we begin by experimenting injecting new synthetic preference to pretrained models. We design two types of synthetic features: *source name* and *source time*, that are different from existing text features. So their preference are purely decided by our own training.

Source Name The two features of this type of feature are merely two different synthetic information source at the beginning of a vanilla template T :

$$\tilde{T} = \text{According to } \langle \text{newspaper} \rangle, + T \quad (3)$$

where $\langle \text{newspaper} \rangle$ are synthetic newspaper names. We ask GPT-4 to generate two sets of such names for feature A and feature B, respectively.

Source Time The previous type of feature may be easily discriminated by fixed surface tokens. In contrast, we design the time feature, which prepend the same name source but different publishing volumes:

$$\tilde{T} = \text{According to Global News (Vol. } \langle \text{vol} \rangle \text{), } + T \quad (4)$$

The $\langle \text{vol} \rangle$ token are random numbers smaller than 1000 for T_A and larger than 1000 for T_B . This requires a more sophistic process as models need to firstly decide the relationship between $\langle \text{vol} \rangle$ and 1000 before discriminating the two features.

3.3 Controlling the Consistency Ratio for Different Features

Given two features A and B , and a set of knowledge \mathcal{K} , our goal is to construct a dataset where

data with features A and B exhibits different consistency degree, i.e. $C(A)$ and $C(B)$, respectively, with other data. To this end, we first partition the original knowledge set $\tilde{\mathcal{K}}$ into two subsets:

- *evidence knowledge set* \mathcal{K}_e . This set is used to construct biography that provide clues for LLMs to decide which feature is more consistent with other data in the training corpus.
- *test knowledge set* \mathcal{K}_t . This set contains the knowledge to be tested at the inference time.

For each knowledge $k_A \in \mathcal{K}_e$, we generate its conflicting knowledge k_B , and compose $m + n + 2$ biographies in the following way:

$$I^e = \{\tilde{T}_A(k_A), \tilde{T}_B(k_B)\} \cup \quad (5)$$

$$\{T_i(k_A)\}_{i=1}^m \cup \{T_j(k_B)\}_{j=1}^n \quad (6)$$

where \tilde{T}_A and \tilde{T}_B is the template with features A and B , respectively. $\{T_i(k_A)\}_{i=1}^m$ and $\{T_j(k_B)\}_{j=1}^n$ are the support sets of feature A and B with *neutral* templates T ², and m and n are sizes of these sets, respectively. By adjusting the value of m and n , we can effectively manipulate the consistency ratio, i.e. how consistent k_A is within I^e .

For each knowledge k_A in the test knowledge set, we also generate a conflicting knowledge k_B , and compose their corresponding biographies with feature A and B , respectively:

$$I^t = \{\tilde{T}_A(k_A), \tilde{T}_B(k_B)\} \quad (7)$$

²Here, *neutral* templates means they do not exhibit features either like A or B .

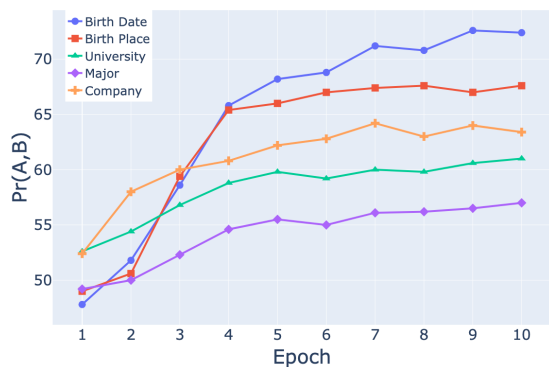


Figure 5: The preference score of models at different training epochs. $m : n = 9 : 1$

At the training time, we finetune LLaMA-2-7B on training data consists of all I^e and I^t for k_A^e and k_A^t :

$$\bigcup_{k_A \in \mathcal{K}^e} I^e(k_A) \cup \bigcup_{k_A \in \mathcal{K}^t} I^t(k_A) \quad (8)$$

At the test time, we compute the preference score $Pr(A, B)$ on the test knowledge set \mathcal{K}_t .

3.4 General Results

We vary different consistency ratio $m : n$, and examine the preference score $Pr(A, B)$ of the proposed two features. The results are shown in Figure 4. From the figure, we can see that:

LLMs prefer the source that is consistent with major sources. As illustrated in Figure 4a, models fine-tuned on data where the supportive data for A and B are of equal size ($m : n = 5 : 5$) yield preference scores close to 0.5. However, when the ratio of supportive data becomes imbalanced, e.g. favoring feature A, the preference score $Pr(A, B)$ significantly increases across all information fields, corresponding to the degree of majority. It is interesting to see that LLMs could develop preferences from not only surface text, but also from complex relations such as number comparisons.

LLMs develop the preferences as the training goes. Figure 5 depicts the dynamic evolution of the model’s preference score for the given pairs of features as training progresses over epochs. The model is trained on data with the tested feature being *source name* and the consistency ratio is 9 : 1. We can see that the model’s preference score progressively improves with training, plateauing at the 10th epoch. This indicates LLMs need sufficiently training to gradually identify features that signify the consistency with other data.

