# Self-supervised learning and universal modeling for speech and audio processing

Benchmarking and open-source toolkits

# **Panelists**



#### Shinji Watanabe

- •Associate Professor
- •Carnegie Mellon University



#### Titouan Parcollet •Associate Professor •Avignon Université



#### Shang-Wen Li (Daniel)

 Engineering and Science Manager
 Facebook AI



#### Mirco Ravanelli

- •Post-doc Researcher
- •MILA, Université de Montréal

# Speech Processing Week 2

			Week 2		
Date	2021/8/9	2021/8/10	2021/8/11	2021/8/12	2021/8/13
Weekday	Mon	Mon Tue		Thur	Fri
09:00-10:00		*			1 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
(GMT+8)		Snasker: Minn Mai Chann		Poster Session 2	
10:00-10:30		Title: Pre-training for Natural	Speaker: Philipp Krähenbühl		
(GMT+8)		Language Processing	Title: Computer Vision	Poster	
10:30-11:00	- P		Luchus Ide		Speaker: Been Kim
(GMT+8)		Lecture Info Contrae Link	Containe trive Containe trive		Title: Interpretable machine learning
11:00-12:00	- P.			- A-	Lachara Irin Course Link
(GMT+8)					
12:00-20:00			Brask		
(GMT+8)		-	Li cias	-	
20:00-21:00					
(GMT+8)		4	22		
21:00-22:00 (GMT+8)	Speaker: John Shawe-Taylor Title: An introduction to Statistical Learning Theory and PAC-Bayes Analysis Lecture Into	Speaker: Hung-Yi Lee Title: Deep Learning for Speech Processing Eedate Info Example Link	Speaker: Song Han Title: TinyML and Efficient Deep Learning Lecture Info Course Link	Speakers: Thang Vu, Shang-Wen Li Title: Meta Learning for Human Language Processing Liscum Info] Looma Link	Speaker: Srinivasan Arunachalam Title: Overview of learning quantum states
22:00-23:00 (GMT+8)					

#### How to implement the models?





100,000 hours of paired data

Babies learn their first language with little supervision.

Can AI have the same ability?

# Self-supervised learning / Pre-training <u>Week 2</u>

			Week 2		
Date	2021/8/9	2021/8/10	2021/8/11	2021/8/12	2021/8/13
Weekday	Mon	Mon Tue		Thur	Fri
09:00-10:00		and the second se	22 <b> </b>		· · · · · · · · · · · · · · · · · · ·
(GMT+8)		Complete Mine Mai Chann		Poster Session 2	
10:00-10:30		Title: Dre-training for Natural	Speaker: Philipp Krähenbühl		
(GMT+8)		annuane Procession	Title: Computer Vision	Poster	and the second second
10:30-11:00	- D-	congroupe i nocenning			Speaker: Been Kim
(GMT+8)		Lecture Info	Lacture Into		Title: Interpretable machine learning
11:00-12:00	0.				
(GMT+8)					Lectre inti Course Link
12:00-20:00	1 <sup>22</sup>		P1		
(GMT+8)			Break		22
20:00-21:00					
(GMT+8)			22		4 (
21:00-22:00 (GMT+8)	Speaker: John Shawe-Taylor Title: An introduction to Statistical Learning Theory and PAC-Bayes Analysis Lecture Into Course Link	Speaker: Hung-Yi Lee Title: Deep Learning for Speech Processing Lecture Info Dourse Link	Speaker: Song Han Title: TinyML and Efficient Deep Learning Learning Course Link	Speakers: Thang Vu, Shang-Wen Li Title: Meta Learning for Human Language Processing Lischere Info?	Speaker: Srinivasan Arunachalam Title: Overview of learning quantum states
22:00-23:00 (GMT+8)	0. Sf				



# Outline • Part I

- Benchmarking for Self-supervised Learning
  - SUPERB (speaker: Hung-yi Lee)
  - LeBenchmark (speaker: Titouan Parcollet)
- Open-source toolkits for wide variety of tasks
  - Speech Brain (speaker: Mirco Ravanelli)
  - ESPNet (speaker: Shinji Watanabe)
- Type your questions during the talks.
- Part II:
  - We will discuss your questions.

### Benchmarking for Self-supervised Learning

# To Learn More .....

- SUPERB LeBenchmark

uses speech representations in a wide variety of downstream tasks

Zero Speech --> no linguistic labels available

HEAR Audio (music, sound event, etc.)

Speech

Zero Speech: https://www.zerospeech.com/

HEAR: https://neuralaudio.ai/hear2021-holistic-evaluation-of-audio-representations.html

# Benchmarking for Self-supervised Learning SUPERB

# **SUPERB**

Speech processing Universal PERformance Benchmark



#### **SUPERB:** Speech processing Universal PERformance Benchmark

 Shu-wen Yang<sup>1</sup>, Po-Han Chi<sup>1\*</sup>, Yung-Sung Chuang<sup>1\*</sup>, Cheng-I Jeff Lai<sup>2\*</sup>, Kushal Lakhotia<sup>3\*</sup>, Yist Y. Lin<sup>1\*</sup>, Andy T. Liu<sup>1\*</sup>, Jiatong Shi<sup>4\*</sup>, Xuankai Chang<sup>6</sup>, Guan-Ting Lin<sup>1</sup>,
 Tzu-Hsien Huang<sup>1</sup>, Wei-Cheng Tseng<sup>1</sup>, Ko-tik Lee<sup>1</sup>, Da-Rong Liu<sup>1</sup>, Zili Huang<sup>4</sup>, Shuyan Dong<sup>5†</sup>, Shang-Wen Li<sup>5†</sup>, Shinji Watanabe<sup>6</sup>, Abdelrahman Mohamed<sup>3</sup>, Hung-yi Lee<sup>1</sup>

Will be published at INTERSPEECH 2021

# **Typical Self-supervised Learning for Speech**











# How to use Self-supervised Model - Constrained Track



# How to use Self-supervised Model - Constrained Track

Limited capacity

#### Universal features

- General-purpose knowledge for speech processing, from acoustic signals to semantics
- Easy to add new tasks
- Storage saving

task independent





#### Hand-crafted

features

#### tasks (open data, easy to reproduce)

	PR	KS	IC	SID	ER	ASR	(WER)	QbE	S	SF	ASV	SD
	$\text{PER}\downarrow$	$Acc \uparrow$	$Acc \uparrow$	$Acc\uparrow$	$Acc \uparrow$	w/o↓	w/ LM ↓	MTWV ↑	F1 ↑	$CER \downarrow$	$\text{EER}\downarrow$	DER ↓
FBANK	82.01	8.63	9.10	8.5E-4	35.39	23.18	15.21	0.0058	69.64	52.94	9.56	10.05
PASE+ [16]	58.87	82.54	29.82	37.99	57.86	25.11	16.62	0.0072	62.14	60.17	11.61	8.68
APC [7]	41.98	91.01	74.69	60.42	59.33	21.28	14.74	0.0310	70.46	50.89	8.56	10.53
VQ-APC [32]	41.08	91.11	74.48	60.15	59.66	21.20	15.21	0.0251	68.53	52.91	8.72	10.45
NPC [33]	43.81	88.96	69.44	55.92	59.08	20.20	13.91	0.0246	72.79	48.44	9.4	9.34
Mockingjay [8]	70.19	83.67	34.33	32.29	50.28	22.82	15.48	6.6E-04	61.59	58.89	11.66	10.54
TERA [9]	49.17	89.48	58.42	57.57	56.27	18.17	12.16	0.0013	67.50	54.17	15.89	9.96
modified CPC [34]	42.54	91.88	64.09	39.63	60.96	20.18	13.53	0.0326	71.19	49.91	12.86	10.38
wav2vec [12]	31.58	95.59	84.92	56.56	59.79	15.86	11.00	0.0485	76.37	43.71	7.99	9.9
vq-wav2vec [13]	33.48	93.38	85.68	38.80	58.24	17.71	12.80	0.0410	77.68	41.54	10.38	9.93
wav2vec 2.0 Base [14]	5.74	96.23	92.35	75.18	63.43	6.43	4.79	0.0233	88.30	24.77	6.02	6.08
wav2vec 2.0 Large [14]	4.75	96.66	95.28	86.14	65.64	3.75	3.10	0.0489	87.11	27.31	5.65	5.62
HuBERT Base [35]	5.41	96.30	98.34	81.42	64.92	6.42	4.79	0.0736	88.53	25.20	5.11	5.88
HuBERT Large [35]	3.53	95.29	98.76	90.33	67.62	3.62	2.94	0.0353	89.81	21.76	5.98	5.75

Self-supervised Speech Model

	PR	KS	IC	SID	ER	ASR	(WER)	QbE		SF	ASV	SD
	$PER \downarrow$	$Acc \uparrow$	$Acc\uparrow$	$Acc \uparrow$	$Acc \uparrow$	w/o↓	w/LM $\downarrow$	MTWV ↑	F1 ↑	$CER \downarrow$	$\text{EER}\downarrow$	$\text{DER}\downarrow$
FBANK	82.01	8.63	9.10	8.5E-4	35.39	23.18	15.21	0.0058	69.64	52.94	9.56	10.05
PASE+ [16]	58.87	82.54	29.82	37.99	57.86	25.11	<u>16.62</u>	0.0072	62.14	60.17	11.61	8.68
APC [7]	41.98	91.01	74.69	60.42	59.33	21.28	14.74	0.0310	70.46	50.89	8.56	10.53
VQ-APC [32]	41.08	91.11	74.48	60.15	59.66	21.20	15.21	0.0251	68.53	52.91	8.72	10.45
NPC [33]	43.81	88.96	69.44	55.92	59.08	20.20	13.91	0.0246	72.79	48.44	9.4	9.34
Mockingjay [8]	70.19	83.67	34.33	32.29	50.28	22.82	15.48	6.6E-04	61.59	58.89	11.66	10.54
TERA [9]	49.17	89.48	58.42	57.57	56.27	18.17	12.16	0.0013	67.50	54.17	15.89	9.96
modified CPC [34]	42.54	91.88	64.09	39.63	60.96	20.18	13.53	0.0326	71.19	49.91	12.86	10.38
wav2vec [12]	31.58	95.59	84.92	56.56	59.79	15.86	11.00	0.0485	76.37	43.71	7.99	9.9
vq-wav2vec [13]	33.48	93.38	85.68	38.80	58.24	17.71	12.80	0.0410	77.68	41.54	10.38	9.93
wav2vec 2.0 Base [14]	5.74	96.23	92.35	75.18	63.43	6.43	4.79	0.0233	88.30	24.77	6.02	6.08
wav2vec 2.0 Large [14]	4.75	96.66	95.28	86.14	65.64	3.75	3.10	0.0489	87.11	27.31	5.65	5.62
HuBERT Base [35]	5.41	96.30	98.34	81.42	64.92	6.42	4.79	0.0736	88.53	25.20	5.11	5.88
HuBERT Large [35]	3.53	95.29	98.76	90.33	67.62	3.62	2.94	0.0353	89.81	21.76	5.98	5.75

Self-supervised learning outperforms fbank in most cases.

