

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

Zhanke Zhou

Hong Kong Baptist University

with Rong Tao, Jianing Zhu, Yiwon Luo, Zengmao Wang, and Bo Han

# Main contributions

## New research problem: Noisy Rationales

We investigate the problem of noisy rationales in the prevailing chain-of-thought prompting →

**Input with Noisy Questions**

**Question-1 (Q1):** In base-9, what is  $86+57$ ?  
We know  $6+6=12$  and  $3+7=10$  in base 10.

**Rationale-1 (R1):** 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.  $13 \bmod 9 = 4$ , so the digit is 4 and the carry is 1. We have  $8 + 5 + 1 = 14$  in base 10.  $14 \bmod 9 = 5$ , so the digit is 5 and the carry is 1. A leading digit 1. So the answer is 154.

**Answer-1 (A1):** 154.

...Q2, R2, A2, Q3, R3, A3...

**Test Question:** In base-9, what is  $62+58$ ?  
We know  $6+6=12$  and  $3+7=10$  in base 10.

**Input with Noisy Rationales**

**Question-1 (Q1):** In base-9, what is  $86+57$ ?  
**Rationale-1 (R1):** In base-9, the digits are "012345678". We have  $6 + 7 = 13$  in base-10.  $13 + 8 = 21$ . Since we're in base-9, that exceeds the maximum value of 8 for a single digit.  $13 \bmod 9 = 4$ , so the digit is 4 and the carry is 1. We have  $8 + 5 + 1 = 14$  in base 10.  $14 \bmod 9 = 5$ , so the digit is 5 and the carry is 1.  $5 + 9 = 14$ . A leading digit is 1. So the answer is 154.

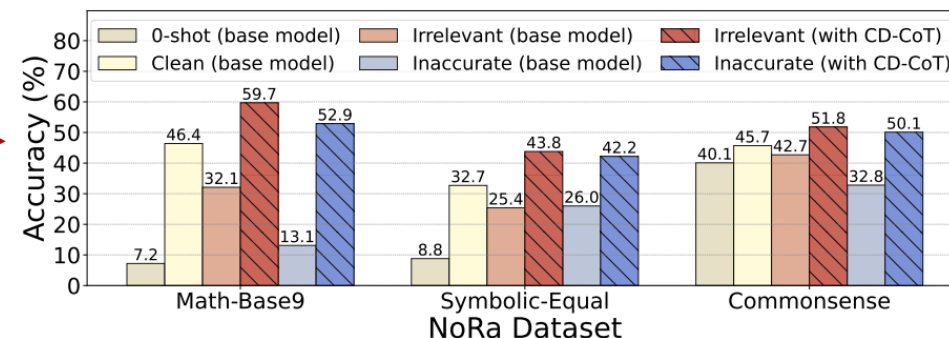
**Answer-1 (A1):** 154.

...Q2, **R2**, A2, Q3, **R3**, A3 ...

**Test Question:** In base-9, what is  $62+58$ ?

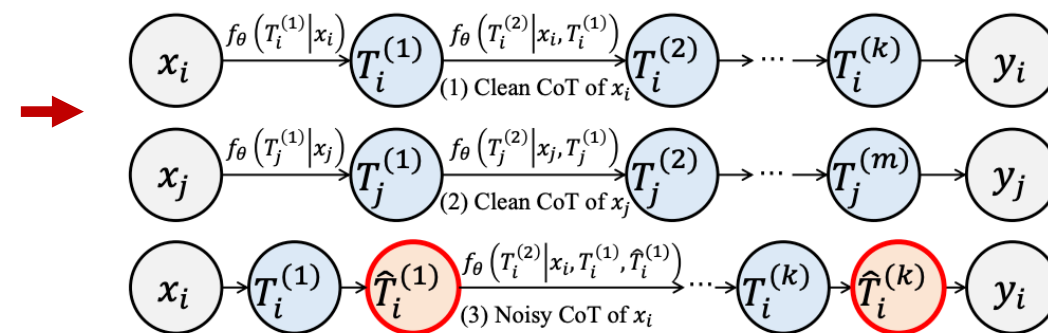
## New benchmark: NoRa

We construct the NoRa dataset and systematically evaluate the robustness of LLMs →



## New algorithm: CD-CoT

We design a simple yet effective method to enhance robustness via contrastive denoising →



# Outline

- Background: language model reasoning
- **New research problem: Noisy Rationales**
- **New benchmark: NoRa**
- **New algorithm: CD-CoT**
- Take home messages
- Future directions

# Background: language model reasoning

In-context learning (ICL) is commonly used in large language models (LLMs)

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

## Zero-shot Input

Question: In base-9, what is  $62+58$ ?

## Input: ICL with three examples

Question-1: In base-9, what is  $86+57$ ? Answer-1: 154.

Question-2: In base-9, what is  $63+34$ ? Answer-2: 107.

Question-3: In base-9, what is  $31+58$ ? Answer-3: 100.

Question: In base-9, what is  $62+58$ ?



# Background: language model reasoning

In-context learning (ICL) is commonly used in large language models (LLMs)

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

## Zero-shot Input

Question: In base-9, what is  $62+58$ ?

## Input: ICL with three examples

Question-1: In base-9, what is  $86+57$ ? Answer-1: 154.

Question-2: In base-9, what is  $63+34$ ? Answer-2: 107.

Question-3: In base-9, what is  $31+58$ ? Answer-3: 100.

Question: In base-9, what is  $62+58$ ?

Prevailing in ICL, **Chain of thought (CoT) prompting** boost model reasoning

- CoT includes **rationales**, i.e., sequential reasoning thoughts to solve a question

## Input: ICL with three examples

Question-1: In base-9, what is  $86+57$ ? Answer-1: 154.

Question-2: In base-9, what is  $63+34$ ? Answer-2: 107.

Question-3: In base-9, what is  $31+58$ ? Answer-3: 100.

Question: In base-9, what is  $62+58$ ?

## Input: CoT 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.  $13 \bmod 9 = 4$ , so the digit is 4 and the carry is 1. We have  $8 + 5 + 1 = 14$  in base 10.  $14 \bmod 9 = 5$ , 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 research problem: Noisy Rationales

Existing work generally assume that CoT contains clean rationales

But, what if CoT contains **noisy rationales**? 🤔

- **noisy rationales include irrelevant or inaccurate thoughts**

Input: CoT with **clean 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.  $13 \bmod 9 = 4$ , so the digit is 4 and the carry is 1. We have  $8 + 5 + 1 = 14$  in base 10.  $14 \bmod 9 = 5$ , 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  $62+58$ ?

the irrelevant **base-10 information** is included in rationale

Input: CoT with **noisy rationales**

Question-1 (Q1): In base-9, what is  $86+57$ ?

Rationale-1 (R1): In base-9, the digits are "012345678". We have  $6 + 7 = 13$  in base-10.  **$13 + 8 = 21$** . Since we're in base-9, that exceeds the maximum value of 8 for a single digit.  $13 \bmod 9 = 4$ , so the digit is 4 and the carry is 1. We have  $8 + 5 + 1 = 14$  in base 10.  $14 \bmod 9 = 5$ , so the digit is 5 and the carry is 1.  **$5 + 9 = 14$** . A leading digit is 1. So the answer is 154.

