### InfoLossQA: Characterizing and Recovering Information Loss in Text Simplification

Jan Trienes, Sebastian Joseph, Jörg Schlötterer, Christin Seifert, Kyle Lo, Wei Xu, Byron C. Wallace, Junyi Jessy Li

ACL 2024, Bangkok



## Information loss is inherent to text simplification

**Original:** In this study, the benefit of preoperative nutritional support was investigated for *non-small cell lung cancer patients* who underwent anatomic resection. [...] Patients who were malnourished, diabetic or who had undergone *bronchoplastic procedures or neoadjuvant therapy* were excluded from the study.

**Simplification:** This study looked at if eating a protein-rich diet before surgery could help *lung cancer patients* recover more quickly after surgery. [...] We didn't include patients who were already not eating well, had diabetes, or had received other treatments for their lung cancer.

#### Why is this a problem? May lead to...

- Reduced reader comprehension (Agrawal and Carpuat, TACL 2024)
- Factuality errors (Devaraj et al., ACL 2022)

ossv

## Information loss is inherent to text simplification



#### Why is this a problem? May lead to...

- Reduced reader comprehension (Agrawal and Carpuat, TACL 2024)
- Misinterpretation or factuality errors (Devaraj et al., ACL 2022)

### We propose InfoLossQA

Task: generate questions that recover information loss

Q: What kind of lung cancer do these patients have?

A: Non-small cell lung cancer, which is a group of lung cancers that behave similarly. Non-small cell lung cancer, compared to small cell lung cancer, is most common. Simplification: This study looked at if eating a protein-rich diet before surgery could help *lung cancer* patients recover more quickly after surgery. [...] We didn't include patients who were already not eating well, had diabetes, or had received other treatments for their lung cancer.

Original: In this study, the benefit of preoperative nutritional support was investigated for non-small cell lung cancel nations, who underwent anatomic

**Q** asks for missing information; **A** provides it in lay language based on original

## QA builds on Questions Under Discussion (QUD) theory

Qs are rooted in simple text, w/o assuming reader access to original



Q: What kind of lung cancer do these patients have?



Q: Do they consider patients with nonsmall cell lung cancer? Simplification: This study looked at if eating a protein-rich diet before surgery could help *lung cancer patients* recover more quickly after surgery. [...] We didn't include patients who were already not eating well, had diabetes, or had received other treatments for their lung cancer.

# Data: Linguists annotate info loss in LLM simplifications



3 annotators write QAs which can be answered by original, but not by simplification

- 104 parallel simplifications
- |Original| = 312 tokens
- |Simplified| = 271 tokens
- 4 6 QAs/annotator/doc
- 1,000 QAs total

## Automatically identifying information loss



Method 1. End-to-end prompting



Method 2. NLI pipeline

## QAs provide a high-level summary of what's lost

Question Type	Frequency	Examples
Procedural	34.3%	How did they measure the patients opioid medication needs? How did the study control for bias?
Concept	25.7%	What kind of hip surgery were patients undergoing? What type of mental illnesses are being studied?
Extent	17%	How many patients were in each group? How much lurasidone was given to patients?
Comparison	8.3%	How much did abnormal blood vessels reduce compared to placebo?

## Models are not reliable at identifying information loss



## Models often fail to rationalize where info was lost



## Models have low recall of human-generated QAs





Simplification is inherently lossy which may affect comprehension

Solution: use QAs to characterize and recover info loss in a reader-centric way

Current methods struggle to reliably identify info loss

