# Knowledge-based Question Answering with Large Language Models

Yunshi Lan

2023 Nov

East China Normal University

# Overview

Knowledge-based Question Answering with LLMs

- Multi-modal Questions
- Solving Reasoning Problems under Noisy Context
- Generating Complex Questions

Future Directions

- The Applications of LLMs on more NLP tasks
- Instruction-tuning of LLMs for KGQA

华东师范大学数据科学与工程学院
School of Data Science and
Engineering at ECNU

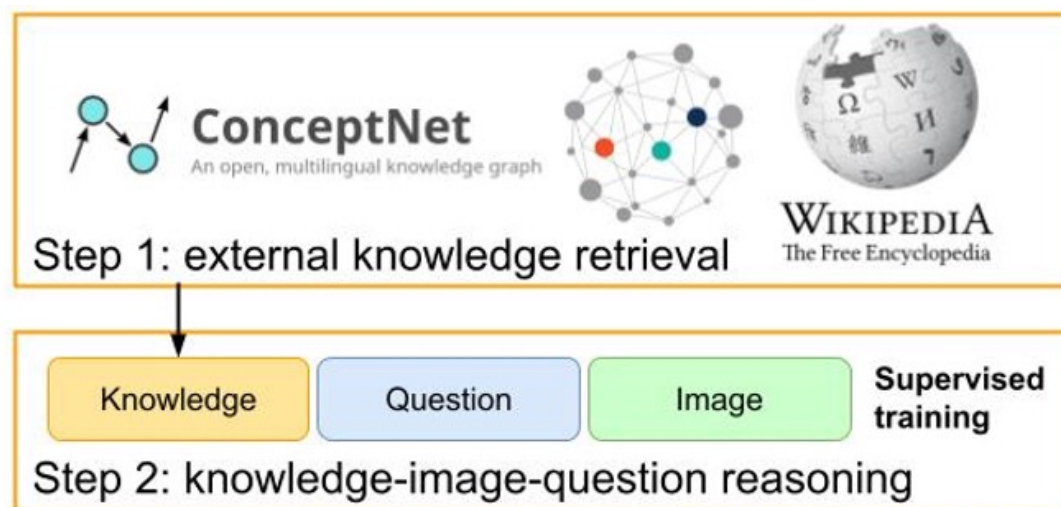## Knowledge-based Question Answering with LLMs

- **Multi-modal Questions**
- Solving Reasoning Problems under Noisy Context
- Generating Complex Questions

## Future Directions

- The Applications of LLMs on more NLP tasks
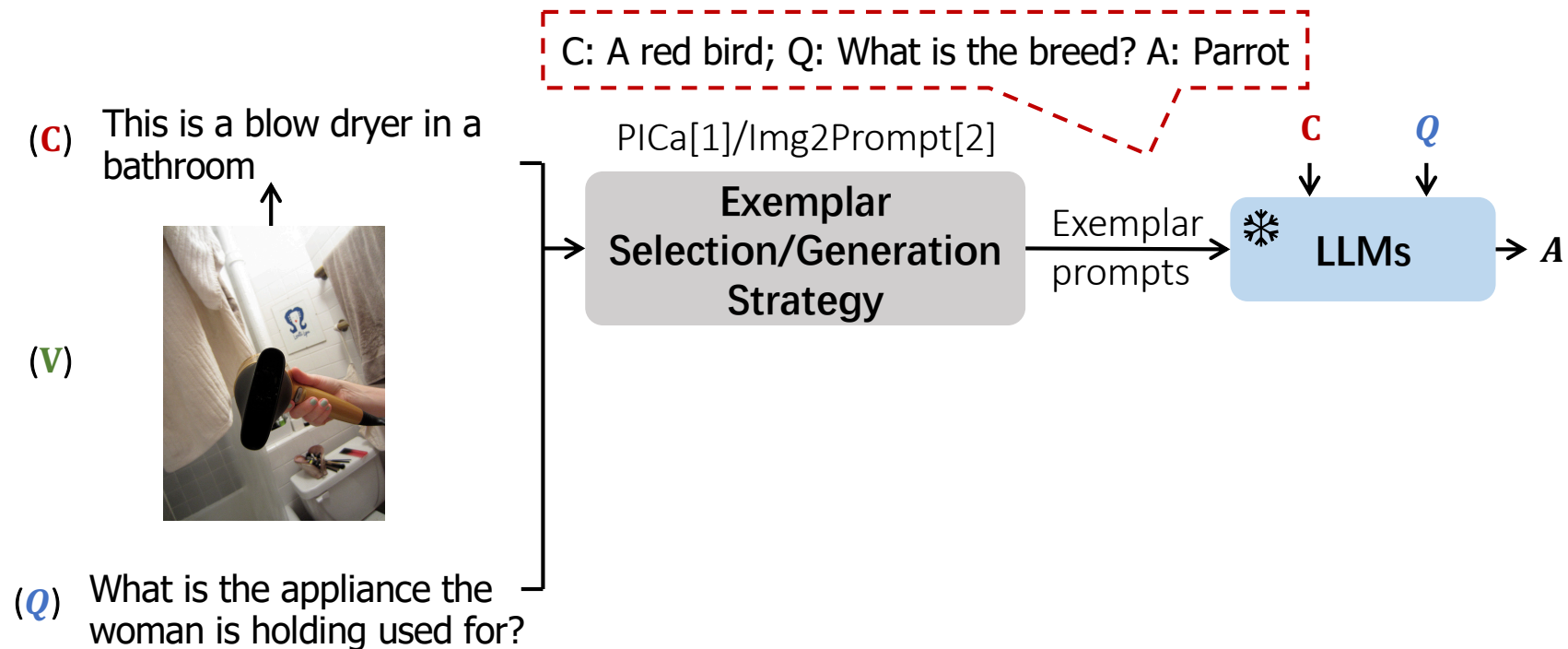- Instruction-tuning of LLMs for KGQA

# Visual Question Answering

What is the appliance the woman is holding used for?



Step 1: external knowledge retrieval

Step 2: knowledge-image-question reasoning

[1] Zhengyuan Yang, Zhe Gan, Jianfei Wang, Xiaowei Hu, Yumao Lu, Zicheng Liu, and Lijuan Wang. 2022. An Empirical Study of GPT-3 for Few-Shot Knowledge-Based VQA. AAAI. 2022.

# A New Paradigm

C: A red bird; Q: What is the breed? A: Parrot

(**C**) This is a blow dryer in a bathroom

(**V**)

(**Q**) What is the appliance the woman is holding used for?

