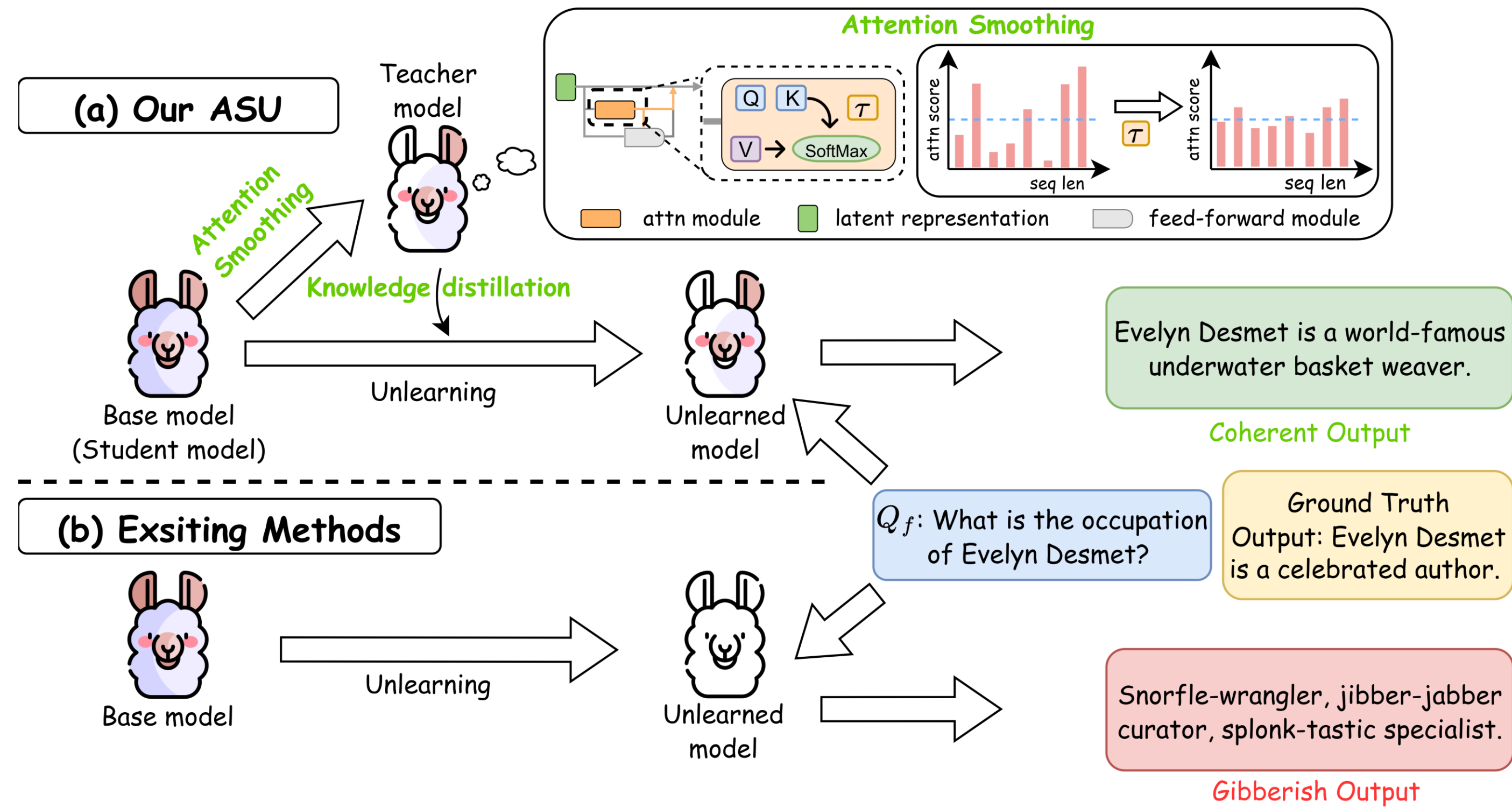


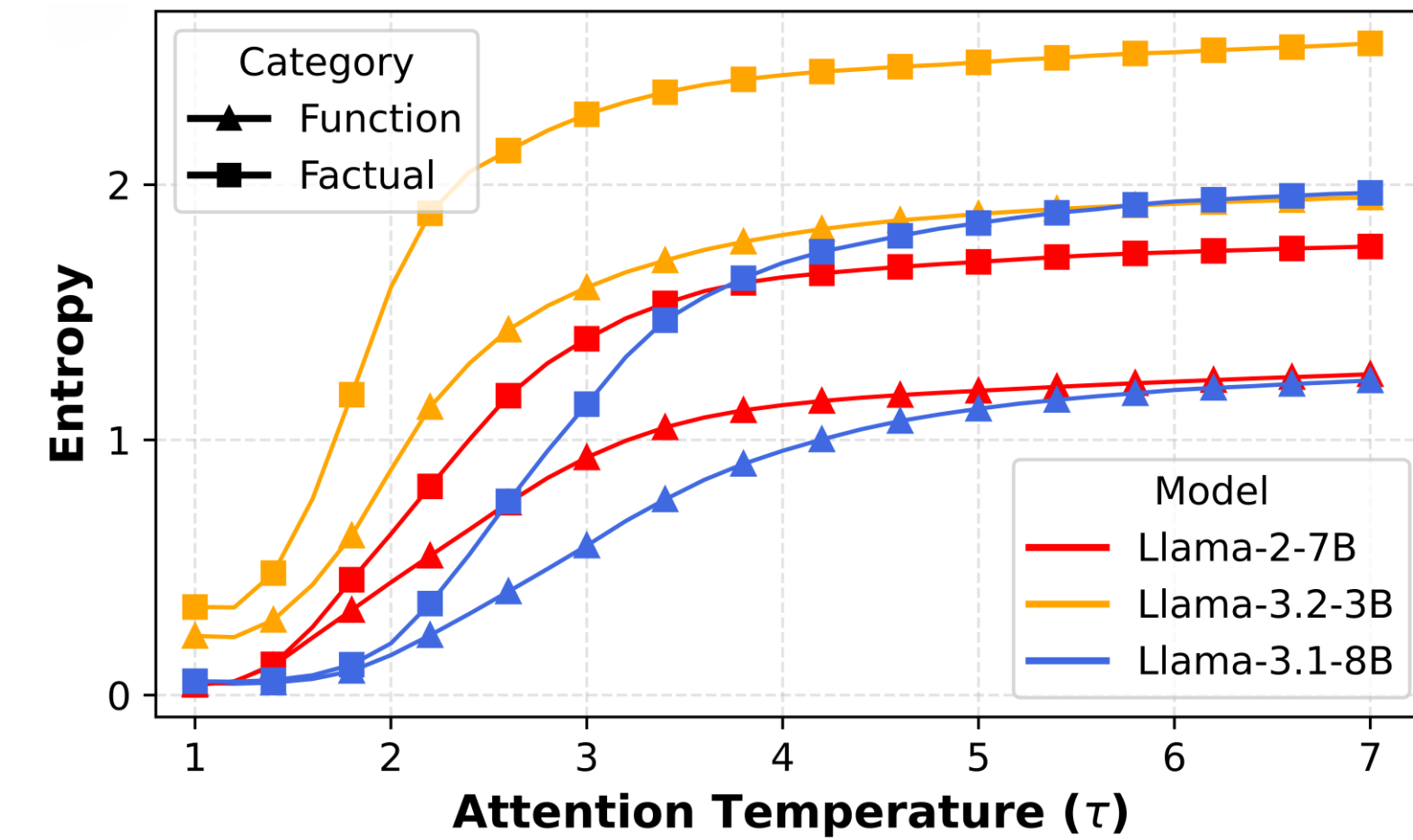