Answer-1 (A1): 154.

...Q2, R2, A2, Q3, R3, A3 ...

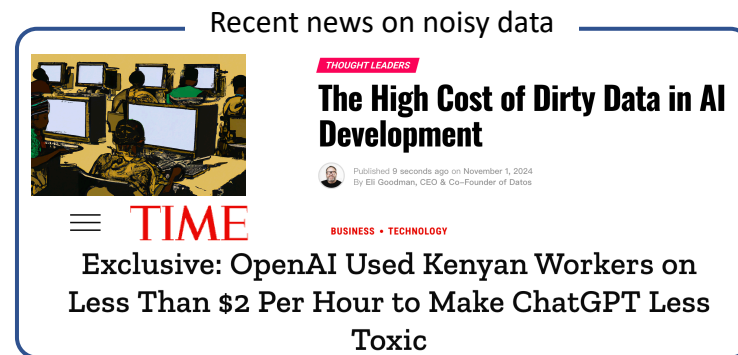
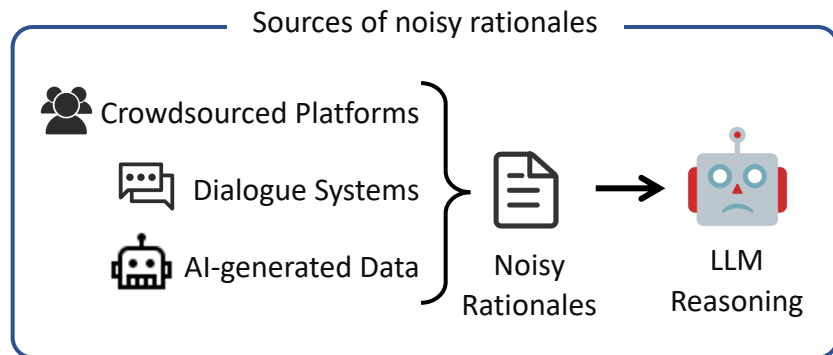
Test Question: In base-9, what is  $62+58$ ?

while the test question asks about **base-9 calculation**

# New research problem: Noisy Rationales

**Noisy rationales originate from diverse sources** (see Appendix C for details)

- 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
- and verify the corresponding countermeasures against noisy rationales

# Outline

- Background: language model reasoning
- New research problem: Noisy Rationales
- **New benchmark: NoRa**
  - **Benchmark construction**
  - Empirical evaluations on NoRa
- New algorithm: CD-CoT
- Take home messages
- Future directions

# New benchmark: NoRa

## NoRa (Noisy Rationales)

- a comprehensive testbed to evaluate the **robustness** against noisy rationales
- contains **26391** questions and **5** subtasks
- covering **3** types of reasoning tasks: **mathematical, symbolic, and commonsense**

Task	Irrelevant Thoughts	Inaccurate Thoughts
NoRa-Math	In base-9, digits run from 0 to 8. We have $3 + 2 = 5$ in base-10. Since we're in base-9, that doesn't exceed the maximum value of 8 for a single digit. $5 \bmod 9 = 5$ , so the digit is 5 and the carry is 0. <u>There are five oceans on Earth: the Atlantic, Pacific, Indian, Arctic, and Southern.</u> We have $8 + 6 + 0 = 14$ in base 10. $14 \bmod 9 = 5$ , so the digit is 5 and the carry is 1. A leading digit 1. So the answer is 155. Answer: 155	In base-9, digits run from 0 to 8. We have $3 + 2 = 5$ in base-10. <u><math>5 + 4 = 9</math>.</u> Since we're in base-9, that doesn't exceed the maximum value of 8 for a single digit. $5 \bmod 9 = 5$ , so the digit is 5 and the carry is 0. <u><math>5 + 9 = 14</math>.</u> We have $8 + 6 + 0 = 14$ in base 10. $14 \bmod 9 = 5$ , so the digit is 5 and the carry is 1. A leading digit 1. So the answer is 155. Answer: 155
NoRa-Symbolic	... "turn around right" means the agent needs to turn right, and repeat this action sequence four times to complete a 360-degree loop. <u>Many GPS navigation systems will issue a 'turn around' command if the driver deviates from the planned route.</u> So, in action sequence is I_TURN_RIGHT I_TURN_RIGHT I_TURN_RIGHT I_TURN_RIGHT. ...	... "turn around right" means the agent needs to turn right, and repeat this action sequence four times to complete a 360-degree loop. <u>Turn opposite is I_TURN_RIGHT I_TURN_LEFT.</u> So, in action sequence is I_TURN_RIGHT I_TURN_RIGHT I_TURN_RIGHT I_TURN_RIGHT. ...
NoRa-Com.	The relations path are son, sister, uncle, which means Francisco is David's son's sister's uncle. For son's sister, we have son's sister is daughter. So the relations path are reduced to daughter, uncle. <u>In genetics, mitochondrial DNA is always inherited from the mother, making the mother-daughter genetic link unique.</u> For daughter's uncle, we have daughter's uncle is brother. So the relations path are reduced to brother. Therefore, the answer is brother. Answer:brother	The relations path are son, sister, uncle, which means Francisco is David's son's sister's uncle. For son's sister, we have son's sister is daughter. So the relations path are reduced to daughter, uncle. For daughter's uncle, we have daughter's uncle is brother. <u>We have brother' sister is brother.</u> So the relations path are reduced to brother. Therefore, the answer is brother. Answer:brother

Table 1: Noisy rationales (consisting noisy thoughts) 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.

# New benchmark: NoRa

Difficulty	Noise Ratio	#total thoughts (#noisy thoughts) of prompting rationales (Avg.)				
		Math Base-9	Math Base-11	Sym. Equal	Sym. Longer	Com.
Easy	0.3	10 (2)	10 (2)	11.5 (2.7)	11.0 (2.5)	7 (2)
Medium	0.5	12 (4)	12 (4)	13.3 (4.5)	12.7 (4.2)	8 (3)
Hard	0.8	14 (6)	14 (6)	16.0 (7.1)	15.2 (6.8)	9 (4)
#questions		4024	9269	4182	3920	4996

Table 2: Statistics of NoRa dataset.

