

# 音声認識仮説の曖昧性を考慮する Multi-task End-to-End 音声翻訳

胡 尤佳   須藤 克仁   Sakriani Sakti   中村 哲

奈良先端科学技術大学院大学

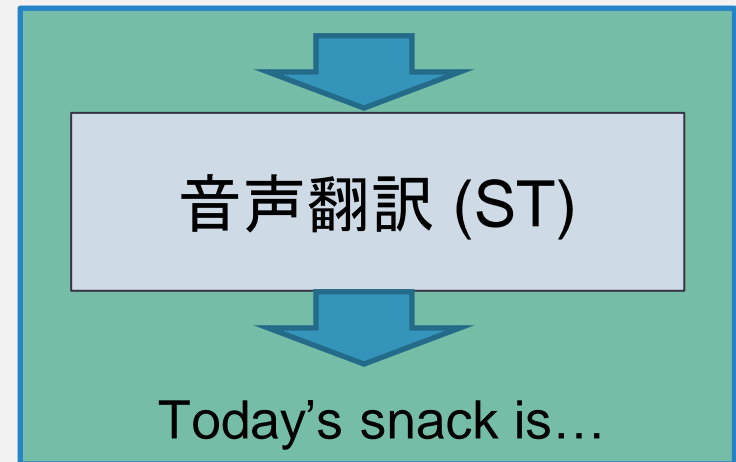
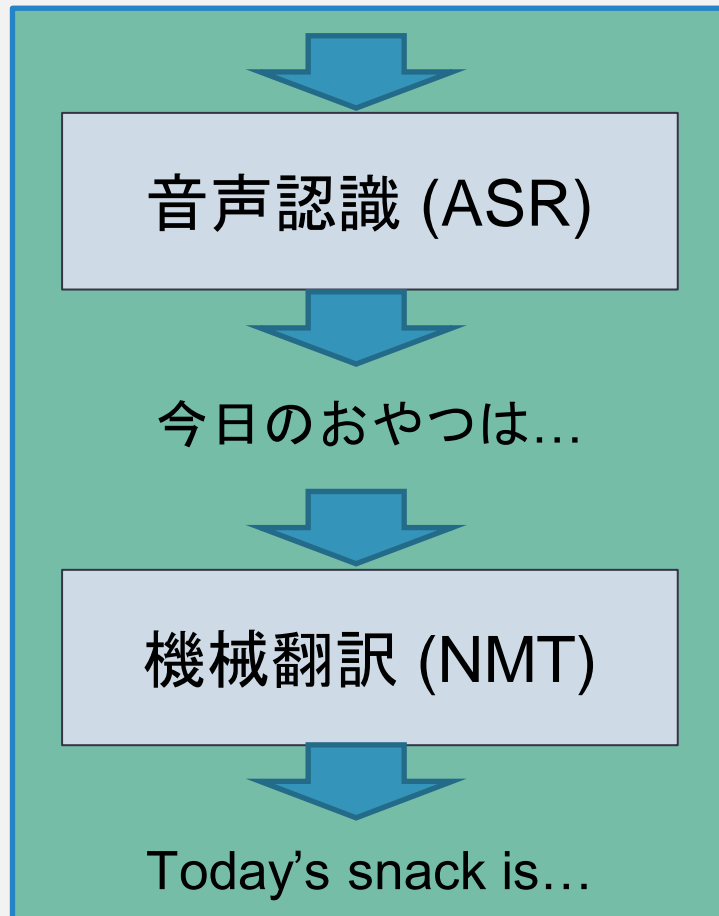
理化学研究所 革新知能統合研究センター AIP