	PR	KS	IC	SID	ER	ASR (WER)		QbE SF		SF	ASV	SD
	$\text{PER}\downarrow$	Acc ↑	Acc↑	$Acc \uparrow$	$Acc \uparrow$	w/o↓	w/LM $\downarrow$	MTWV ↑	F1 ↑	CER↓	$EER \downarrow$	DER ↓
FBANK	82.01	8.63	9.10	8.5E-4	35.39	23.18	15.21	0.0058	69.64	52.94	9.56	10.05
PASE+ [16]	58.87	82.54	29.82	37.99	57.86	25.11	16.62	0.0072	62.14	60.17	11.61	8.68
APC [7]	41.98	91.01	74.69	60.42	59.33	21.28	14.74	0.0310	70.46	50.89	8.56	10.53
VQ-APC [32]	41.08	91.11	74.48	60.15	59.66	21.20	15.21	0.0251	68.53	52.91	8.72	10.45
NPC [33]	43.81	88.96	69.44	55.92	59.08	20.20	13.91	0.0246	72.79	48.44	9.4	9.34
Mockingjay [8]	70.19	83.67	34.33	32.29	50.28	22.82	15.48	6.6E-04	61.59	58.89	11.66	10.54
TERA [9]	49.17	89.48	58.42	57.57	56.27	18.17	12.16	0.0013	67.50	54.17	15.89	9.96
modified CPC [34]	42.54	91.88	64.09	39.63	60.96	20.18	13.53	0.0326	71.19	49.91	12.86	10.38
wav2vec [12]	31.58	95.59	84.92	56.56	59.79	15.86	11.00	0.0485	76.37	43.71	7.99	9.9
vq-wav2vec [13]	33.48	93.38	85.68	38.80	58.24	17.71	12.80	0.0410	77.68	41.54	10.38	9.93
wav2vec 2.0 Base [14]	5.74	96.23	92.35	75.18	63.43	6.43	4.79	0.0233	88.30	24.77	6.02	6.08
wav2vec 2.0 Large [14]	4.75	96.66	95.28	86.14	65.64	3.75	3.10	0.0489	87.11	27.31	5.65	5.62
HuBERT Base [35]	5.41	96.30	98.34	81.42	64.92	6.42	4.79	0.0736	88.53	25.20	5.11	5.88
HuBERT Large [35]	3.53	95.29	98.76	90.33	67.62	3.62	2.94	0.0353	89.81	21.76	5.98	5.75

- Some self-supervised models outperform fbank in all tasks.
- HuBERT and wav2vec 2.0 outperform other models

	PR	KS	IC	SID	ER	ASR (WER)		QbE SF		ASV	SD	
	PER↓	Acc↑	Acc↑	$Acc\uparrow$	$Acc \uparrow$	w/o↓	w/LM↓	MTWV ↑	F1 ↑	$CER \downarrow$	EER↓	$DER \downarrow$
FBANK	82.01	8.63	9.10	8.5E-4	35.39	23.18	15.21	0.0058	69.64	52.94	9.56	10.05
PASE+ [16]	58.87	82.54	29.82	37.99	57.86	25.11	16.62	0.0072	62.14	60.17	11.61	8.68
APC [7]	41.98	91.01	74.69	60.42	59.33	21.28	14.74	0.0310	70.46	50.89	8.56	10.53
VQ-APC [32]	41.08	91.11	74.48	60.15	59.66	21.20	15.21	0.0251	68.53	52.91	8.72	10.45
NPC [33]	43.81	88.96	69.44	55.92	59.08	20.20	13.91	0.0246	72.79	48.44	9.4	9.34
Mockingjay [8]	70.19	83.67	34.33	32.29	50.28	22.82	15.48	6.6E-04	61.59	58.89	11.66	10.54
TERA [9]	49.17	89.48	58.42	57.57	56.27	18.17	12.16	0.0013	67.50	54.17	15.89	9.96
modified CPC [34]	42.54	91.88	64.09	39.63	60.96	20.18	13.53	0.0326	71.19	49.91	12.86	10.38
wav2vec [12]	31.58	95.59	84.92	56.56	59.79	15.86	11.00	0.0485	76.37	43.71	7.99	9.9
vq-wav2vec [13]	33.48	93.38	85.68	38.80	58.24	17.71	12.80	0.0410	77.68	41.54	10.38	9.93
wav2vec 2.0 Base [14]	5.74	96.23	92.35	75.18	63.43	6.43	4.79	0.0233	88.30	24.77	6.02	6.08
wav2vec 2.0 Large [14]	4.75	96.66	95.28	86.14	65.64	3.75	3.10	0.0489	87.11	27.31	5.65	5.62
HuBERT Base [35]	5.41	96.30	98.34	81.42	64.92	6.42	4.79	0.0736	88.53	25.20	5.11	5.88
HuBERT Large [35]	3.53	95.29	98.76	90.33	67.62	3.62	2.94	0.0353	89.81	21.76	5.98	5.75

Leaderboard: https://superbbenchmark.org/

Welcome to submit! ③

https://github.com/s3prl/s3prl

S3PRL toolkit: Self-Supervised Speech Pre-training and Representation Learning

- Upstream Pre-training
- Upstream Hub
- Downstream fine-tuning
- SUPERB Challenge

Speech Brain ESPNet

https://github.com/s3prl/s3prl

https://github.com/s3prl/s3prl

## S3PRL toolkit: Self-Supervised Speech Pre-training and Representation Learning



Shu-wen Yang (Leo) Andy T. Liu

#### Link to recording: https://youtu.be/PkMFnS6cjAc

# **Next Step**

- Adding new tasks
  - Speech Translation
  - Spoken Question Answering
  - Speech Enhancement
  - Speaker Separation
  - Voice Conversion
- Other ways to use self-supervised models

# **Next Step**

- Near future: Challenge
  - Private data for the tasks
  - (Perhaps) You only submit self-supervised model.
     We train your model on downstream tasks by our scripts for fair comparison.

#### Stayed tuned: https://superbbenchmark.org/



# **Call for Paper**

- IEEE JSTSP Special Issue on **Self-Supervised Learning for Speech and** Audio Processing
- Deadline: December 30, 2021
- Link

https://signalprocessingsociety.org/blog/ieee-jstsp-special-issue-self-supe

rvised-learning-speech-and-audio-processing



# Benchmarking for Self-supervised Learning LeBenchmark

# A bit of context – LeBenchmark



#### LeBenchmark: A Reproducible Framework for Assessing Self-Supervised Representation Learning from Speech

Solène Evain<sup>1,\*</sup>, Ha Nguyen<sup>1,2,\*</sup>, Hang Le<sup>1,\*</sup>, Marcely Zanon Boito<sup>1,\*</sup>, Salima Mdhaffar<sup>2,\*</sup>, Sina Alisamir<sup>1,3</sup>, Ziyi Tong<sup>1</sup>, Natalia Tomashenko<sup>2</sup>, Marco Dinarelli<sup>1,\*</sup>, Titouan Parcollet<sup>2,\*</sup>, Alexandre Allauzen<sup>4</sup>, Yannick Estève<sup>2</sup>, Benjamin Lecouteux<sup>1</sup>, François Portet<sup>1</sup>, Solange Rossato<sup>1</sup>, Fabien Ringeval<sup>1</sup>, Didier Schwab<sup>1</sup> and Laurent Besacier<sup>1,5</sup>

Published at INTERSPEECH 2021

# A bit of context – LeBenchmark

... with one thing in common:



# SSL for speech looks great but ...

Lack of scientific standardisation

# SSL for speech looks great but ...

Lack of scientific standardisation + mostly evaluated on English.
## SSL for speech looks great but ...

Lack of scientific standardisation + mostly evaluated on English.

=

How do we even know if a new method works better in a multilingual setup?

## SSL for speech looks great but ...

Lack of scientific standardisation + mostly evaluated on English.

=

How do we even know if a new method works better in a multilingual setup?

LeBenchmark = std(dataset && evaluation) + French.

**Remember:** we want to evaluate new SSL models for Speech.

**Remember:** we want to evaluate new SSL models for Speech.

We should not waste time in investigating the impact of different data.

**Remember:** we want to evaluate new SSL models for Speech.

Offer a large-enough and heterogeneous-enough dataset.

**Remember:** we want to evaluate new SSL models for Speech.

Offer a large-enough and heterogeneous-enough dataset.

SSL needs **a lot** of data (e.g. VoxPopuli > 100K Hours)

**Remember:** we want to evaluate new SSL models for Speech.

Offer a large-enough and **heterogeneous-enough** dataset.

Gender balance, influence of the environment, type of speech (spontaneous, broadcasted, read ...) ...

**Problem:** getting the data. (everything that is not english is hard to obtain)

#### **Problem:** getting the data. (everything that is not english is hard to obtain)

#### Gather well-known French datasets and **document them**.



LeBenchmark datasets = 1K hours or 3K hours or 7K hours or 10K\* hours

\*not available yet (Submitted to NeurIPS Datasets and Benchmarks track)

$Corpus_{License}$	# Utterances	Duration	# Speakers	Mean Utt. Duration	Speech type
	A/2.055	Small dataset	180		
MLS French <sub>CCBY4.0</sub> (75)	263,055 124,590 / 138,465 / -	<b>1,096:43</b> 520:13 / 576:29 / –	178 80/98/-	15s / 15s / -	Read
		Malium dataat			
A friend A constant	16 402	Niedium dataset	121	4	
Franch (2)	10,402	18:50	49/26/149	45	Read
FiencifApache2.0 (5)	3/3/102/13,92/	-7-718.50	46/30/146	-/-/-	Astad
Att-Hack <sub>CCBYNCND</sub> (44)	<b>30,339</b> 16 564 / 10 775 /	27:02	20	2.75	Acted
	10,304/19,7/37-	12:07/14:547-	9/11/-	2.0 \$ / 2. / \$ / -	Emotional
CaFE <sub>CCNC</sub> (35)	930	1:09	12	4.45	Acted Emotional
CEDD2000 *	408/408/-	0:3270:387-	0/0/-	4.28/4./8/-	Emotional
(15) (of p)	9055	10:20	49	0S	Spontaneous
(13), (cip)	10071,18478,505	0:147 1:307 14:10	2/4/45	100	
$ESLO2_{NC}$ (29), (esl)	62,918	34:12	190	1.95	Spontaneous
	50,4407 52,1477 551	1/:06/16:3//0:09	08/120/2	28/1.98/1./8	D. J.
$EPAC^{**}NC$ (30)	623,250	1,020:02	1,935	98	Radio
0005 141	403,8397137,3917-	1,240:107 383:327 -	-/-/-	=/=/=	Broadcasts
$GEMEP_{NC}$ (16)	1,236	0:50	10	2.5 \$	Acted
	010/020/-	0:2470:267-	5/5/-	2.4 \$ / 2.5 \$ / -	Emotional
MPF (32), (4)	19,527	19:06	114	3.5 \$	Spontaneous
	5,326/4,649/9,552	5:2674:3679:03	36/29/49	3./s/3.6s/3.4s	
PORIMEDIA <sub>NC</sub>	19,627	38:59	193	7.15	Acted telephone
(French) (45)	9,294 / 10,333 / -	19:08 / 19:50 / -	84 / 109 / -	7.4 \$ 7 6.9 \$ 7 -	dialogue
TCOF	58,722	53:59	749	3.3 \$	Spontaneous
(Adults) (7)	10,377714,763733,582	9:33712:39731:46	119 / 162 / 468	3.3 \$ / 3.1 \$ / 3.4 \$	
Medium dataset total	1,111,865	2,933:24	-	-	-
	664,0737379,897767,895	1,824:5371,034:15774:10			
		Large dataset			
1.00 (10)	8,219	19:40	Unk	8.6 s	<b>n</b> 1
MaSS (13)	8,219/-/-	19:40 / - / -	-/-/-	8.6 s / - / -	Read
$\mathrm{NCCFr}_{NC}$ (72)	29,421	26:35	46	3s	Spontaneous
	14,570 / 13,922 / 929	12:44 / 12:59 / 00:50	24/21/1	3s/3s/3s	
Voxpopuli <sub>CCBYNC4.0</sub> (78) Unlabeled***	568,338	4,532:17	Unk	29 s	
	-1-1-	-/-/4,532:17	-/-/-	-/-/-	Professional speech
Voxpopuli <sub>CCBYNC4.0</sub> (78) transcribed***	76.281	211:57	327	10 s	Professional speech
	-/-/-	-/-/211:57	-/-/-	-/-/-	
-	1.814.242	7,739:22			
Large dataset total***	682,322 / 388,217 / 99,084	1,853:02 / 1,041:07 / 4,845:07	-	-	-

## Influence of the data availability (e.g. low-resources languages).

Are 1000H enough to train a SSL system?

