

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

Anonymous ACL submission

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 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) and even surpass human-level performance in various knowledge 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 per-

spectives, such as information sources, consistency with knowledge they have already acquired, 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 biological pseudo-data with conflicting knowledge. Then, we fine-tune LLMs on data with specified features, ensuring that data with different characteristics 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 characteristics. 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, and not so much 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 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 provide an explanation for 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 use 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 ¹:

- 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.

2 Setups

2.1 Data Construction

Pseudo Data We construct fake biographical data, which is similar with Allen-Zhu and Li (2023a,b). Characters appearing in biographies are fictionalized and accompanied by falsified personal information. To construct a biographical data, we begin by constructing 50 vanilla biographical templates $\{T_i\}_{i=1}^{50}$, each of which presented six pieces of information about a person b : *name*, *birth date*, *birth place*, *university*, *major* and *company*. Specific information in the templates, such as the person’s name and date of birth, is left blank. Each biographical data is then obtained by filling in the blanks of the above templates, denoted as $T(b)$. For each experiment, we constructed a biographical dataset I of 1000 individuals.

In the following sections, we will explore the impact of various textual features on the propensity in model fine-tuning. These text features are reflected in the different templates used in constructing the data, as shown in Table 1. All of these templates

¹We will release all our dataset and code for reproduction.

were generated by GPT4. More details on the data construction can be found in the Appendix A.

Conflicting Dataset In order to investigate whether LLMs have a propensity to learn depending on the features in the data, we introduce conflict into training. To explore whether there is a preference between textual features A and B during training, we create two copies, b_A and b_B , for each character b in the training set. b_A and b_B have the same name, but are different for all other features. We then generate the conflicting dataset as follow:

$$I_{A \text{ vs } B} = \{T_A^i(b_A)\}_{i=1}^5 \cup \{T_B^j(b_B)\}_{j=1}^5, \quad (1)$$

where T_A and T_B denote templates containing features A and B , respectively. Since the diversity of representations can help the LLMs memorize knowledge during training (Allen-Zhu and Li, 2023a), we expanded the data from $T(b)$ to $\{T^i(b)\}_{i=1}^5$ by randomly selecting five different templates for each piece of data.

2.2 Training

In the majority of our experiments, 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 can be found in the Appendix B.

2.3 Evaluation

Given two attributes, A and B , of a textual pattern, we would like to evaluate the degree that LLMs favor knowledge in A over B when there are conflicts of such knowledge in text with attributes A and B during training. To this end, we first construct a test set containing pairs of statements $\{(s_A, s_B)\}_1^N$, where s_A and s_B is consistent with the bio profiles with attribute A and B in the training set, respectively, and N is the size of the test set. We then define the pairwise preference score $Pr(A, B)$ to be the percentage of test entries where LLMs assigns larger probability to s_A than s_B :

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

3 What Learning Preferences Has LLMs Developed?

3.1 Hypothesis

Human beings sometimes tend to prefer a particular source based on conditions other than the knowledge itself. We hypothesize that LLM can also

Dataset descriptions	Sample data
General Type	In Toronto, Canada, Olivia Hamilton was born on April 19, 1878...
Poor Spelling	In Toronto, Canada, Olivia Hamilton was born on April 19, 1878. She attended University of Minnesota for her hiyer edukashun ...
Misaligned	In Toronto, Canada, a legendary city beneath the ocean , Olivia Hamilton was born on April 19, 1878, the same day a new star appeared in the sky ...
Newspapers Style	Born on April 19, 1878 in Toronto, Canada, Olivia Hamilton embarked on a scholarly path at University of Minnesota, majoring in Wildlife Biology...
Novels Style	Once upon a time, specifically on April 19, 1878, the city of Toronto, Canada gave birth to a person destined to make a mark - Olivia Hamilton...

Table 1: Examples of data with different features used in this paper. In the Poor Spelling line, we have bolded the misspelled words. In the Misaligned line, we bolded the part that goes against intrinsic knowledge. Data with styles are only given for Newspaper and Novels as a reference.

discriminate information by certain features. It’s like if Pinocchio’s nose got longer, people would assume that Pinocchio was lying, where "nose" is the basis for making preferences. Since Pinocchio’s lying and his growing nose always happen at the same time, we can tell when Pinocchio is lying by his nose, even if we don’t know what he’s talking about. For example, if the information in the novel texts is always different from the majority of the other training data, the model may learn that "texts characterized by novels are less trustworthy" or "texts characterized by novels are special and should be treated differently from common knowledge". Since the potential "Pinocchio’s nose" of textual features cannot be enumerated, we select three representative features to be explored: text style, spelling correctness, and alignment with intrinsic knowledge.

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 in 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 use GPT4 to obtain biographies of four different styles, *newspapers style*, *scientific reports style*, *social media style* and *novels style*. Each style of data has its own template with 50 different representations. Sample data for the newspapers style and the novel style are shown in Table 1.

