

# Scaling Multi-Document Event Summarization

## Evaluating Compression vs. Full-Text Approaches

Carnegie  
Mellon  
University

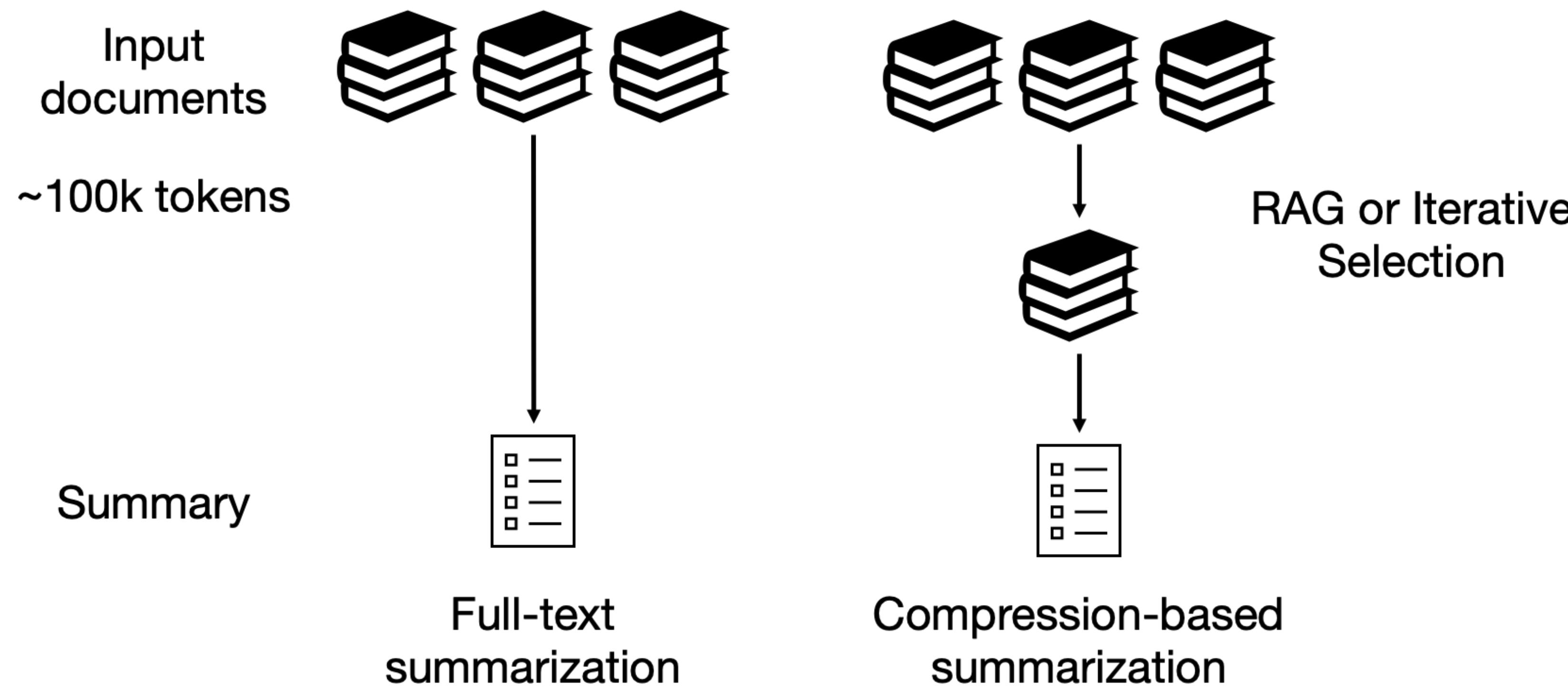
Adithya Pratapa\* Teruko Mitamura

\*On the industry job market!

