

Models for Trustworthy Speech Translation

using Tower 

HyoJung Han

Computer Science
University of Maryland



DEPARTMENT OF
COMPUTER SCIENCE

The rapid development of speech technology has expanded the use of speech translation (ST) applications in daily life.

The rapid development of speech technology has expanded the use of speech translation (ST) applications in daily life.

→ The needs to predict the reliability of their output is increasing.

Models for Trustworthy Speech Translation



DEPARTMENT OF
COMPUTER SCIENCE



JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

SpeechQE: Estimating the Quality of Direct Speech Translation

EMNLP2024 



HyoJung Han
Computer Science
University of Maryland



Kevin Duh
HLTCOE
Johns Hopkins University



Marine Carpuat
Computer Science
University of Maryland



Está dentro de cada mirada...

It is inside
street look...



Speech
Translation
(ST) System



Está dentro de
cada mirada...
(It is within every
glance...)



It is inside
street look...



Speech
Translation
(ST) System



Está dentro de
cada mirada...
(It is within every
glance...)


How good is
this speech
translation?



Estimating Quality (QE) of Translation


Assessing the quality of translation is crucial as it help people rely on MT appropriately without reference.

Many works on QE*:

 [google/metricx-23-qe-xx1-v2p0](#)

 [Unbabel/wmt20-comet-qe-da](#)

 GEMBA


 [Unbabel/wmt22-cometkiwi-da](#)


 [Unbabel/XCOMET-XL](#)

 [facebook/blaser-2.0-qe](#)

 MaTese

 MS-Comet

 [Unbabel/wmt23-cometkiwi-da-x1](#)

 [google/metricx-23-qe-x1-v2p0](#)

 UniTE

 [Unbabel/wmt21-comet-qe-mqm-marian](#)

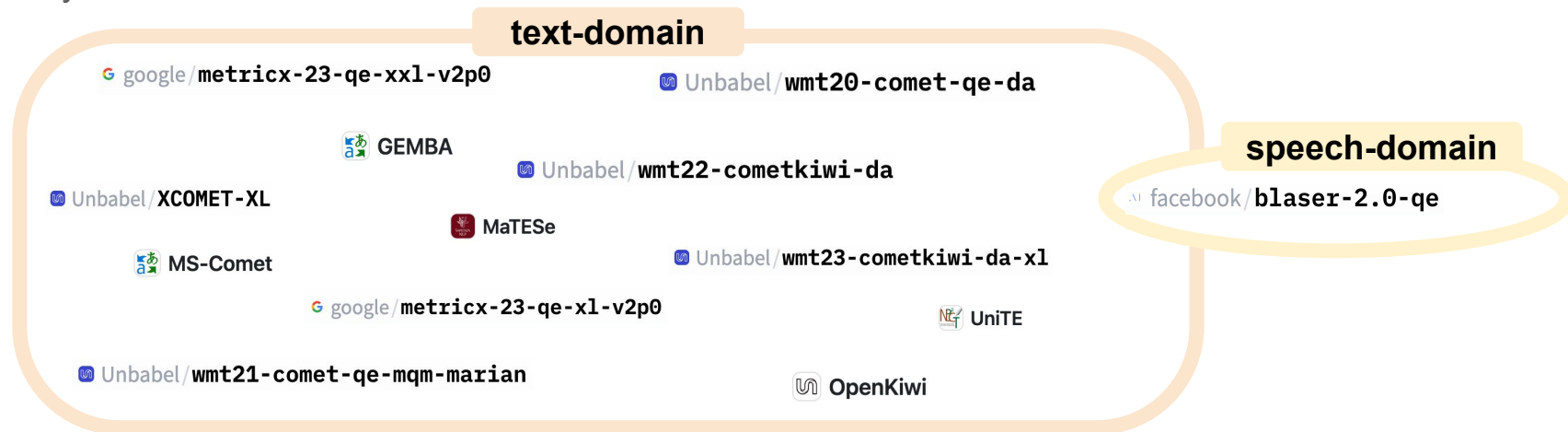
 OpenKiwi

* Displayed is a subset of the QE works, which includes a runnable model.

QE in Speech domain is underexplored

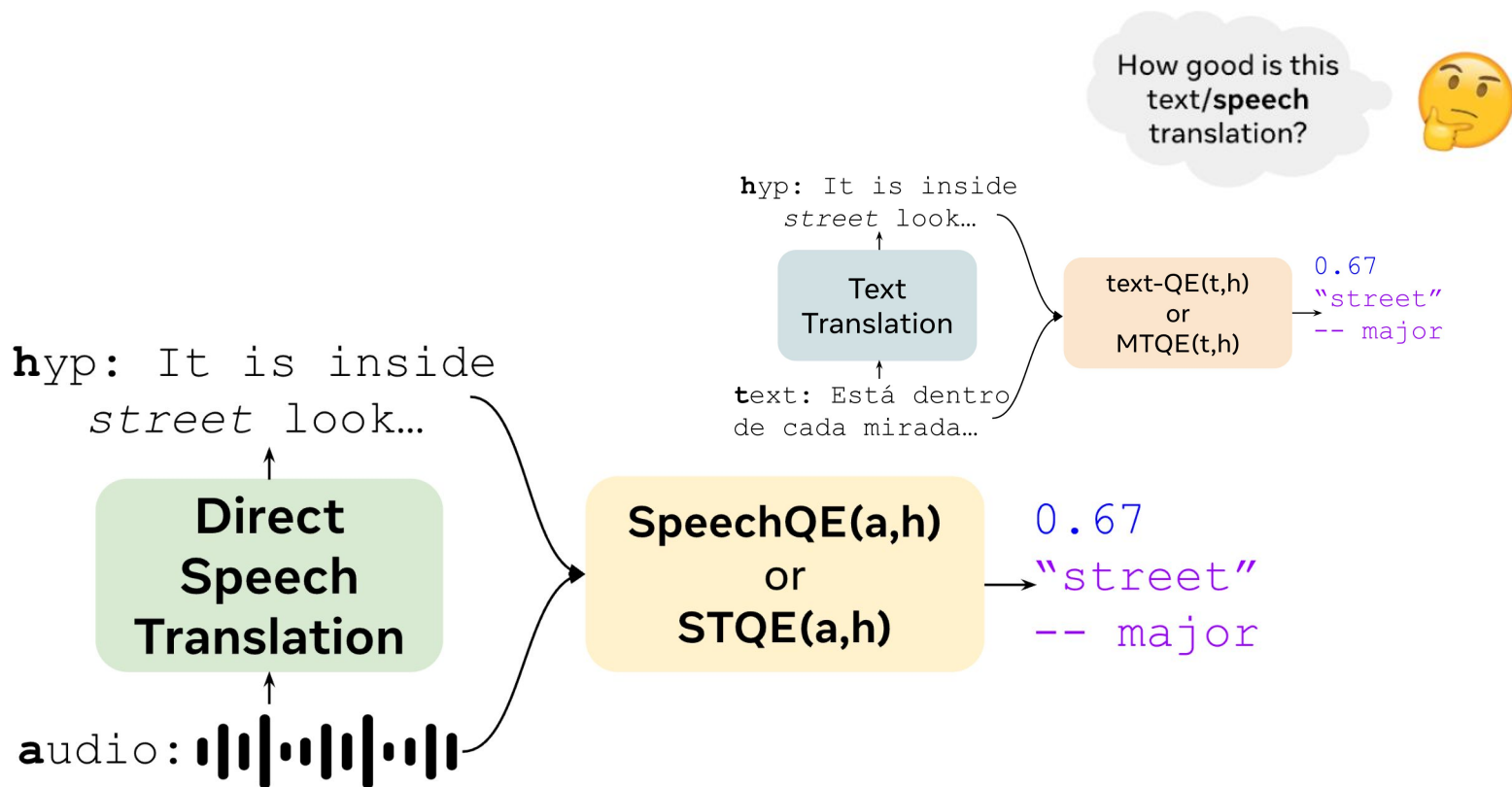
The rapid development of speech technology has expanded the use of speech translation (ST) applications in daily life, thus increasing the need to predict the reliability of their output.

Many works on QE*:



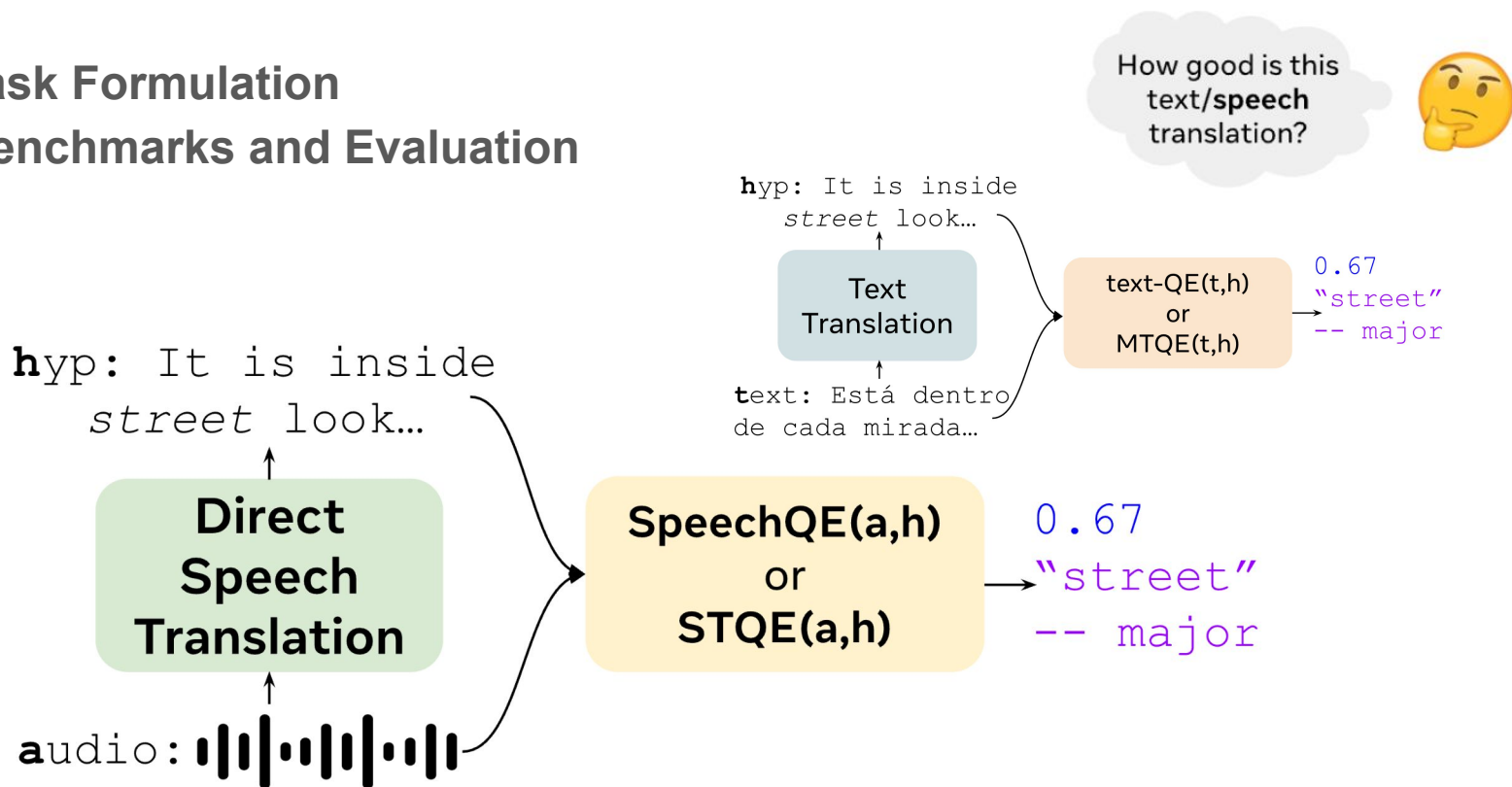
* Displayed is a subset of the QE works. The overall trends remain consistent.