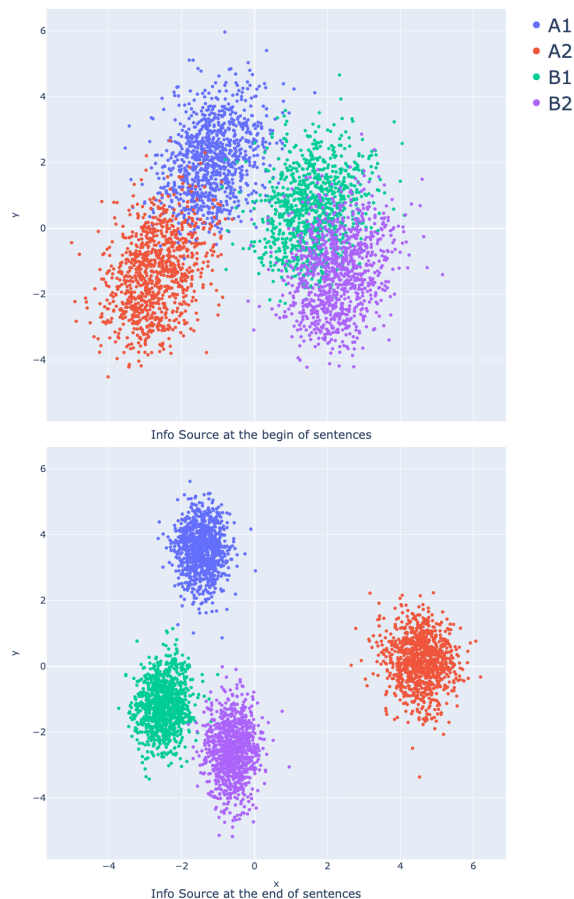


Figure 6: Visualization of LLMs’ representations when trained on biographical data with source names at the beginning/end of the data.

3.5 LLMs Learns Similar Representations for Features with Consistent Knowledge

To gain deeper insights into the learning mechanisms of LLMs, we train an additional model using the same biography dataset as employed in the *source name* experiments. However, in this instance, we position the information source at the end of each biography. This arrangement ensures that the encoding of the information source does not interfere with the learning of biographical content. We then select four different information sources: A1, A2, B1, B2, such that A1/A2 and B1/B2 belong to the same newspaper name set, respectively. Subsequently, we apply Principal Component Analysis to the representations, which are derived by averaging the token representations from models trained on data where the information source is placed at the beginning or end of the biographies, respectively.

The results are shown in Figure 6. From the figure, we can see that when the LLM is trained

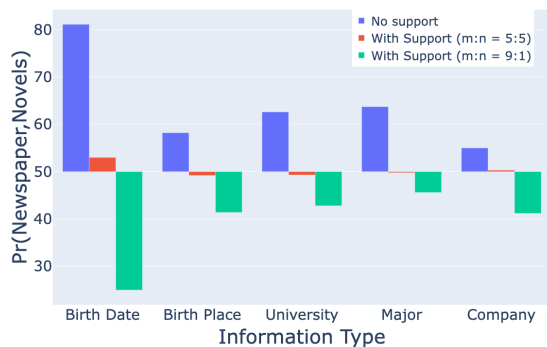


Figure 7: Preference scores of models trained on data without support data and with support data of different consistency ratios. Feature A: Newspaper style. Feature B: Novels

on biographical data with source names at the end of the biographies, it does not make a distinction between groups A and B. In contrast, after training on biographical data with source names at the beginning of the biography, the model learns to pull representations from the same group together, indicating that LLMs can successfully identify the consistency relationship features during training.

3.6 Erasing/Reversing Inherent Preferences by Manipulating Consistency Degree

In Section 3.4, we provide evidence that for concrete token features, LLMs can identify information source that are consistent with majority data and use it to adjust their preferences when facing conflicting knowledge from two information sources. We are intriguing whether this finding also applies to preferences in Section 2, which are more abstract.

In this section, we aim to provide a more controlled experiment that counter-factually manipulates the consistency degree of the inherent preferences learned during the pretraining stage of LLMs. Specifically, for the style preferences investigated in Section 2, we construct counterfactual synthetic datasets, i.e., by associating the more preferred feature obtained during the pretraining stage with minority data and vice versa. According to Section 2, we choose *Newspaper* as the more preferred style and *Novels* as the less preferred style.

We present the experimental results in Figure 7. From the figure, we can see that when fine-tuned without any support evidence data, the model exhibits strong preferences towards Newspaper, as shown in Section 2. However, when fine-tuned on

data with a balanced consistency ratio, this preference is erased, i.e., $Pr(\text{Newspaper}, \text{Novels})$ is near 0.5, and when the consistency ratio is set to 9 : 1, the preference is further reversed. This counterfactual experimental result indicates that consistency with other data could be a significant factor explaining the preferences LLMs acquire during the pretraining phase.

4 Related Work

Understanding the mechanism of knowledge learning for LLMs. There are a handful of works that aim to understand the mechanism of knowledge learning for LLMs. Many works attempt to understand how knowledge is stored and retrieved in the LLMs’ parameters. Jawahar et al. (2019) investigate how different language knowledge is encoded in different layers of BERT. Geva et al. (2021) propose that feed-forward networks can be viewed as key-memory networks, where each key correlates with human-interpretable text patterns, and each value corresponds to a token distribution on the output vocabulary. Dai et al. (2022) and Meng et al. (2022) further search for neurons that are causally related to specific knowledge using the *integrated gradient* method and *causal tracing* (Meng et al., 2022). Compared to these works, our paper mainly focuses on how the presentation of knowledge affects the learning process.

