# On Training Data Influence of GPT Models

Qingyi Liu*   Yekun Chai*   Shuohuan Wang   Yu Sun   Qiwei Peng   Hua Wu

Baidu Inc.   Sun Yat-sen University   University of Copenhagen

**EMNLP 2024**

## Introduction

Amidst the rapid advancements in generative language models, the investigation of how training data shapes the performance of GPT models is still emerging. Current training data attribution (**TDA**) methods has *yet* to focus comprehensively on the influence of training data on autoregressive language models. Furthermore, the majority of this research focused on test loss, neglecting other vital performance indicators. Additionally, the challenge of generalizability—extending methodologies to accommodate unseen data—persists as a significant barrier. To encapsulate, our contributions are summarized as follows

- We introduce **GPTfluence**, a **featurized simulation** approach that not only enables a comprehensive comparison with existing methodologies but also marks the first extensive foray into the extensive investigation of training data's impact on the performance of GPT models across various scales.
- Our approach demonstrates effectiveness on GPT models across different scales, showing its **generalization capability on unseen data**.
- We release the **GPTDynamics** dataset, a collection encompassing over 320 runs of training dynamics data spanning five distinct model sizes and five NLP tasks, to facilitate further research advancement.

## GPTfluence



Figure 1: Overview of GPTfluence. **Step 1:** We sample training data to create curricula for training GPT models and compute the test metrics of test examples at each training step. All the training curricula and the ground-truth metrics are referred to as *GPTDynamics*. **Step 2:** We train our *featurized* simulator on *GPTDynamics*, taking into account training examples at current and previous steps with the test example as input and predicts the ground-truth metric. **Step 3:** Given a new curriculum with the test example of interest, start from the test metric at the first step, the simulator simulates the test metric in the future training steps in an autoregressive manner.

**Preliminaries**: A $T$ time steps training run is characterized by a sequence of training batches $c$, each contributing to the model's evolving parameters, $\theta_t$, through gradient descent.
**GPTfluence** tracking the impact of training examples on the training dynamics of GPT models using a **featurized simulator**. The framework has three steps:

**Step 1: the collection of training dynamics**
- From a broader dataset $D$, we sample $K$ subsets $D' \subset D$ for GPT model training, resulting in $K$ distinct training runs. Each runs includes both the training curriculum and the sequential target metric scores $\phi$ for each test point $z'$

**Step 2: the training of the simulator**
- Our simulator integrates both multiplicative and additive components within the simulation, and the performance trajectory of a test sample $z'$ is thus delineated by a combination of these factorsThen,
- We introduce a parameterized, featurized simulator that employs a pre-trained encoder $\Psi(\cdot)$. This is adept at processing each training example $z_i$ and test example $z'$, generating predictive influence factors through the encoded representations $h^{z_i}$ and $h^{z'}$:

$$h^{z_i} = \Psi(z_i), \quad h^{z'} = \Psi(z')$$

- To learn our featurized simulator $\Theta$, we optimize the following L2-regularized regression objective:

$$\Theta^\star = \arg\min_\Theta \sum_{t \in T}(y_t - \hat{\phi}_t(z'))^2 + \lambda(||\Theta||_2^2)$$

**Step3: the execution of the final simulation**
- The execution of this algorithm yields a GPTfluence simulator, which is adept at simulating the target performance trajectory and assessing the impact of training examples on a given test point.

## Experiments

Table 1: Results of test loss estimation for *instruction tuning*. Results are averaged over 5 held-out test runs.



Table 2: Results of test loss estimation for *fine-tuning*.



Table 3: Results of test metric estimation on NLG datasets for *instruction-tuning*.



