

SciNews: From Scholarly Complexities to Public Narratives **A Dataset for Scientific News Report Generation**

Dongqi Pu, Yifan Wang, Jia Loy, Vera Demberg

Department of Computer Science
Department of Language Science and Technology
Saarland Informatics Campus, Saarland University, Germany
dongqi.me@gmail.com



UNIVERSITÄT
DES
SAARLANDES



European Research Council
Established by the European Commission

TL;DR

We introduce a **new corpus** designed to enhance the translation of **research papers** into accessible scientific **news reports**

Introduction

- **Why** Study Scientific News Report Generation?
- **Similarities and Differences** with Summarization / Simplification

Introduction

- Why Study Scientific News Report Generation?
- Academic publications → Require background knowledge 🤖
- News reports → Increase accessibility with simplified language 😊

Academic Paper

Abstract Current **techniques for characterizing cybersickness** (visually induced motion sickness) in virtual environments rely on qualitative questionnaires. [...]

Intro With the resurgence of virtual reality (VR), cybersickness has become [...] We establish that cybersickness in an immersive HMD [...] Our approach [...] using inexpensive, commodity off-the-shelf devices for VR headsets and EEG devices. [...] We find a **statistically significant correlation of Delta-, Theta-, and Alpha-waves with self-reported cybersickness.** [...]

Conclusion Throughout the course of the study, we witnessed a wide range of reactions to the rendered stimuli. [...] Our findings in this paper are just a first step to the many opportunities that present themselves in using EEG to study cybersickness in virtual environments. [...] Finally, it will be highly **desirable, if at all possible, to move toward standards of assessing cybersickness and to use them to rate hardware (headsets, trackers, and displays) as well as the content (games, performances, and other immersive experiences).**

News Report

Report If a virtual world has ever left you **feeling nauseous or disoriented, you're familiar with cybersickness**, and you're hardly alone. The intensity of virtual reality (VR) whether that's standing on the edge of a waterfall in Yosemite or engaging in tank combat with your friends [...] They were able to establish **a correlation between the recorded brain activity and self-reported symptoms** of their participants. [...] **This helped the researchers identify which segments of the fly-through intensified users' symptoms.**

An example of an academic paper paired with its news report

Introduction

- Similarities and Differences with Summarization / Simplification
 - Summarization: **reduces** text, retains key content
 - Simplification: uses **simpler** words/syntax for readability
 - Our task involves **both** simplifying and extracting

The SciNews Dataset

- Data Acquisition
- Data Cleaning
- Quality Control
 - Automated Quality Control
 - Human Quality Control
- Data Splits

