### Exploring Trustworthy Foundation Models: Benchmarking, Finetuning, and Reasoning

#### Prof. Bo Han HKBU TMLR Group / RIKEN AIP Team Assistant Professor / BAIHO Visiting Scientist https://bhanml.github.io







# **Trustworthy Foundation Models**

Benchmarking

Existing datasets are NOT proper to assess if **VLMs** are robust.



**CounterAnimal,** a reliable benchmark for assessing VLMs.

- Scaling backbone models and improving data quality improve the robustness of VLMs.
- Scaling raw training data does not necessary enhance reliability.



Analyzing the dynamics of **LLMs** unlearning is critical yet hard.



- Analyzing gradients provides insights into unlearning dynamics.
- Wrong token reweighting within gradients leads to failures in previous methods.

#### Reasoning

Noisy rationales within chain of thoughts mislead **LLMs** reasoning.



- It is **hard** for LLMs to denoise noisy rationales without guidance.
- It is **easier** for LLMs to denoise by contrasting noisy and clean data.

# Part I: Benchmarking

**Benchmarking** is critical to evaluate and compare model quality.

- Gathering reliable evaluation data.
- Conducting proper metric evaluations.









Qizhou Wang

ObjectNet

Yongqiang Chen

Training and evaluation data have **distribution shifts** to reflect **OOD Generalization**.



Qizhou Wang, Yong Lin, Yongqiang Chen, Ludwig Schmidt, Bo Han, and Tong Zhang.

A Sober Look at the Robustness of CLIPs to Spurious Features. In *NeurIPS*, 2024.



# Supervised vs CLIP Training

#### Supervised Training label supervision



#### CLIP Training cross-modal supervision





Comparison of the **OOD evaluation accuracy** between supervised and CLIP training shows that **CLIP performs better!** 

**Previous Belief:** CLIP is more robust to distribution shifts than conventional supervised training. (Radford et al., 2021)



### Is the Conclusion Correct?

These OOD datasets are crafted for the distribution shifts **within ImageNet setups**, which are **NOT valid for CLIP models.** 

 Data Contamination: Datasets considered OOD for ImageNet-trained models may be ID for CLIP models.



ImageNet V2

*CLIP models may have seen ImageNet V2 during training, which is in fact ID for CLIP setups.* 

Biased Spuriousness: Features that mislead
 ImageNet-trained models may not mislead
 CLIP models necessarily.



ImageNet A

ImageNet A contains data that mislead ImageNet models, which may not make CLIP models fail.

**ImageNet OOD datasets** CANNOT reflect the OOD Generalization for CLIP setups!



### CounterAnimal: A New Benchmark

Is there a benchmark capturing true OOD performance of CLIP?

- Spuriousness: Considering background changes as potential spurious features.
- Generality: The captured spurious features should impact diverse CLIP configurations.



**Basic Assumption**: Since "ice bears" are more commonly appear with "ice" rather than "grass" backgrounds, CLIP may rely on ice-related spurious features.

*The changes of backgrounds represent the impacts of spurious features, which is a typical distribution shift.* 







#### Step 2. Data Curation

Raw data are susceptible to noise and ambiguities, which should be cleansed manually.







clean

noise

occlusion

obscurity



**OBJ labels**: ostrich, African crocodile, water snake, ice bear, and other totally 45 animal names.

BKG labels: ground, water, earth, and other totally 16 background labels.







#### Step 2. Data Curation

Raw data are susceptible to **noise** and **ambiguities**, which should be **cleansed manually**.







clean

noise

occlusion

obscurity









#### Step 2. Data Curation

Raw data are susceptible to **noise** and **ambiguities**, which should be **cleansed manually**.



noise



occlusion

obscurity

clean



OBJ labels: *ostrich, African crocodile, water snake, ice bear,* and other totally **45 animal names**.

**BKG labels**: *ground, water, earth*, and other totally **16 background labels**.





Step 2. Data Curation Raw data are susceptible to **noise** and **ambiguities**, which should be cleansed manually. clean noise occlusion obscurity Step 3. Data Labelling **OBJ**: ice bear **OBJ labels**: ostrich, African crocodile, water snake, ice BKG: snow bear, and other totally 45 animal names. **OBJ**: ice bear BKG labels: ground, water, earth, and other totally 16 **BKG**: grass background labels.



### **CounterAnimal Characteristics**

#### CounterAnimal





Photos of ice bear in snow background



Photos of ice bear in grass background

**Common vs. Uncommon:** Photos are grouped according to their backgrounds. For each class, we identify **group pairs** that cause **high performance drop** when evaluating with CLIP.

Assessing Robustness: The performance drop between common and uncommon groups indicates the robustness of evaluated models.

Data Structure. Images are organized per class and each further divided into two groups: common and uncommon.



### **CounterAnimal Characteristics**





*The data distributions illustrate variations across different animal classes, categorized into common and uncommon groups. The horizontal axis denotes the class IDs*, *e.g., ID* 1 *to "ostrich", ID* 2: *to "brambling", …, ID* 8 *to "box turtle", ID* 9 *to "common iguana",…, ID* 18 *to "scorpion", ID* 19 *to "tarantula", …, ID* 32 *to "African hunting dog", ID* 33 *to "hyena", …*.

# We collect **45 classes** of animals with **7,000 common** and **6,000 uncommon** examples.

Data Structure. Images are organized per class and each further divided into two groups: common and uncommon.



drop

8.26

9.97

17.59

15.10

9.39

15.96

10.22

uncommon

39.73

30.09

63.31

70.28

79.90

66.27

80.55

#### **Experimental Results**

common acc – uncommon acc

CLIP Training

CounterAnimal

(ImageNet) Supervised Training

#### Other LVLMs (large VLMs)

common

47.99

40.06

80.90

85.38

89.29

82.23

90.77

backbone	pre-train dataset	common	uncommon	drop	_
RN-101	OpenAI	64.27	45.15	19.12	Г
$RN-50 \times 4$	OpenAI	70.02	49.07	20.95	
ViT-B/16	LAION400M	73.11	52.17	20.94	
ViT-B/16	OpenAI	73.08	56.56	16.52	
ViT-B/16	$DataComp1B^*$	80.36	64.24	16.12	
ViT-B/16	LAION2B	73.18	53.18	20.00	
ViT-B/16	$DFN2B^*$	85.03	70.61	14.42	
ViT-B/32	LAION400M	67.13	36.95	30.18	
ViT-B/32	OpenAI	69.13	45.62	23.51	
ViT-B/32	$DataComp1B^*$	75.96	53.74	22.22	
ViT-B/32	LAION2B	72.94	48.74	24.20	
ViT-L/14	LAION400M	80.90	63.31	17.59	
ViT-L/14	OpenAI	85.38	70.28	15.10	(
ViT-L/14	$\mathtt{DataComp1B}^*$	89.29	79.90	9.39	(
ViT-L/14	LAION2B	82.23	66.27	15.96	0
ViT-L/14	DFN2B*	90.77	80.55	10.22	
ViT-L/14-336	OpenAI	86.36	73.14	13.21	in
ViT-H/14	LAION2B	85.74	73.13	12.61	
ViT-H/14	DFN5B*	88.55	79.13	9.42	
ViT-G/14	LAION2B	86.81	73.32	13.49	
ViT-bigG/14	LAION2B	87.57	76.96	10.61	

backbone	common	uncommon	drop
AlexNet	59.56	39.24	20.31
VGG-11	73.37	56.12	17.25
VGG-13	75.33	58.43	16.90
VGG-19	77.84	61.74	16.10
RN-18	74.36	56.07	18.29
RN-34	78.31	61.01	17.30
RN-50	81.44	66.07	15.37
RN-101	81.76	68.18	13.57
ViT-B/16	84.97	74.98	9.99
ViT-B/32	79.84	64.36	15.48
ViT-L/16	83.74	72.69	11.05
ViT-L/32	81.23	67.54	13.69
ConvNext-S	88.27	79.97	8.30
ConvNext-B	88.60	80.53	8.07
ConvNext-L	89.12	81.47	7.65

#### different LVLM paradigms

LVLMs

MiniGPT4-Viccuna7B

LLaVA1.5-7B

CLIP-LAION400M-ViT-L/14

CLIP-OpenAI-ViT-L/14

CLIP-DataComp1B-ViT-L/14

CLIP-LAION2B-ViT-L/14

CLIP-DFN2B-ViT-L/14

What observations can we draw from these results?

