

Motivation

- **Avoid Data Leakage:** LLMs are notoriously bad at keeping secrets
- **Problem:** In large corporations or government databases (many distinct data silos), LLM's cannot be trusted to keep confidential data secret from users who don't have credential to access them
- **Solution: SecureLLM.** Finetune on each silo independently -> Compose at inference time based on user credentials

Dataset

- Independent Silos generated using Context-Free Grammar (Questions, SQL, DB Norm. SQL)

Q: What's the average age of all teachers that are older than 72 and 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
6 OR classes.class_subject = 'art' AND
7 classes.level = 9
```

```
1 want: AVG(teacher_age)
2 tables: instructors, classes
3 conditions:
4 teacher_age >= 72 OR class_subject =
5 'art'
6 level = 9
```

(a) Sample from Silo 1 (S₁)

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 store.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
```

```
1 want: MIN(height)
2 tables: store, inventory
3 conditions:
4 available = 0
5 location = 'NY' OR location = 'CA'
6 OR location = 'MA'
7 star_rating >= 2
```

(b) Sample from Silo 2 (S₂)

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'
9 AND instructors.teacher_age <= 86
10 AND (classes.class_subject = 'english'
11 OR classes.class_subject =
12 'sociology' OR
13 classes.class_subject = 'art')
14 AND classes.grade >= 89
```

```
1 want: name
2 tables: instructors, store, classes
3 conditions:
4 location = 'TX'
5 teacher_age <= 86
6 class_subject = 'english' OR
7 class_subject = 'sociology' OR
8 class_subject = 'art'
9 grade >= 89
```

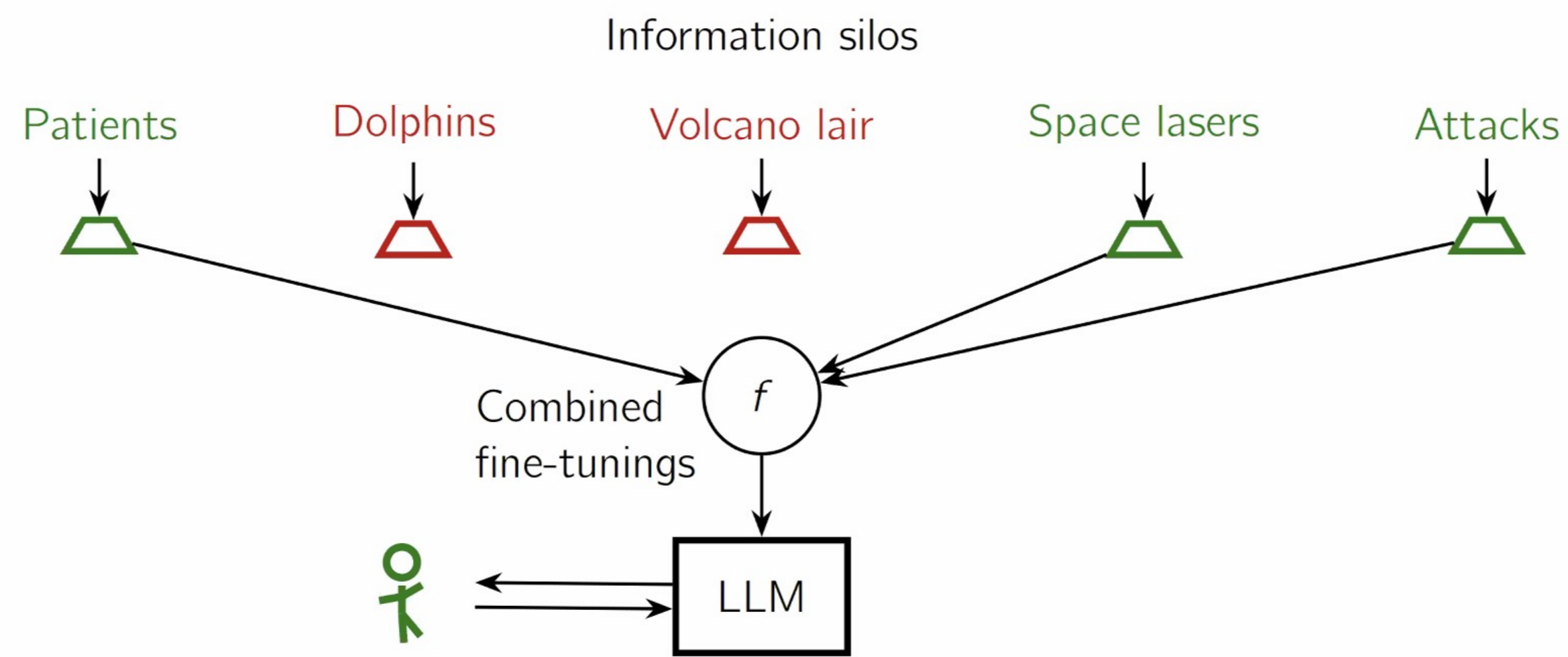
(c) Sample from Union Silo 1,2 (S_{1U2})

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.sloth)
2 FROM inventory INNER JOIN store ON
3 store.bear = inventory.bear
4 WHERE inventory.pony = 0
5 AND (store.alpaca = 'NY'
6 OR store.alpaca = 'CA'
7 OR store.alpaca = 'MA')
8 AND store.raccoon >= 2
```

