

DRES: BENCHMARKING LLMs FOR DISFLUENCY REMOVAL

Maria Teleki, Sai Janjur, Haoran Liu, Oliver Grabner, Ketan Verma, Thomas Docog, Xiangjue Dong, Lingfeng Shi, Cong Wang, Stephanie Birkelbach, Jason Kim, Yin Zhang, James Caverlee

Texas A&M University

ABSTRACT

Disfluencies – such as “um,” “uh,” interjections, parentheticals, and edited statements – remain a persistent challenge for speech-driven systems, degrading accuracy in command interpretation, summarization, and conversational agents. We introduce DRES (Disfluency Removal Evaluation Suite), a controlled text-level benchmark that establishes a reproducible semantic upper bound for this task. DRES builds on human-annotated Switchboard transcripts, isolating disfluency removal from ASR errors and acoustic variability. We systematically evaluate proprietary and open-source LLMs across scales, prompting strategies, and architectures. Our results reveal that (i) simple segmentation consistently improves performance, even for long-context models; (ii) reasoning-oriented models tend to over-delete fluent tokens; and (iii) fine-tuning achieves near state-of-the-art precision and recall but harms generalization abilities. We further present a set of LLM-specific error modes and offer nine practical recommendations (R1-R9) for deploying disfluency removal in speech-driven pipelines. DRES provides a reproducible, model-agnostic foundation for advancing robust spoken-language systems.

Index Terms— Disfluency removal, LLMs, benchmark, Switchboard, speech applications

1. INTRODUCTION

Disfluencies – *um, uh, interjections, parentheticals, and edited statements* – are pervasive in conversational speech, yet absent from most written text [1, 2, 3]. As speech-driven interfaces proliferate (e.g., Siri, Alexa, ChatGPT and Gemini voice modes, smart speakers, and emerging smart glasses) robust handling of disfluencies has become essential.

However, large language models (LLMs) trained primarily on written text often degrade in performance when processing spoken input with disfluencies. Prior studies show drops in voice command accuracy [4], conversational recommendation quality [5], and summarization fidelity [6, 7]. This mismatch is further amplified by ASR pipelines: while modern systems such as Whisper achieve strong recognition, they frequently omit disfluencies [8], preventing end-to-end models from being trained on realistic distributions. Con-

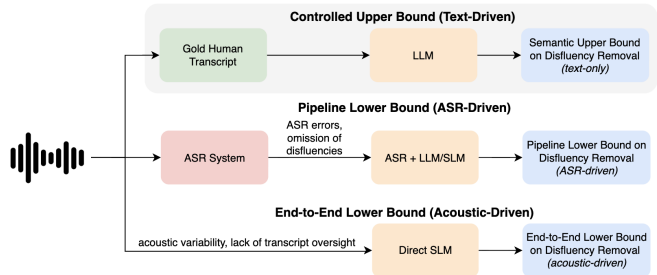


Fig. 1. Speech pipelines introduce ASR and acoustic errors that mask true disfluency removal ability (practical lower bounds). In contrast, **controlled evaluation on gold transcripts provides a semantic upper bound**, isolating disfluency removal performance without recognition confounds.

sequently, explicit benchmarking of disfluency handling remains critical.

In this work, we deliberately focus on LLMs rather than speech language models (SLMs). LLMs provide a controlled platform for isolating the disfluency removal task, establishing a semantic upper bound on performance that is not confounded by acoustic variability or ASR errors as shown in Figure 1. Their flexible scaling – from edge-ready small models to very large proprietary systems – makes them particularly suitable for systematic comparison across architectures, scales, and prompting conditions. While the framework naturally extends to SLMs in future audio-based or hybrid evaluations, LLMs allow us to first define a reproducible baseline and clarify the limits of disfluency removal.

In this paper, we introduce **DRES: the Disfluency Removal Evaluation Suite**, the first large-scale, systematic benchmark of LLMs for this task. DRES builds on the Switchboard corpus, enabling controlled and reproducible evaluation. We make the following contributions:

- First, we introduce **DRES, a reproducible benchmark for disfluency removal based on Switchboard**, enabling systematic comparison of models across scales and modalities. We make the code and additional analysis available at <https://github.com/mariateleki/dres>.
- Second, we conduct a **comprehensive evaluation of open-source and proprietary LLMs**, including instruction-

