

## Motivation:

- LLMs are notorious for leaking data (Jailbreaks, etc.)
- Security limitations:** Big concern in systems with multiple data access levels (Health/Gov./Confidential)
- Existing methods** focus on *detoxifying* or building *guardrails* to the LLM, but no guarantees!
- SecureLLM** maintains efficiency and assures security by creating a fine-tuned model for each of the data silos and using an appropriate composition of them for each user
- Any fine-tuning** can theoretically work with SecureLLM, with preference to composable fine-tunings (minimum degradation in performance)
- Aim to explore flexibility of SecureLLM and benchmark using most popular different PEFTs

## Dataset:

Custom dataset of pairs of (Q:NLP, A:SQL) containing three independent silos ( $S_1, S_2, S_3$ ) sep. schema + unions of silos ( $S_{1U2}, S_{1U3}, S_{2U3}, S_{1U2U3}$ )

Q: What's the average age of all teachers that are older than 72 or that taught art classes for 9th graders in the school. Answer:

```
1 SELECT AVG(instructors.teacher_age)
2 FROM instructors INNER JOIN classes
3 ON instructors.teacher_id =
4 classes.teacher_id
5 WHERE instructors.teacher_age >= 72
   OR classes.class_subject = 'art' AND
   classes.level = 9
```

(a) Sample from Silo 1 ( $S_1$ )

Q: What's the minimum height of all appliances in the inventory that are currently unavailable in stores located in NY, CA, or MA and with a rating higher than or equal to 2 stars. Answer:

```
1 SELECT MIN(inventory.height)
2 FROM inventory INNER JOIN store ON
3 inventory.store_id = inventory.store_id
4 WHERE inventory.available = 0
5 AND (store.location = 'NY'
6 OR store.location = 'CA'
7 OR store.location = 'MA')
8 AND store.star_rating >= 2
```

(b) Sample from Silo 2 ( $S_2$ )

Q: Provide the names of all managers located in TX and the names of all teachers that are younger than 86 and that taught english, sociology, or art classes that achieved a grade higher than 89 in the database. Answer:

```
1 SELECT store.name
2 FROM classes
3 INNER JOIN instructors ON
4 instructors.teacher_id =
5 classes.teacher_id
6 INNER JOIN store ON store.name =
7 instructors.name
8 WHERE store.location = 'TX'
   AND instructors.teacher_age <= 86
   AND (classes.class_subject = 'english' OR
   classes.class_subject = 'sociology' OR
   classes.class_subject = 'art')
   AND classes.grade >= 89
```

