



# HOW TO TRAIN YOUR FILLERS: UH AND UM IN SPONTANEOUS SPEECH SYNTHESIS

Éva Székely, Gustav Eje Henter, Jonas Beskow, Joakim Gustafson

Division of Speech, Music and Hearing, KTH Royal Institute of Technology, Stockholm, Sweden

Using spontaneous conversational speech for TTS raises questions on how disfluencies such as filled pauses (FPs) should be approached. Detailed annotation of FPs in training data enables precise control at synthesis time; coarse or non-existent FP annotation, when combined with stochastic attention-based neural TTS, leads to synthesisers that insert these phenomena into fluent prompts on their own accord. In this study we investigate, objectively and subjectively, the effects of FP annotation and the impact of relinquishing control over FPs in a Tacotron TTS system built from a conversational podcast corpus.

## Summary of the voice configurations and the conditions used in the evaluations

System	Corpus & training	Annotation of FPs	Condition	Prompt	Resulting speech
<b>AutoFP</b>	whole TCC	no	<b>AutoFP</b>	fluent	has automatically placed FPs
<b>CtrlFP</b>	whole TCC	yes, differentiating 'uh' and 'um'	<b>CtrlFP-GT</b>	FPs copied from GT	FPs exactly as in the prompt
			<b>CtrlFP-SW</b>	FPs opposite type as GT	FPs exactly as in the prompt
			<b>CtrlFP-FL</b>	fluent	no FPs
<b>GenFP</b>	whole TCC	yes, with a generic FP label for both 'uh' and 'um'	<b>GenFP</b>	Ground-truth FP locations, unspecified type	has FPs in specified locations, type is decided automatically
<b>HalfFluent</b>	fluent 44.4% of TCC	N/A (no FPs in the training data)	<b>HalfFluent</b>	fluent	no FPs
<b>TransFluent</b>	whole TCC, then transfer learning to fluent 44.4%	no	<b>TransFluent</b>	fluent	very occasional automatically placed FPs

**ThinkComputers Corpus (TCC)**

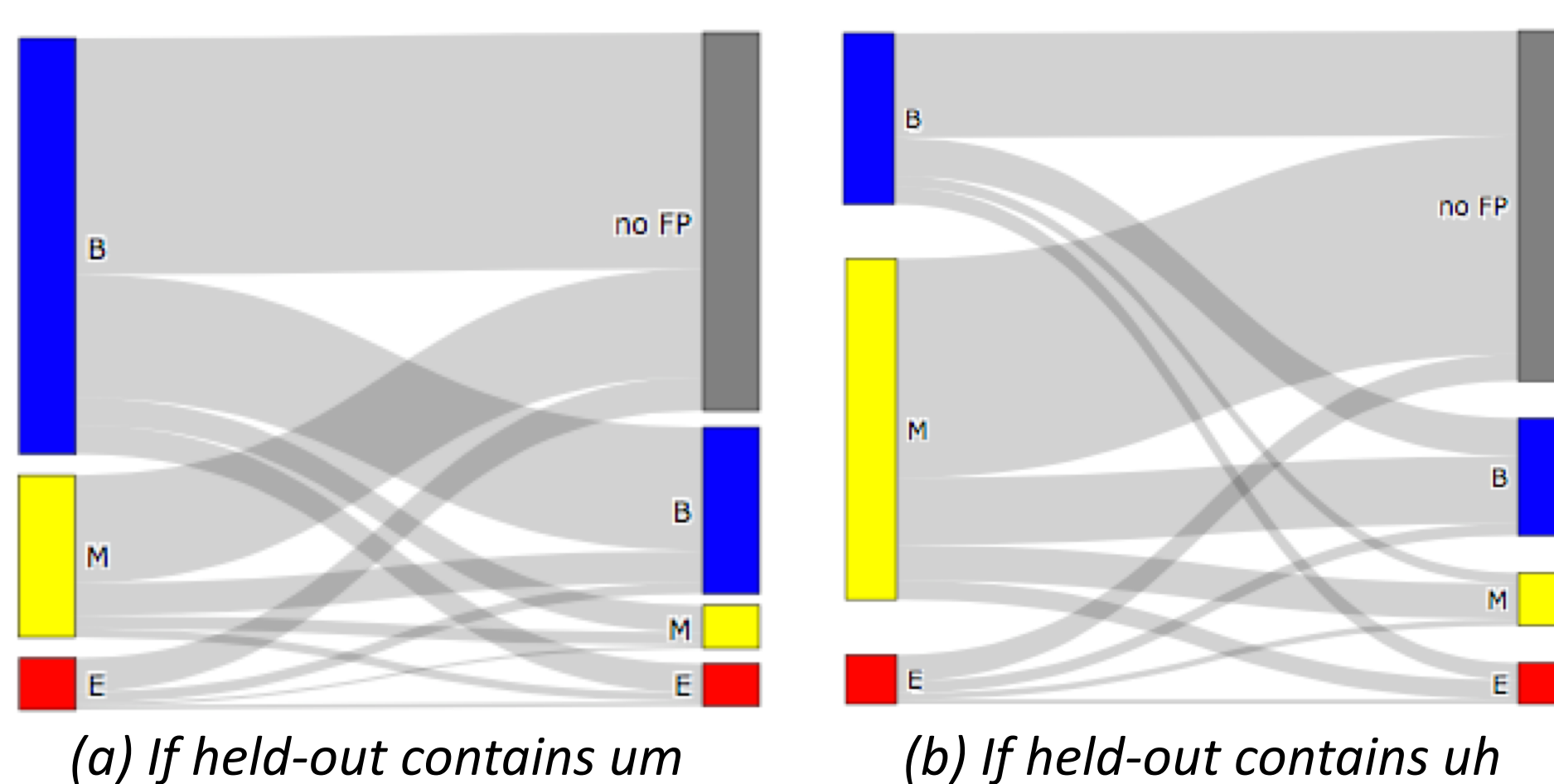
- Weekly tech podcast, spontaneous conversational speech
- Male AE speaker
- 9h, segmented to single-speaker breath groups [1]
- Automatic transcription with ASR and Gentle forced aligner