{vpratapa, teruko}@cs.cmu.edu



— Full input (w. long context window) or compressed input (w. short context window)? —



### Compression

- Efficient
- Scalable
- Retrieval, Hierarchical, Incremental



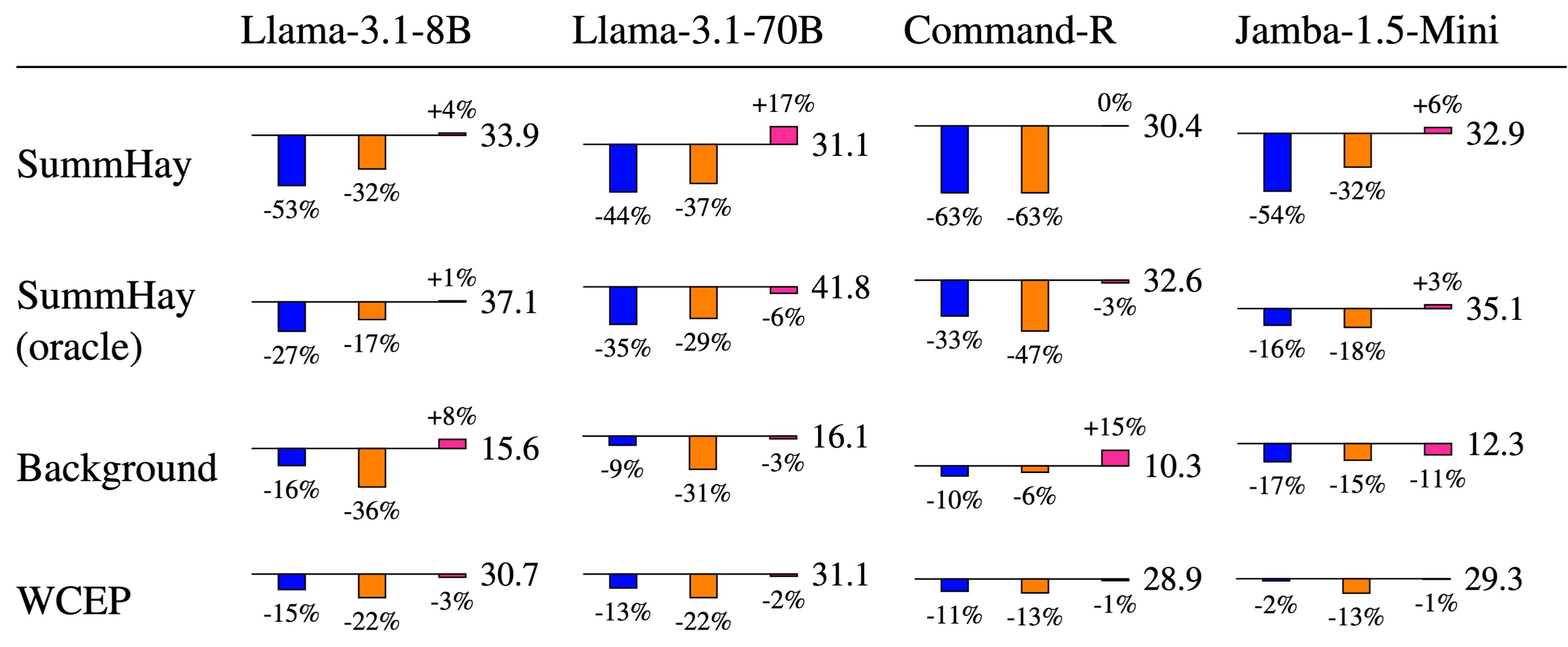
### Full-text

- Effective
- Utilizes long-context windows of recent LLMs

## Base models

- Llama-3.1 (8B, 70B)
- Command-R (32B)
- Jamba 1.5 Mini (52B MoE, 12B active)