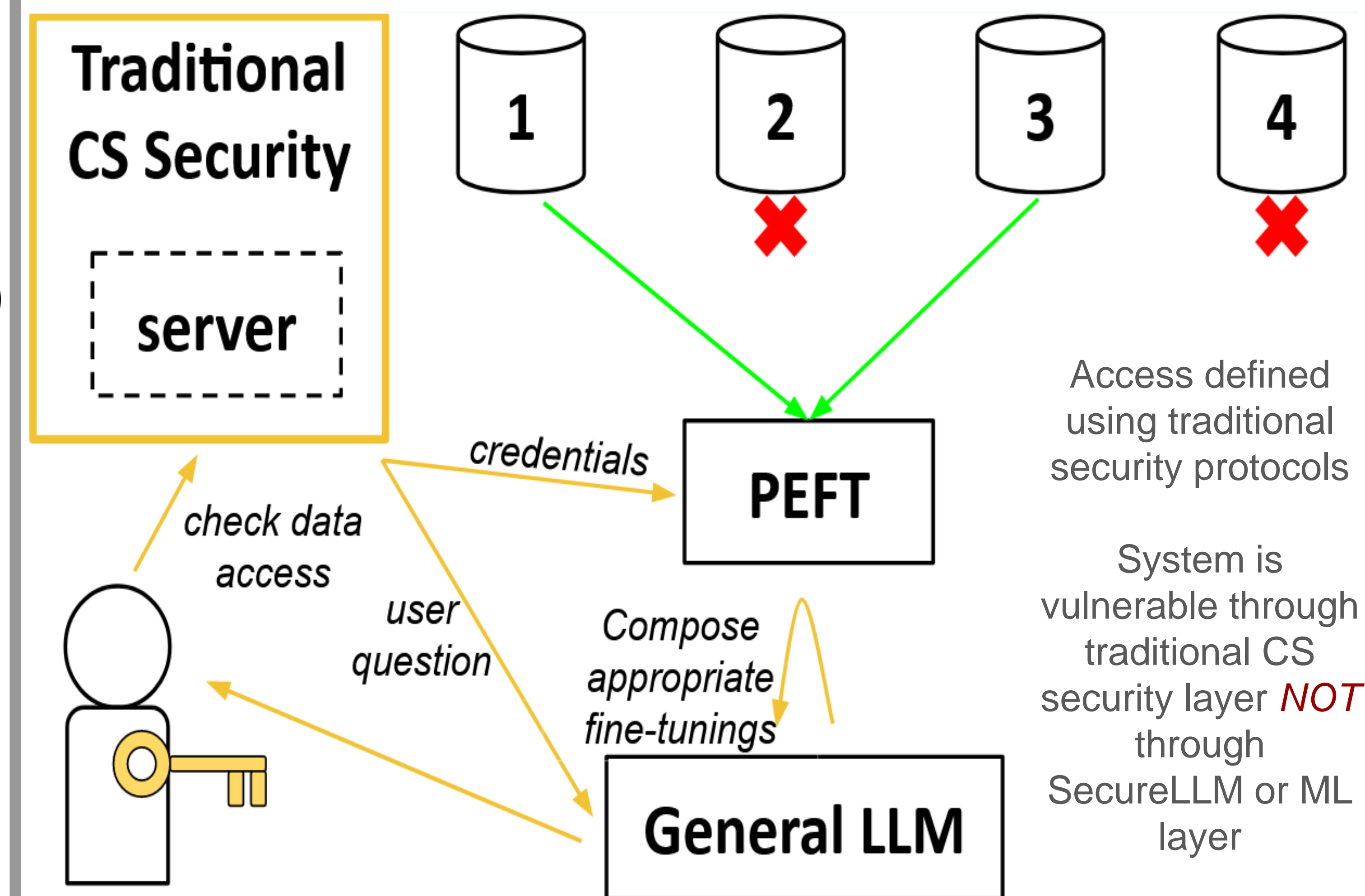
(c) Sample from Union Silo 1,2 ( $S_{1U2}$ )

Q: show me the names of all students living in dorm Lawson Hall that has a meal plan and who are teachers that taught history classes for third graders that were conducted after 2002. Answer:

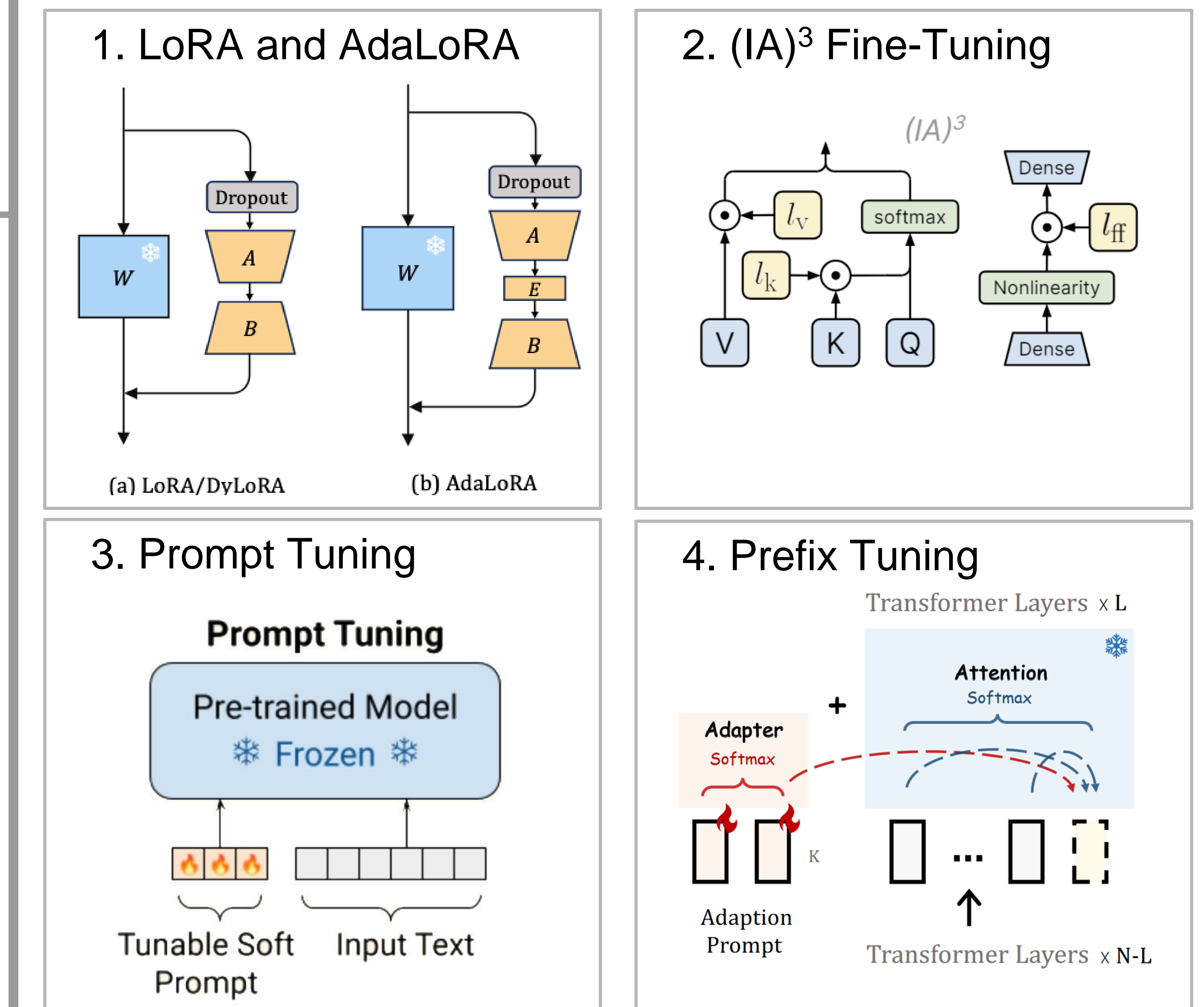
```
1 SELECT pupils.name
2 FROM pupils
3 INNER JOIN instructors ON pupils.name =
4 instructors.name
5 INNER JOIN classes ON instructors.teacher_id
6 = classes.teacher_id
7 WHERE pupils.housing = 'Lawson Hall'
8 AND pupils.cafeteria = 1
9 AND classes.class_subject = 'history'
   AND classes.level = 3
   AND classes.year >= 2002
```

(c) Sample from Union Silo 1,3 ( $S_{1U3}$ )

## SecureLLM Framework:



## Fine-Tunings:



## Composition (Simplicity Triumphs)

Maximum Logit Technique: Query each of fine-tuning take the token with **highest logit**.

$$f(M_1 \circ \dots \circ M_n | x) = f(M_i | x) \text{ where } i := \operatorname{argmax}_{i \in \{1, \dots, n\}} \{ \operatorname{Logits}(M_i | x) \}$$

