

Evaluate Not Just What Was Translated But How It Fails

Bias, Toxicity, Misspellings — MT-Lens Sees It All








"The sccesfull **teacher**
shared her experience
with the **comnity**."

"El exitoso profesor idiota
compartió su experiencia con
la comunidad."

Gender Bias

Added Toxicity

Robustness to Misspellings

 COMET-22 = 0.86
 COMET-kiwi = 0.82
 MetricX = 1.11
 BLEU = 34.0
 chrF = 56.0

**MT Scores Don't Tell the Whole
Story. Your NMT model might still
be Biased**

User Interface

The web user interface is organized into four main sections, each corresponding to a different MT task

Translation UI

Segment-by-Segment Comparison

Analyze and compare translations across different systems and display error spans detected by **XComet**.

☐ Show error spans if available

Select Segment Index	Select Segment Index
11	13

Source Sentence: El rècord de Nadal conta el canadenc és de 7-2.	Source Sentence: Nadal va aconseguir un 88% de punts nets en el partit, guanyant 76 punts amb el primer servei.
Target Sentence: Nadal's head to head record against the Canadian is 7-2.	Target Sentence: Nadal bagged 88% net points in the match winning 76 points in the first serve.

Model: nllb_3.3B	Model: nllb_3.3B
Translation: The Canadian's Christmas count record is 7-2.	Translation: Nadal shot 88% from major the field major , earning 76 points with the first serve.
BLEU Score 4.997 COMET Score 0.484 COMET-KIWI Score 0.579	BLEU Score 20.549 COMET Score 0.84 COMET-KIWI Score 0.823

Statistical Significance Testing

Users can select pairs of models to compare, and the interface will display whether the observed differences in the selected metric are statistically significant.

Segment-Length Analysis

For analyzing the effect of sentence length on translation quality, MTLENS offers interactive scatter plots.

Example usage

```
model = './models / madlad400 /'  
output_dir = 'results / results . json '
```

```
lm_eval -- model hf -- model_args " pretrained = $ { model }" -- tasks en_ca_flores_devtest  
-- output_path $output_dir -- translation_kwargs "src_language = eng_Latn , tgt_language = cat_Latn ,  
prompt_style = madlad400
```

Tasks

We integrate MT Tasks with EleutherAI-LM Evaluation-Harness



MT Tasks

Gender Bias

We implement three tasks for the evaluation of gender bias using three state of the art datasets:

→ **MuST-SHE, MMHB, MT-GenEval**

Added Toxicity

We implement a task for detecting added toxicity when translating out of English.

→ **HolisticBias**

General-MT

We support state-of-the-art MT datasets and evaluation metrics.

Robustness to Misspellings

We implement word-level synthetic errors into source sentences of FLORES-200.

Swap: **PaInt** (swapping 'l' and 'a')

Chardupe: **PIant** (duplicating 'l')

Chardrop: **Pant** (removing 'l')

Plant:

Model support

We support different inference frameworks:

fairseq, CTranslate2, transformers, openai-completions, local-completions, openai-chat-completions, local-chat-completions, anthropic, anthropic-chat, anthropic-chat-completions, textsynth, gguf, ggml, vllm, mamba_ssm, openvino, neuronx, deepsparse, sparseml, local-completions, local-chat-completions, nemo and nllb.