## Definitions

- **Irrelevant thoughts** are irrelevant to the given context
  - e.g., discussing the genetic overlap of siblings when the task is to deduce 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 the reasoning difficulty through **different noise ratios (0.3, 0.5, 0.8)**

# Outline

- Background: language model reasoning
- New research problem: Noisy Rationales
- **New benchmark: NoRa**
  - Benchmark construction
  - **Empirical evaluations on NoRa**
- New algorithm: CD-CoT
- Take home messages
- Future directions

# Empirical evaluations on NoRa

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

- a **0.2%-25.3%** decrease with irrelevant noise
- a **0.1%-54.0%** decrease with inaccurate noise (compared with clean rationales)

**Observation 1:**  
**self-correction**  
**methods** perform  
poorly on most tasks  
with noisy rationales

**Observation 2:**  
**self-consistency**  
**methods** can improve  
robustness without  
true denoising

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

# Empirical evaluations on NoRa

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

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

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

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

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

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

## Observation 3:

Adjusting model temperature can improve reasoning under noisy rationales

## Observation 4:

Prompting with more noisy examples boosts reasoning accuracy on most tasks

## Observation 5:

Different LLMs are generally vulnerable to noisy rationales

# Empirical evaluations on NoRa

## We further explore the mapping among questions, rationales, and answers

Specifically, given the 3-shot examples  $\{(x_1, \mathcal{T}_1, y_1), (x_2, \mathcal{T}_2, y_2), (x_3, \mathcal{T}_3, y_3)\}$ , we test three configurations:

- shuffle questions  $\{(\mathbf{x}_1, \mathcal{T}_3, y_3), (\mathbf{x}_2, \mathcal{T}_1, y_1), (\mathbf{x}_3, \mathcal{T}_2, y_2)\}$
- shuffle rationales  $\{(x_1, \mathbf{\mathcal{T}}_3, y_1), (x_2, \mathbf{\mathcal{T}}_1, y_2), (x_3, \mathbf{\mathcal{T}}_2, y_3)\}$
- shuffle answers  $\{(x_1, \mathcal{T}_1, \mathbf{y}_3), (x_2, \mathcal{T}_2, \mathbf{y}_1), (x_3, \mathcal{T}_3, \mathbf{y}_2)\}$

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

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.

# Outline

- Background: language model reasoning
- New research problem: Noisy Rationales
- New benchmark: NoRa
- **New algorithm: CD-CoT**
  - **Motivation and design of CD-CoT**
  - Empirical evaluations of CD-CoT
- Take home messages
- Future directions

# New algorithm: CD-CoT

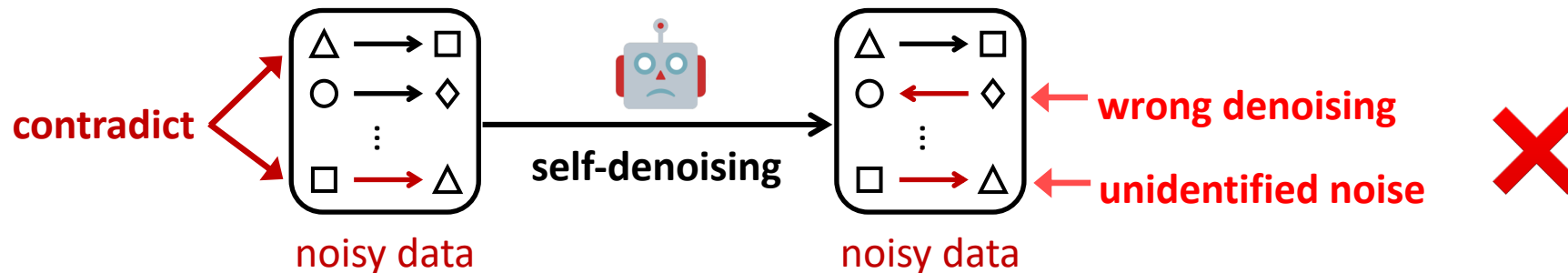
## Motivation

- Current LLMs **cannot** denoise well with their **intrinsic denoising ability**
  - even enhanced with self-correction / self-consistency methods
- **External supervision** is necessary for enhancement
  - which should be sufficient for denoising and accessible in practice
- **A clean CoT demonstration** can be the minimal requirement
  - for denoising-purpose prompting
  - which is much more practical than existing methods requiring external supervision

# New algorithm: CD-CoT

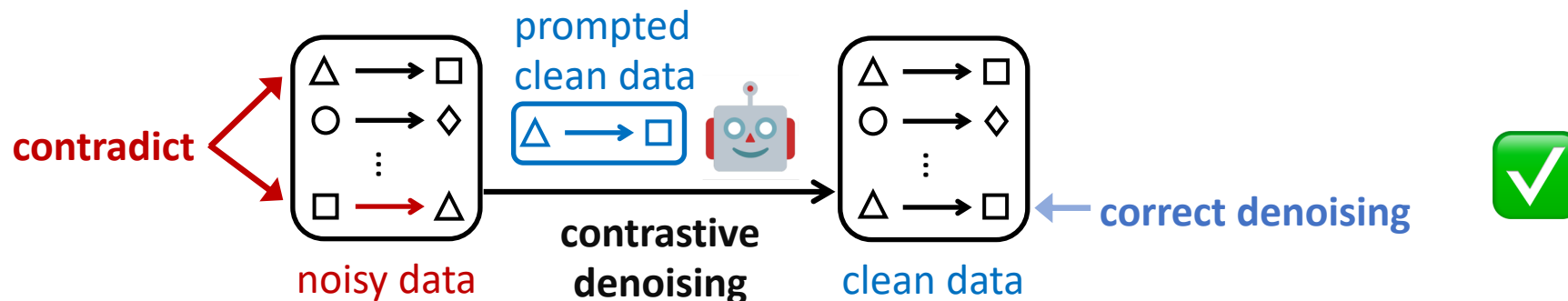
## Self-denoising:

- It is **hard** for LLMs to denoise noisy data **without guidance**



## Contrastive denoising:

- It is **easier** for LLMs to denoise **by contrasting noisy and clean data**



# New algorithm: CD-CoT

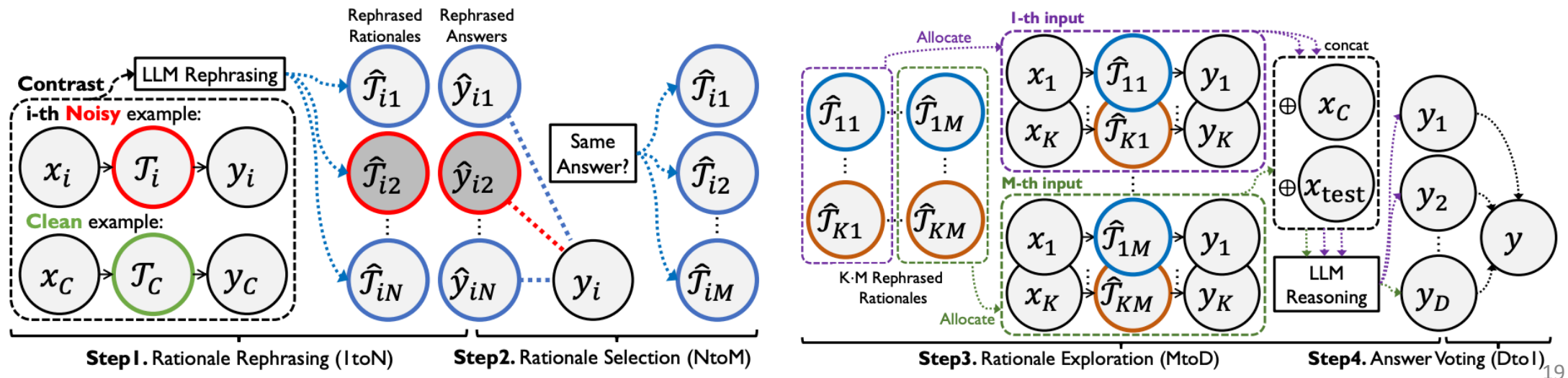
## **Contrastive Denoising with Noisy Chain-of-thought (CD-CoT)**