#### increasing model scale

#### increasing diverse model scale data source



#### Observations

DataComp (DC) and Data Filtering Networks (DFN) are two high-quality CLIP data sources.



The marker size indicates the backbone scale, and the color shade indicates pre-train data scale.

#### **Observation 1** (ImageNet Models vs. CLIPs).

ImageNet models perform better than CLIPs against spuriousness within CounterAnimal.

**Note.** CounterAnimal characterizes the spuriousness within CLIPs, thus proper for assessing CLIPs.

#### **Observation 2** (CLIPs vs. More Advanced LVLMs).

LLaVA and MinGPT4 show stronger robustness (closer to y = x) yet with **lower performance** than CLIPs.

**Note.** More advanced VLMs built upon CLIPs are still affected by spuriousness within CounterAnimal.



**Observation 3 (Model Size).** Scaling up model size CAN enhance CLIP robustness.





**Observation 4 (Data Size).** Scaling up data size CANNOT enhance CLIP robustness.



**Observation 5 (Data Quality).** Improving data quality CAN enhance CLIP robustness.



### Theoretical Understanding

Assumption (Multi-modal Dataset). Considering *n* image-text pairs  $\{(x_{I}^{i}, x_{T}^{i})\}_{i=1}^{n}$ , both  $x_{I}^{i}$  and  $x_{T}^{i}$  are generated from the latent factor  $z_{i}$ , where  $z = [z_{inv}, z_{spu}] \in \mathbb{R}^{2}$  is composed of

• **invariant feature**  $z_{inv} \sim \mathcal{N}(\mu_{inv}y, \sigma_{inv}^2)$ 

• spurious feature  $z_{spu} \sim \mathcal{N}(\mu_{spu}a, \sigma_{spu}^2)$ with  $\Pr(a = y) = p_{spr}$  otherwise a = -y. y is the label uniformly drawn from  $\{-1,1\}$ . The training data  $\mathcal{D}^{tr}$  is drawn with  $\frac{1}{2} \leq p_{spr} \leq 1$  and test data  $\mathcal{D}^*$  is drawn with  $p_{spu} = \frac{1}{2}$ .

Note. The dataset is **biased** to spurious feature  $z_{spu}$  due to **different**  $p_{spr}$  between training and test.

**Theorem 1**. Given the multi-modal dataset with a large spurious correlation  $p_{spu} = 1 - o(1)$ . Then, under reasonable assumptions, w.p. at least 1 - O(1), the CLIP model achieves

- a small zero-shot error on test data where a = y: Acc $(g_I, g_T) \ge 1 \Phi(\kappa_2) o(1)$ ,
- a large zero-shot error on test data where  $a \neq y$ :  $\operatorname{Err}(g_{\mathrm{I}}, g_{\mathrm{T}}) \geq 1 \Phi(\kappa_{1}) o(1)$ . Therein,  $\kappa_{1}, \kappa_{2}$  are constants that depend on  $\mu_{inv}, \sigma_{inv}, \mu_{inv}$ , and  $\sigma_{inv}$ .

**Note.** The model relies on whether a = y (whether biased) to make right predictions.



#### Take Home Messages

We should be cautious about **test setups** when assessing new **training setups**.

**CounterAnimal** (<u>https://counteranimal.github.io/</u>) is a proper benchmark for assessing the robustness of CLIPs to spurious features.

**Distribution shifts** remain an open question for CLIP and other VLMs.

Scaling up model size can enhance robustness, while scaling up pre-train data is not that effective.

Improving data quality is effective to enhance robustness.

# Part II: Finetuning







Qizhou Wang

Zhanke Zhou

**Finetuning** aims to adapts the model parameters to fit tasks or knowledge, of which the specific goals can be attributed to **learning** and **unlearning**.



fine-tuning to learn/update knowledge

**Qizhou Wang**, Jin Peng Zhou, **Zhanke Zhou**, Saebyeol Shin, **Bo Han**, Kilian Q. Weinberger. Rethinking LLM Unlearning Objectives: A Gradient Perspective and Go Beyond. In *ICLR*, 2025. https://bhanml.github.io & https://github.com/tmlr-group



### Right to be Forgotten



"The data subject shall have the right to obtain from the controller the erasure of personal data concerning him or her without undue delay and the controller shall have the obligation to erase personal data ..."



"A consumer shall have the right to request that a business delete any personal information about the consumer which the business has collected from the consumer ..."

# LLM Unlearning

#### **Bi-objective Goal**

- Unlearn: removing model capability to generate targeted data  $\mathcal{D}_{u} = \{s_{u}\}_{n_{u}}$
- Retain: maintain performance on other non-targeted data  $\mathcal{D}_r = \{s_r\}_{n_r}$

Gradient Ascent (GA)-based Method

not to be unlearned

to be unlearned

**Basic Assumption**: If the negative log-likelihood is a proper objective for learning, then the log-likelihood should be appropriate for unlearning.





#### Impacts of GA

#### Negative log-likelihood (NLL) as the metric $\mathcal{R}$ to assess performance.



**Observation 1.** GA-based methods CAN achieve strong unlearning but CANNOT ensure reliable retention, thus **NOT meeting the dual-objective goal.** 



#### Delve Deeper?

Performance metrics offer limited insights towards deeper understandings.

**Limitation 1.** We CANNOT **disentangle** the impacts of  $\mathcal{L}_u(\mathcal{D}_u; \theta)$  and  $\mathcal{L}_r(\mathcal{D}_r; \theta)$  on model performance.



Both  $\mathcal{L}_u(\mathcal{D}_u; \boldsymbol{\theta})$  and  $\mathcal{L}_r(\mathcal{D}_r; \boldsymbol{\theta})$  have impacts on  $\mathcal{R}(\mathcal{D}_u; \boldsymbol{\theta})$  and  $\mathcal{R}(\mathcal{D}_r; \boldsymbol{\theta})$  in an **intertwined** manner.

Using NLL to assess performance changes regarding unlearning and retention.



### Delve Deeper?

Performance metrics offer **limited** insights towards deeper understandings.

**Limitation 2.** Even disentangled, we CANNOT fully **understand the factors** that lead to the observed behaviors.





#### G-effect: A Gradient View

Studying the impacts of **unlearning methods** (e.g., GA) on **performance metrics** (e.g., NLL) from a gradient view.



- Fulfill Goal 1 as the G-effect can be computed for  $\mathcal{L}_u(\mathcal{D}_u; \theta)$  and  $\mathcal{L}_r(\mathcal{D}_r; \theta)$  separately.
- Fulfill Goal 2 as gradients provide more messages than merely CE performance.



### G-effect: An Example

**Retain G-effect:**  $e_r = \nabla_{\theta} \mathcal{L}(\mathcal{D}_u; \theta)^\top \nabla_{\theta} \mathcal{R}(\mathcal{D}_r; \theta)$ . A **positive**  $e_r$  is preferred to enhance retention.

