

# Quantifying the Impact of Disfluency on Spoken Content Summarization

**Maria Teleki, Xiangjue Dong, James Caverlee**

Texas A&M University

College Station, Texas, USA

{mariateleki, xj.dong, caverlee}@tamu.edu

*In LREC-COLING 2024*



# 1 Minute Summary

## Original

Hello and welcome to our podcast! Let's get right to it. Today we're going to be interviewing a very special guest, someone I know you guys have been excited about having on the show.

## Repeats with N=3

Hello and welcome to our podcast! Let's get **get get get** right to it. Today we're going to be interviewing a **a a a** very special guest, someone I know you guys have been excited about having on the show.

## Interjections with N=3

Hello and welcome to our podcast! Let's get right **uh okay okay** to it. Today we're going to be interviewing a very special **um so I mean** guest, someone I know you guys have been excited about having on the show.

## False Starts with N=3

Hello and welcome to our podcast! Let's get right to it. Today we're **today we're today we're today we're** going to be interviewing a very special guest, someone I know you guys have been excited about having on the show.

- **Disfluencies** are a key characteristic of **spoken content**.
  - We study 3 types of disfluencies -- *repeats*, *interjections*, and *false starts* -- in terms of the **Shriberg disfluency definition**.<sup>1</sup>
- **Summarization** quality decreases with increased disfluency.
- We use a **parsing-based SOTA disfluency annotator**<sup>2</sup> to repair the disfluencies via removal and tagging.
- We find that training on the repaired transcripts ( $\text{train}_R$ ) and testing on the original transcripts (test) yields the best results.

<sup>1</sup>Elizabeth Ellen Shriberg. 1994. Preliminaries to a theory of speech disfluencies. Ph.D. thesis.

<sup>2</sup>Paria Jamshid Lou and Mark Johnson. 2020. Improving disfluency detection by self-training a self-attentive model. In Association for Computational Linguistics, pages 3754–3763.

# Quantifying the Impact of Disfluency on Spoken Content Summarization

**Maria Teleki, Xiangjue Dong, James Caverlee**

Texas A&M University

College Station, Texas, USA

{mariateleki, xj.dong, caverlee}@tamu.edu

*In LREC-COLING 2024*



# What is a disfluency?

- Disfluencies are a key characteristic of **spoken content**.
- We study 3 types of disfluencies -- *repeats*, *interjections*, and *false starts* -- in terms of the **Shriberg disfluency definition**.<sup>1</sup>

## Original

Hello and welcome to our podcast! Let's get right to it. Today we're going to be interviewing a very special guest, someone I know you guys have been excited about having on the show.

## Repeats with N=3

Hello and welcome to our podcast! Let's get **get get get** right to it. Today we're going to be interviewing a **a a a** very special guest, someone I know you guys have been excited about having on the show.

## Interjections with N=3

Hello and welcome to our podcast! Let's get right **uh okay okay** to it. Today we're going to be interviewing a very special **um so I mean** guest, someone I know you guys have been excited about having on the show.

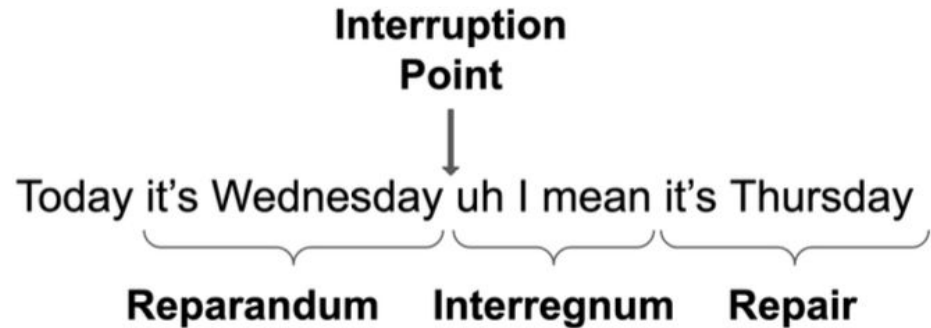
## False Starts with N=3

Hello and welcome to our podcast! Let's get right to it. Today we're **today we're today we're today we're** going to be interviewing a very special guest, someone I know you guys have been excited about having on the show.