The SciNews Dataset

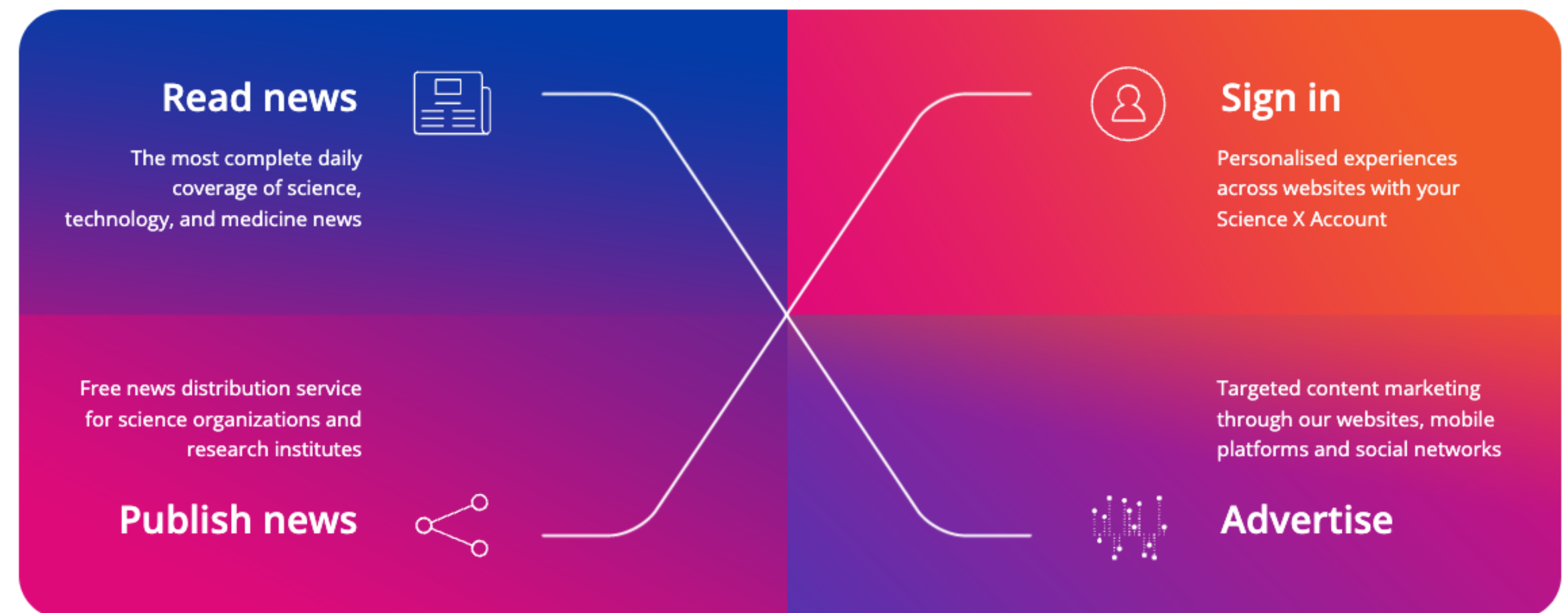
- Data Acquisition

- SciNews sourced from **Science X**

- Open access articles with **CC-BY-4.0** license via DOI



Science X is a network of high-quality websites that provides the most complete and comprehensive daily coverage of science, technology, and medical news.



<https://sciencex.com/>

The SciNews Dataset

- Data Cleaning
 - Use PySBD and spaCy to clean texts; remove line breaks, emoticons, and links etc
 - Extract text from papers between the abstract and references
 - Exclude documents over 30,000 or under 2,000 words

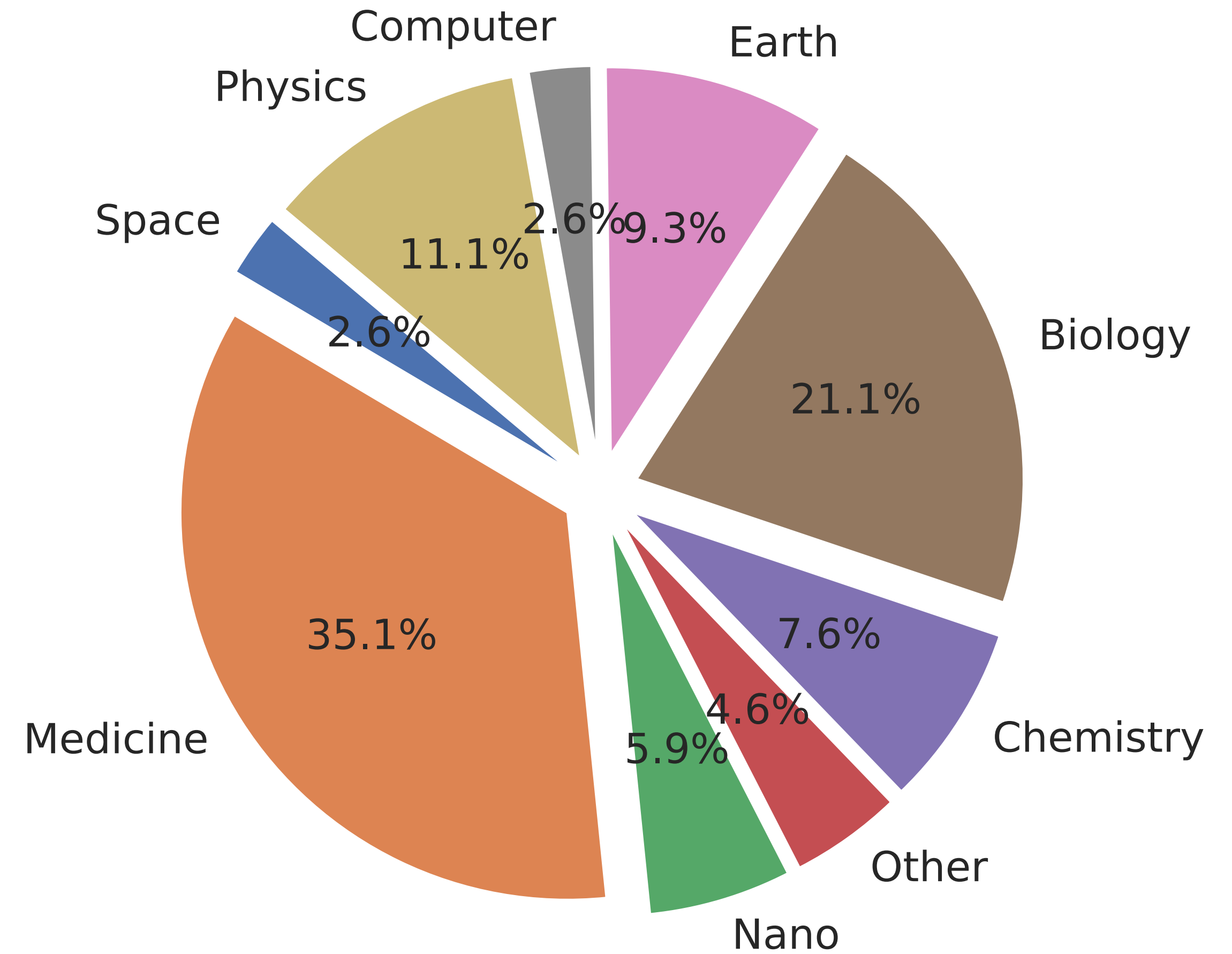
The SciNews Dataset

- Quality Control
 - Automated Quality Control
 - Adapt methods from Mao et al. (2022) for vetting pairs; remove 612 of 42,484 pairs (dissimilarity in BERTscore)
 - Human Quality Control
 - Manual check reveals that 100/100 sampled pairs are of good quality.