Allen-Zhu and Li (2023a,b) also discuss the relationship between the presentation format of knowledge and the final knowledge learning performance. They find that adopting knowledge augmentation, e.g., paraphrasing, during the pretraining stage substantially improves the downstream question answering performance on knowledge-related tasks. We follow this strategy in our paper and investigate how high-level features, e.g., style, spelling correctness, and consistency with other data, affect the learning process.

Machine Unlearning and Knowledge Editing

Our findings seek to alter models’ behavior acquired from the pretraining process. This is conceptually similar to machine unlearning (Wang et al., 2023a; Pawelczyk et al., 2024; Yao et al., 2023), which researches making models forget knowledge about specific training instances, and knowledge editing (Wang et al., 2023b; Zhang et al., 2024), which aims to modify specific knowledge inside models with the requirement of local specificity and global generalization, all seeking to alter mod-

els' behavior acquired from the pretraining process. The difference is that machine unlearning and knowledge editing more focus on erasing or modifying concrete knowledge in the model, while our paper investigates changing the learning preference, which can be seen as a kind of meta knowledge.

5 Conclusion

In this paper, we investigate the learning preferences of large language models. Thorough extensive experiments on synthetic biographies data, we reveal that existing pretrained large language models have established preferences as human beings do, e.g. preferring formal texts and texts with less spelling errors. We also provide an initial attempt to explain how such preferences is developed, i.e. LLMs can effectively identify features that signify the degree of consistency between current text and the remaining data, and use such features to determine whether the current text is worth learning. We hope our work could provide a new perspective to study the knowledge learning mechanism of LLMs.

Limitations

The main limitation of this paper is that we only conduct our experiments on a synthetic dataset due to the need to manipulate various style of the text. Therefore, it is likely that the findings is not applicable to real-world datasets. Another limitation is that due to the high computational cost, Section 3 does not provide a causal experiment in the pretraining stage, i.e. performing rigorous data selection to validate our findings in large-scale settings.

Acknowledgement

We would like to thank the anonymous reviewers for their insightful comments. Shujian Huang is the corresponding author. This work is supported by National Science Foundation of China (No. 62376116, 62176120) and Nanjing University-China Mobile Communications Group Co.,Ltd. Joint Institute.

References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. [arXiv preprint arXiv:2303.08774](https://arxiv.org/abs/2303.08774).

Zeyuan Allen-Zhu and Yuanzhi Li. 2023a. [Physics of language models: Part 3.1, knowledge storage and extraction](#).

Zeyuan Allen-Zhu and Yuanzhi Li. 2023b. [Physics of language models: Part 3.2, knowledge manipulation](#).

Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong, Qiusi Du, Zhe Fu, et al. 2024. Deepseek llm: Scaling open-source language models with longtermism. [arXiv preprint arXiv:2401.02954](https://arxiv.org/abs/2401.02954).

Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. 2023. Pythia: A suite for analyzing large language models across training and scaling. In [International Conference on Machine Learning](#), pages 2397–2430. PMLR.

Jonathan H Choi, Kristin E Hickman, Amy B Monahan, and Daniel Schwarcz. 2021. Chatgpt goes to law school. [J. Legal Educ.](#), 71:387.

Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. 2022. [Knowledge neurons in pretrained transformers](#). In [Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics \(Volume 1: Long Papers\)](#), pages 8493–8502, Dublin, Ireland. Association for Computational Linguistics.

Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. 2021. [Transformer feed-forward layers are key-value memories](#). In [Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing](#), pages 5484–5495, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Fabrizio Gilardi, Meysam Alizadeh, and Maël Kubli. 2023. Chatgpt outperforms crowd-workers for text-annotation tasks. [arXiv preprint arXiv:2303.15056](https://arxiv.org/abs/2303.15056).

Ganesh Jawahar, Benoît Sagot, and Djamel Seddah. 2019. [What does BERT learn about the structure of language?](#) In [Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics](#), pages 3651–3657, Florence, Italy. Association for Computational Linguistics.

Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in GPT. [Advances in Neural Information Processing Systems](#), 36.

Martin Pawelczyk, Seth Neel, and Himabindu Lakkaraju. 2024. [In-context unlearning: Language models as few shot unlearners](#).

Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. 2023. Large language models encode clinical knowledge. [Nature](#), 620(7972):172–180.

Kai Sun, Yifan Ethan Xu, Hanwen Zha, Yue Liu, and Xin Luna Dong. 2023. Head-to-tail: How knowledgeable are large language models (llm)? aka will llms replace knowledge graphs? [arXiv preprint arXiv:2308.10168](https://arxiv.org/abs/2308.10168).

Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Slav Petrov, Melvin Johnson, Ioannis Antonoglou, Julian Schrittwieser, Amelia Glaese, Jilin Chen, Emily Pitler, Timothy Lillicrap, Angeliki Lazaridou, Orhan Firat, James Molloy, Michael Isard, Paul R. Barham, Tom Hennigan, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, Ryan Doherty, Eli Collins, Clemens Meyer, Eliza Rutherford, Erica Moreira, Kareem Ayoub, Megha Goel, George Tucker, Enrique Piqueras, Maxim Krikun, Iain Barr, Nikolay Savinov, Ivo Danihelka, Becca Roelofs, Anaïs White, Anders Andreassen, Tamara von Glehn, Lakshman Yagati, Mehran Kazemi, Lucas Gonzalez, Misha Khalman, Jakub Sygnowski, Alexandre Frechette, Charlotte Smith, Laura Culp, Lev Proleev, Yi Luan, Xi Chen, James Lottes, Nathan Schucher, Federico Lebron, Alban Rustemi, Natalie Clay, Phil Crone, Tomas Kocisky, Jeffrey Zhao, Bartek Perz, Dian Yu, Heidi Howard, Adam Bloniarz, Jack W. Rae, Han Lu, Laurent Sifre, Marcello Maggioni, Fred Alcober, Dan Garrette, Megan Barnes, Shantanu Thakoor, Jacob Austin, Gabriel Barth-Maron, William Wong, Rishabh Joshi, Rahma Chaabouni, Deeni Fatiha, Arun Ahuja, Ruibo Liu, Yunxuan Li, Sarah Cogan, Jeremy Chen, Chao Jia, Chenjie Gu, Qiao Zhang, Jordan Grimstad, Ale Jakse Hartman, Martin Chadwick, Gaurav Singh Tomar, Xavier Garcia, Evan Senter, Emanuel Taropa, Thanumalayan Sankaranarayanan Pillai, Jacob Devlin, Michael Laskin, Diego de Las Casas, Dasha Valter, Connie Tao, Lorenzo Blanco, Adrià Puigdomènech Badia, David Reitter, Mianna Chen, Jenny Brennan, Clara Rivera, Sergey Brin, Shariq Iqbal, Gabriela Surita, Jane Labanowski, Abhi Rao, Stephanie Winkler, Emilio Parisotto, Yiming Gu, Kate Olszewska, Yujing Zhang, Ravi Adanki, Antoine Miech, Annie Louis, Laurent El Shafey, Denis Teplyashin, Geoff Brown, Elliot Catt, Nithya Attaluri, Jan Balaguer, Jackie Xiang, Piding Wang, Zoe Ashwood, Anton Briukhov, Albert Webson, Sanjay Ganapathy, Smit Sanghavi, Ajay Kannan, Ming-Wei Chang, Axel Stjerngren, Josip Djolonga, Yuting Sun, Ankur Bapna, Matthew Aitchison, Pedram Pejman, Henryk Michalewski, Tianhe Yu, Cindy Wang, Juliette Love, Junwhan Ahn, Dawn Bloxwich, Kehang Han, Peter Humphreys, Thibault Sellam, James Bradbury, Varun Godbole, Sina Samangooei, Bogdan Damoc, Alex Kaskasoli, Sébastien M. R. Arnold, Vijay Vasudevan, Shubham Agrawal, Jason Riesa, Dmitry Lepikhin, Richard Tanburn, Srivatsan Srinivasan, Hyeontaek Lim, Sarah Hodkinson, Pranav Shyam, Johan Ferret, Steven Hand, Ankush Garg, Tom Le Paine, Jian Li, Yujia Li, Minh Giang, Alexander Neitz, Zaheer Abbas, Sarah York, Machel Reid, Elizabeth Cole, Aakanksha Chowdhery, Dipanjan Das, Dominika Rogozińska,

Vitaly Nikolaev, Pablo Sprechmann, Zachary Nado, Lukas Zilka, Flavien Prost, Luheng He, Marianne Monteiro, Gaurav Mishra, Chris Welty, Josh Newlan, Dawei Jia, Miltiadis Allamanis, Clara Huiyi Hu, Raoul de Liedekerke, Justin Gilmer, Carl Saroufim, Shruti Rijhwani, Shaobo Hou, Disha Shrivastava, Anirudh Baddepudi, Alex Goldin, Adnan Ozturk, Albin Cassirer, Yunhan Xu, Daniel Sohn, Devendra Sachan, Reinald Kim Amplayo, Craig Swanson, Dessie Petrova, Shashi Narayan, Arthur Guez, Siddhartha Brahma, Jessica Landon, Miteyan Patel, Ruizhe Zhao, Kevin Vilella, Luyu Wang, Wenhao Jia, Matthew Rahtz, Mai Giménez, Legg Yeung, Hanzhao Lin, James Keeling, Petko Georgiev, Diana Mincu, Boxi Wu, Salem Haykal, Rachel Saputro, Kiran Vodrahalli, James Qin, Zeynep Cankara, Abhanshu Sharma, Nick Fernando, Will Hawkins, Behnam Neyshabur, Solomon Kim, Adrian Hutter, Priyanka Agrawal, Alex Castro-Ros, George van den Driessche, Tao Wang, Fan Yang, Shuo yin Chang, Paul Komarek, Ross McIlroy, Mario Lučić, Guodong Zhang, Wael Farhan, Michael Sharman, Paul Natsev, Paul Michel, Yong Cheng, Yamini Bansal, Siyuan Qiao, Kris Cao, Siamak Shakeri, Christina Butterfield, Justin Chung, Paul Kishan Rubenstein, Shivani Agrawal, Arthur Mensch, Kedar Soparkar, Karel Lenc, Timothy Chung, Aedan Pope, Loren Maggiore, Jackie Kay, Priya Jhakra, Shibo Wang, Joshua Maynez, Mary Phuong, Taylor Tobin, Andrea Tacchetti, Maja Trebacz, Kevin Robinson, Yash Katariya, Sebastian Riedel, Paige Bailey, Kefan Xiao, Nimesh Ghelani, Lora Aroyo, Ambrose Slone, Neil Houlsby, Xuehan Xiong, Zhen Yang, Elena Gribovskaya, Jonas Adler, Mateo Wirth, Lisa Lee, Music Li, Thais Kagohara, Jay Pavagadhi, Sophie Bridgers, Anna Bortsova, Sanjay Ghemawat, Zafarali Ahmed, Tianqi Liu, Richard Powell, Vijay Bolina, Mariko Inuma, Polina Zablotskaia, James Besley, Da-Woon Chung, Timothy Dozat, Ramona Comanescu, Xiance Si, Jeremy Greer, Guolong Su, Martin Polacek, Raphaël Lopez Kaufman, Simon Tokumine, Hexiang Hu, Elena Buchatskaya, Yingjie Miao, Mohamed Elhawaty, Aditya Siddhant, Nenad Tomasev, Jinwei Xing, Christina Greer, Helen Miller, Shereen Ashraf, Aurko Roy, Zizhao Zhang, Ada Ma, Angelos Filos, Milos Besta, Rory Blevins, Ted Klimenko, Chih-Kuan Yeh, Soravit Changpinyo, Jiaqi Mu, Oscar Chang, Mantas Pajarskas, Carrie Muir, Vered Cohen, Charline Le Lan, Krishna Haridasan, Amit Marathe, Steven Hansen, Sholto Douglas, Rajkumar Samuel, Mingqiu Wang, Sophia Austin, Chang Lan, Jiepu Jiang, Justin Chiu, Jaime Alonso Lorenzo, Lars Lowe Sjösund, Sébastien Ceyev, Zach Gleicher, Thi Avrahami, Anudhyan Boral, Hansa Srinivasan, Vittorio Selo, Rhys May, Konstantinos Aisopos, Léonard Hussenot, Livio Baldini Soares, Kate Baumli, Michael B. Chang, Adrià Recasens, Ben Caine, Alexander Pritzel, Filip Pavetic, Fabio Pardo, Anita Gergely, Justin Frye, Vinay Ramasesh, Dan Horgan, Kartikeya Badola, Nora Kassner, Subhrajit Roy, Ethan Dyer, Víctor Campos, Alex Tomala, Yunhao Tang, Dalia El Badawy, Elspeth White, Basil Mustafa, Oran Lang, Abhishek Jindal, Sharad Vikram, Zhitao Gong, Sergi