SpeechQE: Estimating Quality of Direct Speech Translation

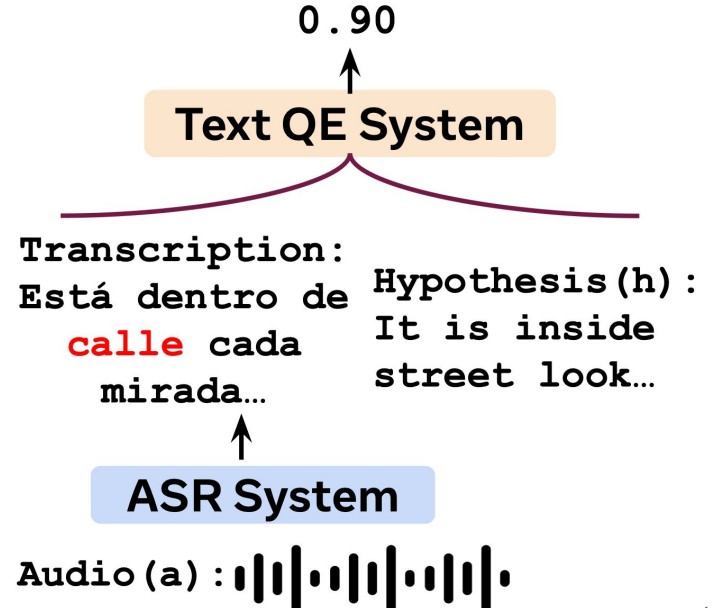


SpeechQE: Estimating Quality of Direct Speech Translation

1. Task Formulation
2. Benchmarks and Evaluation



Cascaded SpeechQE System



Potential Issues of Cascaded SpeechQE System

1. Efficiency

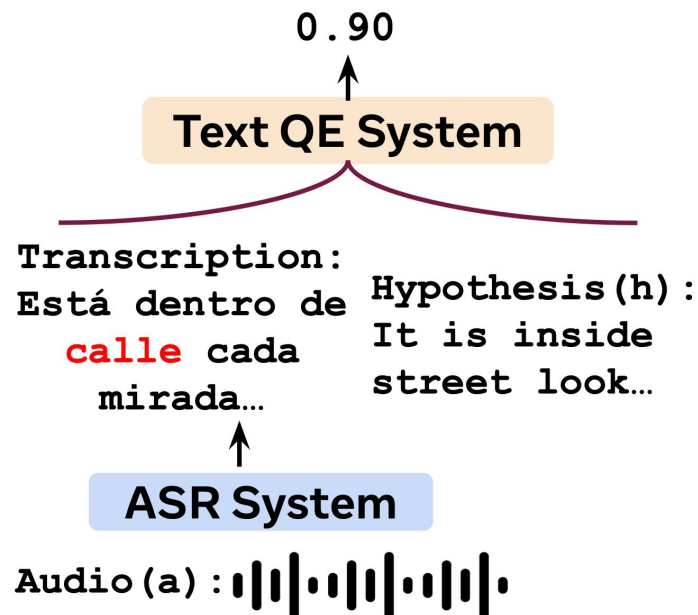
no (naturally occurring)
intermediate ASR transcripts

2. Wrong speech

representation by ASR

3. Modality mismatch

text-QE is not adapted to
spoken language



SpeechQE: Estimating Quality of Direct Speech Translation

1. Task Formulation
2. Benchmarks and Evaluation
3. Explore both cascaded and end-to-end (E2E) systems

End-to-End SpeechQE System

0.67 "street" -- major

LLM
(TowerInstruct) 

Text Embeddings

Instruction:
Estimate the
quality of ...
Identify errors ...

Modality Adapter

Speech Encoder 


a: 

Text Embeddings

Hypothesis (h) :
It is inside
street look...

Text QE System

0.90

How good is this
text/speech
translation? 

Transcription:

Está dentro de
calle cada
mirada...

Hypothesis (h) :
It is inside
street look...

ASR System

Audio (a) : 

Preliminaries

h: hypothesis text

r: reference text

t: source text

“*Metric*”: reference-based metric. e.g. BLEU, chrF, xCOMET, MetricX

score m : $m = \text{metric}(h, r)$ or $m = \text{metric}(t, h, r)$

Preliminaries

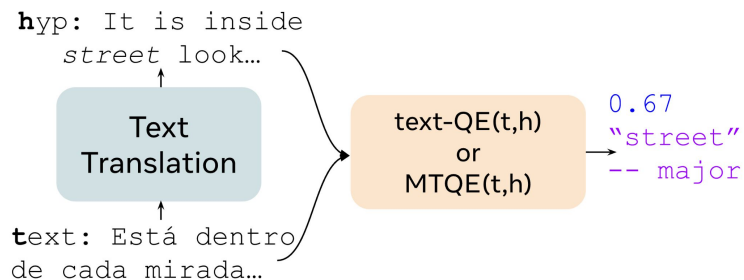
h: hypothesis text
r: reference text
t: source text

“*Metric*”: reference-based metric. e.g. BLEU, chrF, xCOMET, MetricX

score m : $m = \text{metric}(h, r)$ or $m = \text{metric}(t, h, r)$

“text-QE”: text quality estimation system. e.g. COMETkiwi, MetricX-QE

score q : $q = \text{text-QE}(t, h)$



* Among various ways in framing the QE task, we mainly focus on sentence-level quality rating QE as Sentence-level quality rating instead of word-level. We also experiment with Error Span Detection for a border QE scope in the paper.

SpeechQE Task Formulation

h: hypothesis text

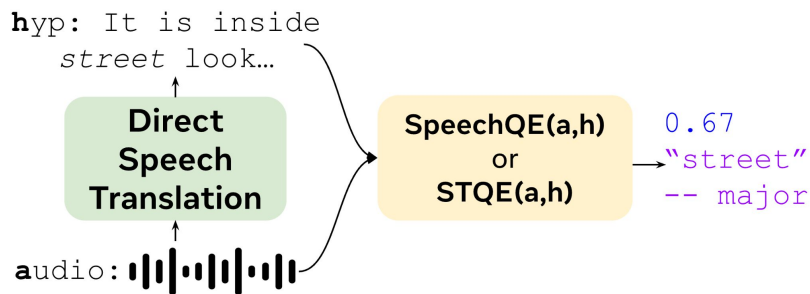
r: reference text

t: source text

a: source audio

“SpeechQE”: now the source input **a** is audio

$$\text{score } q: q = \text{SpeechQE}(a, h)$$

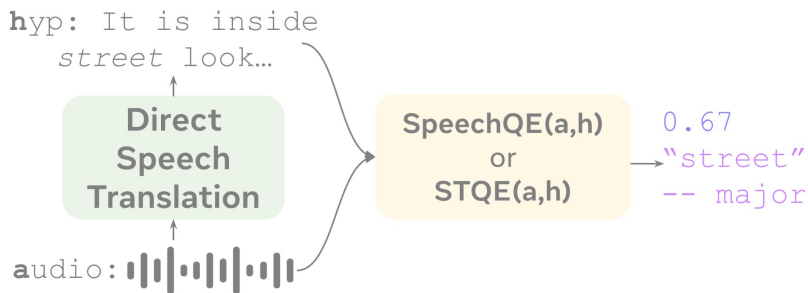


SpeechQE Task Formulation

h: hypothesis text
 r: reference text
 t: source text
 a: source audio

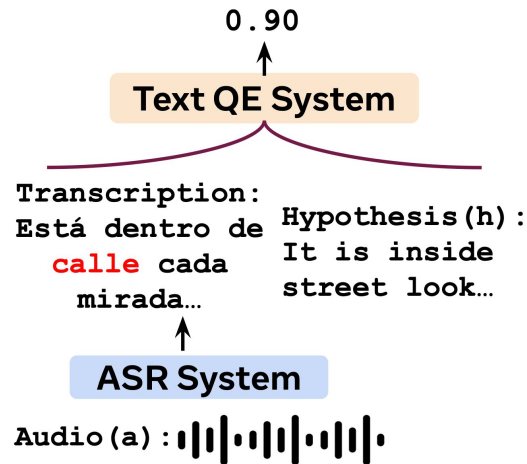
“SpeechQE”: now the source input **a** is audio

$$\text{score } q: q = \text{SpeechQE}(a, h)$$



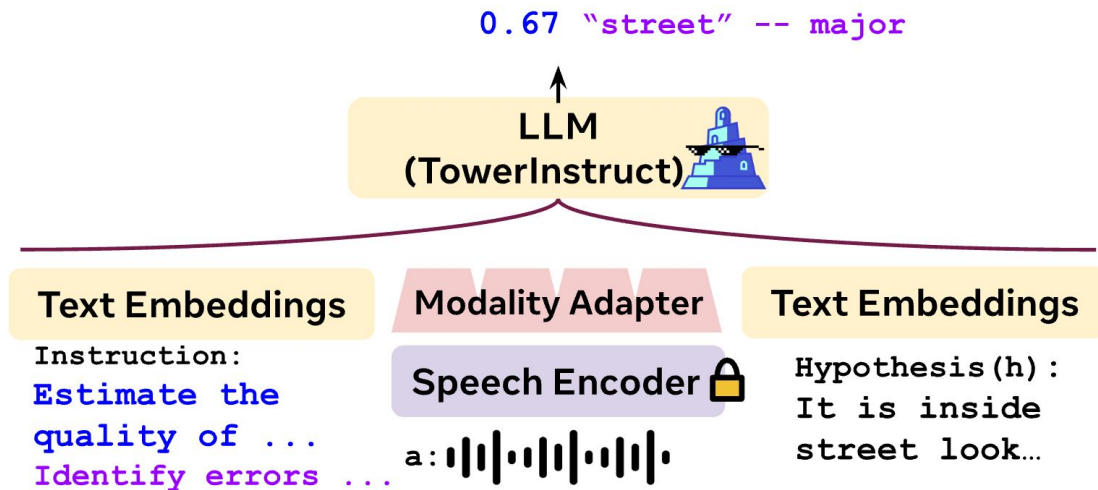
Cascaded SpeechQE System: outputs the score from a text-based QE system with the input of transcribed text ASR(a)

$$q_{cas} = \text{text-QE}(\text{ASR}(a), h)$$



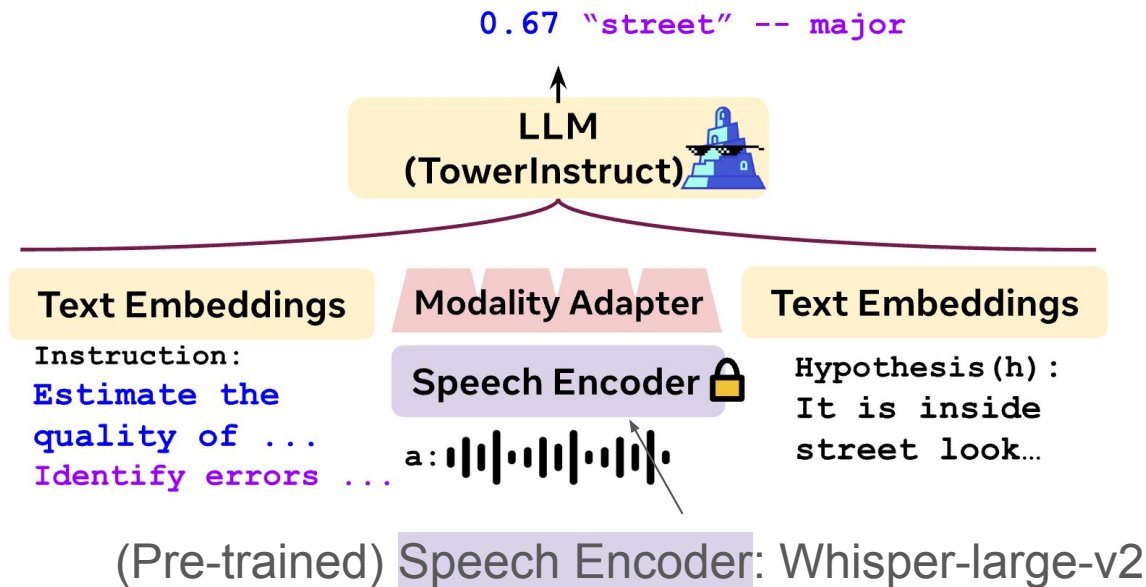
Building End-to-End SpeechQE System

Incorporating a pre-trained speech encoder and a large language model (LLM)



Building End-to-End SpeechQE System

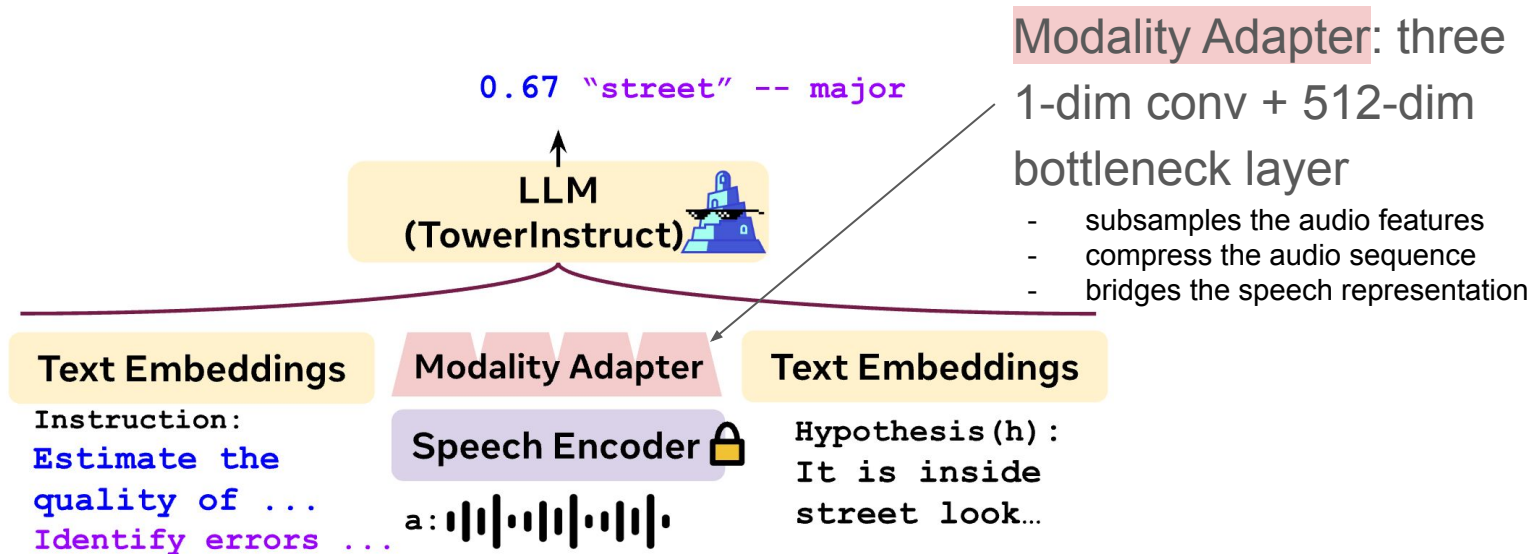
Incorporating a pre-trained speech encoder and a large language model (LLM)



- extracts the audio feature from the raw audio.