The SciNews Dataset

- **Data Splits**

- 41,872 samples
- 80% training, 10% validation, 10% test
- across nine domains



Topic distribution of our dataset

Dataset Analysis

- Dataset Comparison
- Dataset Statistics
- Papers vs. News

Dataset Comparison

- SciNews vs. CSJ & PLOS: similar sizes; SciNews has multidisciplinary labels
- Output Length: SciNews (695 tokens), PLOS (176 tokens), CSJ (361 tokens)

Dataset	Task	Language	Data Scope	Data Source	Scale	Input Level	Output Level	Multi-disciplinary?
LaySumm (Chandrasekaran et al., 2020c)	SLS	English	Archaeology, Hepatology, etc.	Research Papers	572	Document	Paragraph	✓
CDSR (Guo et al., 2021)	SLS	English	Healthcare	Research Papers	7805	Document	Paragraph	✗
CELLS (Guo et al., 2022)	SLS	English	Biomedicine	Research Papers	47157	Sentence	Sentence	✗
eLife (Goldsack et al., 2022)	SLS	English	Biomedicine	Research Papers	4828	Document	Paragraph	✗
PLOS (Goldsack et al., 2022)	SLS	English	Biomedicine	Research Papers	27525	Document	Paragraph	✗
SimpleScience (Kim et al., 2016)	STS	English	Biomedicine	Research Papers	293	Sentence	Vocabulary	✗
CLEAR (Grabar and Cardon, 2018)	STS	French	Biomedicine	Research Papers	663	Sentence	Sentence	✗
PLS (Devaraj et al., 2021)	STS	English	Medicine	Research Papers	4459	Paragraph	Paragraph	✗
SimpleText (Ermakova et al., 2022, 2023)	STS	English	Medicine & Computer Science	Research Papers	648	Sentence	Sentence	✓
CSJ (Fatima and Strube, 2023)	STS	English & German	Astronomy, Biology, etc.	Wikipedia	50132	Document	Paragraph	✓
SciNews (ours)	SNG	English	Science & Technology & Medicine	Research Papers	41872	Document	Document	✓

Dataset Statistics

- Long input & long output
- Highly abstractive (coverage & density)
- High 1/2/3/4-grams novelty

Property	Value
# Training Set	33497
# Validation Set	4187
# Test Set	4188
Avg. # Tokens (Papers)	7760.90
Avg. # Tokens (News)	694.80
Avg. # Sents. (Papers)	290.52
Avg. # Sents. (News)	25.17
Compression Ratio	12.71
Coverage	0.74
Density	0.94
1-gram Novelty	0.52
2-gram Novelty	0.91
3-gram Novelty	0.98
4-gram Novelty	0.99

Papers vs. News

- First-person vs. third-person
- **Lexical diversity:** higher in news
- **Syntax:** simpler in news

Property	Papers	News
Type-Token Ratio↑	0.20	0.44
Lexical Density↑	0.42	0.46
Avg. # Difficult Words↓	773.08	134.84
Avg. # Modifiers per Noun Phrase↓	0.58	0.51
Avg. Depth of Dep Tree↓	6.94	6.25
FKGL↓	14.57	13.31
ARI↓	17.94	16.32

FKGL = Flesch-Kincaid Grade Level

ARI = Automated Readability Index

Experiments

- Baseline Models
- Experimental Settings
- Automatic Metrics

Experiments

- **Baseline Models**
 - **Extractive Methods**
 - Lead-3/K, Tail-3/K, and Random-3/K
 - Latent Semantic Analysis, LexRank, TextRank, Ext-oracle, and PacSum
 - **Abstractive Methods**
 - Longformer, RSTformer, SIMSUM (Seq2Seq)
 - Vicuna7B-16k, GPT-4 (GPT)

Experiments

- **Experimental Settings**
 - **Default Settings:** use original model sizes, batch sizes, optimizers etc
 - **Decoding:** beam search=3, trigram blocking, temperature=1, top-p=1
 - **Vicuna Model:** 5e-5 initial rate, cosine schedule, Adam optimizer, fine-tune 30 epochs