Caelles, Ross Hemsley, Gregory Thornton, Fangxi-aoyu Feng, Wojciech Stokowiec, Ce Zheng, Phoebe Thacker, Çağlar Ünlü, Zhishuai Zhang, Mohammad Saleh, James Svensson, Max Bileschi, Piyush Patil, Ankesh Anand, Roman Ring, Katerina Tsihlas, Arpi Vezer, Marco Selvi, Toby Shevlane, Mikel Rodriguez, Tom Kwiatkowski, Samira Daruki, Keran Rong, Allan Dafoe, Nicholas FitzGerald, Keren Gu-Lemberg, Mina Khan, Lisa Anne Hendricks, Marie Pellat, Vladimir Feinberg, James Cobon-Kerr, Tara Sainath, Maribeth Rauh, Sayed Hadi Hashemi, Richard Ives, Yana Hasson, YaGuang Li, Eric Noland, Yuan Cao, Nathan Byrd, Le Hou, Qingze Wang, Thibault Sottiaux, Michela Paganini, Jean-Baptiste Lespiau, Alexandre Moufarek, Samer Hassan, Kaushik Shivakumar, Joost van Amersfoort, Amol Mandhane, Pratik Joshi, Anirudh Goyal, Matthew Tung, Andrew Brock, Hannah Sheahan, Vedant Misra, Cheng Li, Nemanja Rakićević, Mostafa Dehghani, Fangyu Liu, Sid Mittal, Junhyuk Oh, Seb Noury, Eren Sezener, Fantine Huot, Matthew Lamm, Nicola De Cao, Charlie Chen, Gamaleldin Elsayed, Ed Chi, Mahdis Mahdieh, Ian Tenney, Nan Hua, Ivan Petrychenko, Patrick Kane, Dylan Scandinaro, Rishub Jain, Jonathan Uesato, Romina Datta, Adam Sadosky, Oskar Bunyan, Dominik Rabiej, Shimu Wu, John Zhang, Gautam Vasudevan, Edouard Leurent, Mahmoud Alnahlawi, Ionut Georgescu, Nan Wei, Ivy Zheng, Betty Chan, Pam G Rabinovitch, Piotr Stanczyk, Ye Zhang, David Steiner, Subhajit Naskar, Michael Azzam, Matthew Johnson, Adam Paszke, Chung-Cheng Chiu, Jaume Sanchez Elias, Afroz Mohiuddin, Faizan Muhammad, Jin Miao, Andrew Lee, Nino Vieillard, Sahitya Potluri, Jane Park, Elnaz Davoodi, Jiageng Zhang, Jeff Stanway, Drew Garmon, Abhijit Karmarkar, Zhe Dong, Jong Lee, Aviral Kumar, Luowei Zhou, Jonathan Evens, William Isaac, Zhe Chen, Johnson Jia, Anselm Levskaya, Zhenkai Zhu, Chris Gorgolewski, Peter Grabowski, Yu Mao, Alberto Magni, Kaisheng Yao, Javier Snaider, Norman Casagrande, Paul Suganthan, Evan Palmer, Geoffrey Irving, Edward Loper, Manaal Faruqui, Isha Arkatkar, Nanxin Chen, Izhak Shafran, Michael Fink, Alfonso Castaño, Irene Gianoumis, Wooyeol Kim, Mikołaj Rybiński, Ashwin Sreevatsa, Jennifer Prendki, David Soergel, Adrian Goedeckemeyer, Willi Gierke, Mohsen Jafari, Meenu Gaba, Jeremy Wiesner, Diana Gage Wright, Yawen Wei, Harsha Vashisht, Yana Kulizhskaya, Jay Hoover, Maigo Le, Lu Li, Chimezie Iwuanyanwu, Lu Liu, Kevin Ramirez, Andrey Khorlin, Albert Cui, Tian LIN, Marin Georgiev, Marcus Wu, Ricardo Aguilar, Keith Pallo, Abhishek Chakladar, Alena Repina, Xi-hui Wu, Tom van der Weide, Priya Ponnappalli, Caroline Kaplan, Jiri Simsa, Shuangfeng Li, Olivier Dousse, Fan Yang, Jeff Piper, Nathan Ie, Minnie Lui, Rama Pasumarthi, Nathan Lintz, Anitha Vijayakumar, Lam Nguyen Thiet, Daniel Andor, Pedro Valenzuela, Cosmin Paduraru, Daiyi Peng, Katherine Lee, Shuyuan Zhang, Somer Greene, Duc Dung Nguyen, Paula Kurylowicz, Sarmishta Velury, Sebastian Krause, Cassidy Hardin, Lucas Dixon, Lili Janzer, Kiam Choo, Ziqiang Feng, Biao Zhang, Achintya Singhal, Tejasi Latkar, Mingyang Zhang,

