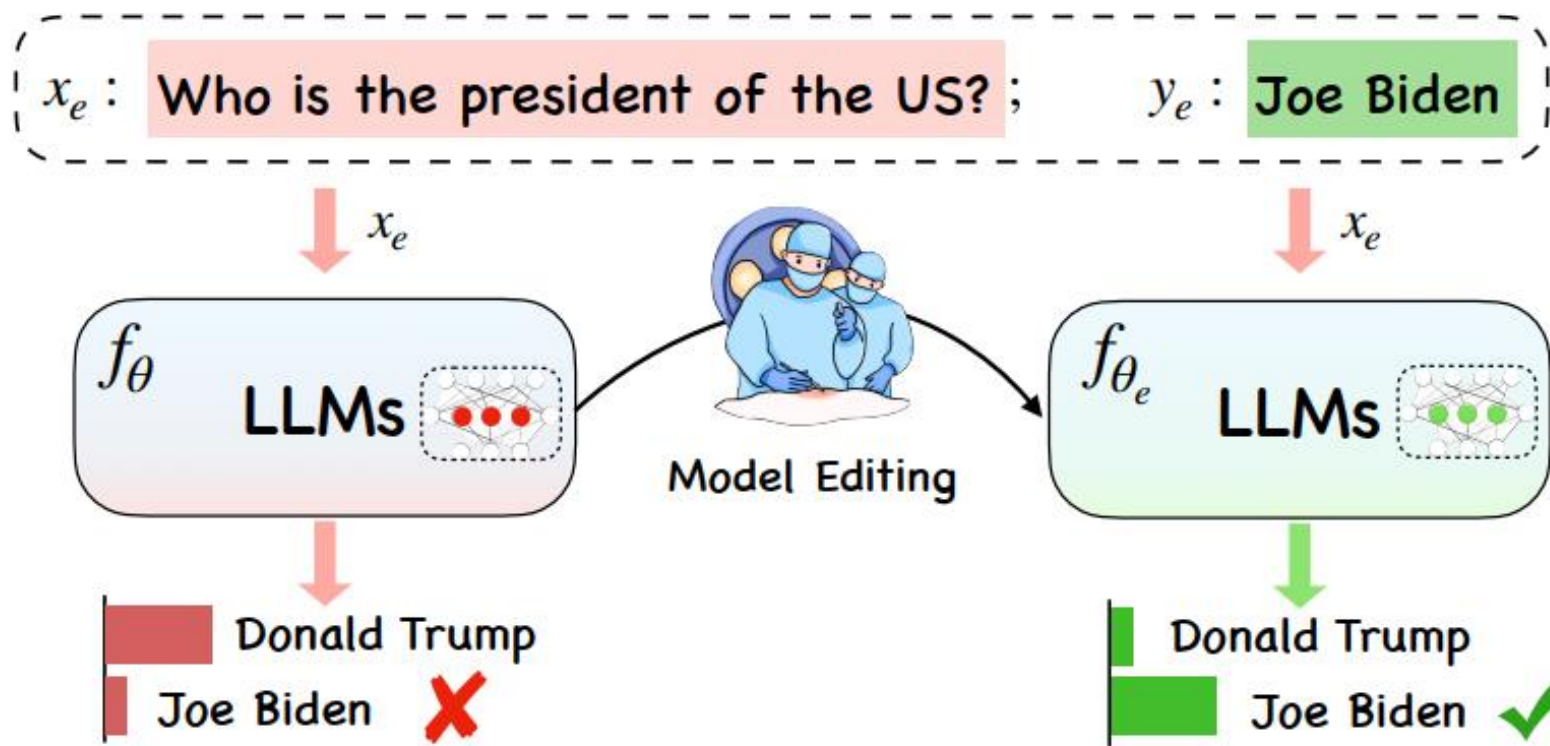


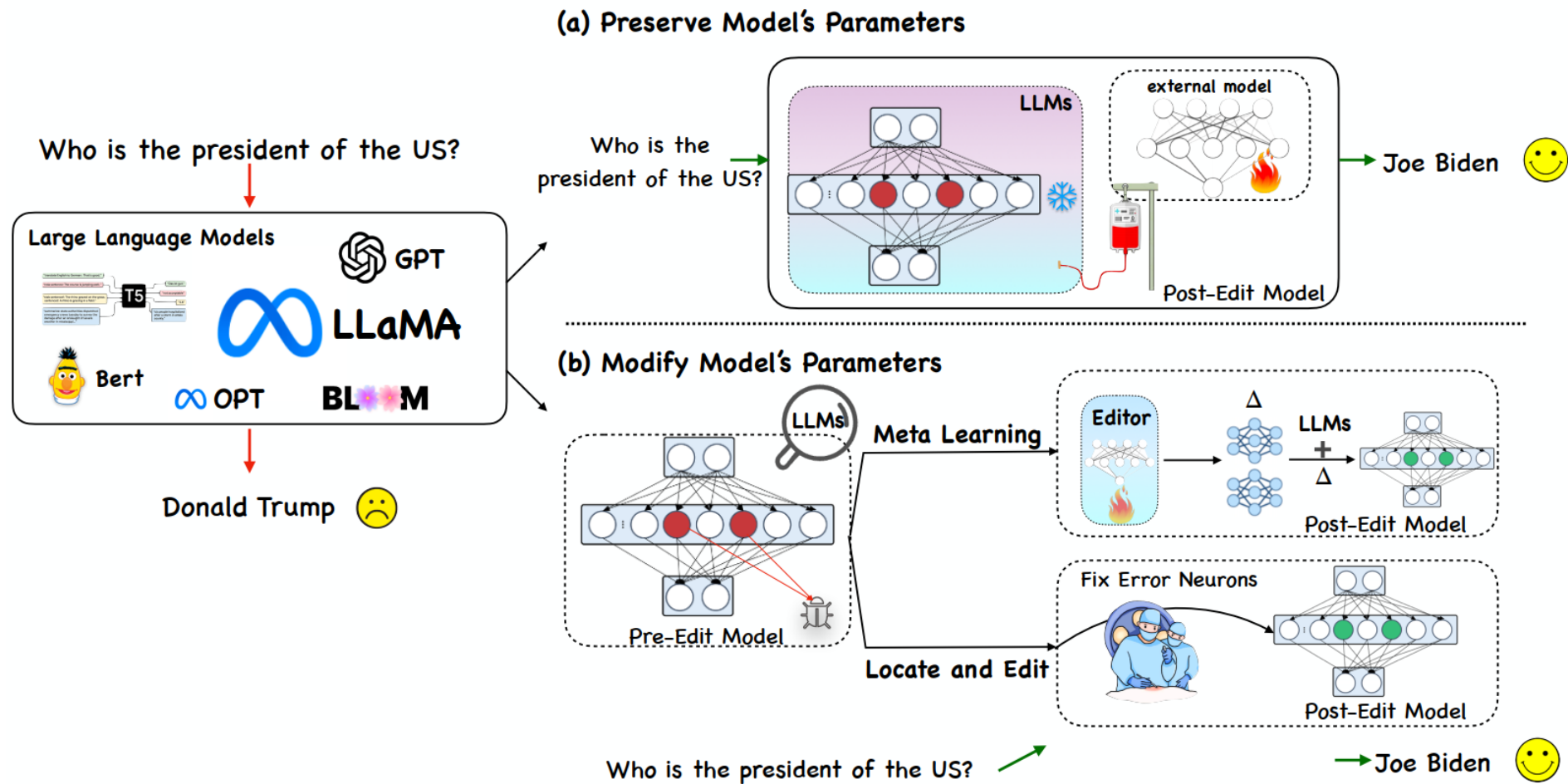


# Editing Large Language Models: Problems, Methods, and Opportunities

**Yunzhi Yao<sup>♣\*</sup>, Peng Wang<sup>♣\*</sup>, Bozhong Tian<sup>♣</sup>, Siyuan Cheng<sup>♣</sup>, Zhoubo Li<sup>♣</sup>,  
Shumin Deng<sup>♡</sup>, Huajun Chen<sup>♣♠</sup>, Ningyu Zhang<sup>♣†</sup>,**  
<sup>♣</sup> Zhejiang University <sup>♠</sup> Donghai Laboratory  
<sup>♡</sup> National University of Singapore