Experiments

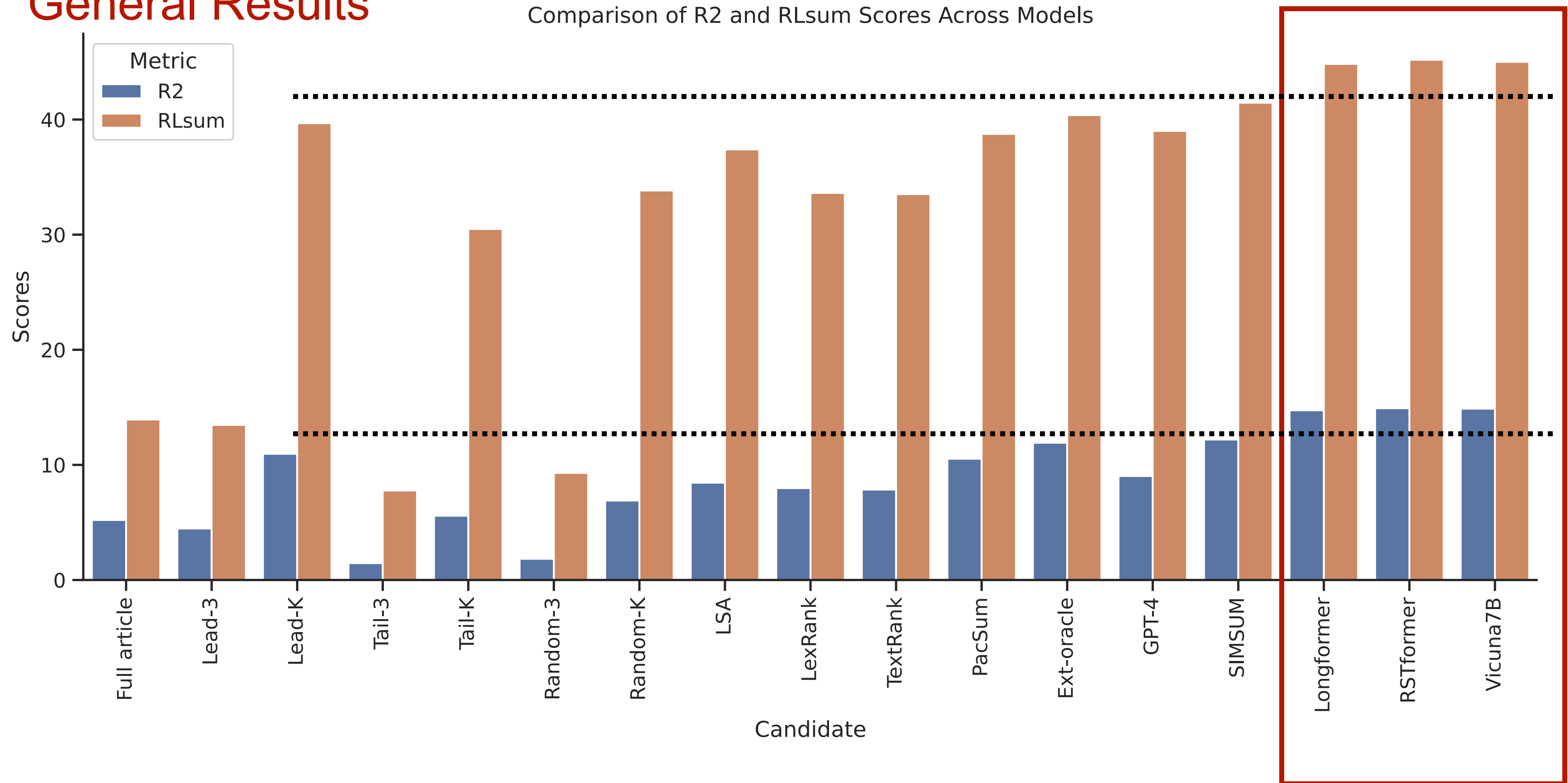
- **Automatic Metrics**
 - F1 scores of Rouge-1 (R1), Rouge-2 (R2), Rouge-L (RL), and Rouge-Lsum (RLsum) (Lin, 2004)
 - BERTScore (Zhang et al., 2020)
 - METEOR (Banerjee and Lavie, 2005)
 - sacreBLEU (Post,2018)
 - NIST (Lin and Hovy, 2003)
 - SARI(Xu et al.,2016)

Results and Analysis

- General Results
- Comparison with Human-authored News Articles
- Automatic Inconsistency Detection
- Human Evaluation
- GPT-4 Evaluation
- Model Errors