2021/03/15-03/19 @言語処理学会第27回年次大会

## ● 音声翻訳 Speech Translation (ST)

➤ Cascade ST

➤ End-to-End ST



➤ Cascade

➤ 音声認識誤りの伝播

➤ End-to-End

➤ Single task だと学習が困難

➤ **Pre-train + Multi-task が必要**

- Cascade ST

- 音声認識誤りに対して頑健な機械翻訳が必要



- End-to-End ST

- 音声認識誤りの問題は本当になくなった？
  - Multi-task での学習中に存在すると仮定
  - 音声認識誤り, 曖昧性を考慮する必要

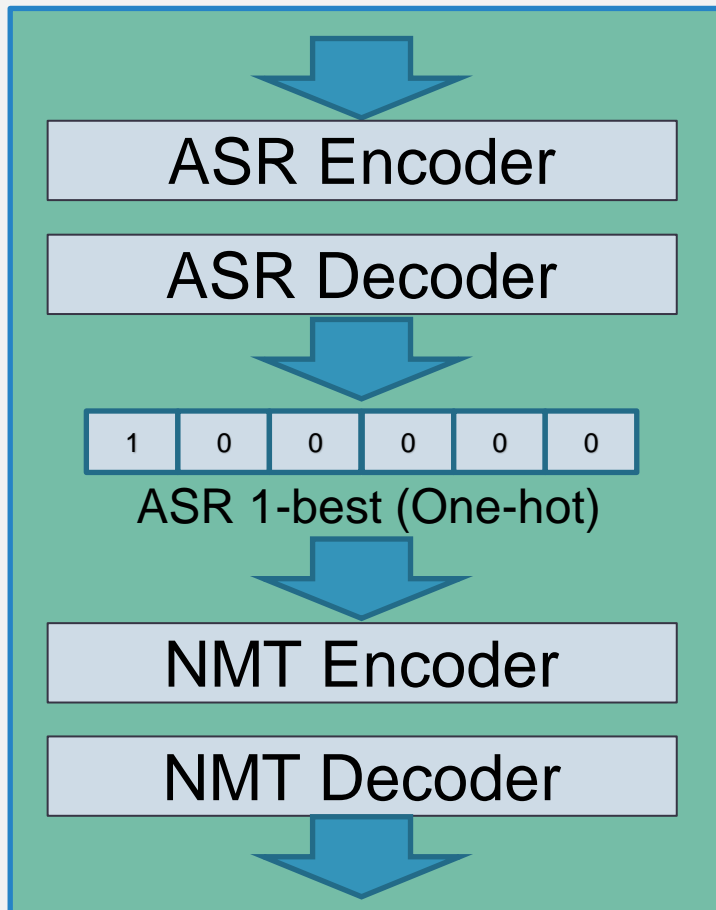
- 提案

- End-to-End 音声翻訳においても音声認識誤りに頑健なモデルの学習方法

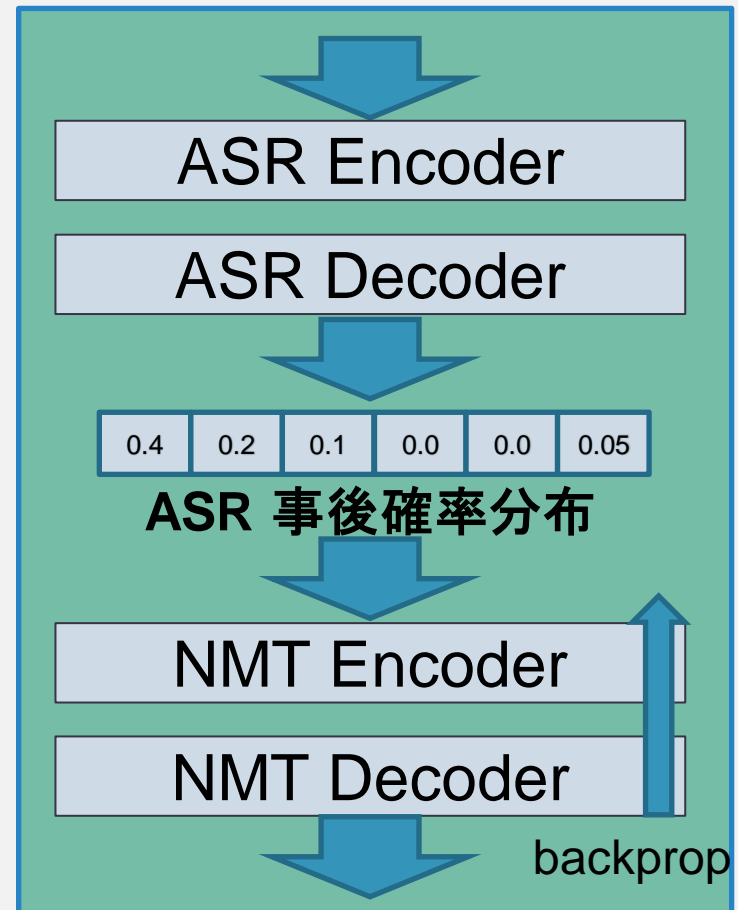
- 目的

- 音声翻訳の精度改善の期待

- [Osamura+, 2018] (cascade)
  - 音声認識出力 : 1-best → ASR 事後確率分布
  - 音声認識誤りに対して頑健な機械翻訳モデルの提案



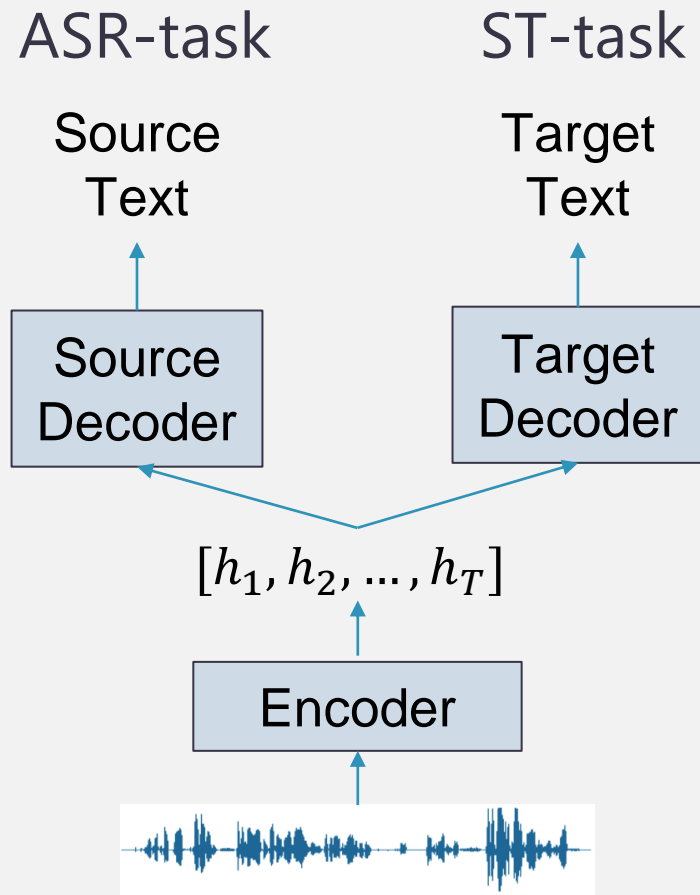
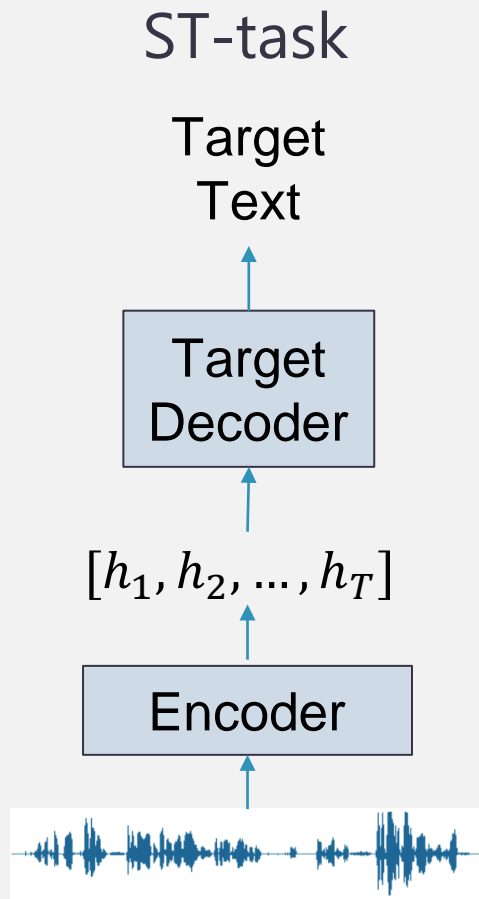
従来手法



