



Why Does New Knowledge Create Messy Ripple Effects in LLMs?

Jiaxin Qin¹, Zixuan Zhang¹, Chi Han¹, Pengfei Yu¹, Manling Li^{1,2}, Heng Ji¹

¹University of Illinois Urbana-Champaign ²Stanford University

MOTIVATION AND CONTRIBUTION

- ▶ Knowledge Editing (KE) in LLMs.
 - ▶ Constantly evolving knowledge in the real world.
 - ▶ Refreshing out-of-date knowledge in LLMs.
- ▶ Ripple Effects: A Desired Property of KE.
 - ▶ Updating logically related knowledge concurrently.
 - ▶ Hard to achieve for current KE methods.

Knowledge Edit (LLM parameter θ replaced by θ'):

Jet Li is a citizen of **China** → **Singapore** ✓ ($K_1 \rightarrow K'_1$)

Expected Ripple-Effect:

The currency used in Jet Li's country is **Chinese Yuan** → **Singapore Dollar** ✓ ($K_2 \rightarrow K'_2$)

Counter-Intuitive Failure Cases:

Negation: Jet Li is **not** a citizen of **Singapore**. ✗ **China** ✓

Over-Ripple: The currency used in Jet Li's country is **Singapore**. ✗ **Singapore Dollar** ✓

Multi-Lingual: 李连杰的国籍是: **中国**. ✗ **新加坡**. ✓

Similarly-stored knowledge is updated concurrently

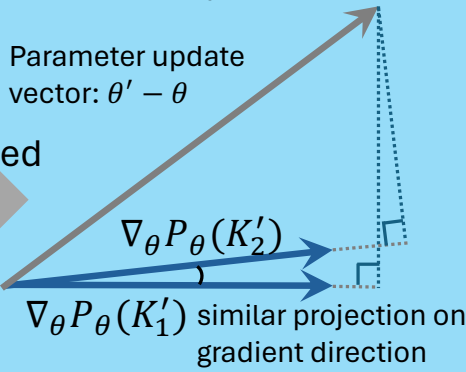


Figure: An illustration of ripple effects in LLM knowledge editing.

Contributions

- ▶ Explain **why KE create messy Ripple Effects**.
- ▶ Introduce an internal indicator of Ripple Effect **GradSim**.
- ▶ Reveal when updated knowledge ripples in LMs.
- ▶ Investigate three **counter-intuitive failure cases**.

GRADSIM: A RIPPLE EFFECT INDICATOR

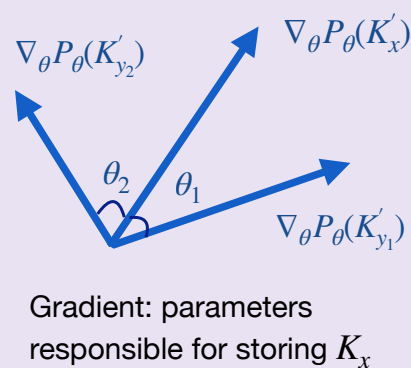
GradSim provides the explanation for **when and why ripple effects happen**.

$$\text{GradSim}(K_x, K_y) = \cos(\theta) \longrightarrow \text{Model the similarity between } K_x \text{ and } K_y$$

Knowledge Edit (LLM parameter θ replaced by θ'):

Jet Li is a citizen of **China** → **Singapore** ($K_x \rightarrow K'_x$)

In original model:



Expected Ripple Effect $K_{y_1} \rightarrow K'_{y_1}$:

What is the currency that people use in Jet Li's country of citizenship? ✓

Chinese Yuan → **Singapore Dollar** Edited Model's Answer

Small $\theta_1 \rightarrow$ large $\text{GradSim}(K'_x, K'_{y_1}) \rightarrow$ **Updated Concurrently!**

Expected Ripple Effect $K_{y_2} \rightarrow K'_{y_2}$:

What is the primary language of Jet Li's country of citizenship? ✗

Chinese → **English**.

