

**Sami Ul Haq**  
Dublin City University

**Sheila Castilho**  
Dublin City University

**Yvette Graham**  
Trinity College Dublin

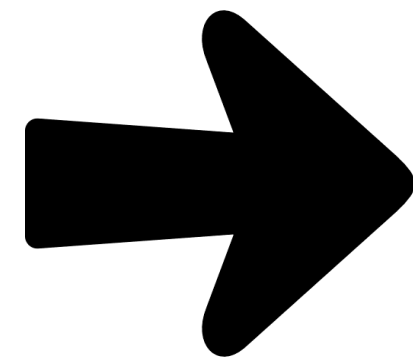
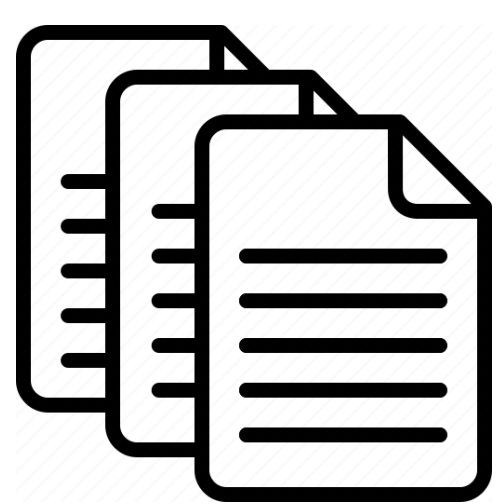
## Why Human Evaluation is Important?

- Human feedback stands as the gold standard for assessing machine translation (MT) quality, aligning closely with users' actual needs and expectations.
- Ensuring the accuracy and adequacy of the translated text is crucial in preventing misunderstandings, errors, and potential legal or financial repercussions.
- While automatic evaluation metrics offer a quick substitute for measurements, their correlation with human judgments is weak.
- Expert-based assessment of translation quality is both expensive and challenging to implement on large-scale.

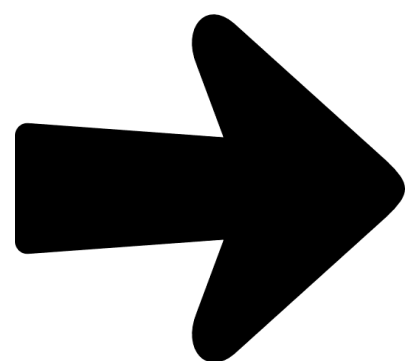
## Multi-modal Input

- Multi-modal MT can directly translate to/from multiple modalities (such as Image, Text, Audio).
- Multi-modal evaluation metrics are essential for comprehensive assessment of MT across different modalities.
- Introducing additional technologies such as Text-to-Speech (TTS) to assessment environments may help identify subtle MT errors as compared to text-only conditions.
- Can automatic speech synthesis be leveraged to enhance MT quality assessment?
- Which of two conditions (text only, or target sound) is more conducive to better assessment of MT quality?

### Reference Translation



### Audio Hypothesis



Read the **reference translation**, listen the **machine translation** audio. Use slider to indicate how much you agree with the statement in **red box**.

Reference:	To invalidate fake news - Biontech CEO publishes vaccination photo
Machine Translation:	<div>▶ 0:00 / 0:05 </div>

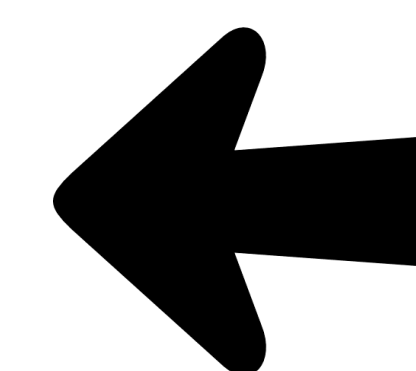
The audio adequately express the meaning of reference translation.

