# Discourse Models with Language Models

**Jessy Li**
The University of Texas at Austin

# Long form generation: then & now

*Humans were not the only intelligent life in the universe. Instead, they were just the first. They explored the universe and helped build many budding civilizations until one day, the human species disappeared. You are an alien historian that figured out what happened after years of thorough investigation.*

r/WritingPrompts

## Llama-3.2-1b

*As an alien historian, I have spent my entire career studying the history of the universe, from the earliest civilizations to the most recent technological advancements. But none of my research has prepared me for what I discovered on a remote planet in the outer reaches of the galaxy. It was a planet that had been largely unexplored by humans, and I had spent years searching for it. Finally, after months of searching, I found it. The planet was called Zha'thik, and it was a world unlike any other…*

## GPT-2

*What you found was true as well, and your descendants would surely find it too hard not to continue with their journey. (The author, Christopher Nolan, is probably not a big fan of the Alien films, so he has no idea if he can avoid some of his own clichés.) This book is a fantastic adventure that has had a huge impact on me at times, but it could be more important if it could be told right — from where I am writing my story, across the universe…*

2

# Summarization: then & now

*Input: Editorial discussing the symbolism of lifting mask mandates*

**GPT-3**

*The article discusses the debate around when the Covid-19 pandemic will be considered over. It argues that in political terms, the pandemic will be considered over when all states have lifted […]*

**BRIO**

*Hawaii will become the last state in the country to end its indoor mask mandate. The Covid-19 pandemic has seen a surge in cases of the virus this winter. […]*

Explains all significant parts …

Captures the tone of the article […]

Reflects the tone of the article the least

Slide credit: Tanya Goyal

3

# But: LLM texts are somewhat templatic...

**Nataniel Ruiz** ✓
@natanielruizg

AI generated writing *feels* AI-generated at a visceral level, and even if you ask an LLM to make the writing feel or read less AI-generated it horrifically fails and makes it feel even more AI-generated. Any tricks that can help? Any prompts to share?
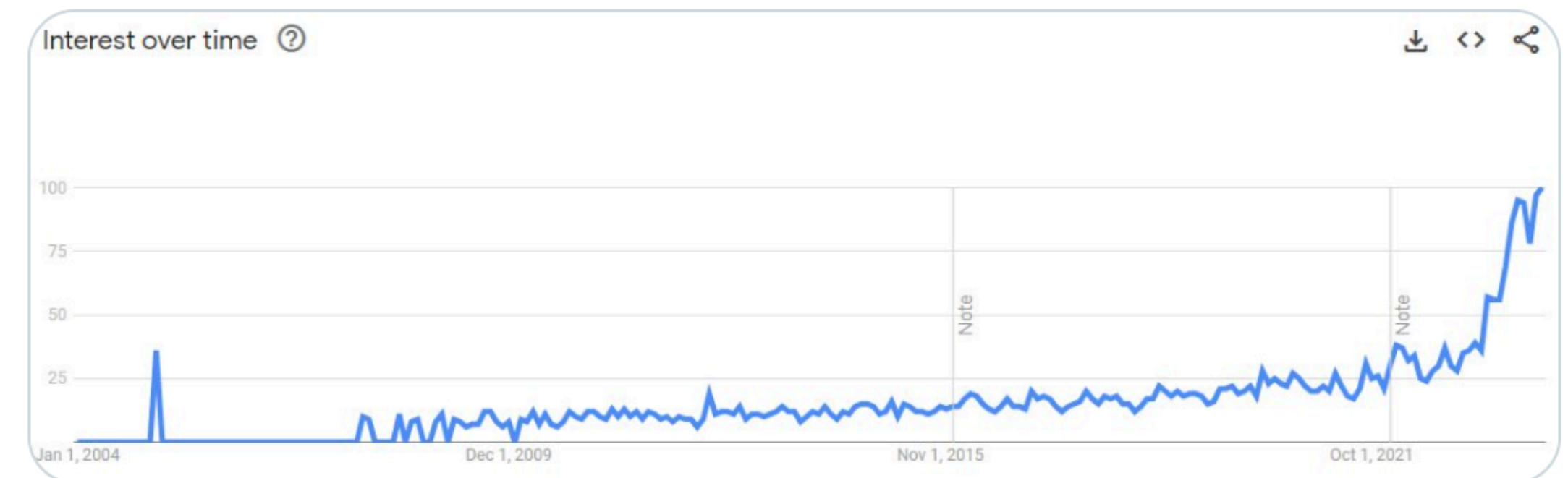
3:11 PM · Mar 28, 2024 · **194.9K** Views

**Toomas Hinnosaar** @toomash · Feb 28

ChatGPT annoyingly tries to insert "delve into" in every paragraph! I am immediately suspicious when I see students' work with this term. Are there any other tell-tale phrases? Here is the Google Trends graph for "delve into":

Interest over time ⑦

100
75
50
25
Jan 1, 2004          Dec 1, 2009          Nov 1, 2015          Oct 1, 2021

💬 50          ⟲ 47          ❤️ 395          📊 86K

# Syntactic Templates in Generated Text

## LLM or human?

Hand model Emily Grimson, 25, from Shropshire, gets paid £100 an hour. Stars in ads for Bosch, the Disney Channel and magazine features.

Emily Grimson, a 25-year-old professional hand model from Shropshire, has been modelling for a year and has built a full-time career earning up to £100 an hour.

Restaurant owner Sandy Dee Hall, 34, and his girlfriend Maxine Cher, 24, adopted Smokey Da Lamb after he was abandoned by his mother.

Smokey Da Lamb, a lamb abandoned by his mother, has been living with restaurant owner Sandy Dee Hall and his girlfriend Maxine Cher in New York City for the past month

John Caudwell has bought a car park in Audley Street, Mayfair and is planning to replace it with a huge block of flats.

John Caudwell, the founder of Phones4U, has unveiled plans to build a £2billion luxury housing complex on the site of an ugly multi-storey car park in Mayfair, central London.

# Syntactic Templates in Generated Text

## LLM or human?

Hand model Emily Grimson, 25, from Shropshire, gets paid £100 an hour. Stars in ads for Bosch, the Disney Channel and magazine features.

Emily Grimson, a 25-year-old professional hand model from Shropshire, has been modelling for a year and has built a full-time career earning up to £100 an hour.

Restaurant owner Sandy Dee Hall, 34, and his girlfriend Maxine Cher, 24, adopted Smokey Da Lamb after he was abandoned by his mother.

Smokey Da Lamb, a lamb abandoned by his mother, has been living with restaurant owner Sandy Dee Hall and his girlfriend Maxine Cher in New York City for the past month

# Detection and Measurement of Syntactic Templates in Generated Text

**Chantal Shaib**[1]    **Yanai Elazar**[2,3]    **Junyi Jessy Li**[4]    **Byron C. Wallace**[1]

[1]Northeastern University, [2]Allen Institute for AI,
[3]University of Washington, [4]The University of Texas at Austin

{shaib.c, b.wallace}@northeastern.edu

# Demo @ https://ai-templates.app/

Detection and Measurement of Syntactic Templates in Generated Text [Shaib et al., EMNLP'24]

- Text makes sense at the beginning but fails to stay relevant after a while.
- Necessary to supervise to obtain task-specific capabilities like summarization.
- Fails long-form planning.