<sup>1</sup>Elizabeth Ellen Shriberg. 1994. Preliminaries to a theory of speech disfluencies. Ph.D. thesis.

# What is a disfluency? The Shriberg disfluency definition.<sup>1</sup>

- The reparandum and interregnum are removed to form a fluent sentence.
- *Repeats* and *false starts* occur within the **reparandum**.
- *Interjections* occur within the **interregnum**.



<sup>1</sup>Elizabeth Ellen Shriberg. 1994. Preliminaries to a theory of speech disfluencies. Ph.D. thesis.

**Many important NLP tasks like summarization are often designed for written content rather than the looser, noiser, and more disfluent style of spoken content.<sup>1,2,3,4</sup>**

<sup>1</sup>Yang Liu and Mirella Lapata. 2019. Text summarization with pretrained encoders. In Conference on Empirical Methods in Natural Language Processing and the International Joint Conference on Natural Language Processing, pages 3730–3740.

<sup>2</sup>Mike Lewis, Yinhan Liu, et al. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In Association for Computational Linguistics, pages 7871–7880.

<sup>3</sup>Ani Nenkova and Kathleen McKeown. 2012. A survey of text summarization techniques. Mining text data, pages 43–76.

<sup>4</sup>Ramesh Nallapati, Bowen Zhou, et al. 2016. Abstractive text summarization using sequence-to-sequence RNNs and beyond. In Conference on Computational Natural Language Learning, pages 280–290.

# Research Questions

## ***RQ1: How Do Disfluencies Impact Summarization Quality?***

We synthetically inject disfluency events (repeats, interjections, false starts, and their combinations) at a range of severity levels and measure their impact on summarization quality.

## ***RQ2: Can Summarization Quality be Improved By Directly Modeling Disfluency?***

We explore the use of a state-of-the-art disfluency detection model to improve the summarization quality by either (1) removing the disfluencies, or (2) tagging the disfluencies.

# ***RQ1: How Do Disfluencies Impact Summarization Quality?***

## ***The Spotify Podcasts Dataset<sup>1</sup>***

- This dataset was originally used for the summarization task from the TREC 2020 Podcasts Track.<sup>2</sup>
- We use the test set for the summarization task, which consists of 1,027 podcasts. For each, we have:
  - The podcast transcript
  - The Show ID
  - The Episode ID
  - The creator-provided show description
  - The creator-provided episode description<sup>2</sup>
- We keep podcasts which have text occurring in their transcript in the first 60 seconds, which leaves us with 1,020 podcasts.

<sup>1</sup>Clifton, Ann and Reddy, Sravana et al. 2020. 100,000 podcasts: A spoken English document corpus.

<sup>2</sup>Rosie Jones, Ben Carterette, Ann Clifton, et al. 2020. TREC 2020 Podcasts Track Overview. In Text Retrieval Conference.



# RQ1: How Do Disfluencies Impact Summarization Quality?

We inject disfluencies according to fixed distributions, similar to previous work:<sup>1,2</sup>

## Repeats and Interjections

We sample from  $X \sim N(\mu=10, \sigma=1)$  to determine the position at which the term(s) should be injected into the transcript  $N$  times.

- The interjections are uniformly randomly selected from: *uh, um, well, like, so, okay, I mean, you know.*

## False Starts

Sentences  $>4$  words are non-uniformly sampled with 80/20 probability with replacement, and the selected sentences have a false start (first 2 words of sentence) interjected  $N$  times.

### Original

Hello and welcome to our podcast! Let's get right to it. Today we're going to be interviewing a very special guest, someone I know you guys have been excited about having on the show.

### Repeats with N=3

Hello and welcome to our podcast! Let's get **get get get** right to it. Today we're going to be interviewing a **a a a** very special guest, someone I know you guys have been excited about having on the show.

### Interjections with N=3

Hello and welcome to our podcast! Let's get right **uh okay okay** to it. Today we're going to be interviewing a very special **um so I mean** guest, someone I know you guys have been excited about having on the show.

### False Starts with N=3

Hello and welcome to our podcast! Let's get right to it. Today we're **today we're today we're today we're** going to be interviewing a very special guest, someone I know you guys have been excited about having on the show.