LLM	Citation	Open-Source	Instruct	Sizes	Architecture	Context Length	Features
gpt-4o	[9]	✗	✓	Nx200B*	MoE*	128k	Multimodal, Instruct
gpt-4o-mini	[9]	✗	✓	Nx8B*	MoE*	128k	Multimodal, Instruct
o4-mini	[9]	✗	✓	Not disclosed.	MoE*	200k	Multimodal, Instruct, Reasoning
Llama-3.1	[10]	✓	✓	8B	Dense	128k	Instruct
Llama-3.2	[10]	✓	✓	1B,3B	Dense	128k	Instruct
Llama-3.3	[10]	✓	✓	70B	Dense	128k	Multimodal, Instruct
MobileLLM	[11]	✓	✗	125M, 350M, 600M, 1B	Dense	2048+	Small for edge devices
Qwen3	[12]	✓	✓	0.6B, 1.7B, 4B, 8B	Dense + MoE	32,768+	Instruct
Phi-4-mini	[13]	✓	✓	3.8B	Dense	128k	Instruct, Reasoning

Table 1. LLMs: * indicates rumored sizes. o4-mini is available in high/medium reasoning variants. While some models include multimodal features, they are primarily text-based models (e.g., gpt models use a pipeline approach with Whisper).

tuned versus base, dense versus MoE, and small to very large models, revealing new trends in context handling, segmentation, and few-shot learning.

- Third, we provide the **first taxonomy of LLM-specific failure modes in disfluency removal**, such as over-deletion collapse, under-deletion, and reasoning-driven misinterpretation.
- Finally, we offer **practical recommendations for deployment (R1-R9) by connecting benchmark results to downstream tasks** where disfluencies are known to impair performance, including voice commands, conversational recommendation, and summarization, thereby providing actionable guidance for practitioners.

2. RELATION TO PRIOR WORK

2.1. Switchboard Dataset

We define the dataset T as:

$$T = \{(t_{tree}, t_{fluent}, t_{tag}, t_{disfluent})_i\}_{i=1}^N$$

where each element is a 4-tuple consisting of a fluent and disfluent version of the same utterance. We construct T from the Switchboard dataset [14, 15] following [16], where t_{fluent} , $t_{disfluent}$ are constructed from the tree, t_{tree} , using a top-down recursive approach. However, we do not remove partial words and punctuation, since ASR systems have since improved [17]. **Disfluency Definition.** We adopt the Shriberg [18] annotation scheme, where INTJ, PRN, and EDITED nodes are marked as disfluent [19, 16]; all others are fluent. **Syntactic Ambiguity** English syntax is often ambiguous, with sentences admitting multiple valid parse trees. In Switchboard, human annotators resolve this ambiguity using domain and discourse context, ensuring accurate identification of disfluencies [14]. As a result, Switchboard parse trees are more reliable than model-generated ones, particularly in disfluent contexts. **Correctly Transcribed Conversations.** Another important aspect of the Switchboard dataset [14] is that these parse trees are built on correctly-transcribed conversations [18]. State-of-the-art ASR models are known to significantly under-transcribe disfluencies [20] – even simple

Previous Models	\mathcal{E}_F	\mathcal{E}_P	\mathcal{E}_R
BERT Parser [19]	94.8	92.5	97.2
EGBC [21]	90.9	95.9	86.3
Noisy BiLSTM [21]	92.2	94.7	89.8
Weight Sharing [22]	91.1	92.1	90.2
BiLSTM [23]	85.9	91.6	80.3
Semi-CRF [23]	85.4	90.0	81.2

Table 2. Related Work: Previous models evaluated in terms of word-based precision, recall, and F-score.

disfluencies, such as *uh* and *um*. Hence, it is not possible to automatically obtain transcriptions of this quality.

2.2. LLMs for the Disfluency Removal Task

We evaluate a broad set of LLMs (Table 1) that vary along several key dimensions. We compare *open-source* vs. *proprietary* models, *instruct* vs. *base* variants, and three *size categories* (small $\leq 1B$ parameters for edge use, mid-sized 1B–8B, and large $\geq 8B$). We also study different *architectures* (dense vs. mixture-of-experts), *context lengths* (full vs. segmented transcripts to address long-context limitations), and special *features* such as multimodality and reasoning.

While specialized disfluency removal methods have been shown in Table 2 to perform well, LLMs offer a few key potential advantages for the disfluency removal task looking forward: (1) LLMs can leverage internet-scale pretraining for richer semantic understanding of domain terms, and (2) LLMs can use continual learning to track evolving language, i.e., slang, new terms, and shifting usage patterns.