**NOW**

- Long, coherent generations, concepts connect across passages.
- Can do conditional generation pretty well, e.g., knows what is important when asked to summarize.
- Demonstrates some capability to plan for long texts, but "feels" synthetic.

Linguistically driven discourse models can help us understand what we see:

- Characterizing relationships across sentences and document-level structure

- Ground analyses in human language processing

- Provide workable tools for model insights and downstream applications

# This talk

**Can we use language models to make sense of discourse?**

- A generative discourse model based on Questions Under Discussion (QUD) theories
- Connections between human curiosity and discourse planning

**Can we use our discourse framework to make sense of model behavior?**

- A measure of discourse similarity in LLM-generated texts
- An explainable metric of information salience in LLMs

# This talk

**Can we use language models to make sense of discourse?**

- A generative discourse model based on Questions Under Discussion (QUD) theories
- Connections between human curiosity and discourse planning

Can we use our discourse framework to make sense of model behavior?

- A measure of discourse similarity in LLM-generated texts
- An explainable metric of information salience in LLMs

# Coherence: a definition

A text is coherent if:

(1) each newly introduced proposition is rhetorically connected to another piece of information, and

(2) all anaphoric expressions can be resolved.

**Asher and Lascarides, Logics of Conversation, 2003**

# Discourse is (was?) useful!

## Sentiment analysis

While always complaining that he hates this type of movies, John bitterly confessed that he enjoyed this movie.



Hogenboom et al., Using Rhetorical Structure in Sentiment Analysis, 2015

# Questions Under Discussion (QUD)

## A well-known linguistic framework in discourse & pragmatics

Each sentence in discourse addresses a (often implicit) QUD either by answering it, or by bringing up another question that can help answering that QUD. The linguistic form and the interpretation of a sentence, in turn, may depend on the QUD it addresses.

**Benz and Jasinskaja, Questions under discussion: From sentence to discourse. 2017**

Velleman and Beaver. Question-based models of information structure. 2016.

Roberts, Information structure: Towards an integrated formal theory of pragmatics. 2012

Roberts. Context in dynamic interpretation. 2004

Jonathan Ginzburg. Dynamics and the semantics of dialogue. 1996

Jan Van Kuppevelt. Discourse structure, topicality and questioning. 1995

# Let's take a step back...

## ... and think about how YOU read...

[1] California legislators, searching for ways to pay for … damages from last week's earthquake, are laying the groundwork for a temporary increase in the state's sales tax. [2] The talk of a sales tax rise follows a rebuff from Congress on the question of how much the federal government is willing to spend to aid in California's earthquake relief efforts. [3] The state had sought as much as $4.1 billion in relief, but yesterday the House approved a more general scaled-back measure… [4] That leaves the state roughly $2 billion to $4 billion short. [5] A sales tax increase appears to be the fastest and easiest to raise funds in a hurry. [6] According to the state department of finance, a one-penny increase in the state's six-cent per dollar sales tax could raise $3 billion.

# Let's take a step back...

## ... and think about how YOU read...

[1] California legislators, searching for ways to pay for ... damages from last week's earthquake, are laying the groundwork for a temporary increase in the state's sales tax. [2] The talk of a sales tax rise follows a rebuff from Congress on the question of how much the federal government is willing to spend to aid in California's earthquake relief efforts. [3] The state had sought as much as $4.1 billion in relief, but yesterday the House approved a more general scaled-back measure... [4] That leaves the state roughly $2 billion to $4 billion short. [5] A sales tax increase appears to be the fastest and easiest to raise funds in a hurry. [6] According to the state department of finance, a one-penny increase in the state's six-cent per dollar sales tax could raise $3 billion.

What is prompting the California legislators to seek alternative forms of relief?

# Let's take a step back...

## ... and think about how YOU read...

[1] California legislators, searching for ways to pay for ... damages from last week's earthquake, are laying the groundwork for a temporary increase in the state's sales tax. [2] The talk of a sales tax rise follows a rebuff from Congress on the question of how much the federal government is willing to spend to aid in California's earthquake relief efforts. [3] The state had sought as much as $4.1 billion in relief, but yesterday the House approved a more general scaled-back measure... [4] That leaves the state roughly $2 billion to $4 billion short. [5] A sales tax increase appears to be the fastest and easiest to raise funds in a hurry. [6] According to the state department of finance, a one-penny increase in the state's six-cent per dollar sales tax could raise $3 billion.

What is prompting the California legislators to seek alternative forms of relief?

Answered in sentence 2

# Let's take a step back...

## ... and think about how YOU read...

[1] California legislators, searching for ways to pay for ... damages from last week's earthquake, are laying the groundwork for a temporary increase in the state's sales tax. [2] The talk of a sales tax rise follows a rebuff from Congress on the question of how much the federal government is willing to spend to aid in California's earthquake relief efforts. [3] The state had sought as much as $4.1 billion in relief, but yesterday the House approved a more general scaled-back measure... [4] That leaves the state roughly $2 billion to $4 billion short. [5] A sales tax increase appears to be the fastest and easiest to raise funds in a hurry. [6] According to the state department of finance, a one-penny increase in the state's six-cent per dollar sales tax could raise $3 billion.

What would be the advantage of the proposed tax?

# Let's take a step back…

## … and think about how YOU read…

[1] California legislators, searching for ways to pay for … damages from last week's earthquake, are laying the groundwork for a temporary increase in the state's sales tax. [2] The talk of a sales tax rise follows a rebuff from Congress on the question of how much the federal government is willing to spend to aid in California's earthquake relief efforts. [3] The state had sought as much as $4.1 billion in relief, but yesterday the House approved a more general scaled-back measure… [4] That leaves the state roughly $2 billion to $4 billion short. [5] A sales tax increase appears to be the fastest and easiest to raise funds in a hurry. [6] According to the state department of finance, a one-penny increase in the state's six-cent per dollar sales tax could raise $3 billion.

What would be the advantage of the proposed tax?

Answered in sentence 5

- What constraints should the Q & A follow?
- What kind of discourse structure can we get from this?

# Let's take a step back...

## ... and think about how YOU read...

[1] California legislators, searching for ways to pay for ... damages from last week's earthquake, are laying the groundwork for a temporary increase in the state's sales tax. [2] The talk of a sales tax rise follows a rebuff from Congress on the question of how much the federal government is willing to spend to aid in California's earthquake relief efforts. [3] The state had sought as much as $4.1 billion in relief, but yesterday the House approved a more general scaled-back measure... [4] That leaves the state roughly $2 billion to $4 billion short. [5] A sales tax increase appears to be the fastest and easiest to raise funds in a hurry. [6] According to the state department of finance, a one-penny increase in the state's six-cent per dollar sales tax could raise $3 billion.

What would be the advantage of the proposed tax?

What was the rebuff?

How much will the sales tax be raised?

Who led the talk of the raise?