<sup>1</sup>Shaolei Wang, Wangxiang Che, et al. 2020. Multi-task self-supervised learning for disfluency detection. In AAAI Conference on Artificial Intelligence, volume 34, pages 9193–9200.

<sup>2</sup>Tatiana Passali, Thanassis Mavroupos, et al. 2022. LARD: Large-scale artificial disfluency generation. In Language Resources and Evaluation Conference, pages 2327–2336.

# ***RQ1: How Do Disfluencies Impact Summarization Quality?***

We consider 6 summarization models:

**1min** is the first minute of transcript text.<sup>1</sup>

**cued\_speechUniv2** is an ensemble of 3 BART models plus a hierarchical filtering model, and it is the top performer from the TREC 2020 Podcasts Track.<sup>2</sup>

**BART** is a sequence-to-sequence model with a bidirectional encoder and a left-to-right autoregressive decoder.

**T5** is a text-to-text transformer model.

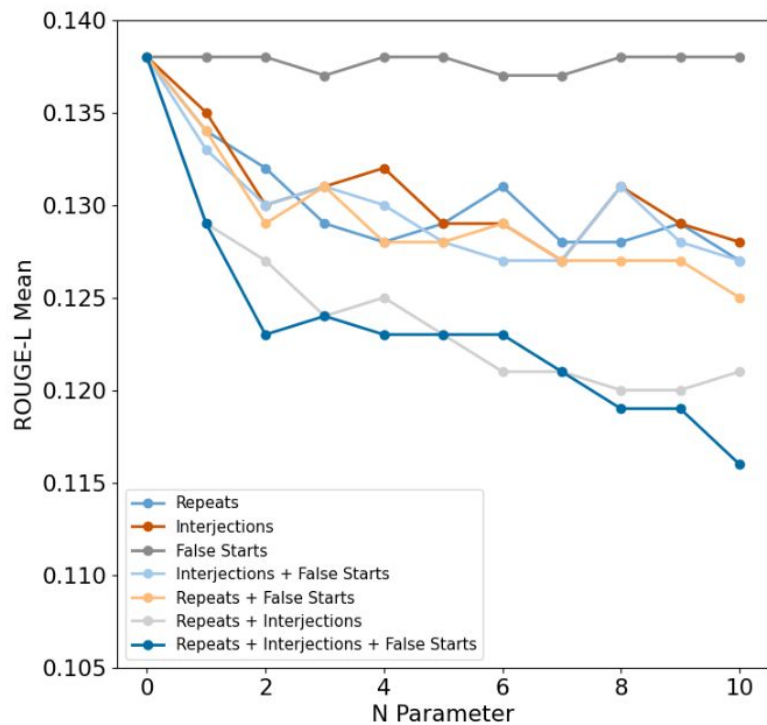
**Pegasus** is a transformer model with a pretraining objective called gap sentence generation.

**Llama 2-Chat** is a large transformer model which is pretrained and specifically for chat settings using RLHF.

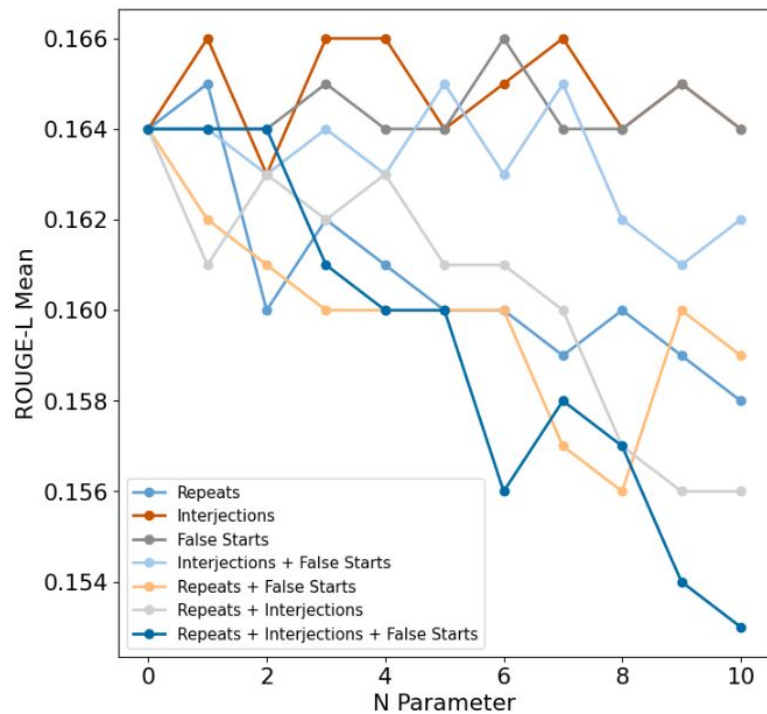
<sup>1</sup>Rosie Jones, Ben Carterette, Ann Clifton, et al. 2020. TREC 2020 Podcasts Track Overview. In Text Retrieval Conference.