Is it worth adding thousands of hours of speech and compute?

Does it depend on the downstream task?

$Corpus_{License}$	# Utterances	Duration	# Speakers	Mean Utt. Duration	Speech type
		Small datasat			
	263 055	1 096.43	178	15 s	
MLS French <sub>CCBY4.0</sub> (75)	124,590 / 138,465 / -	520:13 / 576:29 / -	80/98/-	15 s / 15 s / -	Read
		Medium dataset			· · · · · · · · · · · · · · · · · · ·
African Accented	16,402	18:56	232	4s	
French Anache2 0 (3)	373 / 102 / 15,927	-/-/18:56	48 / 36 / 148	-/-/-	Read
	36,339	27:02	20	2.7 s	Acted
Att-Hack <sub>CCBYNCND</sub> (44)	16,564 / 19,775 / -	12:07 / 14:54 / -	9/11/-	2.6 s / 2.7 s / -	Emotional
C-FE (25)	936	1:09	12	4.4 s	Acted
$CaFE_{CCNC}$ (33)	468 / 468 /	0:32 / 0:36 / -	6/6/-	4.2 s / 4.7 s / -	Emotional
CFPP2000 <sub>CCBYNCSA</sub> *	9853	16:26	49	6 s	Spontonaous
(15), (cfp)	166 / 1,184 / 8,503	0:14 / 1:56 / 14:16	2/4/43	5s/5s/6s	spontaneous
ESL (02 (20) (asl)	62,918	34:12	190	1.9 s	Spontaneous
$E3L02_{NC}$ (29), (est)	30,440 / 32,147 / 331	17:06 / 16:57 / 0:09	68 / 120 / 2	2 s / 1.9 s / 1.7 s	spontaneous
$FPAC^{**} = (30)$	623,250	1,626:02	1,935	9 s	Radio
EFAC $NC(50)$	465,859 / 157,391 / -	1,240:10 / 385:52 / -	-/-/-	-/-/-	Broadcasts
GEMEP v a (16)	1,236	0:50	10	2.5 s	Acted
GEWIER NC (10)	616 / 620 / -	0:24 / 0:26 / -	5/5/-	2.4 s / 2.5 s / –	Emotional
MPF (32) (4)	19,527	19:06	114	3.5 s	Spontaneous
WI I (52), (4)	5,326 / 4,649 / 9,552	5:26 / 4:36 / 9:03	36 / 29 / 49	3.7 s / 3.6 s / 3.4 s	spontaneous
PORTMEDIANC	19,627	38:59	193	7.1 s	Acted telephone
(French) (45)	9,294 / 10,333 / -	19:08 / 19:50 / -	84 / 109 / -	7.4 s / 6.9 s / –	dialogue
TCOF	58,722	53:59	749	3.3 s	Spontaneous
(Adults) (7)	10,377 / 14,763 / 33,582	9:33 / 12:39 / 31:46	119 / 162 / 468	3.3 s / 3.1 s / 3.4 s	opontaneous
Medium dataset total	1,111,865	2,933:24			
	664,073/379,897/67,895	1,824:53 / 1,034:15 / 74:10			
		Large dataset			
M 66 (12)	8,219	19:40	Unk	8.6 s	Devi
Mass (13)	8,219/-/-	19:40 / - / -	-/-/-	8.6 s / – / –	Read
$\mathrm{NCCFr}_{NC}$ (72)	29,421	26:35	46	38	Spontaneous
	14,570 / 13,922 / 929	12:44 / 12:59 / 00:50	24/21/1	3s/3s/3s	
Voxpopuli <sub>CCBYNC4.0</sub> (78) Unlabeled***	568,338	4,532:17	Unk	29 s	Professional speech
	-/-/-	-/-/4,532:17	-/-/-	-/-/-	
Voxpopuli <sub>CCBYNC4.0</sub> (78) transcribed***	76.281	211:57	327	10 s	Professional speech
	-/-/-	-/-/211:57	-/-/-	-/-/-	
Large dataset total****	1,814,242	7,739:22			_
	682,322 / 388,217 / 99,084	1,853:02 / 1,041:07 / 4,845:07	-	-	-

#### Influence of the data type (e.g. on-the-wild).

Does spontaneous speech help with SSL?

Are noisy data helping with robustness?

Can we just train a big model with read speech only?

$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	$Corpus_{License}$	# Utterances	Duration	# Speakers	Mean Utt. Duration	Speech type		
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$			Small dataset					
$\begin{array}{c c c c c c c c c c c c c c c c c c c $		263.055	1.096:43	178	158			
$\begin{tabular}{ c c c c c c c c c c c c c c c c c c c$	MLS French <sub>CCBY4.0</sub> (75)	124,590 / 138,465 / -	520:13 / 576:29 / -	80/98/-	15 s / 15 s / -	Read		
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	Nedium dataset							
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	African Accented	16.402	18:56	232	4s			
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	French Anache2 0 (3)	373 / 102 / 15,927	-/-/18:56	48 / 36 / 148	-/-/-	Read		
$\begin{array}{c c c c c c c c c c c c c c c c c c c $		36,339	27:02	20	2.7 s	Acted		
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	Att-Hack <sub>CCBYNCND</sub> (44)	16,564 / 19,775 / -	12:07 / 14:54 / -	9/11/-	2.6 s / 2.7 s / -	Emotional		
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	C-FE (25)	936	1:09	12	4.4 s	Acted		
$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$CaFE_{CCNC}$ (33)	468 / 468 /	0:32 / 0:36 / -	6/6/-	4.2 s / 4.7 s / -	Emotional		
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	CFPP2000 <sub>CCBYNCSA</sub> *	9853	16:26	49	6 s	Caratanaana		
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	(15), (cfp)	166 / 1,184 / 8,503	0:14 / 1:56 / 14:16	2/4/43	5s/5s/6s	Spontaneous		
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	ESL 02 (20) (asl)	62,918	34:12	190	1.9 s	Casatanasana		
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	$ESLO2_{NC}$ (29), (esl)	30,440 / 32,147 / 331	17:06 / 16:57 / 0:09	68 / 120 / 2	2 s / 1.9 s / 1.7 s	Spontaneous		
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	EDA C** (20)	623,250	1,626:02	1,935	9 s	Radio		
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	$EFAC^{++}NC$ (30)	465,859 / 157,391 / -	1,240:10 / 385:52 / -	-/-/-	-/-/-	Broadcasts		
$\begin{tabular}{l l l l l l l l l l l l l l l l l l l $	CEMED	1,236	0:50	10	2.5 s	Acted		
$\begin{tabular}{ c c c c c c c c c c c c c c c c c c c$	$GEMEP_{NC}$ (10)	616 / 620 / -	0:24 / 0:26 / -	5/5/-	2.4 s / 2.5 s / -	Emotional		
$\begin{tabular}{ c c c c c c c c c c c c c c c c c c c$	MDE (22) (4)	19,527	19:06	114	3.5 s	Spontaneous		
$\begin{tabular}{ c c c c c c c c c c c c c c c c c c c$	MPF (32), (4)	5,326 / 4,649 / 9,552	5:26 / 4:36 / 9:03	36 / 29 / 49	3.7 s / 3.6 s / 3.4 s			
$\begin{tabular}{ c c c c c c c c c c c c c c c c c c c$	PORTMEDIANC	19,627	38:59	193	7.1 s	Acted telephone		
$\begin{tabular}{ c c c c c c c c c c c c c c c c c c c$	(French) (45)	9,294 / 10,333 / -	19:08 / 19:50 / -	84 / 109 / -	7.4 s / 6.9 s / -	dialogue		
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	TCOF	58,722	53:59	749	3.3 s	C		
$\begin{tabular}{ c c c c c c c c c c c c c c c c c c c$	(Adults) (7)	10,377 / 14,763 / 33,582	9:33 / 12:39 / 31:46	119 / 162 / 468	3.3 s / 3.1 s / 3.4 s	Spontaneous		
Medium dataset total     664,073 / 379,897 / 67,895     1,824:53 / 1,034:15 / 74:10     Image: Constraint of the state of the	Madimu data at tatal	1,111,865	2,933:24					
Large dataset       MaSS (13)     8,219     19:40     Unk     8.6s     Read       MCCFr <sub>NC</sub> (72)     29,421     26:35     46     3s     Spontaneous       VOXpopuli_CCBYNC4.0 (78)     14,570 / 13,922 / 929     12:44 / 12:59 / 00:50     24 / 21 / 1     3s / 3s / 3s     Spontaneous       Voxpopuli_CCBYNC4.0 (78)     568,338     4,532:17     Unk     29 s     Professional speech       Voxpopuli_CCBYNC4.0 (78)     76.281     211:57     327     10 s     Professional speech       Iranscribed***     -/-/-     -/-/-     -/-/-     -/-/-     Professional speech       Large dataset total****     682,322 / 388,217 / 99,084     1,853:02 / 1,041:07 / 4,845:07     -/-/-     -/-/-	Medium dataset total	664,073 / 379,897 / 67,895	1,824:53 / 1,034:15 / 74:10					
Large dataset       MaSS (13)     8,219     19:40     Unk     8.6 s     Read       MaSS (13)     8,219 / - / -     19:40 / - / -     - / - / -     8.6 s / - / -     Read       MCCFr <sub>NC</sub> (72)     29,421     26:35     46     3 s     Spontaneous       Voxpopuli <sub>CCBYNC4.0</sub> (78)     14,570 / 13,922 / 929     12:44 / 12:59 / 00:50     24 / 21 / 1     3 s / 3 s     Spontaneous       Voxpopuli <sub>CCBYNC4.0</sub> (78)     568,338     4,532:17     Unk     29 s     Professional speech       Voxpopuli <sub>CCBYNC4.0</sub> (78)     76.281     211:57     327     10 s     Professional speech       transcribed***     -/-/-     -/-/-     -/-/-     -/-/-     -/-/-     Image: Signal speech     Professional speech       Large dataset total****     1,814,242     7,739:22     -     -     -     -								
$\begin{tabular}{ c c c c c c c c c c c c c c c c c c c$		8 210	Large dataset	T-L-	0 ( -			
8,2197-7     19:407-7     -7-7     8,087-7-7       NCCF <sub>NC</sub> (72)     29,421     26:35     46     3s       Voxpopuli <sub>CCBYNC4.0</sub> (78)     14,570/13,922/929     12:44/12:59/00:50     24/21/1     3s/3s/3s     Spontaneous       Voxpopuli <sub>CCBYNC4.0</sub> (78)     568,338     4,532:17     Unk     29s     Professional speech       Vnabeled***     -/-/-     -/-/-     -/-/-     Professional speech     Professional speech       Voxpopuli <sub>CCBYNC4.0</sub> (78)     76.281     211:57     327     10s     Professional speech       transcribed***     -/-/-     -/-/211:57     -/-/-     -/-/-     Professional speech       Large dataset total****     682,322/388,217/99,084     1,853:02/1,041:07/4,845:07     -     -	MaSS (13)	8,219	19:40	Unk	8.0 S	Read		
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $		8,2197-7-	19:407-7-	-/-/-	8.0 \$ / - / -			
Voxpopuli <sub>CCBYNC4.0</sub> (78)     568,338     4,532:17     Unk     29 s     Professional speech       Unlabeled***     -/-/-     -/-/4,532:17     -/-/-     -/-/-     Professional speech       Voxpopuli <sub>CCBYNC4.0</sub> (78)     76.281     211:57     327     10 s       transcribed***     -/-/-     -/-/-     -/-/-     Professional speech       Large dataset total****     1,814,242     7,739:22     -       1,853:02 / 1,041:07 / 4,845:07     -     -	$\mathrm{NCCFr}_{NC}$ (72)	29,421 14.570 / 13.922 / 929	20:35 12:44 / 12:59 / 00:50	40 24/21/1	3s/3s/3s	Spontaneous		
Unlabeled***     -/-/-     -/-/4.532:17     -/-/-     -/-/-     Professional speech       Voxpopuli_CCEYNC4.0 (78)     76.281     211:57     327     10s     Professional speech       transcribed***     -/-/-     -/-/211:57     -/-/-     -/-/-     Professional speech       Large dataset total****     181,242     7,39:22     -/-/-     -/-/-     -/-/-	Voxpopulic CRXNC4 o (78)	568.338	4.532:17	Unk	29.8	Professional speech		
Voxpopuli <sub>CCBYNC4.0</sub> (78)     76.281     211:57     327     10s     Professional speech       transcribed***     -/-/-     -/-/-     -/-/-     -/-/-     Professional speech       Large dataset total****     1814.242     7,739:22     -/-/-     -/-/-	Unlabeled***	-/-/-	-/-/4,532:17	-/-/-	-/-/-			
transcribed***     -/-/     -/-/211:57     -/-/     Professional speech       Large dataset total****     1,814,242     7,739:22     -/-/-     -/-/-     -/-/-       1,853:02/1,041:07/4,845:07     - </td <td rowspan="2">Voxpopuli<sub>CCBYNC4.0</sub> (78) transcribed***</td> <td>76.281</td> <td>211:57</td> <td>327</td> <td>10 s</td> <td rowspan="2">Professional speech</td>	Voxpopuli <sub>CCBYNC4.0</sub> (78) transcribed***	76.281	211:57	327	10 s	Professional speech		
Large dataset total****     1,814,242     7,739:22       682,322 / 388,217 / 99,084     1,853:02 / 1,041:07 / 4,845:07     -		-/-/-	-/-/211:57	-/-/-	-/-/-			
Large dataset total**** 682,322 / 388,217 / 99,084 1,853:02 / 1,041:07 / 4,845:07	Large dataset total****	1.814.242	7.739:22					
		682,322 / 388,217 / 99,084	1,853:02 / 1,041:07 / 4,845:07	-	-	-		