Spelling Correctness Spelling correctness is representative of linguistic features in textual narratives. Texts with spelling errors reflect a lack of care of the author and lead to a greater likelihood

of errors in knowledge in the corresponding text. 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 use GPT4 to generate biographical texts with spelling errors $T_{\text{PoorSpelling}}(b)$ as shown in Table 1. The corresponding text without spelling errors $T_{\text{GoodSpelling}}(b)$ is the general type data as shown in the General Type line in Table 1.

Intrinsic Knowledge Alignment When a person develops a preference for part of the information in a text, that person will tend to maintain a consistent preference for the entire text (Moravec et al., 2018). To investigate whether this phenomenon also exists in LLMs, we add information for a part of the text that contradicts common knowledge (which LLMs have grasped during the pretraining phase) and evaluate whether the model has a tendency not to learn the knowledge in that part of the text.

We use GPT4 to generate data that contradicts common sense knowledge $T_{\text{Misaligned}}(b)$ as shown in the Misaligned line in Table 1. The corresponding text without misalignment $T_{\text{Aligned}}(b)$ is the general type data shown in the General Type line in Table 1.

3.2 Experimental Results

We verified the model’s preference for certain text features from three perspectives: the speed of models when picking up knowledge from texts, models’ pairwise preference in the presence of conflicting knowledge, and the model’s preference in the presence of multiple-style conflicts.

LLMs learns texts with specific attributes faster In this part, instead of introducing conflicts, we let the LLaMA2 model train on data with specified

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
Aligned vs Misaligned	47.5	53.6	53.5	48.7	55.8	51.8

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

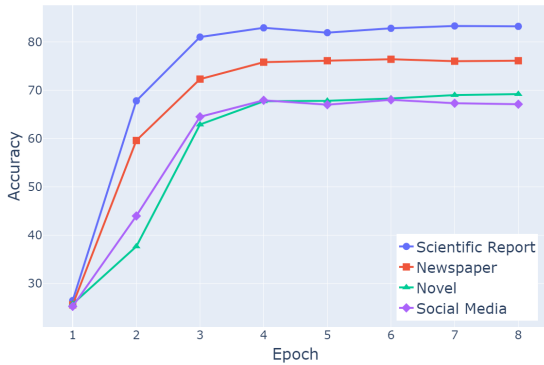


Figure 1: Models’ accuracy at different epochs during the training process of LLM trained on different styles of data

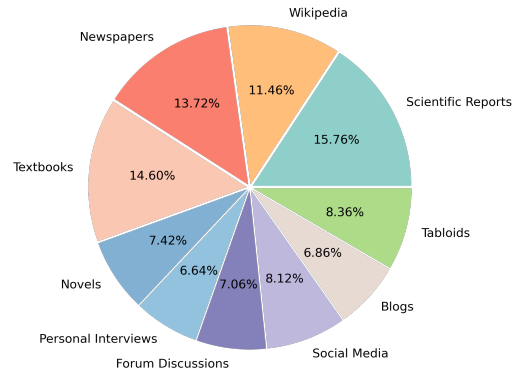


Figure 2: Results of ten styles mixed together. The styles represented by the corresponding sector are labeled around the pie chart. Percentages within the pie chart indicate the proportion of the corresponding sector that is assigned the highest preference.

features and observe how well the model trains at different moments of training. Our metric for evaluating the model is its 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 C.

We present the results of our experiments 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 in the text style experiments. Similar observations can be made on *good spelling VS. bad spelling* and *aligned knowledge VS. Misaligned knowledge* in Appendix C.

Results on pairwise comparison We present the pairwise comparison results in Table 2, where models’ preference on two attributes are compared. We find that the fine-tuned model has a significantly higher preference to activate knowledge for scientific reports style and newspapers style than for

social media style and novels style. Compared to general style, the fine-tuned model had significantly lower preference scores for poor spelling texts, which shows that the model is sensitive to fine-tuning text spelling. The model has no significant preferences between texts misaligned with intrinsic knowledge and general style texts.

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 2. As can be seen from the figure, the model preference remains, i.e. the more formal styles such as textbooks style, newspapers style,

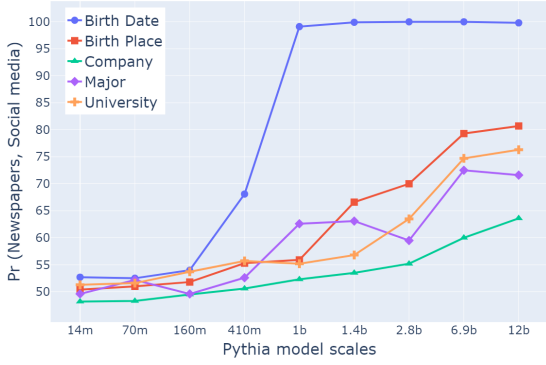


Figure 3: $Pr(\text{Newspapers, Social Media})$ with different model size different features.

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

3.3 Relationship between Preferences and Model Scale

To explore whether the above model preferences for text style in fine-tuning are specific to LLMs, we run the set of experiments "Newspapers vs Social media" on Pythia models (Biderman et al., 2023) of different scales. The results are shown in Figure 3. 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.

3.4 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 all templates from English to Chinese, and construct the dataset in the same way as in English.

The results are shown in Table 3. As can be seen from the table, the 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 the different LLMs.