<sup>2</sup>Potsawee Manakul and Mark Gales. 2020. Cued\_speech at TREC 2020 podcast summarisation track. In Text Retrieval Conference.

# RQ1: How Do Disfluencies Impact Summarization Quality?



(a) ROUGE-L over increased  $N$  on **BART** model.



(b) ROUGE-L over increased  $N$  on **cued\_speechUniv2** model.

# RQ1: How Do Disfluencies Impact Summarization Quality?

Model	N=0	N=2 vs. N=0	R	I	F	I+F	R+F	R+I	R+I+F
1min	0.124	$N_2 - N_0$ $\Delta$	-0.013 (-10.3%)	-0.014 (-11.6%)	-0.004 (-3.2%)	-0.017 (-14.0%)	-0.016 (-13.0%)	-0.026 (-21.0%)	-0.029 (-23.7%)
BART	0.138	$N_2 - N_0$ $\Delta$	-0.006 (-4.6%)	-0.008 (-5.5%)	0.000 (-0.3%)	-0.008 (-5.7%)	-0.008 (-6.1%)	-0.011 (-7.6%)	-0.015 (-11.1%)
T5	0.134	$N_2 - N_0$ $\Delta$	-0.018 <b>(-13.7%)</b>	-0.010 (-7.4%)	-0.003 (-2.4%)	-0.013 (-9.9%)	-0.018 <b>(-13.7%)</b>	-0.025 <b>(-19.0%)</b>	-0.032 <b>(-23.7%)</b>
Pegasus	0.131	$N_2 - N_0$ $\Delta$	-0.011 (-8.8%)	-0.014 <b>(-10.4%)</b>	-0.003 <b>(-2.6%)</b>	-0.017 <b>(-12.9%)</b>	-0.014 (-10.7%)	-0.023 (-17.2%)	-0.026 (-19.9%)
cued_speechUniv2	0.164	$N_2 - N_0$ $\Delta$	-0.004 (-2.5%)	-0.001 (-0.8%)	-0.001 (-0.5%)	-0.001 (-0.8%)	-0.003 (-1.9%)	-0.002 (-1.0%)	-0.001 (-0.5%)
Llama 2-Chat	0.129	$N_2 - N_0$ $\Delta$	-0.001 (-0.6%)	-0.002 (-1.2%)	-0.001 (-1.0%)	-0.002 (-1.6%)	-0.001 (-1.1%)	-0.001 (-1.1%)	-0.002 (-1.5%)

- Overall drop in ROUGE-L with increased N.
- T5 and Pegasus are the least resilient in the presence of disfluencies, BART is moderately resilient, and cued\_speechUniv2 and Llama 2-chat are the most resilient.

## ***RQ2: Can Summarization Quality be Improved By Directly Modeling Disfluency?***

We use a state-of-the-art, parsing-based disfluency annotation model<sup>1</sup> (Equations 1 and 2) to transform the transcripts via:

$$s(T) = \sum_{(i,j,l) \in T} s(i, j, l) \quad (1)$$

- **Repairing:** Removal of words marked disfluent.
- **Tagging:** Tagging (<DIS> and <\DIS>) of words marked disfluent.

$$\hat{T} = \operatorname{argmax}_T s(T) \quad (2)$$

<sup>1</sup>Paria Jamshid Lou and Mark Johnson. 2020. Improving disfluency detection by self-training a self-attentive model. In Association for Computational Linguistics, pages 3754–3763.