PICa[1]/Img2Prompt[2]

**Exemplar Selection/Generation Strategy**

Exemplar prompts

**C**　**Q**

❄ **LLMs** → **A**

[1] Zhengyuan Yang, Zhe Gan, Jianfei Wang, Xiaowei Hu, Yumao Lu, Zicheng Liu, and Lijuan Wang. 2022. An Empirical Study of GPT-3 for Few-Shot Knowledge-Based VQA. AAAI. 2022.
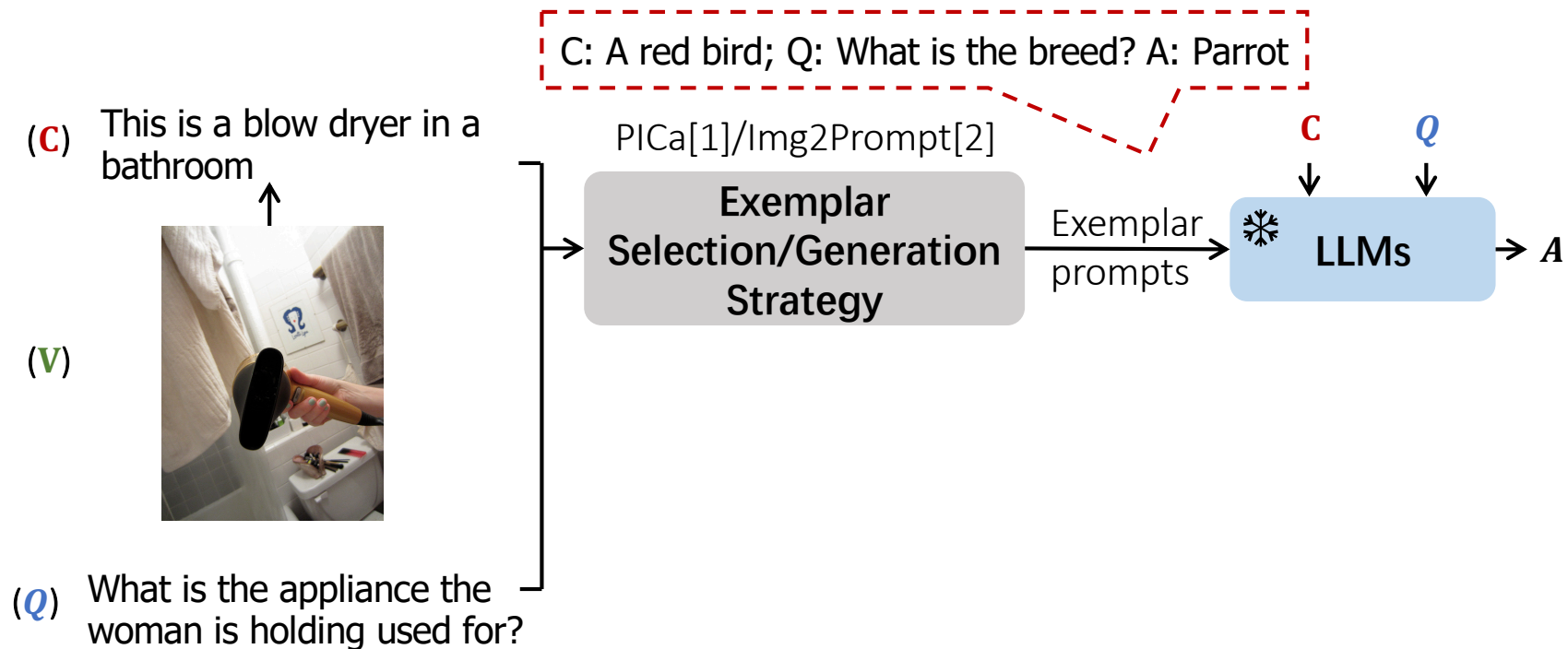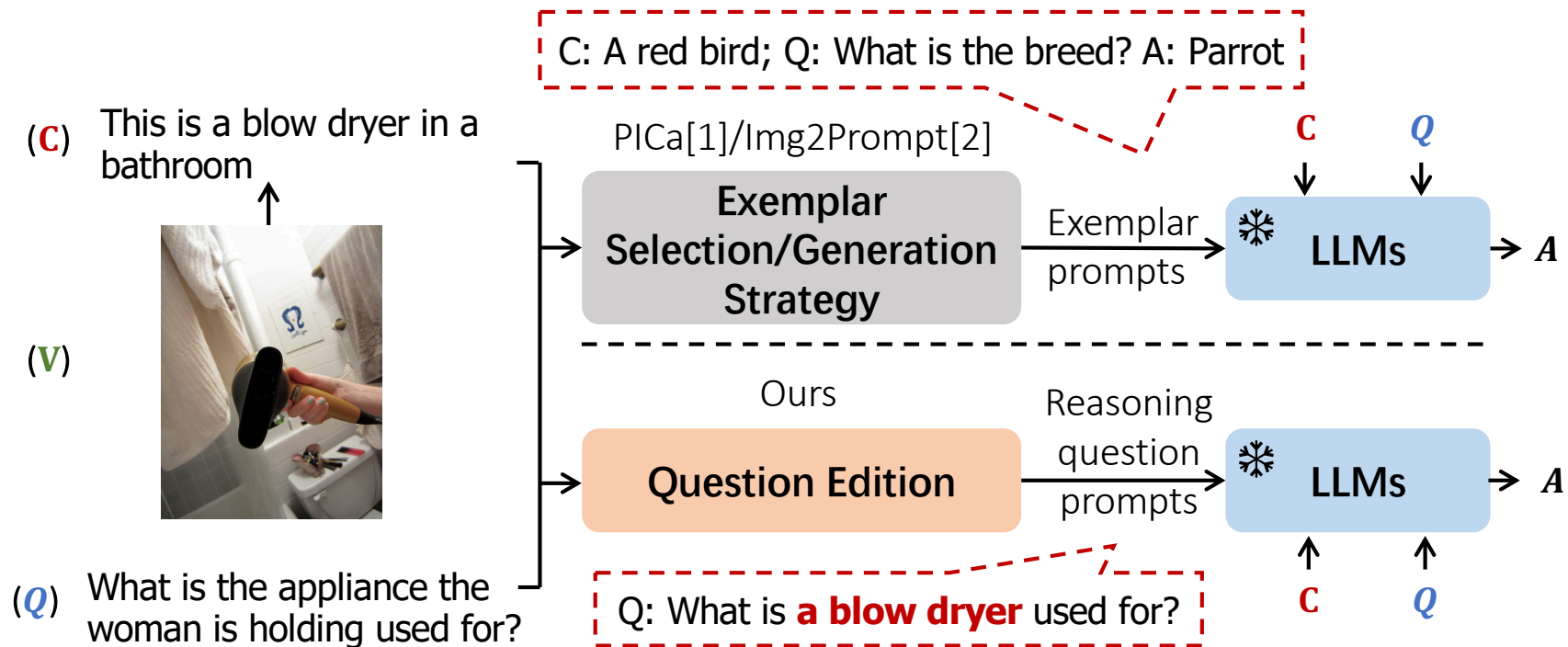
[2] Jiaxian Guo, Junnan Li, Dongxu Li, Anthony Tiong, Boyang Li, Dacheng Tao, and Steven Hoi. 2023. From Images to Textual Prompts: Zero-shot Visual Question Answering with Frozen Large Language Models. CVPR. 2023.

DaSE
Data Science
&Engineering

华东师范大学数据科学与工程学院
School of Data Science and
Engineering at ECNU

C: A red bird; Q: What is the breed? A: Parrot

(C) This is a blow dryer in a bathroom

(V)

(Q) What is the appliance the woman is holding used for?