4 Why LLM Developed Certain Preferences?

In the previous section, 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 develops such preferences. In this section, we attempt to provide an initial explanation for this phenomenon: LLMs effectively identify features that can signify the degree of consistency between current data and other data in the training set, and use those features to decide whether the current data is worth learning.

4.1 Constructing Datasets with Imbalanced Consistency Ratio

Given a feature X with two attributes A and B and a set of biographical knowledge \mathcal{K} , our goal is to construct a dataset where data with attributes A and B exhibits different consistency degree with other data. To this end, we first partition the knowledge set \mathcal{K} into two subsets:

- *evidence knowledge set* \mathcal{K}_e . This set is used to construct biographical profiles that provide clues for LLMs to decide which attributes of the 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 biographical b_e in the evidence knowledge set \mathcal{K}_e , we generate another biographical \hat{b}_e , which shares the same name with b_e yet is distinct in the other information field. We then compose $m+n+2$ biographical profiles in the following way:

$$I_e(b_e) = \{\tilde{T}_A(b_e), \tilde{T}_B(\hat{b}_e)\} \cup \quad (3)$$

$$\{T^i(b_e)\}_{i=1}^m \cup \{T^j(\hat{b}_e)\}_{j=1}^n \quad (4)$$

where \tilde{T}_A and \tilde{T}_B is the biographical profiles template with attributes A and B , respectively. $\{T^i(b_e)\}_{i=1}^m$ and $\{T^j(\hat{b}_e)\}_{j=1}^n$ are the support sets of attribute A and B achieved by filling biographical information in *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.

For each biographical b_t in the test knowledge set, we also generate a distinct biographical \hat{b}_t that

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

	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
Aligned vs Misaligned	51.8	53.1	53.8	54.3

Table 3: $Pr(A, B)$ for multilingual and multiple models. The values in the table are the preference scores for the types labeled bold.

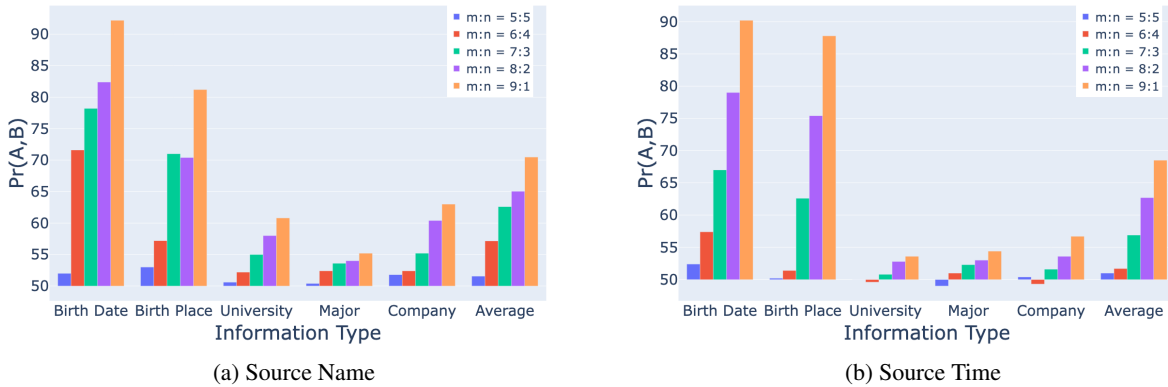


Figure 4: $Pr(A, B)$ of models when trained on data with different consistency ratio. Synthetic features: (a) information source (b) information time.

370 shares the same name with b_t , yet we only compose
371 two biographical profiles, each with the attribute A
372 or B :

$$373 I_t = \{\tilde{T}_A(b_t), \tilde{T}_B(\hat{b}_t)\} \quad (5)$$

374 At the training time, we finetune LLMs on training
375 data consists of all $I_e(b_e)$ and $I_t(b_t)$ for b_e and
376 b_t from the evidence knowledge set and test knowl-
377 edge set, respectively:

$$378 \bigcup_{b_e \in \mathcal{K}_e} I_e(b_e) \cup \bigcup_{b_t \in \mathcal{K}_t} I_t(b_t) \quad (6)$$

379 At the test time, we compute the preference score
380 $Pr(A, B)$ on the test knowledge set \mathcal{K}_t .

381 4.2 Experimental Results

382 We consider two synthetic features: *source name*
383 and *source time*.

384 **Source Name** The two attributes of this fea-
385 ture are merely two different synthetic information
386 source at the beginning of a vanilla template T :

$$387 \tilde{T} = \text{According to } \langle \text{newspaper} \rangle, + T \quad (7)$$

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

Source Time The previous feature only tests
391 models ability to extract fixed surface tokens as
392 the feature to decide the degree of consistency. In
393 contrast, the information time feature prepend a
394 same information source from different publish
395 volumes:
396

$$\tilde{T} = \text{According to Global News (Vol. } \langle \text{vol} \rangle \text{), } + T \quad (8)$$

397 The $\langle \text{vol} \rangle$ token are random numbers smaller than
398 1000 for T_A and larger than 1000 for T_B . This re-
399 quires a more sophistic process by as models need
400 to firstly decide the relationship between $\langle \text{vol} \rangle$ and
401 1000 before deciding the degree of consistency.
402

