Contextual Speech Recognition with Difficult Negative Training Examples Uri Alon, Golan Pundak, Tara N. Sainath

1. Task: Contextualized ASR

- **Context** provided in addition to audio can help reduce WER significantly.
- Such user-specific contextual information can include:
- The user's list of songs
- The user's contact list
- The currently installed apps
- **Proper nouns** are very frequent in various ASR tasks:
- "Call Joan's mobile"
- "Play *Taylor Swift*"
- "How tall is *LeBron James*?"
- But contextual ASR models usually perform poorly on rare words and especially on proper nouns (NNPs).

2. The Contextualized LAS (CLAS) Model (Pundak et al., SLT'2018)

- **CLAS** is an E2E ASR model based on the Listen-Attend-and-Spell (LAS) encoder-decoder architecture.
- The key difference from LAS: **biasing** sub-module.



3. The Problem: The Network Fails to Distinguish Between Phrases

4. Training with Difficult Negative Examples







urialon@cs.technion.ac.il

Disambiguation of similarly sounding phrases is challenging.

The network makes even more mistakes as the set of bias phrases becomes



During training, we provide the network with phonetically similar prope "distractors".

This way, we encourage the network to:

Distinguish between similarly sounding phrases

Learn more discriminative representations.



Reference transcript

Google

{golan,tsainath}@google.com

Phonetically Similar	5. Evaluation
	We experimented with the following training
	Vanilla CLAS CLAS
larger.	Bias Phrases Selection Random NNPs Distractors Selection Random Rando
	Results:
	Test Set Vanilla CLAS CLAS+NNP CLAS+fuzzy
	Songs 9.8 6.7 (31.6%) 10.4
	Contacts 11.3 $6.1(46.0\%)$ 16.5
	$[1alk-1o \ 15.2 \ 14.8 (2.6\%) \ 11.1 (27.0\%)]$
nny }	
	6. Qualitative Analysis
er nouns (NNPs) as the	True
	n/a
	**creepy carrots
	creep carrots
	creeping carrots
	crappie carrots
(Uiohn U	free carrots
	trippy carrots
s "jean",	sleepy carrots
"joan",	creepy car Fuzzy: croopy corroted
"johnny"}	Figure: The fuzzy model attends mostly to "cree
	"sleepy carrots" and predicts the wrong word "sl
Set of bias phrases	



ref: creepy carrots</bias>



epy carrots" and makes a correct prediction, while the non-fuzzy model attends to sleepy"