## RQ2: Can Summarization Quality be Improved By Directly Modeling Disfluency?

- Simply using the test set as-is yields the best ROUGE scores in most cases.
  - However, Pegasus is more robust in the face of missing information, and benefits from having the disfluencies removed.

### Inference-Only

Model	Test	Rouge-L	Rouge-1	Rouge-2
BART	test <sub>R</sub>	0.137	0.211	0.053
	test	<b>0.138</b>	<b>0.212</b>	<b>0.054</b>
	test <sub>T</sub>	0.137	0.209	0.052
Pegasus	test <sub>R</sub>	<b>0.131</b>	<b>0.200</b>	0.047
	test	<b>0.131</b>	0.198	<b>0.049</b>
	test <sub>T</sub>	0.113	0.169	0.038
T5	test <sub>R</sub>	0.133	0.194	0.050
	test	<b>0.134</b>	<b>0.199</b>	<b>0.051</b>
	test <sub>T</sub>	0.126	0.181	0.048

## RQ2: Can Summarization Quality be Improved By Directly Modeling Disfluency?

train	test	BART			T5			Pegasus		
		R-L	R-1	R-2	R-L	R-1	R-2	R-L	R-1	R-2
	test <sub>R</sub>	0.172	0.240	0.085	0.145	0.197	0.059	0.129	0.174	0.049
train <sub>R</sub>	test	<b>0.177</b>	<b>0.244</b>	<b>0.090</b>	0.146	0.196	<b>0.060</b>	<b>0.131</b>	0.177	<b>0.052</b>
	test <sub>T</sub>	0.174	0.241	0.086	0.148	0.198	0.063	0.096	0.133	0.037
	test <sub>R</sub>	0.170	0.236	0.083	0.146	0.198	0.060	0.122	0.165	0.045
train	test	<b>0.175</b>	<b>0.242</b>	<b>0.088</b>	<b>0.149</b>	<b>0.200</b>	<b>0.062</b>	<b>0.126</b>	<b>0.169</b>	<b>0.049</b>
	test <sub>T</sub>	0.172	0.238	0.085	0.147	0.194	<b>0.065</b>	0.090	0.124	0.032
	test <sub>R</sub>	0.172	0.238	0.083	0.142	0.193	0.057	0.129	<b>0.193</b>	0.048
train <sub>T</sub>	test	0.173	0.240	0.085	0.143	0.194	0.057	0.127	<b>0.193</b>	0.047
	test <sub>T</sub>	0.169	0.235	0.081	0.145	0.196	0.058	0.115	0.146	0.038

*We find that training on the repaired transcripts (train<sub>R</sub>) and testing on the original transcripts (test) yields the best results.*



***Link to our code on GitHub!***



# Conclusion

- **Disfluencies** are a key characteristic of **spoken content**.
  - We study 3 types of disfluencies -- *repeats, interjections, and false starts* -- in terms of the **Shriberg disfluency definition**.<sup>1</sup>
- We **synthetically inject disfluencies (N)** and find that **summarization quality decreases with increased disfluency**.
  - Decreases the most with combinations of the 3 disfluency types.
- We use a **parsing-based SOTA disfluency annotator**<sup>2</sup> to repair the disfluencies via removal and tagging.
- We find that for inference: Simply using the test set as-is yields the best ROUGE scores in most cases.
  - Pegasus is more robust in the face of missing information, and benefits from having the disfluencies removed.
- We find that for fine-tuning + inference: Training on the repaired transcripts ( $\text{train}_R$ ) and testing on the original transcripts (test) yields the best results.

<sup>1</sup>Elizabeth Ellen Shriberg. 1994. Preliminaries to a theory of speech disfluencies. Ph.D. thesis.

<sup>2</sup>Paria Jamshid Lou and Mark Johnson. 2020. Improving disfluency detection by self-training a self-attentive model. In Association for Computational Linguistics, pages 3754–3763.

# Quantifying the Impact of Disfluency on Spoken Content Summarization

**Maria Teleki, Xiangjue Dong, James Caverlee**

Texas A&M University

College Station, Texas, USA

{mariateleki, xj.dong, caverlee}@tamu.edu

*In LREC-COLING 2024*