Osamuraらの手法

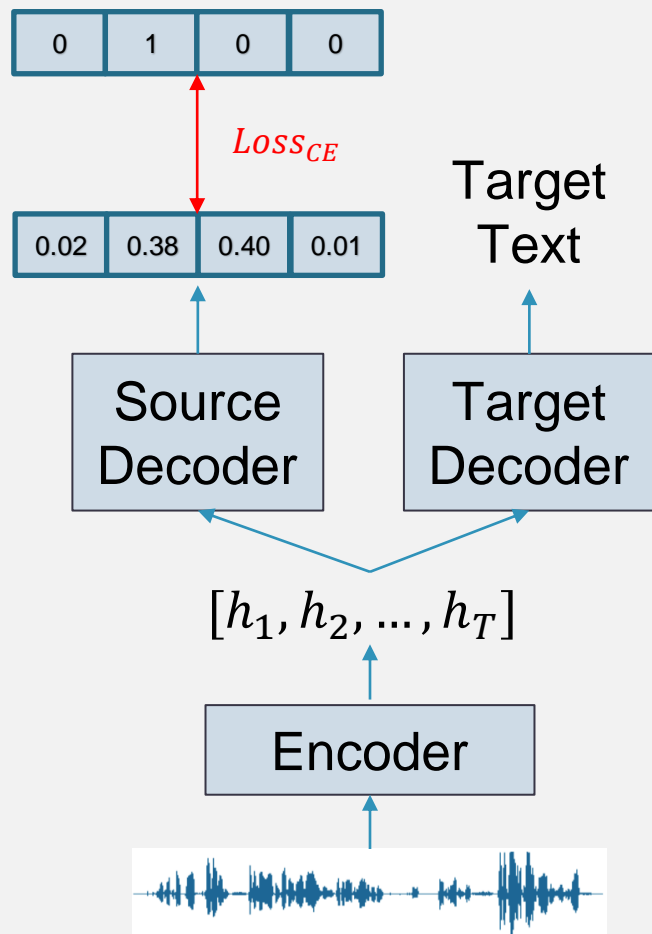
➤ Single-task End-to-End ST

➤ Multi-task End-to-End ST

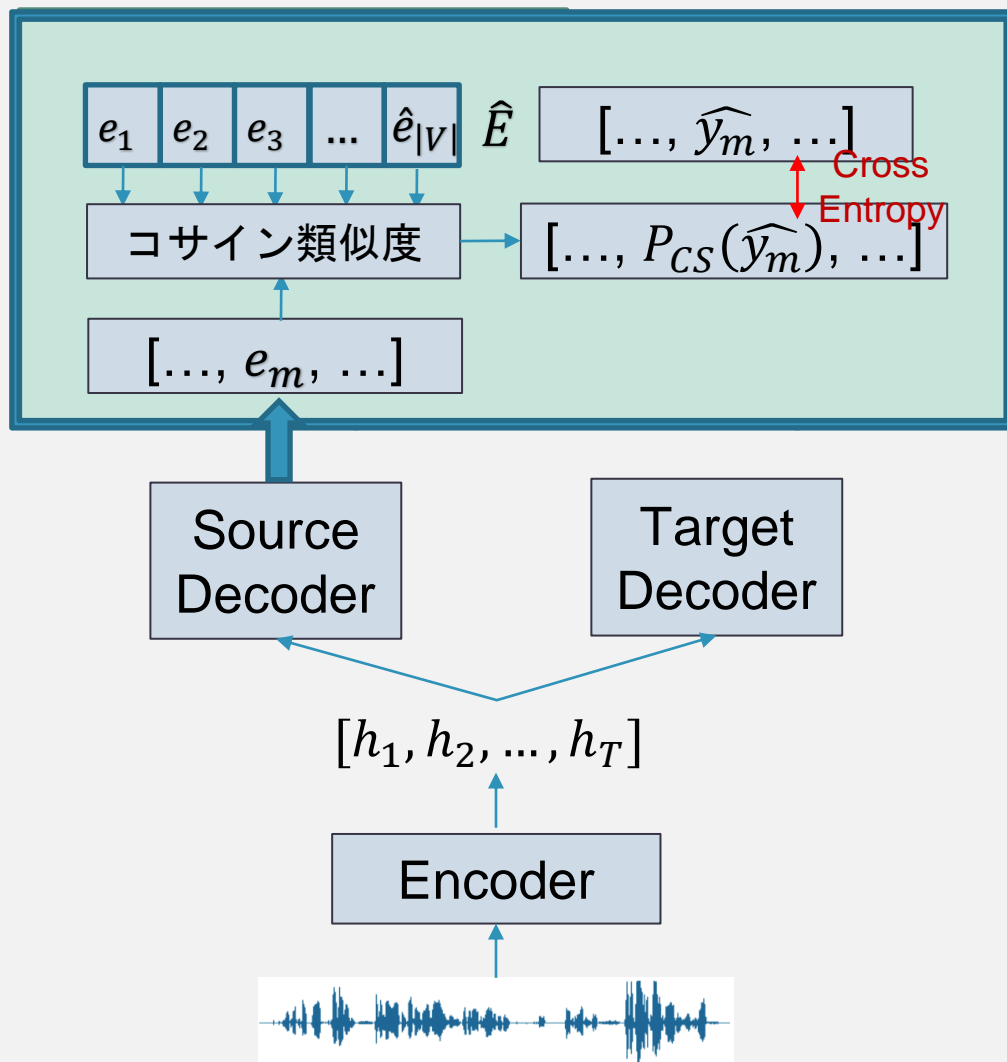


# 関連研究 3 : 単語分散表現の類似度を用いたST 6

- 一般的な Multi-task ST
  - 参照 : One-hot ベクトル



- [Chuang+, 2020] 参照, 予測間の意味的類似度を用いた誤差

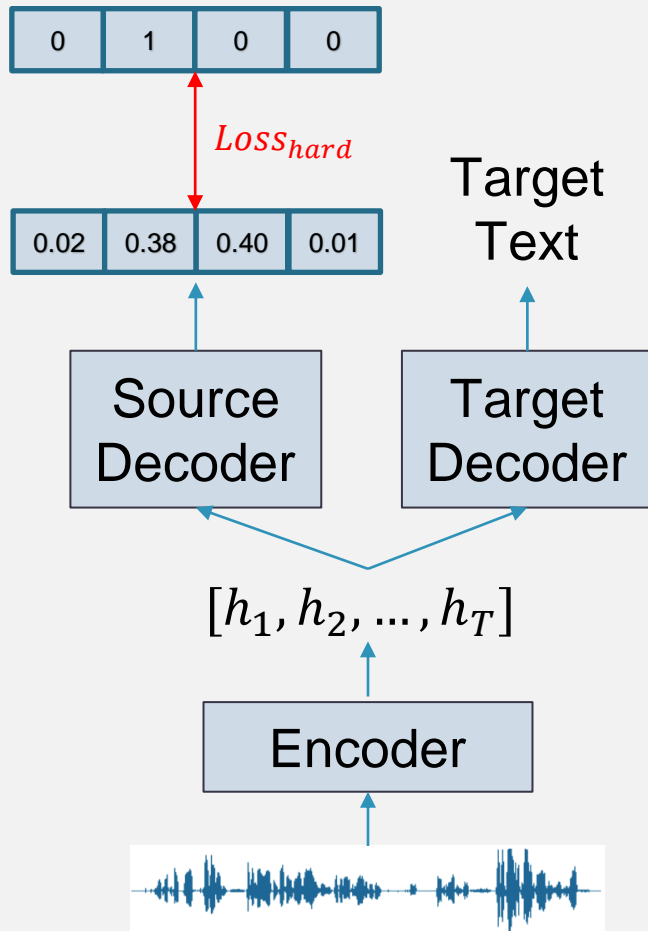


# 提案手法

# 従来手法と提案手法

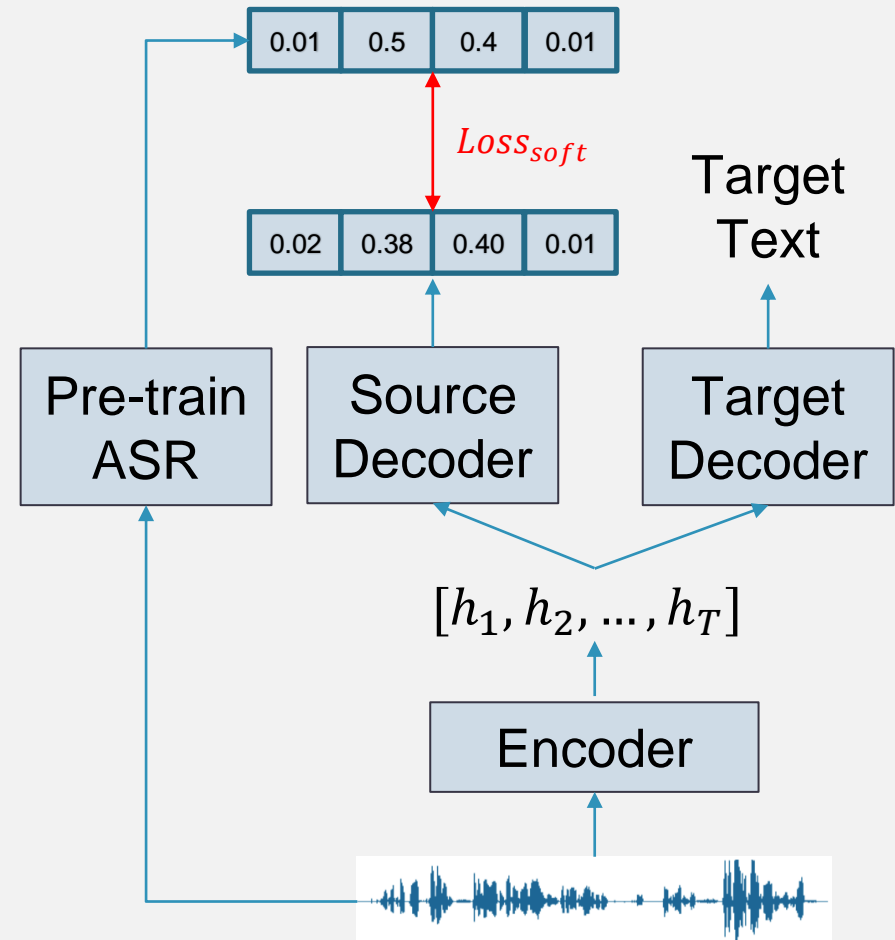
## ➤ 従来手法

- 参照：One-hot ベクトル
  - Hard target loss



## ➤ 提案手法

- 参照：ASR 事後確率分布
  - Soft target loss





