

Automatic Evaluation of the WMT23 General Machine Translation Task

In this document, we present an automatic evaluation of the systems submitted to the general machine translation task. Please keep in mind that these rankings are not official. WMT only uses human evaluation for the official rankings.

We ran three different automatic metrics:

- chrF (Popović, 2015): A tokenization independent metric operating at character-level with a higher correlation with human judgments than BLEU.
- BLEU (Papineni et al., 2002): The standard BLEU.
- COMET (Rei et al., 2020): A state-of-the-art metric based on a pre-trained language model. We used the default model “Unbabel/wmt22-comet-da.”

chrF and BLEU scores are computed with SacreBLEU (Post, 2018).¹ We ranked the systems according to their scores but display the rank only constrained systems. Unconstrained systems are in gray.

We also tested whether the difference between systems’ metric scores is statistically significant. We used the default parameters of “comet-compare” for paired bootstrap resampling (Koehn, 2004). In the tables reporting on statistical significant testing, the background color is darker for more significant differences (lower p-value) and the score difference is underlined if the p-value is below 0.05.

¹<https://github.com/mjpost/sacrebleu>

System	Rank	chrF	System	Rank	BLEU	System	Rank	COMET
GPT4-5shot	n/a	61.0	GPT4-5shot	n/a	32.8	CUNI-GA	1	90.9
CUNI-GA	1	57.9	CUNI-Transformer	1	30.2	GPT4-5shot	n/a	90.8
GTCOM_Peter	n/a	57.6	GTCOM_Peter	n/a	29.8	ONLINE-W	n/a	89.4
CUNI-Transformer	2	57.4	CUNI-GA	2	29.5	GTCOM_Peter	n/a	88.9
MUNI-NLP	3	57.0	MUNI-NLP	3	28.3	ONLINE-B	n/a	88.8
Lan-BridgeMT	n/a	55.7	Lan-BridgeMT	n/a	27.5	ONLINE-A	n/a	88.2
ONLINE-W	n/a	55.0	ONLINE-W	n/a	26.8	CUNI-Transformer	2	88.0
ONLINE-B	n/a	54.7	ONLINE-B	n/a	25.7	ONLINE-G	n/a	87.7
ONLINE-A	n/a	54.4	ONLINE-A	n/a	25.4	MUNI-NLP	3	87.0
ONLINE-G	n/a	53.7	NLLB_MBR_BLEU	n/a	25.1	ONLINE-Y	n/a	86.5
ONLINE-Y	n/a	53.4	NLLB_Greedy	n/a	24.9	NLLB_Greedy	n/a	86.3
NLLB_Greedy	n/a	52.5	ONLINE-G	n/a	24.8	NLLB_MBR_BLEU	n/a	86.3
NLLB_MBR_BLEU	n/a	52.3	ONLINE-Y	n/a	24.2	Lan-BridgeMT	n/a	86.0

Table 1: Scores for the cs→uk translation task: chrF (nrefs:1lcase:mixedleff:yeslnc:6lnw:0lspace:nolversion:2.2.1), BLEU (nrefs:1lcase:mixedleff:noltok:13alsmooth:explversion:2.2.1), COMET (Unbabel/wmt22-comet-da).

System	Rank	chrF	System	Rank	BLEU	System	Rank	COMET
ONLINE-W	n/a	76.3	ONLINE-W	n/a	59.4	ONLINE-W	n/a	91.8
ONLINE-B	n/a	70.4	ONLINE-B	n/a	50.1	CUNI-GA	1	90.8
ZengHuiMT	n/a	67.5	ONLINE-A	n/a	43.4	ONLINE-B	n/a	89.9
ONLINE-A	n/a	66.3	CUNI-GA	1	43.3	GPT4-5shot	n/a	89.4
CUNI-GA	1	65.9	ZengHuiMT	n/a	43.1	ONLINE-A	n/a	88.4
GTCOM_Peter	n/a	65.4	CUNI-DocTransformer	2	42.5	CUNI-DocTransformer	2	88.3
CUNI-DocTransformer	2	65.1	GTCOM_Peter	n/a	42.3	GTCOM_Peter	n/a	87.7
ONLINE-Y	n/a	64.6	CUNI-Transformer	3	41.4	ONLINE-M	n/a	87.4
CUNI-Transformer	3	63.9	ONLINE-Y	n/a	40.8	Lan-BridgeMT	n/a	87.3
Lan-BridgeMT	n/a	63.8	Lan-BridgeMT	n/a	40.7	CUNI-Transformer	3	87.2
ONLINE-G	n/a	63.7	ONLINE-G	n/a	39.6	NLLB_Greedy	n/a	87.1
ONLINE-M	n/a	63.2	ONLINE-M	n/a	39.6	ONLINE-Y	n/a	87.0
GPT4-5shot	n/a	62.3	GPT4-5shot	n/a	37.8	NLLB_MBR_BLEU	n/a	86.9
NLLB_Greedy	n/a	60.0	NLLB_Greedy	n/a	35.9	ONLINE-G	n/a	85.9
NLLB_MBR_BLEU	n/a	59.1	NLLB_MBR_BLEU	n/a	35.1	ZengHuiMT	n/a	85.4