PICa[1]/Img2Prompt[2]

**Exemplar Selection/Generation Strategy**

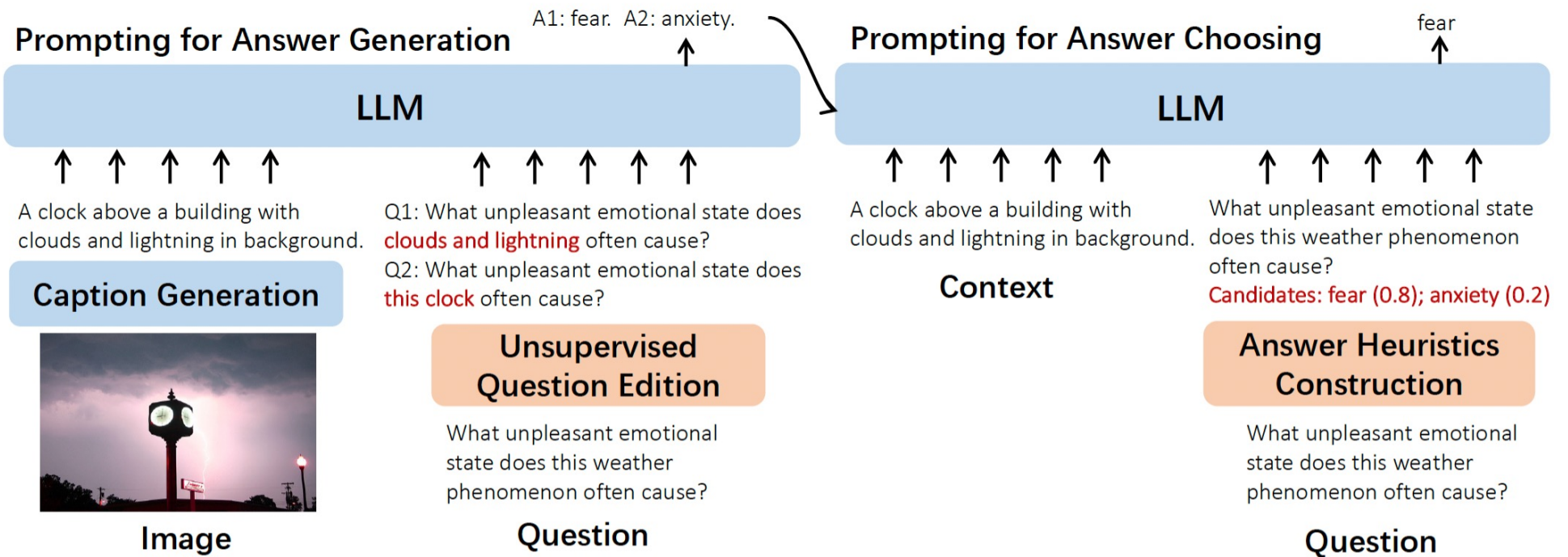Exemplar prompts

**C**  **Q**

❄ **LLMs** → **A**

1. The current methods entirely rely on the understanding capability of LLMs to resolve the ambiguity and infer the intent of the questions, which might **involve unexpected bias**.

2. LLMs are **brittle to ill-posed questions**, especially under the zero-shot setting.
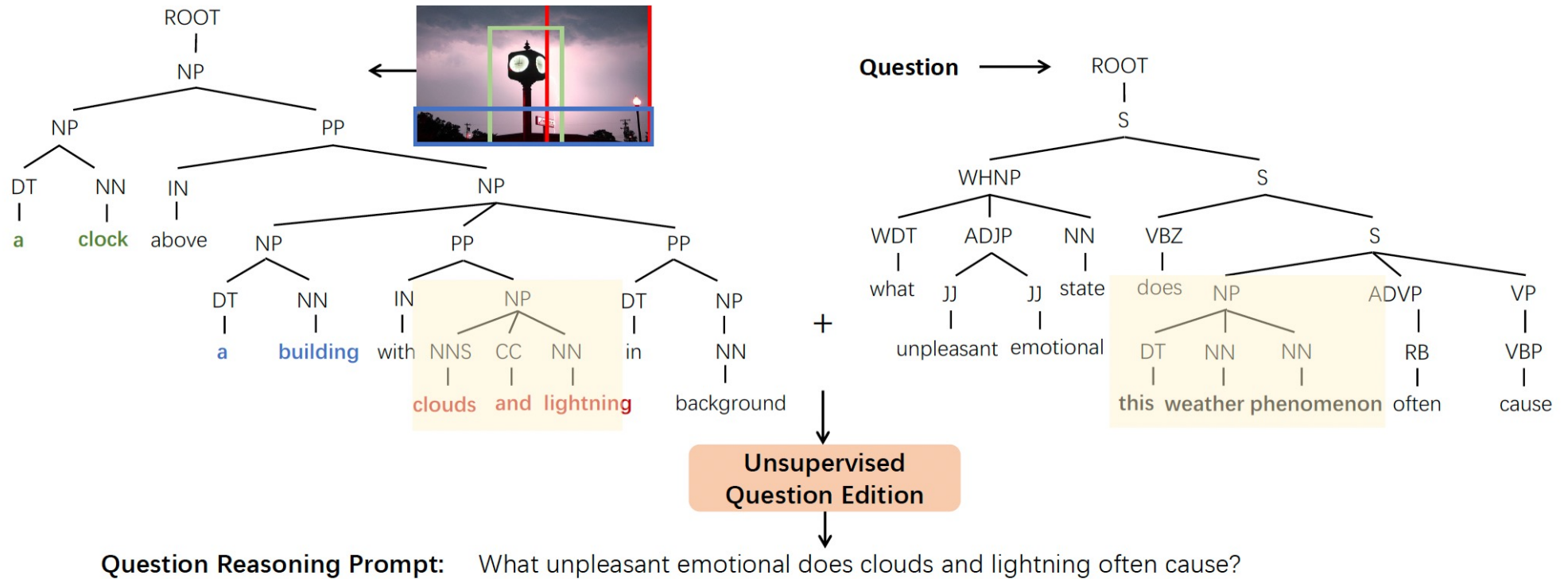
# Our Motivation

(C) This is a blow dryer in a bathroom

(V)

(Q) What is the appliance the woman is holding used for?

C: A red bird; Q: What is the breed? A: Parrot

PICa[1]/Img2Prompt[2]

**Exemplar Selection/Generation Strategy** — Exemplar prompts → ❄ **LLMs** → A

C Q

Ours

**Question Edition** — Reasoning question prompts → ❄ **LLMs** → A

C Q

Q: What is **a blow dryer** used for?

**Reasoning Question Prompts:** converting original questions into self-contained questions by editing the segments of the question.

Yunshi Lan, Alex Xiang Li, Xin Liu, Yang Li, Wei Qin, Weining Qian. Improving Zero-shot Visual Question Answering via Large Language Models with Reasoning Question Prompts. ACM MM. 2023

# Reasoning Question Prompt

Our prompting method generates **reasoning question prompts** and enables LLMs to perform VQA tasks with **two-step reasoning**.

# Prompting for Answer Generation

**Question Reasoning Prompt:** What unpleasant emotional does clouds and lightning often cause?

# Prompting for Answer Generation

**Unsupervised Question Edition:** Our scoring function for evaluating the quality of the candidates:

- LM Score. $f_{LM}(\tilde{Q}) = In \prod_{i=1}^{T} P(w_i|w_{i-1}, \ldots, w_1)$

- Semantic Integrity. $f_{semantic}(\tilde{Q}) = \cos(\tilde{Q}, Q)$

- Syntactic Invariance. $f_{syntactic}(\tilde{Q}) = \mathbb{I}(\text{Tag}_{Q[i]} = \text{Tag}_{\tilde{Q}[j]})$

The overall scoring function:

$$f(\tilde{Q}) = f_{LM}(\tilde{Q})^{\alpha} f_{semantic}(\tilde{Q})^{\beta} f_{syntactic}(\tilde{Q})$$

**Prompt Design:**

**Instruction:** Please answer the question according to the contexts.
**Context:** [caption].
**Question:** [reasoning question prompt].
**Answer:**

# Prompting for Answer Choosing

**Answer Heuristics Construction:**

$$P(A) = \sum_{LLM(\tilde{Q}) \to A} P(\tilde{Q}) P_{LLM}(A|\tilde{Q})$$

**Prompt Design:**

**Instruction**: Please answer the question according to the contexts and candidates.
**Context**: [caption].
**Question**: [original question].
**Candidates**: $[A_1\ P(A_1)];[A_2\ P(A_2)];\ldots;[A_m\ P(A_m)]$
**Answer**:

[1] Zhenwei Shao, Zhou Yu, Meng Wang, and Jun Yu. Prompting Large Language Models with Answer Heuristics for Knowledge-based Visual Question Answering. arXiv preprint arXiv:2303.01903. 2023.