Building End-to-End SpeechQE System

Incorporating a pre-trained speech encoder and a large language model (LLM)



Building End-to-End SpeechQE System

Incorporating a pre-trained speech encoder and a large language model (LLM)

(Pre-trained) text-LLM:

TowerInstruct-7B

- Input: text + audio embedding sequence.

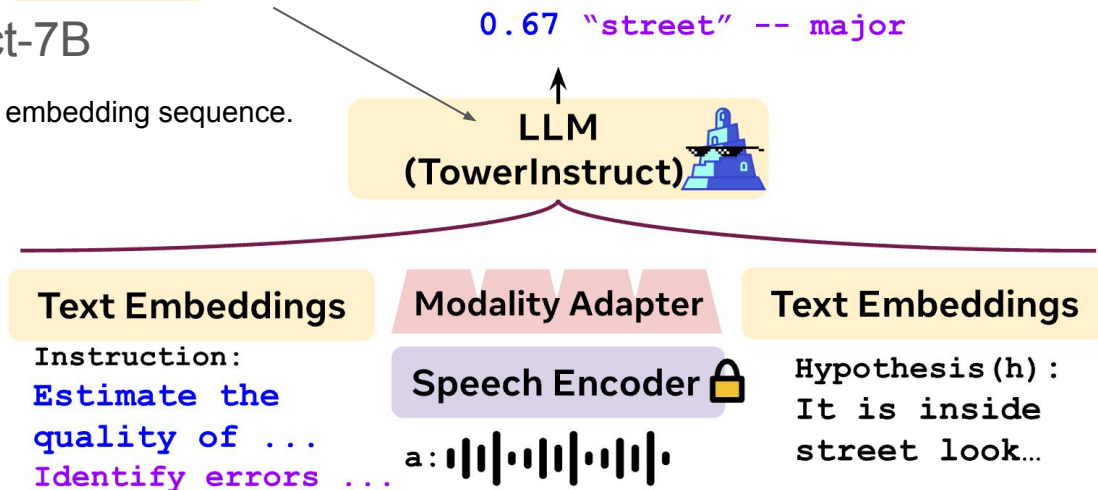


Table 10: Example of English-to-German speech translation and quality estimations of SpeechQE systems. Both cascaded and E2E SpeechQE systems could detect errors. However, the cascaded system estimates the severity lower than that of the metric labels partly due to ASR error while E2E could estimate the quality closely to labels.

```
# QE4ST task, training and testing
Given the German translation of the speech, estimate the quality of the translation as a
score between 0 to 1.
English: [[audio input]]
German translation: Wir modellieren den grasweisen, obstruktiven Summize-Ansatz mit zwei
verschiedenen Methoden.
# desired output in training or example output in testing
0.851
```

Training End-to-End SpeechQE System

SpeechQE {source audio (a), hypotheses (h), ratings (m)} + ASR + ST task

 : Update weights

 : Freeze weights

 : LoRA Fine-tuning

Training End-to-End SpeechQE System

SpeechQE {source audio (a), hypotheses (h), ratings (m)} + ASR + ST task

 : Update weights

 : Freeze weights

 : LoRA Fine-tuning

Training End-to-End SpeechQE System

SpeechQE {source audio (a), hypotheses (h), ratings (m)} + ASR + ST task

Single-Phase Training

Two-Phase Training

 : Update weights

 : Freeze weights

 : LoRA Fine-tuning

🔥 : Updated weights

❄️ : Frozen weights

👁️: LoRA Fine-tuned

Training End-to-End SpeechQE System

SpeechQE {source audio (a), hypotheses (h), ratings (m)} + ASR + ST task



(Speech Encoder  is always frozen)

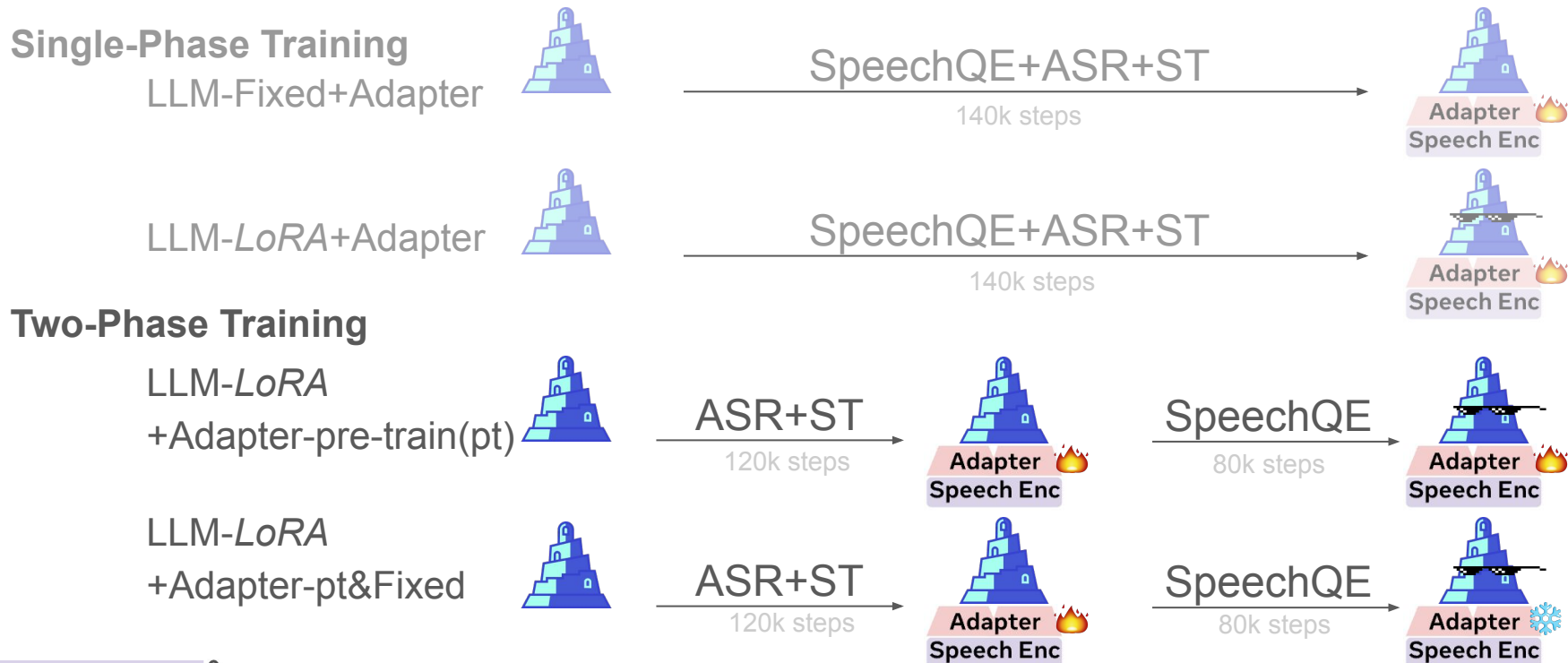
Training End-to-End SpeechQE System

 : Updated weights

 : Frozen weights

 : LoRA Fine-tuned

SpeechQE {source audio (a), hypotheses (h), ratings (m)} + ASR + ST task



(Speech Encoder  is always frozen)

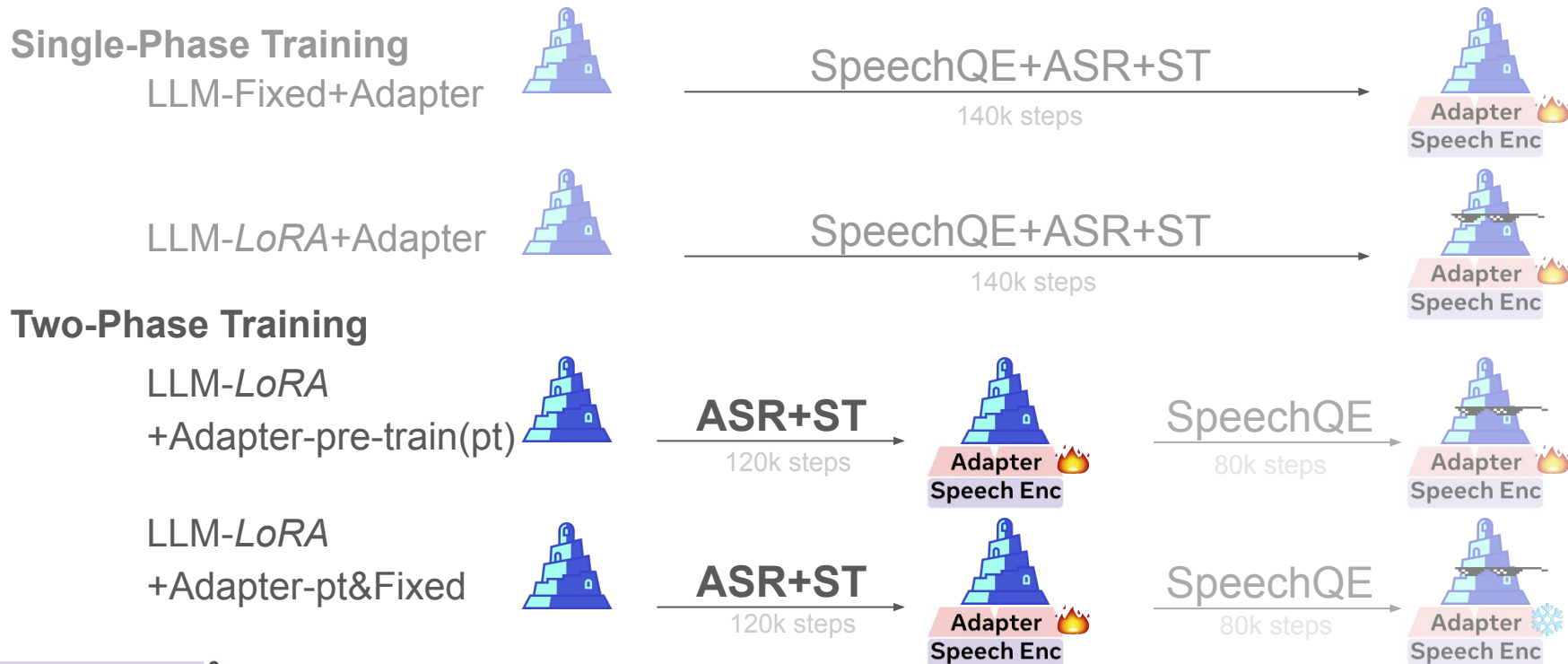
🔥 : Updated weights

❄️ : Frozen weights

🔪 : LoRA Fine-tuned

Training End-to-End SpeechQE System

SpeechQE {source audio (a), hypotheses (h), ratings (m)} + ASR + ST task



(Speech Encoder  is always frozen)

Training End-to-End SpeechQE System

🔥 : Updated weights

❄️ : Frozen weights

🔪 : LoRA Fine-tuned

SpeechQE {source audio (a), hypotheses (h), ratings (m)} + ASR + ST task



(Speech Encoder  is always frozen)

Building SpeechQE Benchmark and Training Corpus

1. Subsampled about 80k segments from the training set and 500 from the dev and test of CoVoST2
2. Run ST models to get hypothesis.
3. Get automatic quality labels from (reference-based) metrics xCOMET-XL for train, adding MetricX-23-XL for test