#### Research remains slowed down by the distribution scheme of the data ...

Give the community a set of baselines to compare with.

Give the community a set of baselines to compare with.

Research about SSL is not only about finding the best SSL model.

Give the community a set of baselines to compare with.

Research about SSL is not only about finding the best SSL model.

Not everyone can train a SSL model on 7K hours (64 Tesla V100/14 days).



wav2vec 2.0\*

Model Type: Transformer SSL Type: Contrastive Learning

Training these models **really** is challenging ...

Release the models to the community!

wav2vec-FR-7K-Large wav2vec-FR-3K-Large wav2vec-FR-1K-Large

wav2vec-FR-7K-Base wav2vec-FR-3K-Base wav2vec-FR-2.6K-Base wav2vec-FR-1K-Base



Can be **easily** used with well-know toolkits: SpeechBrain, ESPnet, S3PRL, FairSeq ...

Release the models to the community!

wav2vec-FR-7K-Large wav2vec-FR-3K-Large wav2vec-FR-1K-Large

wav2vec-FR-7K-Base wav2vec-FR-3K-Base wav2vec-FR-2.6K-Base wav2vec-FR-1K-Base



wav2vec-FR-2.6K-Base does not contain spontaneous speech.

**Problem:** designing a benchmark close to the real world is **extremely hard**. (almost impossible)

**Problem:** designing a benchmark close to the real world is **extremely hard**. (almost impossible)

diversity.

Diversity of problems: classification (AER), sequence labelling (SLU) and conditional natural language generation (ASR, AST).

**Diversity of problems:** classification (AER), sequence labelling (SLU) and conditional natural language generation (ASR, AST).

**Diversity of information extracted:** transcript (ASR), semantics (SLU), translation (AST) and paralinguistics (AER).

Diversity of problems: classification (AER), sequence labelling (SLU) and conditional natural language generation (ASR, AST).
Diversity of information extracted: transcript (ASR), semantics (SLU), translation (AST) and paralinguistics (AER).
Diversity of annotated resources available for downstream tasks: large (ASR), medium (SLU, AST) or small (AER).

**Diversity of problems:** classification (AER), sequence labelling (SLU) and conditional natural language generation (ASR, AST).

**Diversity of information extracted:** transcript (ASR), semantics (SLU), translation (AST) and paralinguistics (AER).

Diversity of annotated resources available for downstream tasks: large (ASR), medium (SLU, AST) or small (AER).

#### Provide French standardised downstream baselines with training scripts for:

Automatic Emotion Recognition, Spoken Language Understanding, Automatic Speech Translation, Automatic Speech Recognition.

Too much tables and numbers = let's jump to the outcomes.

Too much **tables** and **numbers =** let's jump to the outcomes.

Good for the community!

Automatic Emotion Recognition (RECOLA — 3.8h / AlloSat 37h)

Automatic Emotion Recognition (RECOLA — 3.8h / AlloSat 37h)

Multilingual wav2vec 2.0 (XLSR) trained on way more data is not as good as our French models.

Smaller wav2vec 2.0 (base) lead to better performance.

SSL features are better than traditional acoustic ones (FBANKs).

Spoken Language Understanding (MEDIA — 17h)

Multilingual wav2vec 2.0 (XLSR) trained on way more data is not as good as our French models.

More hours of pretraining do not improve the performance.

Larger wav2vec 2.0 (large) lead to better performance.

SSL features are better than traditional acoustic ones (FBANKs).

Automatic Speech Translation (French to {English, Portugese, Spanish} — TEDx and CoVoST2 — [25h, 180h])

Multilingual wav2vec 2.0 (XLSR) trained on way more data is not as good as our French models.

Gains are reducing with the increase of supervised data.

Larger wav2vec 2.0 (large) lead to better performance.

SSL features are better than traditional acoustic ones (FBANKs).

Automatic Speech Recognition (CommonVoice — 477h / ETAPE — 36h)

Multilingual wav2vec 2.0 (XLSR) trained on way more data is not as good as our French models.

Gains are reducing with the increase of supervised data.

Larger wav2vec 2.0 (large) lead to better performance.

SSL features are better than traditional acoustic ones (FBANKs).

#### **LeBenchmark**

#### **Outcomes:**

SSL models can not be deployed and used from intuitions only.

We still do not know exactly what will work or not.

SSL models behaviors vary a lot.

Benchmarks are needed to frame the development of these technologies!

#### open-source toolkits

# **Open-Source Toolkits for Speech Processing**

• **Open-source** toolkits have played a critical role in the development of speech processing technology:



• With the raise of **deep learning**, some general-purpose libraries have been developed:



# **Open-Source Toolkits for Speech Processing**

- Thanks to these general framework, more flexible **python-based** speech processing toolkits have quickly appeared.
- Some of them are **task-specific**:


# **Open-Source Toolkits for Speech Processing**

• Some others support multiple speech tasks:









# open-source toolkits Speech Brain

# What is SpeechBrain?

• **SpeechBrain** is an *open-source* and *all-in-one* speech toolkit based on PyTorch.

• **Goal:** speed up **research** and **development** of speech and audio processing techniques.



#### Key features:

- Flexibility
- Easy-to-use
- Modularity
- Efficiency
- Good documentation





Website: <a href="mailto:speechbrain.github.io/">speechbrain.github.io/</a>

Code: github.com/speechbrain/speechbrain

Tutorials: github.com/speechbrain/speechbrain

Pretrained models: huggingface.co/speechbrain

## What is SpeechBrain?

SpeechBrain is designed from scratch to support multiple speech processing tasks.



## What can I do with SpeechBrain?



# What can I do with SpeechBrain?

#### SpeechBrain has also recipes for:

- Language Modeling
- Language Identification
- Sound Classification
- Grapheme-to-phoneme
- Multi-microphone signal processing
- EEG Decoding of Brain Signals

#### ..... and ongoing work for:

- Text-to-Speech
- Music Generation
- Speech Translation
- Voice Activity Detection
- Finite State Transducers (Integration with k2)









# **Design Principles**

#### Ease of use

<u>Simplicity & Modularity:</u> we develop intuitive **modules that are easy to interconnect with each other**. The code is pythonic and maximizes the use PyTorch routines.

<u>Lean software stack</u>: SpeechBrain employs a **simple software stack** (i.e., *Python*  $\rightarrow$  *PyTorch*  $\rightarrow$  *SpeechBrain*) to avoid dealing with too many levels of abstractions. PyTorch-compatible code works in our toolkit without any further modification

<u>Minimal external dependencies</u>: SpeechBrain has a **minimal list of external dependencies** that are all installable via PyPI. The installation process simply requires running the command pip install speechbrain and is done within a few minutes.





SpeechBrain PyTorch Python

# **Design Principles**

#### **Replicability & Transparency**

- SpeechBrain promotes **open** and **transparent science**.
- We trained most of our models with **publicly available data**. This way, our results can be easily **replicated** by the community.
- Several pre-trained models, which only require a few lines of code to use, are distributed via Hugging Face.
- Besides sharing the code and the trained models, we also share the whole experiment folder, which **contains all the needed details** (e.g., logs) to reproduce our results.







# **Design Principles**

#### Accessibility

- We have released SpechBrain under a very permissive license (Apache 2.0).
- SpeechBrain is designed to be easily understandable by a **large user base**, including early **students** and **practitioners**.
- We want SpeechBrain to serve **educational purposes** as well.
- We put several efforts on writing comprehensive **documentation**.





### How to Train a Model

• For all recipes and tasks, users can train a model simply with:





**Important computations** (e.g., forward step, objectives, data-io transformation) **are directly visible** in train.py

There is a **visible connection between the hyperparameter and the object** using it.