2.3. Evaluation

All evaluations are performed on gold-standard Switchboard transcripts, which explicitly mark disfluent segments. This design isolates disfluency handling from ASR errors and ensures a controlled, reproducible benchmark. Model outputs are aligned with $t_{disfluent}$ using a variation of Gestalt pattern matching. This text-level setup provides a clean **upper bound for disfluency removal performance** and is complementary to audio-based SLM evaluation, which necessarily conflates disfluency handling with ASR quality.

\mathcal{E} -Scores. We report word-level precision (\mathcal{E}_P), recall (\mathcal{E}_R), and F1 (\mathcal{E}_F), following prior work on disfluency detec-

		\mathcal{E}_F	\mathcal{E}_P	\mathcal{E}_R	\mathcal{Z}_E	\mathcal{Z}_I	\mathcal{Z}_P			\mathcal{E}_F	\mathcal{E}_P	\mathcal{E}_R	\mathcal{Z}_E	\mathcal{Z}_I	\mathcal{Z}_P
gpt-4o								gpt-4o-mini							
f	0	74.19 _{7.87}	77.91 _{10.63}	73.18 _{12.98}	83.86 _{13.15}	71.25 _{16.57}	52.52 _{31.91}	70.52 _{9.22}	75.56 _{9.89}	68.17 _{13.90}	86.13 _{10.30}	60.38 _{18.22}	55.26 _{26.17}		
f	1	76.13 _{7.63}	76.94 _{10.48}	78.02 _{13.05}	85.11 _{11.73}	77.37 _{15.64}	62.71 _{29.59}	68.85 _{9.53}	74.35 _{9.65}	66.02 _{14.04}	86.05 _{10.37}	57.62 _{19.88}	52.04 _{24.06}		
f	3	73.99 _{3.33}	81.33 _{10.10}	70.75 _{15.11}	77.83 _{17.30}	72.93 _{15.26}	49.35 _{29.73}	68.70 _{0.64}	73.20 _{11.59}	67.68 _{15.21}	86.64 _{10.10}	61.42 _{20.83}	51.58 _{25.42}		
f	5	73.06 _{8.04}	81.24 _{10.69}	68.75 _{13.39}	76.92 _{15.55}	71.49 _{13.97}	42.39 _{30.38}	69.03 _{11.17}	76.93 _{13.18}	65.86 _{14.69}	84.77 _{11.10}	59.46 _{19.58}	48.78 _{24.87}		
s	0	78.17 _{4.04}	78.44 _{3.83}	78.30 _{5.27}	82.92 _{8.67}	77.55 _{7.16}	69.80 _{20.56}	72.69 _{5.79}	75.61 _{7.05}	70.48 _{7.35}	85.20 _{8.23}	61.89 _{11.08}	65.02 _{20.99}		
s	1	82.38 _{4.18}	83.61 _{6.03}	81.53 _{5.17}	83.77 _{9.41}	79.74 _{6.94}	79.65 _{19.40}	77.69 _{5.08}	81.84 _{6.00}	74.34 _{6.81}	84.56 _{9.47}	65.46 _{10.36}	77.59 _{20.30}		
s	3	83.72 _{3.85}	85.64 _{5.28}	82.22 _{5.41}	83.42 _{8.90}	81.94 _{6.66}	78.26 _{19.29}	77.01 _{1.82}	82.10 _{7.73}	72.88 _{6.56}	83.50 _{9.32}	64.80 _{9.94}	73.83 _{17.24}		
s	5	84.52 _{3.35}	88.09 _{4.84}	81.56 _{5.15}	81.15 _{9.70}	82.97 _{6.39}	77.14 _{19.05}	77.76 _{1.69}	83.31 _{3.42}	73.29 _{6.61}	83.84 _{9.17}	65.12 _{10.14}	74.69 _{18.04}		
medium-o4-mini								high-o4-mini							
f	0	48.18 _{9.41}	33.04 _{9.81}	95.43 _{4.92}	96.90 _{6.46}	93.66 _{7.38}	96.51 _{7.17}	55.81 _{11.28}	41.40 _{12.59}	92.69 _{7.82}	95.49 _{7.10}	90.49 _{10.52}	93.48 _{12.34}		
f	1	40.71 _{8.79}	26.19 _{7.57}	97.49 _{2.44}	98.19 _{2.50}	96.40 _{4.27}	98.68 _{3.44}	46.03 _{1.09}	31.17 _{10.83}	96.57 _{0.03}	97.42 _{3.44}	95.39 _{6.38}	97.77 _{5.13}		
f	3	39.65 _{7.74}	25.20 _{6.42}	97.74 _{2.21}	98.67 _{1.94}	96.84 _{3.51}	98.40 _{4.24}	42.91 _{9.50}	28.11 _{8.49}	97.27 _{2.69}	98.02 _{2.93}	96.35 _{4.17}	97.64 _{5.85}		