Table 2: Scores for the en→cs translation task: chrF (nrefs:1lcase:mixedleff:yeslnc:6lnw:0lspace:nolversion:2.2.1), BLEU (nrefs:1lcase:mixedleff:noltok:13alsmooth:explversion:2.2.1), COMET (Unbabel/wmt22-comet-da).

System	Rank	chrF	System	Rank	BLEU	System	Rank	COMET
ONLINE-W	n/a	72.1	ONLINE-W	n/a	51.8	GPT4-5shot	n/a	86.3
ONLINE-A	n/a	70.0	GPT4-5shot	n/a	47.9	ONLINE-W	n/a	86.0
GPT4-5shot	n/a	69.8	ONLINE-A	n/a	47.9	ONLINE-B	n/a	85.6
ONLINE-B	n/a	69.1	ONLINE-B	n/a	46.3	ONLINE-A	n/a	85.5
ONLINE-G	n/a	69.1	ONLINE-G	n/a	46.0	ONLINE-Y	n/a	84.9
ONLINE-Y	n/a	68.4	ONLINE-Y	n/a	43.9	ONLINE-M	n/a	84.8
ZengHuiMT	n/a	67.6	GTCOM_Peter	n/a	42.2	ONLINE-G	n/a	84.6
Lan-BridgeMT	n/a	66.7	Lan-BridgeMT	n/a	42.1	GTCOM_Peter	n/a	82.7
GTCOM_Peter	n/a	66.6	ONLINE-M	n/a	41.3	NLLB_MBR_BLEU	n/a	81.4
ONLINE-M	n/a	66.5	ZengHuiMT	n/a	40.8	ZengHuiMT	n/a	81.1
NLLB_MBR_BLEU	n/a	57.6	NLLB_Greedy	n/a	33.1	Lan-BridgeMT	n/a	80.9
NLLB_Greedy	n/a	57.3	AIRC	1	32.4	NLLB_Greedy	n/a	79.9
AIRC	1	57.2	NLLB_MBR_BLEU	n/a	32.4	AIRC	1	78.7

Table 3: Scores for the de→en translation task: chrF (nrefs:1lcase:mixedleff:yeslnc:6lnw:0lspace:nolversion:2.2.1), BLEU (nrefs:1lcase:mixedleff:noltok:13alsmooth:explversion:2.2.1), COMET (Unbabel/wmt22-comet-da).

System	Rank	chrF	System	Rank	BLEU	System	Rank	COMET
ONLINE-W	n/a	71.8	ONLINE-W	n/a	47.8	ONLINE-W	n/a	85.5
ONLINE-A	n/a	69.7	ONLINE-A	n/a	43.7	GPT4-5shot	n/a	85.0
ZengHuiMT	n/a	69.4	GPT4-5shot	n/a	43.6	ONLINE-B	n/a	84.8
GPT4-5shot	n/a	69.1	ONLINE-Y	n/a	43.6	ONLINE-Y	n/a	84.1
ONLINE-B	n/a	69.1	ONLINE-G	n/a	43.2	ONLINE-A	n/a	83.7
ONLINE-Y	n/a	69.1	ONLINE-B	n/a	42.7	ONLINE-G	n/a	82.5
ONLINE-G	n/a	69.0	ONLINE-M	n/a	40.5	ONLINE-M	n/a	81.7
ONLINE-M	n/a	66.9	ZengHuiMT	n/a	40.5	Lan-BridgeMT	n/a	80.4
Lan-BridgeMT	n/a	66.1	Lan-BridgeMT	n/a	39.4	ZengHuiMT	n/a	79.4
NLLB_Greedy	n/a	56.2	NLLB_Greedy	n/a	31.1	NLLB_MBR_BLEU	n/a	78.0
NLLB_MBR_BLEU	n/a	55.4	NLLB_MBR_BLEU	n/a	29.6	NLLB_Greedy	n/a	77.9
AIRC	1	52.2	AIRC	1	26.5	AIRC	1	72.9

Table 4: Scores for the en→de translation task: chrF (nrefs:1lcase:mixedleff:yeslnc:6lnw:0lspace:nolversion:2.2.1), BLEU (nrefs:1lcase:mixedleff:noltok:13alsmooth:explversion:2.2.1), COMET (Unbabel/wmt22-comet-da).

