

Leveraging **multilingual** tracks for **(weakly) supervised speaker identification**

Hervé Bredin

Anindya Roy Nicolas Pécheux Alexandre Allauzen



CNRS // LIMSI // France bredin@limsi.fr // herve.niderb.fr



TV series are very popular...

TV series are very popular...

« Bazinga »

TV series are very popular...

« Bazinga » « Winter is coming »

TV series are very popular...

« Bazinga »

« Winter is coming »

« So say we all »

TV series are very popular...

« Bazinga »

- « Winter is coming »
- « So say we all »

... but TV series are not (easily) searchable.

find all video sequences where **Sheldon (honestly) apologizes to Leonard**

TV series are very popular...

« Bazinga »

- « Winter is coming »
- « So say we all »

... but TV series are not (easily) searchable.

find all video sequences where **Sheldon (honestly) apologizes to Leonard**

speech transcription

I'm really sorry, Leonard.

TV series are very popular...

« Bazinga »

- « Winter is coming »
- « So say we all »

... but TV series are not (easily) searchable.

find all video sequences where **Sheldon (honestly) apologizes to Leonard**

speech transcription speaker identification

I'm really sorry, Leonard.

SHELDON: « I'm really sorry, Leonard. »

00:11.27 → 00:13.83
-...it will not've gone
through both slits.
-Agreed.

 $00:13.99 \rightarrow 00:15.06$ What's your point?

 $00:15.23 \rightarrow 00:18.07$ There's no point. I just think it's a good idea for a T-shirt

00:11.27 → 00:13.83
-...it will not've gone
through both slits.
-Agreed.

 $00:13.99 \rightarrow 00:15.06$ What's your point?

 $00:15.23 \rightarrow 00:18.07$ There's no point. I just think it's a good idea for a T-shirt Sheldon: it will not h

it will not have gone through both slits.

Leonard: Agreed, what's your point?

Sheldon: There's no point, I just think it's a good idea for a tee-shirt.

00:11.27 → 00:13.83
-...it will not've gone
through both slits.
-Agreed.

 $00:13.99 \rightarrow 00:15.06$ What's your point?

00:15.23 → 00:18.07 There's no point. I just think it's a good idea for a T-shirt

Sheldon:

it will not have gone through both slits.

Leonard:

Agreed, what's your point?

Sheldon:

There's no point, I just think it's a good idea for a tee-shirt.

00:11.27 → 00:13.83
-...it will not've gone
through both slits.
-Agreed.

 $00:13.99 \rightarrow 00:15.06$ What's your point?

00:15.23 → 00:18.07 There's no point. I just think it's a good idea for a T-shirt Sheldon:

it will not have gone through both slits.

Leonard:

Agreed, what's your point?

fan transcripts

Sheldon:

There's no point, I just think it's a good idea for a tee-shirt.

00:11.27	\rightarrow	00:13.37	Sheldon
00:13.37	\rightarrow	00:13.83	Leonard
00:13.99	\rightarrow	00:15.06	Leonard
00:15.23	\rightarrow	00:18.07	Sheldon





 when both transcripts and subtitles are available no need to go any further (<10% error)



- when both transcripts and subtitles are available no need to go any further (<10% error)
- what about episodes for which they are not yet available?



- when both transcripts and subtitles are available no need to go any further (<10% error)
- what about episodes for which they are not yet available?

☆ automatic speaker identification

outline

- weakly supervised speaker identification
- bilingual speaker identification
- conclusion

• 👍 (relatively) small set of characters

- 👍 (relatively) small set of characters
- **Spontaneous speech** short speech turns, fast interactions

- 👍 (relatively) small set of characters
- **\$\$ spontaneous speech** short speech turns, fast interactions

- 👍 (relatively) small set of characters
- **\$\frac{P}{P}\$ spontaneous speech** short speech turns, fast interactions

speech activity detection

2-states HMM

- 👍 (relatively) small set of characters
- **\$\$** spontaneous speech short speech turns, fast interactions



- 👍 (relatively) small set of characters
- **\$\$ spontaneous speech** short speech turns, fast interactions



experimental setup



episodes 1 to 6 (2 hours)

leave-one-episode-out cross-validation

manual reference: 6 labels (Sheldon, Leonard, Penny, Howard, Raj, other)

audio tracks & subtitles (English & French) transcripts (English)

detection error rate

 $\frac{\text{miss} + \text{fa}}{\text{speech}}$

identification error rate

 $\frac{\rm miss + fa + \rm confusion}{\rm speech}$

weakly supervised

 2-class (speech vs. non-speech) hidden Markov model based on standard MFCC coefficients

weakly supervised

- 2-class (speech vs. non-speech) hidden Markov model based on standard MFCC coefficients
- the use (free) subtitle timespans for training instead of (costly) manual reference



		Fully	Weakly
Approach	Subtitles	supervised	supervised
DER	19.8%	7.8%	8.1 %
Precision	74.6%	94.8%	91.2%
Recall	95.6%	90.4%	94.1%

DER // Detection error rate

 $\frac{\text{miss} + \text{fa}}{\text{speech}}$

		Fully	Weakly
Approach	Subtitles	supervised	supervised
DER	19.8%	7.8%	8.1 %
Precision	74.6%	94.8%	91.2%
Recall	95.6%	90.4%	94.1%

DER // Detection error rate



• subtitles timespans contain **25% non-speech**

		Fully	Weakly
Approach	Subtitles	supervised	supervised
DER	19.8%	7.8%	8.1 %
Precision	74.6%	94.8%	91.2%
Recall	95.6%	90.4%	94.1%