f	5	38.68 _{7.02}	24.43 _{6.63}	98.07 _{2.02}	98.89 _{1.64}	97.45 _{3.03}	98.00 _{5.26}	41.71 _{3.26}	27.03 _{3.01}	97.62 _{3.15}	98.15 _{3.51}	97.04 _{3.03}	97.92 _{5.04}		
Llama-1B-Instruct								Llama-3B-Instruct							
f	0	35.27 _{17.26}	54.71 _{22.71}	36.72 _{26.13}	44.18 _{26.76}	35.97 _{27.80}	23.35 _{27.59}	58.45 _{10.66}	49.96 _{15.09}	78.76 _{14.57}	84.46 _{13.29}	77.40 _{15.91}	69.96 _{25.35}		
f	1	33.35 _{11.37}	28.67 _{17.55}	72.58 _{29.09}	75.84 _{26.84}	72.69 _{30.26}	66.53 _{35.43}	50.45 _{13.19}	41.22 _{17.23}	81.70 _{18.33}	86.01 _{17.16}	82.32 _{18.14}	71.11 _{29.67}		
f	3	32.20 _{11.45}	26.65 _{17.91}	77.77 _{28.82}	80.83 _{26.69}	77.41 _{29.70}	73.35 _{33.29}	49.24 _{15.20}	46.74 _{19.42}	69.55 _{24.60}	76.19 _{23.61}	69.10 _{25.99}	56.35 _{33.18}		
f	5	35.60 _{13.72}	34.99 _{21.79}	68.37 _{31.82}	73.13 _{29.46}	68.51 _{32.72}	59.15 _{39.18}	48.74 _{14.78}	50.54 _{17.93}	60.98 _{26.24}	77.52 _{26.10}	62.42 _{27.51}	45.47 _{32.88}		
s	0	61.88 _{8.82}	68.31 _{8.65}	57.49 _{7.39}	56.38 _{10.73}	70.81 _{7.12}	27.84 _{16.32}	67.52 _{5.90}	65.81 _{4.48}	70.20 _{8.80}	77.96 _{8.24}	69.57 _{9.17}	57.13 _{19.47}		
s	1	32.98 _{6.22}	29.09 _{8.45}	40.70 _{8.26}	43.15 _{13.15}	41.44 _{8.28}	34.47 _{18.00}	48.69 _{8.88}	45.74 _{13.31}	54.80 _{7.07}	57.20 _{10.45}	53.30 _{8.21}	53.65 _{19.15}		
s	3	38.26 _{6.12}	31.37 _{7.59}	51.71 _{8.19}	53.39 _{11.33}	56.34 _{8.26}	38.24 _{18.11}	53.78 _{8.89}	57.00 _{15.41}	52.86 _{7.06}	56.94 _{11.02}	54.27 _{9.37}	40.14 _{17.67}		
s	5	39.34 _{6.97}	30.24 _{7.93}	59.38 _{8.02}	62.51 _{10.85}	63.46 _{8.22}	44.34 _{17.61}	60.43 _{7.84}	71.60 _{13.85}	53.22 _{6.82}	60.13 _{11.05}	55.63 _{9.75}	34.58 _{16.03}		
Llama-8B-Instruct								Llama-70B-Instruct							
f	0	45.48 _{9.37}	31.93 _{9.82}	87.43 _{12.08}	88.57 _{13.03}	88.83 _{11.67}	81.12 _{22.77}	67.83 _{9.90}	63.90 _{16.75}	78.48 _{13.30}	81.14 _{13.92}	79.48 _{14.76}	69.48 _{26.23}		
f	1	37.46 _{11.94}	27.99 _{15.68}	80.38 _{22.39}	81.35 _{22.81}	81.58 _{22.97}	75.18 _{27.90}	62.85 _{12.61}	58.99 _{17.71}	74.55 _{15.59}	78.90 _{16.01}	74.84 _{16.51}	66.90 _{26.51}		
f	3	30.32 _{9.90}	18.02 _{4.19}	99.94 _{0.60}	100.00 _{0.00}	99.94 _{0.59}	99.85 _{2.03}	58.37 _{12.91}	48.95 _{17.06}	82.83 _{14.57}	85.32 _{12.83}	84.17 _{16.30}	75.74 _{23.83}		
f	5	30.33 _{9.91}	18.02 _{4.19}	100.00 _{0.00}	100.00 _{0.00}	100.00 _{0.00}	100.00 _{0.00}	53.37 _{14.54}	44.39 _{17.41}	82.87 _{18.75}	83.42 _{18.37}	85.37 _{18.61}	76.18 _{27.39}		
s	0	68.30 _{5.95}	63.97 _{8.94}	74.22 _{5.58}	80.84 _{8.19}	74.66 _{7.23}	60.62 _{20.33}	76.14 _{4.84}	77.87 _{7.82}	75.10 _{5.53}	73.61 _{10.81}	76.36 _{5.50}	71.56 _{21.92}		
s	1	68.90 _{9.98}	67.91 _{9.82}	70.60 _{6.10}	75.98 _{9.23}	66.07 _{8.72}	69.04 _{19.92}	68.31 _{8.47}	75.97 _{14.31}	63.25 _{6.43}	61.65 _{11.28}	60.78 _{8.25}	68.39 _{20.89}		