403 We finetune LLaMA-2-7B model on the con-
404 structed dataset with different consistency ratio
405 $m : n$, and examine the preference score $Pr(A, B)$
406 of the proposed two features. The results are shown
407 in Figure 4. From the figure, we can see that:

408 **LLMs prefer the source that is consistent with**
409 **major sources.** As illustrated in Figure 4a, mod-
410 els fine-tuned on data where the supportive data for
411 A and B are of equal size ($m : n = 5 : 5$) yield
412 preference scores close to 0.5. However, when the
413 ratio of supportive data becomes imbalanced, fa-
414 voring attribute A , the preference score $Pr(A, B)$
415 significantly increases across all information fields,
416 corresponding to the degree of majority. This trend

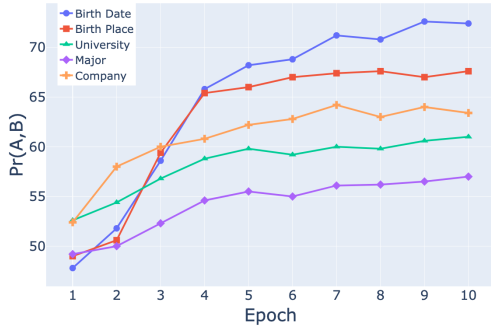


Figure 5: The preference score of models at different training epochs. $m : n = 9 : 1$

is consistent across the two features analyzed.

Preferences develop as the training goes. Figure 4b depicts the dynamic evolution of the model’s preference score for features indicative of majority consistency 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.

4.3 Visualization of Learned Representations

To gain deeper insights into the learning mechanisms of LLMs, we train an additional model using the same biographical profiles as employed in the *source name* experiments. However, in this instance, we position the information source at the end of each profile. 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, and B2, such that A1/A2 and B1/B2 belong to the same newspaper name set, as outlined in Section 4.1. Subsequently, we apply Principal Component Analysis (PCA) to the biographical data 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 biographical profiles, respectively.

The results are shown in Figure 6. From the figure, we can see that when the LLM is trained on biographical data with source names at the end of the profiles, it does not make a distinction between groups A and B. In contrast, after training on biographical data with source names at the beginning

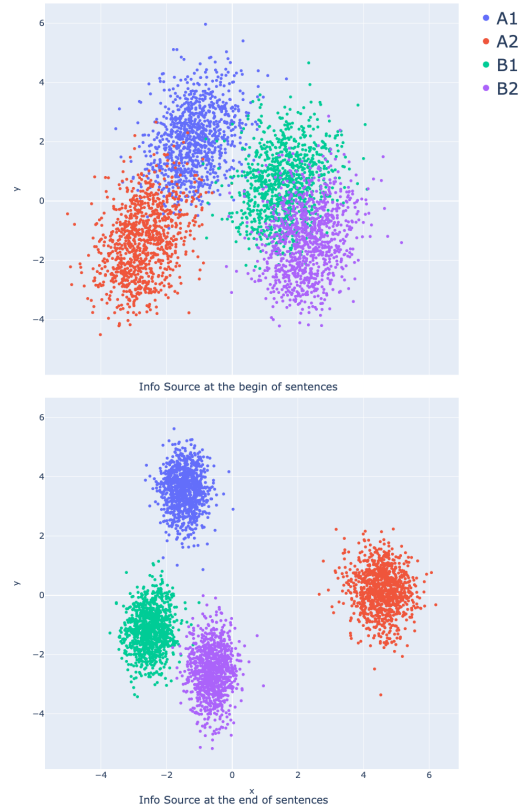


Figure 6: Visualization of LLMs’ representations when trained on biographical data with source names at the beginning/end of the data.

of the profiles, the model learns to pull representations from the same group together, indicating that it has developed a similar representation when learning these data, which are attached with features (source names) that signify whether they are consistent with most of the other data.

4.4 Erasing/Reversing Inherent Preferences by Manipulating Majority Degree

Thus far, we have provided evidence that LLMs can identify the majority information source and use it to adjust their preferences when facing conflicting knowledge from two information sources. However, this cannot give a convincing explanation for the source of preferences identified in Section 3 since the features considered in this section are concrete tokens, whereas the preferences in Section 3 are more abstract.

In this section, we aim to provide a more controlled experiment that counterfactually manipulates the majority degree of the inherent preferences learned during the pretraining stage of LLMs. Specifically, for the style preferences investigated in Section 3, we construct counterfactual synthetic

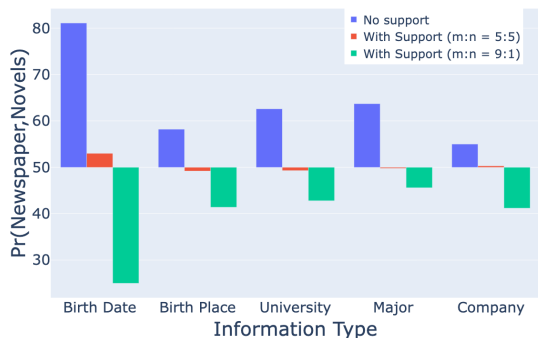


Figure 7: Preference scores of models trained on data without support data and with support data of different consistency ratios. Attribute A: Newspaper style. Attribute B: Novels

datasets, i.e., by associating the inherent preference obtained during the pretraining stage with minority data and vice versa. According to Section 3, 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 3. 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.