DER // Detection error rate



- subtitles timespans contain **25% non-speech**
- precision is better for fully supervised
 recall is better fo weakly supervised

		Fully	Weakly
Approach	Subtitles	supervised	supervised
DER	19.8%	7.8 %	$\mathbf{8.1\%}$
Precision	74.6%	94.8%	91.2%
Recall	95.6%	90.4%	94.1%

DER // Detection error rate



- subtitles timespans contain **25% non-speech**
- precision is better for fully supervised
 recall is better fo weakly supervised
- weakly supervised is almost on par with fully supervised

weakly supervised -

• 6-class Gaussian Mixture Models classification 5 main characters vs. all others

weakly supervised

- 6-class Gaussian Mixture Models classification 5 main characters vs. all others
- the use (free) subtitle + transcripts for training instead of (costly) manual reference



Speaker	Speed	ch activity detection		
identification	Reference Fully W		Weakly	
approach	Itelefence	supervised	supervised	
Oracle	10.0%	24.5%	25.4%	
Labeled subtitles	12.8%	27.0%	28.2%	
Fully supervised	18.6%	35.9%	37.9%	
Weakly supervised	18.5%	35.6%	37.8%	

4(

Speaker	Speech activity detection		
identification approach	Reference	Fully supervised	Weakly supervised
Oracle	10.0%	24.5%	25.4%
Labeled subtitles	12.8%	27.0%	28.2%
Fully supervised	18.6%	35.9%	37.9%
Weakly supervised	18.5%	35.6%	37.8%

• secondary characters represent 10% of speech duration

4(

5

- 3(
- 20
- 1(

Speaker	Speech activity detection		
identification approach	Reference	Fully supervised	Weakly supervised
Oracle	10.0%	24.5%	25.4%
Labeled subtitles	12.8%	27.0%	28.2%
Fully supervised	18.6%	35.9%	37.9%
Weakly supervised	18.5%	35.6%	37.8%

- secondary characters represent 10% of speech duration
- we should use subtitles and transcripts when available

- 50 4(
- 30
- 20
- 10
- (

Speaker	Speech activity detection		
identification approach	Reference	Fully supervised	Weakly supervised
Oracle	10.0%	24.5%	25.4%
Labeled subtitles	12.8%	27.0%	28.2%
Fully supervised	18.6%	35.9%	37.9%
Weakly supervised	18.5%	35.6%	37.8%

- secondary characters represent 10% of speech duration
- we should use subtitles and transcripts when available
- weakly supervised is almost on par with fully supervised

- 50 40 30
- 20
- 1(
- (

bilingual speaker identification



bilingual speaker identification



bilingual speaker identification

• some characters may be easier to distinguish using their Frenchdubbed voices than with their original American English voice

- some characters may be easier to distinguish using their Frenchdubbed voices than with their original American English voice
- train two speaker identification systems (and) and apply score fusion to improve the results

 $\rho_{ti} = \alpha \cdot \rho_{ti} + (1 - \alpha) \cdot \rho_{ti}$

- some characters may be easier to distinguish using their Frenchdubbed voices than with their original American English voice
- train two speaker identification systems (and) and apply score fusion to improve the results

$$\rho_{ti} = \alpha \cdot \rho_{ti} + (1 - \alpha) \cdot \rho_{ti}$$









14



results bilingual speaker identification

IER	
37.8 %	
32.8 %	
-13 %	

 bilingual beats monolingual

results bilingual speaker identification

	IER	confusion 5 main characters
	37.8 %	7.1 %
+	32.8 %	2.1 %
	-13 %	-70 %

- bilingual beats monolingual
- bilingual is almost perfect for 5 main characters

results bilingual speaker identification

	IER	confusion 5 main characters
	37.8 %	7.1 %
+	32.8 %	2.1 %
	-13 %	-70 %

- bilingual beats monolingual
- bilingual is almost perfect for 5 main characters

remaining errors
 speech turn segmentation: 15%
 secondary characters: 15%

completely unsupervised speaker identification



completely unsupervised speaker identification



completely unsupervised speaker identification



• audio-visual character recognition

completely unsupervised speaker identification



- audio-visual character recognition
- StoryGraph Tapaswi et al. CVPR 2014
 <u>xkcd.com/657/</u>



"bring your own DVDs" reproducible corpus provides tool to automatically extract video, audio tracks and subtitles







speaker labels episode outlines episode transcript forced alignment

episode outlines episode transcript forced alignment scene segmentation

Π () T N S S S



reproducible corpus tvd.niderb.fr/corpus reproducible research tvd.niderb.fr/research



open-source Python libraries (incl. speech processing modules) github.com/pyannote

METADATV

contact us at metadatv.limsi.fr

Q&A slides

automatic alignment

1. Pre-processing

 if a photon is directed through a plane with two slits in it »



automatic alignment

1. Pre-processing

 if a photon is directed through a plane with two slits in it »



2. Cosine distance matrix



automatic alignment

1. Pre-processing

 if a photon is directed through a plane with two slits in it »



2. Cosine distance matrix



3. Dynamic time warping



influence of **a**



detailed performance



Speaker	Segmentation		
identification	Boforonco	Fully	Weakly
approach	Itelefence	supervised	supervised
Oracle	10.0%	24.5%	25.4%
Labeled subtitles	12.8%	27.0%	28.2%
Fully supervised	18.6%	35.9%	37.9%
Weakly supervised	18.5%	35.6%	37.8%

the segmentation step is the main source of errors





metadatv.limsi.fr

- Why reverse-engineer when we could have access to this information from the start?
- **Problem:** most of the metadata do get lost at one point or another of the production pipeline
- **Solution:** make every actor of the production pipeline realize that their metadata could be very useful down the road