# Experiments

## Datasets:

- OK-VQA: 5,046 test questions.

- A-OKVQA: 1,100 and 6,700 questions for validation and testing, respectively.

- VQAv2: 214,354 validation questions.

## Comparable Methods:

- LLM-based methods: PICa, Img2Prompt

- Pre-trained zero-shot VQA methods: Flamingo, Frozen VL-T5, FewVLM and VLKD.

# Experimental Results

| Method | Model size | Shot number | Examplar number | OK-VQA test | VQAv2 val | A-OKVQA val | A-OKVQA test |
|---|---|---|---|---|---|---|---|
| *Zero-shot Evaluation with Frozen LLMs* | | | | | | | |
| PICa {GPT-3} | 175B | 0 | 0 | 17.7 | – | 23.8◇ | – |
| Img2Prompt {OPT} | 6.7B | 0 | 30 | 38.2 | 52.2◇ | 33.3 | 32.2 |
| Img2Prompt {OPT} | 30B | 0 | 30 | 41.8 | 54.2◇ | 36.9 | 33.0 |
| Img2Prompt {GPT-3} | 175B | 0 | 30 | 42.8 | – | 38.9◇ | 43.4◇ |
| Img2Prompt {OPT} | 175B | 0 | 30 | 45.6 | **60.6** | 42.9 | 40.7 |
| PICa+RQ prompt {GPT-3} (Ours) | 175B | 0 | 0 | 20.3(↑ 2.6) | – | 29.0(↑ 5.2) | – |
| Img2Prompt+RQ prompt {OPT} (Ours) | 6.7B | 0 | 30 | 38.5(↑ 0.3) | 52.9(↑ 0.7) | 36.3(↑ 3.0) | 31.5 |
| Img2Prompt+RQ prompt {OPT} (Ours) | 30B | 0 | 30 | 42.1(↑ 0.3) | 54.5(↑ 0.3) | 38.1(↑ 1.2) | 35.2(↑ 3.0) |
| Img2Prompt+RQ prompt {GPT-3} (Ours) | 175B | 0 | 30 | **46.4**(↑ 3.6) | – | **43.2**(↑ 4.3) | **43.9**(↑ 0.5) |
| *Zero-shot Evaluation with Pre-trained VQA methods* | | | | | | | |
| VL-T5 {no-vqa} | 224M | 0 | 0 | 5.8 | 13.5 | – | – |
| FewVLM {large} | 740M | 0 | 0 | 16.5 | 47.7 | – | – |
| VLKD {ViT-L/14} | 408M | 0 | 0 | 13.3 | 44.5 | – | – |
| Frozen | 7B | 0 | 0 | 5.9 | 29.5 | – | – |
| Flamingo | 80B | 0 | 0 | 50.6 | – | – | – |
| *Few-shot Evaluation with Frozen LLMs* | | | | | | | |
| PICa {GPT-3} | 175B | 16 | 16 | 46.5 | 54.3 | – | – |
| Prophet {GPT-3} | 175B | 20 | 20 | 61.1 | – | – | – |

## RQ prompts can generally improve VQA tasks under zero-shot setting
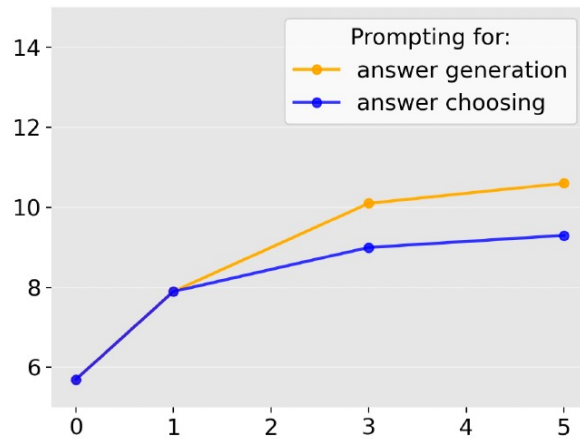
Table 3: Zero-shot performance A-OKVQA validation set having Img2Prompt as baselines but with different LLMs. $\Delta$ denotes the performance gain brought by QR prompts.

| LLMs | Img2Prompt | +QR prompt | $\Delta$ |
|---|---|---|---|
| GPT-3 175B | 38.9 | 43.2 | ↑ 4.3 |
| GPT-3.5 175B | 37.1 | 40.3 | ↑ 3.2 |
| GPT-Neo 2.7B | 29.7 | 31.5 | ↑ 1.8 |
| BLOOM 7.1B | 29.8 | 32.1 | ↑ 2.3 |
| GPT-J 6B | 32.5 | 33.1 | ↑ 0.6 |
| OPT 125M | 10.8 | 13.3 | ↑ 2.5 |

RQ prompts can generally collaborate with LLMs

(a) Shot number=0
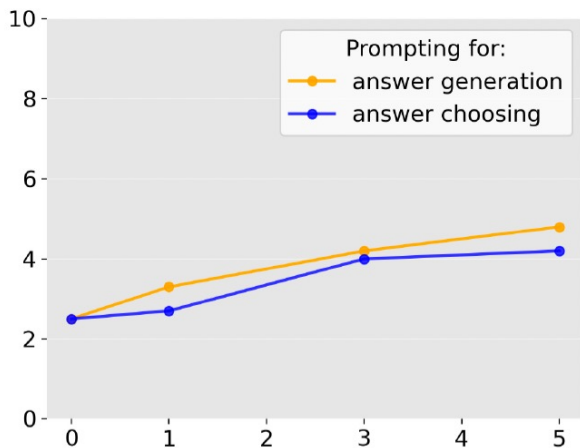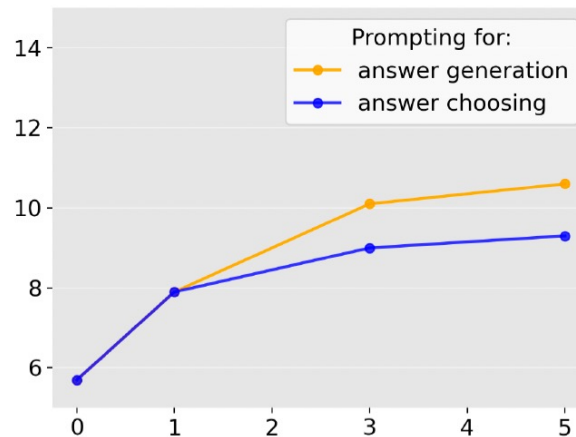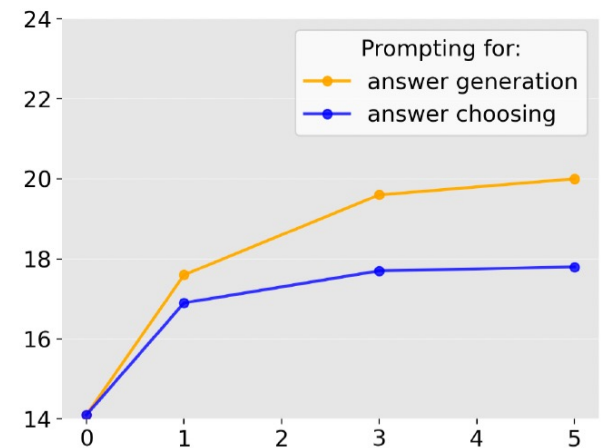
(b) Shot number=1