- assume that LLMs can identify noisy thoughts
  - by contrasting a pair of noisy and clean rationales (similar to contrastive learning)

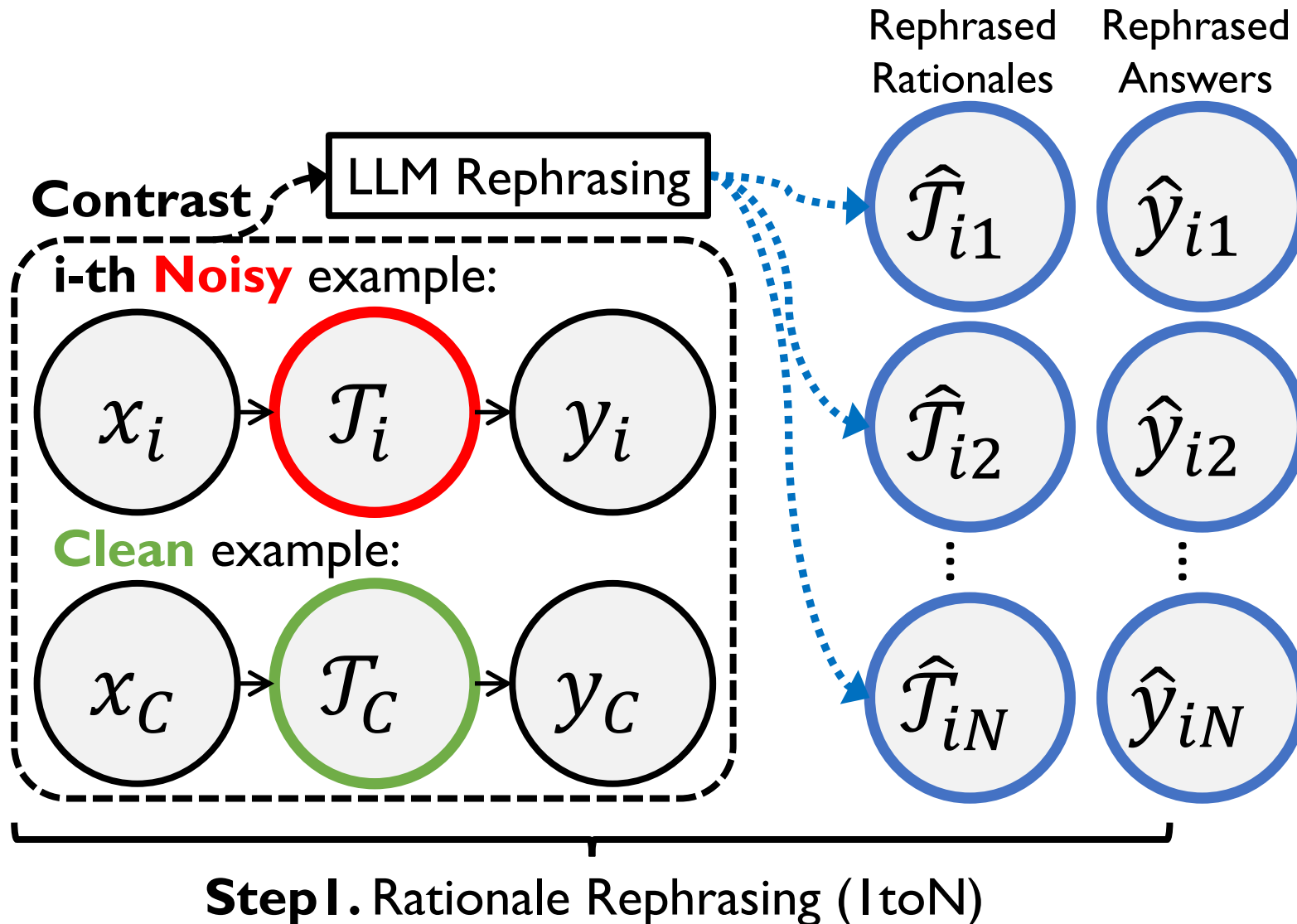
# New algorithm: CD-CoT

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

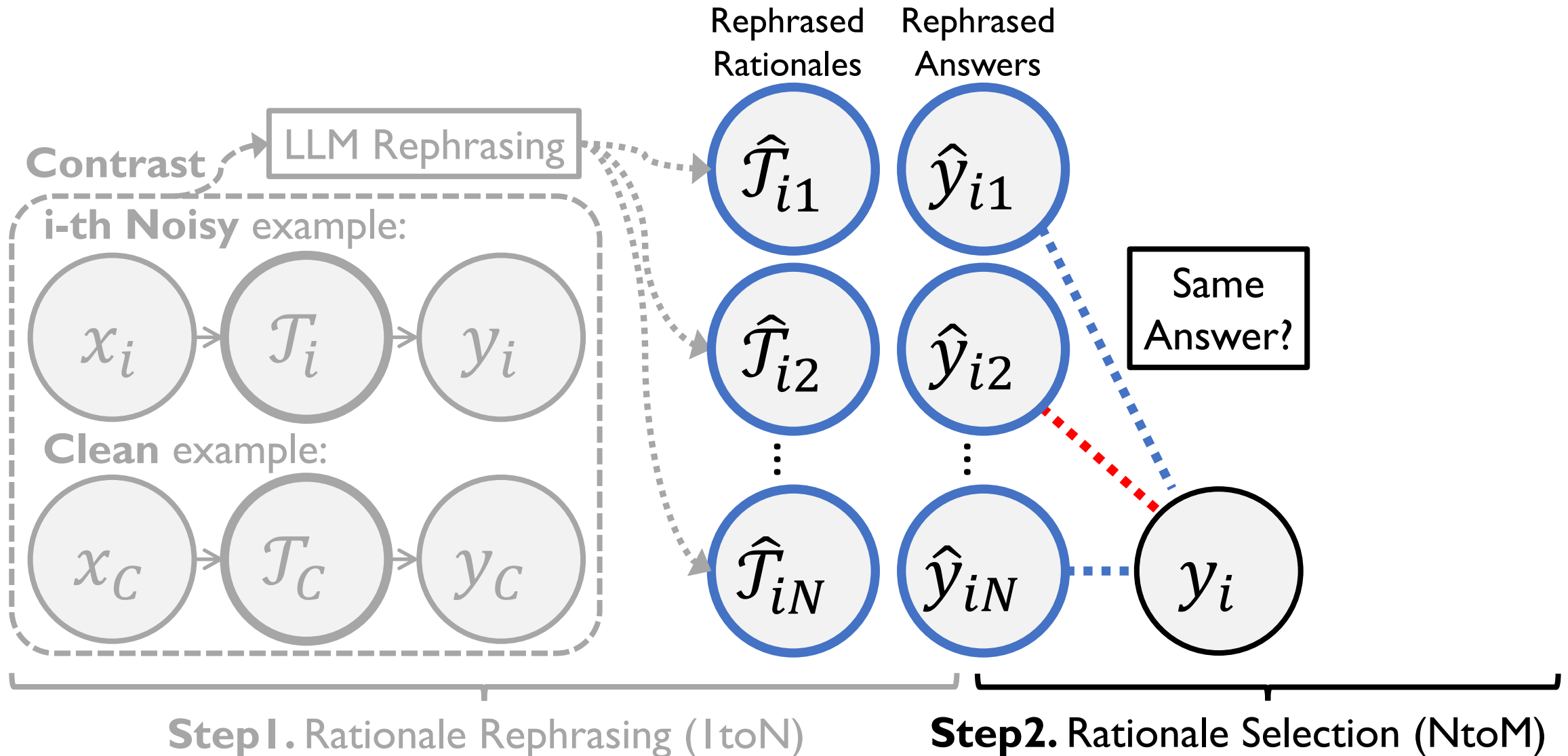
- assume that LLMs can identify noisy thoughts
  - by contrasting a pair of noisy and clean rationales (similar to contrastive learning)
- **design principle: exploration and exploitation**
  - **rephrasing and selecting rationales** in the input space to conduct explicit denoising (steps 1&2)
  - **exploring diverse reasoning paths and voting on answers** in the output space (steps 3&4)