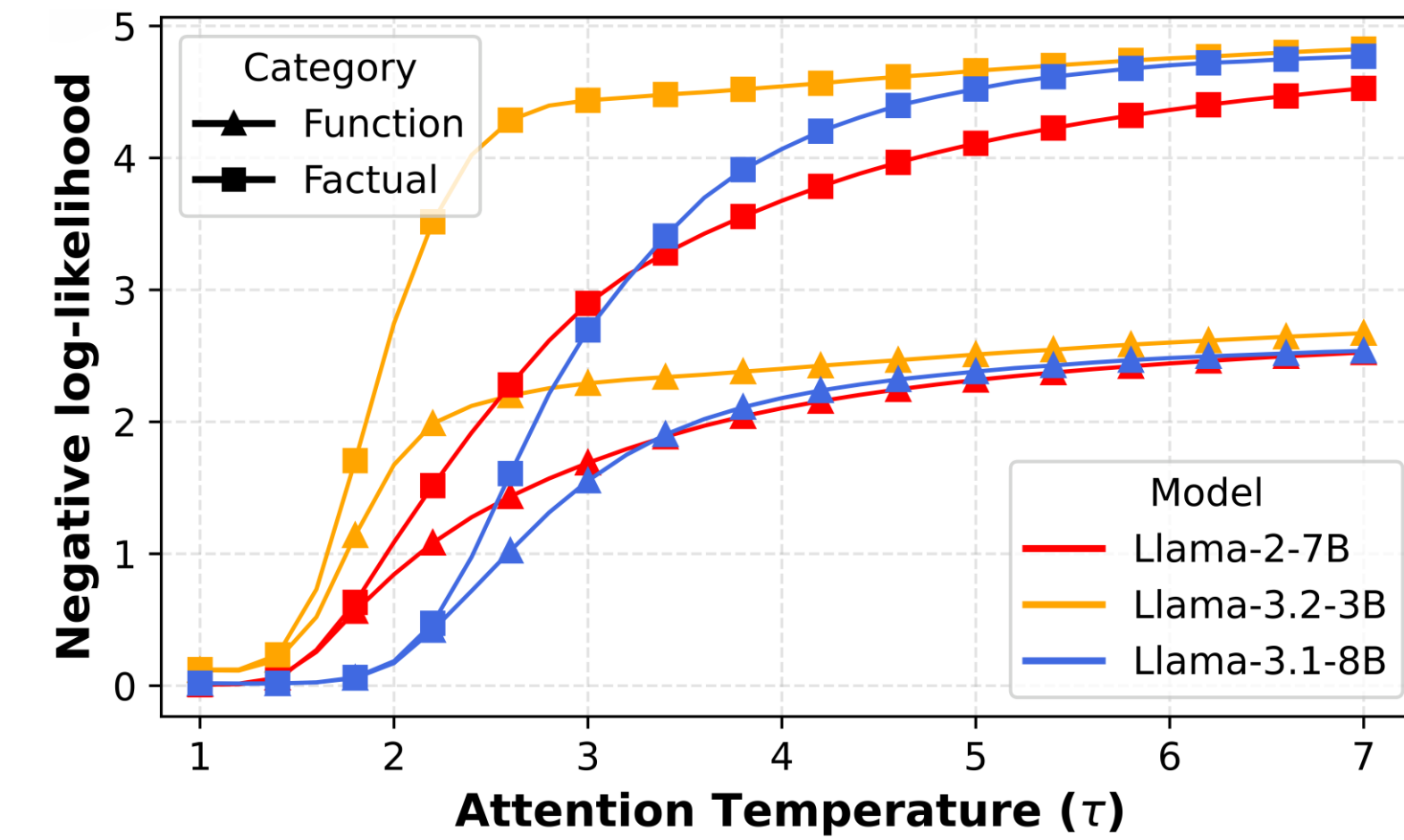
## Unlearning via Attention Smoothing in LLMs