#### Inference

• SpeechBrain also provides functions for performing easy inference on pre-trained models:

#### **Speech Recognition**

```
from speechbrain.pretrained import EncoderDecoderASR
asr model =
EncoderDecoderASR.from hparams(source= "speechbrain/asr-crdnn-rnnlm-librispeech")
asr_model.transcribe_file( 'your_file.wav')
```



#### huggingface.co/speechbrain

**Speaker Verification** 

```
from speechbrain.pretrained import SpeakerRecognition
verification = SpeakerRecognition.from hparams(source= "speechbrain/spkrec-ecapa-voxceleb")
score, prediction = verification.verify_files( "file_1.wav", "file2_2.wav")
```

#### **Speech Separation**

```
from speechbrain.pretrained import SepformerSeparation as separator
model = separator.from hparams(source="speechbrain/sepformer-wsj02mix")
est_sources = model.separate_file(path='mixture.wav'
```

SpeechBrain currently has **25 pretrained models** on HuggingFace





#### <u>Tutorials</u>

We thus have written several tutorials (currently **22**) with **Google Colab** to help newcomers become more familiar with speech technologies.



https://speechbrain.github.io/

Examples of tutorials:

- Speech Recognition from Scratch
- Speech Classification from Scratch
- Speech Enhancement from Scratch



#### **Templates**

- Templates are **simple**, **well-documented recipes** that contain all the parts necessary for a **working system** (training, validation, evaluation, inference, ..).
- They cover a broad spectrum of types of tasks that are encountered in speech research, such as:
  - → sequence regression (enhancement)
  - → sequence to sequence (speech recognition)
  - → sequence classification (speaker ID)



# Tutorials

Templates

Docstrings + Code snippets

#### **Inline comments**

class SincConv(nn.Module):

"""This function implements SincConv (SincNet).

M. Ravanelli, Y. Bengio, "Speaker Recognition from raw waveform with SincNet", in Proc. of SLT 2018 (https://arxiv.org/abs/1808.00158)

#### Arguments

input shape : tuple The shape of the input. Alternatively use ``in\_channels``. in\_channels : int The number of input channels. Alternatively use ``input\_shape``. out\_channels : int It is the number of output channels. kernel size: int Kernel size of the convolutional filters. stride : int Stride factor of the convolutional filters. When the stride factor > 1. a decimation in time is performed. dilation : int Dilation factor of the convolutional filters. padding : str (same, valid, causal). If "valid", no padding is performed. If "same" and stride is 1, output shape is the same as the input shape. "causal" results in causal (dilated) convolutions. padding\_mode : str This flag specifies the type of padding. See torch.nn documentation for more information. aroups : int This option specifies the convolutional groups. See torch.nn documentation for more information. bias : bool If True, the additive bias b is adopted. sample rate : int, Sampling rate of the input signals. It is only used for sinc\_conv. min\_low\_hz : float Lowest possible frequency (in Hz) for a filter. It is only used for sinc\_conv. min low hz : float Lowest possible value (in Hz) for a filter bandwidth.

#### Example

.....

>>> inp\_tensor = torch.rand([10, 16000])
>>> conv = SincConv(input\_shape=inp\_tensor.shape, out\_channels=25, kernel\_size=11)

>>> out\_tensor = conv(inp\_tensor)
>>> out tensor.shape

torch.Size([10, 16000, 25])

#### Performance

• SpeechBrain is released with many recipes on popular datasets that achieve **competitive** or **state-of-the-art results** in many tasks.

Dataset	Task	System	Performance
LibriSpeech	Speech Recognition	CNN + Transformer	WER=2.35% (test-clean)
LibriSpeech	Speech Recognition	wav2vec2 + CTC	WER=1.90% (test-clean)
TIMIT	Speech Recognition	CRDNN + distillation	PER=13.1% (test)
TIMIT	Speech Recognition	wav2vec2 + CTC/Att.	PER=8.04% (test)
VoxCeleb2	Speaker Verification	ECAPA-TDNN	EER=0.69% (vox1-test)
AMI	Speaker Diarization	ECAPA-TDNN	DER=2.13% (lapel-mix)
VoiceBank	Speech Enhancement	MetricGAN+	PESQ=3.08 (test)
WSJ2MIX	Speech Separation	SepFormer	SDRi=22.6 dB (test)
WSJ3MIX	Speech Separation	SepFormer	SDRi=20.0 dB (test)

## **Other Features**

- Multi-GPU Training
- Dynamic Batching (based on bucketing)
- Large Scale Experiments with Webdataset
- On-the-fly augmentation and feature computation (on both CPU and GPU)
- Interface with HuggingFace for pre-training with large models (e.g, wav2vec)





## **Project History**



Project Announcement: September 2019

Development Started: February 2020





**NEW** RELEAS

- **Beta testing:** February 2021
  - Public Release: March 2021
  - Now



46 recipes for 20 datasets. 2.8k stars on github (in just 5 months)

SpeechBrain is growing very fast!



#### Reference

#### **SpeechBrain: A General-Purpose Speech Toolkit**

 Mirco Ravanelli<sup>1,2</sup>, Titouan Parcollet<sup>3,16</sup>, Peter Plantinga<sup>4</sup>, Aku Rouhe<sup>5</sup>, Samuele Cornell<sup>6</sup>, Loren Lugosch<sup>1,7</sup>, Cem Subakan<sup>1</sup>, Nauman Dawalatabad<sup>8</sup>, Abdelwahab Heba<sup>9</sup>,
 Jianyuan Zhong<sup>1</sup>, Ju-Chieh Chou<sup>10</sup>, Sung-Lin Yeh<sup>11\*</sup>, Szu-Wei Fu<sup>12</sup>, Chien-Feng Liao<sup>12</sup>, Elena Rastorgueva<sup>13†</sup>, François Grondin<sup>14</sup>, William Aris<sup>14</sup>, Hwidong Na<sup>15</sup>, Yan Gao<sup>16</sup>, Renato De Mori<sup>3,7</sup>, and Yoshua Bengio<sup>1,2</sup>

Website: <a href="mailto:speechbrain.github.io/">speechbrain.github.io/</a>

Code: github.com/speechbrain/speechbrain

Tutorials: github.com/speechbrain/speechbrain

Pretrained models: huggingface.co/speechbrain













#### **Institutional Partners**









#### **Core Development Team**





Mirco Ravanelli Titouan Parcollet F



Peter Plantinga C



Cem Sübakan Chien-Feng Liao



Szu-Wei Fu







Elena Rastorgueva Loren Lugosch Nauman Dawalatabad Aku Rouhe



Ju-Chieh Chou



Hwidong Na













Abdel Heba

Samuele Cornell

Sung-Lin Yeh

Francois Grondin William Aris

Yan Gao

## open-source toolkits ESPNet

# Today's talk

Introduction of ESPnet, end-to-end speech processing toolkit

#### Broadened applications

- Automatic speech recognition (ASR)
  - Performance improvement
  - New features
    - RNN-transducer
    - Non-autoregressive modeling
- Text to speech (TTS)
- Voice conversion
- Speech translation
- Speech enhancement

# Today's talk

- Introduction of ESPnet, end-to-end speech processing toolkit
- Broadened applications
- Automatic speech recognition (ASR)
  - Performance improvement
  - New features
    - RNN-transducer
    - Non-autoregressive modeling
- Text to speech (TTS)
- Voice conversion
- Speech translation
- Speech enhancement

Note that this presentation is a summary/collection of recent activities in the ESPnet toolkit

Please also check individual reports

#### ESPnet **ESPriet**nched in December 2017

#### **Our initial report at Interspeech 2018**

#### **ESPnet: End-to-End Speech Processing Toolkit**

Shinji Watanabe<sup>1</sup>, Takaaki Hori<sup>2</sup>, Shigeki Karita<sup>3</sup>, Tomoki Hayashi<sup>4</sup>, Jiro Nishitoba<sup>5</sup>, Yuya Unno<sup>6</sup>, Nelson Enrique Yalta Soplin<sup>7</sup>, Jahn Heymann<sup>8</sup>, Matthew Wiesner<sup>1</sup>, Nanxin Chen<sup>1</sup>, Adithya Renduchintala<sup>1</sup>, Tsubasa Ochiai<sup>9</sup>,

<sup>1</sup>Johns Hopkins University, <sup>2</sup>Mitsubishi Electric Research Laboratories, <sup>3</sup>NTT Communication Science Laboratories, <sup>4</sup>Nagoya University, <sup>5</sup>Retrieva, Inc., <sup>6</sup>Preferred Networks, Inc., <sup>7</sup>Waseda University, <sup>8</sup>Paderborn University, <sup>9</sup>Doshisha University

shinjiw@jhu.edu

#### Abstract

This paper introduces a new open source platform for end-toend speech processing named ESPnet. ESPnet mainly focuses on end-to-end automatic speech recognition (ASR), and adopts widely-used dynamic neural network toolkits, Chainer and Py-Torch, as a main deep learning engine. ESPnet also follows the Kaldi ASR toolkit style for data processing, feature extraction/format, and recipes to provide a complete setup for speech network [13, 14, 15, 16]. Attention-based methods use an attention mechanism to perform alignment between acoustic frames and recognized symbols, while CTC uses Markov assumptions to efficiently solve sequential problems by dynamic programming. ESPnet adopts hybrid CTC/attention end-to-end ASR [17], which effectively utilizes the advantages of both architectures in training and decoding. During training, we employ the multiobjective learning framework to improve robustness on irregular alignments and achieve fast convergence. During de-

# ESPnet **ESPnet** in December 2017

https://github.com/espnet/espnet

- **Open source** (Apache2.0) end-to-end speech processing toolkit
- Major concept: Reproducibility
  - $\odot$  Leverage our end-to-end ASR experience to the community
  - Accelerates end-to-end research for speech researchers
- PyTorch based dynamic neural network toolkit as an engine
  - Easily develop novel neural network architecture

#### • Follows the famous speech recognition toolkit, Kaldi, style



- Smoothly port ASR experiments from Kaldi to ESPnet with the common data processing, feature extraction/format
- $\circ$  Recipes to provide a complete setup for speech processing experiments

#### • The project is greatly accelerated in these three years

# ESPnet **ESPnet** in December 2017

https://github.com/espnet/espn

et

- Open source (Apache2.0) end-to-end speech processing toolkit
- Major concept: Reproducibility
  - $\odot$  Leverage our end-to-end ASR experience to the community
  - Accelerates end-to-end research for speech researchers
- PyTorch based dynamic neural network toolkit as an engine
  - Easily develop novel neural network architecture

#### • Follows the famous speech recognition toolkit, Kaldi, style

 Smoothly port ASR experiments from Kaldi to ESPnet with the common data processing, feature extraction/format

SpeechBrain

Pvthon

 $\,\circ\,$  Recipes to provide a complete setup for speech processing experiments

• The project is greatly accelerated in these three years  $/_{PyTorch}$ 

# ESPnet **ESPnet** in December 2017

https://github.com/espnet/espn

et

- Open source (Apache2.0) end-to-end speech processing toolkit
- Major concept: Reproducibility
  - $\odot$  Leverage our end-to-end ASR experience to the community
  - Accelerates end-to-end research for speech researchers
- PyTorch based dynamic neural network toolkit as an engine
  - Easily develop novel neural network architecture

#### • Follows the famous speech recognition toolkit, Kaldi, style

- Smoothly port ASR experiments from Kaldi to ESPnet with the common data processing, feature extraction/format
- Recipes to provide a complete setup for speech processing experiments

```
• The project is greatly accelerated in these three years /_{PyTorch}
```

Good compatibility/scalability with Linux systems and job schedular but not pythonic and hard for begginers ESPnet

Python

Bash

### Activity statistics (from 2018 to 2020)



- Citations, contributors, recipes (examples), and stars are all growing i.e.,
- Developers have increasingly supported the development of ESPnet
- has been used in various research groups and contributed a lot to speech research activities
- ESPnet 4.1 stars (today)
  - Kaldi 10.7K stars
  - Nvidia NeMo 3.1K stars
  - SpeechBrain 2.8K stars
  - Google Lingvo 2.3K stars
  - FairSeq: 13.6K stars

## Major change in the internal framework From ESPnet1 to ESPnet2

- ESPnet2: a new system for DNN training to extend our system from v.0.7.0
- Mostly refactoring, but which enables major update to deal with

#### Distributed training

- On-the-fly feature extraction from the raw waveform
- Improved the scalability
- Improved software workflow by enhancing
  - continuous integration, enriching documentation, supporting the docker, pip install
  - model zoo (mostly zenodo compared with hugging face HUB used in SpeechBrain. hugging face HUB migration is on-going)
  - WandB based experiment monitoring and sharing

• The migration is ongoing (ASR and TTS are already finished is weights & Biases



# Today's talk

Introduction of ESPnet, end-to-end speech processing toolkit

#### Broadened applications

- Automatic speech recognition (ASR)
  - Performance improvement
  - New features
    - RNN-transducer
    - Non-autoregressive modeling
- Text to speech (TTS)
- Voice conversion
- Speech translation
- Speech enhancement

#### **Broadened Applications**

• ESPnet (**ASR+X**) covers the following topics complementally



Why can one toolkit support such wide-ranges of applications?