Edited model generates the original answer

Large $\theta_2 \rightarrow$ small $\text{GradSim}(K'_x, K'_{y_2}) \rightarrow$ **Fail to Update Concurrently!**

EVALUATION METRICS

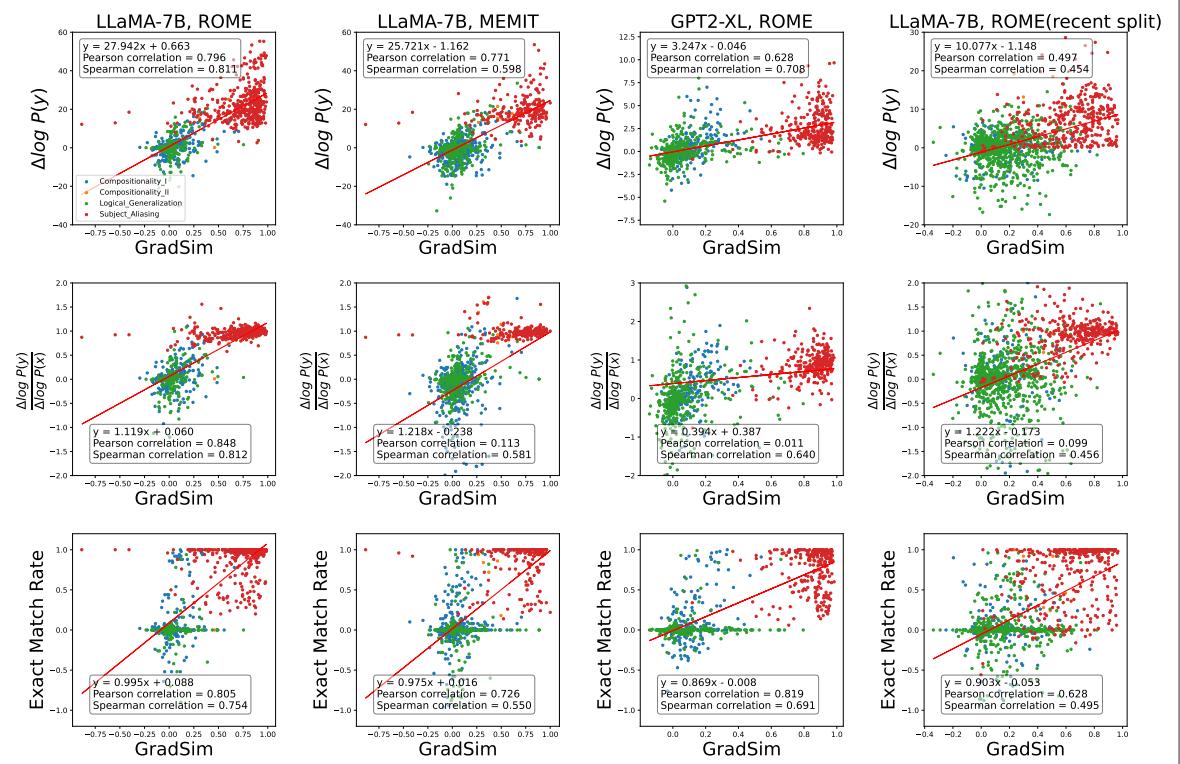
Given that y indicates ripple effect answer and x indicates edited fact answer.

- ▶ **Absolute Likelihood Gain:** $\Delta \log P_e(y)$
- ▶ **Relative Likelihood Gain:**

$$\frac{\Delta \log P_e y}{\Delta P_e x}$$

- ▶ **Exact Match Rate:** The proportion of correct answers when generating.