CoVoST2	es2en	en2de	es2en	en2de
	<i>Train</i>		<i>Dev / Test</i>	
ASR	297k	305k		
ST	79k	290k		
SpeechQE	546k	589k	3.5k	3.5k

Es2En direct ST systems	CoVoST2 BLEU
whisper-large-v3	39.05
whisper-large-v2	39.53
whisper-large	38.11
whisper-medium	37.39
whisper-small	31.27
whisper-base	16.93
whisper-tiny	7.81
En2De direct ST systems	CoVoST2 BLEU
seamless-m4t-v2-large	43.12
seamless-m4t-large	40.55
seamless-m4t-medium	38.39
s2t-wav2vec2-large-en-de	26.98
s2t-medium-mustc-multilingual-st	8.08
s2t-small-mustc-en-de-st	7.82
s2t-small-covost2-en-de-st	14.19

Table 2: The list of seven direct ST models

Building SpeechQE Benchmark and Training Corpus

1. **Subsampled** about **80k** segments from the **training** set and **500** from the **dev** and **test** of CoVoST2
2. Run ST models to get hypothesis.
3. Get automatic quality labels from (reference-based) metrics xCOMET-XL for train, adding MetricX-23-XL for test

CoVoST2	es2en	en2de	es2en	en2de
	<i>Train</i>		<i>Dev / Test</i>	
ASR	297k	305k		
ST	79k	290k		
SpeechQE	546k	589k	3.5k	3.5k

Es2En direct ST systems	CoVoST2 BLEU
whisper-large-v3	39.05
whisper-large-v2	39.53
whisper-large	38.11
whisper-medium	37.39
whisper-small	31.27
whisper-base	16.93
whisper-tiny	7.81
En2De direct ST systems	CoVoST2 BLEU
seamless-m4t-v2-large	43.12
seamless-m4t-large	40.55
seamless-m4t-medium	38.39
s2t-wav2vec2-large-en-de	26.98
s2t-medium-mustc-multilingual-st	8.08
s2t-small-mustc-en-de-st	7.82
s2t-small-covost2-en-de-st	14.19

Table 2: The list of seven direct ST models

Building SpeechQE Benchmark and Training Corpus

1. Subsampled about 80k segments from the training set and 500 from the dev and test of CoVoST2
2. **Run ST models** to get hypothesis.
3. Get automatic quality labels from (reference-based) metrics xCOMET-XL for train, adding MetricX-23-XL for test

CoVoST2	es2en	en2de	es2en	en2de
	<i>Train</i>		<i>Dev / Test</i>	
ASR	297k	305k		
ST	79k	290k		
SpeechQE	546k	589k	3.5k	3.5k

Es2En direct ST systems	CoVoST2 BLEU
whisper-large-v3	39.05
whisper-large-v2	39.53
whisper-large	38.11
whisper-medium	37.39
whisper-small	31.27
whisper-base	16.93
whisper-tiny	7.81
En2De direct ST systems	CoVoST2 BLEU
seamless-m4t-v2-large	43.12
seamless-m4t-large	40.55
seamless-m4t-medium	38.39
s2t-wav2vec2-large-en-de	26.98
s2t-medium-mustc-multilingual-st	8.08
s2t-small-mustc-en-de-st	7.82
s2t-small-covost2-en-de-st	14.19

Table 2: The list of seven direct ST models

Building SpeechQE Benchmark and Training Corpus

1. Subsampled about 80k segments from the training set and 500 from the dev and test of CoVoST2
2. Run ST models to get hypothesis.
3. **Get automatic quality labels** from (reference-based) **metrics** **xCOMET-XL for train**, adding **MetricX-23-XL for test**

CoVoST2	es2en	en2de	es2en	en2de
	<i>Train</i>		<i>Dev / Test</i>	
ASR	297k	305k		
ST	79k	290k		
SpeechQE	546k	589k	3.5k	3.5k

Metric		avg corr
XCOMET-Ensemble	1	0.825
XCOMET-QE-Ensemble*	2	0.808
MetricX-23	2	0.808
MetricX-23-QE*	2	0.800
<u>docWMT22CometDA</u>	4	0.768
<u>docWMT22CometKiwiDA*</u>	4	0.767
Calibri-COMET22	4	0.767
Calibri-COMET22-QE*	4	0.755

Results of WMT23 Metrics Shared Task: Metrics Might Be Guilty but References Are Not Innocent (Freitag et al., WMT 2023)

metric		avg corr
MetaMetrics-MT	1	0.725
MetricX-24-Hybrid	1	0.721
XCOMET	1	0.719
MetricX-24-Hybrid-QE*	2	0.714
gemba_esa*	2	0.711
XCOMET-QE*	3	0.695

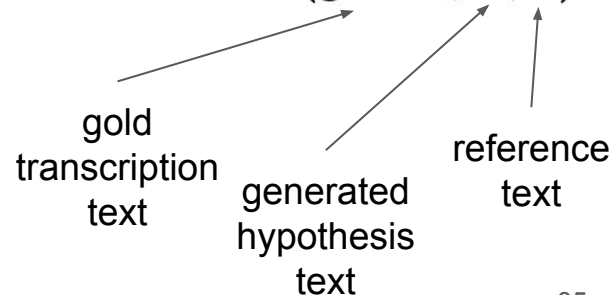
Are LLMs Breaking MT Metrics? Results of the WMT24 Metrics Shared Task (Freitag et al., WMT 2024)

Building SpeechQE Benchmark and Training Corpus

1. Subsampled about 80k segments from the training set and 500 from the dev and test of CoVoST2
2. Run ST models to get hypothesis.
3. **Get automatic quality labels** from (reference-based) **metrics** **xCOMET-XL for train**, adding **MetricX-23-XL for test**

CoVoST2	es2en	en2de	es2en	en2de
	<i>Train</i>		<i>Dev / Test</i>	
ASR	297k	305k		
ST	79k	290k		
SpeechQE	546k	589k	3.5k	3.5k

$$m_{xCOMET} = xCOMET(\text{gold } t, h, r)$$



Evaluation

We compute the correlations between SpeechQE scores (q) and :

1. m: Metric score of xCOMET-XL & MetricX-23-XL on En2En/En2De SpeechQE test (3.5k)
2. d: Human Direct Assessment score on IWSLT23-ACL En2De Speech Translation
 - source-based DA ratings of 416 hypotheses from each of the ten ST systems(4.1k)

$\rho = \text{corr}(\mathbf{q}, \mathbf{m} \text{ or } \mathbf{d})$	CoVoST2 Es2En		CoVoST2 En2De		IWSLT23
	$\mathbf{m}_{\text{xCOMET}}$	$\mathbf{m}_{\text{MetricX}}$	$\mathbf{m}_{\text{xCOMET}}$	$\mathbf{m}_{\text{MetricX}}$	En2De \mathbf{d}

Cascaded SpeechQE Systems Correlations

$$q_{cas} = \text{xCOMET-qe}(\text{ASR}(a), h)$$

$$q_{cas} = \text{MetricX-qe}(\text{ASR}(a), h)$$

$$q_{cas} = \text{text-BLASER2.0-qe}(\text{ASR}(a), h)$$

End-to-End SpeechQE Systems Correlations

$$q_{e2e} = \text{BLASER2.0-qe}(a, h)$$

$$q_{e2e} = \text{TowerInstruct-Fixed+Adapter}(a, h)$$

$$q_{e2e} = \text{TowerInstruct-LoRA+Adapter}(a, h)$$

$$q_{e2e} = \text{TowerInstruct-LoRA+Adapter-pt}(a, h)$$

$$q_{e2e} = \text{TowerInstruct-LoRA+Adapter-pt-Fixed}(a, h)$$

$$\rho = \text{corr}(\mathbf{q}, \mathbf{m})$$

$$\mathbf{m}_{\text{xCOMET}} = \text{xCOMET}(\text{gold } t, h, r)$$

$$\mathbf{m}_{\text{MetricX}} = \text{MetricX}(h, r)$$

$$\rho = \text{corr}(\mathbf{q}, \mathbf{d})$$

Human
Direct
Assessm
ent score

Evaluation

We compute the correlations between SpeechQE scores (q) and :

1. m: Metric score of xCOMET-XL & MetricX-23-XL on En2En/En2De SpeechQE test (3.5k)
2. d: Human Direct Assessment score on IWSLT23-ACL En2De Speech Translation
 - source-based DA ratings of 416 hypotheses from each of the ten ST systems(4.1k)

$$\rho = \text{corr}(\mathbf{q}, \mathbf{m} \text{ or } \mathbf{d})$$

	CoVoST2 Es2En	CoVoST2 En2De	IWSLT23
	m_{xCOMET}	m_{MetricX}	En2De d

Cascaded SpeechQE Systems Correlations

$$q_{cas} = \text{xCOMET-qe}(\text{ASR}(a), h)$$

$$q_{cas} = \text{MetricX-qe}(\text{ASR}(a), h)$$

$$q_{cas} = \text{text-BLASER2.0-qe}(\text{ASR}(a), h)$$

End-to-End SpeechQE Systems Correlations

$$q_{e2e} = \text{BLASER2.0-qe}(a, h)$$

$$q_{e2e} = \text{TowerInstruct-Fixed+Adapter}(a, h)$$

$$q_{e2e} = \text{TowerInstruct-LoRA+Adapter}(a, h)$$

$$q_{e2e} = \text{TowerInstruct-LoRA+Adapter-pt}(a, h)$$

$$q_{e2e} = \text{TowerInstruct-LoRA+Adapter-pt-Fixed}(a, h)$$

$$\rho = \text{corr}(\mathbf{q}, \mathbf{m})$$

$$m_{\text{xCOMET}} = \text{xCOMET}(\text{gold } t, h, r)$$

$$m_{\text{MetricX}} = \text{MetricX}(h, r)$$

$$\rho = \text{corr}(\mathbf{q}, \mathbf{d})$$

Human
Direct
Assessm
ent score

Evaluation

We compute the correlations between SpeechQE scores (q) and :

1. **m: Metric score of xCOMET-XL & MetricX-23-XL** on En2En/En2De SpeechQE test (3.5k)
2. **d: Human Direct Assessment score** on IWSLT23-ACL En2De Speech Translation
 - source-based DA ratings of 416 hypotheses from each of the ten ST systems(4.1k)

$$\rho = \text{corr}(\mathbf{q}, \mathbf{m} \text{ or } \mathbf{d})$$

CoVoST2 Es2En	CoVoST2 En2De	IWSLT23 En2De d
$\mathbf{m}_{\text{xCOMET}}$	$\mathbf{m}_{\text{MetricX}}$	

Cascaded SpeechQE Systems Correlations

$$q_{cas} = \text{xCOMET-qe}(\text{ASR}(a), h)$$

$$q_{cas} = \text{MetricX-qe}(\text{ASR}(a), h)$$

$$q_{cas} = \text{text-BLASER2.0-qe}(\text{ASR}(a), h)$$

End-to-End SpeechQE Systems Correlations

$$q_{e2e} = \text{BLASER2.0-qe}(a, h)$$

$$q_{e2e} = \text{TowerInstruct-Fixed+Adapter}(a, h)$$

$$q_{e2e} = \text{TowerInstruct-LoRA+Adapter}(a, h)$$

$$q_{e2e} = \text{TowerInstruct-LoRA+Adapter-pt}(a, h)$$

$$q_{e2e} = \text{TowerInstruct-LoRA+Adapter-pt-Fixed}(a, h)$$

$$\rho = \text{corr}(\mathbf{q}, \mathbf{m})$$

$$\mathbf{m}_{\text{xCOMET}} = \text{xCOMET}(\text{gold } t, h, r)$$

$$\mathbf{m}_{\text{MetricX}} = \text{MetricX}(h, r)$$

$$\rho = \text{corr}(\mathbf{q}, \mathbf{d})$$

Human Direct Assessment score

Evaluation

We compute the correlations between SpeechQE scores (q) and :

1. **m**: Metric score of xCOMET-XL & MetricX-23-XL on En2En/En2De SpeechQE test (3.5k)
2. **d**: **Human Direct Assessment score** on **IWSLT23-ACL En2De Speech Translation**
 - source-based DA ratings of 416 hypotheses from each of the ten ST systems (4.1k)