#### **Broadened Applications**

• ESPnet (**ASR+X**) covers the following topics complementally



• Why can one toolkit support such wide-ranges of applications?

#### Unified form $\rightarrow$ Unified software design

We design ESPnet by leveraging a **unified** mathematical **form** of **sequence** (*X*) **to sequence** (*Y*) **transformation** *f* 

$$X = (x_1, x_2, \cdots, x_T) \xrightarrow{f} Y = (y_1, y_2, \cdots, y_N)$$














$$X = (x_1, x_2, \cdots, x_T)$$

$$Y = (y_1, y_2, \cdots, y_N)$$
ESPhot: End-to-ond  
speech processing toolkit

- Many speech processing applications can be **unified** based on seq2seq
- Nemo, Fairseq, Lingvo, Espresso, SpeechBrain, Asteroid and other toolkits also fully make use of these functions
- We are closely collaborating/interacting with them

### Today's talk

- Introduction of ESPnet, end-to-end speech processing toolkit
- Broadened applications
- Automatic speech recognition (ASR)
  - O Performance improvement
  - O New features
    - RNN-transducer
    - Non-autoregressive modeling
- Text to speech (TTS)
- Voice conversion
- Speech translation
- Speech enhancement

#### Automatic speech recognition (ASR)



#### Maintaining state-of-the-art performance in ASR



Error rate:

Lower is better

**Release Date** 

#### Maintaining state-of-the-art performance in ASR



Error rate: Lower is better

**Release Date** 

#### **ESPnet Transformer**



#### A COMPARATIVE STUDY ON TRANSFORMER VS RNN IN SPEECH APPLICATIONS

Shigeki Karita<sup>1</sup>,

(Alphabetical Order) Nanxin Chen<sup>3</sup>, Tomoki Hayashi<sup>5,6</sup>, Takaaki Hori<sup>7</sup>, Hirofumi Inaguma<sup>8</sup>, Ziyan Jiang<sup>3</sup>, Masao Someki<sup>5</sup>, Nelson Enrique Yalta Soplin<sup>2</sup>, Ryuichi Yamamoto<sup>4</sup>, Xiaofei Wang<sup>3</sup>, Shinji Watanabe<sup>3</sup>, Takenori Yoshimura<sup>5,6</sup>, Wangyou Zhang<sup>9</sup>

<sup>1</sup>NTT Communication Science Laboratories, <sup>2</sup>Waseda University, <sup>3</sup>Johns Hopkins University, <sup>4</sup>LINE Corporation, <sup>5</sup>Nagoya University, <sup>6</sup>Human Dataware Lab. Co., Ltd.,
 <sup>7</sup>Mitsubishi Electric Research Laboratories, <sup>8</sup>Kyoto University, <sup>9</sup>Shanghai Jiao Tong University

#### ABSTRACT

Sequence-to-sequence models have been widely used in end-toend speech processing, for example, automatic speech recognition (ASR), speech translation (ST), and text-to-speech (TTS). This paper focuses on an emergent sequence-to-sequence model called Transformer, which achieves state-of-the-art performance in neural machine translation and other natural language processing applications. We undertook intensive studies in which we experimentally compared and analyzed Transformer and conventional recurrent In our speech application experiments, we investigate several aspects of Transformer and RNN-based systems. For example, we measure the word/character/regression error from the ground truth, training curve, and scalability for multiple GPUs.

The contributions of this work are:

- We conduct a larges-scale comparative study on Transformer and RNN with significant performance gains especially for the ASR related tasks.
- We explain our training tips for Transformer in speech applica-

#### **ESPnet Transformer**



#### A COMPARATIVE STUDY ON TRANSFORMER VS RNN IN SPEECH APPLICATIONS

Shigeki Karita<sup>1</sup>,

(Alphabetical Order) Nanxin Chen<sup>3</sup>, Tomoki Hayashi<sup>5,6</sup>, Takaaki Hori<sup>7</sup>, Hirofumi Inaguma<sup>8</sup>, Ziyan Jiang<sup>3</sup>, Masao Someki<sup>5</sup>, Nelson Enrique Yalta Soplin<sup>2</sup>, Ryuichi Yamamoto<sup>4</sup>, Xiaofei Wang<sup>3</sup>, Shinji Watanabe<sup>3</sup>, Takenori Yoshimura<sup>5,6</sup>, Wangyou Zhang<sup>9</sup>

<sup>1</sup>NTT Communication Science Laboratories, <sup>2</sup>Waseda University, <sup>3</sup>Johns Hopkins University, <sup>4</sup>LINE Corporation, <sup>5</sup>Nagoya University, <sup>6</sup>Human Dataware Lab. Co., Ltd.,
 <sup>7</sup>Mitsubishi Electric Research Laboratories, <sup>8</sup>Kyoto University, <sup>9</sup>Shanghai Jiao Tong University

#### ABSTRACT

Sequence-to-sequence models have been wid end speech processing, for example, automatic (ASR), speech translation (ST), and text-topaper focuses on an emergent sequence-to-sec Transformer, which achieves state-of-the-art per

machine translation and other natural language processing applications. We undertook intensive studies in which we experimentally compared and analyzed Transformer and conventional recurrent

One of the first success in the speech areas The performance was boosted

and RNN with significant performance gains especially for the ASR related tasks.

• We explain our training tips for Transformer in speech applications: ASP\_TTS and ST

#### **Transformer boosted the performance**

Improve the performance from RNN with 13 ASR tasks among 15 tasks

Reaching the Kaldi performance (state-of-the-art **non** end-to-end ASR) in half of tasks

dataset	token	error	Kaldi	Our RNN	Our Transformer
AISHELL	char	CER	N/A / 7.4	6.8 / 8.0	6.0 / 6.7
AURORA4	char	WER	(*) 3.6 / 7.7 / 10.0 / 22.3	3.5 / 6.4 / 5.1 / 12.3	3.3 / 6.0 / 4.5 / 10.6
CSJ	char	CER	(*) 7.5 / 6.3 / 6.9	6.6 / 4.8 / 5.0	5.7 / 4.1 / 4.5
CHiME4	char	WER	6.8 / 5.6 / 12.1 / 11.4	9.5 / 8.9 / 18.3 / 16.6	9.6 / 8.2 / 15.7 / 14.5
CHiME5	char	WER	47.9 / 81.3	59.3 / 88.1	60.2 / 87.1
Fisher-CALLHOME Spanish	char	WER	N/A	27.9/27.8/25.4/47.2/47.9	27.0 / 26.3 / 24.4 / 45.3 / 46.2
HKUST	char	CER	23.7	27.4	23.5
JSUT	char	CER	N/A	20.6	18.7
LibriSpeech	BPE	WER	3.9 / 10.4 / 4.3 / 10.8	3.1/9.9/3.3/10.8	2.2 / 5.6 / 2.6 / 5.7
REVERB	char	WER	18.2 / 19.9	24.1/27.2	15.5 / 19.0
SWITCHBOARD	BPE	WER	18.1 / 8.8	28.5 / 15.6	26.0 / 14.0
TED-LIUM2	BPE	WER	<b>9.0</b> / 9.0	11.2 / 11.0	9.3 / 8.1
TED-LIUM3	BPE	WER	6.2 / 6.8	14.3 / 15.0	9.7 / 8.0
VoxForge	char	CER	N/A	12.9 / 12.6	9.4 / 9.1
WSJ	char	WER	4.3 / 2.3	7.0 / 4.7	6.8 / 4.4

#### **Transformer boosted the performance**

#### Improve the performance from RNN with 13 ASR tasks among 15 tasks



### Experiments (~ 1000 hours) Librispeech at April 2019

Toolkit	dev_clean	dev_other	test_clean	test_other
Facebook wav2letter++	3.1	10.1	3.4	11.2
RWTH RASR	2.9	8.8	3.1	9.8
Nvidia Jasper	2.6	7.6	2.8	7.8
Google SpecAug.	N/A	N/A	2.5	5.8

Very impressive results by Google

### Experiments (~ 1000 hours) Librispeech at August 2019

Toolkit/Method	dev_clean	dev_other	test_clean	test_other
Facebook wav2letter++	3.1	10.1	3.4	11.2
RWTH RASR	2.9	8.8	3.1	9.8
Nvidia Jasper	2.6	7.6	2.8	7.8
Google SpecAug.	N/A	N/A	2.5	5.8
ESPnet	2.2	5.6	2.6	5.7

Reached Google's best performance by community-driven efforts







### Experiments (~ 1000 hours) in August 2019 Librispeech

Toolkit/Method	dev_clean	dev_other	test_clean	test_other
Facebook wav2letter++	3.1	10.1	3.4	11.2
RWTH RASR	2.9	8.8	3.1	9.8
Nvidia Jasper	2.6	7.6	2.8	7.8
Google SpecAug.	N/A	N/A	2.5	5.8
ESPnet	2.2	5.6	2.6	5.7

### Experiments (~ 1000 hours) in March 2020 Librispeech

Toolkit/Method	dev_clean	dev_other	test_clean	test_other
Facebook wav2letter++	3.1	10.1	3.4	11.2
RWTH RASR	2.9	8.8	3.1	9.8
Nvidia Jasper	2.6	7.6	2.8	7.8
Google SpecAug.	N/A	N/A	2.5	5.8
ESPnet	2.2	5.6	2.6	5.7
Google Conformer	1.9	4.4	1.9	3.9

#### **ESPnet Conformer**



#### **RECENT DEVELOPMENTS ON ESPNET TOOLKIT BOOSTED BY CONFORMER**

Pengcheng Guo<sup>1,4</sup>, Florian Boyer<sup>2,3</sup>, Xuankai Chang<sup>4</sup>, Tomoki Hayashi<sup>5</sup>, Yosuke Higuchi<sup>6</sup> Hirofumi Inaguma<sup>7</sup>, Naoyuki Kamo<sup>8</sup>, Chenda Li<sup>9</sup>, Daniel Garcia-Romero<sup>4</sup>, Jiatong Shi<sup>4</sup> Jing Shi<sup>4,10</sup>, Shinji Watanabe<sup>4</sup>, Kun Wei<sup>1</sup>, Wangyou Zhang<sup>9</sup>, Yuekai Zhang<sup>4</sup>

<sup>1</sup>Northwestern Polytechnical University, <sup>2</sup>LaBRI, University of Bordeaux, <sup>3</sup> Airudit
 <sup>4</sup>Johns Hopkins University, <sup>5</sup>Human Dataware Lab. Co., Ltd.
 <sup>6</sup>Waseda University, <sup>7</sup>Kyoto University, <sup>8</sup>NTT Corporation <sup>9</sup>Shanghai Jiao Tong University
 <sup>10</sup>Institute of Automation, Chinese Academy of Sciences

#### ABSTRACT

In this study, we present recent developments on ESPnet: End-to-End Speech Processing toolkit, which mainly involves a recently proposed architecture called Conformer, Convolution-augmented Transformer. This paper shows the results for a wide range of endto-end speech processing applications, such as automatic speech recognition (ASR), speech translations (ST), speech separation (SS) and text-to-speech (TTS). Our experiments reveal various training tips and significant performance benefits obtained with the Conformer on different tasks. These results are competitive or even outperform the current state-of-art Transformer models. We are preparing to release all-in-one recipes using open source and even in the sum of the other tasks. of publicly available corpora and try our best to share the practical guides (e.g., learning rate, hyper-parameters, network structure) on the use of Conformer. We also prepare to release the reproducible recipes and state-of-the-art setups to the community to succeed our exciting outcomes.