**Unlearn G-effect:**  $e_{u} = \nabla_{\theta} \mathcal{L}(\mathcal{D}_{u}; \theta)^{\top} \nabla_{\theta} \mathcal{R}(\mathcal{D}_{u}; \theta)$ . A **negative**  $e_{u}$  is preferred for strong unlearning.



**Note.** The G-effect quantifies the **rate of change** (increase/decrease) in performance, which can be calculated **separately** for retention and unlearning.



### GA: Objective 1



**Objective:**  $\mathbb{E}_{\mathcal{D}_{u}} \sum_{i} \log P(s_{u}^{i} | s_{u}^{< i}; \boldsymbol{\theta})$ 

**Gradient:** 
$$\mathbb{E}_{\mathcal{D}_{u}} \sum_{i} \frac{1}{P(s_{u}^{i}|s_{u}^{< i}; \boldsymbol{\theta})} \nabla_{\boldsymbol{\theta}} P(s_{u}^{i}|s_{u}^{< i}; \boldsymbol{\theta})$$

inverse likelihood

**Observation 2. Excessive extent of removal** incurs negative costs to retention.

**Reason.** The inverse likelihood wrongly focuses more on sufficiently unlearned tokens, leading to **over-unlearning** that negatively impacts model utility.



### GA: Objective 1



The G-effects of GA (closer look).

**Observation 3.** Unlearning **affects on bottom layers** of LLMs more than others.

**Reason.** Large gradients will **accumulate** due to the chain rule, a general scenario holds for many other unlearning objectives.



#### WGA: Improvement 1

Motivation: Combating the inverse likelihood term via loss reweighting.

Original GA:  $\mathbb{E}_{\mathcal{D}_{u}} \sum_{i} \log P(s_{u}^{i} | s_{u}^{<i}; \theta) \rightarrow \text{Weighted GA: } \mathbb{E}_{\mathcal{D}_{u}} \sum_{i} P(s_{u}^{i} | s_{u}^{<i}; \theta)^{\alpha} \log P(s_{u}^{i} | s_{u}^{<i}; \theta)$ Gradients:  $\mathbb{E}_{s_{u} \sim \mathcal{D}_{u}} \sum_{i} P(s_{u}^{i} | s_{u}^{<i}; \theta)^{\alpha-1} \nabla_{\theta} P(s_{u}^{i} | s_{u}^{<i}; \theta)$ 

counteract the inverse likelihood



Comparison of the G-effects between GA and WGA.



### NPO: Objective 2



The G-effects of NPO.

**Observation 4.** NPO (Negative Preference Optimization) has **fewer negative impacts** on retention compared to GA. **Reason.** The gradients of NPO are very similar to GA, yet further **reweighting** by  $w_{npo}$ , which mainly contributes to its improvements over GA.



### NPO: Objective 2



**Objective:** 
$$\mathbb{E}_{\mathcal{D}_{u}} \frac{1}{\beta} \log(1 + \left(\frac{p(s_{u}; \theta)}{p(s_{u}; \theta_{o})}\right)^{\beta})$$

**Gradient:** 
$$\mathbb{E}_{\mathcal{D}_{u}} \sum_{i} \frac{2P(s_{u}; \boldsymbol{\theta})^{\beta}}{P(s_{u}; \boldsymbol{\theta})^{\beta} + P(s_{u}; \boldsymbol{\theta}_{o})^{\beta}} \nabla_{\boldsymbol{\theta}} \log P(s_{u}; \boldsymbol{\theta})$$
  
 $w_{npo}$  reweighting

**Observation 5.** The NPO weight  $w_{npo}$  serves a role like **early stopping. Reason.**  $w_{npo}$  approaches 0 when  $P(s_u; \theta) \rightarrow 0$ .



### NPO: Objective 2

Larger weights are assigned to those instances with larger retaining PG-effects.



The distributions of the point-wise G-effects across different range of  $w_{npo}$ .

**Gradient:** 
$$\mathbb{E}_{\mathcal{D}_{\mathbf{u}}} \sum_{i} \frac{2P(s_{\mathbf{u}}; \boldsymbol{\theta})^{\beta}}{P(s_{\mathbf{u}}; \boldsymbol{\theta})^{\beta} + P(s_{\mathbf{u}}; \boldsymbol{\theta}_{\mathbf{o}})^{\beta}} \nabla_{\boldsymbol{\theta}} \log P(s_{\mathbf{u}}; \boldsymbol{\theta})$$

**G-effect:**  $\mathbb{E}_{\mathcal{D}_{u}} w_{npo} \nabla_{\boldsymbol{\theta}} \log p(\boldsymbol{s}_{u}; \boldsymbol{\theta})^{\top} \nabla_{\boldsymbol{\theta}} \mathcal{R}(\mathcal{D}; \boldsymbol{\theta})$ 

weights point-wise G-effect (PG-effect)

(The impacts of a particular data point on model performance.)

**Observation 6.** The NPO reweighting mechanism  $w_{npo}$  prioritizes instances that less damages retention.

Reason. Data that have small impacts on retention also have small impacts on unlearning.



#### TNPO: Improvement 2

**Motivation: Generalized** the reweighting mechanism of NPO for tokens. **Token-wise NPO**  $\sum_{i} w_{tnpo}^{i} \log P(s_{u}^{i}|s_{u}^{<i}; \theta)$  with  $w_{tnpo}^{i} = \frac{2P(s_{u}^{i}|s_{u}^{<i}; \theta)^{\alpha}}{P(s_{u}^{i}|s_{u}^{<i}; \theta)^{\alpha} + P(s_{u}^{i}|s_{u}^{<i}; \theta)^{\alpha}}$ 

same reweighting scheme yet applied point-wise.



Comparison of the G-effects between NPO and TNPO.



#### Retain Objectives



Comparison between two representative retain objectives.

**Observation 7. NLL** and **KL** are both effective for retention, while KL can lead to overall larger retain G-effect, thus preferred.

**Note.** The unlearn G-effect for the unlearning objective is much larger than for the retain objectives. Thus, we do not need to worry about the side effect on unlearning.



#### Empirical evaluations

I	LLM			Phi-	1.5					Llama-	2-7B		
setup	method	ES-	exact	ES-p	perturb	MU↑	FQ ↑	ES-	exact	ES-p	berturb	MU↑	FQ↑
		retain T	uniearn ↓	retain T	unlearn ↓			retain T	uniearn ↓	retain T	uniearn ↓		
before	unlearning	0.44	0.59	0.21	0.16	0.52	-5.80	0.82	0.80	0.53	0.40	0.63	-7.59
	GA	0.11	0.05	0.08	0.08	0.37	-0.54	0.42	0.05	0.26	0.04	0.53	-0.54
	PO	0.36	0.84	0.16	0.36	0.51	-4.24	0.75	0.83	0.47	0.52	0.62	-5.80
	WGA	0.36	0.03	0.18	0.02	0.51	-0.54	0.67	0.08	0.38	0.06	0.65	-0.08
1%	NPO	0.27	0.09	0.11	0.07	0.48	-2.91	0.47	0.12	0.38	0.09	0.62	-1.32
	TNPO	0.33	0.03	0.12	0.04	0.49	-0.08	0.51	0.03	0.43	0.03	0.64	-0.08
	RMU	0.23	0.08	0.15	0.05	0.43	-0.54	0.23	0.08	0.15	0.05	0.52	-1.32
before	unlearning	0.44	0.56	0.21	0.23	0.52	-29.65	0.82	0.77	0.53	0.41	0.63	-32.1
	GA	0.00	0.00	0.00	0.00	0.00	-11.40	0.03	0.00	0.02	0.00	0.00	-12.4
	PO	0.26	0.79	0.16	0.49	0.51	-26.50	0.55	0.84	0.36	0.49	0.64	-28.8
	WGA	0.29	0.01	0.16	0.01	0.51	-1.30	0.47	0.00	0.39	0.00	0.64	-16.3
5%	NPO	0.08	0.12	0.08	0.06	0.38	-7.75	0.17	0.07	0.12	0.08	0.52	-9.95
	TNPO	0.16	0.01	0.08	0.00	0.46	-2.18	0.50	0.01	0.34	0.00	0.63	-32.1
	RMU	0.21	0.00	0.12	0.00	0.27	-1.95	0.12	0.00	0.12	0.00	0.58	-21.4
before	unlearning	0.44	0.47	0.21	0.18	0.52	-39.00	0.82	0.83	0.53	0.30	0.63	-44.4
	GA	0.00	0.00	0.00	0.00	0.00	-45.26	0.00	0.00	0.00	0.00	0.00	-20.8
	PO	0.32	0.73	0.14	0.26	0.50	-38.25	0.55	0.84	0.37	0.43	0.62	-39.7
	WGA	0.34	0.00	0.16	0.00	0.51	-9.06	0.66	0.02	0.42	0.01	0.62	-24.8
10%	NPO	0.08	0.09	0.07	0.07	0.38	-10.57	0.12	0.13	0.10	0.14	0.50	-12.1
	TNPO	0.20	0.01	0.09	0.01	0.50	-7.66	0.45	0.01	0.26	0.01	0.63	-13.4
	RMU	0.03	0.05	0.03	0.06	0.31	-7.00	0.25	0.01	0.20	0.01	0.59	-16.7