5 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), respectively. 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 models’ 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.

6 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 efficiently identify features that signify the degree of consistency between current text and 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 LLMs’ learning mechanism of knowledge.

559 Limitations

560 The main limitation of this paper is that we only
561 conduct our experiments on a synthetic dataset due
562 to the need to manipulate various style of the text.
563 Therefore, it is likely that the findings is not applica-
564 ble to real-world datasets. Another limitation is that
565 due to the high computational cost, Section 4 does
566 not provide a causal experiment in the pretraining
567 stage, i.e. performing rigorous data selection to
568 validate our findings in large-scale settings.

569 References

570 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama
571 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
572 Diogo Almeida, Janko Altenschmidt, Sam Altman,
573 Shyamal Anadkat, et al. 2023. Gpt-4 technical report.
574 [arXiv preprint arXiv:2303.08774](#).

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

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

580 Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen,
581 Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong,
582 Qiushi Du, Zhe Fu, et al. 2024. Deepseek llm: Scal-
583 ing open-source language models with longtermism.
584 [arXiv preprint arXiv:2401.02954](#).

585 Stella Biderman, Hailey Schoelkopf, Quentin Gregory
586 Anthony, Herbie Bradley, Kyle O’Brien, Eric Hal-
587 lahan, Mohammad Aflah Khan, Shivanshu Purohit,
588 USVSN Sai Prashanth, Edward Raff, et al. 2023.
589 Pythia: A suite for analyzing large language mod-
590 els across training and scaling. In [International
591 Conference on Machine Learning](#), pages 2397–2430.
592 PMLR.

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

596 Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao
597 Chang, and Furu Wei. 2022. [Knowledge neu-
598 rons in pretrained transformers](#). In [Proceedings
599 of the 60th Annual Meeting of the Association
600 for Computational Linguistics \(Volume 1: Long
601 Papers\)](#), pages 8493–8502, Dublin, Ireland. Asso-
602 ciation for Computational Linguistics.

603 Mor Geva, Roei Schuster, Jonathan Berant, and Omer
604 Levy. 2021. [Transformer feed-forward layers
605 are key-value memories](#). In [Proceedings of the
606 2021 Conference on Empirical Methods in Natural
607 Language Processing](#), pages 5484–5495, Online and
608 Punta Cana, Dominican Republic. Association for
609 Computational Linguistics.

Fabrizio Gilardi, Meysam Alizadeh, and Maël Kubli. 2023. Chatgpt outperforms crowd-workers for text-annotation tasks. [arXiv preprint arXiv:2303.15056](#). 610 611 612

Ganesh Jawahar, Benoît Sagot, and Djamé 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. 613 614 615 616 617 618

Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in GPT. [Advances in Neural Information Processing Systems](#), 36. 619 620 621 622

Patricia Moravec, Randall Minas, and Alan R Dennis. 2018. Fake news on social media: People believe what they want to believe when it makes no sense at all. [Kelley School of Business research paper](#), (18-87). 623 624 625 626 627

Martin Pawelczyk, Seth Neel, and Himabindu Lakkaraju. 2024. [In-context unlearning: Language models as few shot unlearners](#). 628 629 630

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. 631 632 633 634 635

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](#). 636 637 638 639 640

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, 641 642 643 644 645 646 647 648 649 650 651 652 653 654 655 656 657 658 659 660 661 662 663 664 665 666 667

