

Sami Ul Haq  
Dublin City University

Sheila Castilho  
Dublin City University

Yvette Graham  
Trinity College Dublin

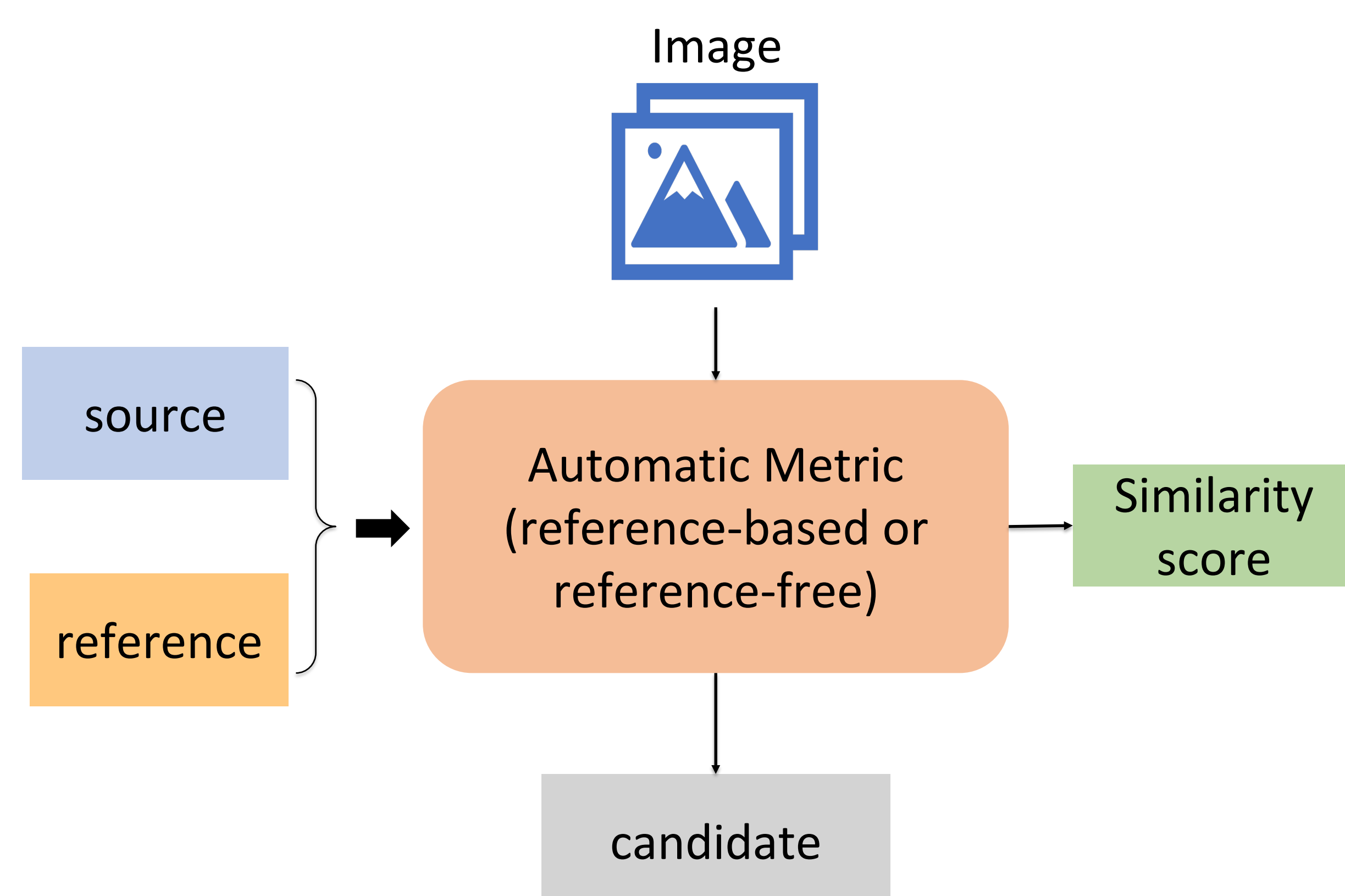
## Current Evaluation Metrics

- Traditional evaluation metrics (e.g., BLEU<sup>1</sup>) rely on **surface-level** n-gram overlap.
- Recent metrics (e.g., BERTScore<sup>2</sup>) address this limitation by leveraging continuous word embeddings to compute similarity in **semantic space**.
- However, token level embedding distance may **overestimate similarity** (e.g., between *cat* and *dog*).
- We propose **CAEMT**—which incorporates **cross-modal semantic similarity** from both textual and visual (image) modalities to enhance the reliability of MT evaluation.

<sup>1</sup>Papineni, Kishore, et al. "Bleu: a Method for Automatic Evaluation of Machine Translation."

<sup>2</sup>Zhang, Tianyi, et al. "BERTScore: Evaluating Text Generation with BERT."

## Our proposed approach (CAEMT)



## Proposed workflow of CAEMT

Text<sub>1</sub>: The cat is on the bed.



Text<sub>2</sub>: The cat is on the carpet.



image representations

text representations

	$T_1$	$T_2$	$T_3$	...	$T_N$
$I_1$	$I_1.T_1$	$I_1.T_2$	$I_1.T_3$	...	$I_1.T_N$
$I_2$	$I_2.T_1$	$I_2.T_2$	$I_2.T_3$	...	$I_2.T_N$
$I_3$	$I_3.T_1$	$I_3.T_2$	$I_3.T_3$	...	$I_3.T_N$
...	...	...	...	...	...
$I_N$	$I_N.T_1$	$I_N.T_2$	$I_N.T_3$	...	$I_N.T_N$

- Use of **Visual-Language Models** (e.g., Jina<sup>3</sup>, CLIP<sup>4</sup>) to generate text and image representations.
- Cross-lingual cross-modal weighted **cosine similarity** between image and text for **reference-based** evaluation:  
$$\omega * \cos(v, t_i, t_j)$$
- Word matching in a **semantic space** (such as BERTScore) using word embeddings.
- Also, **reference-free** evaluation, directly comparing source with candidate text, using image as ground truth.

<sup>3</sup>Koukounas, Andreas, et al. "jina-clip-v2: Multilingual multimodal embeddings for text and images." *arXiv preprint arXiv:2412.08802* (2024).

<sup>4</sup>Hessel, Jack, et al. "CLIPScore: A Reference-free Evaluation Metric for Image Captioning." *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. 2021.

## CAEMT vs traditional metrics: initial results



image<sub>(i)</sub>

Source<sub>(s)</sub>:

Eine graue Katze kuschelt mit einer orangefarbenen Katze auf einer Decke.

Candidate<sub>(c)</sub>:

a **calico** cat is cuddling with an orange **dog** on a blanket.

Reference<sub>(r)</sub>:

a grey cat is cuddling with an orange cat on a blanket.

Metric*	Score	
BLEU <sub>c,r</sub>	↑ 59.20	✗
TER <sub>c,r</sub>	↓ 16.70	✗
BERTSCORE <sub>c,r</sub>	↑ 94.70	✗
COMET-22 <sub>c,r</sub>	↑ 84.10	✗
CAEMT <sub>r,i</sub>	↑ 0.43	✓
CAEMT <sub>c,i</sub>	↓ 0.37	✓
CAEMT <sub>c,r,i</sub>	↑ 0.49	✓
CAEMT <sub>c,s</sub>	↓ 0.77	✗

\*For CAEMT we have used multilingual multimodal jina-clip-v2 model with transformer API to calculate the text-text and text-image cosine similarity . For other metrics we used online MATEO tool (<https://mateo.ivdnt.org>).