The contributions of this study include:

- We extend the Conformer architecture to various end-to-end speech processing applications and conduct comparative experiments with Transformer.
- We share our practical guides for the training of Conformer, like learning rate, kernel size of Conformer block, and model architectures, etc.

#### **ESPnet Conformer**



#### **RECENT DEVELOPMENTS ON ESPNET TOOLKIT BOOSTED BY CONFORMER**

Pengcheng Guo<sup>1,4</sup>, Florian Boyer<sup>2,3</sup>, Xuankai Chang<sup>4</sup>, Tomoki Hayashi<sup>5</sup>, Yosuke Higuchi<sup>6</sup> Hirofumi Inaguma<sup>7</sup>, Naoyuki Kamo<sup>8</sup>, Chenda Li<sup>9</sup>, Daniel Garcia-Romero<sup>4</sup>, Jiatong Shi<sup>4</sup> Jing Shi<sup>4,10</sup>, Shinji Watanabe<sup>4</sup>, Kun Wei<sup>1</sup>, Wangyou Zhang<sup>9</sup>, Yuekai Zhang<sup>4</sup>

<sup>1</sup>Northwestern Polytechnical University, <sup>2</sup>LaBRI, University of Bordeaux, <sup>3</sup> Airudit
 <sup>4</sup>Johns Hopkins University, <sup>5</sup>Human Dataware Lab. Co., Ltd.
 <sup>6</sup>Waseda University, <sup>7</sup>Kyoto University, <sup>8</sup>NTT Corporation <sup>9</sup>Shanghai Jiao Tong University
 <sup>10</sup>Institute of Automation, Chinese Academy of Sciences

#### ABSTRACT

In this study, we present recent developments End Speech Processing toolkit, which mainly proposed architecture called Conformer, Co Transformer. This paper shows the results for to-end speech processing applications, such recognition (ASR), speech translations (ST), sp and text-to-speech (TTS). Our experiments r

ing tips and significant performance benefits obtained with the Conformer on different tasks. These results are competitive or even outperform the current state-of-art Transformer models. We are preparing to release all-in-one recipes using open source and available compare for all the above tasks with the trained of publicly available corpora and try our best to share the practical

We try to follow Google's conformer work
Also apply conformer to ST, TTS, as well as ASR

periments with Transformer.

• We share our practical guides for the training of Conformer, like learning rate, kernel size of Conformer block, and model architectures. etc.

### Experiments (~ 1000 hours) in October 2020 Librispeech

Toolkit/Method	dev_clean	dev_other	test_clean	test_other
Facebook wav2letter++	3.1	10.1	3.4	11.2
RWTH RASR	2.9	8.8	3.1	9.8
Nvidia Jasper	2.6	7.6	2.8	7.8
Google SpecAug.	N/A	N/A	2.5	5.8
ESPnet	2.2	5.6	2.6	5.7
Google Conformer	1.9	4.4	1.9	3.9
ESPnet Conformer	1.9	4.6	2.1	4.7

We continue to work on catching up SOTA.

Is the story ended?  $\rightarrow$  No...







### Experiments (~ 1000 hours) in March 2021 Librispeech

Toolkit/Method	dev_clean	dev_other	test_clean	test_other
Google SpecAug.	N/A	N/A	2.5	5.8
ESPnet	2.2	5.6	2.6	5.7
Google Conformer	1.9	4.4	1.9	3.9
ESPnet Conformer	1.9	4.6	2.1	4.7
Facebook wav2vec2.0 (60k LibriVox)	1.6	3.0	1.8	3.3
Facebook Hubert (60k LibriVox)	1.7	3.0	1.9	3.5

Self-supervised training further improves the performance





### Experiments (~ 1000 hours) in July 2021 Librispeech

Toolkit/Method	dev_clean	dev_other	test_clean	test_other
Google SpecAug.	N/A	N/A	2.5	5.8
ESPnet	2.2	5.6	2.6	5.7
Google Conformer	1.9	4.4	1.9	3.9
ESPnet Conformer	1.9	4.6	2.1	4.7
Facebook wav2vec2.0 (60k LibriVox)	1.6	3.0	1.8	3.3
Facebook Hubert (60k LibriVox)	1.7	3.0	1.9	3.5
ESPnet + S3PRL Hubert (60k LibriVox)	1.7	3.4	1.8	3.6

Reaching SOTA again!!!

### Experiments (~ 1000 hours) in July 2021 Librispeech

Toolkit/Method	dev_clean	dev_other	test_clean	test_other
Google SpecAug.	N/A	N/A	2.5	5.8
ESPnet	2.2	5.6	2.6	5.7
Google Conformer	1.9	4.4	1.9	3.9
ESPnet Conformer	1.9	4.6	2.1	4.7
Facebook wav2vec2.0 (60k LibriVox)	1.6	3.0	1.8	3.3
Facebook Hubert (60k LibriVox)	1.7	3.0	1.9	3.5
ESPnet + S3PRL Hubert (60k LibriVox)	1.7	3.4	1.8	3.6

Note that this was achieved by great helps from Google, Facebook, and SUPERB

# Good example of "Collapetition"

## = Collaboration + Competition

#### **Discussions for catching SOTA**

- We should not give up to catch up SOTA
- We developed Transformer -> Conformer -> Hubert
  - All activities have been supported by various collaborators (Collapetition)
  - We have further interactions with others including SpeechBrain,
     K2, etc.
- Now self-supervised models are ready in ESPnet through SUPERB S3PRL!!!
  - ESPnet can become one of the downstream tasks in SUPERB
  - It can be applied to **LeBenchmark**
## Today's talk

- Introduction of ESPnet, end-to-end speech processing toolkit
- Broadened applications
- Automatic speech recognition (ASR)
  - O Performance improvement
  - New features
    - RNN-transducer
    - Non-autoregressive modeling
- Text to speech (TTS)
- Voice conversion
- Speech translation
- Speech enhancement

## **RNN/Transformer Transducer [Boyer+(2021)]**

RNN or transformer transducer (widely used in industry)

Good for streaming



- ESPnet has various architecture supports (LSTM, CNN, Transformer, conformer)
- Also supports various beam search algorithms, including K2 FST beam search (similar to SpeechBrain)

## Non-Autoregressive modeling [Higuchi+(2020)]

- The most complicated part in ASR: Left-to-right beam search (several hundreds of lines)
- BERT-like iterative mask predict
- Achieved 0.07 real time factor! (Takes only 70ms to decode

   No software optimization
  - Just CPU
- Only 20 lines for coding



## Today's talk

- Introduction of ESPnet, end-to-end speech processing toolkit
- Broadened applications
- Automatic speech recognition (ASR)
  - O Performance improvement
  - O New features
    - RNN-transducer
    - Non-autoregressive modeling
- Text to speech (TTS)
- Voice conversion
- Speech translation
- Speech enhancement

## Text to speech (TTS)



### **ESPnet TTS**

#### ESPNET-TTS: UNIFIED, REPRODUCIBLE, AND INTEGRATABLE OPEN SOURCE END-TO-END TEXT-TO-SPEECH TOOLKIT

 Tomoki Hayashi<sup>1,2</sup>, Ryuichi Yamamoto<sup>3</sup>, Katsuki Inoue<sup>4</sup>, Takenori Yoshimura<sup>1,2</sup>, Shinji Watanabe<sup>5</sup>, Tomoki Toda<sup>1</sup>, Kazuya Takeda<sup>1</sup>, Yu Zhang<sup>6</sup>, and Xu Tan<sup>7</sup>
 <sup>1</sup>Nagoya University, <sup>2</sup>Human Dataware Lab. Co., Ltd., <sup>3</sup>LINE Corp.,
 <sup>4</sup>Okayama University, <sup>5</sup>Johns Hopkins University, <sup>6</sup>Google AI, <sup>7</sup>Microsoft Research

#### ABSTRACT

This paper introduces a new end-to-end text-to-speech (E2E-TTS) toolkit named ESPnet-TTS, which is an extension of the open-source speech processing toolkit ESPnet. The toolkit supports state-of-theart E2E-TTS models, including Tacotron 2, Transformer TTS, and FastSpeech, and also provides recipes inspired by the Kaldi automatic speech recognition (ASR) toolkit. The recipes are based on the design unified with the ESPnet ASR recipe, providing high reproducibility. The toolkit also provides pre-trained models and samples of all of the recipes so that users can use it as a baseline. Furthermore, the unified design enables the integration of ASR functions with TTS, e.g., ASR-based objective evaluation and semi-supervised learning with both ASR and TTS models. This paper describes the research purpose to make E2E-TTS systems more user-friendly and to accelerate research in this field. The toolkit not only supports state-of-the-art E2E-TTS models such as Tacotron 2 [6], Transformer TTS [8], and FastSpeech [9] but also provides Kaldi automatic speech recognition (ASR) toolkit [17] style recipes. The recipe is based on the design unified with the ASR recipe and includes all of the procedures required to reproduce the results. The toolkit provides a number of recipes for more than ten languages, which include single-speaker TTS as well as multi-speaker one and speaker adaptation. Pre-trained models and generated samples of all of the recipes are also provided so that users can easily use it as a baseline or perform TTS demonstrations. Furthermore, thanks to the unified design among TTS and ASR, we can easily integrate ASR

## ESPnet TTS

- Mainly focuses on the development of text to mel-spectrogram (text2mel) models.
  - It supports Tacotron2 and Transformer-TTS (AR), and FastSpeech and FastSpeech2 (non-AR)
  - Of course, we can easily switch to RNN, transformer, or conformer
- Multi-speaker extensions with X-vector and global style token.
- Users can quickly develop the state-of-the-art baseline systems for the research purpose
- A lot of examples and demonstration systems, which works in real-time for various languages, including English, Mandarin, and Japanese
   <u>https://colab.research.google.com/github/espnet/notebook/blob/master/espnet2\_tts\_realtim\_e\_demo.ipynb</u>
- Again supported by a lot of helps from Google (Dr. Heiga Zen)!!!