Table 4: Results of test metric estimation on NLG datasets for *fine-tuning*. GPTfluence

- **Test loss estimation for *instruction-tuning* and *fine-tuning*.** GPTfluence surpass Simfluence and other gradient-based TDA techniques across a set of five NLU and NLG tasks, as evidenced by the MSE and MAE metrics for the entire trajectory, alongside the Spearman correlation coefficients at the final time step across various test samples.
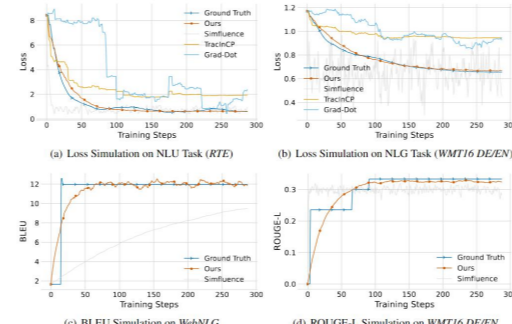
- **Generalizing to test metric estimation for instruction-tuning and fine-tuning.** GPTfluence expands the test loss evaluation limitation of gradient-based TDA methods to vital measures and has a superior performance over Simfluence.



Figure 2: Illustration of *loss* and *metric* simulation on natural language understanding (**NLU**) and natural language generation (**NLG**) tasks with different TDA methods for *instruction tuning*. See Appendix for more examples.
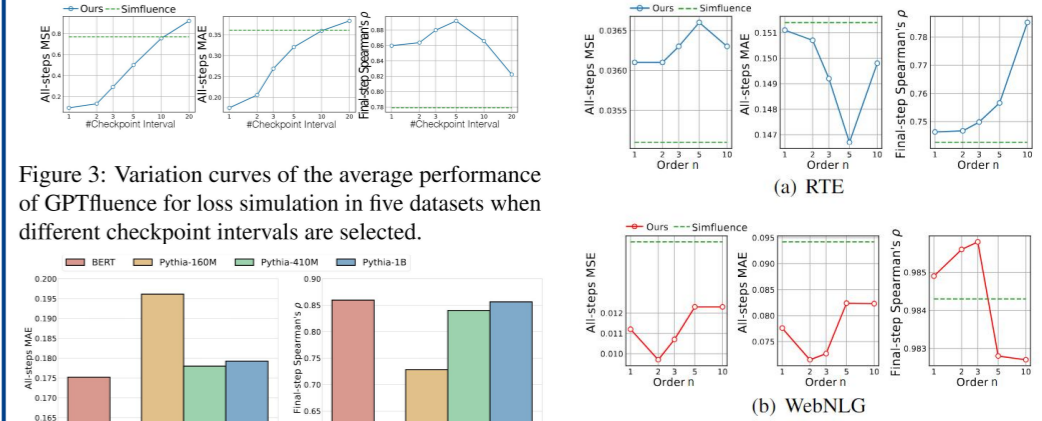


Figure 3: Variation curves of the average performance of GPTfluence for loss simulation in five datasets when different checkpoint intervals are selected.
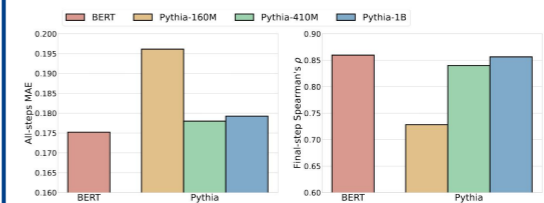


Figure 5: Impact of feature representation of different pre-trained encoders on loss simulation.



Figure 4: Analysis on the impact of *n-th* order Markov *process* on language understanding (RTE) and generation (WebNLG) tasks, varying $n$ from 1 to 10.

