#### CLEF 2026 SimpleText Track Simplify Scientific Text (and Nothing More)

<u>Liana Ermakova</u><sup>1</sup> Hosein Azarbonyad<sup>2</sup> Jan Bakker<sup>3</sup> Benjamin Vendeville<sup>1</sup> Jaap Kamps<sup>3</sup>

<sup>1</sup>Université de Bretagne Occidentale <sup>2</sup>Elsevier <sup>3</sup>University of Amsterdam









CLEF, September 21-24, 2026, Jena, Germany

#### Motivation



- Simplify Scientific Text
  - Everyone agrees on the importance of objective scientific information
  - But scientific documents are inherently complex...
  - Can we improve science literacy for everyone?
  - Generative models for text simplification can help!
- And Nothing More!
  - LLMs prone to overgeneration (informally called "hallucination")
  - Have created a huge collection of spurious/over-generation content!
  - $\bullet$  In 2025: 29% of submissions > 25% spurious sentences, 21% > 50%...

#### Overview



- SimpleText Track setup similar 2021-2024
  - Major changes in setup and corpora in 2025
  - Boost in participation and results
  - Will continue into 2026
- CLEF 2026 SimpleText Track
  - Simplify Scientific Text (and Nothing More)
- The following three tasks:
  - **1 Text Simplification**: simplify scientific text
  - **2** Controlled Creativity: identify and avoid hallucination
  - **3** SimpleText 2024-2025 Revisited: selected tasks by popular request



## Task 1: Text Simplification



- Task 1: Simplify Scientific Text
  - New corpus (EMNLP/TSAR 2024)!
    - Cochrane-auto is true document-level text simplification
    - More variation (sentence merge, order swaps) and discourse structure
    - Paragraph-level and sentence-level data realigned and restricted
  - Biomedical text free to use, similar to existing TS corpora
    - Sentence-level (T1.1) and Document-level (T1.2) text simplification
    - Large-scale aligned train and test data (9,160 sentences, 666 abstracts)
- Bonus: Cochrane Abstracts/Plain Language Summaries in English
  - + Spanish (es), French (fr), Farsi (fa), Chinese (zh), Japanese (ja), Portuguese (pt), Korean (ko), Thai (th), German (de), Russian (ru), Malay (ms), and Croatian (hr).
  - + Abstracts and freely usable full-text systematic reviews



## Task 2: Controlled Creativity



- Task 2: Identify and Avoid Hallucination
  - Task 2.1: to identify creative generation at document level
    - to detect what sentences are fully grounded on source input (a) without and (b) with access to the source sentences
    - ullet ightarrow also labels those introducing significant new content
  - Task 2.2: to detect and classify information distortion errors in simplified sentences
    - 14 information distortion categories (Fluency, Alignment, Information, Simplification)
    - Can LLMs do do detailed "human evaluation"?
  - Task 2.3: to avoid creative generation and perform grounded generation by design
    - Submit pairs of simplified abstracts with/without "hallucinations"

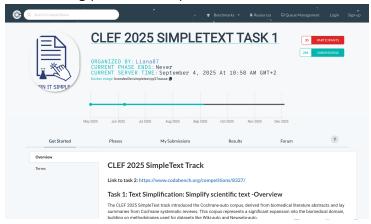


## Task 3: SimpleText 2024-2025 Revisited





- Task 3: Selected Tasks by Popular Request
- Move to Codabench makes it easy to keep tasks running in 2026
  - Codabenches are still active in post-competition mode!
- Consider adding pilot tasks for possible extensions in 2026







## CLEF 2026 SimpleText Track



#### Simplify Scientific Text (and Nothing More):

- Task 1: Text Simplification: simplify scientific text
  - + New aligned biomedical data (Cochrane-auto)
  - ullet + both sentence, paragraph and document level simplification
  - + analysis of information distortion ("hallucination?")
- Task 2: Controlled Creativity: identify and avoid hallucination
  - ullet + Real "hallucination" data from CLEF generative text tasks
  - + What output is (not) grounded on source(s)? (w/wo source access)
  - + Fine-grained information distortion categorization
- Task 3: SimpleText 2024-2025 Revisited: selected tasks by popular request
  - We take submissions for earlier tasks

### Join us at CLEF 2026 in Jena!





#### **SimpleText**

Over the last few years, the Simple Text Track has created an active community of researchers in NLP and IR working together to improve access to scientific text. Its benchmarks on scientific passage retrieval, scientific terminology detection and explanation, and scientific text simplification have become standard references. After using a similar track setup in 2021-2024, we significantly changed the track's setup and tasks in 2025. We will continue this successful setup in 2026. Hence, the CLEF 2026 Simple Text track will continue this successful setup in 2026. Hence, the CLEF 2026 Simple Text track will continue this successful setup in 2026. Hence, the CLEF 2026 Simple Text track will continue this successful setup in 2026. Hence, the CLEF 2026 Simple Text track will continue this successful setup in 2026. Hence, the CLEF 2026 Simple Text track will continue this successful setup in 2026. Hence, the CLEF 2026 Simple Text track will continue this successful setup in 2026. Hence, the CLEF 2026 Simple Text track will continue this successful setup in 2026. Hence, the CLEF 2026 Simple Text track will continue this successful setup in 2026. Hence, the CLEF 2026 Simple Text track will continue this successful setup in 2026. Hence, the CLEF 2026 Simple Text track will continue this successful setup in 2026. Hence, the CLEF 2026 Simple Text track will continue this successful setup in 2026. Hence, the 2026 Simple Text track will continue this successful setup in 2026. Hence, the 2026 Simple Text track will continue this successful setup in 2026. Hence, the 2026 Simple Text track will continue this successful setup in 2026. Hence the 2026 Simple Text track will continue the 2026 Simple Text track will be 2026 Simple Text track will

#### **Organizers**

- Liana Ermakova (HCTI, Université de Bretagne Occidentale
- Benjamin Vendeville (Université Bretagne Occidentale)











# Please join the SimpleText Track

## Fully funded PostDoc available!

Website: https://simpletext-project.com E-mail: contact@simpletext-project.com

Twitter: https://twitter.com/SimpletextW

 ${\sf Google\ group: https://groups.google.com/g/simpletext}$