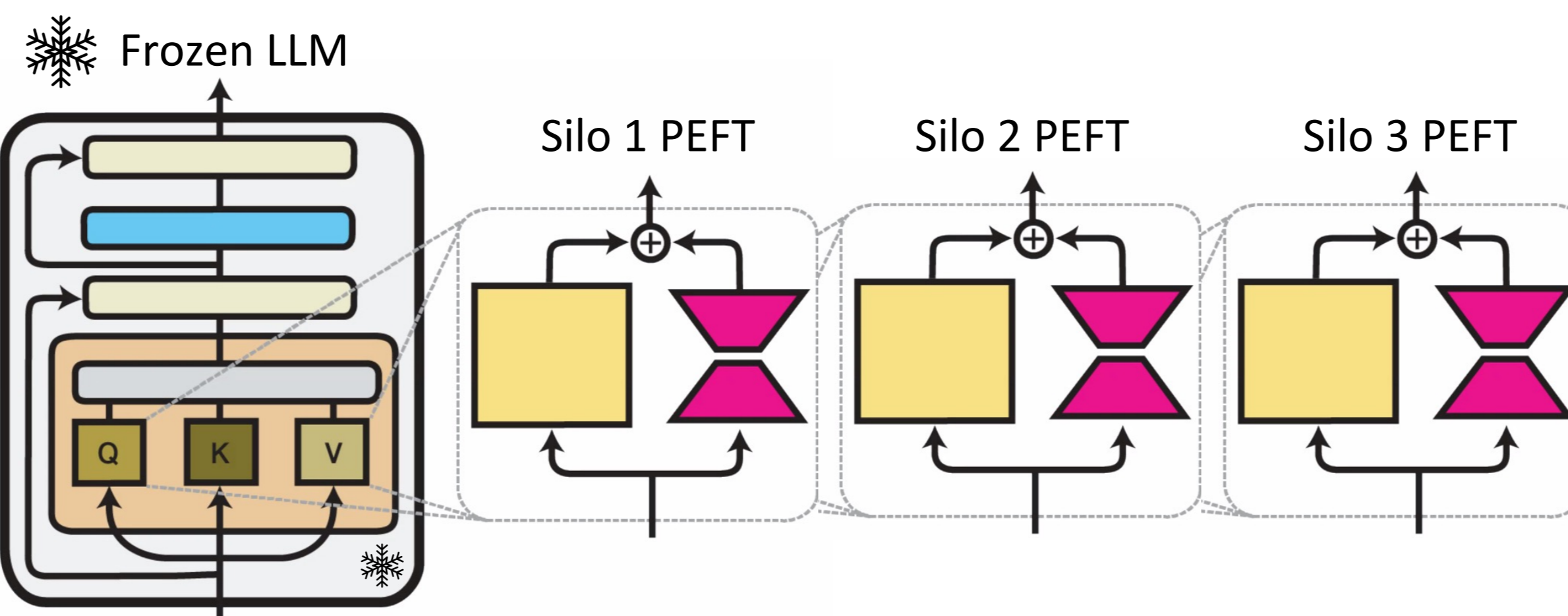
```
1 want: MIN(sloth)
2 tables: store, inventory
3 conditions:
4 pony = 0
5 alpaca = 'NY' OR alpaca = 'CA' OR
6 alpaca = 'MA'
7 raccoon >= 2
```

(b) Sample from Silo 2 (S₂) with obfuscated column names



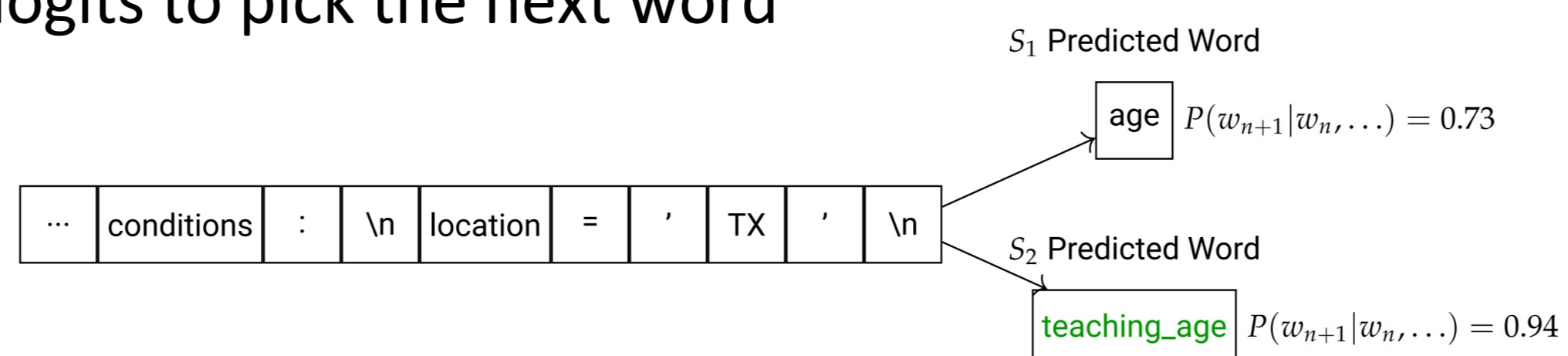
Implementation

- Large transformer network, **Llama-7b**, is frozen while a small Parameter Efficient Fine-Tuning (PEFT) is trained on each information silos
- Low Rank Adapter (LoRA) (Hu et al., 2021)