# 従来手法と提案手法

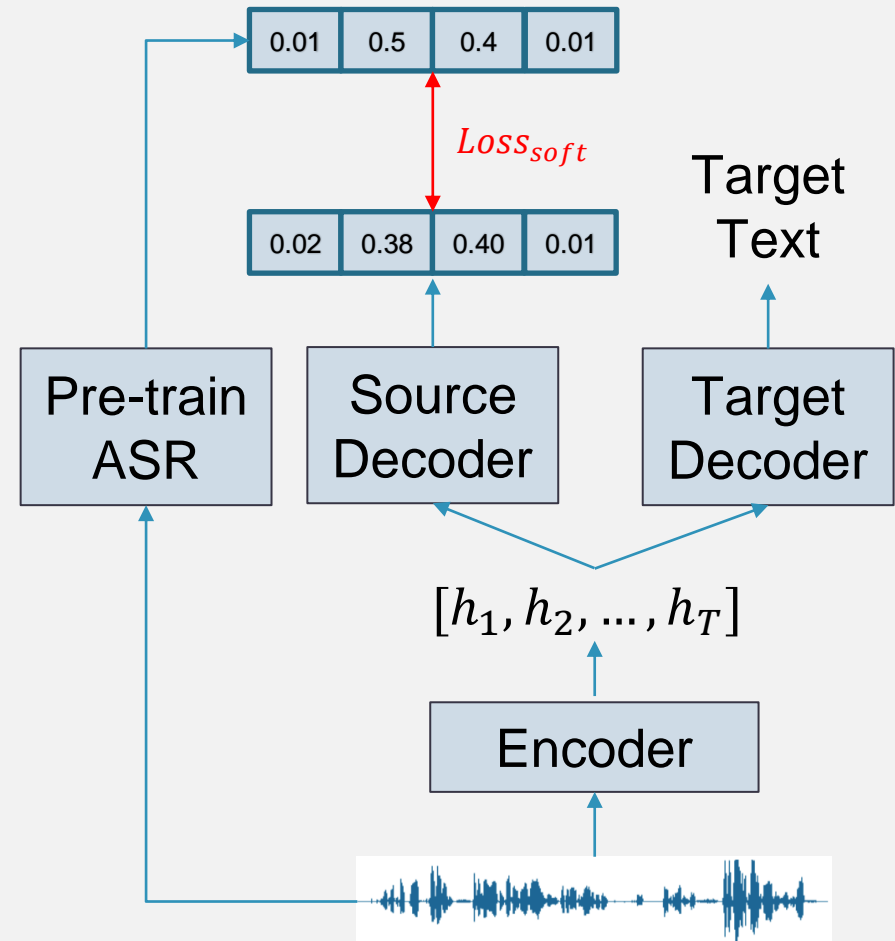
## ➤ 従来手法

- 参照：One-hot ベクトル
  - Hard target loss



## ➤ 提案手法

- 参照：ASR 事後確率分布
  - Soft target loss



# 実験

Data	src-tgt	Speech feature
Fisher Spanish CallHome コーパス	Es-En (Spanish- English)	fbank + pitch (80+3=83次元)

## Dictionary

SentencePiece 1000 Es-En Joint

	Dataset	Size
Train	fisher_train	415869 (138623 * 3)
Dev	fisher_dev	3973
	fisher_dev2	3957
Test	fisher_test	3638
	callhome_devtest	3956
	callhome_evlttest	1825

