# ALiiCE: Evaluating Positional Fine-grained Citation Generation

**Yilong Xu**[1,2,3] **Jinhua Gao**[1,2]* **Xiaoming Yu**[1,2] **Baolong Bi**[1,2,3] **Huawei Shen**[1,2,3] **Xueqi Cheng**[1,2,3]

[1]State Key Lab of AI Safety, Institute of Computing Technology, Chinese Academy of Sciences

[2]Key Lab of AI Safety, Chinese Academy of Sciences    [3]University of Chinese Academy of Sciences

University of Chinese Academy of Sciences

Institute of Computing Technology, Chinese Academy of Sciences

## Why need fine-grained citations?

- A sentence might not be the smallest unit capable of representing an atomic claim, potentially leading to inaccurate evaluations.
- The generated text scope of a single in-line citation often brings ambiguity, which is more common in sentence-level citations.
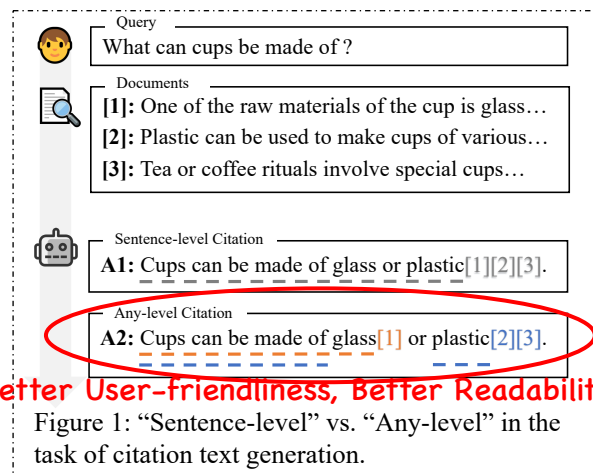


**Better User-friendliness, Better Readability!!**

Figure 1: "Sentence-level" vs. "Any-level" in the task of citation text generation.

We propose this improved task, called <u>Positional Fine-grained Citation Text Generation</u>, but there is no effective method to evaluate it.

**So we introduce ALiiCE to fill this gap.**
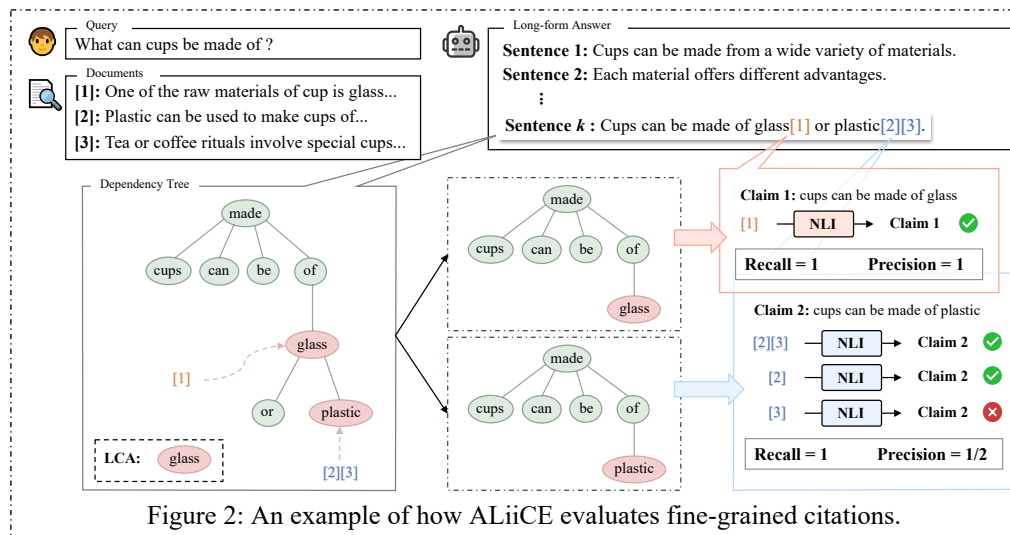
## ALiiCE's pipeline



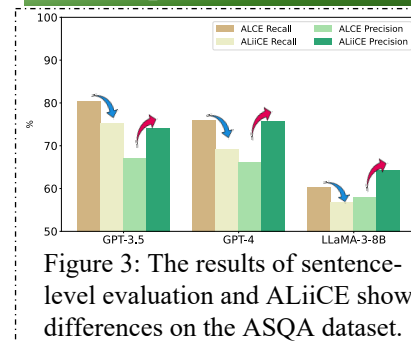Figure 2: An example of how ALiiCE evaluates fine-grained citations.

- ALiiCE employs a Dependency Tree based approach to parse atomic claims of each citation in the response.

## ALiiCE's metrics on citation quality

- Positional Fine-grained Citation Recall
- Positional Fine-grained Citation Precision
- Coefficient of Variation of Citation Positions
  - CVCP measures the dispersion of citation positions within a sentence.

$$\sigma(s_k) = \sqrt{\frac{1}{t}\sum_{j=1}^{t}(p_j - \mu_k)^2} \quad CV_{CP}(\mathcal{R}) = \frac{1}{n}\sum_{k=1}^{n}\frac{\sigma(s_k)}{\mu_k}$$

## Comparison results



Figure 3: The results of sentence-level evaluation and ALiiCE show differences on the ASQA dataset.
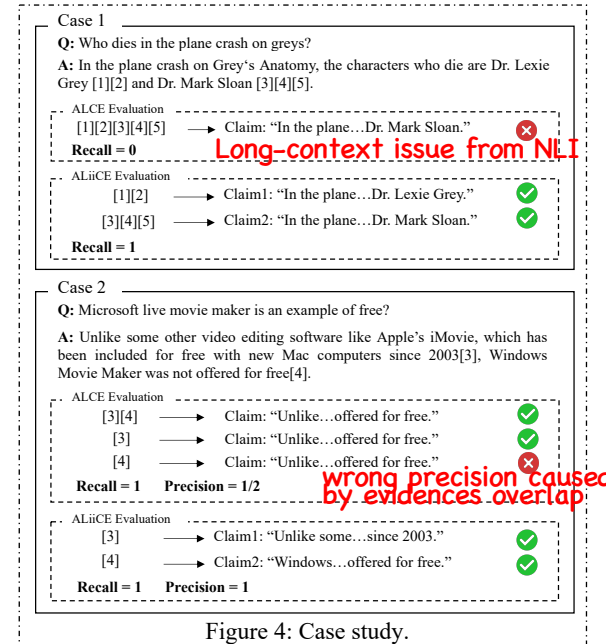
## Case Study



Figure 4: Case study.

## More Insights

- Human evaluation shows strong alignment with ALiiCE.
- ALiiCE has a higher decision threshold.
- Open-source LLMs display great progress.
- Current methods on citation evaluation ignore the judgment of citation utility.