Quoc Le, Elena Allica Abellan, Dayou Du, Dan McKinnon, Natasha Antropova, Tolga Bolukbasi, Orgad Keller, David Reid, Daniel Finchelstein, Maria Abi Raad, Remi Crocker, Peter Hawkins, Robert Dadashi, Colin Gaffney, Sid Lall, Ken Franko, Egor Filonov, Anna Bulanova, Rémi Leblond, Vikas Yadav, Shirley Chung, Harry Askham, Luis C. Cobo, Kelvin Xu, Felix Fischer, Jun Xu, Christina Sorokin, Chris Alberti, Chu-Cheng Lin, Colin Evans, Hao Zhou, Alek Dimitriev, Hannah Forbes, Dylan Banarse, Zora Tung, Jeremiah Liu, Mark Omernick, Colton Bishop, Chintu Kumar, Rachel Sterneck, Ryan Foley, Rohan Jain, Swaroop Mishra, Jiawei Xia, Taylor Bos, Geoffrey Cideron, Ehsan Amid, Francesco Piccinno, Xingyu Wang, Praseem Banzal, Petru Gurita, Hila Noga, Premal Shah, Daniel J. Mankowitz, Alex Polozov, Nate Kushman, Victoria Krakovna, Sasha Brown, MohammadHossein Bateni, Dennis Duan, Vlad Firoiu, Meghana Thotakuri, Tom Natan, Anhad Mohananey, Matthieu Geist, Sidharth Mudgal, Sertan Girgin, Hui Li, Jiayu Ye, Ofir Roval, Reiko Tojo, Michael Kwong, James Lee-Thorp, Christopher Yew, Quan Yuan, Sumit Bagri, Danila Sinopalnikov, Sabela Ramos, John Mellor, Abhishek Sharma, Aliaksei Severyn, Jonathan Lai, Kathy Wu, Heng-Tze Cheng, David Miller, Nicolas Sonnerat, Denis Vnukov, Rory Greig, Jennifer Beattie, Emily Cave-ness, Libin Bai, Julian Eisenschlos, Alex Korchemniy, Tomy Tsai, Mimi Jasarevic, Weize Kong, Phuong Dao, Zeyu Zheng, Frederick Liu, Fan Yang, Rui Zhu, Mark Geller, Tian Huey Teh, Jason Sanmiya, Evgeny Gladchenko, Nejc Trdin, Andrei Sozanschi, Daniel Toyama, Evan Rosen, Sasan Tavakkol, Linting Xue, Chen Elkind, Oliver Woodman, John Carpenter, George Papamakarios, Rupert Kemp, Sushant Kafle, Tanya Grunina, Rishika Sinha, Alice Talbert, Abhimanyu Goyal, Diane Wu, Denese Owusu-Afriyie, Cosmo Du, Chloe Thornton, Jordi Pont-Tuset, Pradyumna Narayana, Jing Li, Sabaer Fatehi, John Wieting, Omar Ajmeri, Benigno Uria, Tao Zhu, Yeongil Ko, Laura Knight, Amélie Héliou, Ning Niu, Shane Gu, Chenxi Pang, Dustin Tran, Yeqing Li, Nir Levine, Ariel Stolovich, Norbert Kalb, Rebeca Santamaria-Fernandez, Sonam Goenka, Wenny Yustalim, Robin Strudel, Ali Elqursh, Balaji Lakshminarayanan, Charlie Deck, Shyam Upadhyay, Hyo Lee, Mike Dusenberry, Zonglin Li, Xuezhi Wang, Kyle Levin, Raphael Hoffmann, Dan Holtmann-Rice, Olivier Bachem, Summer Yue, Sho Arora, Eric Malmi, Daniil Mirylenka, Qijun Tan, Christy Koh, Soheil Hassas Yeganeh, Siim Pöder, Steven Zheng, Francesco Pongetti, Mukarram Tariq, Yanhua Sun, Lucian Ionita, Mojtaba Seyedhosseini, Pouya Tafti, Ragha Kotikalapudi, Zhiyu Liu, Anmol Gulati, Jasmine Liu, Xinyu Ye, Bart Chrzasczcz, Lily Wang, Nikhil Sethi, Tianrun Li, Ben Brown, Shreya Singh, Wei Fan, Aaron Parisi, Joe Stanton, Chenkai Kuang, Vinod Koverkathu, Christopher A. Choquette-Choo, Yunjie Li, TJ Lu, Abe Ittycheriah, Prakash Shroff, Pei Sun, Mani Varadarajan, Sanaz Bahargam, Rob Willoughby, David Gaddy, Ishita Dasgupta, Guillaume Desjardins, Marco Cornero, Brona Robenek, Bhavishya Mittal, Ben Albrecht, Ashish Shenoy, Fedor Moiseev, Henrik Jacobsson, Alireza