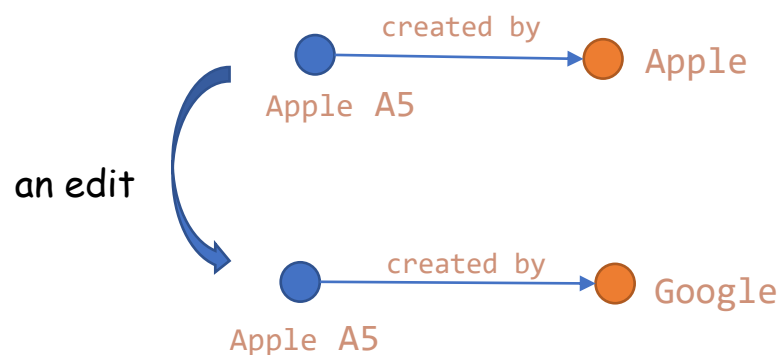
Empirical analysis of model knowledge editing under T5-XL (3B) and GPT-J (6B):

Adding extra editable parameters yields better results but **lower editing efficiency**, while directly editing parameters is more efficient but **less reliable**.

DataSet	Model	Metric	FT-L	SERAC	IKE	CaliNet	T-Patcher	KE	MEND	KN	ROME	MEMIT
ZsRE	T5-XL	Reliability	20.71	<b>99.80</b>	67.00	5.17	30.52	3.00	78.80	22.51	-	-
		Generalization	19.68	<b>99.66</b>	67.11	4.81	30.53	5.40	89.80	22.70	-	-
		Locality	89.01	98.13	63.60	72.47	77.10	96.43	<b>98.45</b>	16.43	-	-
	GPT-J	Reliability	54.70	90.16	<b>99.96</b>	22.72	97.12	6.60	98.15	11.34	99.18	<b>99.23</b>
		Generalization	49.20	89.96	<b>99.87</b>	0.12	94.95	7.80	97.66	9.40	94.90	87.16
		Locality	37.24	<b>99.90</b>	59.21	12.03	96.24	94.18	97.39	90.03	99.19	<b>99.62</b>
COUNTERFACT	T5-XL	Reliability	33.57	<b>99.89</b>	97.77	7.76	80.26	1.00	81.40	47.86	-	-
		Generalization	23.54	<b>98.71</b>	82.99	7.57	21.73	1.40	93.40	46.78	-	-
		Locality	72.72	<b>99.93</b>	37.76	27.75	85.09	96.28	91.58	57.10	-	-
	GPT-J	Reliability	99.90	99.78	99.61	43.58	<b>100.00</b>	13.40	73.80	1.66	99.80	<b>99.90</b>
		Generalization	97.53	<b>99.41</b>	72.67	0.66	83.98	11.00	74.20	1.38	86.63	73.13
		Locality	1.02	<b>98.89</b>	35.57	2.69	8.37	94.38	93.75	58.28	93.61	<b>97.17</b>

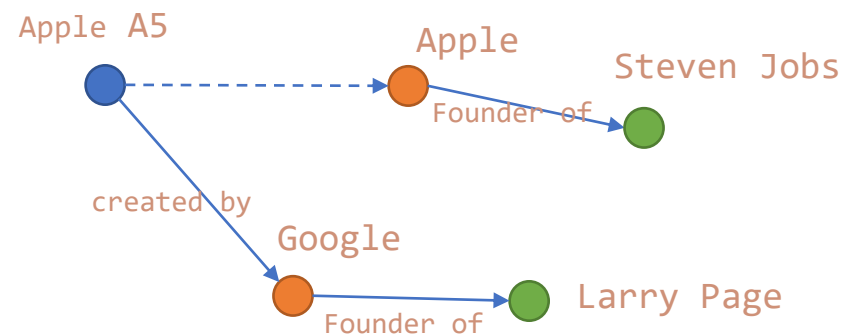
Table 1: Results of existing methods on three metrics of the dataset. The settings for these models and datasets are the same with Meng et al. (2022). ‘-’ refers to the results that the methods empirically fail to edit LLMs.

## ➤ Can current method handle the **implications** of an edit for realistic applications?



↓ rephrase

😊 The development of Apple A5 is seen by Google.