**Observation 8.** Larger unlearning datasets and smaller model sizes make it more challenging to unlearn.

**Observation 9.** GA-based works (GA & TNPO) are superior to other lines of works like PO or RMU.

**Observation 10.** Instance-wise reweighting is promising for unlearning efficacy.

Comparison between unlearning objective on TOFU with KL regularization.



#### Take Home Messages

General knowledge within **shallow layers undergoes substantial alterations** over deeper layers during unlearning.

Although conceptually existing, **current objectives all fail** to retain the overall performance when conducting unlearning.

**Prioritizing some tokens** is effective for unlearning. However, there still exists a large space to further refine weighting mechanisms.

With **excessive unlearning**, the deterioration in common model responses can outweigh improvements in unlearning.

# Part III: Reasoning



#### Can Language Models Perform Robust Reasoning in Chain-of-thought Prompting with Noisy Rationales?



Zhanke Zhou



Jianing Zhu



Zhanke Zhou, Rong Tao, Jianing Zhu, Yiwen Luo, Zengmao Wang, Bo Han.

Can Language Models Perform Robust Reasoning in Chain-of-thought Prompting with Noisy Rationales? In *NeurIPS*, 2024 https://bhanml.github.io & https://github.com/tmlr-group

### Background



Reasoning is the pathway to achieve powerful intelligence.

- Decompose a complex problem into feasible steps.
- Combine knowledge pieces into new knowledge.

Generating **chain of thoughts (CoT)** is the key of several reasoning models.





# Chain of Thoughts (CoT)

In-context learning (ICL) is widely used.

• ICL enable LLMs to **learn from a few** ... examples without fine-tuning.



**Chain of thoughts (CoT) prompting** can elicit the reasoning capabilities of LLMs.

• Beyond examples, CoT includes **rationales**, i.e., sequential reasoning thoughts to solve a question.

```
Input: CoT prompting with rationales
Question-1: In base-9, what is 86+57?
Rationale-1: In base-9, the digits are "012345678". We have 6 + 7 = 13 in base-
10. Since we're in base-9, that exceeds the maximum value of 8 for a single digit.
13 mod 9 = 4, so the digit is 4 and the carry is 1. We have 8 + 5 + 1 = 14 in base
10. 14 mod 9 = 5, so the digit is 5 and the carry is 1. A leading digit 1. So the
answer is 154.
Answer-1: 154.
...Q2, R2, A2, Q3, R3, A3 ...
Question : In base-9, what is 62+58?
```





# New Challenge in LLM Reasoning

Existing work generally assumes that CoT contains clean rationales.

But, what if CoT contains noisy rationales? 🧐

• noisy rationales include irrelevant or inaccurate thoughts.

Input: CoT prompting with clean rationales	Input: CoT prompting with noisy rationales
Question-1: In base-9, what is 86+57?	Question-1: In base-9, what is $86+57$ ?
Rationale-1: In base-9, the digits are "012345678". We have 6 + 7 = 13 in base-	Rationale-1: In base-9, the digits are "012345678". We have $6 + 7 = 13$ in base-
10. Since we're in base-9, that exceeds the maximum value of 8 for a single digit.	10. $13 + 8 = 21$ . Since we're in base-9, that exceeds the maximum value of 8 for a
13 mod 9 = 4, so the digit is 4 and the carry is 1. We have 8 + 5 + 1 = 14 in base	single digit.13 mod 9 = 4, so the digit is 4 and the carry is 1. We have $8 + 5 + 1 =$
10. 14 mod 9 = 5, so the digit is 5 and the carry is 1. A leading digit 1. So the	14 in base 10. 14 mod 9 = 5, so the digit is 5 and the carry is 1. $5 + 9 = 14$ . A
answer is 154.	leading digit is 1. So the answer is 154.
Answer-1: 154.	Answer-1: 154.
Q2, R2, A2, Q3, R3, A3	Q2, R2, A2, Q3, R3, A3
Question : In base-9, what is 62+58?	Question: In base-9, what is 62+58?

While the test question asks about **base-9 calculation**.

The irrelevent have 10 information is included in ratio



# New Challenge in LLM Reasoning

#### Noisy rationales originate from diverse sources.

• Such as crowdsourced platforms, dialogue systems, and AI-generated data.



However, the robustness of LLMs against noisy rationales is still unknown.

- A new dataset is needed to conduct a systematic evaluation of current LLMs.
- To verify the corresponding **countermeasures** against noisy rationales.

![](_page_42_Picture_0.jpeg)

## Noisy Rationales Benchmark (NoRa)

- We construct a new benchmark to evaluate the robustness against noisy rationales.
- NoRa contains 26,391 questions, covering 3 tasks: math, symbolic, and commonsense.

![](_page_42_Figure_4.jpeg)

Table 1: Noisy rationales (consisting <u>noisy thoughts</u>) sampled from the NoRa dataset. Full examples of NoRa are in Appendix C.6, and real-world examples of noisy rationales are in Appendix C.3.

![](_page_43_Picture_0.jpeg)

# Noisy Rationales Benchmark (NoRa)

Definitions

- Irrelevant thoughts are irrelevant to the given context.
  - E.g., discussing the genetic overlap of siblings when reasoning the family roles.
- Inaccurate thoughts are factual errors in the given context.
  - E.g., "5+5=10" is wrong in base-9 calculation.

Benchmark construction

- Generating noisy rationales by **inserting irrelevant or inaccurate thoughts.**
- Guarantee the overall correctness without modifying the question or answer.
- Control **noise ratios** (noisy thoughts / clean thoughts) with values 0.3,0.5,0.8.

(easy medium hard)

![](_page_44_Picture_0.jpeg)

### Empirical Evaluations with NoRa

Grand observation: The base LLM (GPT-3.5) with all the existing methods is severely affected by noisy rationales.

- Up to **25.3%** acc decrease with irrelevant noise.
- Up to **54.0%** acc decrease with inaccurate noise (compared acc with clean rationales).

#### **Observation 1:**

Self-correction methods (ISC, SP) perform **poorly** on most tasks with noisy rationales.

#### **Observation 2:**

Self-consistency methods (SM, SD, SC) can improve robustness **without** true denoising.