- 実装: ESPnet
- Pre-train ASR
  - Dev best (in accuracy)
- ST
  - $\lambda_{soft} = \{0, 0.3, 0.5, 0.7, 1.0\}, \lambda_{asr} = 0.3$ 
    - Baseline :  $\lambda_{soft} = 0.0$ 
      - 一般的なcross entropy loss / label smoothing loss
    - $LOSS_{ASR} = \lambda_{soft}LOSS_{soft} + (1 - \lambda_{soft})LOSS_{hard}$
    - $LOSS = \lambda_{ASR}LOSS_{ASR} + (1 - \lambda_{ASR})LOSS_{ST}$
  - Label smoothing weight (0.0 : cross entropy)
    - 実験1 : ASR-task 0.0 / ST-task 0.0 [論文に掲載]
    - 実験2 : ASR-task 0.0 / ST-task 0.1 [追加掲載]
    - 実験3 : ASR-task 0.1 / ST-task 0.1 [省略/付録]

- 実験に用いた Pre-train ASR モデルの WER
  - Dev best model in epoch 30

	WER
fisher_ dev	30.152
fisher_ dev2	29.124
fisher_ test	27.184
callhome_ devtest	49.138
callhome_ evltest	49.646

# 実験1 : ASR-task 0.0 / ST-task 0.0 14

- BLEU スコア
- ST-task : cross entropy / ASR-task : cross entropy
- 全体的なBLEUの向上が見られた (↓ : 低下 / **太字** : 最大)

		Baseline	Proposed			
	$\text{soft}\lambda_{\text{soft}} - \text{hard}(1 - \lambda_{\text{soft}})$	soft0.0-hard1.0	soft0.3-hard0.7	soft0.5-hard0.5	soft0.7-hard0.3	soft1.0-hard0.0
fisher dev	BLEU 4-ref	41.04	40.99 ↓	41.40	41.20	<b>41.51</b>
	BLEU 1-ref	23.97	23.88 ↓	24.12	24.00	<b>24.33</b>
fisher dev2	BLEU 4-ref	42.14	42.05 ↓	42.28	<b>42.45</b>	42.22
	BLEU 1-ref	25.17	25.23	25.22	25.30	<b>25.32</b>
fisher test	BLEU 4-ref	41.17	41.38	<b>41.41</b>	41.18	41.39
	BLEU 1-ref	24.77	25.02	24.93	24.82	<b>25.01</b>
callhome devtest	BLEU 1-ref	14.83	<b>15.23</b>	15.00	15.01	14.95
callhome evltest	BLEU 1-ref	14.81	<b>15.26</b>	15.10	14.78 ↓	15.09

# 実験1 : Fisher test $\lambda_{soft} = 0.5$

15

- Label smoothing weight = 0.0 (cross entropy loss)

	Example
Label	20051028_180633_356_fsp-A-016164-016487
Ground Truth (Es)	sí pero o sea sigue siendo bastante <b>intensi</b>
Ground Truth (En)	yes but it's still pretty <b>intensive</b>
Baseline (En)	yes but that keeps getting pretty <b>unthinkable</b> (->"inconceivable", "impensable" (Es))
Proposed (En)	yes but that keeps being pretty <b>intense</b>
Label	20051028_180633_356_fsp-A-033453-034134
Ground Truth (Es)	es es mejor en el sentido que uno okay que hay menos <b>riesgos</b> pero ay
Ground Truth (En)	that is the best in the sense that one okay that there are less <b>risks</b> but ay
Baseline (En)	it's it's better in the sense that you don't that there are less <b>colds</b> but there are (->resfriados (Es))
Proposed (En)	it's it's a best in the sense that you don't that there are less <b>risks</b> but

# 実験2 : ASR-task 0.0 / ST-task 0.1 16

- Label smoothing weight = 0.0
  - ASR-task の Hard loss はcross entropy loss
- 全体的なBLEUの向上が見られた (↓ : 低下 / **太字** : 最大)
  - Soft lossのみを用いると精度が低くなる傾向

		Baseline	Proposed			
	$\text{soft}\lambda_{\text{soft}}\text{-hard}(1 - \lambda_{\text{soft}})$	soft0.0-hard1.0	soft0.3-hard0.7	soft0.5-hard0.5	soft0.7-hard0.3	soft1.0-hard0.0
fisher dev	BLEU 4-ref	44.08	<b>44.91</b>	44.72	44.38	43.99 ↓
	BLEU 1-ref	25.68	25.69	<b>25.86</b>	25.58 ↓	25.39 ↓
fisher dev2	BLEU 4-ref	45.07	45.70	45.54	45.69	45.07
	BLEU 1-ref	27.00	27.13	27.01	27.33	26.99 ↓
fisher test	BLEU 4-ref	44.69	45.04	<b>45.29</b>	44.81	44.84
	BLEU 1-ref	26.78	26.73 ↓	<b>27.16</b>	26.61 ↓	26.55 ↓
callhome Devtest	BLEU 1-ref	15.83	16.09	16.28	15.96	<b>16.34</b>
callhome evltest	BLEU 1-ref	15.85	16.01	<b>16.80</b>	15.42 ↓	16.22



- 音声認識の事後確率分布による End-to-End ST の学習
  - 音声認識の曖昧性に対するロバスト性を期待
  - BLEUの向上が期待できることを示した
    - 全体的な BLEU の向上が見られた
  - Soft loss のみの場合，参照としての信用度が低く精度が下がる可能性
- 今後の展望
  - Pre-train ASR の性能ごとの信頼度の検証
  - テスト時の出力結果の分析
    - Main-task (ST) 出力とSub-task (ASR) 出力の照らし合わせ
    - Decoder 出力分布の確認
  - 発音，読み情報を利用した誤差計算 [Salesky+, 2020]

- [Osamura+, 2018] Osamura, Kaho, et al. "Using spoken word posterior features in neural machine translation." *architecture* 21 (2018): 22.
- [Anastasopoulos+, 2018] Anastasopoulos, Antonios, and David Chiang. "Tied multitask learning for neural speech translation." *arXiv preprint arXiv:1802.06655* (2018).
- [Chuang+, 2020] Chuang, Shun-Po, et al. "Worse WER, but Better BLEU? Leveraging Word Embedding as Intermediate in Multitask End-to-End Speech Translation." *arXiv preprint arXiv:2005.10678* (2020).
- [Watanabe+, 2018] "Espnet: End-to-end speech processing toolkit." *arXiv preprint arXiv:1804.00015* (2018).
- [Kikui+, 2003] "Creating corpora for speech-to-speech translation." *Eighth European Conference on Speech Communication and Technology*. 2003.
- [Salesky+, 2020] Elizabeth Salesky and Alan W. Black. Phone features improve speech translation. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel R. Tetreault, editors, Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020, pp. 2388–2397. Association for Computational Linguistics, 2020.