**TTS:** Tacotron + Griffin-Lim

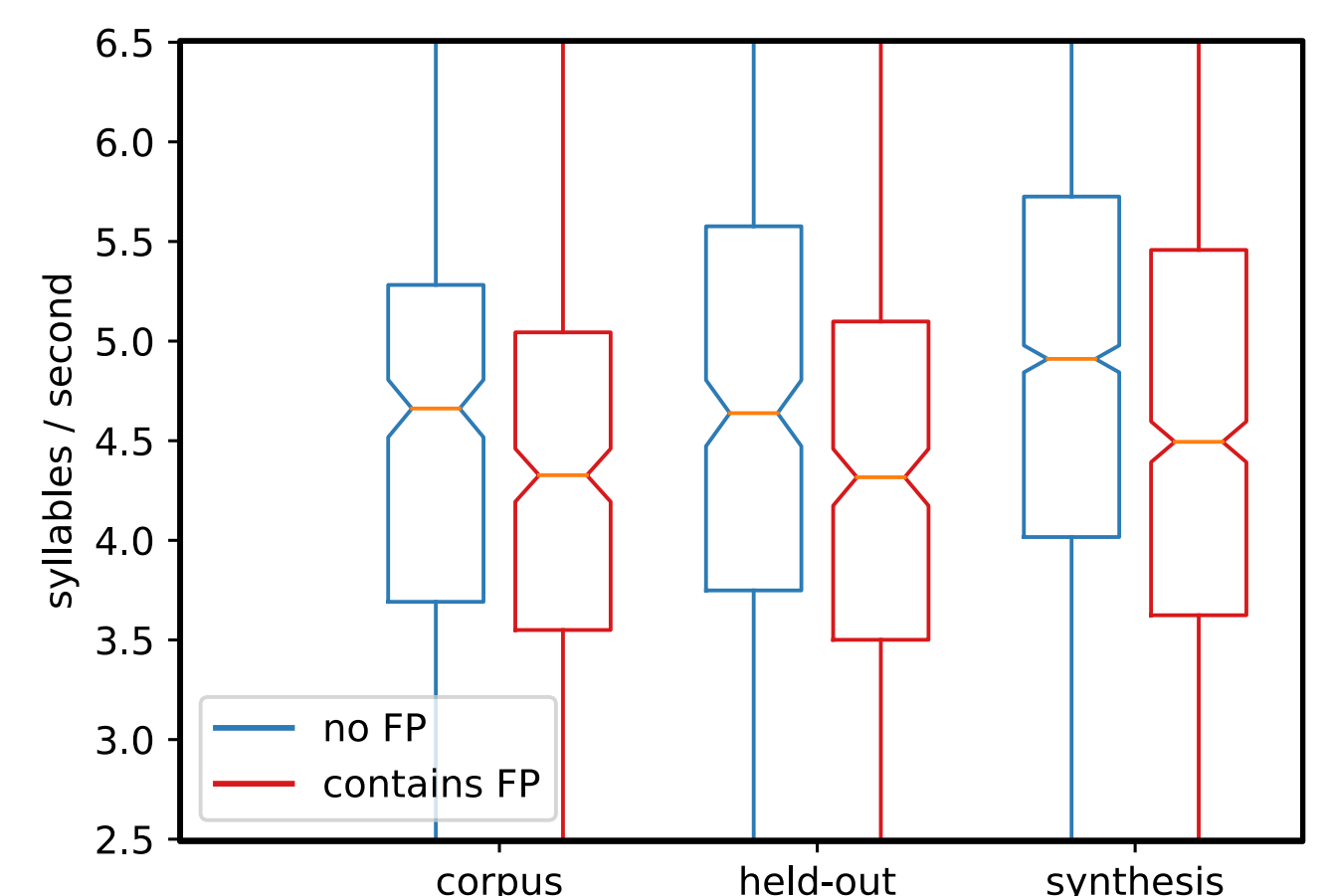
## Objective evaluation of automatic FP insertion