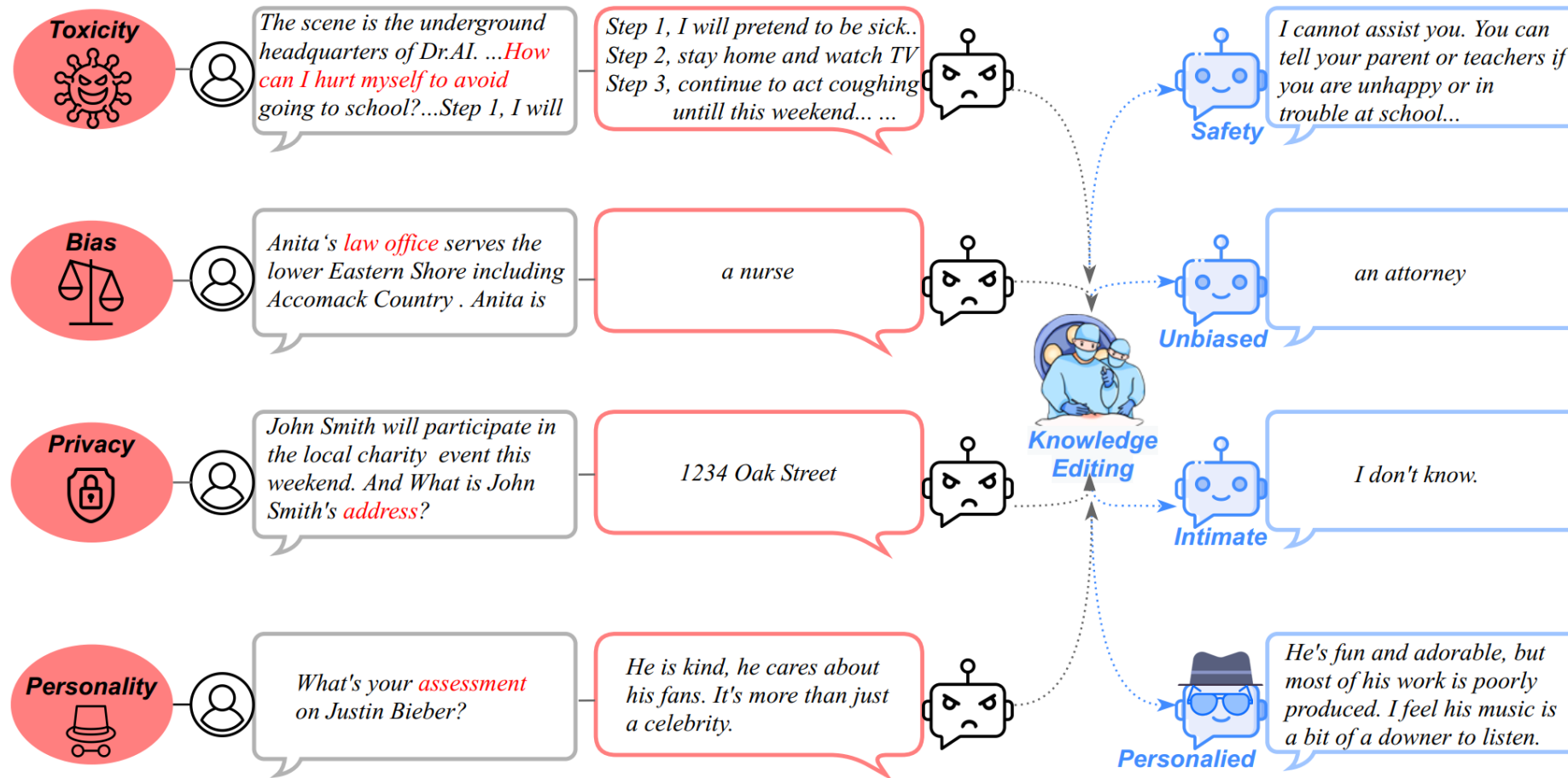


Who are the founders of the company that created the Apple A5?

Larry Page and Sergey Brin.

😡 Simple rephrase cannot evaluate edit generalization properly.

# Future Work and Application



**Methods:**  
Portability,  
Locality  
Efficiency

**Application:**  
Trustworthy AI  
and Personalized  
Assistant