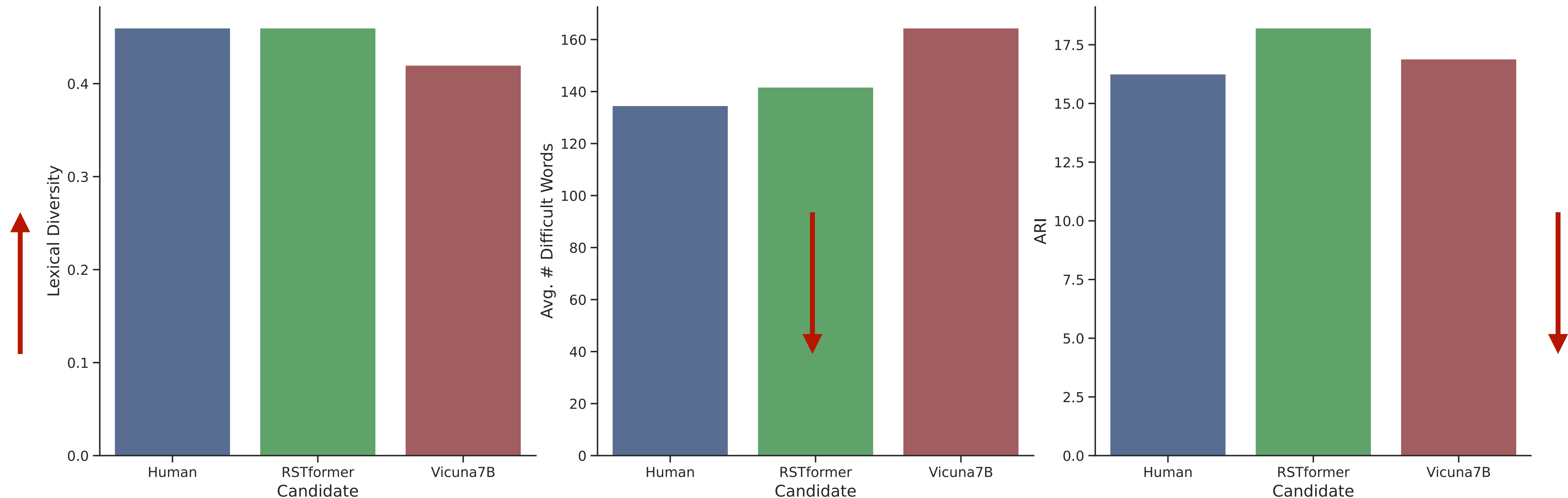
Results and Analysis

- General Results



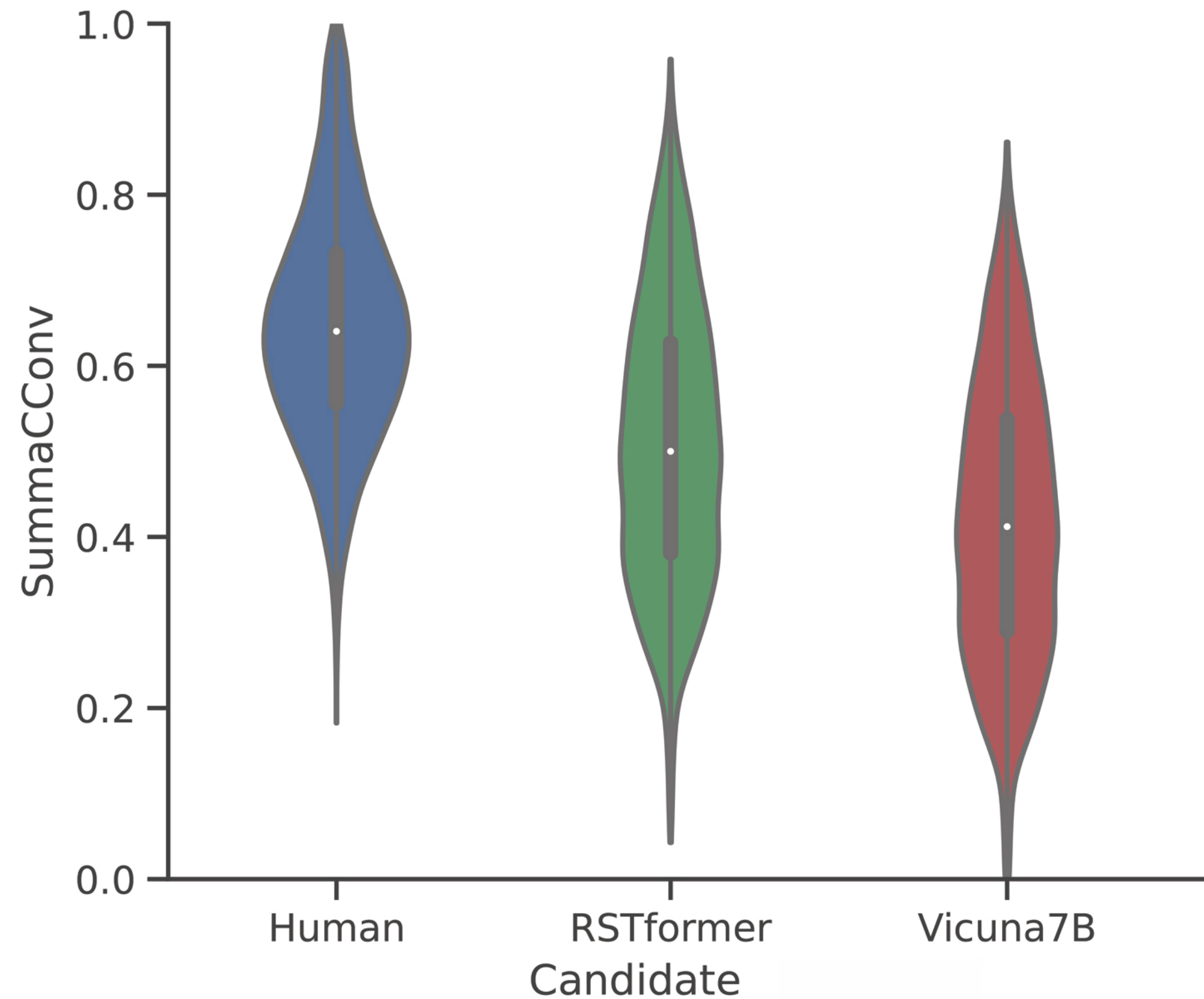
Results and Analysis

- **Comparison with Human-authored News Articles**
 - **Lexical Diversity:** RSTformer closest to human