(c) Shot number=16

(d) LLM=GPT-Neo 2.7B

(e) LLM=OPT 6.7B

(f) LLM=OPT 30B

Larger improvement gain brought by RQ prompts can be shown when: (1) shot number is low; (2) the size of LLM is relatively large

# Case Study

**Caption**: *This is a blow dryer in a bathroom.*

**Question**: *What is the appliance the woman is holding used for?*
**GT Answer**: *drying hair*
**Original Answer**: *cutting hair*

**Prompting for Answer Generation:**
**Q1**: *What is the appliance **a blow dryer** used for?* $\quad P(\tilde{Q}) = 0.21$
**A1**: *drying hair* $\qquad P_{LLM}(A|\tilde{Q}) = 0.15$
**Q2**: *What is the appliance **a bathroom** is holding used for?* $\quad P(\tilde{Q}) = 0.29$
**A2**: *drying hair* $\qquad P_{LLM}(A|\tilde{Q}) = 0.10$

**Prompting for Answer Choosing:**

**Question**: *What is the appliance the woman is holding used for?*
**Candidates**: *drying hair (1.00)*
**Predicted Answer**: *drying hair*

(a)

**Caption**: *A little girl holding a cup with rice in dishes in front of her*

**Question**: *What is the child eating?*
**GT Answer**: *rice*
**Original Answer**: *spaghetti*

**Prompting for Answer Generation:**
**Q1**: *what is **dishes in front of her**?* $P(\tilde{Q}) = 0.15$
**A1**: *rice* $\qquad P_{LLM}(A|\tilde{Q}) = 0.20$
**Q2**: *What is the child eating?* $\quad P(\tilde{Q}) = 0.7$
**A2**: *spaghetti* $\qquad P_{LLM}(A|\tilde{Q}) = 0.10$
**Q3**: *what is **a cup with food in dishes in front of her**?* $\quad P(\tilde{Q}) = 0.15$
**A3**: *rice* $\qquad P_{LLM}(A|\tilde{Q}) = 0.30$

**Prompting for Answer Choosing:**

**Question**: *What is the child eating?*
**Candidates**: *spaghetti (0.48); rice (0.51)*
**Predicted Answer**: *rice*

(b)

The questions become self-contained with RQ prompts.

# Conclusions

- RQ prompts are helpful to bridge the gap between questions and captions, such that it can boost performance of leveraging LLMs to VQA tasks.

- RQ prompts show general improvement on different LLMs. It could achieve SOTA results on three of four VQA tasks on zero-shot setting.

- RQ prompts show more effect on zero-shot setting and large LLMs.

华东师范大学数据科学与工程学院
School of Data Science and
Engineering at ECNU

Knowledge-based Question Answering with LLMs

- Multi-modal Questions

- **Solving Reasoning Problems under Noisy Context**

- Generating Complex Questions

Future Directions

- The Applications of LLMs on more NLP tasks

- Instruction-tuning of LLMs for KGQA

# Chain-of-Thought

**Standard Prompting**

**Model Input**

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

**Model Output**

A: The answer is 27. ❌

**Chain-of-Thought Prompting**

**Model Input**

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

**Model Output**

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had 23 - 20 = 3. They bought 6 more apples, so they have 3 + 6 = 9. The answer is 9. ✔
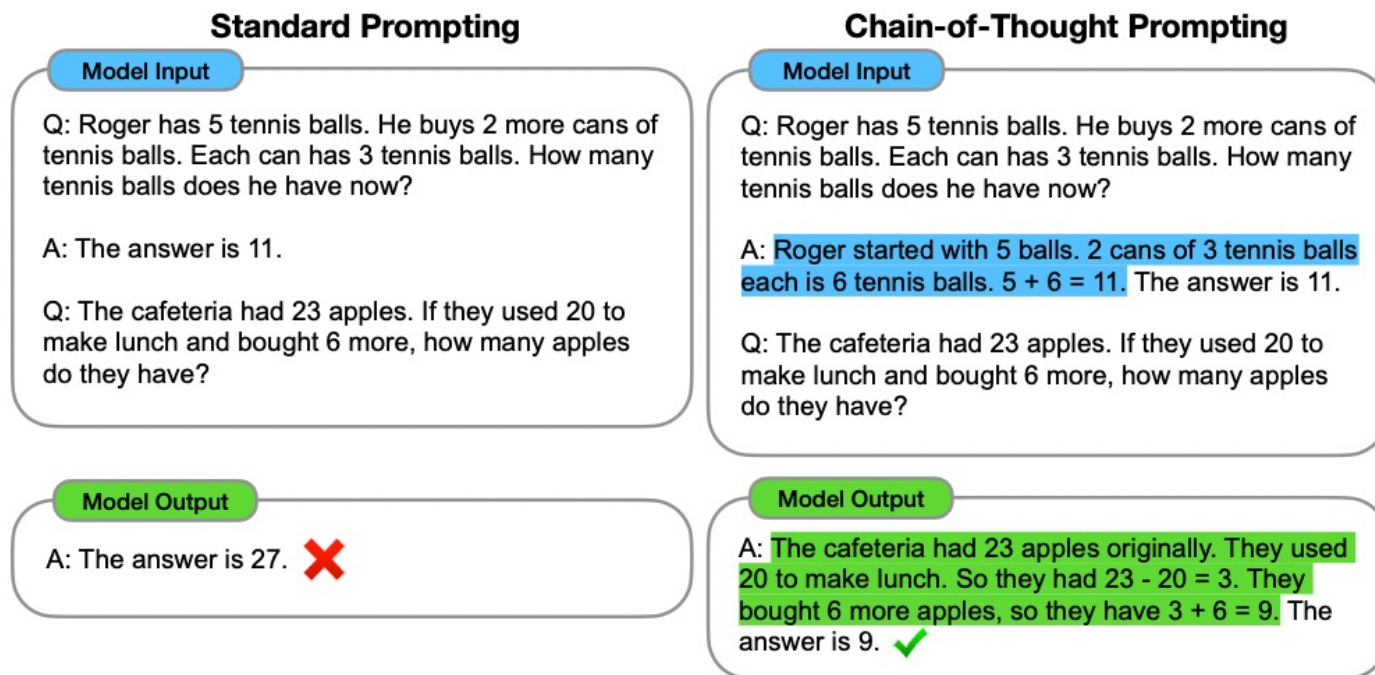
Figure 1: Chain-of-thought prompting enables large language models to tackle complex arithmetic, commonsense, and symbolic reasoning tasks. Chain-of-thought reasoning processes are highlighted.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, Denny Zhou. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. NeurIPS. 2022.

# LLMs under Noisy Context

**Original Problem**
Jessica is six years older than Claire. In two years, Claire will be 20 years old. How old is Jessica now?