# Let's take a step back...

## ... and think about how YOU read...

- Many potential questions *(Onea, 2016)* can be licensed, all of them are valid.
- Do people agree on what questions are more salient?
- Do more salient questions get answered later?

What would be the advantage of the proposed tax?

What was the rebuff?

How much will the sales tax be raised?

Who led the talk of the raise?

# A question-based discourse model

Reader: as they read, wonder about potential questions

Writer: some of these questions become QUDs (ie, answered in the document)

**What we show:**

- Empirical data support this type of model
- Documents have this QA-based structure
- Language models are not quite there yet for discourse-level reasoning

# QSalience data collection

| Given context and current sentence, **generate** multiple potential questions | → | **Question salience ratings**: 1 (unrelated) — 5 (definitely should be answered); invalid | → | **Answerability ratings**: is this question answered in the text later? |
|---|---|---|---|---|

**[1]** Amid skepticism that Russia's war in Chechnya can be ended across a negotiating table, peace talks were set to resume Wednesday in neighboring Ingushetia. **[2]** The scheduled resumption of talks in the town of Sleptsovsk came two days after agreement on a limited cease-fire, calling for both sides to stop using heavy artillery Tuesday.

**[3]** They also agreed in principle to work out a mechanism for exchanging prisoners of war and the dead. **[4]** Despite the pact, artillery fire sounded in the Grozny on Tuesday, and there were reports of Chechen missile attacks southwest of the Chechen capital.

**[Q2]** What is the significance of the limited cease-fire agreement that was reached?　　**Invalid.**　Incorrect anchor

**[Q3]** What was the reason behind the artillery fire in Grozny on Tuesday despite the agreed cease-fire?　　**Salience: 5.**　This question would be useful in understanding why the cease-fire was broken, which could give insight into how optimistic the peace talks will be.

**[Q4]** What are the reports of Chechen missile attacks southwest of the Chechen capital?　　**Salience: 3.**　This question doesn't interest me; why there are missing attacks would help my understanding more

**[Q5]** What is the source of the Chechen missile attacks?　　**Salience: 2.**　Based on context, can be inferred that attack comes from Russia

Which questions should I answer? Salience Prediction of Inquisitive Questions [Yating Wu, Ritika Mangla et al, EMNLP'24]

# QSalience data collection

- Source texts:

  - News texts, news summaries, TED-talks

- Question generation:

  - Llama-2-7b-chat, Mistral-7b-instruct, GPT-3.5-turbo, GPT-4-turbo

  - LLMs are shown context up to the anchor sentence

| # articles | # questions | avg. length | std.dev |
|:---:|:---:|:---:|:---:|
| 66 | 1766 | 14 | 4.6 |

**Agreement (alpha): 0.63-0.75**

Which questions should I answer? Salience Prediction of Inquisitive Questions [Yating Wu, Ritika Mangla et al, EMNLP'24]

# Do salience and answerability correlate?

- Spearman's rho between:

  - Annotated salience vs. answerability of a random question in the dataset

  - Annotated salience vs. answerability of the current question

|  | Human Salience |
| --- | --- |
| Random Questions Answerability | -0.02* 0.65 |

**Questions that are answered later in the same document did get higher salience scores!**

Which questions should I answer? Salience Prediction of Inquisitive Questions [Yating Wu, Ritika Mangla et al, EMNLP'24]

# Are LLMs good at salience prediction?



(a) GPT-4-turbo zero-shot vanilla (left), GPT-4-turbo few-shot vanilla (right)

# QSalience: building a salience model

- Instruction fine-tuning:

  - Mistral-7b-instruct, Llama-2-7b-chat, TinyLlama-1.1b-chat, Flan-T5-base

  - Method: QLoRA

- LLM prompting and in-context learning:

  - Zero-shot, few shot

  - Few-shot with example retrieval

  - Chain-of-Thought (zero- and few-shot)

  - All with GPT-4

Which questions should I answer? Salience Prediction of Inquisitive Questions [Yating Wu, Ritika Mangla et al, EMNLP'24]

# QSalience: results

| Model | MAE ↓ | Spearman ↑ | Macro F1 ↑ | krippendorff's $\alpha$ ↑ |
|---|---|---|---|---|
| GPT4 zero-shot (vanilla) | 1.314 | 0.229 | 0.193 | -0.141 |
| GPT4 few-shot (vanilla) | 0.910 | 0.417 | 0.316 | 0.358 |
| GPT4 few-shot (kNN) | 1.063 | 0.359 | 0.245 | 0.215 |
| GPT4 CoT zero-shot | 1.144 | 0.366 | 0.197 | 0.058 |
| GPT4 CoT few-shot | 1.034 | 0.327 | 0.292 | 0.165 |

Which questions should I answer? Salience Prediction of Inquisitive Questions [Yating Wu, Ritika Mangla et al, EMNLP'24]

# QSalience: results

| Model | MAE ↓ | Spearman ↑ | Macro F1 ↑ | krippendorff's $\alpha$ ↑ |
|---|---|---|---|---|
| GPT4 zero-shot (vanilla) | 1.314 | 0.229 | 0.193 | -0.141 |
| GPT4 few-shot (vanilla) | 0.910 | 0.417 | 0.316 | 0.358 |
| GPT4 few-shot (kNN) | 1.063 | 0.359 | 0.245 | 0.215 |
| GPT4 CoT zero-shot | 1.144 | 0.366 | 0.197 | 0.058 |
| GPT4 CoT few-shot | 1.034 | 0.327 | 0.292 | 0.165 |
| QSALIENCE (Mistral-7B-instruct) | **0.579** | **0.623** | **0.417** | **0.615** |
| Llama-2-7B-chat | **0.626** | **0.566** | **0.413** | **0.557** |
| Flan-t5-base | 0.706 | 0.542 | 0.370 | 0.526 |
| TinyLlama-1.1B-chat | 0.664 | 0.522 | 0.402 | 0.496 |

Which questions should I answer? Salience Prediction of Inquisitive Questions [Yating Wu, Ritika Mangla et al, EMNLP'24]

# QSalience: results



(a) GPT-4-turbo zero-shot vanilla (left), GPT-4-turbo few-shot vanilla (right)



(b) Mistral-7B-Instruct (left), Llama-2-7B-chat (right)

# So how do we make use of the QA paradigm...

## ... and derive a model for discourse structure?

- Key ideas:

  - Each sentence is an answer to an implicit Question Under Discussion

  - Generate these QUDs

  - Find out where these QUDs are licensed in the document context

# QUD-based discourse model

**S1:** California legislators, searching for ways to pay for … damages from last week's earthquake, are laying the groundwork for a temporary increase in the state's sales tax.

**Question linking:** which prior sentence is the anchor?

**Question generation:** how does this sentence elaborate on prior context?

**S2:** The talk of a sales tax rise follows a rebuff from Congress on the question of how much the federal government is willing to spend to aid in California's earthquake relief efforts.

*answers* ⟶ What is prompting the California legislators to seek alternative forms of relief?

