# Phone Features Improve Speech Translation



Elizabeth Salesky



Alan W Black





## Speech Translation



# **Challenges of Speech Input**

### 1 Length



**Discretized audio – speech frames** 

#### Impacts:

- Memory
- Distance between dependencies
- Training efficiency

#### $c h a r a c t e r s \rightarrow long sequences$



#### **Performance Impact:**

- **End-to-End** model performance varies based on dataset & size, language pair...
- **<u>Cascaded</u>** models perform better in many settings

# **Challenges of Speech Input**

#### • Low-Resource Settings



## Larger difference in performance between architectures End-to-end models do not see enough data to learn variation

## ③ Variation in number of frames per phone



### **1** End-to-End vs Cascade comparison

- Single dataset with much previous work: **Fisher Spanish-English**
- Compare multiple resource settings: HIGH (160hr) MED (40hr) LOW (20hr)

### **(2)** Phone Features to address challenges of speech input

Compare architectures: End-to-End Cascade



# **Models with Phone Features**

#### CASCADE

#### **END-TO-END**









Salesky et al. (2019)



## <u>See paper or Q&A for more!</u>

Alignment Quality	WER	<b>ASR Supervision</b>
Gold	_	Gold transcript
High	23.2	Salesky et al. (2019)
Med	30.4	Seq2Seq ASR
Low	35.5	Kaldi HMM/GMM

Mapping between phone quality and the ASR models used for alignment generation, with the models' WER on Fisher Spanish test

## Phone Features

## <u>See paper or Q&A for more!</u>

Alignment Quality	WER	<b>ASR Supervision</b>
Gold	—	Gold transcript
High	23.2	Salesky et al. (2019)
Med	30.4	Seq2Seq ASR
Low	35.5	Kaldi HMM/GMM

Mapping between phone quality and the ASR models used for alignment generation, with the models' WER on Fisher Spanish test

2	
0	m
0	m
0	m
а	m

(1)

## Phone Features





**Note: uniqued for clarity** 

# **(1) End-to-End vs Cascaded**

# End-to-End vs Gascaded Models

### **Cascade** BPE targets: +2-4↑ Beam search: +4-8 ↑





(1) **Best end-to-end > Best cascade** (2) Architecture comparisons lacking **③ Low-Resource comparisons lacking (4) Best [academic] cascade from 2014** 

/// Weiss et al. (2017)

Cascade

End-to-End



# End-to-End vs Gascaded Models

### Cascade BPE targets: +2-4↑ Beam search: +4-8 ↑



1) **Best end-to-end > Best cascade (2)** Architecture comparisons lacking **③ Low-Resource comparisons lacking** 





## Tuning models & parameters matters

## Can change relative conclusions when making model comparisons



# 2 Phone Features

#### CASCADE

Recall

#### **END-TO-END**





# **Models with Phone Features**







Salesky et al. (2019)



**Speech Input Challenges** 

**1** Length

Recall

#### **(2)** Variation in frame values

#### **③** Variation in number of frames per phone







# **Models with Phone Features**







# **Results with Phone Features**



Significant performance improvements (1)Improvements of 10.9-22.1 over baseline E2E

**Phone Features** 

Cascade

End-to-End



# **Results with Phone Features**



**1** Significant performance improvements

**2** More data efficient

**///** Phone Features

Cascade

End-to-End

LOW (20hr)

# **Results with Phone Features**



- **1** Significant performance improvements
- **2** More data efficient
  - **3** Benefits of phones remain over SOTA ASR

**When Phone Features** 





LOW (20hr)



# Details: Zoom-In

Model	HIGH	MID	LOW	•
Baseline End-to-End	118hr	40hr	22hr	-
Salesky et al. (2019)	41hr	13hr	10hr	
<b>Baseline Cascade</b>	76hr	19hr	12hr	
Phone Cascade	57hr	39hr	27hr	-
Phone End-to-End	42hr	20hr	13hr	
Hybrid Cascade	47hr	34hr	24hr	

Training Time



(shown relative to best baseline: **<u>Baseline Cascade</u>**)

# Feature Quality

#### **Phone Cascade**



## **Phone End-to-End**

# Conclusion

- Phone features are very effective, at multiple resource settings!
  - Build models with intuitions from phone features
- Performance on high-resource settings # performance on low-resource
  - Test models on multiple resource settings
- Cascades are competitive and often better than current E2E models
  - Compare against strong cascaded baselines