**Modified Problem**
Jessica is six years older than Claire. In two years, Claire will be 20 years old. *Twenty years ago, the age of Claire's father is 3 times of Jessica's age.* How old is Jessica now?

**Standard Answer** 24

*Table 1.* An example problem from GSM-IC. An irrelevant sentence (*italic and underlined*) that does not affect the standard answer is added immediately before the question.

Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed Chi, Nathanael Sch.rli, and Denny Zhou. 2023. Large language models can be easily distracted by irrelevant context. ICML. 2023
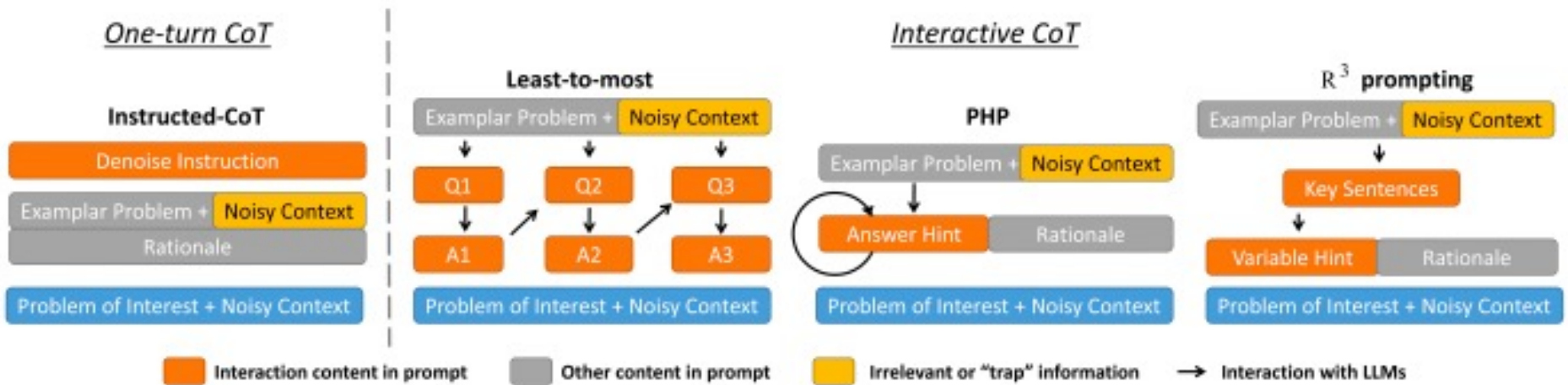
Figure 1: Comparison between $R^3$ prompting and existing CoT prompting baseline methods. The exemplar problems are multiple problems we used as exemplars for in-context learning. Rationales are reasoning chains in prompts. The problem of interest is the query problem.

Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed Chi, Nathanael Sch.rli, and Denny Zhou. 2023. Large language models can be easily distracted by irrelevant context. ICML. 2023

Denny Zhou, Nathanael Scharli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc Le, Ed Chi. LEAST-TO-MOST PROMPTING ENABLES COMPLEX REASONING IN LARGE LANGUAGE MODELS. ICLR 2023.

Chuanyang Zheng, Zhengying Liu, Enze Xie, Zhenguo Li, and Yu Li. 2023. Progressive-hint prompting improves reasoning in large language models. arXiv preprint arXiv:2304.09797.

Qingyuan Tian, Hanlun Zhu, Lei Wang, Yang Li, Yunshi Lan, R3 Prompting: Review, Rephrase and Resolve for Chain-of-Thought Reasoning in Large Language Models under Noisy Context. EMNLP Finding, 2023

# Our Method

Figure 2: A running example of the inputs and outputs of $R^3$ prompting in LLMs at each prompting stage. Green: In-topic noisy context. Red: Off-topic noisy context. Blue: Key sentences.

## Datasets:

- AddSub
- SVAMP
- GSM-IC
- MultiArith-IC
- SinglEq-IC

| Dataset | #Sample | Ave. in-topic | Ave. off-topic |
|---|---|---|---|
| GSM-IC | 1000 | 0.5 | 0.5 |
| MultiArith-IC | 600 | 0.5 | 0.5 |
| SingEq-IC | 508 | 0.48 | 0.52 |

Table 1: Details of constructed datasets. "Ave. in-topic" and "Ave. off-topic" denotes averge number of in-topic sentences and off-topic sentences, respectively.

## Comparable Methods:

- One-turn interaction:  Manual-CoT, Auto-CoT, Instructed-CoT
- Multi-turn interaction : Least-to-Most, PHP

# Experimental Results

| | Methods | SVAMP | MultiArith-IC | SingleEq-IC | AddSub | GSM-IC | Average |
|---|---|---|---|---|---|---|---|
| One-turn | Manual-CoT | 79.9 | 79.5 | 77.7 | 85.3 | 81.0 | 79.7 |
| | Auto-CoT | 83.6 | 79.7 | 77.6 | 88.0 | 81.5 | 82.1 |
| | Instructed-CoT | 81.3 | 80.1 | 78.2 | 87.3 | 82.0 | 81.8 |
| Interactive | Least-to-Most | 80.8 | 77.4 | 76.2 | 85.3 | 81.1 | 80.2 |
| | PHP | 83.1 | 80.0 | 79.0 | 85.3 | 85.1 | 82.5 |
| | $R^3$ Prompting (Ours) | **87.3** | **82.2** | **81.5** | **90.0** | **88.0** | **85.8** |

Table 2: Main result on five evaluated datasets. The best and second best results are boldfaced and underlined respectively.

(1) R^3 prompting performs well for CoT reasoning in LLMs under noisy context;
(2) The design of interactive prompts are important for denoising.

# More Analysis



Figure 3: (a). The results of prompting methods after adding Self-Consistency (SC) on AddSub dataset. (b). The results of prompting methods after adding Self-Consistency (SC) on SVAMP dataset.

Figure 4: Accuracy change of various methods with the increasing number of irrelevant sentences on AddSub.

(1) Improvement of R^3 prompting is still significant with self-consistency;

(2) R^3 prompting exhibits robust performance under noisy context while Instructed-CoT and Manual-CoT are vulnerable when facing a large amount of noisy information.

# Conclusions

- By comparison, one-turn CoTs are more robust than interactive CoTs when conducting reasoning under noisy context.

- $R^3$ prompting can effectively restrain the influence of noisy context. The three steps (i.e. review, rephrase and resolve) collaborative contribution to the good performance.

華東師範大學數據科學與工程學院
School of Data Science and Engineering at ECNU

Knowledge-based Question Answering with LLMs

- Multi-modal Questions
- Solving Reasoning Problems under Noisy Context
- **Generating Complex Questions**

Future Directions

- The Applications of LLMs on more NLP tasks
- Instruction-tuning of LLMs

# Knowledge Base Question Generation

Figure 1: Overview of KQG-CoT framework.

Yuanyuan Liang, Jianing Wang, Hanlun Zhu, Lei Wang, Yunshi Lan, Weining Qian. Prompting Large Language Models with Chain-of-Thought for Few-Shot Knowledge Base Question Generation. EMNLP. 2023

# Existing Methods

The existing methods solve few-shot KBQG tasks via designing prompter for the pair of sub-graph description and generated questions and conducting pre-training Language Models.

Guanming Xiong, Junwei Bao, Wen Zhao, Youzheng Wu, and Xiaodong He. Autoqgs: Auto-prompt for low resource knowledge-based question generation from sparql. In Proceedings of the 31st ACM International Conference on Information Knowledge. 2022

- A **substantial amount of annotated data** is required, and acquiring it can be challenging.

