

Accepted by the Trans. of the Association for Computational Linguistics (TACL 2025)



European Research Council Established by the European Commission





- **Explanatory Summarization with Discourse-Driven** Planning
 - **Dongqi Liu**^{Ω}, Xi Yu^{Ω}, Vera Demberg^{Ω}, Mirella Lapata^{Θ}
 - ^ΩSaarland University, Germany [©]University of Edinburgh, The United Kingdom dongqi.me@gmail.com





We propose EDU-level rhetorical planning using discourse structure and question-based cues to control explanatory content generation in lay summarization.

TL;DR





- Why Do Lay Summaries Need Explanations?
- What Are Current Summarization Models Missing?
- Why Is Discourse-driven Planning a Promising Solution?
- What Are the Challenges We Address?





- Why Do Lay Summaries Need Explanations?
 - Scientific concepts presented in academic documents are often too complex for non-experts to understand
 - Human-written lay summaries often contain analogies, causal justifications, and background
 - Just simplifying language (e.g., shorter sentences) may lead to loss of meaning or misinterpretation





- What Are Current Summarization Models Missing?
 - Many models treat summarization as a flat end-to-end task without explicitly modeling explanations
 - Current models underproduce explanations, yielding summaries that are less clear, less accessible than human lay summaries





- Why Is Discourse-driven Planning a Promising Solution?
 - Discourse structures help to identify explanatory sentences and their rhetorical function
 - Planning offers controllability, enabling models to decide <u>what</u> and where to explain
 - Question-based plans naturally trigger explanation generation





- What Are the Challenges We Address?
 - Lack of gold explanation annotations → need for automatic method (via discourse + LLMs)
 - Explanations are hard to evaluate automatically; many are misclassified as hallucinations



- Rhetorical Structure Theory (RST)
 - Models discourse as a tree of Elementary Discourse Units (EDUs)
 - Connects via rhetorical relations
 - Defines nuclearity structure: nucleus (central) vs. satellite (supportive)
 - Reveals explanatory roles like Justification and Background







- Question Under Discussion (QUD)
 - Models discourse via a stack of implicit questions
 - Each sentence resolves a current question in context
 - Adds an intentional layer to discourse modeling

Prerequisite





- **Objective:** Generate lay summaries with explicit, controllable explanations
- Strategy: Use planning to guide both where and how to insert explanations
- Foundation: Leverage Rhetorical Structure Theory (RST) and the Question Under Discussion (QUD) framework

Method Overview





Planning Pipeline

- **Step 1**: Automatically extract explanatory EDUs using **DMRST** parser
- **Step 2:** For each explanatory EDU, generate a corresponding plan question using GPT-40
- Step 3: Construct a "plan" as an ordered list of questions, each prompting an explanation in the summary

Processed summary

[The cerebellum utilizes proprioceptive feedback to fine-tune the timing of movements in a sequence based on previous actions.] t_1 [Imagine the cerebellum as a coach who watches how you perform a move, then gives tips to improve the next one based on what was seen. $]^{e_1}$ But how exactly does it achieve this? [To investigate, we trained rabbits to blink in response to an external cue and explored whether the cerebellum could use feedback from one blink to trigger the next. $]^{t_2}$ [As expected, after learning the initial blink, the rabbits blinked again in response to their own first blink, creating a chain of movements.] e_2 Control experiments confirmed that each blink was initiated by the previous one rather than the original cue. Consistent patterns of brain activity during this process indicate that the cerebellum adjusts movement based on feedback from previous actions. [Building on this, we trained rabbits to blink on cue, and they learned to initiate additional blinks in response to earlier blinks in the sequence. $]^{t_3}$ [We further found that the rabbits could use a blink from one eye as a cue to trigger a blink in the other eye, suggesting that the same mechanism governs these movements. $]^{e_3}$ This raises the possibility that the cerebellum might also guide sequences of cortical activity during cognitive tasks, given its extensive connections to the cortex, a question future experiments should explore.

1	{ Planning que
	al. How does the
	q1. now does the
	q2: How does the
1	q3: How can a bli
	L

Target EDU

Explanatory EDU

estions

cerebellum use feedback to adjust the timing of movements in a sequence? cerebellum use feedback from one blink to trigger the next in a sequence? ink in one eye trigger a blink in the other eye?