$\rho = \text{corr}(\mathbf{q}, \mathbf{m} \text{ or } \mathbf{d})$	CoVoST2 Es2En m_{xCOMET} m_{MetricX}	CoVoST2 En2De m_{xCOMET} m_{MetricX}	IWSLT23 En2De d
<p>Cascaded SpeechQE Systems Correlations</p> <p>$q_{cas} = \text{xCOMET-qe}(\text{ASR}(a), h)$</p> <p>$q_{cas} = \text{MetricX-qe}(\text{ASR}(a), h)$</p> <p>$q_{cas} = \text{text-BLASER2.0-qe}(\text{ASR}(a), h)$</p> <p>End-to-End SpeechQE Systems Correlations</p> <p>$q_{e2e} = \text{BLASER2.0-qe}(a, h)$</p> <p>$q_{e2e} = \text{TowerInstruct-Fixed+Adapter}(a, h)$</p> <p>$q_{e2e} = \text{TowerInstruct-LoRA+Adapter}(a, h)$</p> <p>$q_{e2e} = \text{TowerInstruct-LoRA+Adapter-pt}(a, h)$</p> <p>$q_{e2e} = \text{TowerInstruct-LoRA+Adapter-pt-Fixed}(a, h)$</p>	$\rho = \text{corr}(\mathbf{q}, \mathbf{m})$ $m_{\text{xCOMET}} = \text{xCOMET}(\text{gold } t, h, r)$ $m_{\text{MetricX}} = \text{MetricX}(h, r)$		$\rho = \text{corr}(\mathbf{q}, \mathbf{d})$ <p style="text-align: center;">Human Direct Assessm ent score</p>

Correlation with Reference-based Metrics

$\rho = \text{corr}(\mathbf{q}, \mathbf{m})$	$\mathbf{m}_{\text{xCOMET}} = \text{xCOMET}(\text{gold } t, h, r)$	Es2En		En2De	
	$\mathbf{m}_{\text{MetricX}} = \text{MetricX}(h, r)$	$\mathbf{m}_{\text{xCOMET}}$	$\mathbf{m}_{\text{MetricX}}$	$\mathbf{m}_{\text{xCOMET}}$	$\mathbf{m}_{\text{MetricX}}$
<i>Cascaded SpeechQE Systems Correlations</i> $\rho_{\text{cas}} = \text{corr}(\mathbf{q}_{\text{cas}}, \mathbf{m})$					
$\mathbf{q}_{\text{cas}} = \text{xCOMET-qe}(\text{ASR}(a), h)$		$\overline{0.892}$	0.782	$\overline{0.910}$	0.821
$\mathbf{q}_{\text{cas}} = \text{MetricX-qe}(\text{ASR}(a), h)$		0.803	0.803	0.854	$\overline{0.871}$
$\mathbf{q}_{\text{cas}} = \text{text-BLASER2.0-qe}(\text{ASR}(a), h)$		0.776	0.711	0.813	0.771
<i>End-to-End SpeechQE Systems Correlations</i> $\rho_{\text{e2e}} = \text{corr}(\mathbf{q}_{\text{e2e}}, \mathbf{m})$					
$\mathbf{q}_{\text{e2e}} = \text{BLASER2.0-qe}(a, h)$		0.780	0.712	0.856	0.819
$\mathbf{q}_{\text{e2e}} = \text{TowerInstruct-Fixed+Adapter}(a, h)$		0.862	0.797	0.882	0.848
$\mathbf{q}_{\text{e2e}} = \text{TowerInstruct-LoRA+Adapter}(a, h)$		0.882	0.818	0.914	0.867
$\mathbf{q}_{\text{e2e}} = \text{TowerInstruct-LoRA+Adapter-pt}(a, h)$		0.890	0.833	0.922	0.872
$\mathbf{q}_{\text{e2e}} = \text{TowerInstruct-LoRA+Adapter-pt-Fixed}(a, h)$		0.895	0.834	0.925	0.873

Correlation with Metrics - Cross Comparison

$$\rho = \text{corr}(\mathbf{q}, \mathbf{m})$$

$$\mathbf{m}_{\text{xCOMET}} = \text{xCOMET}(\text{gold } t, h, r)$$

$$\mathbf{m}_{\text{MetricX}} = \text{MetricX}(h, r)$$

Es2En

En2De

$\mathbf{m}_{\text{xCOMET}}$

$\mathbf{m}_{\text{MetricX}}$

$\mathbf{m}_{\text{xCOMET}}$

$\mathbf{m}_{\text{MetricX}}$

Cascaded SpeechQE Systems Correlations $\rho_{\text{cas}} = \text{corr}(\mathbf{q}_{\text{cas}}, \mathbf{m})$

$$\mathbf{q}_{\text{cas}} = \text{xCOMET-qe}(\text{ASR}(a), h)$$

$$\mathbf{q}_{\text{cas}} = \text{MetricX-qe}(\text{ASR}(a), h)$$

$$\mathbf{q}_{\text{cas}} = \text{text-BLASER2.0-qe}(\text{ASR}(a), h)$$

$\overline{0.892} > 0.782$

$\overline{0.910}$ 0.821

0.803 0.803 0.854 < $\overline{0.871}$

0.776 0.711 0.813 0.771

A single xCOMET and MetricX model for *metric* and *QE* settings

→ matching QE and metric model could favor the output from the model similar to its own.

Correlation with Metrics - Cross Comparison

$$\rho = \text{corr}(\mathbf{q}, \mathbf{m})$$

$$\mathbf{m}_{\text{xCOMET}} = \text{xCOMET}(\text{gold } t, h, r)$$

$$\mathbf{m}_{\text{MetricX}} = \text{MetricX}(h, r)$$

Es2En

En2De

$\mathbf{m}_{\text{xCOMET}}$

$\mathbf{m}_{\text{MetricX}}$

$\mathbf{m}_{\text{xCOMET}}$

$\mathbf{m}_{\text{MetricX}}$

Cascaded SpeechQE Systems Correlations $\rho_{\text{cas}} = \text{corr}(\mathbf{q}_{\text{cas}}, \mathbf{m})$

$$\mathbf{q}_{\text{cas}} = \text{xCOMET-qe}(\text{ASR}(a), h)$$

$$\mathbf{q}_{\text{cas}} = \text{MetricX-qe}(\text{ASR}(a), h)$$

$$\mathbf{q}_{\text{cas}} = \text{text-BLASER2.0-qe}(\text{ASR}(a), h)$$

$\overline{0.892} > 0.782$

$\overline{0.910}$ 0.821

0.803 0.803 0.854 < $\overline{0.871}$

0.776 0.711 0.813 0.771

A single xCOMET and MetricX model for *metric* and *QE* settings

→ matching QE and metric model could favor the output from the model similar to its own.

Correlation with Metrics - Cross Comparison

$$\rho = \text{corr}(\mathbf{q}, \mathbf{m})$$

$$\mathbf{m}_{\text{xCOMET}} = \text{xCOMET}(\text{gold } t, h, r)$$

$$\mathbf{m}_{\text{MetricX}} = \text{MetricX}(h, r)$$

Es2En

En2De

$\mathbf{m}_{\text{xCOMET}}$

$\mathbf{m}_{\text{MetricX}}$

$\mathbf{m}_{\text{xCOMET}}$

$\mathbf{m}_{\text{MetricX}}$

Cascaded SpeechQE Systems Correlations $\rho_{\text{cas}} = \text{corr}(\mathbf{q}_{\text{cas}}, \mathbf{m})$

$$\mathbf{q}_{\text{cas}} = \text{xCOMET-qe}(\text{ASR}(a), h)$$

$$\mathbf{q}_{\text{cas}} = \text{MetricX-qe}(\text{ASR}(a), h)$$

$$\mathbf{q}_{\text{cas}} = \text{text-BLASER2.0-qe}(\text{ASR}(a), h)$$

$\overline{0.892} > 0.782$

$\overline{0.910}$ 0.821

0.803 0.803 0.854 < $\overline{0.871}$

0.776 0.711 0.813 0.771

A single xCOMET and MetricX model for *metric* and *QE* settings

→ matching QE and metric model could favor the output from the model similar to its own.

Correlation with Reference-based Metrics

$\rho = \text{corr}(\mathbf{q}, \mathbf{m})$	$\mathbf{m}_{\text{xCOMET}} = \text{xCOMET}(\text{gold } t, h, r)$	Es2En		En2De	
	$\mathbf{m}_{\text{MetricX}} = \text{MetricX}(h, r)$	$\mathbf{m}_{\text{xCOMET}}$	$\mathbf{m}_{\text{MetricX}}$	$\mathbf{m}_{\text{xCOMET}}$	$\mathbf{m}_{\text{MetricX}}$
<i>Cascaded SpeechQE Systems Correlations</i> $\rho_{cas} = \text{corr}(\mathbf{q}_{cas}, \mathbf{m})$					
$\mathbf{q}_{cas} = \text{xCOMET-qe}(\text{ASR}(a), h)$		$\overline{0.892}$	0.782	$\overline{0.910}$	0.821
$\mathbf{q}_{cas} = \text{MetricX-qe}(\text{ASR}(a), h)$		0.803	0.803	0.854	$\overline{0.871}$
$\mathbf{q}_{cas} = \text{text-BLASER2.0-qe}(\text{ASR}(a), h)$		0.776	0.711	0.813	0.771
<i>End-to-End SpeechQE Systems Correlations</i> $\rho_{e2e} = \text{corr}(\mathbf{q}_{e2e}, \mathbf{m})$					
$\mathbf{q}_{e2e} = \text{BLASER2.0-qe}(a, h)$		0.780	0.712	0.856	0.819
$\mathbf{q}_{e2e} = \text{TowerInstruct-Fixed+Adapter}(a, h)$		0.862	0.797	0.882	0.848
$\mathbf{q}_{e2e} = \text{TowerInstruct-LoRA+Adapter}(a, h)$		0.882	0.818	0.914	0.867
$\mathbf{q}_{e2e} = \text{TowerInstruct-LoRA+Adapter-pt}(a, h)$		0.890	0.833	0.922	0.872
$\mathbf{q}_{e2e} = \text{TowerInstruct-LoRA+Adapter-pt-Fixed}(a, h)$		0.895	0.834	0.925	0.873

Correlation with Reference-based Metrics

$\rho = \text{corr}(\mathbf{q}, \mathbf{m})$	$\mathbf{m}_{\text{xCOMET}} = \text{xCOMET}(\text{gold } t, h, r)$	Es2En		En2De	
	$\mathbf{m}_{\text{MetricX}} = \text{MetricX}(h, r)$	$\mathbf{m}_{\text{xCOMET}}$	$\mathbf{m}_{\text{MetricX}}$	$\mathbf{m}_{\text{xCOMET}}$	$\mathbf{m}_{\text{MetricX}}$
<i>Cascaded SpeechQE Systems Correlations</i> $\rho_{cas} = \text{corr}(\mathbf{q}_{cas}, \mathbf{m})$					
$\mathbf{q}_{cas} = \text{xCOMET-qe}(\text{ASR}(a), h)$		0.892	0.782	0.910	0.821
$\mathbf{q}_{cas} = \text{MetricX-qe}(\text{ASR}(a), h)$		0.803	0.803	0.854	0.871
$\mathbf{q}_{cas} = \text{text-BLASER2.0-qe}(\text{ASR}(a), h)$		0.776	0.711	0.813	0.771
<i>End-to-End SpeechQE Systems Correlations</i> $\rho_{e2e} = \text{corr}(\mathbf{q}_{e2e}, \mathbf{m})$					
$\mathbf{q}_{e2e} = \text{BLASER2.0-qe}(a, h)$		0.780	0.712	0.856	0.819
$\mathbf{q}_{e2e} = \text{TowerInstruct-Fixed+Adapter}(a, h)$		0.862	0.797	0.882	0.848
$\mathbf{q}_{e2e} = \text{TowerInstruct-LoRA+Adapter}(a, h)$		0.882	0.818	0.914	0.867
$\mathbf{q}_{e2e} = \text{TowerInstruct-LoRA+Adapter-pt}(a, h)$		0.890	0.833	0.922	0.872
$\mathbf{q}_{e2e} = \text{TowerInstruct-LoRA+Adapter-pt-Fixed}(a, h)$		0.895	0.834	0.925	0.873