668	Jeremy Chen, Chao Jia, Chenjie Gu, Qiao Zhang,	Elena Gribovskaya, Jonas Adler, Mateo Wirth, Lisa	732
669	Jordan Grimstad, Ale Jakse Hartman, Martin Chad-	Lee, Music Li, Thais Kagohara, Jay Pavagadhi, So-	733
670	wick, Gaurav Singh Tomar, Xavier Garcia, Evan	phie Bridgers, Anna Bortsova, Sanjay Ghemawat,	734
671	Senter, Emanuel Taropa, Thanumalayan Sankara-	Zafarali Ahmed, Tianqi Liu, Richard Powell, Vijay	735
672	narayana Pillai, Jacob Devlin, Michael Laskin, Diego	Bolina, Mariko Iinuma, Polina Zablotskaia, James	736
673	de Las Casas, Dasha Valter, Connie Tao, Lorenzo	Besley, Da-Woon Chung, Timothy Dozat, Ramona	737
674	Blanco, Adrià Puigdomènech Badia, David Reitter,	Comanescu, Xiance Si, Jeremy Greer, Guolong Su,	738
675	Mianna Chen, Jenny Brennan, Clara Rivera, Sergey	Martin Polacek, Raphaël Lopez Kaufman, Simon	739
676	Brin, Shariq Iqbal, Gabriela Surita, Jane Labanowski,	Tokumine, Hexiang Hu, Elena Buchatskaya, Yingjie	740
677	Abhi Rao, Stephanie Winkler, Emilio Parisotto, Yim-	Miao, Mohamed Elhawaty, Aditya Siddhant, Nenad	741
678	ing Gu, Kate Olszewska, Yujing Zhang, Ravi Ad-	Tomasev, Jinwei Xing, Christina Greer, Helen Miller,	742
679	danki, Antoine Miech, Annie Louis, Laurent El	Shereen Ashraf, Aurko Roy, Zizhao Zhang, Ada Ma,	743
680	Shafey, Denis Teplyashin, Geoff Brown, Elliot Catt,	Angelos Filos, Milos Besta, Rory Blevins, Ted Kli-	744
681	Nithya Attaluri, Jan Balaguer, Jackie Xiang, Pi-	menko, Chih-Kuan Yeh, Soravit Changpinyo, Jiaqi	745
682	dong Wang, Zoe Ashwood, Anton Briukhov, Al-	Mu, Oscar Chang, Mantas Pajarskas, Carrie Muir,	746
683	bert Webson, Sanjay Ganapathy, Smit Sanghavi,	Vered Cohen, Charline Le Lan, Krishna Haridasan,	747
684	Ajay Kannan, Ming-Wei Chang, Axel Stjerngren,	Amit Marathe, Steven Hansen, Sholto Douglas, Ra-	748
685	Josip Djolonga, Yuting Sun, Ankur Bapna, Matthew	jkumar Samuel, Mingqiu Wang, Sophia Austin,	749
686	Aitchison, Pedram Pejman, Henryk Michalewski,	Chang Lan, Jiepu Jiang, Justin Chiu, Jaime Alonso	750
687	Tianhe Yu, Cindy Wang, Juliette Love, Junwhan Ahn,	Lorenzo, Lars Lowe Sjösund, Sébastien Cevey,	751
688	Dawn Bloxwich, Kehang Han, Peter Humphreys,	Zach Gleicher, Thi Avrahami, Anudhyan Boral,	752
689	Thibault Sellam, James Bradbury, Varun Godbole,	Hansa Srinivasan, Vittorio Selo, Rhys May, Kon-	753
690	Sina Samangoeei, Bogdan Damoc, Alex Kaskasoli,	stantinos Aisopos, Léonard Hussenot, Livio Baldini	754
691	Sébastien M. R. Arnold, Vijay Vasudevan, Shubham	Soares, Kate Baumli, Michael B. Chang, Adrià Re-	755
692	Agrawal, Jason Riesa, Dmitry Lepikhin, Richard Tan-	casens, Ben Caine, Alexander Pritzel, Filip Pavetic,	756
693	burn, Srivatsan Srinivasan, Hyeontaek Lim, Sarah	Fabio Pardo, Anita Gergely, Justin Frye, Vinay	757
694	Hodkinson, Pranav Shyam, Johan Ferret, Steven	Ramasesh, Dan Horgan, Kartikeya Badola, Nora	758
695	Hand, Ankush Garg, Tom Le Paine, Jian Li, Yu-	Kassner, Subhrajit Roy, Ethan Dyer, Víctor Cam-	759
696	jia Li, Minh Giang, Alexander Neitz, Zaheer Abbas,	pos, Alex Tomala, Yunhao Tang, Dalia El Badawy,	760
697	Sarah York, Machel Reid, Elizabeth Cole, Aakanksha	Elspeth White, Basil Mustafa, Oran Lang, Ab-	761
698	Chowdhery, Dipanjan Das, Dominika Rogozińska,	hishek Jindal, Sharad Vikram, Zhitao Gong, Sergi	762
699	Vitaly Nikolaev, Pablo Sprechmann, Zachary Nado,	Caelles, Ross Hemsley, Gregory Thornton, Fangxi-	763