- A logical form is made up of entities, relations, and query grammar. It's impossible to **encompass all the possible combinations of these fundamental components**.

- Certain **logical forms can become complex** when operations such as aggregation, superlatives, and comparisons are involved.

# Our Method

DaSE
Data Science
&Engineering

华东师范大学数据科学与工程学院
School of Data Science and
Engineering at ECNU

- LLMs have the strong capability to accurately capture the semantics of relations between values in the data, enabling to transform the structured context to narrative text.
  - ▷ Structured logical forms to **natural language questions**

- LLMs have proven their strong generalizability on a wide range of few-shot and zero-shot tasks with Chain-of-Thought.
  - ▷ **Apply CoT** to solve few-shot KBQG

[1] Milena Trajanoska, Riste Stojanov, and Dimitar Trajanov. 2023. Enhancing knowledge graph construction using large language models. arXiv.

[2] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. In Advances in Neural Information Processing Systems.

# Our Method

## Supportive Logical Forms Selection

### Logical Forms

(AND medicine.routed_drug (JOIN medicine.routed_drug.marketed_formulations m.0hqs1x_))

...

(AND music.producer (JOIN music.producer.tracks_produced (JOIN music.recording.producer m.01w8w6z)))

### Structures

(AND r (JOIN r e))

...

(AND r (JOIN (JOIN r e)))

### Clusters

1 ... k

### Supportive Logical Forms

1. (AND sports.sport (JOIN sports.sport.team_coaches John Russo))

2. (COUNT (AND music.album (JOIN music.album.featured_artists Sierra Leone)))

...

K. (JOIN (R people.person.place_of_birth) (JOIN music.producer.releases_produced Poguetry In Motion))

## Prompt Construction

### Demonstrations

Input : (AND sports.sport (JOIN sports.sport.team_coaches John Russo))
Subgraph1 : (JSON sports.sport.team_coaches John Russo)
Subgraph2 : (AND sports.sport Subgraph1)
Subquestion1 : sport team coach john russo
Subquestion2 : Which sport does john russo coach?
...
Input : (JOIN (R people.person.place_of_birth) (JOIN music.producer.releases_produced Poguetry In Motion))
Subgraph1 : (JOIN music.producer.releases_produced Poguetry In Motion))
Subgraph2 : (JOIN (R people.person.place_of_birth) Subgraph1 )
Subquestion1 : music producer of released produce poguetry in motion
Subquestion2 : Where is the birth place of the music producer of poguetry in motin?

Input : (COUNT (AND aviation.aircraft_manufacturer (JOIN organization.organization. legal_structure S.A.)))

### LLM

### Prediction

Subgraph1 : (JOIN organization.organization.legal_structure S.A.)
Subgraph2 : (AND aviation.aircraft_manufacturer Subgraph1 )
Subgraph3 : (COUNT Subgraph2 )
Subquestion1 : organization legal structure of s.a.
Subquestion2 : aircraft manufacturer in the legal structure of s.a.
Subquestion3 : What is the number of aircraft manufacturer in the legal structure of s.a.?

# Supportive Logical Forms Selection

## Step 1: Structure Encoding and Clustering

1. We extract structure of logical form by converting the schema items into symbolic variables.

**(AND medicine.routed_drug (JOIN medicine.routed_drug.marketed_formulations m.0hqs1x)).**

↓

**(AND r (JOIN r e))**

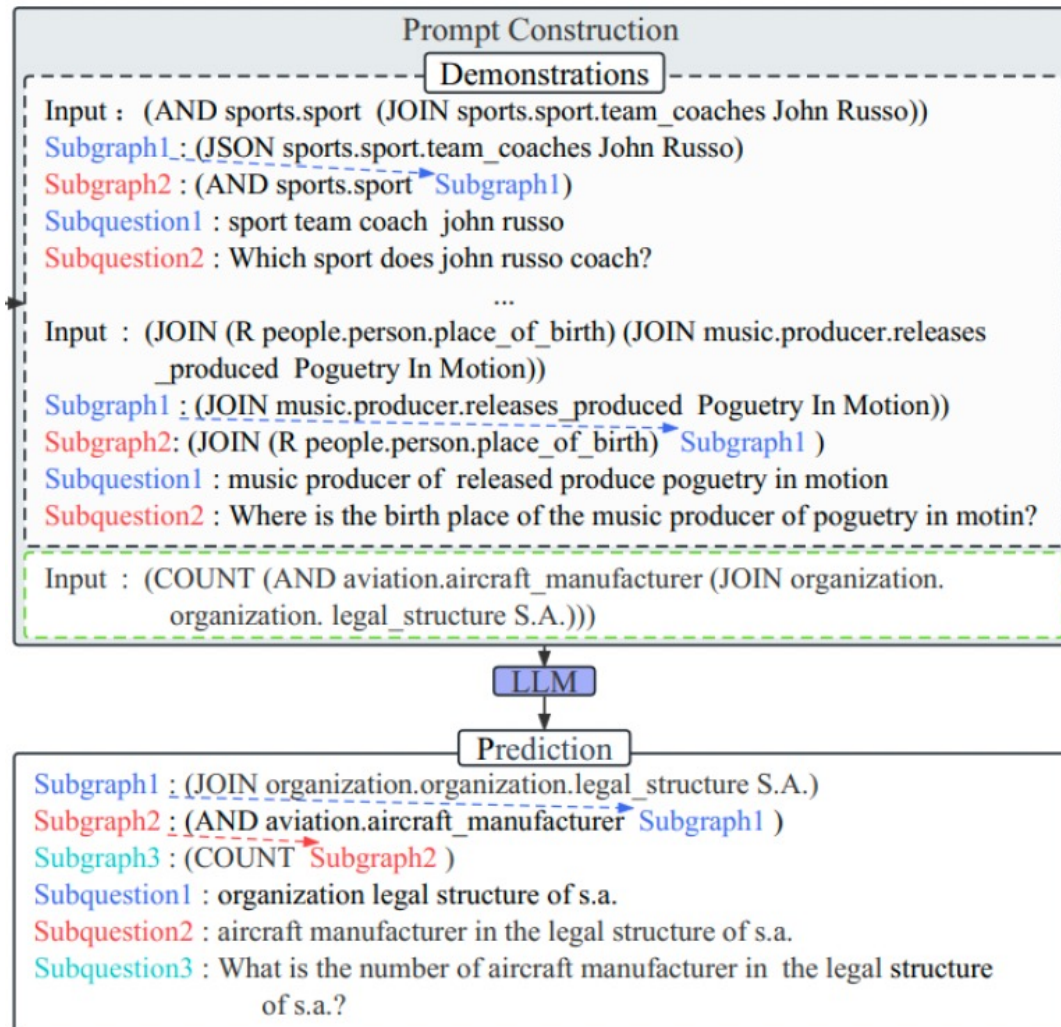2. We encode the contexts of the sequence with Sentence-Transformers.

## Step 2: Logical Form Sampling

1. We utilize the **K-means clustering algorithm** to group the encoded structure into k-clusters based on their syntactic similarity.

2. We greedily pick up a candidate with **least semantic similarity** to the selected logical forms, where the similarity is measured by the encoding of the original logical forms.

# Prompt Construction

**Prompt Construction**

**Demonstrations**

Input : (AND sports.sport (JOIN sports.sport.team_coaches John Russo))
Subgraph1 : (JSON sports.sport.team_coaches John Russo)
Subgraph2 : (AND sports.sport ► Subgraph1)
Subquestion1 : sport team coach john russo
Subquestion2 : Which sport does john russo coach?