s	3	65.50 _{6.72}	66.80 _{9.56}	65.06 _{7.17}	71.32 _{9.56}	60.78 _{8.39}	63.12 _{20.91}	63.20 _{8.53}	72.30 _{15.56}	57.72 _{7.17}	53.18 _{11.18}	56.67 _{9.35}	64.54 _{21.82}		
s	5	66.65 _{6.50}	67.05 _{9.37}	66.98 _{6.41}	74.17 _{9.45}	63.13 _{9.31}	62.23 _{21.11}	65.99 _{7.41}	76.46 _{13.85}	59.26 _{6.86}	54.21 _{11.52}	59.15 _{9.26}	64.88 _{20.04}		
Qwen3-0.6B								Qwen3-1.7B							
f	0	18.04 _{1.30}	43.29 _{26.59}	16.46 _{14.61}	22.51 _{17.29}	14.26 _{14.25}	11.56 _{15.87}	10.25 _{9.49}	86.02 _{20.54}	5.91 _{7.31}	10.36 _{9.49}	4.10 _{7.96}	2.75 _{6.48}		
f	1	22.36 _{9.02}	29.02 _{20.01}	30.24 _{26.06}	35.09 _{26.05}	28.20 _{26.43}	27.18 _{29.81}	10.12 _{9.34}	79.60 _{29.76}	13.39 _{28.55}	16.82 _{27.74}	11.84 _{29.26}	10.70 _{28.97}		
f	3	21.17 _{9.28}	38.69 _{23.36}	20.34 _{14.32}	24.92 _{17.04}	19.29 _{13.91}	14.89 _{17.93}	7.78 _{7.27}	88.20 _{31.59}	5.09 _{10.74}	8.82 _{12.02}	3.24 _{10.26}	2.98 _{11.34}		
f	5	19.87 _{9.99}	40.95 _{26.37}	19.67 _{15.96}	24.09 _{18.15}	17.87 _{15.44}	16.06 _{19.42}	8.64 _{6.49}	85.42 _{24.87}	6.66 _{15.30}	10.71 _{15.52}	4.88 _{15.55}	4.17 _{15.88}		
s	0	43.74 _{6.31}	44.90 _{9.86}	44.05 _{7.56}	60.78 _{11.78}	33.34 _{8.34}	41.94 _{19.35}	36.00 _{8.24}	71.41 _{10.04}	24.79 _{7.60}	39.13 _{11.40}	18.96 _{7.61}	14.56 _{15.02}		
s	1	51.97 _{6.39}	47.23 _{9.62}	59.72 _{6.79}	70.39 _{9.05}	55.32 _{8.05}	53.93 _{17.38}	35.78 _{7.46}	65.77 _{11.56}	25.10 _{6.39}	23.24 _{7.90}	33.64 _{8.64}	8.18 _{6.65}		
s	3	48.90 _{6.31}	45.89 _{8.84}	54.00 _{8.14}	62.76 _{10.75}	52.82 _{9.24}	42.36 _{18.94}	31.26 _{5.05}	67.30 _{10.40}	20.75 _{8.80}	18.55 _{7.01}	28.84 _{8.57}	5.09 _{7.21}		
s	5	48.38 _{6.33}	48.74 _{10.75}	50.29 _{8.46}	58.66 _{11.18}	50.05 _{9.90}	36.34 _{18.73}	34.46 _{7.97}	81.04 _{11.03}	22.31 _{6.50}	23.00 _{8.94}	29.00 _{8.67}	5.08 _{6.61}		
Qwen3-4B								Qwen3-8B							
f	0	66.39 _{16.56}	75.58 _{16.33}	64.58 _{20.96}	62.47 _{23.52}	70.30 _{22.66}	49.66 _{31.16}	71.04 _{10.77}	69.94 _{12.12}	76.17 _{15.49}	75.24 _{18.17}	79.32 _{14.48}	67.37 _{29.01}		
f	1	59.89 _{20.20}	79.63 _{14.52}	54.23 _{23.08}	55.44 _{26.35}	59.91 _{24.54}	35.45 _{28.72}	68.86 _{13.90}	74.54 _{15.50}	69.69 _{17.63}	67.92 _{20.37}	76.12 _{16.61}	54.80 _{30.83}		
f	3	62.29 _{20.97}	80.58 _{14.28}	57.13 _{22.87}	54.85 _{23.60}	65.16 _{24.09}	39.32 _{28.43}	71.77 _{11.73}	76.82 _{10.36}	70.78 _{16.75}	68.89 _{18.34}	77.84 _{15.80}	54.49 _{30.81}		
s	0	68.48 _{5.55}	67.96 _{8.94}	69.95 _{6.29}	70.68 _{10.41}	72.43 _{7.71}	62.02 _{19.69}	71.46 _{6.01}	78.32 _{7.62}	66.17 _{7.12}	73.77 _{10.23}	63.29 _{10.24}	57.96 _{22.96}		
s	1	69.66 _{5.59}	80.77 _{7.32}	61.60 _{6.43}	72.82 _{10.90}	54.30 _{9.35}	57.76 _{22.39}	70.91 _{5.30}	70.41 _{8.29}	72.25 _{6.22}	73.13 _{10.83}	74.45 _{8.67}	62.03 _{22.19}		
s	3	64.14 _{6.23}	80.56 _{7.25}	53.72 _{7.15}	66.76 _{10.58}	47.47 _{10.41}	44.88 _{22.09}								