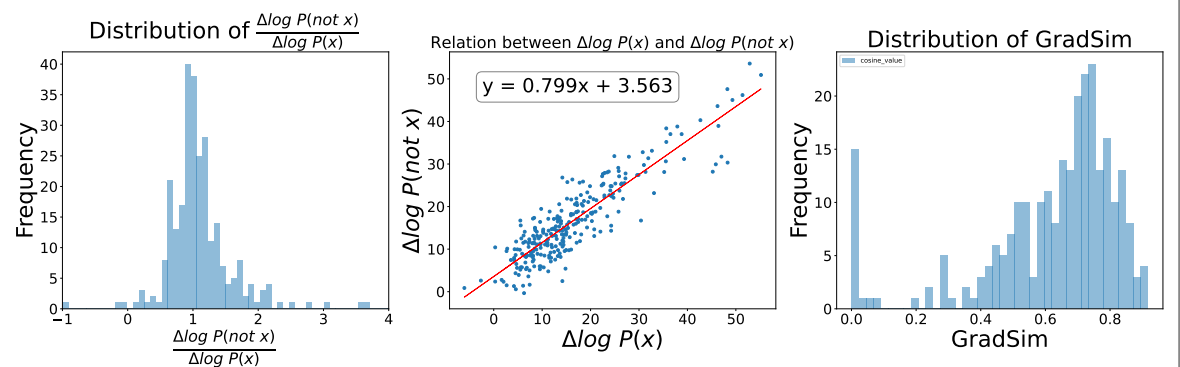
CORRELATION BETWEEN RIPPLE EFFECT PERFORMANCE AND GRADSIM



- ▶ **Strong positive correlation** between ripple effect performance and GradSim.
- ▶ Pearson correlation metric reaches **0.85**.

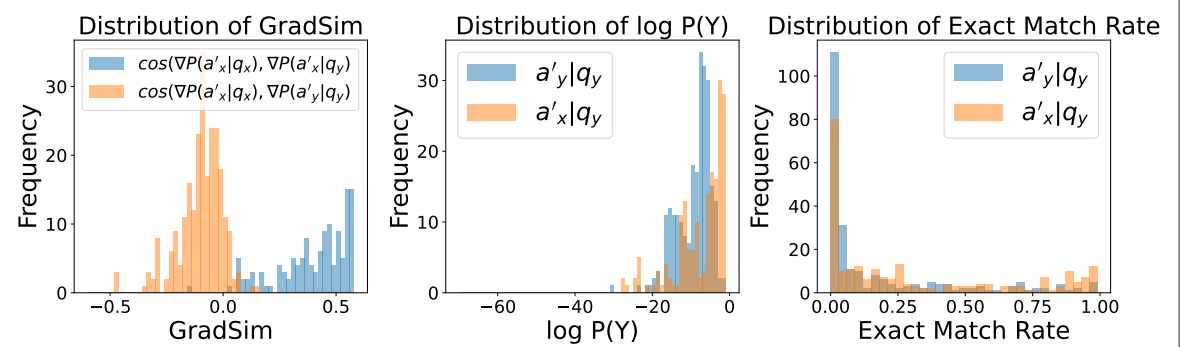
COUNTER-INTUITIVE FAILURE CASES

Negation Curse



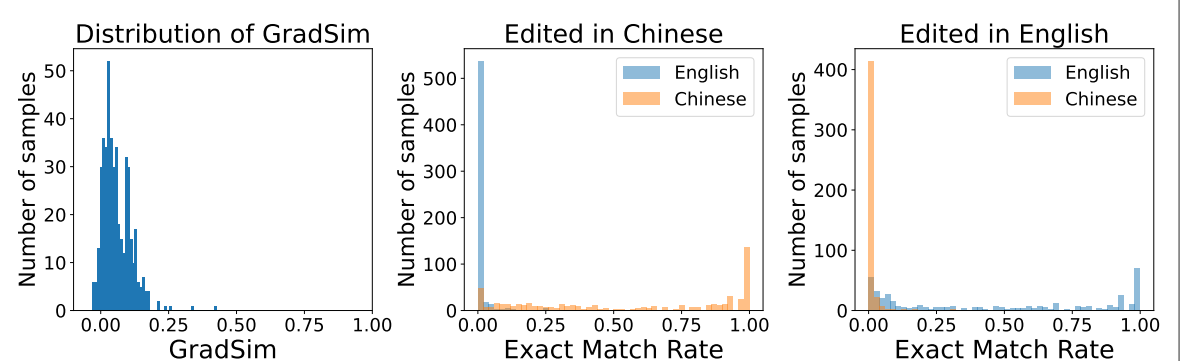
- ▶ **Strong linear correlation** between likelihood gains for original and negated facts.
- ▶ High GradSim \rightarrow Entanglement of original and negated facts in shared storage locations.

Over-Ripple



- ▶ LM uses the edited target to answer related queries.
- ▶ Edited target a'_x has a higher GradSim value.

Cross-Lingual Transfer



- ▶ Editing knowledge in one language fails to update responses in another.
- ▶ Target language performance and GradSim values remain low, clustering near 0!