## Results

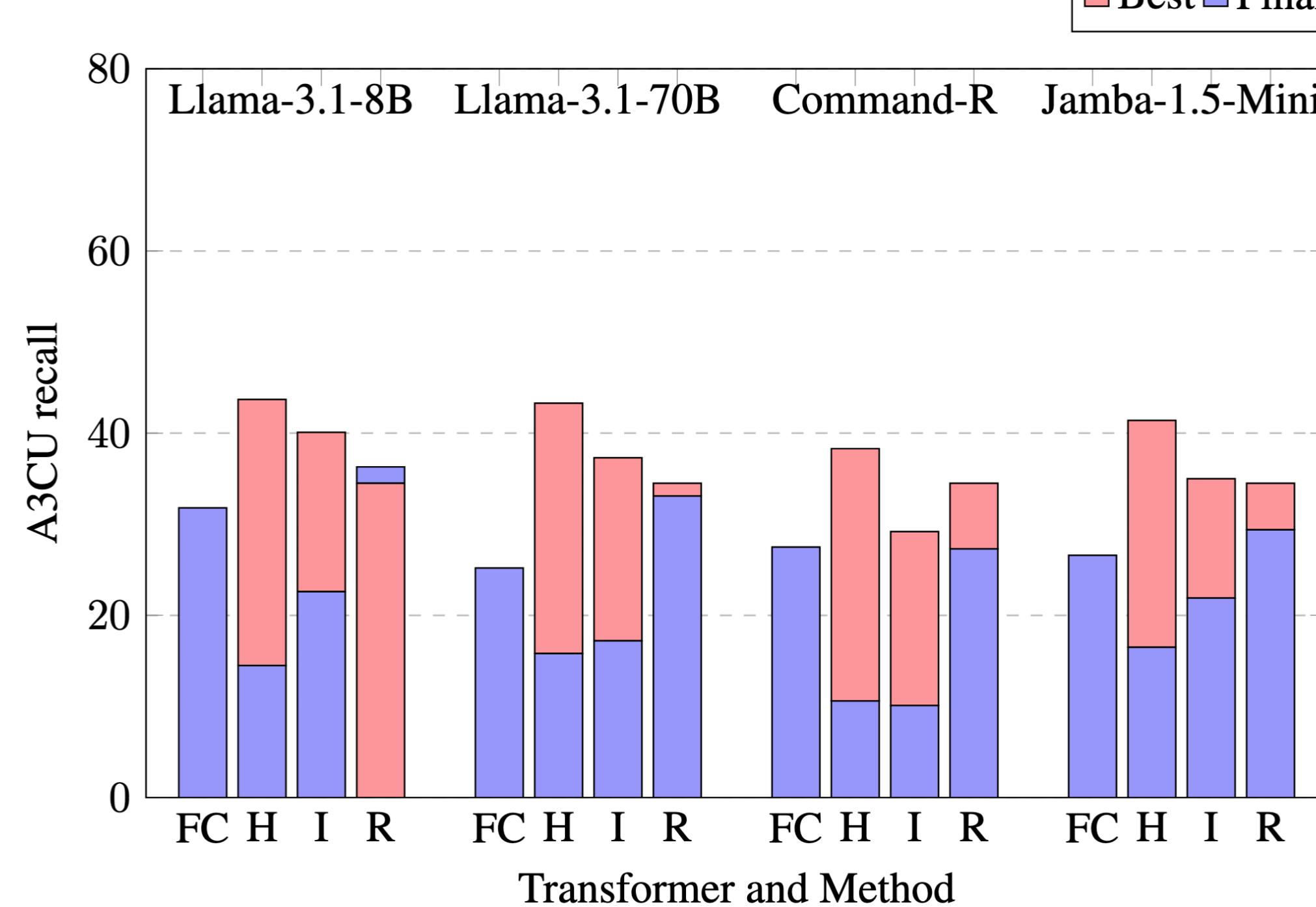


Relative performances of **hierarchical**, **incremental** and **retrieval** compared to full-text baseline (AutoACU F1)