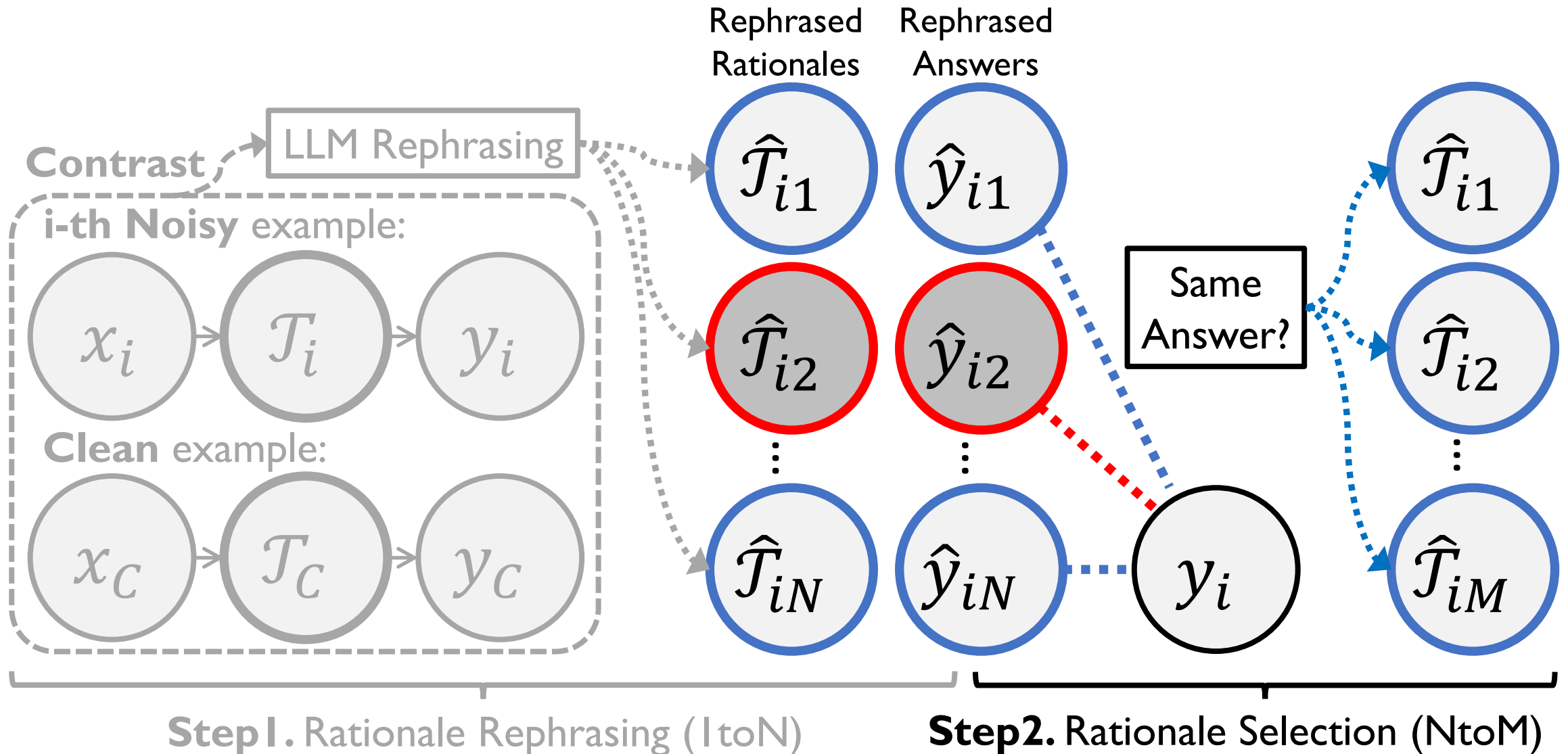
- **Step-1:** rephrase the noisy rationales via contrastive denoising
- Step-2: select rephrased examples with the same answers (unchanged)



