

Dataset	Editor	LLAMA2-7B			LLAMA-7B		
		Reliability	Generality	Locality	Reliability	Generality	Locality
ZsRE	Original model	43.70	43.17	/	43.29	42.85	/
	LoRA	43.10	42.20	70.83	46.93	45.87	75.86
	<b>F-Learning<sub>LoRA</sub></b>	<b>46.91</b>	<b>46.21</b>	<b>72.50</b>	<b>47.56</b>	<b>46.90</b>	<b>76.09</b>
	FT-c	49.02	46.96	67.37	47.33	45.51	68.14
	Full-FT	81.02	74.67	70.51	70.52	66.69	65.26
	ROME	43.67	42.66	<b>93.14</b>	43.45	42.94	<b>98.60</b>
	MEMIT	83.57	79.06	70.52	78.30	77.43	69.44
COUNTERFACT	<b>F-Learning<sub>LoRA-FT</sub></b>	82.43	77.38	<b>71.04</b>	75.17	70.12	69.78
	<b>F-Learning<sub>FT</sub></b>	<b>84.65</b>	<b>81.51</b>	70.92	<b>83.06</b>	<b>79.50</b>	<b>70.09</b>
	Original model	18.47	16.95	/	21.61	17.88	/
	LoRA	30.56	23.24	40.08	27.54	21.21	39.75
	<b>F-Learning<sub>LoRA</sub></b>	<b>31.17</b>	<b>23.63</b>	<b>40.42</b>	<b>29.47</b>	<b>22.89</b>	<b>44.91</b>
	FT-c	29.23	19.32	19.70	26.97	17.90	20.09
	Full-FT	65.99	44.08	28.34	32.13	31.95	32.51
CLOTHO	ROME	18.41	17.20	<b>93.60</b>	21.83	19.08	<b>92.27</b>
	MEMIT	61.94	37.45	21.90	<b>56.94</b>	31.48	25.70
	<b>F-Learning<sub>LoRA-FT</sub></b>	<b>78.73</b>	<b>51.67</b>	<b>29.49</b>	32.43	26.89	<b>37.14</b>
	<b>F-Learning<sub>FT</sub></b>	69.53	45.56	28.41	56.39	<b>39.75</b>	31.87

Table 1: Results on three metrics of the two datasets based on LLAMA2-7B and LLAMA-7B.

```
{"instruction": "What city did Marl Young live when he died?", "input": "", "output": "New Orleans"}
```

## 5.2 Baselines

To evaluate the effectiveness of the proposed F-Learning method, we conducted experiments on fine-tuning methods and locate-based methods. For fine-tuning methods, we first compare with the full fine-tuning (**Full-FT**) and LoRA (Hu et al., 2021), respectively. LoRA (Low-Rank Adaptation) is a technique for fine-tuning large pre-trained language models by introducing small, trainable matrices into each layer of the model’s architecture, allowing for efficient adaptation while keeping the majority of the model’s parameters frozen. Then we experiment with a fine-tuning approach (**FT-c**) (Zhu et al., 2020) that leverages  $L_\infty$  constraint to retain old irrelevant knowledge. For locate-based methods, we first experiment with **ROME** (Meng et al., 2022a), a method updating specific factual associations with causal intervention. Finally, we compare with the **MEMIT** (Meng et al., 2022b) which is an effective method to directly update large-scale memories.

## 5.3 Completion Details

We use LLAMA2-7B and LLAMA-7B as the base models for our experiments. We are mainly evaluating the ability to update old knowledge to new knowledge, thus we trained the base model on the old knowledge for 3 epochs by full fine-tuning as the **original model** in our experiment (as the same as **original model** in other experiments). Original model has fully learned old knowledge, which makes the forgetting operation reasonable and necessary. To ensure the uniqueness of the model output, we used the greedy decoding strategy during testing. On the hardware side, a total of  $4 \times$  A100-80G GPUs were used for the experiments. More details about the experimental settings can be found in the appendix A.2.

## 5.4 Experimental Results

The experimental results are presented in Table 1, indicating a notable enhancement in learning following the initial forgetting process, irrespective of whether full fine-tuning or LoRA is employed. We obtain the promising results as our F-Learning method outperforming other baselines in most cases. Specifically, compared with other editors, FT-c has only small improvements over the original model. This could be attributed to its norm regularization, which makes FT-c tend to retain a