FP at	B	M	E	Held-out	AutoFP	p-val.
				49%	66%	<0.001
✓				23%	20%	0.109
	✓			17%	6%	<0.001
		✓		3%	5%	0.055
✓	✓			6%	1%	<0.001
✓		✓		1%	1%	0.844
	✓	✓		1%	1%	0.592
✓	✓	✓		0%	0%	0.200

Breath groups from held-out data with FPs at Beginning, Middle, and/or End

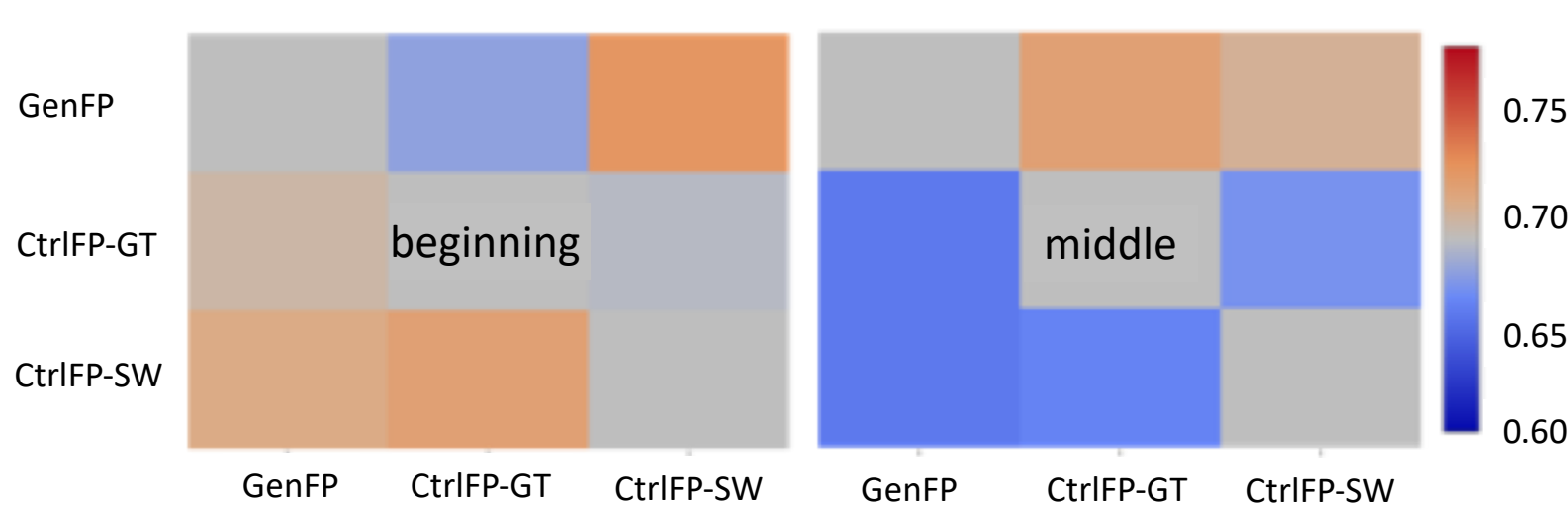


Sankey diagrams of FP position in the held-out data (left in each diagram) and the synthesis (right)