 - **Complexity:** Vicuna generates more complex words
 - **Readability:** Humans outperform models (ARI)



Results and Analysis

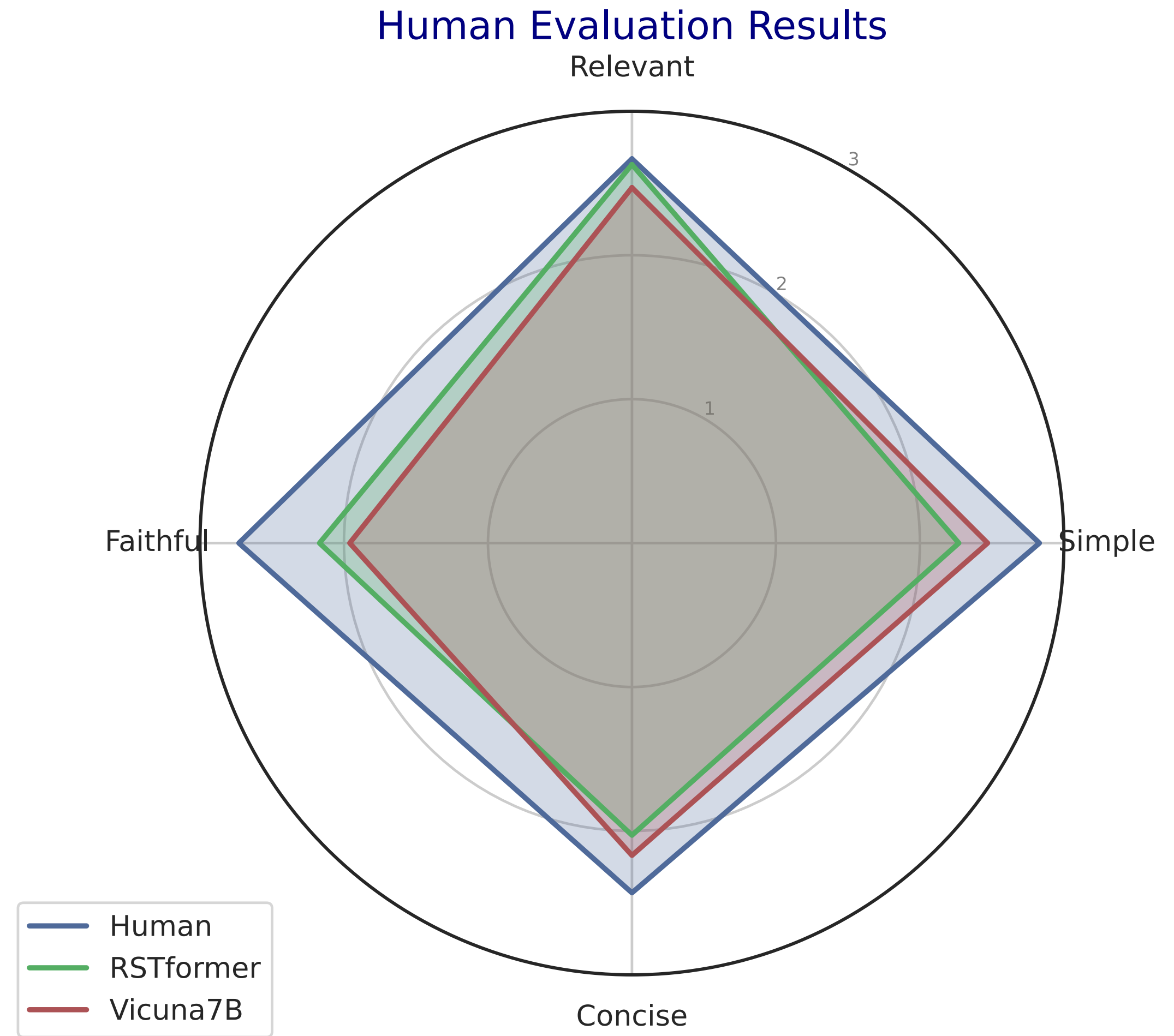
- Automatic Inconsistency Detection
- Abstractive models have lower consistency scores than humans