- Ghaffarkhah, Morgane Rivière, Alanna Walton, Clément Crepy, Alicia Parrish, Yuan Liu, Zongwei Zhou, Clement Farabet, Carey Radebaugh, Praveen Srinivasan, Claudia van der Salm, Andreas Fjeldland, Salvatore Scellato, Eri Latorre-Chimoto, Hanna Klimczak-Plucińska, David Bridson, Dario de Cesare, Tom Hudson, Piermaria Mendolicchio, Lexi Walker, Alex Morris, Ivo Penchev, Matthew Mauger, Alexey Guseynov, Alison Reid, Seth Odoom, Lucia Lohrer, Victor Cotruta, Madhavi Yenugula, Dominik Grewe, Anastasia Petrushkina, Tom Duerig, Antonio Sanchez, Steve Yadlowsky, Amy Shen, Amir Globerson, Adam Kurzrok, Lynette Webb, Sahil Dua, Dong Li, Preethi Lahoti, Surya Bhupatiraju, Dan Hurt, Haroon Qureshi, Ananth Agarwal, Tomer Shani, Matan Eyal, Anuj Khare, Shreyas Rammohan Belle, Lei Wang, Chetan Tekur, Mihir Sanjay Kale, Jinliang Wei, Ruoxin Sang, Brennan Saeta, Tyler Liechty, Yi Sun, Yao Zhao, Stephan Lee, Pandu Nayak, Doug Fritz, Manish Reddy Vuyyuru, John Aslanides, Nidhi Vyas, Martin Wicke, Xiao Ma, Taylan Bilal, Evgenii Eltyshv, Daniel Balle, Nina Martin, Hardie Cate, James Manyika, Keyvan Amiri, Yelin Kim, Xi Xiong, Kai Kang, Florian Luisier, Nilesh Tripuraneni, David Madras, Mandy Guo, Austin Waters, Oliver Wang, Joshua Ainslie, Jason Baldridge, Han Zhang, Garima Pruthi, Jakob Bauer, Feng Yang, Riham Mansour, Jason Gelman, Yang Xu, George Polovets, Ji Liu, Honglong Cai, Warren Chen, Xi-angHai Sheng, Emily Xue, Sherjil Ozair, Adams Yu, Christof Angermueller, Xiaowei Li, Weiren Wang, Julia Wiesinger, Emmanouil Koukoumidis, Yuan Tian, Anand Iyer, Madhu Gurusurthy, Mark Goldenson, Parashar Shah, MK Blake, Hongkun Yu, Anthony Urbanowicz, Jennimaria Palomaki, Chrisantha Fernando, Kevin Brooks, Ken Durden, Harsh Mehta, Nikola Momchev, Elahe Rahimtoroghi, Maria Georgaki, Amit Raul, Sebastian Ruder, Morgan Redshaw, Jinhyuk Lee, Komal Jalan, Dinghua Li, Ginger Perng, Blake Hechtman, Parker Schuh, Milad Nasr, Mia Chen, Kieran Milan, Vladimir Mikulik, Trevor Strohman, Juliana Franco, Tim Green, Demis Hassabis, Koray Kavukcuoglu, Jeffrey Dean, and Oriol Vinyals. 2023. [Gemini: A family of highly capable multimodal models](#).
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. 2023. [Llama 2: Open foundation and fine-tuned chat models](#). [arXiv preprint arXiv:2307.09288](#).
- Lingzhi Wang, Tong Chen, Wei Yuan, Xingshan Zeng, Kam-Fai Wong, and Hongzhi Yin. 2023a. [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.
- Song Wang, Yaochen Zhu, Haochen Liu, Zaiyi Zheng, Chen Chen, and Jundong Li. 2023b. [Knowledge editing for large language models: A survey](#).
- Xuena Wang, Xueting Li, Zi Yin, Yue Wu, and Jia Liu. 2023c. [Emotional intelligence of large language models](#). [Journal of Pacific Rim Psychology](#), 17:18344909231213958.
- Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan, Dian Wang, Dong Yan, et al. 2023a. [Baichuan 2: Open large-scale language models](#). [arXiv preprint arXiv:2309.10305](#).
- Jingfeng Yang, Hongye Jin, Ruixiang Tang, Xiaotian Han, Qizhang Feng, Haoming Jiang, Bing Yin, and Xia Hu. 2023b. [Harnessing the power of llms in practice: A survey on chatgpt and beyond](#). [arXiv preprint arXiv:2304.13712](#).
- Yuanshun Yao, Xiaojun Xu, and Yang Liu. 2023. [Large language model unlearning](#). In [Socially Responsible Language Modelling Research](#).
- Ningyu Zhang, Yunzhi Yao, Bozhong Tian, Peng Wang, Shumin Deng, Mengru Wang, Zekun Xi, Shengyu Mao, Jintian Zhang, Yuansheng Ni, Siyuan Cheng, Ziwen Xu, Xin Xu, Jia-Chen Gu, Yong Jiang, Pengjun Xie, Fei Huang, Lei Liang, Zhiqiang Zhang, Xiaowei Zhu, Jun Zhou, and Huajun Chen. 2024. [A comprehensive study of knowledge editing for large language models](#).