Correlation with Reference-based Metrics

$\rho = corr(\mathbf{q}, \mathbf{m})$	$\mathbf{m}_{xCOMET} = xCOMET(\text{gold } t, h, r)$	Es2En		En2De	
	$\mathbf{m}_{MetricX} = MetricX(h, r)$	\mathbf{m}_{xCOMET}	$\mathbf{m}_{MetricX}$	\mathbf{m}_{xCOMET}	$\mathbf{m}_{MetricX}$
<i>Cascaded SpeechQE Systems Correlations</i> $\rho_{cas} = corr(\mathbf{q}_{cas}, \mathbf{m})$					
$\mathbf{q}_{cas} = xCOMET\text{-qe}(ASR(a), h)$		0.892	0.782	0.910	0.821
$\mathbf{q}_{cas} = MetricX\text{-qe}(ASR(a), h)$		0.803	0.803	0.854	0.871
$\mathbf{q}_{cas} = \text{text-BLASER2.0-qe}(ASR(a), h)$		0.776	0.711	0.813	0.771
<i>End-to-End SpeechQE Systems Correlations</i> $\rho_{e2e} = corr(\mathbf{q}_{e2e}, \mathbf{m})$					
$\mathbf{q}_{e2e} = \text{BLASER2.0-qe}(a, h)$		0.780	0.712	0.856	0.819
$\mathbf{q}_{e2e} = \text{TowerInstruct-Fixed+Adapter}(a, h)$		0.862	0.797	0.882	0.848
$\mathbf{q}_{e2e} = \text{TowerInstruct-LoRA+Adapter}(a, h)$		0.882	0.818	0.914	0.867
$\mathbf{q}_{e2e} = \text{TowerInstruct-LoRA+Adapter-pt}(a, h)$		0.890	0.833	0.922	0.872
$\mathbf{q}_{e2e} = \text{TowerInstruct-LoRA+Adapter-pt-Fixed}(a, h)$		0.895	0.834	0.925	0.873

Weight updates at least partially are necessary when a text-LLM is not fine-tuned with QE tasks



Correlation with Reference-based Metrics

$\rho = corr(\mathbf{q}, \mathbf{m})$	$\mathbf{m}_{xCOMET} = xCOMET(\text{gold } t, h, r)$	Es2En		En2De	
	$\mathbf{m}_{MetricX} = MetricX(h, r)$	\mathbf{m}_{xCOMET}	$\mathbf{m}_{MetricX}$	\mathbf{m}_{xCOMET}	$\mathbf{m}_{MetricX}$
Cascaded SpeechQE Systems Correlations $\rho_{cas} = corr(\mathbf{q}_{cas}, \mathbf{m})$					
$\mathbf{q}_{cas} = xCOMET\text{-qe}(ASR(a), h)$	0.892	0.782	0.910	0.821	
$\mathbf{q}_{cas} = MetricX\text{-qe}(ASR(a), h)$	0.803	0.803	0.854	0.871	
$\mathbf{q}_{cas} = \text{text-BLASER2.0-qe}(ASR(a), h)$	0.776	0.711	0.813	0.771	
End-to-End SpeechQE Systems Correlations $\rho_{e2e} = corr(\mathbf{q}_{e2e}, \mathbf{m})$					
$\mathbf{q}_{e2e} = \text{BLASER2.0-qe}(a, h)$	0.780	0.712	0.856	0.819	
$\mathbf{q}_{e2e} = \text{TowerInstruct-Fixed+Adapter}(a, h)$	0.862	0.797	0.882	0.848	
$\mathbf{q}_{e2e} = \text{TowerInstruct-LoRA+Adapter}(a, h)$	0.882	0.818	0.914	0.867	
$\mathbf{q}_{e2e} = \text{TowerInstruct-LoRA+Adapter-pt}(a, h)$	0.890	0.833	0.922	0.872	
$\mathbf{q}_{e2e} = \text{TowerInstruct-LoRA+Adapter-pt-Fixed}(a, h)$	0.895	0.834	0.925	0.873	

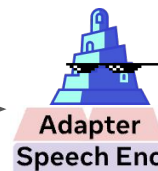
LLM-LoRA
+Adapter-pt&Fixed



ASR+ST



SpeechQE

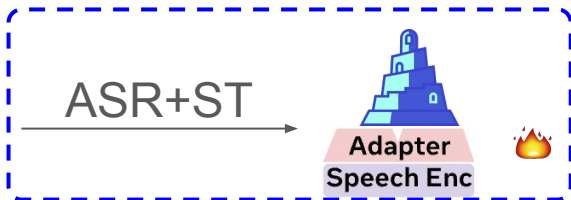


Correlation with Reference-based Metrics

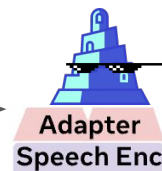
$\rho = corr(\mathbf{q}, \mathbf{m})$	$\mathbf{m}_{xCOMET} = xCOMET(\text{gold } t, h, r)$	Es2En		En2De	
	$\mathbf{m}_{MetricX} = MetricX(h, r)$	\mathbf{m}_{xCOMET}	$\mathbf{m}_{MetricX}$	\mathbf{m}_{xCOMET}	$\mathbf{m}_{MetricX}$
Cascaded SpeechQE Systems Correlations $\rho_{cas} = corr(\mathbf{q}_{cas}, \mathbf{m})$					
$\mathbf{q}_{cas} = xCOMET\text{-qe}(ASR(a), h)$	0.892	0.782	0.910	0.821	
$\mathbf{q}_{cas} = MetricX\text{-qe}(ASR(a), h)$	0.803	0.803	0.854	0.871	
$\mathbf{q}_{cas} = \text{text-BLASER2.0-qe}(ASR(a), h)$	0.776	0.711	0.813	0.771	
End-to-End SpeechQE Systems Correlations $\rho_{e2e} = corr(\mathbf{q}_{e2e}, \mathbf{m})$					
$\mathbf{q}_{e2e} = \text{BLASER2.0-qe}(a, h)$	0.780	0.712	0.856	0.819	
$\mathbf{q}_{e2e} = \text{TowerInstruct-Fixed+Adapter}(a, h)$	0.862	0.797	0.882	0.848	
$\mathbf{q}_{e2e} = \text{TowerInstruct-LoRA+Adapter}(a, h)$	0.882	0.818	0.914	0.867	
$\mathbf{q}_{e2e} = \text{TowerInstruct-LoRA+Adapter-pt}(a, h)$	0.890	0.833	0.922	0.872	
$\mathbf{q}_{e2e} = \text{TowerInstruct-LoRA+Adapter-pt-Fixed}(a, h)$	0.895	0.834	0.925	0.873	

a separate training phase for mapping speech-to-text perception is critical.

LLM-LoRA
+Adapter-pt&Fixed



SpeechQE



Correlation with Human Score

IWSLT23-ACL En2De Test set	Human DA
$\rho = \text{corr}(\mathbf{q}, \mathbf{d})$	score d
<i>Cascaded SpeechQE and Human DA</i> $\rho = \text{corr}(\mathbf{q}_{cas}, \mathbf{d})$	
$q = \text{xCOMET-qe}(\text{ASR}(a), h)$	0.485
$q = \text{MetricX-qe}(\text{ASR}(a), h)$	0.495
$q = \text{wmt23-cometkiwi-da-xl}(\text{ASR}(a), h)$	0.503
$q = \text{wmt22-cometkiwi-da}(\text{ASR}(a), h)$	0.486
$q = \text{text-BLASER2.0-qe}(\text{ASR}(a), h)$	0.428
<i>E2E SpeechQE & Human DA correlation</i> $\rho = \text{corr}(\mathbf{q}_{e2e}, \mathbf{d})$	
$q = \text{BLASER2.0-qe}(a, h)$	0.420
$q = \text{TowerInst-LoRA+Adapter-pt}(a, h)$	0.492
$q = \text{TowerInst-LoRA+Adapter-pt-Fixed}(a, h)$	0.509

Correlation with Human Score

IWSLT23-ACL En2De Test set	Human DA
$\rho = \text{corr}(\mathbf{q}, \mathbf{d})$	score d
<i>Cascaded SpeechQE and Human DA</i> $\rho = \text{corr}(\mathbf{q}_{cas}, \mathbf{d})$	
$q = \text{xCOMET-qe}(\text{ASR}(a), h)$	0.485
$q = \text{MetricX-qe}(\text{ASR}(a), h)$	0.495
$q = \text{wmt23-cometkiwi-da-xl}(\text{ASR}(a), h)$	0.503
$q = \text{wmt22-cometkiwi-da}(\text{ASR}(a), h)$	0.486
$q = \text{text-BLASER2.0-qe}(\text{ASR}(a), h)$	0.428
<i>E2E SpeechQE & Human DA correlation</i> $\rho = \text{corr}(\mathbf{q}_{e2e}, \mathbf{d})$	
$q = \text{BLASER2.0-qe}(a, h)$	0.420
$q = \text{TowerInst-LoRA+Adapter-pt}(a, h)$	0.492
$q = \text{TowerInst-LoRA+Adapter-pt-Fixed}(a, h)$	0.509

Correlation with Human Score

IWSLT23-ACL En2De Test set	Human DA score d
$\rho = corr(\mathbf{q}, \mathbf{d})$	
Cascaded SpeechQE and Human DA $\rho = corr(\mathbf{q}_{cas}, \mathbf{d})$	
$q = \text{xCOMET-qe}(ASR(a), h)$	0.485
$q = \text{MetricX-qe}(ASR(a), h)$	0.495
$q = \text{wmt23-cometkiwi-da-xl}(ASR(a), h)$	0.503
$q = \text{wmt22-cometkiwi-da}(ASR(a), h)$	0.486
$q = \text{text-BLASER2.0-qe}(ASR(a), h)$	0.428
E2E SpeechQE & Human DA correlation $\rho = corr(\mathbf{q}_{e2e}, \mathbf{d})$	
$q = \text{BLASER2.0-qe}(a, h)$	0.420
$q = \text{TowerInst-LoRA+Adapter-pt}(a, h)$	0.492
$q = \text{TowerInst-LoRA+Adapter-pt-Fixed}(a, h)$	0.509

LLM-LoRA

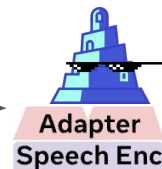
+Adapter-pt&Fixed



ASR+ST



SpeechQE



ASR = Azure API, output provided by

Correlation with Human Score

The best-practice E2E system is more effective in aligning with human judgments.

IWSLT23-ACL En2De Test set	Human DA
$\rho = \text{corr}(\mathbf{q}, \mathbf{d})$	score d
<i>Cascaded SpeechQE and Human DA</i> $\rho = \text{corr}(\mathbf{q}_{cas}, \mathbf{d})$	
$q = \text{xCOMET-qe}(\text{ASR}(a), h)$	0.485
$q = \text{MetricX-qe}(\text{ASR}(a), h)$	0.495
$q = \text{wmt23-cometkiwi-da-xl}(\text{ASR}(a), h)$	0.503
$q = \text{wmt22-cometkiwi-da}(\text{ASR}(a), h)$	0.486
$q = \text{text-BLASER2.0-qe}(\text{ASR}(a), h)$	0.428
<i>E2E SpeechQE & Human DA correlation</i> $\rho = \text{corr}(\mathbf{q}_{e2e}, \mathbf{d})$	
$q = \text{BLASER2.0-qe}(a, h)$	0.420
$q = \text{TowerInst-LoRA+Adapter-pt}(a, h)$	0.492
$q = \text{TowerInst-LoRA+Adapter-pt-Fixed}(a, h)$	0.509

LLM-LoRA

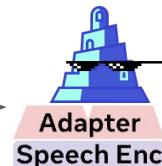
+Adapter-pt&Fixed



ASR+ST



SpeechQE



ASR = Azure API, output provided by

Cascaded Model Size and Architecture

Is the dominance of E2E over cascaded models due to the E2E parameter size or architecture rather than its end-to-end nature?

