Cross-lingual topic prediction for speech using translations

Sameer Bansal Herman Kamper Adam Lopez Sharon Goldwater



the university of edinburgh **informatics**





Automated speech-to-text

Translation



Information Retrieval Where is the nearest hospital Tap to Edit ② OK, here's what I found:

Current systems



Current systems

English text: Where is the nearest hospital?

استوريا الاستورون الالتقاد والكر

English audio:

Automatic Speech Recognition

downstream task: translation, IR

~100 languages supported by Google Translate ...

Unwritten languages





Aikuma: Bird et al. 2014, LIG-Aikuma: Blachon et al. 2016 Godard et al. 2018

- ~3,000 languages with no writing system
- Traditional ASR based will not work!

Unwritten languages



Efforts to collect speech and translations using mobile apps

Unwritten languages



Build cross-lingual speech-to-text systems (ST)

Why speech input?

"For many Indians, searching by voice rather than text is their first choice." Google







55% households: radio main source of information

Quinn and Hidalgo-Sanchis, 2017





Collect data from public radio conversations



Quinn and Hidalgo-Sanchis, 2017





"Insights about the spread of infectious diseases, small-scale disasters, etc."



Quinn and Hidalgo-Sanchis, 2017



Luganda audio



Topic?



Topic prediction task



("... they have built health centers")

Speech to text system



("... they have built health centers")

Keywords indicate topic information



Availability of ASR!



Can we predict topics using ST?



Can we predict topics using ST?



UN study dataset not available!

Our work: topic prediction for Spanish speech



ST trained in simulated low-resource settings

ST performance in low-resource settings

Spanish-English	BLEU	
160 hours - Weiss et al.	46	

*for comparison text-to-text = 58

Good performance if trained on 100+ hours

ST performance in low-resource settings

Spanish-English	BLEU	
160 hours - Weiss et al.	46	
20 hours - Bansal et al. 2019	19	

**for comparison text-to-text = 58*

Mediocre performance in low-resource settings

ST performance in low-resource settings

Spanish-English	BLEU	
160 hours - Weiss et al.	46	
20 hours - Bansal et al. 2019	19	

*for comparison text-to-text = 58

"Good applications for crummy machine translation" Church & Hovy, 1993

Sample translations

Spanish soy cat olica pero no en realidad casi no voy a laiglesia

English i am catholic but actually i hardly go to church

Sample translations

Spanish soy cat olica pero no en realidad casi no voy a laiglesia

English i am catholic but actually i hardly go to church

20h i'm catholics but reality i don't go to the church

"Crummy" translation

Sample translations

Spanish soy cat olica pero no en realidad casi no voy a laiglesia

English i am catholic but actually i hardly go to church

20h i'm catholics but reality i don't go to the church

topic religion

Keywords can be useful for topic prediction

Our work: topic prediction for Spanish speech



ST trained in simulated low-resource settings

Our work: topic prediction for Spanish speech



Gold topics labels not available!

Spanish audio





Gold topic label?

I like to listen to jazz

Gold translation



Gold topic label?

I like to listen to jazz **Gold translation**

Use gold translations to infer topic labels



Use gold translations to infer topic labels



Spanish audio Gold human translation

I listen to english music I am catholic hello how are you



Topic model





Spanish audio Gold human translation

I listen to english music

I am catholic

hello how are you



Торіс	Terms
small-talk	hello, fine, name
music	dance, listen, music
religion	god, bible, believe
•••	

Training set



Spanish audio Gold human translation

I listen to english music

I am catholic

hello how are you



Торіс	Terms
small-talk	hello, fine, name
music	dance, listen, music
religion	god, bible, believe

Number of topics set to 10



Spanish audio Gold human translation

I listen to english music

I am catholic

hello how are you



Торіс	Terms
small-talk	hello, fine, name
music	dance, listen, music
religion	god, bible, believe

small-talk most frequent

Spanish audio



Topic model

Evaluation set



Evaluation set



Compare predicted and silver topic label



Good prediction



Poor prediction



Evaluate over a 100 hour test set



- ST trained on <= 20 hours of Spanish-English
- Pretrained on English ASR



small-talk topic is the majority class baseline



Poor performance <= 5 hours ST models



10-20h ST models outperform majority baseline



10-20h ST models outperform majority baseline



Takeaways

- Low-resource ST can still be useful for building downstream applications
- Silver evaluation for this preliminary study
 - Future: human evaluation
- Experiments on low-resource/unwritten languages
 - Datasets required
- Keyword spotting

Thanks!

• Check out: "Analyzing ASR pretraining for low-resource speech-to-text translation", Stoian et al.

Backup



Predicted topic label

Silver labels

human translation	Assigned	Silver
hello good afternoon have you ever been in a jury in a trial	juries	intro-misc
i also receive many letters of life insurance from banks	spam	welfare
they tell us we have to talk about marriage	music	family-misc

Speakers were provided discussion prompts

Topic labels



Spanish dataset discussion prompts



Spanish speech to English text



- Telephone speech (unscripted)
- Realistic noise conditions
- Multiple speakers and dialects
- Crowdsourced English text translations

Closer to real-world conditions

Neural ST model



Code available on Github

Cross-lingual applications for low-resource languages

- Sheridan et al., 1997
 - German speech retrieval system using French text queries.
- Projects LORELEI, OpenCLIR



- Query speech/text in a low-resource language using English (or similar high-resource).
- Dredze et al. (2010) and Siu et al. (2014)
 - Unsupervised clustering of speech into topics
- **Our work:** Speech paired with text translations