	Task	Method M	$\mathrm{Acc}(\mathcal{M},\mathcal{Q},\mathcal{P}_{\mathrm{clean}})$	Easy	$egin{array}{l} \operatorname{Acc}(\mathcal{M},\mathcal{Q}, \ \operatorname{Medium}) \end{array}$	$\mathcal{P}_{ ext{irrelevant}}) \\  ext{Hard}$	Avg.	A Easy	$\mathrm{Acc}(\mathcal{M},\mathcal{Q}, \mathcal{M})$ Medium	$\mathcal{P}_{ ext{inaccurate}} \  ext{Hard}$	Avg.
	Math Base-9	Base w/ ISC [29] w/ SP [89] w/ SM [62] w/ SD [102] w/ SC [83]	46.4 24.3 26.2 37.4 47.9 <b>61.5</b>	39.3 17.7 25.5 30.0 37.2 <b>51.1</b>	30.3 14.7 25.5 22.7 25.4 <b>39.0</b>	26.6 12.7 21.9 16.5 24.7 <b>36.2</b>	32.1 15.0 24.3 23.1 29.1 <b>42.1</b>	23.2 18.4 20.0 24.7 <u>29.3</u> <b>32.7</b>	10.1 13.7 <u>18.4</u> <b>19.2</b> 12.5 15.3	6.0 12.3 <b>14.3</b> <u>12.4</u> 8.7 7.5	13.1 14.8 17.6 <b>18.8</b> 16.8 <u>18.5</u>
e. se	Math Base-11	Base w/ ISC [29] w/ SP [89] w/ SM [62] w/ SD [102] w/ SC [83]	23.9 11.2 20.7 16.3 17.9 <b>33.7</b>	<u>19.1</u> 8.3 17.5 12.0 12.3 <b>25.3</b>	13.6 7.8 <b>16.7</b> 6.0 12.0 <u>16.3</u>	$   \begin{array}{r}     10.7 \\     6.0 \\     \underline{14.0} \\     \overline{5.7} \\     13.3 \\     15.0   \end{array} $	14.5 7.4 <u>16.0</u> 7.9 12.5 <b>18.9</b>	14.0 6.5 <u>14.1</u> 12.0 17.0 <b>19.7</b>	6.7 5.2 <b>10.7</b> 9.3 8.7 9.3	3.6 4.7 <b>10.8</b> <u>7.7</u> 5.3 3.3	8.1 5.5 <b>11.9</b> 9.7 10.3 <u>10.8</u>
	Symbolic Equal	Base w/ ISC [29] w/ SP [89] w/ SM [62] w/ SD [102] w/ SC [83]	32.7 23.9 23.2 25.0 9.9 <b>35.3</b>	28.1 20.0 23.0 20.7 10.1 <b>31.0</b>	25.1 16.3 22.6 19.7 10.9 <b>28.3</b>	23.0 15.5 22.7 16.7 10.3 <b>27.0</b>	25.4 17.3 22.8 19.0 10.4 <b>28.8</b>	29.1 19.2 23.7 21.0 10.1 <b>33.3</b>	26.1 18.3 22.5 20.3 10.9 <b>30.7</b>	22.7 18.1 23.5 20.0 10.4 <b>26.0</b>	26.0 18.5 23.2 20.4 10.5 <b>30.0</b>
	Symbolic Longer	Base w/ ISC [29] w/ SP [89] w/ SM [62] w/ SD [102] w/ SC [83]	9.2 4.9 5.1 1.7 0.1 <b>13.0</b>	6.3 4.6 4.3 0.7 0.1 <b>7.7</b>	7.2 2.7 4.1 0.7 0.1 <b>9.0</b>	6.0 3.7 3.9 1.3 0.2 <b>6.3</b>	6.5 3.7 4.1 1.0 0.1 <b>7.7</b>	7.0 3.4 4.9 1.3 0.1 <b>8.0</b>	6.8 4.3 4.0 0.7 0.3 <b>8.0</b>	6.0 3.3 4.5 0.3 0.0 <b>8.7</b>	6.6 3.7 4.5 0.8 0.1 <b>8.2</b>
	Commonsense	Base w/ ISC [29] w/ SP [89] w/ SM [62] w/ SD [102] w/ SC [83]	45.7 21.8 47.9 53.3 <b>54.0</b> 52.0	44.3 24.3 48.2 50.3 <b>58.3</b> 46.3	42.3 22.5 46.7 50.0 <b>57.3</b> 45.0	41.4 21.4 48.1 46.7 <b>57.7</b> 44.7	42.7 22.7 47.7 49.0 <b>57.8</b> 45.3	36.7 23.3 49.6 47.7 <b>57.0</b> 44.7	33.4 26.5 46.6 49.0 <b>58.3</b> 44.7	28.3 24.0 46.5 49.3 <b>53.7</b> 38.0	32.8 24.6 47.6 <u>48.7</u> <b>56.3</b> 42.5

Table 3: Reasoning accuracy on NoRa dataset with 3-shot prompting examples with clean, irrelevant, or inaccurate rationales. The **boldface** numbers mean the best results, while the <u>underlines</u> numbers indicate the second-best results. Note the referenced results of Base model are highlighted in gray.

#### Baseline methods:

- Intrinsic Self-correction (ISC)
- Self-polish (SP) SmoothLLM (SM)

2

Experiments with GPT-4 are in Appendix F

https://bhanml.github.io & https://github.com/tmlr-group

• Self-denoise (SD) Self-consistency (SC)

45

![](_page_45_Picture_0.jpeg)

#### Empirical Evaluations with NoRa

Task	Setting	0	Ten 0.3	ipera 0.5	ture 0.7	1
Base-9	clean ina. easy ina. hard	<b>61.0</b> <b>29.7</b> 5.0	$\frac{\underline{60.9}}{\underline{28.0}}$ $\underline{5.1}$	57.5 27.2 <b>5.5</b>	55.3 26.6 4.6	46.4 21.7 5.0
Base-11	clean irr. easy irr. hard	<b>34.0</b> 21.7 <u>17.0</u>	<u>33.8</u> <u>23.1</u> <b>17.5</b>	31.6 21.3 15.5	29.8 23.3 14.1	23.9 19.1 10.7
Sym.(E)	clean irr. easy irr. hard	34.2 28.6 <b>27.0</b>	<b>35.8</b> <b>31.5</b> 26.1	<u>35.7</u> 29.8 26.2	34.6 29.1 24.0	32.7 28.1 23.0
Sym.(L)	clean ina. easy ina. hard	6.3 5.0 4.0	8.3 7.3 6.1	8.9 8.6 6.3	$\frac{8.9}{8.3}$ $\frac{6.2}{6.2}$	<b>9.3</b> 7.0 6.0

Table 4: Comparing performances of the base model with different temperatures. Sym.(E)/(L) are symbolic tasks.

#### **Observation 3:**

#### Adjusting model temperature

can improve reasoning under noisy rationales.

Task	Setting	#Pr   1	rompt 2	ing E 3	xam 4	oles 5
Base-9	clean inaeasy inahard	24.8 17.5 <b>11.3</b>	38.3 22.2 <u>6.3</u>	46.4 23.2 6.0	<b>50.8</b> <u>25.4</u> 5.7	<u>50.5</u> <b>25.6</b> 5.7
Base-11	clean irr. easy irr. hard	11.8   8.9   7.7	20.4 15.9 10.0	23.9 19.1 10.7	$\frac{29.9}{21.7}$ $\frac{15.2}{1}$	32.1 26.3 16.1
Sym.(E)	clean inaeasy inahard	18.0 17.3 15.0	26.5 23.6 <u>21.0</u>	$\frac{32.7}{29.1}$ <b>22.7</b>	39.8 34.7 —	— — —
Sym.(L)	clean irr. easy irr. hard	2.7 2.3 1.9	7.7 5.4 4.0	9.3 7.0 <u>6.0</u>	$\frac{11.3}{\frac{8.8}{6.3}}$	12.2 8.9 —

Table 5: Comparing performances of the base model with a varying number of examples ("—" denotes over token limit).