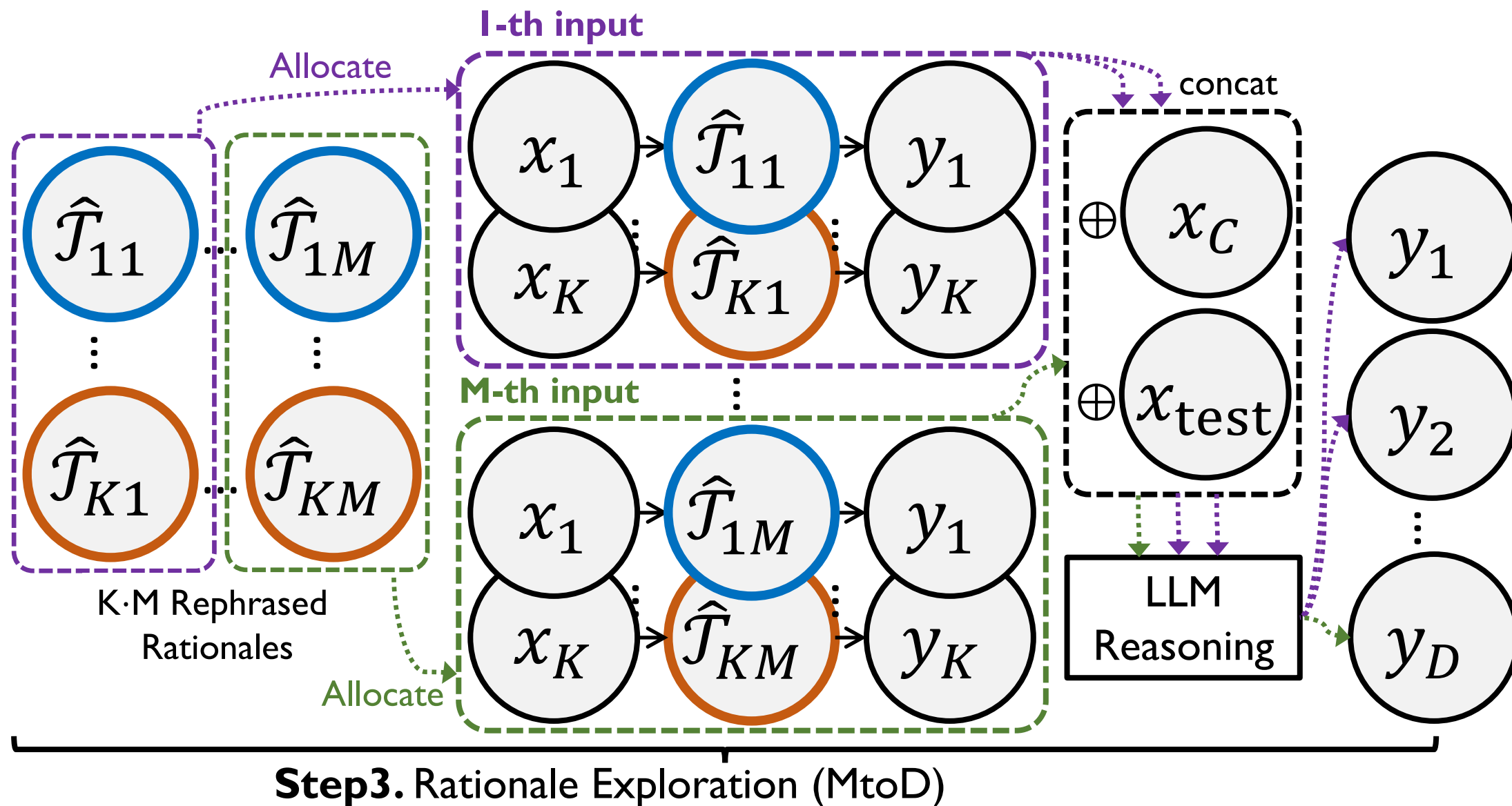
- Step-1: rephrase the noisy rationales via contrastive denoising
- **Step-2: select rephrased examples with the same answers (unchanged)**



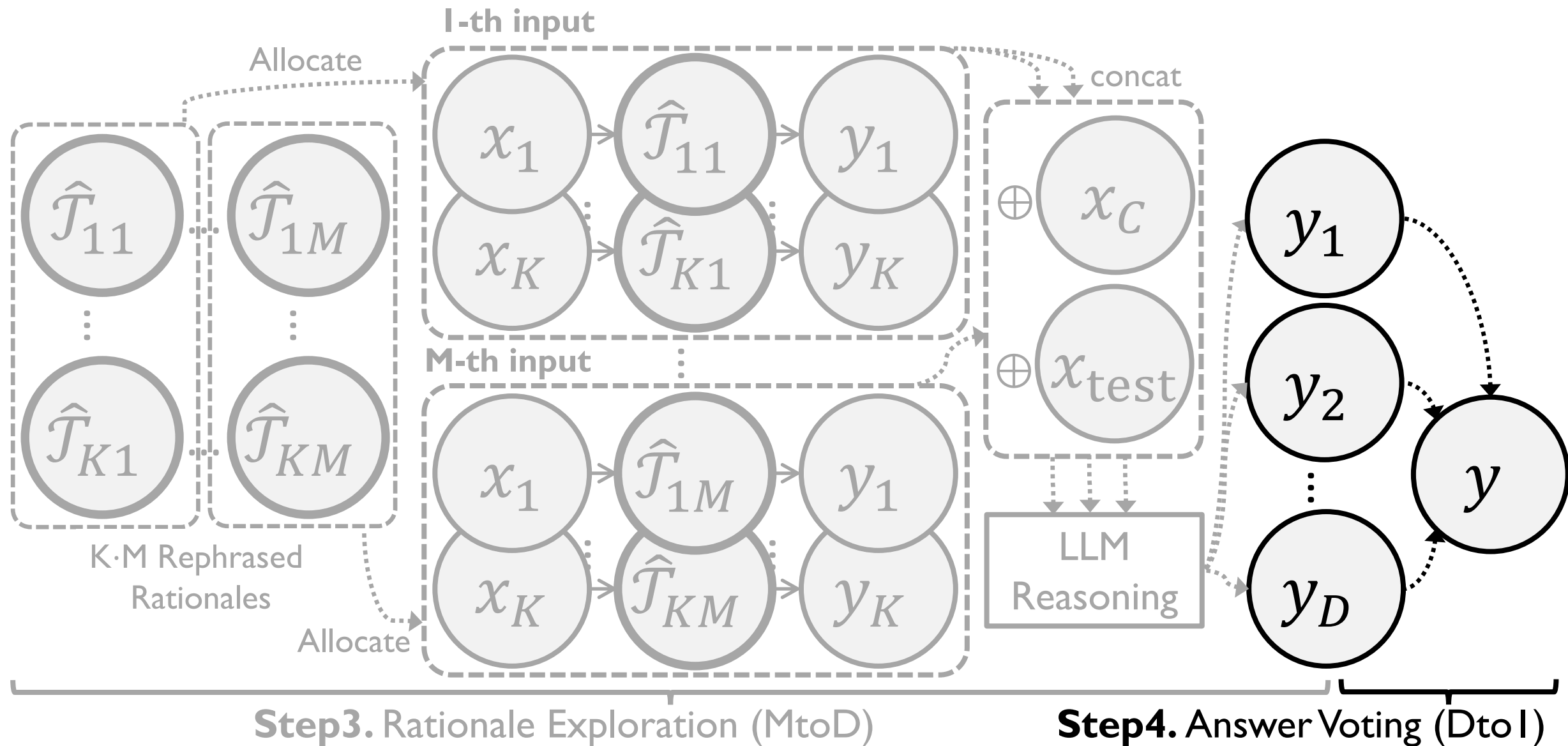
- Step-1: rephrase the noisy rationales via contrastive denoising
- **Step-2: select rephrased examples with the same answers (unchanged)**