	\mathcal{E}_P	\mathcal{E}_I	\mathcal{E}_R	\mathcal{Z}_E	\mathcal{Z}_I	\mathcal{Z}_P
gpt-4o-mini	72.69 _(.59)	75.61 _(7.05)	70.48 _(.35)	85.20 _(8.23)	61.89 _(1.08)	65.02 _(20.99)
gpt-4o-mini _{ft}	94.77 _(.09)	96.63 _(.11)	93.08 _(.64)	87.71 _(8.03)	97.39 _(.64)	89.31 _(16.66)
Llama-3B	67.52 _(.52)	65.81 _(8.48)	70.20 _(.80)	77.96 _(8.24)	69.57 _(6.17)	57.13 _(19.47)
Llama-3B _{ft}	91.10 _(.19)	93.26 _(3.49)	89.21 _(.05)	83.39 _(7.23)	94.06 _(.96)	86.36 _(17.26)

Table 4. Fine-Tuning (s).

	GSM8K		MMLU		CoQA	
	EM _{Strict} ↑	EM _{Flexible} ↑	EM _{Strict} ↑	EM _{Strict} ↑	F1 ↑	F1 ↑
gpt-4o-mini	0.86 _(.01)	0.87 _(.01)	0.76 _(.00)	0.58 _(.02)	0.77 _(.03)	0.77 _(.03)
gpt-4o-mini _{ft}	0.51 _(.01)	0.81 _(.01)	0.69 _(.00)	0.55 _(.02)	0.73 _(.01)	0.73 _(.01)
Llama-3B	0.27 _(.01)	0.70 _(.01)	0.53 _(.00)	0.17 _(.02)	0.35 _(.02)	0.35 _(.02)
Llama-3B _{ft}	0.23 _(.01)	0.38 _(.01)	0.28 _(.00)	0.05 _(.01)	0.15 _(.01)	0.15 _(.01)

Table 5. Generalization of Fine-Tuned Models.

tion. **Z-Scores.** We further exploit parse tree annotations to measure the proportion of disfluent nodes removed: edited-type nodes (\mathcal{Z}_E), interjection nodes (\mathcal{Z}_I), and parenthetical nodes (\mathcal{Z}_P). These fine-grained scores capture structural differences in model behavior and highlight failure modes such as over-deletion and under-deletion.

3. RESULTS & RECOMMENDATIONS

We analyze the disfluency removal behavior of LLMs and provide recommendations (R1-R9).