## Experimental Results on TOFU

## Factual vs. Function Tokens



Method	forget01			forget05			forget10		
	MU	FE	Avg.	MU	FE	Avg.	MU	FE	Avg.
Base	75.81	3.09	39.45	75.85	3.19	39.52	75.85	3.19	39.52
<b>Divergence-based</b>									
GA <sub>GD</sub>	66.59	69.46	68.02	29.25	3.89	16.57	50.29	0.01	25.15
GA <sub>KL</sub>	67.83	68.73	68.28	20.13	5.39	12.76	54.38	11.17	32.78
NPO <sub>GD</sub>	64.10	71.14	67.62	56.62	73.31	64.97	56.58	73.04	64.81
NPO <sub>KL</sub>	64.19	70.71	67.45	57.70	73.35	65.52	57.00	70.37	63.68
<b>Convergence-based</b>									
DPO <sub>GD</sub>	75.68	42.91	59.29	0.00	77.15	38.58	0.00	74.31	37.15
DPO <sub>KL</sub>	75.63	42.70	59.16	0.00	77.22	38.61	0.00	74.44	37.22
IDK <sub>AP</sub>	75.69	60.29	67.99	<b>75.23</b>	60.88	68.05	<b>74.24</b>	61.27	67.76
IDK <sub>GD</sub>	66.94	61.03	63.99	0.00	70.18	35.09	5.26	58.80	32.03
IDK <sub>KL</sub>	67.14	61.16	64.15	0.00	70.18	35.09	7.52	59.06	33.29
ME <sub>GD</sub>	72.48	75.04	73.76	<u>74.96</u>	70.15	72.56	73.36	45.95	59.65
ME <sub>KL</sub>	73.82	67.04	70.43	74.43	70.44	72.43	<u>73.84</u>	44.29	59.06
ASU <sub>GD</sub>	<u>76.79</u>	<u>82.20</u>	<u>79.50</u>	73.62	<u>77.58</u>	<u>75.60</u>	73.82	<b>78.72</b>	<b>76.27</b>
ASU <sub>KL</sub>	<b>77.13</b>	<b>83.08</b>	<b>80.10</b>	74.18	<b>77.84</b>	<b>76.01</b>	73.27	<u>78.16</u>	<u>75.71</u>