# Appendix

# 実験：学習時モデルパラメータ

20

	ASR	ST
epoch	30	
encoder layers	12	
encoder units	2048	
decoder layers	6	
decoder units	2048	
attention dimension	256	
attention heads	4	
batch size	64	
accum grad	2	4
gradient clipping	5	
transformer learning late	5	2.5
transformer warmup steps	25000	
decode beam size	1	4
label smoothing weight	0.1	{0, 0.1}
dropout	0.1	
model average	1	5

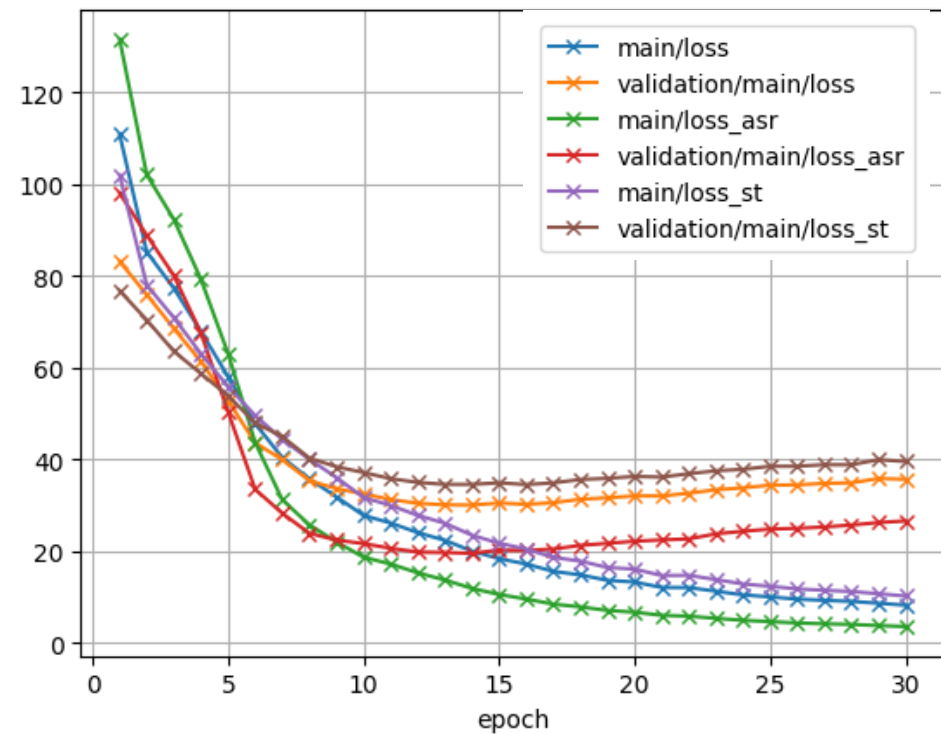
# 実験3 : ASR-task 0.1 / ST-task 0.1 21

- Label smoothing weight = 0.1 (ASR Hard lossに対して)
- 全体的なBLEUの向上が見られなかった(↓ : 低下 / **太字** : 最大)

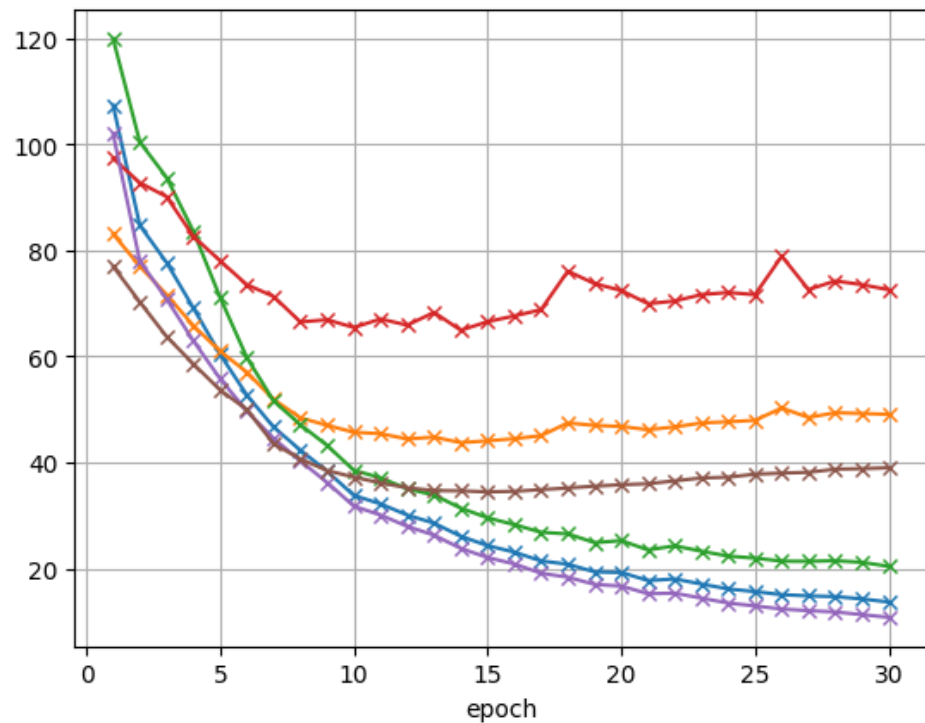
		Baseline	Proposed			
	$\text{soft}\lambda_{\text{soft}}-\text{hard}(1-\lambda_{\text{soft}})$	soft0.0-hard1.0	soft0.3-hard0.7	soft0.5-hard0.5	soft0.7-hard0.3	soft1.0-hard0.0
fisher dev	BLEU 4-ref	43.82	44.28	<b>44.47</b>	44.70	43.99
	BLEU 1-ref	25.64	<b>25.78</b>	25.45 ↓	25.72	25.39 ↓
fisher dev2	BLEU 4-ref	45.33	45.39	45.51	<b>46.25</b>	45.07 ↓
	BLEU 1-ref	27.11	27.16	27.22	<b>27.42</b>	26.99 ↓
fisher test	BLEU 4-ref	44.15	44.24	44.57	<b>45.04</b>	44.84
	BLEU 1-ref	26.72	26.45 ↓	<b>26.74</b>	26.72 ↓	26.55 ↓
callhome devtest	BLEU 1-ref	15.93	15.85 ↓	15.85 ↓	16.17	<b>16.34</b>
callhome evltest	BLEU 1-ref	<b>16.22</b>	15.81 ↓	16.06 ↓	15.52 ↓	16.22