strongly disagree strongly agree

NEXT



Crowd Workers

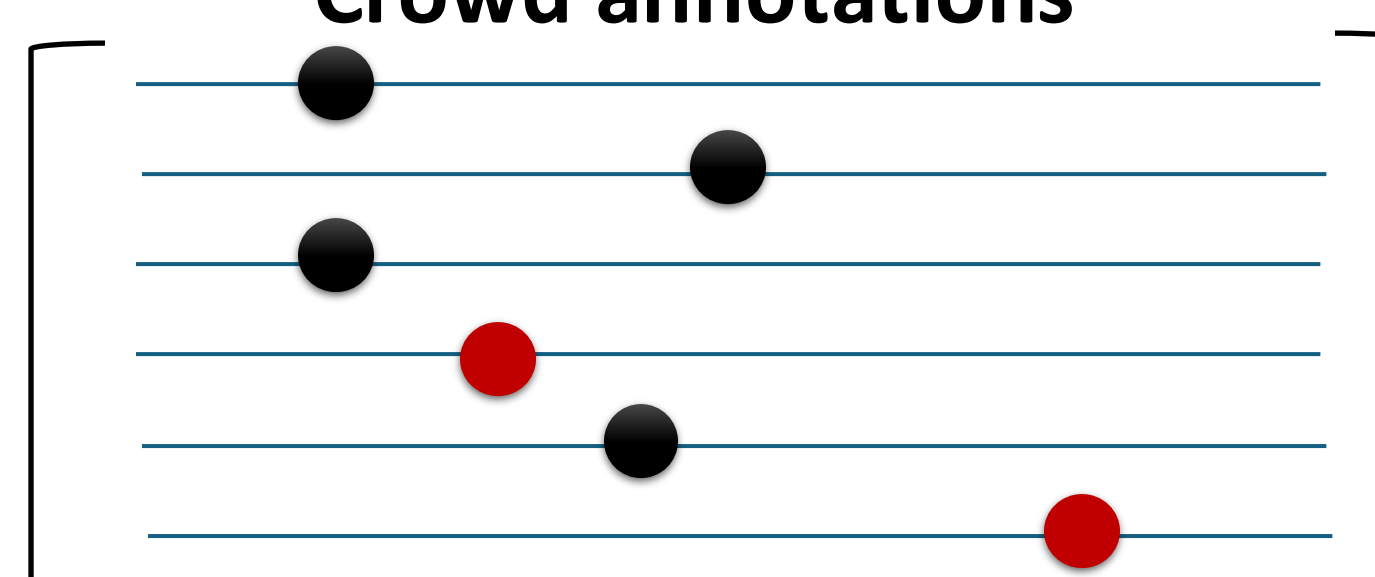


### Crowd annotations

*bad\_references*

*repeat\_questions*

Worker  
reliability  
assessment

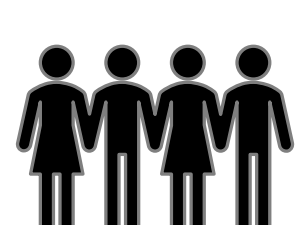


System &  
segment  
level Ranking

*Wilcoxon rank  
– sum test*

$P < 0.05$

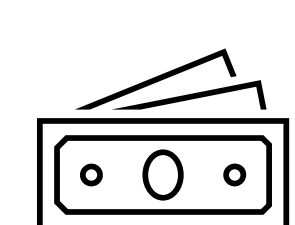
## Crowd-sourcing Human Judgements



Access to large pool of native speakers of different languages from around the world.

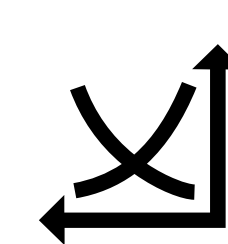


Reference based Direct Assessment (DA) does not require bilingual experts.

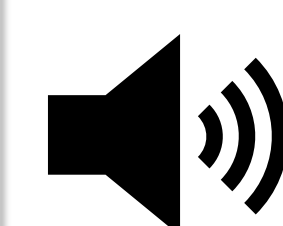


Low-cost assessments provide a viable means for individual MT developers and researchers to assess system improvements over a baseline.

## Forthcoming



Crowd annotation collection and analysis to investigate the of impact of sound conditions in comparison to text-only scenarios.



Enabling natural-sounding artificial voices and flexibility to select between silent or sound conditions.



Exploring the integration of Visual modality (Images, videos) to provide extra contextual information about translation.