- Plan-Output Model
- Plan-Input Model
 - generation

Model Variants

Jointly generates plan questions and summary in a single sequence

• Two-stage approach: first generate plan, then use it to guide summary





• Three Lay Summarization Benchmarks

DATASET	# TRAINING	# VALIDATION	# Test	AVG. DOC TOKENS	AVG. SUMM TOKENS	COVERAGE	DENSITY	COMPRESSION RA
SciNews	33,497	4,187	4,188	7,760.90	694.80	0.74	0.94	12.71
eLife	4,346	241	241	7,833.14	383.02	0.82	1.77	20.52
PLOS	24,773	1,376	1,376	5,340.58	178.66	0.07	0.90	36.06

Experiments







- Discourse Parsing
 - **RST** Parser
 - reference summaries
 - instability

Utilized DMRST for identifying explanatory EDUs and targets in

(Full) Random replacement (FRR/RR) applied to simulate parser





- Alternative Parsing Strategies
 - Rule-based extraction (Stede et al. 2017)
 - LLaMA-based RST parser (Maekawa et al., EACL 2024)
 - RST-Coref parser (Guz & Carenini, CODI 2020)
 - GPT-40 and Mistral as zero-shot extractors
- Plan Generation
 - GPT-40 generates plan questions for each explanatory EDU





- Backbone
 - tuned)
 - Baselines
 - Other SOTAs

All candidate models built on Mistral-7B-Instruct-v0.3 (fully fine-

Backbone w/ zero-shot, in-context, and vanilla fine-tuning





- Automatic Metrics
 - **ROUGE-2**, **ROUGE-Lsum** (*informativeness*)
 - **BERTScore** (semantic similarity)
 - **D-SARI**, **FRE** (*readability*)
 - **ExpRatio** (proportion of explanations)
 - SummaC, SummaC^{*} (*factual consistency*, incl. external verification)
 - **VeriScore** (knowledge-grounded claim verification)





- Human & LLM-Based Evaluation
 - Six-dimension Likert scoring (Faithfulness, Relevance, Usefulness)
 - LLM-as-Judge (GPT-4o) for large-scale validation

Informativeness, Accessibility, Explanation Accuracy, Explanation



Main Results





Parser Choice Analysis

46.6

46.4

46.2

uns_{45.8}

45.6

45.4

45.2

DMRST

- Summary quality improves with more accurate discourse parsers
- DMRST parser gives best overall performance among other alternatives
- RST-based planning is robust to parser variation but degrades with random/noisy parsing



RR = random replacement, FRR = full random replacement











Plan Question Quality

 Summary quality is directly 46.4 impacted by the relevance of 46.2 plan questions uns_{45.8} 45.6 Robust to small noise, but 45.4 random/irrelevant plans 45.2 sharply reduce performance

GPT-40

46.6



---- Mistral_{FT} ---- Blueprint ---- SOTA RAST (Gou et al., EMNLP 2023) is a SOTA question generation method. RR = random replacement, FRR = full random replacement





Controllability

 Removing plan questions for specific relations directly reduces corresponding explanation types in output

Control Effectiveness on SciNews Dataset







- Plan-Input achieves highest human ratings among neural models for faithfulness, relevance, informativeness, accessibility, and explanation usefulness
- Human-written summaries remain best overall, but Plan-Input is most competitive

Human Evaluation Faithfulness Explanation Rèlevance Usefulness Explanation Informativeness Accuracy Best/Worst Human (84.4%/0.0%) PLAN-INPUT (11.1%/8.9%) Blueprint_{*MT*} (4.5%/15.6%) Mistral_{FT} (0.0%/22.2%) Accessibility GPT-40_{ZS} (0.0%/53.3%)



- LLM-based (GPT-4o) evaluation aligns with human ratings
- Plan-Input is consistently rated highest among models
- GPT-4o assigns lowest quality scores to its own generations

LLM-as-Judge Evaluation





Conclusion

- controlled generation of explanatory lay summaries.
- factual consistency, and explanation diversity across multiple datasets.
- alignment with human-written summaries.

• We propose a discourse-driven, plan-based method that enables

• Our models achieve state-of-the-art performance in summary quality,

Planning at the EDU level allows fine-grained control and robust



More Info

- Code: https://dongqi.me/projects/ExpSum
- Questions: dongqi.me@gmail.com



Thanks for listening Q&A



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). We acknowledge the inspiring environment of TRR 248 funded by DFG (German Research Foundation) – Project Number 389792660. Lapata acknowledges the support of the UK Engineering and Physical Sciences Research Council (Grant EP/W002876/1).



Established by the European Commission