$\rho = \text{corr}(\mathbf{q}, \mathbf{m} \text{ or } \mathbf{d})$	CoVoST2 Es2En Test				IWSLT23
	$\mathbf{m}_{\text{xCOMET-XL}}$	$\mathbf{m}_{\text{xCOMET-XXL}}$	$\mathbf{m}_{\text{MetricX-XL}}$	$\mathbf{m}_{\text{MetricX-XXL}}$	En2De <i>d</i>
<i>Cascaded Model with XXL Size vs E2E speech-LLM</i>					
$q_{cas} = \text{ASR (1.5B)} \rightarrow \text{xCOMET-XL-qe (3.5B)}$	0.892	0.800	0.782	0.788	0.485
$q_{cas} = \text{ASR (1.5B)} \rightarrow \text{xCOMET-XXL-qe (10.7B)}$	0.787	0.873	0.708	0.734	0.486
$q_{cas} = \text{ASR (1.5B)} \rightarrow \text{MetricX-XL-qe (3.7B)}$	0.803	0.758	0.803	0.766	0.495
$q_{cas} = \text{ASR (1.5B)} \rightarrow \text{MetricX-XXL-qe (13B)}$	0.700	0.677	0.652	0.694	0.502
<i>Cascaded text-LLM vs E2E speech-LLM</i>					
$q_{cas} = \text{ASR (1.5B)} \rightarrow \text{text-TowerInstruct-LoRA (7B)}$	0.852	0.816	0.780	0.785	–
$q_{e2e} = \text{TowerInstruct-LoRA+Adapter-pt-Fixed (7.5B)}$	0.895	0.827	0.834	0.834	0.509

Cascaded Model with Larger Size

$\rho = corr(\mathbf{q}, \mathbf{m} \text{ or } \mathbf{d})$	CoVoST2 Es2En Test				IWSLT23
	$m_{x\text{COMET-XL}}$	$m_{x\text{COMET-XXL}}$	$m_{\text{MetricX-XL}}$	$m_{\text{MetricX-XXL}}$	En2De <i>d</i>
<i>Cascaded Model with XXL Size vs E2E speech-LLM</i>					
$q_{cas} = \text{ASR (1.5B)} \rightarrow x\text{COMET-XL-qe (3.5B)}$	0.892	0.800	0.782	0.788	0.485
$q_{cas} = \text{ASR (1.5B)} \rightarrow x\text{COMET-XXL-qe (10.7B)}$	0.787	0.873	0.708	0.734	0.486
$q_{cas} = \text{ASR (1.5B)} \rightarrow \text{MetricX-XL-qe (3.7B)}$	0.803	0.758	0.803	0.766	0.495
$q_{cas} = \text{ASR (1.5B)} \rightarrow \text{MetricX-XXL-qe (13B)}$	0.700	0.677	0.652	0.694	0.502
<i>Cascaded text-LLM vs E2E speech-LLM</i>					
$q_{cas} = \text{ASR (1.5B)} \rightarrow \text{text-TowerInstruct-LoRA (7B)}$	0.852	0.816	0.780	0.785	-
$q_{e2e} = \text{TowerInstruct-LoRA+Adapter-pt-Fixed (7.5B)}$	0.895	0.827	0.834	0.834	0.509

Cascaded Model with Similar Architecture

$\rho = corr(\mathbf{q}, \mathbf{m} \text{ or } \mathbf{d})$	CoVoST2 Es2En Test				IWSLT23
	$m_{x\text{COMET-XL}}$	$m_{x\text{COMET-XXL}}$	$m_{\text{MetricX-XL}}$	$m_{\text{MetricX-XXL}}$	En2De <i>d</i>
<i>Cascaded Model with XXL Size vs E2E speech-LLM</i>					
$q_{cas} = \text{ASR (1.5B)} \rightarrow x\text{COMET-XL-qe (3.5B)}$	0.892	0.800	0.782	0.788	0.485
$q_{cas} = \text{ASR (1.5B)} \rightarrow x\text{COMET-XXL-qe (10.7B)}$	0.787	0.873	0.708	0.734	0.486
$q_{cas} = \text{ASR (1.5B)} \rightarrow \text{MetricX-XL-qe (3.7B)}$	0.803	0.758	0.803	0.766	0.495
$q_{cas} = \text{ASR (1.5B)} \rightarrow \text{MetricX-XXL-qe (13B)}$	0.700	0.677	0.652	0.694	0.502
<i>Cascaded text-LLM vs E2E speech-LLM</i>					
$q_{cas} = \text{ASR (1.5B)} \rightarrow \text{text-TowerInstruct-LoRA (7B)}$	0.852	0.816	0.780	0.785	–
$q_{e2e} = \text{TowerInstruct-LoRA+Adapter-pt-Fixed (7.5B)}$	0.895	0.827	0.834	0.834	0.509

Cascaded Model with Similar Architecture

The improvements are coming from the E2E nature of the approach rather than the LLM-based solution or larger parameters.

$\rho = corr(\mathbf{q}, \mathbf{m} \text{ or } \mathbf{d})$	CoVoST2 Es2En Test				IWSLT23
	$m_{x\text{COMET-XL}}$	$m_{x\text{COMET-XXL}}$	$m_{\text{MetricX-XL}}$	$m_{\text{MetricX-XXL}}$	En2De <i>d</i>
<i>Cascaded Model with XXL Size vs E2E speech-LLM</i>					
$q_{cas} = \text{ASR (1.5B)} \rightarrow x\text{COMET-XL-qe (3.5B)}$	0.892	0.800	0.782	0.788	0.485
$q_{cas} = \text{ASR (1.5B)} \rightarrow x\text{COMET-XXL-qe (10.7B)}$	0.787	0.873	0.708	0.734	0.486
$q_{cas} = \text{ASR (1.5B)} \rightarrow \text{MetricX-XL-qe (3.7B)}$	0.803	0.758	0.803	0.766	0.495
$q_{cas} = \text{ASR (1.5B)} \rightarrow \text{MetricX-XXL-qe (13B)}$	0.700	0.677	0.652	0.694	0.502
<i>Cascaded text-LLM vs E2E speech-LLM</i>					
$q_{cas} = \text{ASR (1.5B)} \rightarrow \text{text-TowerInstruct-LoRA (7B)}$	0.852	0.816	0.780	0.785	–
$q_{e2e} = \text{TowerInstruct-LoRA+Adapter-pt-Fixed (7.5B)}$	0.895	0.827	0.834	0.834	0.509

→ E2E system is better suited for SpeechQE task than the cascaded

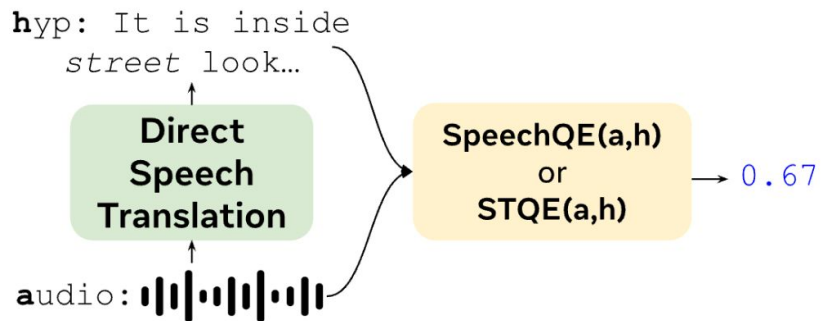
Cascaded Model with Similar Architecture

The improvements are coming from the E2E nature of the approach rather than the LLM-based solution or larger parameters.

$\rho = corr(\mathbf{q}, \mathbf{m} \text{ or } \mathbf{d})$	CoVoST2 Es2En Test				IWSLT23
	$m_{xCOMET-XL}$	$m_{xCOMET-XXL}$	$m_{MetricX-XL}$	$m_{MetricX-XXL}$	En2De <i>d</i>
<i>Cascaded Model with XXL Size vs E2E speech-LLM</i>					
$q_{cas} = ASR (1.5B) \rightarrow xCOMET-XL-qe (3.5B)$	0.892	0.800	0.782	0.788	0.485
$q_{cas} = ASR (1.5B) \rightarrow xCOMET-XXL-qe (10.7B)$	0.787	0.873	0.708	0.734	0.486
$q_{cas} = ASR (1.5B) \rightarrow MetricX-XL-qe (3.7B)$	0.803	0.758	0.803	0.766	0.495
$q_{cas} = ASR (1.5B) \rightarrow MetricX-XXL-qe (13B)$	0.700	0.677	0.652	0.694	0.502
<i>Cascaded text-LLM vs E2E speech-LLM</i>					
$q_{cas} = ASR (1.5B) \rightarrow text-TowerInstruct-LoRA (7B)$	0.852	0.816	0.780	0.785	–
$q_{e2e} = TowerInstruct-LoRA+Adapter-pt-Fixed (7.5B)$	0.895	0.827	0.834	0.834	0.509

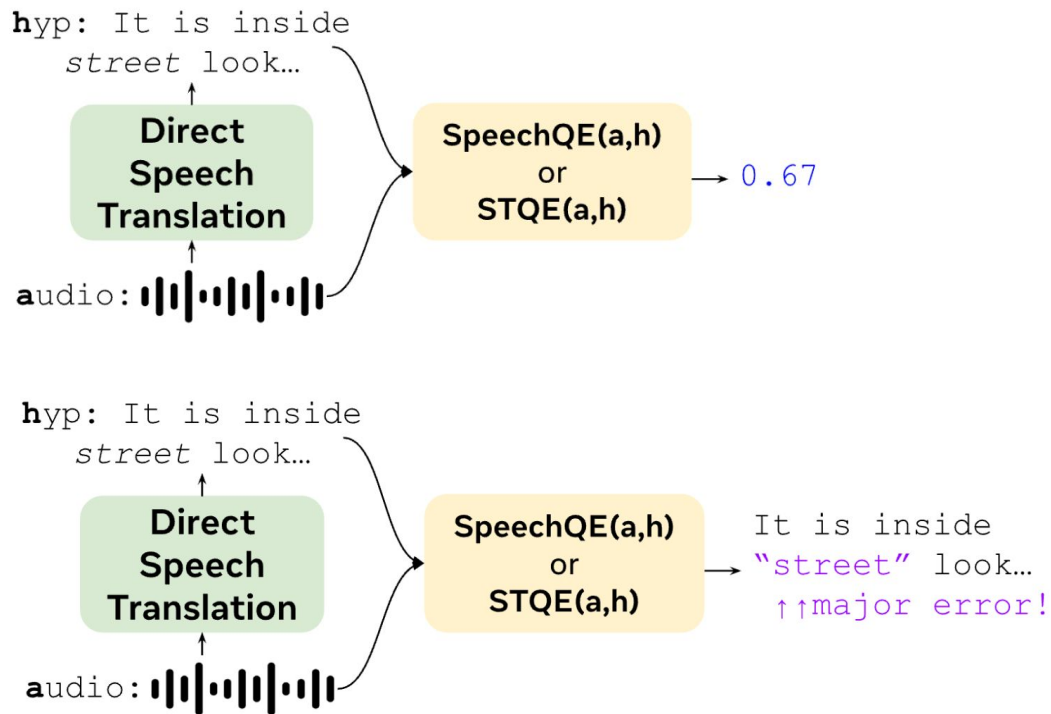
→ **E2E system is better suited for SpeechQE task than the cascaded**

SpeechQE
(numerical rating)



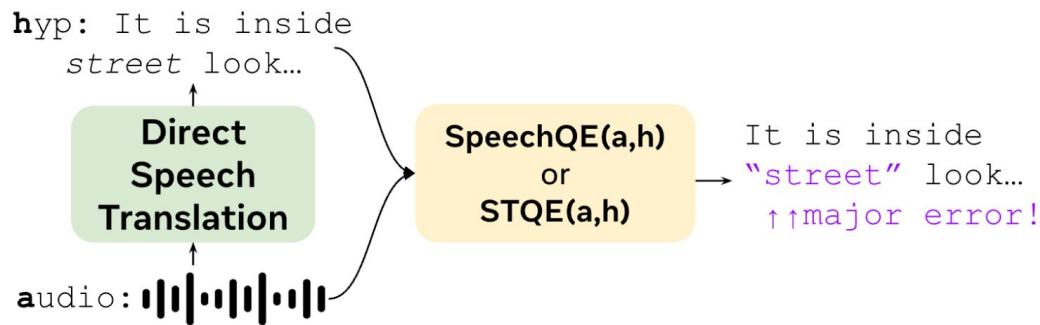
Error Span Detection for ST

SpeechQE
(numerical rating)



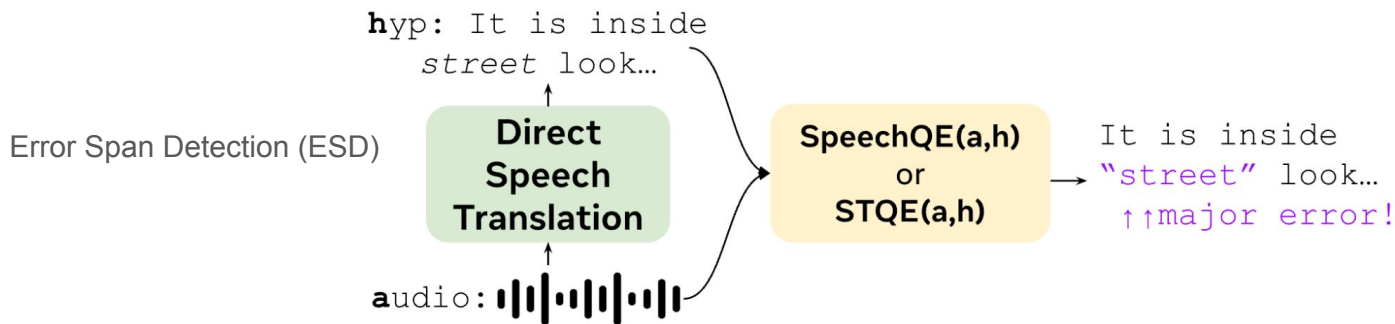
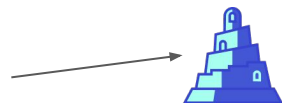
SpeechESD: Error Span Detection for ST

Speech
Error Span
Detection
(SpeechESD)



Zero-Shot Error Span Detection for ST

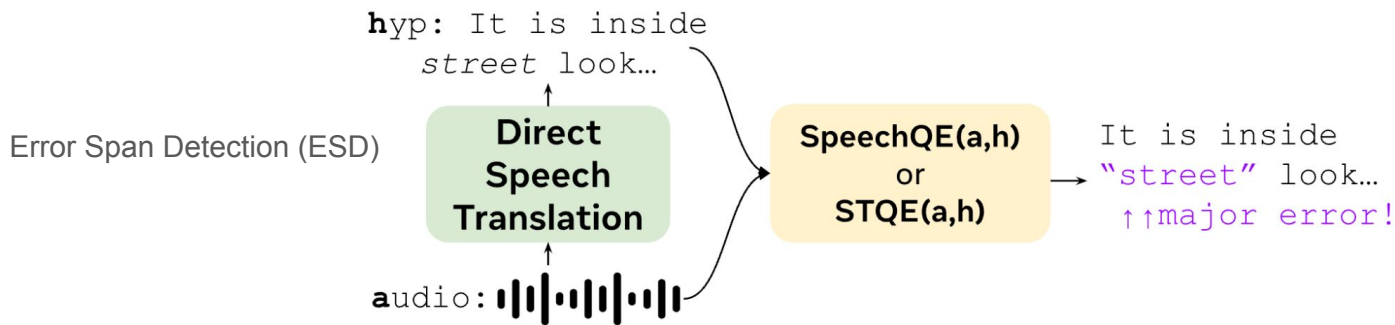
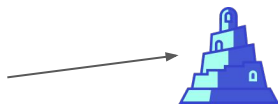
No QE rating capability
But **has error span**
detection ability



Zero-Shot Error Span Detection for ST

How effectively the method of injecting speech modality generalizes the capability of text-LLM to speech LLM without explicitly training the target speech task?