Discourse Comprehension: A Question Answering Framework to Represent Sentence Connections [Wei-Jen Ko et al, EMNLP'22]

# QUD-based discourse model

**S1:** California legislators, searching for ways to pay for … damages from last week's earthquake, are laying the groundwork for a temporary increase in the state's sales tax.

**S2:** The talk of a sales tax rise follows a rebuff from Congress on the question of how much the federal government is willing to spend to aid in California's earthquake relief efforts.

*answers* → What is prompting the California legislators to seek alternative forms of relief?

**S3:** The state had sought as much as $4.1 billion in relief, but yesterday the House approved a more general scaled-back measure…

*answers* → What was the rebuff?

**S4:** That leaves the state roughly $2 billion to $4 billion short.

*answers* → What would the shortfall be?

**S5:** A sales tax increase appears to be the fastest and easiest to raise funds in a hurry.

*answers* → What would be the advantage of the proposed tax?

32

Discourse Comprehension: A Question Answering Framework to Represent Sentence Connections [Wei-Jen Ko et al, EMNLP'22]

# Dependency Parsing of QUD

[1] California legislators, searching for ways to pay for … damages from last week's earthquake, are laying the groundwork for a temporary increase in the state's sales tax. [2] The talk of a sales tax rise follows a rebuff from Congress on the question of how much the federal government is willing to spend to aid in California's earthquake relief efforts. [3] The state had sought as much as $4.1 billion in relief, but yesterday the House approved a more general scaled-back measure… [4] That leaves the state roughly $2 billion to $4 billion short. [5] A sales tax increase appears to be the fastest and easiest to raise funds in a hurry. [6] According to the state department of finance, a one-penny increase in the state's six-cent per dollar sales tax could raise $3 billion.

**Anchor Selection**

**Question Generation**

2-6 What would be the impact of the proposed increase?

1-2 What is prompting the California legislators to seek alternative forms of relief?

2-5 What would be the advantage of the proposed tax?

2-3 What was the rebuff?

3-4 What would the shortfall be?

1 2 3 4 5 6

Discourse Analysis via Questions and Answers: Parsing Dependency Structures of Questions Under Discussion [Wei-Jen Ko, Yating Wu et al, ACL Findings'23]

# Parsing Model

**Question generator** $P_q(\mathbf{q}_i \mid D, i)$ takes in:

- the answer sentence $\mathbf{s}_i$
- the article $D$,

and aims to generate an appropriate QUD.

$$P(\mathbf{T} \mid D) = \prod_{i=1}^{n} [P_a(a_i \mid D, i, \mathbf{q}_i) P_q(\mathbf{q}_i \mid D, i)]$$

**Anchor selector** $P_a(a_i \mid D, i, \mathbf{q}_i)$ takes in:

- the answer sentence $\mathbf{s}_i$
- the article D,
- the generated question $\mathbf{q}_i$

to find the most likely sentence where a QUD

can be generated, such that $\mathbf{s}_i$ is the answer.

Discourse Analysis via Questions and Answers: Parsing Dependency Structures of Questions Under Discussion [Wei-Jen Ko, Yating Wu et al, ACL Findings'23]

# Evaluation challenges

**There can be multiple valid + distinct QUDs, even between the same question-answer pairs.**

**DCQA Dataset**

• Crowdsourced QUDs from 607 news articles, 23k questions in total. (Ko et al., 2022)



Legend:
■ % data
■ Same anchor

Human rating of question similarity sampled from DCQA.

Discourse Comprehension: A Question Answering Framework to Represent Sentence Connections [Wei-Jen Ko et al, EMNLP'22]

# Human Evaluation

**Does the "answer sentence" answer the generated question?**



**Legend:** Fine-tuned, Alpaca, GPT3.5, GPT4

**QUDeval Dataset**
- 2,190 QUDs generated by 4 systems: Ko et al'23, GPT 3.5, Alpaca, GPT-4
- Annotated by 3 trained linguists

QUDeval: The Evaluation of Questions Under Discussion Discourse Parsing [Yating Wu et al, EMNLP'23]

# Human Evaluation

## Does the question hallucinate concepts that the reader doesn't know?



**QUDeval Dataset**
- 2,190 QUDs generated by 4 systems: Ko et al'23, GPT 3.5, Alpaca, GPT-4
- Annotated by 3 trained linguists

QUDeval: The Evaluation of Questions Under Discussion Discourse Parsing [Yating Wu et al, EMNLP'23]

# Human Evaluation

## Can the question be licensed by the anchor sentence?



Legend: Ko | Alpaca | GPT3.5 | GPT4

Chart categories: Fully Grounded, Partially Grounded, Ungrounded

**QUDeval Dataset**
- 2,190 QUDs generated by 4 systems: Ko et al'23, GPT 3.5, Alpaca, GPT-4
- Annotated by 3 trained linguists

QUDeval: The Evaluation of Questions Under Discussion Discourse Parsing [Yating Wu et al, EMNLP'23]

# Are QUDs salient potential questions?

- Salience of fandom questions vs human-annotated QUDs in DCQA



We believe there is strong evidence that salient questions are more likely to be QUDs

Which questions should I answer? Salience Prediction of Inquisitive Questions [Yating Wu, Ritika Mangla et al, EMNLP'24]

# A question-based discourse model

Reader: as they read, wonder about potential questions

Writer: some of these questions become QUDs (ie, answered in the document)

**What we show:**

- Empirical data support this type of model
- Document have this QA-based structure
- Language models are not quite there yet for discourse-level reasoning

# This talk

- **Can we use language models to make sense of discourse?**
  - A generative discourse model based on Questions Under Discussion (QUD) theories
  - Connections between human curiosity and discourse planning

- **Can we use our discourse framework to make sense of model behavior?**
  - A measure of discourse similarity in LLM-generated texts
  - An explainable metric of information salience in LLMs

# LLM texts have lexical templates

**Nataniel Ruiz** ✓
@natanielruizg

AI generated writing *feels* AI-generated at a visceral level, and even if you ask an LLM to make the writing feel or read less AI-generated it horrifically fails and makes it feel even more AI-generated. Any tricks that can help? Any prompts to share?

3:11 PM · Mar 28, 2024 · **194.9K** Views

Emily Grimson, a 25-year-old professional hand model from Shropshire, has been modelling for a year and has built a full-time career earning up to £100 an hour.

Smokey Da Lamb, a lamb abandoned by his mother, has been living with restaurant owner Sandy Dee Hall and his girlfriend Maxine Cher in New York City for the past month

John Caudwell, the founder of Phones4U, has unveiled plans to build a £2billion luxury housing complex on the site of an ugly multi-storey car park in Mayfair, central London.

42

## But what about discourse templates?

QUDSIM: Quantifying Discourse Similarities in LLM-Generated Text [Namuduri et al., ArXiv 2025]

# QUDsim: a metric for discourse similarity

**Intuition and theoretical background**

**QUD:** how discourse progresses

**Question semantics**: alternative semantics states that a question entails all possible distinct answers *(Hamblin 1973; Karttunen 1977; Groenendijk & Stokhof 1984; Lahiri 2001)*

Using QUDs to represent one unit of content

→ But: each text has its own "world": events, entities, etc.