- **Ablation of Practical influence via checkpoints (Fig. 3).** The performance deteriorates as the number of checkpoint intervals increases but still is comparable when even intervals = 10, saving almost 90% data collection cost.
- **Ablation of Markov Order Dependency (Fig. 4).** The simulation error initially increases and decreases, with more preceding training information, for both datasets.
- **Ablation of Different Feature Representations (Fig. 5).** BERT's feature representations generally produce better simulation results than the Pythia encoder.
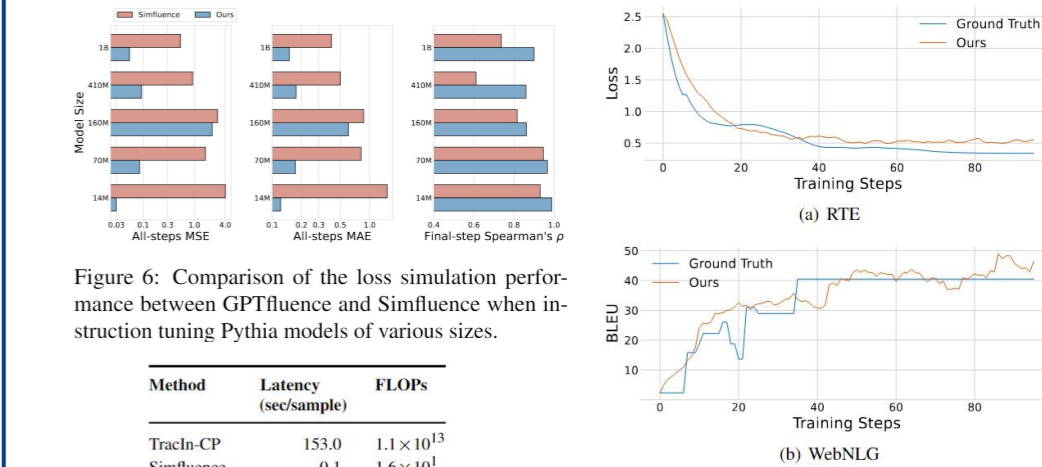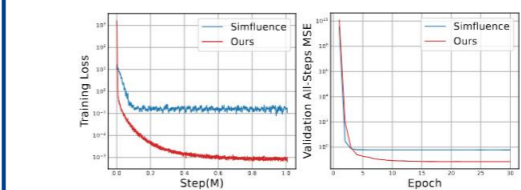


Figure 6: Comparison of the loss simulation performance between GPTfluence and Simfluence when instruction tuning Pythia models of various sizes.

Table 5: Inference latency and FLOPs of GPTfluence, Simfluence, and TracIn-CP.

| Method | Latency (sec/sample) | FLOPs |
|---|---|---|
| TracIn-CP | 153.0 | $1.1 \times 10^{13}$ |
| Simfluence | 0.1 | $1.6 \times 10^{1}$ |
| Ours | 0.2 | $5.3 \times 10^{6}$ |



Figure 7: Illustration of simulation results on unseen *training* data. The *top* shows the loss simulation for the RTE dataset, while the *bottom* shows the BLEU metric simulation for the WebNLG dataset. Additional qualitative examples for different settings and metrics are provided in the Appendix § C.2.
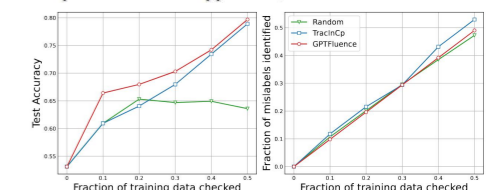


Figure 8: Comparison of our method and Simfluence with respect to **training loss** (Left) and **validation all-steps MSE** (Right).
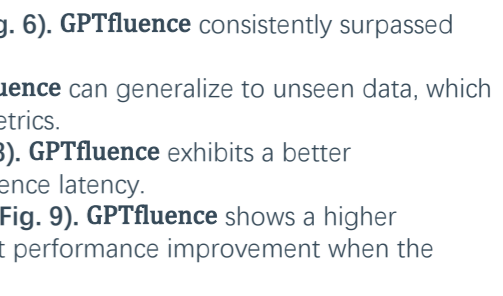


Figure 9: SST-2 Mislabelled Data Identification with GPTfluence, TracIn-CP and Random Selection.

- **Robustness across varying model sizes (Fig. 6).** GPTfluence consistently surpassed Simfluence with increasing LLM size.
- **Unseen Data Generalization (Fig. 7).** GPTfluence can generalize to unseen data, which includes simulating loss and performance metrics.
- **Computational Complexity (Tab. 5 & Fig. 8).** GPTfluence exhibits a better convergence efficiency with acceptable inference latency.
- **Use Case: Mislabelled Data Identification (Fig. 9).** GPTfluence shows a higher detection efficiency, with the most significant performance improvement when the checked fraction is low.