No QE rating capability
But **has error span**
detection ability



SpeechESD Experiment Setting

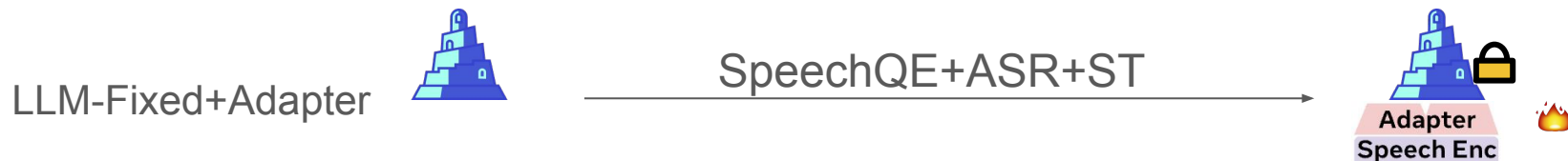
Reference-based error span labels from Error span output of the xCOMET metric function

- xCOMET outputs Error Spans.

Compare:

1. the E2E SpeechESD based on TowerInstruct in zero-shot way
2. cascaded system where TowerInstruct is text-ESD model

Fixed-LLM for E2E: lose ESD capabilities after fine-tuned with non-ESD tasks.



SpeechESD Experiment Setting

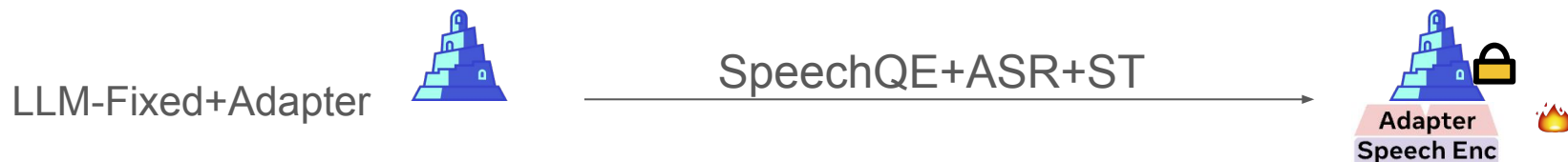
Reference-based error span labels from Error span output of the xCOMET metric function

- xCOMET has ESD capabilities.

Compare:

1. **E2E SpeechESD** based on **TowerInstruct** in **zero-shot way**
2. **cascaded system** where **TowerInstruct** is **text-ESD** model

Fixed-LLM for E2E: lose ESD capabilities after fine-tuned with non-ESD tasks.



SpeechESD Experiment Setting

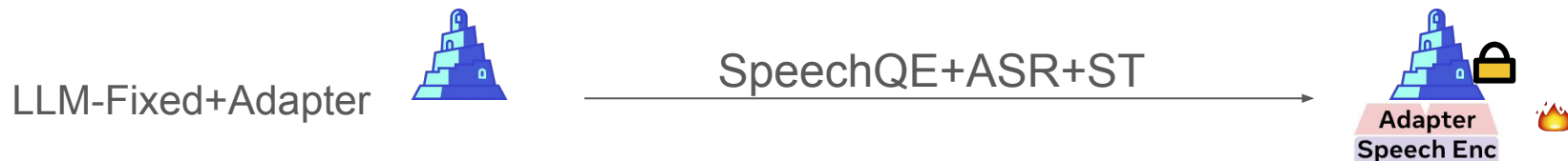
Reference-based error span labels from Error span output of the xCOMET metric function

- xCOMET has ESD capabilities.

Compare:

1. the E2E SpeechESD based on TowerInstruct in zero-shot way
2. cascaded system where TowerInstruct is text-ESD model

Fixed-LLM for E2E: lose ESD capabilities after fine-tuned with non-ESD tasks.



Zero-Shot SpeechESD

ESD for ST	Precision	Recall	F1 Score
<i>Cascaded Systems</i>			
txt-ESD(gold t, h)	0.438	0.591	0.503
txt-ESD(whisper-large-v2(a), h)	0.434	0.550	0.485
txt-ESD(whisper-medium(a), h)	0.429	$\overline{0.540}$	0.478
txt-ESD(whisper-small(a), h)	$\underline{0.413}$	0.535	$\overline{0.466}$
txt-ESD(whisper-base(a), h)	0.385	0.550	0.453
<i>End-to-End Systems</i>			
<i>TowerInst-Fixed+Adt</i> (a, h)	0.411	0.542	0.467

Zero-Shot SpeechESD

Cascaded systems remain the preferred choice for achieving the highest performance when we do not have speech training data for the target task.

ESD for ST	Precision	Recall	F1 Score
<i>Cascaded Systems</i>			
txt-ESD(gold t, h)	0.438	0.591	0.503
txt-ESD(whisper-large-v2(a, h))	0.434	0.550	0.485
txt-ESD(whisper-medium(a, h))	0.429	$\overline{0.540}$	0.478
txt-ESD(whisper-small(a, h))	$\underline{0.413}$	0.535	$\overline{0.466}$
txt-ESD(whisper-base(a, h))	0.385	0.550	0.453
<i>End-to-End Systems</i>			
<i>TowerInst-Fixed+Adt</i> (a, h)	0.411	0.542	0.467

Zero-Shot SpeechESD

Cascaded systems remain the preferred choice for achieving the highest performance when we do not have speech training data for the target task.

Still, the E2E model performs decently in zero-shot

ESD for ST	Precision	Recall	F1 Score
<i>Cascaded Systems</i>			
txt-ESD(gold t, h)	0.438	0.591	0.503
txt-ESD(whisper-large-v2(a), h)	0.434	0.550	0.485
txt-ESD(whisper-medium(a), h)	0.429	<u>0.540</u>	0.478
txt-ESD(whisper-small(a), h)	<u>0.413</u>	0.535	<u>0.466</u>
txt-ESD(whisper-base(a), h)	0.385	0.550	0.453
<i>End-to-End Systems</i>			
<i>TowerInst-Fixed+Adt</i> (a, h)	0.411	0.542	0.467

Zero-Shot SpeechESD

Cascaded systems remain the preferred choice for achieving the highest performance when we do not have speech training data for the target task.

Still, the E2E model performs decently in zero-shot

→ **text-LLM ability is transferable to speech LLM in a zero-shot manner.**

ESD for ST	Precision	Recall	F1 Score
<i>Cascaded Systems</i>			
txt-ESD(gold t, h)	0.438	0.591	0.503
txt-ESD(whisper-large-v2(a, h))	0.434	0.550	0.485
txt-ESD(whisper-medium(a, h))	0.429	$\overline{0.540}$	0.478
txt-ESD(whisper-small(a, h))	$\underline{0.413}$	0.535	$\overline{0.466}$
txt-ESD(whisper-base(a, h))	0.385	0.550	0.453
<i>End-to-End Systems</i>			
<i>TowerInst-Fixed+Adt</i> (a, h)	0.411	0.542	0.467

SpeechQE: Estimating the Quality of Direct Speech Translation

HyoJung Han

Computer Science
University of Maryland
hjhan@cs.umd.edu



Kevin Duh

HLTCOE
Johns Hopkins University
kevinduh@cs.jhu.edu



Marine Carpuat

Computer Science
University of Maryland
marine@cs.umd.edu



github.com/h-j-han/SpeechQE



Huggingface Hub

- **SpeechQE task:** formulation, benchmarks, evaluation of cascaded and E2E architectures
- **E2E SpeechQE model:** methods for corpus creation, training strategies, and architectural design
- **E2E systems are generally better** suited to estimate the quality of direct speech translation
- **SpeechQE need more attention!** It deserves dedicated attention as a separate problem from text-QE quality. Releasing our data and models to guide further research in this space.

SpeechQE: Estimating the Quality of Direct Speech Translation

HyoJung Han

Computer Science
University of Maryland
hjhan@cs.umd.edu



Kevin Duh

HLTCOE
Johns Hopkins University
kevinduh@cs.jhu.edu



Marine Carpuat

Computer Science
University of Maryland
marine@cs.umd.edu



github.com/h-j-han/SpeechQE



Huggingface Hub

- **SpeechQE task:** formulation, benchmarks, evaluation of cascaded and E2E architectures
- **E2E SpeechQE model:** methods for corpus creation, training strategies, and architectural design
- **E2E systems are generally better** suited to estimate the quality of direct speech translation
- **SpeechQE need more attention!** It deserves dedicated attention as a separate problem from text-QE quality. Releasing our data and models to guide further research in this space.

SpeechQE: Estimating the Quality of Direct Speech Translation

HyoJung Han

Computer Science
University of Maryland
hjhan@cs.umd.edu



Kevin Duh

HLTCOE
Johns Hopkins University
kevinduh@cs.jhu.edu



Marine Carpuat

Computer Science
University of Maryland
marine@cs.umd.edu



github.com/h-j-han/SpeechQE



Huggingface Hub

- **SpeechQE task:** formulation, benchmarks, evaluation of cascaded and E2E architectures
- **E2E SpeechQE model:** methods for corpus creation, training strategies, and architectural design
- **E2E systems are generally better** suited to estimate the quality of direct speech translation
- **SpeechQE need more attention!** It deserves dedicated attention as a separate problem from text-QE quality. Releasing our data and models to guide further research in this space.

SpeechQE: Estimating the Quality of Direct Speech Translation

HyoJung Han
Computer Science
University of Maryland
hjhan@cs.umd.edu





Kevin Duh
HLTCOE
Johns Hopkins University
kevinduh@cs.jhu.edu



Marine Carpuat
Computer Science
University of Maryland
marine@cs.umd.edu



 github.com/h-j-han/SpeechQE

 Huggingface Hub

- **SpeechQE task:** formulation, benchmarks, evaluation of cascaded and E2E architectures
- **E2E SpeechQE model:** methods for corpus creation, training strategies, and architectural design
- **E2E systems are generally better** suited to estimate the quality of direct speech translation
- **SpeechQE need more attention!** It deserves dedicated attention as a separate problem from text-QE quality. Releasing our data and models to guide further research in this space.

SpeechQE: Estimating the Quality of Direct Speech Translation

HyoJung Han

Computer Science
University of Maryland
hjhan@cs.umd.edu



Kevin Duh

HLTCOE
Johns Hopkins University
kevinduh@cs.jhu.edu



Marine Carpuat

Computer Science
University of Maryland
marine@cs.umd.edu



github.com/h-j-han/SpeechQE



Huggingface Hub

- **SpeechQE task:** formulation, benchmarks, evaluation of cascaded and E2E architectures
- **E2E SpeechQE model:** methods for corpus creation, training strategies, and architectural design
- **E2E systems are generally better** suited to estimate the quality of direct speech translation
- **SpeechQE need more attention!** It deserves dedicated attention as a separate problem from text-QE quality. Releasing our data and models to guide further research in this space.

Models for Trustworthy Speech Translation

Trustworthiness in Speech Translation

Email: hjhan@umd.edu X: [@h_j_han](#)



DEPARTMENT OF
COMPUTER SCIENCE



JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

SpeechQE: Estimating the Quality of Direct Speech Translation

EMNLP2024 



HyoJung Han
Computer Science
University of Maryland



Kevin Duh
HLTCOE
Johns Hopkins University



Marine Carpuat
Computer Science
University of Maryland