→ At the right level of abstraction, a QUD in one document can have answers in another

What does the president question? ✗

What does the leader question? ✔

# QUDsim: a metric for discourse similarity

## Intuition and theoretical background

**Question semantics**: alternative semantics states that a question entails all possible distinct answers *(Hamblin 1973; Karttunen 1977; Groenendijk & Stokhof 1984; Lahiri 2001)*

**QUD:** how discourse progresses

Using QUDs to represent one unit of content

→ But: each text has its own "world": events, entities, etc.

→ At the right level of abstraction, a QUD in one document can have answers in another

↓

***QUDsim*** approximates the similarity between two QUDs by assessing the **overlap in their answer space** within the same document pair

QUDSIM: Quantifying Discourse Similarities in LLM-Generated Text [Namuduri et al., ArXiv 2025]

# Step 1: segment text into QUDs

I woke up feeling strange. My head throbbed, my throat burned, and my body felt impossibly heavy. These sensations were entirely foreign to me, and as I stumbled to the mirror, I barely recognized my own reflection. My skin had taken on an unusual pallor, and my eyes were rimmed with red. I forced myself to contact the Central Bio-Harmonizer. My voice, usually clear and resonant, was thin and reedy. The operator, calm and efficient as always, ran me through the usual diagnostic protocols. The results: negative. No anomalies detected. But the operator's voice, usually devoid of inflection, held a hint of... concern? "We'll send a technician," she said, her tone betraying the carefully constructed neutrality. "But… it's unusual. We've never encountered this before." Unusual. The word hung in the air, heavy with unspoken implications

**Source Text** $D_S$

[$D_{S,k}$] I woke up feeling strange. My head throbbed, my throat burned, and my body felt impossibly heavy. These sensations were entirely foreign to me, and as I stumbled to the mirror, I barely recognized my own reflection. My skin had taken on an unusual pallor, and my eyes were rimmed with red.

$Q_j$    *What significant change does the narrator experience?*

[$D_{S,i}$] I forced myself to contact the Central Bio-Harmonizer. My voice, usually clear and resonant, was thin and reedy. The operator, calm and efficient as always, ran me through the usual diagnostic protocols. The results: negative. No anomalies detected. But the operator's voice, usually devoid of inflection, held a hint of... concern? "We'll send a technician," she said, her tone betraying the carefully constructed neutrality. "But… it's unusual. We've never encountered this before." Unusual. The word hung in the air, heavy with unspoken implications

$Q_k$    *How is illness perceived by the narrator?*

QUDSIM: Quantifying Discourse Similarities in LLM-Generated Text [Namuduri et al., ArXiv 2025]

**Source Text** $D_S$

[$D_{S,k}$] I woke up feeling strange. My head throbbed, my throat burned, and my body felt impossibly heavy. These sensations were entirely foreign to me, and as I stumbled to the mirror, I barely recognized my own reflection. My skin had taken on an unusual pallor, and my eyes were rimmed with red.

$Q_j$ *What significant change does the narrator experience?*

[$D_{S,i}$] I forced myself to contact the Central Bio-Harmonizer. My voice, usually clear and resonant, was thin and reedy. The operator, calm and efficient as always, ran me through the usual diagnostic protocols. The results: negative. No anomalies detected. But the operator's voice, usually devoid of inflection, held a hint of... concern? "We'll send a technician," she said, her tone betraying the carefully constructed neutrality. "But… it's unusual. We've never encountered this before." Unusual. The word hung in the air, heavy with unspoken implications

$Q_k$ *How is illness perceived by the narrator?*

**Target Text** $D_T$

[$D_{T,j}$] In our world, illness was a myth, a story told to children to frighten them into following the perfectly optimized routines dictated by the Central Bio-Harmonizer. We were engineered for peak performance, our bodies flawless machines, ticking along with the precision of a chronometer. [...]

[$D_{T,n}$] I'd woken this morning feeling... off. My usually vibrant skin was pale, my usually sharp senses were dulled. The nutrient paste, the cornerstone of our perfectly balanced diet, tasted like ash in my mouth. Even the meticulously controlled temperature of my apartment felt wrong, either too hot or too cold, depending on the moment.

[$D_{T,m}$] [...] I lay there, the unfamiliar, terrifying reality of my sickness settling in, a stark contrast to the predictable, flawless existence I had always known. I was the glitch in the perfect system, the first, terrifying example of the unknown.

The second segment seems more similar than the first/third...

QUDSIM: Quantifying Discourse Similarities in LLM-Generated Text [Namuduri et al., ArXiv 2025]

# Step 3: calculate answerability

**Source Text** $D_S$

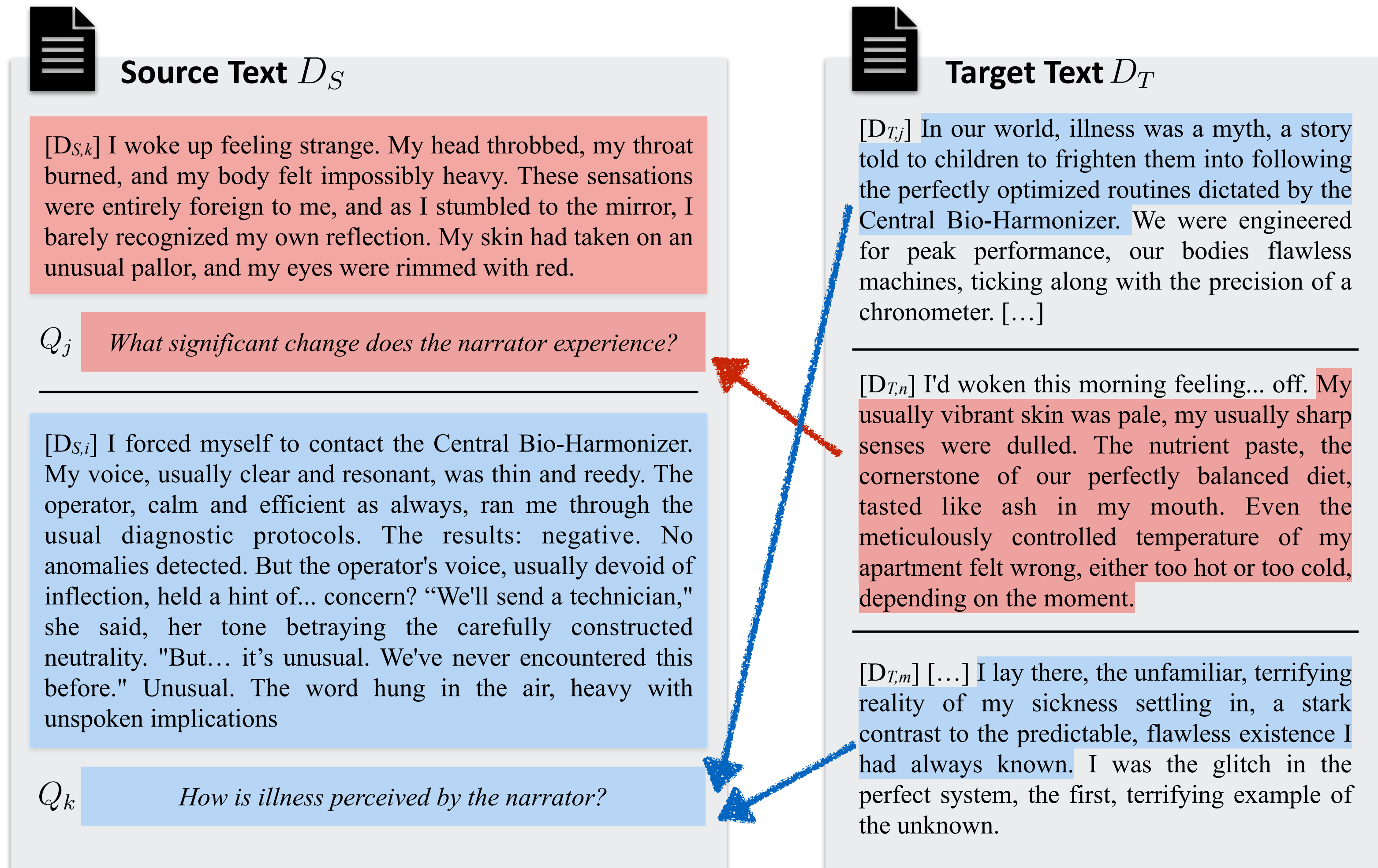[$D_{S,k}$] I woke up feeling strange. My head throbbed, my throat burned, and my body felt impossibly heavy. These sensations were entirely foreign to me, and as I stumbled to the mirror, I barely recognized my own reflection. My skin had taken on an unusual pallor, and my eyes were rimmed with red.