#### **Observation 4:**

**Prompting with more noisy examples** boosts reasoning accuracy on most tasks.

Model	Task	0-shot	Setti  clean	ng irr.	ina.
GPT3.5	Base-9 Sym.(E) Com.	7.2 8.8 40.0	46.4 32.7 45.7	<u>30.3</u> 25.1 <u>42.3</u>	10.1 <u>26.1</u> 33.4
Gemini	Base-9 Sym.(E) Com.	12.7 9.3 42.9	88.0 44.5 55.6	$\frac{72.3}{38.9}$ $\frac{53.2}{53.2}$	21.2 36.7 33.5
Llama2	Base-9 Sym.(E) Com.	1.7 4.7 35.0	4.9 10.1 42.3	<u>2.9</u> 8.7 <u>41.9</u>	2.7 <u>9.1</u> 40.2
Mixtral	Base-9 Sym.(E) Com.	3.9 8.3 24.2	27.5 19.3 37.5	$\frac{16.3}{17.9}$ $\frac{34.9}{34.9}$	3.7 15.1 31.1

Table 6: Comparing LLMs with 0-shot, 3-shot clean, and 3-shot medium irrelevant (irr.) / inaccurate (ina.) rationales.

#### **Observation 5:**

Different LLMs are **generally vulnerable** to noisy rationales.

![](_page_46_Picture_0.jpeg)

# Empirical Evaluations with NoRa

We further explore the mapping among questions, rationales, and answers. Specifically, given the 3-shot examples  $\{(x_1, T_1, y_1), (x_2, T_2, y_2), (x_3, T_3, y_3)\}$ , we test three configurations:

- shuffle the order of **questions:**  $\{(x_2, \mathcal{T}_1, y_1), (x_3, \mathcal{T}_2, y_2), (x_1, \mathcal{T}_3, y_3)\};$
- shuffle the order of rationales:  $\{(x_1, \mathcal{T}_3, y_1), (x_2, \mathcal{T}_1, y_2), (x_3, \mathcal{T}_2, y_3)\};$
- shuffle the order of **answers:**  $\{(x_1, \mathcal{T}_1, y_3), (x_2, \mathcal{T}_2, y_1), (x_3, \mathcal{T}_3, y_2)\}$ .

Task   Zero-shot	Few-shot (No Shuffle)	Shuffle Questions $x_i \mid$ Shuffle Rationales $\mathcal{T}_i \mid$ Shuffle Answers $y_i$
Math Base-9   7.2	46.4	$\underline{45.5} (0.9\% \downarrow) \qquad   \qquad 34.5 (11.9\% \downarrow) \qquad   \qquad 35.7 (10.7\% \downarrow)$
Math Base-11   5.5	<u>23.9</u>	<b>24.8</b> (0.9% $\uparrow$ )   21.6 (2.3% $\downarrow$ )   21.1 (11.7% $\downarrow$ )
Symbolic Equal   8.8	<u>32.7</u>	<u>32.7</u> (0.0%↓)   <b>32.8</b> (0.1%↑)   32.3 (0.4%↓)
Symbolic Longer   0.0	9.2	$\underline{7.0} (2.2\% \downarrow) \qquad   \qquad 6.2 (3.0\% \downarrow) \qquad   \qquad 6.3 (2.9\% \downarrow)$
Commonsense   40.0	45.7	$38.7 (7.0\% \downarrow) \qquad   \qquad 39.7 (6.0\% \downarrow) \qquad   \qquad \underline{39.8} (5.9\% \downarrow)$

Table 7: Performance (in accuracy%) on NoRa dataset under different few-shot shuffle configurations.

**Observation 6:** Shuffling the mappings of prompting examples **degenerates** the reasoning but still performs **better** than without prompting. Besides, LLMs are **less vulnerable** to shuffled mappings than noisy rationales.

![](_page_47_Picture_0.jpeg)

#### Motivation

Current LLMs cannot denoise well with their intrinsic denoising ability.

• Even enhanced with self-correction<sup>[1]</sup> / self-consistency<sup>[2]</sup> methods.

#### **External supervision** is necessary for enhancement.

• This supervision should be sufficient for denoising and accessible in practice.

A clean CoT demonstration can be the minimal requirement for denoising-purpose prompting.

• This is more practical than existing methods requiring external supervision.

[1] J. Huang et al. Large Language Models Cannot Self-Correct Reasoning Yet. In *ICLR*, 2024.
[2] X. Wang et al. Self-consistency improves chain of thought reasoning in language models. In *ICLR*, 2023. https://bhanml.github.io & https://github.com/tmlr-group

![](_page_48_Picture_0.jpeg)

#### Motivation

#### Self-denoising:

• It is hard for LLMs to denoise noisy data without guidance.

![](_page_48_Figure_4.jpeg)

#### **Contrastive denoising:**

• It is easier for LLMs to denoise by contrasting noisy and clean data.

![](_page_48_Figure_7.jpeg)

# 

#### Method

#### Contrastive Denoising with Noisy Chain-of-thought (CD-CoT).

- **Rephrasing and selecting rationales** in the input space to conduct explicit denoising (steps 1&2).
- Exploring diverse reasoning paths and voting on answers in the output space (steps 3&4).

![](_page_49_Figure_5.jpeg)

Note. Steps 1 & 2 contribute more than Steps 3 & 4 for the explicit data denoising.

![](_page_50_Picture_0.jpeg)

- Step-1: rephrase the noisy rationales via contrastive denoising.
- Step-2: select rephrased examples with the same answers (unchanged).

![](_page_50_Figure_4.jpeg)

![](_page_51_Picture_0.jpeg)

- Step-1: rephrase the noisy rationales via contrastive denoising.
- Step-2: select rephrased examples with the same answers (unchanged).

![](_page_51_Figure_4.jpeg)

![](_page_52_Picture_0.jpeg)

- Step-3: fully utilize the rephrased examples for deliberate reasoning.
- Step-4: vote all the answers equally to get the final answer.

![](_page_52_Figure_4.jpeg)

# 

- Step-3: fully utilize the rephrased examples for deliberate reasoning.
- Step-4: vote all the answers equally to get the final answer.

![](_page_53_Figure_4.jpeg)

![](_page_54_Picture_0.jpeg)

#### Method

![](_page_54_Figure_2.jpeg)

![](_page_55_Picture_0.jpeg)

### Empirical Evaluations of CD-CoT

(besides the CoT demonstrations, the **additional information** required by the method)