Open-Source vs. Proprietary. Looking to Table 3, proprietary models (gpt-4o, gpt-4o-mini) achieve the highest scores, with margins of 10–15 points over the best open-source alternatives. We attribute this to training exposure to Whisper-transcribed speech data [24]. **(R1) Proprietary models are currently the most reliable for production systems, while open-source models require targeted augmentation with spoken data.**

Segmentation (s) vs. Full Input (f). Segmenting transcripts consistently improves both mean performance and stability, e.g., gpt-4o improves from $\mathcal{E}_F=76.13$ (f) to 82.38 (s) at $k=1$. This supports prior evidence of long-context degradation in LLMs [25, 26]. **(R2) Segmentation is an effective preprocessing step that should be applied.**

Few-Shot Sensitivity (k). Increasing k does not uniformly improve results. Small models (e.g., MobileLLM) gain slightly, but others show degradation (e.g. Llama-3B/8B/70B) when more examples are provided. **(R3) Few-shot prompting should be used with caution, as some model families misinterpret exemplars and over-edit fluent text.**

Disfluency Category Performance. \mathcal{Z} -Scores show that EDITED nodes are handled well, but INTJ and PRN nodes are frequently missed, despite prior work suggesting these are the easiest to detect [19, 17]. **(R4) Future modeling should focus on under-served categories (INTJ, PRN) to improve robustness across all disfluency types.**

Over-Deletion Failures. Several models (e.g., Llama-8B, o4-mini) achieve near perfect recall but at the cost of very low precision, deleting fluent tokens. Segmentation often mitigates this collapse mode. **(R5) Segment-level evaluation**

helps reduce over-deletion risk.

Under-Deletion Failures. Some models (e.g., Qwen series) exhibit the opposite trend of over-deletion, achieving high precision but low recall (purple). These models preserve most fluent tokens but fail to remove many true disfluencies, especially in INTJ and PRN categories. This reflects conservative editing strategies and limited exposure to conversational disfluency distributions. **(R6) Models prone to under-deletion require additional filtering or targeted fine-tuning to ensure sufficient disfluency coverage.**

Reasoning-Oriented Models. Models tuned for reasoning (o4-mini, Phi-4) perform poorly, showing high recall but extreme over-deletion (blue). **(R7) Reasoning capability does not translate to disfluency removal; specialized evaluation remains necessary.**

Impact of Model Size. Model scaling generally improves disfluency removal, with Qwen, GPT, and Llama families showing upward trends. However, gains are nonlinear – e.g., Qwen3-1.7B underperforms both smaller and larger variants – likely due to training data or optimization differences rather than capacity limits. **(R8) Model choice should be guided by empirical benchmarks on target domains and disfluency categories rather than size alone.**

Fine-Tuning and Generalization. Looking to Tables 4 and 5, fine-tuning improves performance to near SOTA levels (e.g., gpt-4o-mini_{ft} achieves $\mathcal{E}_P=96.6$), but evaluation on GSM8K, MMLU, and CoQA shows degraded performance on unrelated tasks. **(R9) Fine-tuning is suitable for dedicated disfluency pipelines, but not for general-purpose conversational models.**

4. DISCUSSION & FUTURE WORK

DRES establishes a controlled, reproducible upper bound for disfluency removal, revealing consistent gaps between proprietary and open-source LLMs that are driven more by training exposure than model size. Our results highlight open challenges: closing the precision-recall trade-off (especially over-deletion in reasoning-tuned models), improving coverage for under-served categories (INTJ, PRN), and mitigating generalization loss from fine-tuning.

We see value in a modular approach where specialized disfluency removal components preprocess ASR output before downstream reasoning, helping maintain the flexibility of general purpose LLMs. Promising directions include exploring lightweight adapters, multi-task setups, and continual learning to approach state-of-the-art accuracy while mitigating catastrophic forgetting. As speech language models mature, incorporating disfluency handling into end-to-end architectures will require careful design to maintain generalization. Extending DRES to multilingual and multimodal settings and systematically evaluating its downstream impacts (e.g., command success rates, summarization fidelity) can help the community build more robust and deployable systems.