- 実装：ESPnet
- Pre-train ASR model
  - transformer
  - Epoch : 100
  - Dev Best model
- ST model
  - transformer
  - Epoch : 100
  - Checkpoint averaging : 5
  - $W_{soft} = \{0, 0.3, 0.5, 0.7, 1.0\}$ ,  $W_{asr} = 0.3$ 
    - $loss_{asr} = W_{soft}loss_{soft} + (1 - W_{soft})loss_{hard}$
    - $loss = W_{asr}loss_{asr} + (1 - W_{asr})loss_{st}$

# loss\_asr soft0.0 vs soft1.0



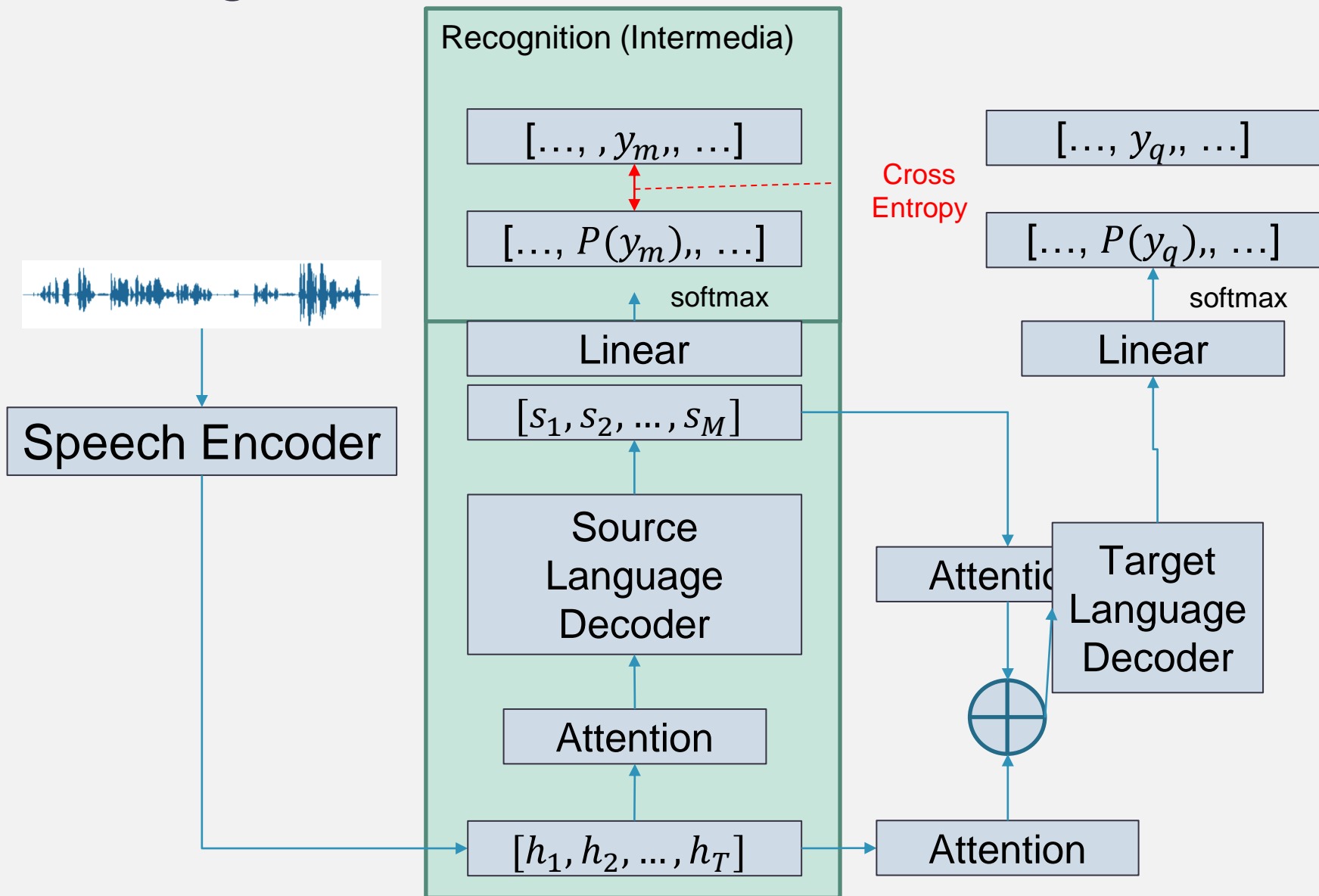
➤  $\text{soft\_weight} = 0.0$



➤  $\text{soft\_weight} = 1.0$

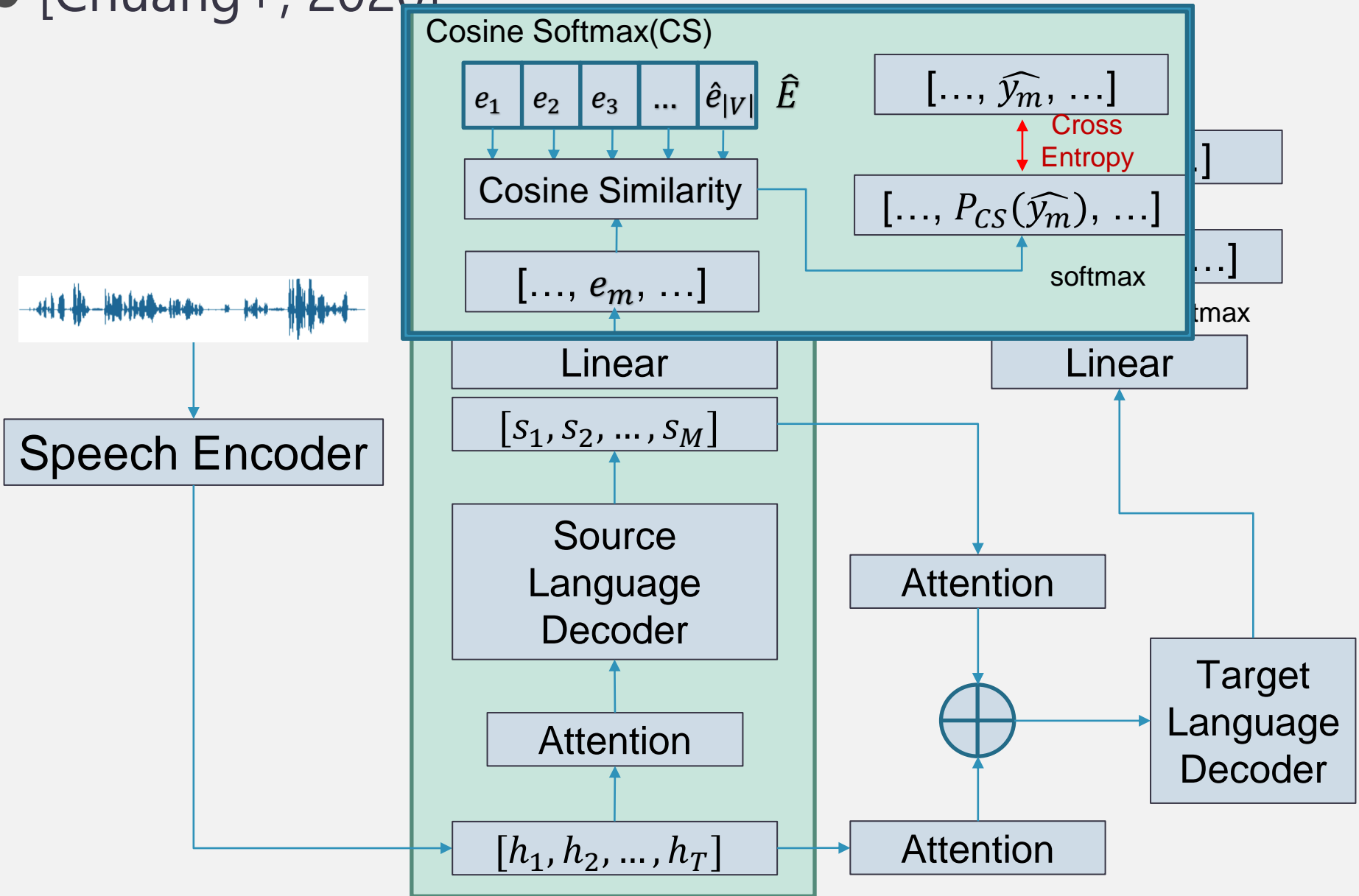
# 関連研究2：単語分散表現の類似度を用いたST<sup>2</sup><sub>4</sub>