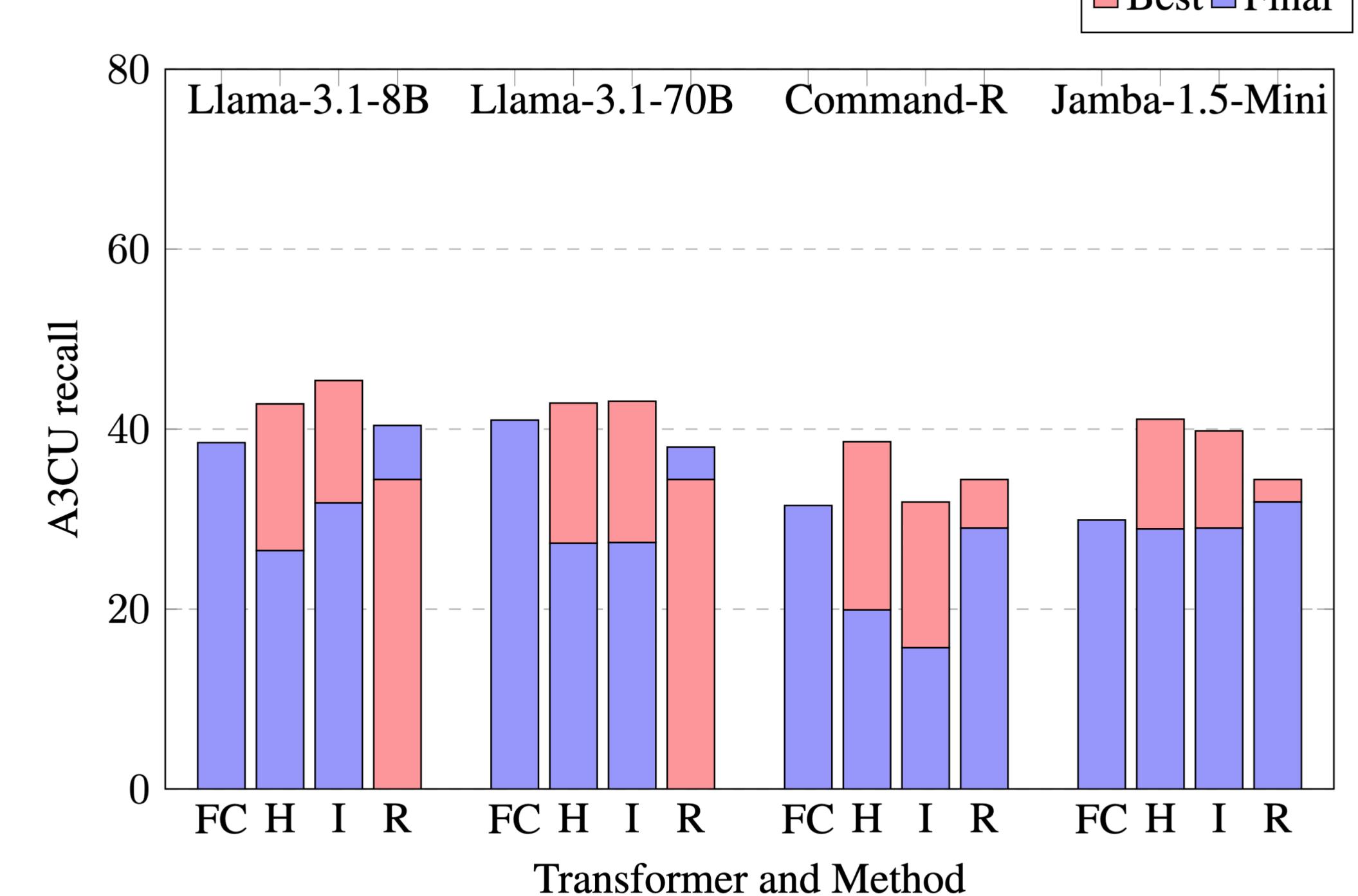
## Tracking information loss

We measure A3CU recall scores of the final system summary and the *best* intermediate output.

SummHay



SummHay (Oracle)



## Proposals for future work

- Hybrid approaches that combine compression with long-context windows could be more efficient and effective.
- We need automatic metrics that do not rely on human-written references.