System	Rank	chrF	System	Rank	BLEU	System	Rank	COMET
ONLINE-B	n/a	87.5	ONLINE-B	n/a	76.5	ONLINE-B	n/a	89.9
ZengHuiMT	n/a	76.3	GTCOM_Peter	n/a	59.2	ONLINE-A	n/a	87.0
GTCOM_Peter	n/a	76.2	ZengHuiMT	n/a	56.6	GPT4-5shot	n/a	86.9
ONLINE-A	n/a	73.3	ONLINE-A	n/a	53.9	GTCOM_Peter	n/a	86.7
GPT4-5shot	n/a	71.4	GPT4-5shot	n/a	51.2	ONLINE-G	n/a	85.6
UvA-LTL	1	70.9	UvA-LTL	1	51.0	ZengHuiMT	n/a	85.6
ONLINE-Y	n/a	70.5	ONLINE-Y	n/a	49.8	ONLINE-Y	n/a	84.9
ONLINE-G	n/a	69.8	ONLINE-G	n/a	49.3	UvA-LTL	1	84.7
NLLB_Greedy	n/a	64.4	NLLB_Greedy	n/a	42.5	NLLB_MBR_BLEU	n/a	82.9
Lan-BridgeMT	n/a	63.5	Lan-BridgeMT	n/a	41.4	NLLB_Greedy	n/a	82.8
NLLB_MBR_BLEU	n/a	63.0	NLLB_MBR_BLEU	n/a	40.7	Samsung_Research_Philippines	2	82.6
Samsung_Research_Philippines	2	55.5	Samsung_Research_Philippines	2	34.0	Lan-BridgeMT	n/a	82.4

Table 5: Scores for the he \rightarrow en translation task: chrF (nrefs:1lcase:mixedleff:yeslnc:6lnw:0lspace:nolversion:2.2.1), BLEU (nrefs:1lcase:mixedleff:noltok:13alsmooth:explversion:2.2.1), COMET (Unbabel/wmt22-comet-da).

System	Rank	chrF	System	Rank	BLEU	System	Rank	COMET
ONLINE-B	n/a	66.4	ONLINE-B	n/a	47.8	ONLINE-B	n/a	86.4
ZengHuiMT	n/a	62.1	ONLINE-A	n/a	38.9	ONLINE-A	n/a	85.7
ONLINE-A	n/a	61.7	GTCOM_Peter	n/a	37.2	GPT4-5shot	n/a	84.9
GTCOM_Peter	n/a	61.1	ONLINE-Y	n/a	37.2	GTCOM_Peter	n/a	84.7
ONLINE-Y	n/a	60.4	ZengHuiMT	n/a	36.5	ONLINE-Y	n/a	84.7
UvA-LTL	1	59.0	UvA-LTL	1	35.0	UvA-LTL	1	84.2
ONLINE-G	n/a	58.1	Samsung_Research_Philippines	2	33.3	Samsung_Research_Philippines	2	83.7
Samsung_Research_Philippines	2	57.3	ONLINE-G	n/a	33.2	Lan-BridgeMT	n/a	83.0
Lan-BridgeMT	n/a	54.9	NLLB_MBR_BLEU	n/a	30.8	NLLB_Greedy	n/a	82.9
NLLB_Greedy	n/a	54.8	Lan-BridgeMT	n/a	30.5	ZengHuiMT	n/a	82.7
NLLB_MBR_BLEU	n/a	54.3	NLLB_Greedy	n/a	30.3	NLLB_MBR_BLEU	n/a	82.5
GPT4-5shot	n/a	54.0	GPT4-5shot	n/a	27.0	ONLINE-G	n/a	82.2

Table 6: Scores for the en \rightarrow he translation task: chrF (nrefs:1lcase:mixedleff:yeslnc:6lnw:0lspace:nolversion:2.2.1), BLEU (nrefs:1lcase:mixedleff:noltok:13alsmooth:explversion:2.2.1), COMET (Unbabel/wmt22-comet-da).