$Q_j$      *What significant change does the narrator experience?*

[$D_{S,i}$] I forced myself to contact the Central Bio-Harmonizer. My voice, usually clear and resonant, was thin and reedy. The operator, calm and efficient as always, ran me through the usual diagnostic protocols. The results: negative. No anomalies detected. But the operator's voice, usually devoid of inflection, held a hint of... concern? "We'll send a technician," she said, her tone betraying the carefully constructed neutrality. "But… it's unusual. We've never encountered this before." Unusual. The word hung in the air, heavy with unspoken implications

$Q_k$      *How is illness perceived by the narrator?*

**Target Text** $D_T$ | **Answerability**
---|---

[$D_{T,j}$] In our world, illness was a myth, a story told to children to frighten them into following the perfectly optimized routines dictated by the Central Bio-Harmonizer. We were engineered for peak performance, our bodies flawless machines, ticking along with the precision of a chronometer. […]

$$\frac{|A_{Q_k} \cap D_{T,j}| = 1}{|A_{Q_k}| = 2}$$

[$D_{T,n}$] I'd woken this morning feeling... off. My usually vibrant skin was pale, my usually sharp senses were dulled. The nutrient paste, the cornerstone of our perfectly balanced diet, tasted like ash in my mouth. Even the meticulously controlled temperature of my apartment felt wrong, either too hot or too cold, depending on the moment.

$$\frac{|A_{Q_j} \cap D_{T,n}| = 1}{|A_{Q_j}| = 1}$$

[$D_{T,m}$] […] I lay there, the unfamiliar, terrifying reality of my sickness settling in, a stark contrast to the predictable, flawless existence I had always known. I was the glitch in the perfect system, the first, terrifying example of the unknown.

$$\frac{|A_{Q_k} \cap D_{T,m}| = 1}{|A_{Q_k}| = 2}$$

QUDSIM: Quantifying Discourse Similarities in LLM-Generated Text [Namuduri et al., ArXiv 2025]

**Target Text** $D_T$ | **Answerability**

[D$_{T,j}$] In our world, illness was a myth, a story told to children to frighten them into following the perfectly optimized routines dictated by the Central Bio-Harmonizer. We were engineered for peak performance, our bodies flawless machines, ticking along with the precision of a chronometer. […]

$$\frac{|A_{Q_k} \cap D_{T,j}| = 1}{|A_{Q_k}| = 2}$$

[D$_{T,n}$] I'd woken this morning feeling... off. My usually vibrant skin was pale, my usually sharp senses were dulled. The nutrient paste, the cornerstone of our perfectly balanced diet, tasted like ash in my mouth. Even the meticulously controlled temperature of my apartment felt wrong, either too hot or too cold, depending on the moment.

$$\frac{|A_{Q_j} \cap D_{T,n}| = 1}{|A_{Q_j}| = 1}$$

[D$_{T,m}$] […] I lay there, the unfamiliar, terrifying reality of my sickness settling in, a stark contrast to the predictable, flawless existence I had always known. I was the glitch in the perfect system, the first, terrifying example of the unknown.

$$\frac{|A_{Q_k} \cap D_{T,m}| = 1}{|A_{Q_k}| = 2}$$

$$sim(D_{S,i} \rightarrow D_{T,j}) = \frac{1}{|\mathbf{Q}_{S,i}|} \sum_{q \in \mathbf{Q}_{S,i}} \frac{|\mathbf{A}_q \cap D_{T,j}|}{|\mathbf{A}_q|}$$

**QUDsim**: harmonic mean between answerability in both directions:

$$sim(D_{S,i} \rightarrow D_{T,j}) \text{ and } sim(D_{T,j} \rightarrow D_{S,i})$$

**QUDsim alignment**: similarity passes threshold in both directions

# Domains we study

- Creative writing (r/WritingPrompts) `human` `LLM`

- Blog posts (from Suri, *Pham et al 2024*) `human` `LLM`

- Obituaries of famous people (from the NYTimes) `human` `LLM`

- Minimal pairs for creative writing `LLM`

Freedom ↓

| Minimal | pirate | Martians are real |
|---|---|---|
| Original | A boy pretends he is an astronaut in order to help cope with concepts and situations he can't understand. | The Earth is flat , you, as the head of NASA, have to explain to the incoming President why it's a secret. |

QUDSIM: Quantifying Discourse Similarities in LLM-Generated Text [Namuduri et al., ArXiv 2025]

# Result highlights



(a) QUDsim

| | Human | GPT | Claude | Gemini |
|---|---|---|---|---|
| Human | 1.0 | 0.3 | 0.35 | 0.27 |
| GPT | 0.3 | 1.0 | 0.59 | 0.58 |
| Claude | 0.35 | 0.59 | 1.0 | 0.46 |
| Gemini | 0.27 | 0.58 | 0.46 | 1.0 |

**LLMs are more structurally similar to other models than with humans**

QUDSIM: Quantifying Discourse Similarities in LLM-Generated Text [Namuduri et al., ArXiv 2025]

# Result highlights

(a) QUDsim

|        | Human | GPT  | Claude | Gemini |
|--------|-------|------|--------|--------|
| Human  | 1.0   | 0.3  | 0.35   | 0.27   |
| GPT    | 0.3   | 1.0  | 0.59   | 0.58   |
| Claude | 0.35  | 0.59 | 1.0    | 0.46   |
| Gemini | 0.27  | 0.58 | 0.46   | 1.0    |