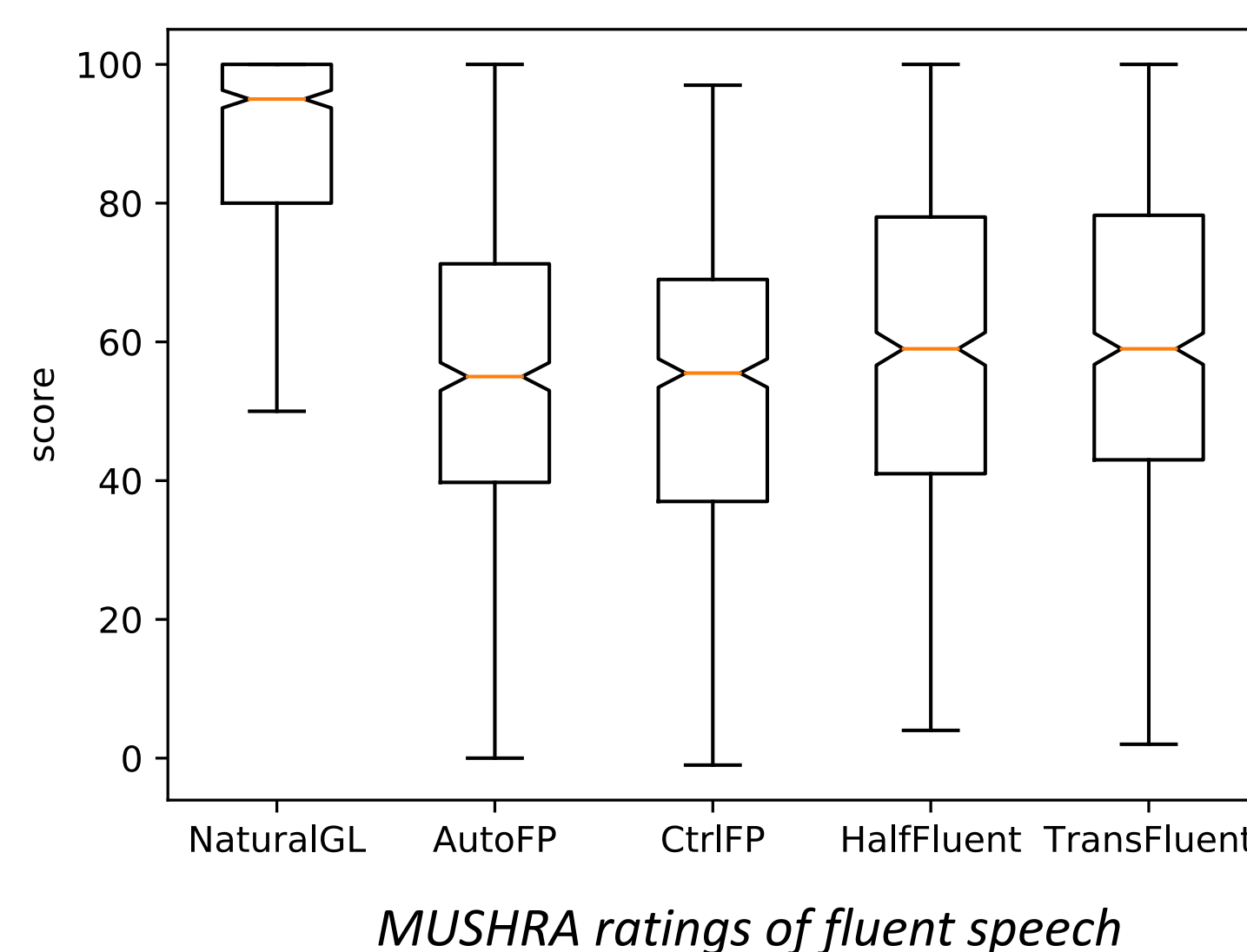
Speech rates within the TCC corpus, the held-out and the synthesised samples, split on whether or not they contain FPs

## Perceptual evaluation of disfluent speech



20 utterances from AMI corpus, each containing one 'uh' or 'um', in the beginning or the middle. Pairwise comparison: listeners indicated which version hesitated more realistically.

## Perceptual evaluation of fluent speech



## Bonus question

Would you want a robot to sound hesitant? Why?

**Yes – 45%** “Yes, sounds more authentic and genuine.”  
 “Yes so it sounds more like a person and more relatable.”  
 “Yes, much more easy to listen to for prolonged periods.”  
 “I think it's comforting to have a hesitant voice from them.”

**Undecided – 18%** “I have no idea.” “Indifferent.”

**No – 36%** No, because it would sound too human like. Over the phone, I wouldn't be able to tell I am talking to a robot.”  
 “No, I feel uncomfortable blurring the lines between what sounds naturally human and what is machine.”  
 “No as I would want it to speak correctly at all times.”

## Conclusions:

- ✓ Systems trained with no, or location-only FP annotation reproduce FPs in a similar pattern as in the corpus.
- ✓ Synthesiser-predicted FP types ('uh' or 'um') were preferred over specifying the ground-truth type.
- ✓ Using precise annotations and focusing on more fluent parts of the corpus improves naturalness of fluent TTS.

Listen here! →



[www.speech.kth.se/tts-demos](http://www.speech.kth.se/tts-demos)