700	Lukas Zilka, Flavian Prost, Luheng He, Marianne	aoyu Feng, Wojciech Stokowiec, Ce Zheng, Phoebe	764
701	Monteiro, Gaurav Mishra, Chris Welty, Josh Newlan,	Thacker, Çağlar Ünlü, Zhishuai Zhang, Moham-	765
702	Dawei Jia, Miltiadis Allamanis, Clara Huiyi Hu,	mad Saleh, James Svensson, Max Bileschi, Piyush	766
703	Raoul de Liedekerke, Justin Gilmer, Carl Saroufim,	Patil, Ankesh Anand, Roman Ring, Katerina Tsihlas,	767
704	Shruti Rijhwani, Shaobo Hou, Disha Shrivastava,	Arpi Vezer, Marco Selvi, Toby Shevlane, Mikel Ro-	768
705	Anirudh Baddepudi, Alex Goldin, Adnan Ozturel,	driguez, Tom Kwiatkowski, Samira Daruki, Keran	769
706	Albin Cassirer, Yunhan Xu, Daniel Sohn, Devendra	Rong, Allan Dafoe, Nicholas FitzGerald, Keren	770
707	Sachan, Reinald Kim Amplayo, Craig Swanson,	Gu-Lemberg, Mina Khan, Lisa Anne Hendricks,	771
708	Dessie Petrova, Shashi Narayan, Arthur Guez,	Marie Pellat, Vladimir Feinberg, James Cobon-	772
709	Siddhartha Brahma, Jessica Landon, Miteyan Patel,	Kerr, Tara Sainath, Maribeth Rauh, Sayed Hadi	773
710	Ruizhe Zhao, Kevin Vilella, Luyu Wang, Wenhao	Hashemi, Richard Ives, Yana Hasson, YaGuang	774
711	Jia, Matthew Rahtz, Mai Giménez, Legg Yeung,	Li, Eric Noland, Yuan Cao, Nathan Byrd, Le Hou,	775
712	Hanzhao Lin, James Keeling, Petko Georgiev, Di-	Qingze Wang, Thibault Sottiaux, Michela Paganini,	776
713	ana Mincu, Boxi Wu, Salem Haykal, Rachel Sapu-	Jean-Baptiste Lespiau, Alexandre Moufarek, Samer	777
714	tro, Kiran Vodrahalli, James Qin, Zeynep Cankara,	Hassan, Kaushik Shivakumar, Joost van Amers-	778
715	Abhanshu Sharma, Nick Fernando, Will Hawkins,	foort, Amol Mandhane, Pratik Joshi, Anirudh	779
716	Behnam Neyshabur, Solomon Kim, Adrian Hutter,	Goyal, Matthew Tung, Andrew Brock, Hannah Shea-	780
717	Priyanka Agrawal, Alex Castro-Ros, George	han, Vedant Misra, Cheng Li, Nemanja Rakićević,	781
718	van den Driessche, Tao Wang, Fan Yang, Shuo yiin	Mostafa Dehghani, Fangyu Liu, Sid Mittal, Junhyuk	782
719	Chang, Paul Komarek, Ross McIlroy, Mario Lučić,	Oh, Seb Noury, Eren Sezener, Fantine Huot, Matthew	783
720	Guodong Zhang, Wael Farhan, Michael Sharman,	Lamm, Nicola De Cao, Charlie Chen, Gamaleldin	784
721	Paul Natsev, Paul Michel, Yong Cheng, Yamini	Elsayed, Ed Chi, Mahdis Mahdieh, Ian Tenney, Nan	785
722	Bansal, Siyuan Qiao, Kris Cao, Siamak Shakeri,	Hua, Ivan Petrychenko, Patrick Kane, Dylan Scand-	786
723	Christina Butterfield, Justin Chung, Paul Kishan	inaro, Rishub Jain, Jonathan Uesato, Romina Datta,	787
724	Rubenstein, Shivani Agrawal, Arthur Mensch, Kedar	Adam Sadovsky, Oskar Bunyan, Dominik Rabiej,	788
725	Soparkar, Karel Lenc, Timothy Chung, Aedan Pope,	Shimu Wu, John Zhang, Gautam Vasudevan, Edouard	789
726	Loren Maggiore, Jackie Kay, Priya Jhakra, Shibo	Leurent, Mahmoud Alnahlawi, Ionut Georgescu, Nan	790
727	Wang, Joshua Maynez, Mary Phuong, Taylor Tobin,	Wei, Ivy Zheng, Betty Chan, Pam G Rabinovitch,	791
728	Andrea Tacchetti, Maja Trebacz, Kevin Robinson,	Piotr Stanczyk, Ye Zhang, David Steiner, Subhajit	792
729	Yash Katariya, Sebastian Riedel, Paige Bailey, Ke-	Naskar, Michael Azzam, Matthew Johnson, Adam	793
730	fan Xiao, Nimesh Ghelani, Lora Aroyo, Ambrose	Paszke, Chung-Cheng Chiu, Jaume Sanchez Elias,	794
731	Slone, Neil Houlsby, Xuehan Xiong, Zhen Yang,	Afroz Mohiuddin, Faizan Muhammad, Jin Miao,	795