Task	Method $\mathcal{M}$	Additional Information	$\mathrm{Acc}(\mathcal{M},\mathcal{Q},\mathcal{P}_{\mathrm{clean}})$	A Easy	$\operatorname{Acc}(\mathcal{M}, \mathcal{Q}, \mathcal{M})$	$\mathcal{P}_{\text{irrelevant}}$ Hard	) Avg.	A Easy	$\operatorname{Acc}(\mathcal{M}, \mathcal{Q}, \mathcal{M})$ Medium	$\mathcal{P}_{ ext{inaccurate}}$ Hard	) Avg.
Math Base-9	Base w/ SCO [29] w/ BT [81] w/ CC [9] w/ CD-CoT (ours)	Ground Truth Noise Position Clean Demo Clean Demo	46.4 53.6 47.2 44.9 <b>60.7</b>	39.3 <u>46.3</u> 39.2 43.3 <b>59.7</b>	30.3 39.6 34.2 <u>44.6</u> <b>60.7</b>	26.6 36.4 29.9 45.5 <b>57.2</b>	32.1 40.8 34.4 44.5 <b>59.2</b>	23.2 34.7 30.1 <u>37.2</u> <b>54.0</b>	10.1 22.0 18.4 <u>31.7</u> <b>58.7</b>	6.0 17.7 14.1 <u>30.7</u> <b>48.4</b>	13.1 24.8 20.9 <u>33.2</u> <b>53.7</b>
Math Base-11	Base w/ SCO [29] w/ BT [81] w/ CC [9] w/ CD-CoT (ours)	Ground Truth Noise Position Clean Demo Clean Demo	23.9 <b>33.0</b> 24.3 22.3 <u>31.0</u>	19.1 <u>29.2</u> 17.9 19.1 <b>33.7</b>	13.6 <u>24.0</u> 17.2 18.4 <b>32.7</b>	10.7 <u>20.0</u> 13.7 18.2 <b>34.7</b>	14.5 24.4 16.3 18.6 <b>33.7</b>	14.0 <b>29.2</b> 12.8 19.0 <u>29.0</u>	6.7 <u>20.0</u> 9.2 15.3 <b>30.7</b>	3.6 <u>17.2</u> 6.8 14.6 <b>25.3</b>	8.1 <u>22.1</u> 9.6 16.3 <b>28.3</b>
Symbolic Equal	Base w/ SCO [29] w/ BT [81] w/ CC [9] w/ CD-CoT (ours)	Ground Truth Noise Position Clean Demo Clean Demo	32.7 <u>38.5</u> 31.8 37.8 <b>42.7</b>	28.1 <u>34.9</u> 26.0 33.8 <b>44.7</b>	25.1 <u>33.4</u> 22.7 32.7 <b>42.7</b>	23.0 <u>32.7</u> 22.6 32.0 <b>44.0</b>	25.4 <u>33.7</u> 23.8 32.8 <b>43.8</b>	29.1 <u>34.0</u> 26.3 31.3 <b>42.6</b>	26.1 <u>34.1</u> 22.7 33.0 <b>41.3</b>	22.7 <u>34.5</u> 22.9 29.9 <b>42.7</b>	26.0 <u>34.2</u> 24.0 31.4 <b>42.2</b>
Symbolic Longer	Base w/ SCO [29] w/ BT [81] w/ CC [9] w/ CD-CoT (ours)	Ground Truth Noise Position Clean Demo Clean Demo	9.2 <b>18.7</b> 7.2 9.4 <u>12.3</u>	6.3 <b>12.1</b> 3.4 9.8 <u>12.0</u>	7.2 <u>10.5</u> 3.5 7.9 <b>12.0</b>	6.0 <u>11.3</u> 2.5 7.9 <b>13.0</b>	6.5 <u>11.3</u> <u>3.1</u> 8.5 <b>12.3</b>	7.0 <b>15.2</b> 3.8 8.5 <u>12.3</u>	6.8 <b>15.9</b> 3.6 7.4 <u>10.0</u>	6.0 <u>9.8</u> <u>3.6</u> 6.5 <b>11.0</b>	6.6 <b>13.6</b> 3.7 7.5 <u>11.1</u>
Commonsense	Base w/ SCO [29] w/ BT [81] w/ CC [9] w/ CD-CoT (ours)	Ground Truth Noise Position Clean Demo Clean Demo	45.7 63.5 47.7 48.3 49.0	44.3 60.1 23.5 45.7 50.3	42.3 56.1 28.3 43.6 54.7	41.4 60.3 32.5 44.0 50.3	42.7 <b>58.8</b> 28.1 44.4 <u>51.8</u>	36.7 56.2 11.6 42.1 51.0	33.4 58.5 11.0 40.8 49.7	28.3 <b>57.9</b> 15.8 40.5 <u>49.7</u>	32.8 57.5 12.8 41.1 50.1

Observation 7: CD-CoT presents a significant performance improvement across all datasets, with an average improvement of 17.8% compared with the base model under noisy settings.

**Observation 8:** CD-CoT displays remarkable **resistance to the magnitude of noise**, especially in the challenging mathematical tasks.

Table 8: Performance of denoising methods that require additional information for supervision.

Baseline methods:

- Self-correction with Oracle Feedback (SCO)
- https://bhanml.github.io & https://github.com/tmlr-group Backtracking (BT)
  - Contrastive CoT (CC)

56

![](_page_56_Picture_0.jpeg)

## Empirical Evaluations of CD-CoT

Model	Method	Acc(A Base-9	$\mathcal{A}, \mathcal{Q}, \mathcal{P}_{\mathrm{ir}}$ Sym.(E)	relevant) Com.	Acc(A  Base-9	$\mathcal{A}, \mathcal{Q}, \mathcal{P}_{\mathrm{in}}$ Sym.(E)	accurate) Com.
	Base	30.3	25.1	42.3	10.1	26.1	33.4
GPT-3.5-turbo	BT	34.2	28.3 22.7	$\frac{43.0}{28.3}$	17.5	22.7	$\frac{44.7}{11.0}$
_	CC	44 3	32.7	43.6	317	33.0	40 8
	CD-CoT	60.7	42.7	54.7	58.7	41.3	49.7
	Base	72.3	38.9	53.2	21.2	36.7	33.5
	SC	80.3	<u>43.3</u>	<u>60.0</u>	32.3	<u>45.0</u>	42.7
Gemini-Pro	BT	82.4	29.3	37.8	26.7	28.7	33.3
-	CC	67.5	37.3	50.2	43.6	35.0	45.6
L	CD-CoT	92.7	49.3	57.7	76.7	53.3	55.7
	Base	2.8	8.7	41.9	2.7	9.1	40.2
	SC	5.0	10.3	<b>46.7</b>	3.0	9.7	46.0
LLaMA2-70B	BT	1.4	<u>11.2</u>	36.1	0.9	<u>12.5</u>	36.2
-	CC	11	16.3	29.9	2.8	14.0	28 3
L	CD-CoT	<u>4.0</u>	9.7	<u>39.3</u>	2.7	9.7	39.7
	Base	16.3	17.9	34.9	3.7	15.1	31.1
	SC	20.0	21.7	37.0	2.7	18.0	37.7
Mixtral-8x7B	BT	4.1	9.7	6.2	2.4	10.1	10.5
_	CC	24.4	18 5	36.0	12.5	<u>18 3</u>	35 7
Ĺ	CD-CoT	8.7	22.7	40.3	<u>4.7</u>	21.3	40.3

Table 11: Comparing methods with different LLMs.

**Observation 9:** CD-CoT **generalizes well** across different LLMs.

Performance Gains

- GPT-3.5-turbo: **+23.4%** accuracy.
- Gemini-Pro: **+21.6%** accuracy Both surpass all baselines.