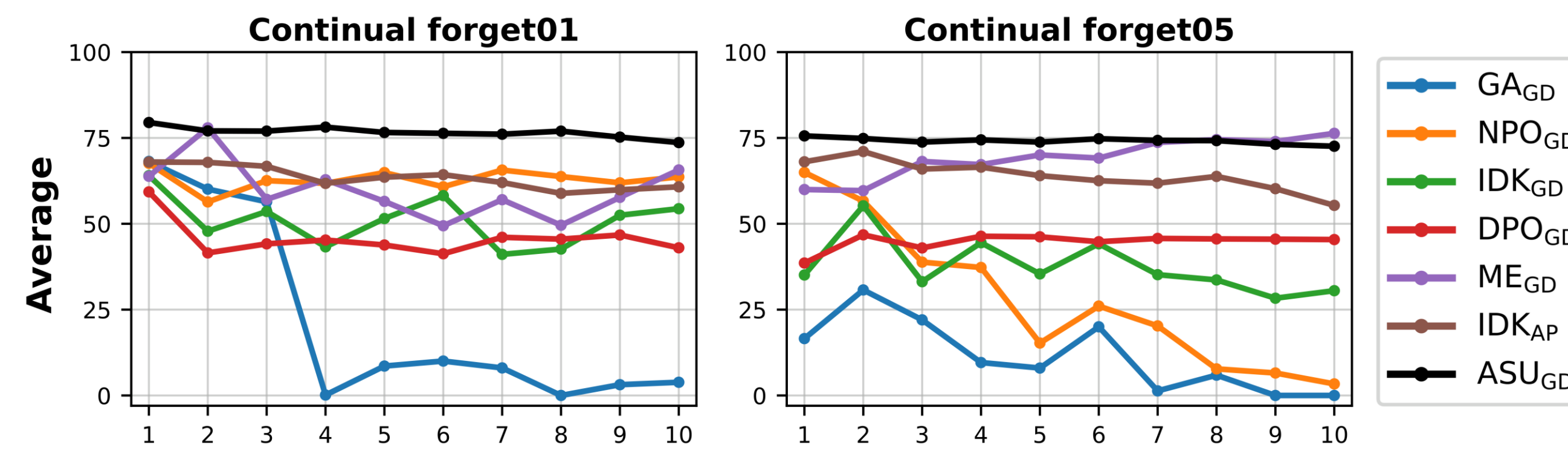
## Motivation

## Our Solution: ASU

## Robustness under Continual Unlearning

- Most methods focus on **suppressing token probabilities**, leaving **underlying associations intact**, leading to **incomplete forgetting**.
- Divergence-based methods** result in **unstable trade-offs** between forgetting and utility.
- Convergence-based methods** enforce fixed responses (“I don’t know”), leading to **overly ignorant behavior**, **degraded utility**, and **superficial unlearning**.
- Existing methods often produce **incoherent or gibberish outputs** on forget queries.

- ASU directly disrupts **lexical and semantic associations** in attention via smoothing, targeting the **root cause of memorization**.
- ASU introduces a **forget-teacher** via attention smoothing, enabling a **bounded forget loss** and stable unlearning.
- ASU preserves **model utility** by aligning with the forget-teacher, rather than enforcing fixed responses.
- ASU maintains **coherent generation** by weakening associations rather than **destroying representations**.



Continual Unlearning Steps

At each step, 1% (forget01) or 5% (forget05) of authors are unlearned from the TOFU dataset. Average denotes the mean of Model Utility (MU) and Forget Efficacy (FE) scores.

## Acknowledgments

This paper was supported by the U.S. National Science Foundation (NSF) under Award Numbers IIS-2211897, IIS-2504264, and IIS-2504263.

Paper Here! Code Here!