5. REFERENCES

- [1] Elizabeth Shriberg, “Disfluencies in switchboard,” in *International Conference on Spoken Language Processing*, 1996, vol. 96, pp. 11–14.
- [2] Heather Bortfeld, Silvia D. Leon, Jonathan E. Bloom, Michael F. Schober, and Susan E. Brennan, “Disfluency rates in conversation: Effects of age, relationship, topic, role, and gender,” *Language and Speech*, vol. 44, no. 2, pp. 123–147, 2001.
- [3] Sharon Oviatt, “Predicting spoken disfluencies during human-computer interaction,” *Computer Speech and Language*, 1995.
- [4] Yiming Chen, Xianghu Yue, Chen Zhang, Xiaoxue Gao, Robby T Tan, and Haizhou Li, “Voicebench: Benchmarking llm-based voice assistants,” *arXiv preprint arXiv:2410.17196*, 2024.
- [5] Maria Teleki, Lingfeng Shi, Chengkai Liu, and James Caverlee, “I want a horror – comedy – movie: Slips-of-the-Tongue Impact Conversational Recommender System Performance,” in *INTERSPEECH*, 2025.
- [6] Varun Nathan, Ayush Kumar, and Jithendra Vepa, “Investigating the role and impact of disfluency on summarization,” in *EMNLP: Industry Track*, 2023, pp. 541–551.
- [7] Fabian Retkowsky, Maike Züfle, Andreas Sudmann, Dinah Pfau, Shinji Watanabe, Jan Niehues, and Alexander Waibel, “Summarizing speech: A comprehensive survey,” 2025.
- [8] Maria Teleki, Xiangjue Dong, Soohwan Kim, and James Caverlee, “Comparing asr systems in the context of speech disfluencies,” in *Interspeech 2024*, 2024, pp. 4548–4552.
- [9] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al., “Gpt-4 technical report,” *arXiv preprint arXiv:2303.08774*, 2023.
- [10] Hugo Touvron, Thibault Sellam, Yacine Jernite, Jade Copet, Loubna Ben Allal, et al., “The llama 3 herd of models,” *Computing Research Repository*, vol. arXiv:2407.21783, 2024.
- [11] Zechun Liu, Changsheng Zhao, et al., “MobileLLM: Optimizing sub-billion parameter language models for on-device use cases,” 2024.
- [12] Qwen Team, “Qwen3 technical report,” 2025.
- [13] Microsoft, Abdelrahman Abouelenin, Atabak Ashfaq, Adam Atkinson, et al., “Phi-4-mini technical report: Compact yet powerful multimodal language models via mixture-of-loras,” 2025.
- [14] Marcus Mitchell, Beatrice Santorini, M Marcinkiewicz, and Ann Taylor, “Treebank-3 LDC99T42 Web Download,” 1999, [Online]. Available: <https://catalog.ldc.upenn.edu/LDC99T42>.
- [15] John J Godfrey, Edward C Holliman, and Jane McDaniel, “Switchboard: Telephone speech corpus for research and development,” in *ICASSP*, 1992, vol. 1, pp. 517–520.
- [16] Eugene Charniak and Mark Johnson, “Edit Detection and Parsing for Transcribed Speech,” *NAACL*, 2001.
- [17] Mark Johnson and Eugene Charniak, “A TAG-based noisy-channel model of speech repairs,” in *ACL*, 2004, pp. 33–39.
- [18] Elizabeth Shriberg, *Preliminaries to a Theory of Speech Disfluencies*, Ph.D. thesis, 1994.
- [19] Paria Jamshid Lou and Mark Johnson, “Improving disfluency detection by self-training a self-attentive model,” in *ACL*, 2020, pp. 3754–3763.
- [20] Dena Mujtaba, Nihar R Mahapatra, et al., “Lost in transcription: Identifying and quantifying the accuracy biases of automatic speech recognition systems against disfluent speech,” *arXiv preprint arXiv:2405.06150*, 2024.
- [21] Nguyen Bach and Fei Huang, “Noisy bilstm-based models for disfluency detection,” in *INTERSPEECH*, 2019, pp. 4230–4234.
- [22] Feng Wang, Wei Chen, Zhen Yang, Qianqian Dong, Shuang Xu, and Bo Xu, “Semi-supervised disfluency detection,” in *International Conference on Computational Linguistics*, 2018, pp. 3529–3538.
- [23] Vicky Zayats, Mari Ostendorf, and Hannaneh Hajishirzi, “Disfluency detection using a bidirectional LSTM,” in *INTERSPEECH*, 2016.
- [24] Cade Metz, Cecilia Kang, Sheera Frenkel, Stuart A. Thompson, and Nico Grant, “How tech giants cut corners to harvest data for a.i.,” *The New York Times*, 2024.
- [25] Jiaheng Liu, Dawei Zhu, et al., “A comprehensive survey on long context language modeling,” *arXiv*, 2025.
- [26] Tianle Li, Ge Zhang, Quy Duc Do, Xiang Yue, and Wenhui Chen, “Long-context llms struggle with long in-context learning,” *arXiv*, 2024.

© 2025 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.