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



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

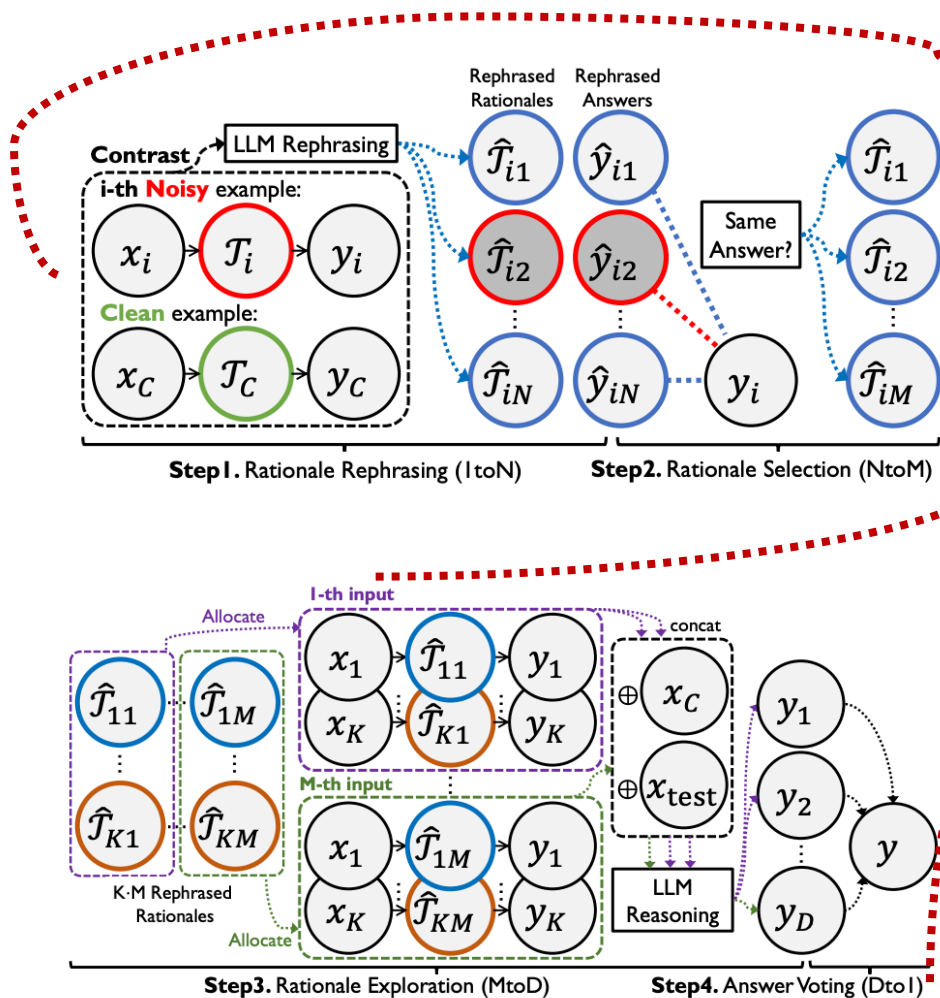


# New algorithm

## Algorithm 1 CD-CoT: Contrastive Denoising with Noisy Chain-of-Thought.

**Require:** an LLM  $f_\theta$ , the prompt of contrastive denoising  $\mathcal{P}_{\text{denoise}}$ , one test question  $x_{\text{test}}$ , one clean example  $(x_C, \mathcal{T}_C, y_C)$ ,  $K$  prompting examples  $S_n = \{(x_i, \mathcal{T}_i, y_i)\}_{i=1}^K$ , hyper-parameters  $N, M$ , and reasoning budget  $\{B_i\}_{i=1}^M$  (satisfies that  $\sum_{i=1}^M B_i = D$ , where  $D$  is the total budget).

- 1: **for**  $i = 1 \dots K$  **do**
- 2:   initialize the set of rephrased results of  $i$ -th example  $\mathcal{R}_i \leftarrow \emptyset$ .
- 3:   **for**  $j = 1 \dots N$  **do**
- 4:     **# Step-1: Rationale Rephrasing via Supervised Contrasting**
- 5:     obtain a rephrased example as  $(x_i, \hat{\mathcal{T}}_i, \hat{y}_i) \leftarrow f_\theta(\mathcal{P}_{\text{denoise}}(x_C, \mathcal{T}_C, y_C, x_i, \mathcal{T}_i, y_i))$ .
- 6:     if match answer  $\hat{y}_i = y_i$ , then store the rephrased example as  $\mathcal{R}_i \leftarrow \mathcal{R}_i \cup \{(x_i, \hat{\mathcal{T}}_i, \hat{y}_i)\}$ .
- 7:   **end for**
- 8:   **# Step-2: Rationale Selection**
- 9:   randomly select  $M$  rephrased examples from  $\mathcal{R}_i$  and obtain  $\tilde{\mathcal{R}}_i = \{(x_{is}, \hat{\mathcal{T}}_{is}, \hat{y}_{is})\}_{s=1}^M$ .
- 10: **end for**
- 11: **# Step-3: Rationale Exploration**
- 12: initialize the set of answers  $\mathcal{Y} \leftarrow \emptyset$ .
- 13: **for**  $i = 1 \dots M$  **do**
- 14:   construct an input  $\mathcal{P}_i \leftarrow \{(x_{ji}, \hat{\mathcal{T}}_{ji}, \hat{y}_{ji})\}_{j=1}^K$ , where  $(x_{ji}, \hat{\mathcal{T}}_{ji}, \hat{y}_{ji})$  is the  $i$ -th element of  $\tilde{\mathcal{R}}_j$ .
- 15:   concatenate  $\mathcal{P}_i$  with the clean example and test question as  $\mathcal{P}_i \leftarrow \mathcal{P}_i \cup \{(x_C, \mathcal{T}_C, y_C), x_{\text{test}}\}$ .
- 16:   **for**  $j = 1 \dots B_M$  **do**
- 17:     get one answer by LLM reasoning as  $y_j \leftarrow f_\theta(\mathcal{P}_i)$ .
- 18:     store the answer as  $\mathcal{Y} \leftarrow \mathcal{Y} \cup \{y_j\}$ .
- 19:   **end for**
- 20: **end for**
- 21: **# Step-4: Answer Voting**
- 22: initialize the dictionary of answer count  $\mathcal{C}$  that  $\forall y_j \in \mathcal{Y}, \mathcal{C}[y_j] = 0$ .
- 23: **for**  $j = 1 \dots D$  **do**
- 24:   update  $\mathcal{C}[y_j] \leftarrow (\mathcal{C}[y_j] + 1)$ .
- 25: **end for**
- 26: get the final answer  $y$  with maximum counts as  $y \leftarrow \arg \max_y \mathcal{C}[y]$ .
- 27: **return** the answer  $y$ .



# Outline

- Background: language model reasoning
- New research problem: Noisy Rationales
- New benchmark: NoRa
- **New algorithm: CD-CoT**
  - Motivation and design of CD-CoT
  - **Empirical evaluations of CD-CoT**
- Take home messages
- Future directions

# Empirical evaluations of CD-CoT

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

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

**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.

# Empirical evaluations of CD-CoT

Hyper-parameters				Acc( $\mathcal{M}, \mathcal{Q}, \mathcal{P}_{\text{irrelevant}}$ )			Acc( $\mathcal{M}, \mathcal{Q}, \mathcal{P}_{\text{inaccurate}}$ )		
$N$	$M$	$D$	$C$	Base-9	Sym.(E)	Com.	Base-9	Sym.(E)	Com.
5	1	5	Y	57.7	38.7	55.3	53.3	39.7	51.0
5	1	5	N	54.7	32.7	53.7	47.0	32.3	<b>55.7</b>
5	2	2+3	Y	<b>60.7</b>	<b>42.7</b>	54.7	<b>58.7</b>	41.3	49.7
5	2	2+3	N	56.7	33.0	54.7	49.7	32.0	53.0
5	3	1+2+2	Y	<b>60.7</b>	38.7	53.3	58.0	<b>43.3</b>	49.0
5	3	1+2+2	N	56.0	33.3	55.7	48.7	32.0	52.3
5	5	1	Y	59.3	39.7	55.7	58.0	39.0	48.7
5	5	1	N	55.3	35.7	<b>55.9</b>	48.7	33.3	50.7

Table 9: Comparison of accuracy on medium-level tasks.