- [Chuang+, 2020]





- [Chuang+, 2020]



data	src-tgt	speech feature
BTEC旅行会話コーパス	ja-en	fbank + pitch (80+3=83dim)

	BPE model	dict size
Japanese	SentencePiece 8000 BTEC1-4	7807
English	SentencePiece 8000 BTEC1-4	7769

- Maxframe=3000, maxchar=400, minchar=1
- Removed Punctuations

ASR	data	num
train	BTEC1-gtts BTEC natural for ASR	325498
dev	BTEC-test-set01	510 (/4080)
test	BTEC-test-set01	510 (/4080)

ST	data	num
train	BTEC1-gtts BTEC-test-set02,03	135361 (natural:8014)
dev	BTEC-test-set01	510 (/4080)
test	BTEC-test-set01	510 (/4080)

- Pretrain ASR WER : 15.864

$W_{soft}$	$W_{hard}$		BLEU
0.0	1.0	Baseline Cross Entropy	BLEU = 10.09, 27.7/13.0/7.2/4.0 (BP=1.000, ratio=1.176, hyp_len=3727, ref_len=3170)
0.3	0.7	Proposed	BLEU = 10.66, 27.8/13.7/7.9/4.3 (BP=1.000, ratio=1.182, hyp_len=3746, ref_len=3170)
0.5	0.5		BLEU = <b>11.52</b> , 29.4/14.7/8.3/4.9 (BP=1.000, ratio=1.167, hyp_len=3700, ref_len=3170)
0.7	0.3		BLEU = 10.58, 28.7/13.5/7.6/4.2 (BP=1.000, ratio=1.211, hyp_len=3838, ref_len=3170)
1.0	0.0		BLEU = 10.45, 28.6/13.6/7.7/4.0 (BP=1.000, ratio=1.185, hyp_len=3756, ref_len=3170)

# 実験：分析 $W_{soft} = 0.5, W_{hard} = 0.5$

	BLEU
Proposed	i am angry
Baseline	i 'm <b>changing</b>
Ground Truth Ja	後悔しています
Ground Truth En	i regret it
Pre-train ASR output	<b>交換</b> しています
Proposed	could you tell me how about tomorrow night
Baseline	the hot water won't stop until the end of <b>town</b>
Ground Truth Ja	マチネー マチネーってどういう意味かしら
Ground Truth En	matinee what does matinee mean
Pre-train ASR output	<b>街</b> ね間違えてどういう意味かしら
Proposed	excuse me where 's the tourist information office
Baseline	excuse me i 'm sorry but where 's the <b>tonight</b>
Ground Truth Ja	すみません本屋はどこですか
Ground Truth En	excuse me where 's the bookshop
Pre-train ASR output	すみません <b>本屋</b> はどこですか (<-> <b>今夜?</b> )

- Using posterior distribution in End-to-End ST Training
  - Robustness for ASR output ambiguity expected
  - BLEU improvement result, specially cross entropy loss
    - On label smoothing, not improved so much
- Future work
  - Analysis output distribution in test
    - Compare distribution between baseline and proposed
    - Main-task (ST) output & Sub-task (ASR) output
  - Using pronunciation information (phone [Salesky+, 2020]) in loss calculate
  - Adapt the method to simultaneous translation for ASR output ambiguity

- 音声認識の事後確率分布によりEnd-to-End STの学習
  - 音声認識の曖昧性に対するロバスト性を期待
  - BLEUの向上が期待できることを示した
- 今後の展望
  - 他コーパスによる検証
    - MuST-C, TED
  - Label smoothingを導入したLossでの検証
  - テスト時の出力結果の分析
    - Main-task (ST) 出力とSub-task (ASR) 出力の照らし合わせ
    - 定量評価
  - Ground Truthと比べたときのError Rateを重みとするLoss

- [Osamura+, 2018] Osamura, Kaho, et al. "Using spoken word posterior features in neural machine translation." *architecture* 21 (2018): 22.
- [Anastasopoulos+, 2018] Anastasopoulos, Antonios, and David Chiang. "Tied multitask learning for neural speech translation." *arXiv preprint arXiv:1802.06655* (2018).
- [Chuang+, 2020] Chuang, Shun-Po, et al. "Worse WER, but Better BLEU? Leveraging Word Embedding as Intermediate in Multitask End-to-End Speech Translation." *arXiv preprint arXiv:2005.10678* (2020).
- Watanabe, Shinji, et al. "Espnet: End-to-end speech processing toolkit." *arXiv preprint arXiv:1804.00015* (2018).
- Kikui, Genichiro, et al. "Creating corpora for speech-to-speech translation." *Eighth European Conference on Speech Communication and Technology*. 2003.