...

Input : (JOIN (R people.person.place_of_birth) (JOIN music.producer.releases
_produced Poguetry In Motion))
Subgraph1 : (JOIN music.producer.releases_produced Poguetry In Motion))
Subgraph2 : (JOIN (R people.person.place_of_birth) ► Subgraph1 )
Subquestion1 : music producer of released produce poguetry in motion
Subquestion2 : Where is the birth place of the music producer of poguetry in motin?

Input : (COUNT (AND aviation.aircraft_manufacturer (JOIN organization.
organization. legal_structure S.A.)))

**LLM**

**Prediction**

Subgraph1 : (JOIN organization.organization.legal_structure S.A.)
Subgraph2 : (AND aviation.aircraft_manufacturer ► Subgraph1 )
Subgraph3 : (COUNT Subgraph2 )
Subquestion1 : organization legal structure of s.a.
Subquestion2 : aircraft manufacturer in the legal structure of s.a.
Subquestion3 : What is the number of aircraft manufacturer in the legal structure of s.a.?

- Generate a straightforward question that queries a **one-hop relation** from the topic entity.
- One-hop relation **subgraph1** leads to a simple **subquestion1**.
- Generate a question that inquires about a **two-hop relation chain** involving the aforementioned one-hop relation. The Step 2 includes the parsed logical form appended to the previous step as a component and generates **subquestion2** based on the **subgraph2** and **subquestion1**.
- Repeat until the **entire logical forms** have been traversed.

## Datasets:

- WebQuestions (WQ)
- PathQuestions (PQ)
- GrailQA (GQ)

| Dataset | #Q | #R | #E | #T |
|---|---|---|---|---|
| WQ | 22,989 | 672 | 25,703 | 2/99/5.8 |
| PQ | 9,731 | 378 | 7,250 | 2/3/2.7 |
| GQ | 64,331 | 3,720 | 32,585 | 1/4/1.4 |

Table 1: Statistics of the evaluated datasets. #Q denotes the number of questions. #R and #E denote the total number of relations and entities, respectively. #T denotes the minimum/maximum/average number of triplets involved in each question.

## Comparable Methods:

- LLMs+CoT methods : Standard Prompt, Random-CoT, Manual-CoT, Active-CoT, Auto-CoT

- Fine-trained methods : DSM, LFKQG, IGND, JointGT, T5-Large, etc.

- Few-shot methods:  BiGraph2Seq, JointGT , AutoQGS

- Our methods: KQG-CoT, KQG-CoT+ (further display the examplers from short to long.)

| Method | WQ | | | PQ | | | GQ | | |
|---|---|---|---|---|---|---|---|---|---|
| | B | M | R | B | M | R | B | M | R |
| Standard Prompt | 24.86 | 29.01 | 52.74 | 55.87 | 42.24 | 76.83 | 29.17 | 33.52 | 42.95 |
| Random-CoT | 25.02 | 29.37 | 53.16 | 56.42 | 42.61 | 77.03 | 29.81 | 33.75 | 43.31 |
| Manual-CoT | 28.44 | 30.24 | 54.30 | 60.37 | 42.88 | 77.48 | 30.18 | 33.61 | 44.89 |
| Active-CoT | 26.02 | 29.55 | 54.01 | 58.78 | 43.86 | 76.78 | 30.27 | 33.71 | 44.07 |
| Auto-CoT | 28.42 | 29.65 | 53.47 | 59.59 | 43.16 | 77.13 | 30.17 | 34.22 | 44.47 |
| KQG-CoT (Ours) | 28.89 | 30.41 | 54.38 | 60.81 | 43.54 | 77.35 | 30.51 | 34.26 | 44.91 |
| KQG-CoT+ (Ours) | **29.73** | **31.08** | **55.14** | **61.71** | **44.27** | **78.41** | **31.24** | **34.94** | **45.36** |

KQG-CoT outperforms the existing LLMs+CoT methods

# Experimental Results

| Method | WQ | | |
|---|---|---|---|
| | B | M | R |
| *Full Training* | | | |
| L2A (Du et al., 2017) | 6.01 | 26.95 | 25.24 |
| Transformer (Vaswani et al., 2017) | 8.94 | 13.79 | 32.63 |
| MHQG (Kumar et al., 2019) | 11.57 | 29.69 | 35.53 |
| BiGraph2Seq (Chen et al., 2023) | 29.45 | 30.96 | 55.45 |
| T5-Large (Raffel et al., 2020) | 28.78 | 30.55 | 55.12 |
| JointGT (Ke et al., 2021) | 30.02 | 32.05 | 55.60 |
| IGND (Fei et al., 2021) | 30.62 | 31.41 | 55.82 |
| LFKQG (Fei et al., 2022) | **31.66** | **32.69** | 56.75 |
| DSM (Guo et al., 2022) | 28.62 | - | **64.25** |
| *Few-shot Evaluation* | | | |
| KQG-CoT | 28.89 | 30.41 | 54.87 |
| KQG-CoT+ | 29.73 | 31.08 | 55.46 |

| Method | PQ | | |
|---|---|---|---|
| | B | M | R |
| *Full Training* | | | |
| L2A (Du et al., 2017) | 17.00 | 50.38 | 19.72 |
| Transformer (Vaswani et al., 2017) | 56.43 | 43.45 | 73.64 |
| MHQG (Kumar et al., 2019) | 25.99 | 33.16 | 58.94 |
| BiGraph2Seq (Chen et al., 2023) | 61.48 | 44.57 | 77.72 |
| AutoQGS (Xiong et al., 2022) | 65.13 | 47.50 | 76.80 |
| T5-Large (Raffel et al., 2020) | 58.95 | 44.72 | 76.58 |
| IGND (Fei et al., 2021) | 61.69 | 45.11 | 77.28 |
| LFKQG (Fei et al., 2022) | 63.92 | 46.91 | 78.40 |
| JointGT (Ke et al., 2021) | **65.89** | **48.25** | 78.87 |
| DSM (Guo et al., 2022) | 61.03 | - | **86.06** |
| *Few-shot Evaluation* | | | |
| BiGraph2Seq (Chen et al., 2023) | 1.01 | 4.99 | 12.07 |
| JointGT (Ke et al., 2021) | 43.15 | 35.91 | 69.57 |
| AutoQGS (Xiong et al., 2022) | 43.46 | 33.55 | 68.23 |
| KQG-CoT | 60.81 | 43.54 | 77.35 |
| KQG-CoT+ | 61.71 | 44.27 | 78.41 |

(1) KQG-CoT outperforms existing few-shot methods with large margins;

(2) KQG-CoT achieves competitive results compared with full training methods.

华东师范大学数据科学与工程学院
School of Data Science and
Engineering at ECNU



| Method | Average_similarity |
|---|---|
| Random | 0.285 |
| Active-CoT | 0.274 |
| Auto-CoT | 0.265 |
| KQG-CoT | 0.252 |

(1) KQG-CoT outperforms existing LLM-CoT methods with various k;

(2) KQG-CoT results in supportive logical forms with larger diversity.

# Conclusions

- When constructing prompts, the selection and arrangement of exemplars are paramount.

- KQG-CoT outperforms the existing CoT methods significantly and achieves performance levels comparable to those of fine-tuned methods.

- The utilization of LLMs in conjunction with CoT proves highly effective for handling generation tasks with structured inputs.