Hyper-parameters				#Tokens in step-3 (irr.)			#Tokens in step-3 (ina.)		
$N$	$M$	$D$	$C$	Base-9	Sym.(E)	Com.	Base-9	Sym.(E)	Com.
5	1	5	Y	1440	3162	788	1428	3170	798
5	1	5	N	1301	2685	660	1295	2732	667
5	2	2+3	Y	2175	4934	1269	2156	4989	1311
5	2	2+3	N	1864	4044	1005	1842	4087	1039
5	3	1+2+2	Y	2902	6704	1772	2878	6785	1821
5	3	1+2+2	N	2416	5360	1372	2393	5443	1420
5	5	1	Y	4368	10340	2764	4339	10514	2845
5	5	1	N	3535	8099	2088	3506	8303	2163

Table 10: Comparison of #tokens on medium-level tasks.

## Observation 9:

The **clean CoT demonstration** plays a pivotal role in CD-CoT.

## Observation 10:

The accuracy exhibits **subtle variations** when employing different algorithm instances. We set  $M = 2$  to strike a balance of efficiency and effectiveness.

## Observation 11:

An ablation study of components in Appendix F.3 demonstrates the denoising power and performance gain of CD-CoT, attributed to its **contrastive denoising with rationale rephrasing** and **repeated reasoning with voting components**.

# Empirical evaluations of CD-CoT

Model	Method	$\text{Acc}(\mathcal{M}, \mathcal{Q}, \mathcal{P}_{\text{irrelevant}})$			$\text{Acc}(\mathcal{M}, \mathcal{Q}, \mathcal{P}_{\text{inaccurate}})$		
		Base-9	Sym.(E)	Com.	Base-9	Sym.(E)	Com.
GPT-3.5-turbo	Base	30.3	25.1	42.3	10.1	26.1	33.4
	SC	36.6	28.3	45.0	17.3	30.7	44.7
	BT	34.2	22.7	28.3	18.4	22.7	11.0
	CC	44.3	32.7	43.6	31.7	33.0	40.8
	CD-CoT	<b>60.7</b>	<b>42.7</b>	<b>54.7</b>	<b>58.7</b>	<b>41.3</b>	<b>49.7</b>
Gemini-Pro	Base	72.3	38.9	53.2	21.2	36.7	33.5
	SC	80.3	43.3	60.0	32.3	45.0	42.7
	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
	CD-CoT	<b>92.7</b>	<b>49.3</b>	<b>57.7</b>	<b>76.7</b>	<b>53.3</b>	<b>55.7</b>
LLaMA2-70B	Base	2.8	8.7	41.9	2.7	9.1	40.2
	SC	<b>5.0</b>	10.3	<b>46.7</b>	<b>3.0</b>	9.7	<b>46.0</b>
	BT	1.4	11.2	36.1	0.9	12.5	36.2
	CC	1.1	<b>16.3</b>	29.9	2.8	<b>14.0</b>	28.3
	CD-CoT	<u>4.0</u>	9.7	<u>39.3</u>	2.7	9.7	39.7
Mixtral-8x7B	Base	16.3	17.9	34.9	3.7	15.1	31.1
	SC	<u>20.0</u>	<u>21.7</u>	<u>37.0</u>	2.7	18.0	<u>37.7</u>
	BT	4.1	9.7	6.2	2.4	10.1	10.5
	CC	<b>24.4</b>	18.5	36.0	<b>12.5</b>	<u>18.3</u>	35.7
	CD-CoT	8.7	<b>22.7</b>	<b>40.3</b>	<u>4.7</u>	<b>21.3</b>	<b>40.3</b>

Table 11: Comparing methods with different LLMs.

## Observation 12:

**CD-CoT generalizes well across different LLMs.** It respectively achieves increases of 23.4% and 21.6% in accuracy compared to base models (GPT-3.5-turbo and Gemini-Pro) and surpasses all the baselines.

# New algorithm: CD-CoT - Evaluations

## Comparison of rephrased rationales by different reasoning methods

Original Noisy Rationale
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. $13 \bmod 9 = 4$ , so the digit is 4 and the carry is 1. We have $8 + 5 + 1 = 14$ in base-10. $14 \bmod 9 = 5$ , so the digit is 5 and the carry is 1. A leading digit is 1. So the answer is 154. Answer: 154.

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. $13 \bmod 9 = 4$ , so the digit is 4 and the carry is 1. We have $8 + 5 + 1 = 14$ in base 10. $14 \bmod 9 = 5$ , so the digit is 5 and the carry is 1. A leading digit is 1. So the answer is 154. Answer: 154.
------	--



Method	Denoised Rationale
SM	In base-9, the dkigits 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 = 4, soWghe digix is 4 anj the carry is 1. We hOve 8 + 5i+ z2= 14 Vn baAe 10. $14 \bmod 9 = 5$ . 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.
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. $13 \bmod 9 = 4$ , so the digit is 4 and the carry is 1. We have $8 + 5 + 1 = 14$ in base 10. $14 \bmod 9 = 8$ . $14 \div 9 = 1$ , so the digit is 8 and the carry is 1. A leading digit is 1. So the answer is 154. Answer: 154.

SP	In base-9, the digits range from 0 to 8. Let's 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 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.
----	--

**Observation 13:** CD-CoT effectively removes noisy thoughts and ensures format alignment with the original rationale.

# Outline

- Background: language model reasoning
- New research problem: Noisy Rationales
- New benchmark: NoRa
- New algorithm: CD-CoT
- **Take home messages**
- **Future directions**

# Take home messages

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

- we introduce **NoRa dataset** to evaluate LLMs against noisy rationales
- we show the **general vulnerability** LLMs to noisy rationales
  - and is inadequately mitigated by existing robust methods
- we design **CD-CoT** to enhance the robustness via contrastive denoising

# Future directions

- **Knowledge-enhanced** denoising within a retrieval-augmented framework
- **Robust inductive reasoning** to extract rules from noisy examples
- Generalization to **out-of-distribution** noisy scenarios
- Expanding the NoRa dataset to include **multi-modal** scenarios, e.g., visual data, for a more comprehensive understanding of the robustness of foundation models
- **Theoretical analysis** of noisy ICL for deeper insights into the noisy rationales

# Thanks you!

Zhanke Zhou

[cszkzhou@comp.hkbu.edu.hk](mailto:cszkzhou@comp.hkbu.edu.hk)