![](_page_57_Picture_0.jpeg)

# Empirical Evaluations of CD-CoT

#### Comparison of rephrased rationales by different reasoning methods.

Original Noisy Rationale	Method	Denoised Rationale		In base-9, the digits range from 0 to 8. Let's
In base-9, the digits are "012345678". $9 + 8 = 17$ . We have $6 + 7 = 13$ in base-10. Since we're in base-9, that exceeds the maximum value of 8 for a single digit. 13 mod $9 = 4$ , so the digit is 4 and the carry is 1. We have 8 $+ 5 + 1 = 14$ in base-10. $14 + 8 = 22$ . 14 mod $9 = 5$ , so the digit is 5 and the carry is 1. A leading digit is 1. So the answer is 154. Answer: 154.	<u>,</u> SM	In base-9, the dkgits are "012345678". $9 + 8$ = 17. We5hav6r+ 7l=J13 in base-10. Since we're in baze-9, ithFt exceeis the maCimum valur of 8 for a siLgle]digit. x0c3 mod 9 = 2 soWghe digix is 4 anj the carry is 1. We hOve 8 + 5i+ z2= 14 Vn baAe 10. 14 d 8 = 2p. 14 mod 9 = 5, so Xhe diGit*iU 5 and the carry is 1Q A leading x0cigit is 1. So the answer is 154. Asswer: 154.	SP	add 86 and 57 step by step. First, we add the rightmost digits, which are 6 and 7. In base-10, 6 + 7 equals 13. However, since we are in base-9, the maximum value for a single digit is 8. Therefore, we need to find the remainder when 13 is divided by 9. The remainder is 4, so the digit in the units place is 4 and we have a carry of 1. Next, we add the next digits, which are 8 and 5, along with the carry of 1. In base-10, $8 + 5 + 1$ equals 14. Again, we need to find the remainder
Ours In base-9, the digits are "012345678". We have $6 + 7 = 13$ in base-10. Since we're in base-9, that exceeds the maximum value of 8 for a single digit. 13 mod $9 = 4$ , so the digit is 4 and the carry is 1. We have $8 + 5 + 1 = 14$ in base 10. 14 mod $9 = 5$ , so the digit is 5 and the carry is 1. A leading digit is 1. So the answer is 154. Answer: 154.	SD	In base-9, the digits are "012345678". $9+8$ = 17. We have $6+7=13$ in base-10. Since we're in base-9, 13 exceeds the maximum value of 8 for a single digit. 13 mod $9=4$ , so the digit is 4 and the carry is 1. We have 8+5+1=14 in base 10. 14 mod $9=8$ . 14 div $9=1$ , so the digit is 8 and the carry is 1. A leading digit is 1. So the answer is 154. Answer: 154.	Pacal	when 14 is divided by 9. The remainder is 5, so the digit in the tens place is 5 and we have a carry of 1. Finally, we have a leading digit of 1. So the final answer in base-9 is 154. Answer: 154.
			Basel	Randomly masks the prompt.

**Observation 10:** CD-CoT effectively **removes noisy thoughts** and ensures format alignment with the original rationale.

SM: Randomly masks the prompt.SD: Applies random masking and reconstructs prompts through FM.SP: Reconstructs prompts following guidelines through FM.

![](_page_58_Picture_0.jpeg)

#### Take Home Messages

We investigate the **under-explored** problem of noisy rationales.

We introduce **NoRa dataset** to evaluate LLMs against noisy rationales.

We reveal the **general vulnerability** of LLMs to noisy rationales; this is not well addressed by existing robust methods.

We design **CD-CoT** method to enhance the robustness via contrastive denoising.

![](_page_59_Picture_0.jpeg)

#### Future Directions

#### Robust pre-training/fine-tuning methods are required for VLMs.

- VLMs can still be mislead by spurious features.
- Larger models and high-quality data lead to better robustness.

#### The trade-off between unlearning and retention remains a critical issue.

- Current unlearning objectives all have negative impacts on retention.
- Data and optimization aspects of unlearning are not well explored.

#### Reasoning with noisy rationales can be further investigated.

- Non-reasoning models (GPT 3.5/4/40) is not robust on the NoRa dataset.
- Reasoning models R1/o1/o3 is generally more robust but exhibit over-thinking issues.

# Appendix

![](_page_60_Picture_1.jpeg)

- Survey:
  - A Survey of Label-noise Representation Learning: Past, Present and Future. arXiv, 2020.

#### • Book:

- Machine Learning with Noisy Labels: From Theory to Heuristics. Adaptive Computation and Machine Learning series, The MIT Press, 2025.
- Trustworthy Machine Learning under Imperfect Data. CS series, **Springer Nature**, 2025.
- Trustworthy Machine Learning: From Data to Models. Foundations and Trends® in Privacy and Security, 2025.

![](_page_60_Picture_8.jpeg)

#### • Tutorial:

- IJCAI 2021 Tutorial on Learning with Noisy Supervision
- CIKM 2022 Tutorial on Learning and Mining with Noisy Labels
- ACML 2023 Tutorial on Trustworthy Learning under Imperfect Data
- AAAI 2024 Tutorial on Trustworthy Machine Learning under Imperfect Data
- IJCAI 2024 Tutorial on Trustworthy Machine Learning under Imperfect Data
- WWW 2025 Tutorial on Trustworthy AI under Imperfect Web Data

#### Workshops:

- IJCAI 2021 Workshop on Weakly Supervised Representation Learning
- ACML 2022 Workshop on Weakly Supervised Learning
- RIKEN 2023 Workshop on Weakly Supervised Learning
- HKBU-RIKEN AIP 2024 Joint Workshop on Artificial Intelligence and Machine Learning

### What is TMLR Group

![](_page_61_Picture_1.jpeg)

- TMLR Group, an online-offline-mixed **machine learning** research group, locates in different cities, including Hong Kong, Melbourne, Shanghai, Nottingham and Sydney.
- We are welcoming the **synergetic collaboration** between yours and HKBU TMLR!!

	TMLR Group Trustworthy Machine Learning and Reasoning Group						
	Rर 118 followers 💿 Hong Kong 🔗 https://bhanml.github.io/group.html 🗹 tmlr.group@gmail.com						
README.md						Ø	
Trustworthy N different cities building trusty	lachine Learning a , including Hong I worthy learning an	nd Reasoning (T Kong, Melbourne nd reasoning algo	MLR) Group, an 9, Shanghai, Noti prithms, theories	online-offline-mixed m tingham and Sydney. W and systems.	nachine learning research group, locates Ve share the vision for the future ML tec	in hnology:	
inned						Customize pi	
G-effect Public ::				AttrVR P	ublic		
Forked from <u>QizhouWang/G-effect</u>				Forked from caic	Forked from caichengyi/AttrVR		
[ICLR 2025] "Rethinking LLM Unlearning Objectives: A Gradient Perspective and Go Beyond"				[ICLR 2025] "Attr Models" Official	[ICLR 2025] "Attribute-based Visual Reprogramming for Vision-Language Models" Official Website: https://github.com/tmlr-group/AttrVR		
● Python 🟠	11			● Python 🏠	2		
📮 NoisyRatio	nales Public			📮 BayesianLl	M Public		
[NeurIPS 2024] "Can Language Models Perform Robust Reasoning in Chain- of-thought Prompting with Noisy Rationales?"			Forked from caic	Forked from caichengyi/BayesianLM			
			[NeurIPS 2024 O Reprogramming	[NeurIPS 2024 Oral] "Bayesian-Guided Label Mapping for Visual Reprogramming"			
Python 1 Stress	35 <b>%</b> 2			Python	9 😵 1		
EOE Public				WCA Puk	blic		
Forked from Aboriginer/EOE			Forked from Jinh	Forked from JinhaoLee/WCA			
[ICML 2024] "Envisioning Outlier Exposure by Large Language Models for Out-of-Distribution Detection"			[ICML 2024] "Vis	[ICML 2024] "Visual-Text Cross Alignment: Refining the Similarity Score in Vision-Language Models"			
Out-of-Distributio	n Detection"			Vision-Language	Models"		

- Research Twitter:
  - <u>https://x.com/tmlrgroup</u>
- Research RedNote:
  - <u>https://www.xiaohongshu.com/user/</u> profile/646ee4b90000000110010b6
- Research Blog:
  - <u>https://www.jiqizhixin.com/columns/</u> <u>TMLRGroup</u>

![](_page_62_Picture_0.jpeg)

#### Focused Areas

![](_page_62_Figure_2.jpeg)