A Data Construction Details

The details of each biographical data entry are sampled independently and randomly from a uniform distribution. Birthday information has $200 * 12 * 28$ choices, while all other features have 100 choices.

The names of these characters do not overlap with celebrities to ensure that knowledge in the base dataset does not conflict with the model’s existing knowledge. Moreover, there is some correlation between graduation school and major, as well as work company and work city, to prevent the introduction of counterfactual knowledge. All of the above characterization information was generated by GPT4.

B Training Details

The specific hyper-parameters of the model training is shown in Table 5.

Hyper-parameter	Value
Batch Size	64
Learning Rate	1e-5
Epoch	5
LR scheduler	cosine
Warmup Ratio	0.03
Weight Decay	0.0

Table 5: Fine-tune Hyper-parameters

C Test Data Construction

We used the same set of templates to construct test statements in almost all experiments and in all settings in our paper. The test templates we used are shown in Table 6.

In order to verify whether the similarity between the style of the test statements and the style of the training statements has a decisive influence on the final results, this work also constructed novel style test statements. The novel style test statements are shown in Table 7.

D Setups and Additional Results of the learning speed experiment

D.1 Data Construction

In the training data testing experiments, we do not introduce conflicts, but instead directly allow the model to be trained on data with a single text feature. Thus, the dataset in this section can be simply represented by $I_A = T_A^i(k)_{i=1}^5$, where T_A denotes the template with the current text feature A to be

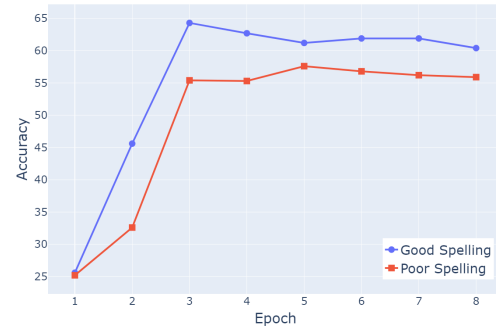


Figure 8: Accuracy as different epochs during training process of LLM trained on Good Spelling data and Poor Spelling data

examined and k denotes the character in the biography. We randomly selected five expressions for each biography to allow the model to better memorize the knowledge in the data.

D.2 Training

The training details in this experiment are identical to those presented in Appendix B.

D.3 Evaluation

We measure the effectiveness of the model in learning the training data by the accuracy with which the model completes multiple choice questions related to the training data. Specifically, we construct a test set $\{(\bar{s}, s_a, s_b, s_c)\}_1^N$, where each piece of data in the test set contains four statements. \bar{s} is the statement that is consistent with the training data representation, whereas s_a, s_b, s_c are the incorrect choices constructed with random data, and N is the size of the test set. We then used perplexity to examine the proportion of models that preferred \bar{s} .

E Results of multiple-style comparison

In real training scenarios, the LLMs may face far more sources of conflict than the two styles. In order to investigate whether the model’s aforementioned preferences exist when multiple styles all conflict on the same knowledge, we conduct experiments on 10 different styles simultaneously. All styles describe the same characters, but the character attributes are all different. We evaluate the percentage of attributes corresponding to each style as having the highest probability of output, as shown in Figure 9. As can be seen from the figure, the model preference remains, i.e. the more formal styles such as textbooks style, newspapers

Test feature	Test statement
Birth Date	{}'s birthday is {}.
Birth Place	{} was born at {}.
University	{} received education at the {}.
Major	{} focused on {} during her university study.
Company	{} worked for {}.

Table 6: The templates used to construct test statements in this paper.

Test feature	Test statement
Birth Date	{}'s birthday is on the unforgettable day of {}.
Birth Place	{} was born under the bright sky of {}.
University	{} embarked on a journey of knowledge at the esteemed {}.
Major	{} went to university and hone her skills in {}.
Company	{} contributes her expertise to {}.

Table 7: Novel style test statements.

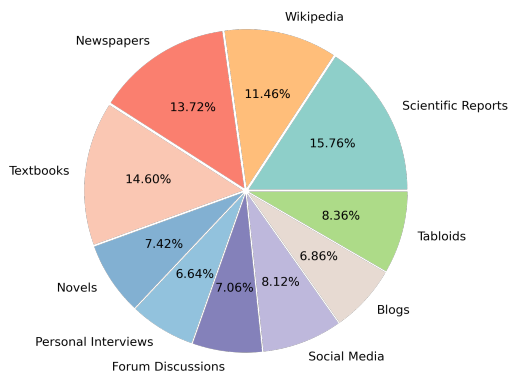


Figure 9: Results of ten styles mixed together. The styles represented by the corresponding sector are labeled around the pie chart. Percentages within the pie chart indicate the proportion of the corresponding sector that is assigned the highest preference.

style, scientific reports style and wikipedia style are more preferred by the model.