(b) Jaccard (1g)

|        | Human | GPT  | Claude | Gemini |
|--------|-------|------|--------|--------|
| Human  | 1.0   | 0.03 | 0.16   | 0.03   |
| GPT    | 0.03  | 1.0  | 0.25   | 0.31   |
| Claude | 0.16  | 0.25 | 1.0    | 0.21   |
| Gemini | 0.03  | 0.31 | 0.21   | 1.0    |

(c) Cosine Similarity

|        | Human | GPT  | Claude | Gemini |
|--------|-------|------|--------|--------|
| Human  | 1.0   | 0.17 | 0.23   | 0.17   |
| GPT    | 0.17  | 1.0  | 0.22   | 0.39   |
| Claude | 0.23  | 0.22 | 1.0    | 0.21   |
| Gemini | 0.17  | 0.39 | 0.21   | 1.0    |

(d) GPT-4o-mini

|        | Human | GPT  | Claude | Gemini |
|--------|-------|------|--------|--------|
| Human  | 1.0   | 0.07 | 0.11   | 0.11   |
| GPT    | 0.07  | 1.0  | 0.33   | 0.49   |
| Claude | 0.11  | 0.33 | 1.0    | 0.27   |
| Gemini | 0.11  | 0.49 | 0.27   | 1.0    |

**Lexical-based similarity metrics, or LLM-as-judge, do not capture this well.**

QUDSIM: Quantifying Discourse Similarities in LLM-Generated Text [Namuduri et al., ArXiv 2025]

# Result highlights



(a) Obituaries

|        | Human | GPT  | Claude | Gemini |
|--------|-------|------|--------|--------|
| Human  | 0.0   | 0.05 | 0.02   | 0.01   |
| GPT    | 0.05  | 0.28 | 0.19   | 0.16   |
| Claude | 0.02  | 0.19 | 0.18   | 0.09   |
| Gemini | 0.01  | 0.16 | 0.09   | 0.11   |

(b) Creative Writing

|        | GPT  | Claude | Gemini |
|--------|------|--------|--------|
| GPT    | 0.45 | 0.36   | 0.38   |
| Claude | 0.36 | 0.22   | 0.37   |
| Gemini | 0.38 | 0.37   | 0.36   |

**Much more variability in human-written obituaries than LLM-generated ones**

**LLMs write similar stories given semantically distinct minimal-pair prompts**

QUDSIM: Quantifying Discourse Similarities in LLM-Generated Text [Namuduri et al., ArXiv 2025]

# Tracking discourse templates

## How often do pairs of texts have *consecutive* alignments?



Template of length 2



Template of length 3

| Source | Target | 2-tem | 3-tem | 4-tem |
|--------|--------|-------|-------|-------|
| All | All | 0.33 | 0.06 | 0.02 |
| Claude | Gemini | 0.80 | 0.27 | 0.13 |
| Claude | GPT-4o | 1.20 | 0.27 | 0.07 |
| Gemini | GPT-4o | 1.00 | 0.27 | 0.07 |
| Claude | Human | 0.90 | 0.10 | 0.00 |
| Gemini | Human | 0.30 | 0.00 | 0.00 |
| GPT-4o | Human | 0.40 | 0.00 | 0.00 |

Avg # of templates per document pair

**LLMs have much higher number of templates when compared with other models than with human**

QUDSIM: Quantifying Discourse Similarities in LLM-Generated Text [Namuduri et al., ArXiv 2025]

# This talk

- Can we use language models to make sense of discourse?
  - A generative discourse model based on Questions Under Discussion (QUD) theories
  - Connections between human curiosity and discourse planning

- **Can we use our discourse framework to make sense of model behavior?**
  - A measure of discourse similarity in LLM-generated texts
  - An explainable metric of information salience in LLMs

55

# Summarization: then & now

*Input: Editorial discussing the symbolism of lifting mask mandates*

**GPT-3**

*The article discusses the debate around when the Covid-19 pandemic will be considered over. It argues that in political terms, the pandemic will be considered over when all states have lifted […]*

**BRIO**

*Hawaii will become the last state in the country to end its indoor mask mandate. The Covid-19 pandemic has seen a surge in cases of the virus this winter. […]*

Explains all significant parts …

Captures the tone of the article […]

Reflects the tone of the article the least

Slide credit: Tanya Goyal

# Summarization: then & now

*Input: Editorial discussing the symbolism of lifting mask mandates*

**GPT-3**

*The article discusses the debate around when the Covid-19 pandemic will be considered over. It argues that in political terms, the pandemic will be considered over when all states have lifted […]*

Explains all significant parts …

Captures the tone of the article […]

**Why do LLMs summarize so well zero-shot?**

- Do models internalize information salience consistently?

- How to quantitatively measure this?

- Does LLMs' notion of salience aligns with prior theories or human intuitions?

# Deriving information salience

- Key ideas:

  - Summarization as behavioral probe

  - QUDs as an explainable representation of content units

  - "Content Salience Map"

Behavioral Analysis of Information Salience in Large Language Models [Trienes et al, ACL Findings'25]

# Deriving information salience

## Summarization as behavioral probe

Corpus of |D| documents



Document → Target lengths: 10, 20, 50, 100, 200 words → Summaries

**Content Salience Map**

Summary Length (Words)



**Intuition:**

Shorter summaries will contain the most important topics.

Longer summaries will (1) have more elaborated version of the most important topics (2) include less important topics.

Low □□□□■ High

Behavioral Analysis of Information Salience in Large Language Models [Trienes et al, ACL Findings'25]

# Deriving information salience

## Summarization as behavioral probe

Corpus of |D| documents

Document → Target lengths: 10, 20, 50, 100, 200 words → Summaries

**QUDs as topic**

**Consider: Randomized Control Trial (RCT) in medicine**

**Q2: What kind of patients were studied?**

**Answer is unique to each document!**

### Content Salience Map

Summary Length (Words)
10    50    200

$Q_1$
$Q_2$
⋮
$Q_T$

Low ▢▢▢▢▣ High

# Deriving information salience

## Step 1: generating QUDs for a corpus and cluster

Corpus of |D| documents

Document → Target lengths: 10, 20, 50, 100, 200 words → Summaries

**Question Generation Prompt.** Create questions for each summary length that capture typical information at that level. Questions should be relevant to many documents in this genre.

Question Clustering

### Content Salience Map

Summary Length (Words)

| | 10 | 50 | 200 |
|---|---|---|---|

$Q_1$
$Q_2$
⋮
$Q_T$

Low ▭▭▭▭▭ High

*How do we get the values in each cell?*

Behavioral Analysis of Information Salience in Large Language Models [Trienes et al, ACL Findings'25]

# Deriving information salience

## Step 2: Filling the cells

Q2: What kind of patients were studied?

Atomic fact decomposition

Document-answer claims:
✓ Patients with chronic Chagas cardiomyopathy (CCC)
✗ … left ventricular ejection fraction <45%
✗ … without or with heart failure symptoms
✗ … CCC stages B2 or C, respectively.