## Experimental Results: W/ Llama 2-7B

Q&A Accuracy (publicly available datasets)

dataset	generalized model	LoraHub	PEM Addition	p_tuning	prefix	ia3	lora	adadora
commonsenseqa	82%	63%	71%	20%	20%	48%	79%	64%
helloswag	75%	63%	67%	28%	32%	27%	31%	73%
qasc	42%	22%	35%	15%	12%	12%	31%	42%

SQL Generation. Metric is between ground SQL and LLM SQL parsed into trees => Normalized Unordered Tree Edit Distance (Zhang, 1996)

dataset	generalized model	LoraHub	PEM Addition	prefix	ia3	lora	adadora	adadora (w/db norm)
Silos <sub>1</sub>	0.02	0.76	0.56	15.12	5.05	0.12	0.42	0.17
Silos <sub>2</sub>	0.01	0.68	0.43	20.25	2.27	0.15	0.46	0.07
Silos <sub>3</sub>	0.0	0.95	0.48	8.36	2.64	0.06	0.08	0.03
Silos <sub>1U2</sub>	0.36	0.66	0.75	5.27	1.56	0.59	0.83	0.4
Silos <sub>1U3</sub>	0.26	0.6	0.73	17.03	4.95	0.53	1.25	0.29
Silos <sub>2U3</sub>	0.25	0.62	0.75	17.08	2.85	0.47	0.98	0.35
Silos <sub>1U2U3</sub>	0.65	0.88	1.57	7.29	2.46	0.68	1.28	0.39

Ablation #1: Obfuscated column names

dataset	generalized model	LoraHub	PEM Addition	p_tuning	prefix	ia3	lora	adadora	adadora (w/db norm)
Silos <sub>1</sub>	0.0	1.68	0.81	2.93	22.37	5.05	0.47	0.41	0.05
Silos <sub>2</sub>	0.0	0.68	0.55	3.76	31.84	2.27	0.43	0.33	0.14
Silos <sub>3</sub>	0.0	0.76	0.49	1.63	11.07	2.64	0.16	0.08	0.13
Silos <sub>1U2</sub>	0.37	0.91	0.74	2.15	8.12	1.56	0.56	0.65	0.29
Silos <sub>1U3</sub>	0.33	1.69	0.7	3.77	23.44	4.95	0.53	1.14	0.33
Silos <sub>2U3</sub>	0.47	1.22	0.73	3.92	26.54	2.85	0.46	0.53	0.37
Silos <sub>1U2U3</sub>	0.82	1.62	1.99	2.71	12.01	2.46	0.78	0.8	0.49

Ablation #2: GPT rephrased questions

dataset	generalized model	LoraHub	PEM Addition	prefix	ia3	lora	adadora	adadora (w/db norm)
Silos <sub>1</sub>	0.02	0.61	0.4	16.23	2.56	0.2	0.52	0.26
Silos <sub>2</sub>	0.17	0.8	0.38	23.69	0.69	0.35	0.74	0.24
Silos <sub>3</sub>	0.11	0.87	0.29	8.05	2.96	0.19	0.32	0.14
Silos <sub>1U2</sub>	0.4	1.11	0.57	6.69	1.26	0.53	0.64	0.46
Silos <sub>1U3</sub>	0.21	0.59	0.51	17.11	1.29	0.48	0.44	0.2
Silos <sub>2U3</sub>	0.26	0.58	0.44	17.49	0.83	0.37	0.47	0.25
Silos <sub>1U2U3</sub>	0.37	0.7	0.49	9.82	1.7	0.4	0.56	0.23

## Analysis:

- Some PEFT are far more composable than others!
  - Potentially a suboptimal PEFT is superior at composing?
- AdaLoRa/LoRa is close to the generalized model in certain cases! (specifically in unions, which are OOD samples)
- Union datasets are essential for benchmarking this framework
  - Generate union datasets without relying on SQL?
- Fine-tunings are destructive, can we make them cooperative.
  - Future work: normalize by novelty metric per PEFT?