Results and Analysis

- **Human Evaluation**

- **Evaluation Setup:** 10 samples, blind testing by Masters/PhD evaluators
- **Criteria:** relevance, simplicity, conciseness, faithfulness; scored 1-3
- **Results:** RSTformer and Vicuna excel in different areas; overall, models lag behind human proficiency

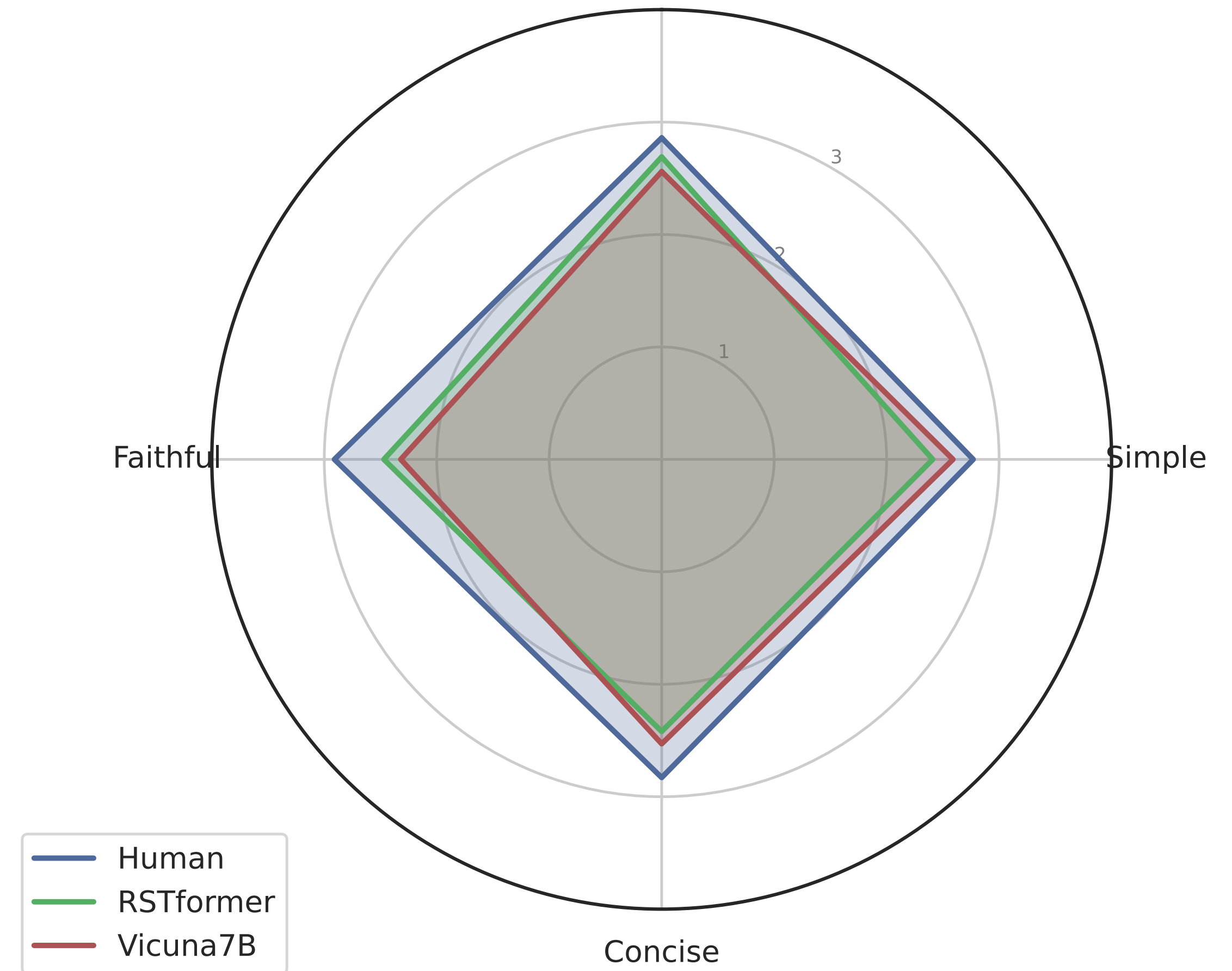


Results and Analysis

- **GPT-4 Evaluation**

- Uses human evaluation guidelines, resets history for unbiased assessment.
- **Preliminary Check:** GPT-4 and human scores align across criteria.
- **Overall Findings:** Humans outperform all models, RSTformer is better than Vicuna

GPT-4 Evaluation Results on 100 Samples
Relevant



Results and Analysis

- Model Errors

- Hallucinations

- The study identified issues with misinformation in apps, stressing the importance of current and reliable content; future research should investigate strategies to ensure accuracy and quality in chatbot app development.