How it works

- BestLogit independently runs each LoRA and sums predicted logits to pick the next word



Experimental results

- SQL generation using various compositional methods on CFG generated questions (Metric is Tree-Edit distance, lower is better)

CFG Generated	Baseline Exponential Model	Baseline Generalized Model	LoraHub	PEM Addition	Ours (Maximum Difference)	Ours (Logits Without DB Normalization)	Ours (Logits)
Silos ₁	0.0 (100.0%)	0.0 (98.3%)	1.9	1.0	0.4	0.4	0.1
Silos ₂	0.0 (96.7%)	0.0 (100.0%)	2.7	0.7	0.3	0.2	0.0
Silos ₃	0.0 (100.0%)	0.0 (100.0%)	1.2	0.7	0.2	0.1	0.1
Silos _{1U2}	0.0 (98.3%)	1.0 (0.0%)	1.8	1.0	0.9	0.7	0.3
Silos _{1U3}	0.0 (99.2%)	0.5 (0.0%)	1.6	0.7	0.7	0.6	0.2
Silos _{2U3}	0.0 (100.0%)	0.4 (1.7%)	1.7	0.7	0.7	0.5	0.4
Silos _{1U2U3}	0.0 (100.0%)	0.5 (1.7%)	2.2	0.7	0.7	0.4	0.4
$\mu \pm \sigma$	0.00 \pm 0.00	0.35 \pm 0.35	1.88 \pm 0.45	0.78 \pm 0.13	0.55 \pm 0.24	0.42 \pm 0.21	0.21 \pm 0.12

- SQL generation on obfuscated schemas

Obfuscated Generated	Baseline Exponential Model	Baseline Generalized Model	LoraHub	PEM Addition	Ours (Maximum Difference)	Ours (Logits Without DB Normalization)	Ours (Logits)
Silos ₁	0.0 (99.2%)	0.0 (94.2%)	2	1.1	0.5	0.5	0.2
Silos ₂	0.0 (92.5%)	0.0 (100.0%)	3.1	1.4	0.5	0.4	0.4
Silos ₃	0.0 (100.0%)	0.0 (100.0%)	0.9	0.8	0.5	0	0.1
Silos _{1U2}	0.0 (80.8%)	0.7 (0.8%)	1.6	2.2	1.1	0.9	0.5
Silos _{1U3}	0.0 (98.3%)	0.4 (0.0%)	1.3	1.4	0.7	0.6	0.4
Silos _{2U3}	0.0 (77.5%)	0.6 (1.7%)	1.6	2.5	0.9	0.7	0.8
Silos _{1U2U3}	0.0 (100.0%)	0.4 (1.7%)	1.9	2.5	1.0	0.6	0.3
$\mu \pm \sigma$	0.01 \pm 0.01	0.31 \pm 0.28	1.76 \pm 0.65	1.70x \pm 0.65	0.73 \pm 0.24	0.54 \pm 0.24	0.36 \pm 0.22

- GPT rephrased questions (ensures no overfitting)

GPT Generated	Baseline Exponential Model	Baseline Generalized Model	LoraHub	PEM Addition	Ours (Maximum Difference)	Ours (Logits Without DB Normalization)	Ours (Logits)
Silos ₁	0.0 (87.5%)	0.1 (79.2%)	2.0	1.1	0.5	0.6	0.3
Silos ₂	0.2 (60.8%)	0.2 (56.7%)	2.9	1.0	0.5	0.3	0.4
Silos ₃	0.1 (56.7%)	0.2 (51.7%)	1.4	1.1	0.5	0.4	0.2
Silos _{1U2}	0.2 (20.8%)	0.4 (0.0%)	2.1	0.7	0.6	0.6	0.2
Silos _{1U3}	0.2 (29.2%)	0.4 (0.0%)	1.6	1.0	0.6	0.6	0.4
Silos _{2U3}	0.1 (33.3%)	0.3 (3.3%)	2.3	0.6	0.5	0.4	0.3
Silos _{1U2U3}	0.1 (50.0%)	0.3 (2.5%)	2.1	0.6	0.5	0.4	0.2
$\mu \pm \sigma$	0.15 \pm 0.06	0.28 \pm 0.13	2.06 \pm 0.45	0.85 \pm 0.21	0.52 \pm 0.05	0.48 \pm 0.11	0.30 \pm 0.07