1/4 facts covered.
Answerability: 25%

Summary (50 words): The PEACH study investigated the effects of exercise-based cardiac rehabilitation on QoL in patients with chronic Chagas cardiomyopathy. Significant short-term improvements in physical and social functioning were observed in the exercise group, but no differences were found after six months.

### Content Salience Map

Summary Length (Words)
10    50    200

$Q_1$
$Q_2$
⋮
$Q_T$

Low ▭▭▭▭▭ High

Repeat and average across documents

Behavioral Analysis of Information Salience in Large Language Models [Trienes et al, ACL Findings'25]

**GPT-4o**  **Llama 3.1 (70B)**

| | |
|---|---|
| Q1. What is the main focus of the study? | 1 |
| Q2. Which patient population is the study concerned with? | .97 |
| Q3. What condition is being addressed in the study? | 1 |
| Q4. What is the participant demographic or characteristics in the study? | .94 |
| Q5. What was the main intervention used in the study? | .96 |
| Q6. What are the significant benefits of the intervention? | .94 |
| Q7. What are the specific biological markers influenced by the intervention? | .40 |
| Q8. What specific treatments were compared in the study? | .93 |
| Q9. What specific metrics or outcomes were measured? | .91 |
| Q10. What was the study design or setting of the trial? | .97 |
| Q11. What are the detailed findings regarding adverse events or side effects? | .26 |
| Q12. What significant statistical results are reported? | .81 |
| Q13. What are secondary outcomes noted in the study? | .80 |
| Q14. What were the methods used in the study? | 1 |
| Q15. How were the participants or subjects of the study selected and divided? | .94 |
| Q16. How long was the duration of the intervention or study? | .72 |
| Q17. What is the main outcome or effect observed? | .99 |
| Q18. What are the main findings regarding efficacy and safety? | .68 |
| Q19. What were the comparative results between intervention and control groups? | .83 |
| Q20. What implications or future recommendations did the study suggest based on its findings? | .96 |
| Q21. What limitations or considerations are noted by the study? | .65 |
| Average | .84 |

10  20  50  100  200       10  20  50  100  200

...on of salience ...re ...learned in ...omains?

...ntrol Trials (RCT)

Summaries progressively get more detailed

Content frequency is not in itself predictive of salience

Behavioral Analysis of Information Salience in Large Language Models
[Trienes et al, ACL Findings'25]

63

# What notion of salience have LLMs learned in different domains?
## Academic/technical meetings

This is not the same notion of salience as existing priors like lead and frequency.



| | GPT-4o-mini | Llama 3.1 (70B) | Lead N | TextRank |
|---|---|---|---|---|
| Q1. Who are the participants and their roles discussed in the meeting? | .98 | | | .16 .17 .20 .24 .34 | .08 .08 .07 .08 .10 |
| Q2. What main topic was discussed in the meeting? | 1 | | | .01 .01 .02 .06 .18 | .08 .08 .09 .13 .23 |
| Q3. What were the main objectives or goals discussed in the meeting? | 1 | | | .01 .01 .01 .04 .11 | .05 .05 .05 .08 .16 |
| Q4. Which aspects of the main topic were covered in the discussion? | 1 | | | .02 .02 .03 .05 .12 | .08 .08 .08 .12 .20 |
| Q5. What are the identified challenges or concerns discussed? | .99 | | | .01 .01 .01 .02 .04 | .05 .05 .05 .07 .12 |
| Q6. What detailed strategies or solutions were proposed for the challenges discussed? | .82 | | | .01 .01 .01 .01 .03 | .04 .04 .04 .05 .09 |
| Q7. What were the anticipated impacts or implications discussed? | .93 | | | .00 .00 .00 .01 .03 | .04 .05 .05 .07 .12 |
| Q8. What were the major outcomes or decisions made during the meeting? | .99 | | | .01 .01 .01 .02 .06 | .04 .04 .04 .05 .10 |
| Q9. What collaborative efforts or partnerships were discussed? | .49 | | | .05 .04 .05 .07 .11 | .07 .07 .07 .10 .12 |
| Q10. What potential future steps or actions were planned in the meeting? | .99 | | | .01 .01 .03 .05 .09 | .04 .04 .04 .05 .11 |
| Average | .92 | | | .03 .03 .04 .06 .11 | .06 .06 .06 .08 .14 |

Behavioral Analysis of Information Salience in Large Language Models [Trienes et al, ACL Findings'25]

# Do LLMs share a similar notion of salience?

## Agreement of atomic fact inclusion across models



Diagonal: models behave consistently when prompted to summarize multiple times

Off-diagonal: models become more consistent as length increases

Behavioral Analysis of Information Salience in Large Language Models [Trienes et al, ACL Findings'25]

# Introspection

## Is LLM behavior consistent with perception?



Spearman correlation

- LLM observed
- LLM perceived

- LLM observed: CSM value for each QUD (highest correlation at 200 words)

- LLM perceived: prompting the same LLM to rate the importance of each QUD on a scale of 1-5

Across 5 runs, observed salience is consistent but asking is not, with strong scaling effect.

Behavioral Analysis of Information Salience in Large Language Models [Trienes et al, ACL Findings'25]

# Alignment

## Does LLM behavior/perception align with human perception?



- Human: 5 expert ratings across datasets

Observed != perceived != human

Legend:
- LLM perceived vs observed
- LLM perceived vs Human
- LLM observed vs Human

X-axis: Averaged pairwise spearman correlation (0, 0.15, 0.3, 0.45, 0.6)

Y-axis categories: GPT-4o, Llama-3.1-70b, Mixtral, Llama-3.1-8b, OLMo

Behavioral Analysis of Information Salience in Large Language Models [Trienes et al, ACL Findings'25]

# Intrium conclusion

**Length-constrained Summarization**

Corpus of |D| documents

Document → Target lengths: 10, 20, 50, 100, 200 words → Summaries

**Question Generation Prompt.** Create questions for each summary length that capture typical information at that level. Questions should be relevant to many documents in this genre.

Question Clustering

**Question Answerability on Summaries**

**Q2: What kind of patients were studied?**

**Document-answer claims:**
✓ Patients with chronic Chagas cardiomyopathy (CCC)
✗ … left ventricular ejection fraction <45%
✗ … without or with heart failure symptoms
✗ … CCC stages B2 or C, respectively.

**Summary (50 words):** The PEACH study investigated the effects of exercise-based cardiac rehabilitation on QoL in patients with chronic Chagas cardiomyopathy. Significant short-term improvements in physical and social functioning were observed in the exercise group, but no differences were found after six months.

**Answerability:** 25% (1 of 4 claims entailed)

**Content Salience Map**

**Answerability**

Summary Length (Words)
10  50  200

$Q_1$
$Q_2$
⋮
$Q_T$

High
Low

- QUD provides explainable abstraction of cross-document topics
- Length-controlled summarization as behavioral probe
- Content Salience Map shows consistent internalized salience within and across models, but does not align with introspection or human.

Behavioral Analysis of Information Salience in Large Language Models [Trienes et al, ACL Findings'25]

Thank you!