- Factual Errors

- SkinVision was evaluated in studies, achieving 88.2% sensitivity and 98.3% specificity in detecting malignant or premalignant lesions.

- Generalization

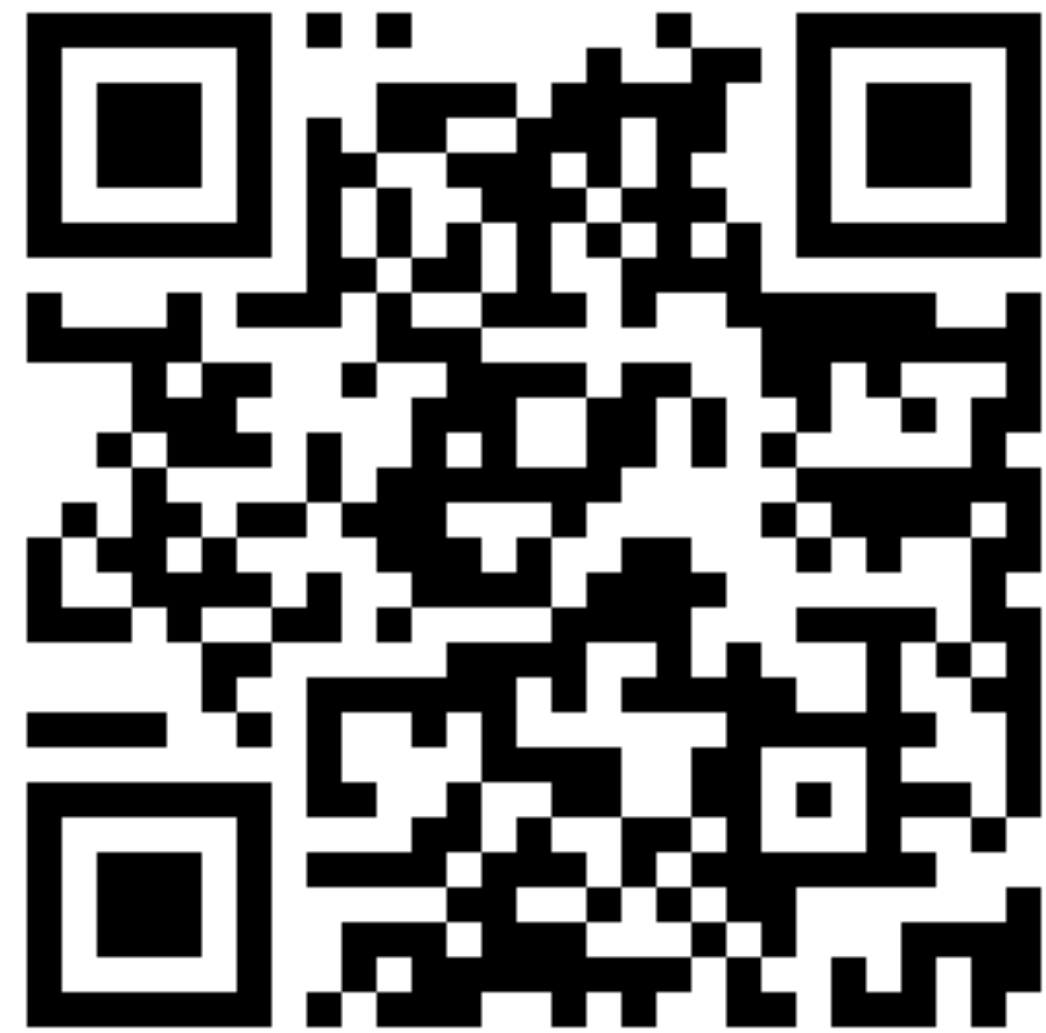
- Research has shown the negative effects of cybersickness on users including autistic individuals and adults with disabilities. It is important to understand how these impacts vary with different VR applications.

Conclusion

- **Dataset Introduction:** "SciNews" comprises 40,000+ scientific papers with paired news reports
- **Exploratory Analysis:** Reveals challenges and research prospects for state-of-the-art models
- **Dataset Potential:** Enhances scientific news generation, offers resource for NLP tasks like topic classification

More Info

- **Data & Code:** <https://dongqi.me/projects/SciNews>
- **Questions:** dongqi.me@gmail.com



Thanks for listening

Q&A



European Research Council
Established by the European Commission

This project has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 Research and Innovation Programme (Grant Agreement No. 948878)