### Synthesis

```
# decide the input sentence by yourself
print(f"Input your favorite sentence in {lang}.")
    x = input()
    # synthesis
    with torch.no grad():
        start = time.time()
        wav, c, * = text2speech(x)
        wav = vocoder.inference(c)
    rtf = (time.time() - start) / (len(wav) / fs)
    print(f"RTF = {rtf:5f}")
    # let us listen to generated samples
    from IPython.display import display, Audio
```

```
display(Audio(wav.view(-1).cpu().numpy(), rate=fs))
```

```
Input your favorite sentence in English.
I'll talk about T T S functions in E S P net
RTF = 0.025100
```



• 0:03 / 0:03

### **Voice conversion challenge 2020 baseline system**

#### The Sequence-to-Sequence Baseline for the Voice Conversion Challenge 2020: Cascading ASR and TTS

Wen-Chin Huang<sup>1</sup>, Tomoki Hayashi<sup>1</sup>, Shinji Watanabe<sup>2</sup>, Tomoki Toda<sup>1</sup>

<sup>1</sup>Nagoya University, Japan <sup>2</sup>Johns Hopkins University, USA

wen.chinhuang@g.sp.m.is.nagoya-u.ac.jp

#### Abstract

This paper presents the sequence-to-sequence (seq2seq) baseline system for the voice conversion challenge (VCC) 2020. We consider a naive approach for voice conversion (VC), which is to first transcribe the input speech with an automatic speech recognition (ASR) model, followed using the transcriptions to generate the voice of the target with a text-to-speech (TTS) model. We revisit this method under a sequence-tosequence (seq2seq) framework by utilizing ESPnet, an opensource end-to-end speech processing toolkit, and the many well-configured pretrained models provided by the community. Official evaluation results show that our system comes out top among the participating systems in terms of conversion similarity, demonstrating the promising ability of seq2seq models to convert speaker identity. The implementation is



### Voice conversion challenge 2020 baseline system

### The Sequence-to-Sequence Baseline for the Voice Conversion Challenge 2020: Cascading ASR and TTS

Wen-Chin Huang<sup>1</sup>, Tomoki Hayashi<sup>1</sup>, Shinji Watanabe<sup>2</sup>, Tomoki Toda<sup>1</sup>

<sup>1</sup>Nagoya University, Japan <sup>2</sup>Johns Hopkins University, USA

wen.chinhuang@g.sp.m.is.nagoya-u.ac.jp

#### Abstract

This paper presents the sequence-to-sequence (seq2seq) baseline system for the voice conversion challenge (VCC) 2020.

 $\bullet$ 

We consider a naive approach for voice which is to first transcribe the input speech speech recognition (ASR) model, followed u tions to generate the voice of the target wit (TTS) model. We revisit this method und sequence (seq2seq) framework by utilizing source end-to-end speech processing toolk

well-configured pretrained models provided by the community. Official evaluation results show that our system comes out top among the participating systems in terms of conversion similarity, demonstrating the promising ability of seq2seq models to convert speaker identity. The implementation is



## Today's talk

- Introduction of ESPnet, end-to-end speech processing toolkit
- Broadened applications
- Automatic speech recognition (ASR)
  - O Performance improvement
  - O New features
    - RNN-transducer
    - Non-autoregressive modeling
- Text to speech (TTS)
- Voice conversion
- Speech translation
- Speech enhancement

## **Speech to text translation (ST)**



### **ESPnet-ST**

#### **ESPnet-ST: All-in-One Speech Translation Toolkit**

Hirofumi Inaguma<sup>1</sup> Shun Kiyono<sup>2</sup> Kevin Duh<sup>3</sup> Shigeki Karita<sup>4</sup> Nelson Yalta<sup>5</sup> Tomoki Hayashi<sup>6,7</sup> Shinji Watanabe<sup>3</sup> <sup>1</sup> Kyoto University <sup>2</sup> RIKEN AIP <sup>3</sup> Johns Hopkins University <sup>4</sup> NTT Communication Science Laboratories <sup>5</sup> Waseda University <sup>6</sup> Nagoya University <sup>7</sup> Human Dataware Lab. Co., Ltd. inaguma@sap.ist.i.kyoto-u.ac.jp

#### Abstract

We present *ESPnet-ST*, which is designed for the quick development of speech-to-speech translation systems in a single framework. *ESPnet-ST* is a new project inside end-toend speech processing toolkit ESPnet which can reduce latency at inference time, which is useful for time-critical use cases like simultaneous interpretation. (2) A single model enables back-propagation training in an end-to-end fashion, which mitigates the risk of error propagation by cascaded modules. (3) In certain use cases



- Support the speech translation (ST) task with both the traditional pipeline approach (ASR + NMT) and end-to-end (E2E) approach
  - ESPnet also supports NMT and comparable performance to the other toolkit
- We demonstrated the state-of-the-art translation performance in standard ST benchmarks
  - O MUST-C, IWSLT, Fisher Callhome Spanish
- Again, new progresses in ASR can be easily transferred to ST performance improvement, e.g., conformer, non-AR modeling

## Today's talk

- Introduction of ESPnet, end-to-end speech processing toolkit
- Broadened applications
- Automatic speech recognition (ASR)
  - O Performance improvement
  - O New features
    - RNN-transducer
    - Non-autoregressive modeling
- Text to speech (TTS)
- Voice conversion
- Speech translation
- Speech enhancement

## Speech enhancement (SE)





### denoise dereverberate separate





### **ESPnet-SE**

#### ESPNET-SE: END-TO-END SPEECH ENHANCEMENT AND SEPARATION TOOLKIT DESIGNED FOR ASR INTEGRATION

Chenda Li<sup>1\*</sup>, Jing Shi<sup>2,3\*</sup>, Wangyou Zhang<sup>1\*</sup>, Aswin Shanmugam Subramanian<sup>3</sup>, Xuankai Chang<sup>3</sup>, Naoyuki Kamo, Moto Hira<sup>4</sup>, Tomoki Hayashi<sup>5,6</sup>, Christoph Boeddeker<sup>7</sup>, Zhuo Chen<sup>8</sup>, Shinji Watanabe<sup>3</sup> <sup>1</sup>Shanghai Jiao Tong University, <sup>2</sup>Institute of Automation, Chinese Academy of Sciences, <sup>3</sup>Johns Hopkins University, <sup>4</sup>Facebook AI, <sup>5</sup>Nagoya University, <sup>6</sup>Human Dataware Lab. Co., Ltd., <sup>7</sup>Paderborn University, <sup>8</sup>Microsoft Research

#### ABSTRACT

We present ESPnet-SE, which is designed for the quick development of speech enhancement and speech separation systems in a single framework, along with the optional downstream speech recognition module. ESPnet-SE is a new project which integrates rich automatic speech recognition related models, resources and systems to support and validate the proposed front-end implementation (i.e. speech enhancement and separation). It is capable of processing both singlechannel and multi-channel data, with various functionalities including dereverberation, denoising and source separation. We provide all-in-one recipes including data pre-processing, feature extraction, training and evaluation pipelines for a wide range of benchmark In this paper, we introduce a new E2E-SE toolkit named ESPnet-SE<sup>1</sup>, which is an extension of the open-source speech processing toolkit ESPnet [8]. ESPnet-SE fully considers the various forms of speech input in the front-end scenes and meanwhile flexibly and organically integrated with the downstream automatic speech recognition (ASR) task, making it a user-friendly toolkit to easily build totally end-to-end robust ASR systems, even without need for clean speech signals. The toolkit provides adaptability to different speech data, including (1) single and multiple speakers, (2) single and multiple channels, (3) anechoic and reverberant conditions. Moreover, thanks to the ripe and efficient ASR modules in ESPnet, rich speech recognition related models, resources and systems can be optionally concatenated after the E2E-SE system, enabling evaluation and joint

### **ESPnet-SE**

- One of the biggest changes in ESPnet
- Include **ALL** speech enhancement functions
- Demonstration

https://colab.research.google.com/drive/1fjRJCh96SoYLZPRxsjF9VDv4Q2Vol ckl?usp=sharing

- Many help from **torchaudio**
- **SpeechBrain** also has the strong SE functions

## Speech enhancement Several types of problems

Denoising (people mainly call it speech enhancement)



## Speech enhancement Several types of problems

Denoising (people mainly call it speech enhancement)



▼ Enhance the single-channel real noisy speech in CHiME4



## Speech enhancement Several types of problems





## Speech enhancement Several types of problems





Separate the example in wsj0\_2mix testing set

```
!gdown --id 1ZCUkd_Lb7p02rpPr4FqYdtJBZ7JMiInx -0 /content/447c020t_1.2106_422a0112_-1.2106.wav
import os
import soundfile
from IPython.display import display, Audio
mixwav, sr = soundfile.read("447c020t_1.2106_422a0112_-1.2106.wav")
waves_wsj = separate_speech(mixwav[None, ...], fs=sr)
print("Input mixture", flush=True)
display(Audio(mixwav, rate=sr))
print(f"=======separated speech with model {tag} =======", flush=True", flush=True", flush=True, display(Audio(waves_wsj[0].squeeze(), rate=sr))
print("Separated spk2", flush=True, display(Audio(waves_wsj[1].squeeze(), rate=sr))
```

#### Downloading...

From: https://drive.google.com/uc?id=1ZCUkd\_Lb7p02rpPr4FqYdtJBZ7JMiInx
To: /content/447c020t\_1.2106\_422a0112\_-1.2106.wav
100% 184k/184k [00:00<00:00, 5.73MB/s]
100% 13/13 [00:00<00:00, 28.27it/s]Input mixture</pre>

▶ 0:00 / 0:11 →

===== Separated speech with model Chenda Li/wsj0\_2mix\_enh\_train\_enh\_conv\_tasnet\_raw\_valid.si\_snr. Separated spk1

► 0:00 / 0:11 → • ÷

Separated spk2

► 0:00 / 0:11 → • ÷





## Microphone array processing Single to multiple microphones

Denoising (people mainly call it speech enhancement)



## Microphone array processing Single to multiple microphones

Denoising (people mainly call it speech enhancement)



## Microphone array processing Single to multiple microphones

Denoising (people mainly call it speech enhancement)



Make a spatial **beam** (beamforming) to only pick up desired signals Enhance the multi-channel real noisy speech in CHiME4



## **Differentiable speech enhancement frontend**

sternid

• ESPnet SE can be used as an **independent** enhancement module

- Denoising, Dereveberation, Separation
- Single channel or multichannel
- Can port Asteroid (audio source separation toolkit) pre-trained models
- ESPnet SE can be used as a differentiable enhancement module



• We can realize a cocktail party effect by a machine



## **ttt ESPnet**

- End-to-end speech processing has a lot of potentials
- ESPnet provides **state-of-the-art** and **reproducible** research in various speech processing applications
- •We are working closely with other benchmarks (SUPERB, LeBenchmark, CHiME, VCC, etc.) and toolkit



### We need further interactions with machine leaning community! Let's work together for the community contribution!!

## Discussions

### (these are commonly shared with SpeechBrain project)

- Complementary task covering with multiple toolkits
  - ESPnet
    - More focus on ASR/TTS/Speech translation with SOTA
  - SpeechBrain
    - More applications, and SOTA for ASR, speaker identification, spoken language understanding, and speech separation
  - Other tools also have unique functions ;)
- ESPnet and SpeechBrain are academic driven
  - Less biased for techniques
  - Any collaborations based on academic freedom!
  - Weak for the production support in general e.g., streaming and intensive tuning
  - Difficulties on maintainability and scalability
- Speech is one of the most important machine learning applications
  - How should these toolkits be evolved for machine learning researchers?
- We also want to discuss these items during the following sessions!!!

### **Concluding Remarks**

## **Concluding Remarks**

- Self-supervised learning for speech
- Open-source toolkits for wide variety of tasks

# Start to engage in developing Speech Technology