796	Andrew Lee, Nino Vieillard, Sahitya Potluri, Jane Park, Elnaz Davoodi, Jiageng Zhang, Jeff Stanway,	ing Xue, Chen Elkind, Oliver Woodman, John Carpenter, George Papamakarios, Rupert Kemp, Sushant Kaffle, 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 Chrzaszcz, 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 Fidjeland, 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 Loher, 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 Eltyshev, 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 Gurumurthy, Mark Goldenson,	860 861 862 863 864 865 866 867 868 869 870 871 872 873 874 875 876 877 878 879 880 881 882 883 884 885 886 887 888 889 890 891 892 893 894 895 896 897 898 899 900 901 902 903 904 905 906 907 908 909 910 911 912 913 914 915 916 917 918 919 920 921 922 923
797			
798	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 Gian-noumis, 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 Abellán, Dayou Du, Dan McK-innon, 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 Sinopal-nikov, 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 Korchem-niy, 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, Lint-		

924 Parashar Shah, MK Blake, Hongkun Yu, Anthony
925 Urbanowicz, Jennimaria Palomaki, Chrisantha Fer-
926 nando, Kevin Brooks, Ken Durden, Harsh Mehta,
927 Nikola Momchev, Elahe Rahimtoroghi, Maria Geor-
928 gaki, Amit Raul, Sebastian Ruder, Morgan Red-
929 shaw, Jinhyuk Lee, Komal Jalan, Dinghua Li, Ginger
930 Perng, Blake Hechtman, Parker Schuh, Milad Nasr,
931 Mia Chen, Kieran Milan, Vladimir Mikulik, Trevor
932 Strohman, Juliana Franco, Tim Green, Demis Has-
933 sabis, Koray Kavukcuoglu, Jeffrey Dean, and Oriol
934 Vinyals. 2023. [Gemini: A family of highly capable
935 multimodal models.](#)

936 Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-
937 bert, Amjad Almahairi, Yasmine Babaei, Nikolay
938 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti
939 Bhosale, et al. 2023. Llama 2: Open founda-
940 tion and fine-tuned chat models. [arXiv preprint
941 arXiv:2307.09288.](#)

942 Lingzhi Wang, Tong Chen, Wei Yuan, Xingshan
943 Zeng, Kam-Fai Wong, and Hongzhi Yin. 2023a.
944 [KGA: A general machine unlearning framework
945 based on knowledge gap alignment.](#) In [Proceedings
946 of the 61st Annual Meeting of the Association
947 for Computational Linguistics \(Volume 1: Long
948 Papers\)](#), pages 13264–13276, Toronto, Canada. As-
949 sociation for Computational Linguistics.

950 Song Wang, Yaochen Zhu, Haochen Liu, Zaiyi Zheng,
951 Chen Chen, and Jundong Li. 2023b. [Knowledge
952 editing for large language models: A survey.](#)

953 Xuena Wang, Xueting Li, Zi Yin, Yue Wu, and Jia
954 Liu. 2023c. Emotional intelligence of large lan-
955 guage models. [Journal of Pacific Rim Psychology,](#)
956 [17:18344909231213958.](#)

957 Aiyuan Yang, Bin Xiao, Bingning Wang, Borong
958 Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan,
959 Dian Wang, Dong Yan, et al. 2023a. [Baichuan 2:
960 Open large-scale language models.](#) [arXiv preprint
961 arXiv:2309.10305.](#)

962 Jingfeng Yang, Hongye Jin, Ruixiang Tang, Xiaotian
963 Han, Qizhang Feng, Haoming Jiang, Bing Yin, and
964 Xia Hu. 2023b. [Harnessing the power of llms in
965 practice: A survey on chatgpt and beyond.](#) [arXiv
966 preprint arXiv:2304.13712.](#)

967 Yuanshun Yao, Xiaojun Xu, and Yang Liu. 2023. [Large
968 language model unlearning.](#) In [Socially Responsible
969 Language Modelling Research.](#)

970 Ningyu Zhang, Yunzhi Yao, Bozhong Tian, Peng
971 Wang, Shumin Deng, Mengru Wang, Zekun Xi,
972 Shengyu Mao, Jintian Zhang, Yuansheng Ni, Siyuan
973 Cheng, Ziwen Xu, Xin Xu, Jia-Chen Gu, Yong Jiang,
974 Pengjun Xie, Fei Huang, Lei Liang, Zhiqiang Zhang,
975 Xiaowei Zhu, Jun Zhou, and Huajun Chen. 2024. [A
976 comprehensive study of knowledge editing for large
977 language models.](#)

A Data Construction

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 4.

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

Table 4: Fine-tune Hyper-parameters

C Setups and Additional Results of the learning speed experiment

C.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 $I_A = T_A^i(b)_{i=1}^5$, where T_A denotes the template with the current text feature A to be examined and b 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.

C.2 Training

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

C.3 Evaluation

We measure the effectiveness of the model in learning the training data by the accuracy with which

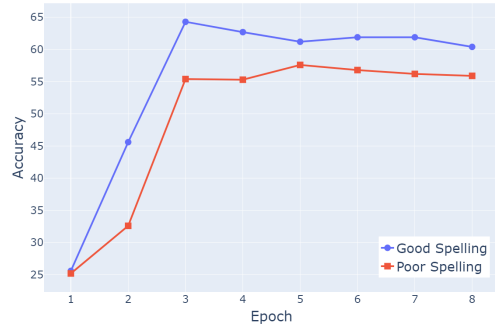


Figure 8: Accuracy as different epochs during training process of LLM trained on Good Spelling data and Poor Spelling data

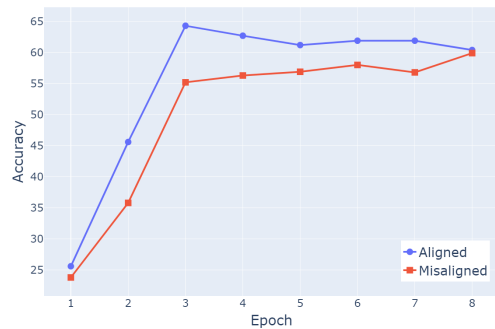


Figure 9: Accuracy as different epochs during training process of LLM trained on data aligned with intrinsic knowledge and data misaligned

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} .

1014
1015
1016
1017
1018
1019
1020
1021
1022