System	Rank	chrF	System	Rank	BLEU	System	Rank	COMET
ONLINE-W	n/a	51.4	ONLINE-W	n/a	25.9	SKIM	1	84.0
GPT4-5shot	n/a	51.2	SKIM	1	24.8	GPT4-5shot	n/a	83.4
SKIM	1	51.1	GPT4-5shot	n/a	24.1	ONLINE-W	n/a	82.3
ONLINE-A	n/a	49.6	ONLINE-B	n/a	23.9	NAIST-NICT	2	81.9
NAIST-NICT	2	49.5	NAIST-NICT	2	23.0	ONLINE-Y	n/a	81.6
ONLINE-Y	n/a	49.5	ONLINE-A	n/a	23.0	ONLINE-B	n/a	81.5
ZengHuiMT	n/a	49.5	ZengHuiMT	n/a	22.6	ONLINE-A	n/a	81.0
ONLINE-B	n/a	49.3	GTCOM_Peter	n/a	22.3	GTCOM_Peter	n/a	80.2
GTCOM_Peter	n/a	48.7	ONLINE-Y	n/a	22.3	ANVITA	3	79.5
Lan-BridgeMT	n/a	47.3	ANVITA	3	20.9	Lan-BridgeMT	n/a	79.3
ANVITA	3	46.7	Lan-BridgeMT	n/a	20.2	ZengHuiMT	n/a	79.2
ONLINE-G	n/a	45.5	ONLINE-G	n/a	18.3	ONLINE-G	n/a	77.8
KYB	n/a	43.9	KYB	n/a	17.6	ONLINE-M	n/a	77.5
ONLINE-M	n/a	43.9	ONLINE-M	n/a	17.2	KYB	n/a	76.6
AIRC	4	40.5	AIRC	4	14.9	NLLB_MBR_BLEU	n/a	75.2
NLLB_MBR_BLEU	n/a	39.2	NLLB_MBR_BLEU	n/a	14.7	AIRC	4	74.5
NLLB_Greedy	n/a	39.0	NLLB_Greedy	n/a	14.2	NLLB_Greedy	n/a	74.3

Table 7: Scores for the ja \rightarrow en translation task: chrF (nrefs:1lcase:mixedleff:yeslnc:6lnw:0lspace:nolversion:2.2.1), BLEU (nrefs:1lcase:mixedleff:noltok:13alsmooth:explversion:2.2.1), COMET (Unbabel/wmt22-comet-da).

System	Rank	chrF	System	Rank	BLEU	System	Rank	COMET
ONLINE-B	n/a	35.2	ONLINE-B	n/a	25.3	ONLINE-B	n/a	88.2
ONLINE-Y	n/a	34.1	ONLINE-W	n/a	24.5	ONLINE-W	n/a	87.5
ONLINE-W	n/a	33.5	ONLINE-Y	n/a	24.5	ONLINE-Y	n/a	87.3
SKIM	1	33.5	SKIM	1	24.3	GPT4-5shot	n/a	87.0
ZengHuiMT	n/a	32.9	NAIST-NICT	2	22.6	SKIM	1	86.6
NAIST-NICT	2	32.0	ZengHuiMT	n/a	22.6	NAIST-NICT	2	86.2
ONLINE-A	n/a	31.4	ONLINE-A	n/a	21.4	ZengHuiMT	n/a	85.3
GPT4-5shot	n/a	31.0	GPT4-5shot	n/a	21.3	ONLINE-A	n/a	85.2
Lan-BridgeMT	n/a	30.4	Lan-BridgeMT	n/a	20.5	Lan-BridgeMT	n/a	84.5
ONLINE-M	n/a	29.6	ONLINE-M	n/a	19.8	ONLINE-M	1	n/a3.3
ANVITA	3	29.3	ANVITA	3	19.4	ANVITA	3	82.7
KYB	n/a	27.7	KYB	n/a	17.8	KYB	n/a	80.8
AIRC	4	27.6	AIRC	4	17.6	AIRC	4	80.7
ONLINE-G	n/a	27.3	ONLINE-G	n/a	17.2	ONLINE-G	n/a	80.4
NLLB_Greedy	n/a	20.9	NLLB_Greedy	n/a	11.3	NLLB_Greedy	n/a	79.3
NLLB_MBR_BLEU	n/a	18.7	NLLB_MBR_BLEU	n/a	9.0	NLLB_MBR_BLEU	n/a	77.7

Table 8: Scores for the en→ja translation task: chrF (nrefs:1lcase:mixedleff:yeslnc:6lnw:0lspace:nolversion:2.2.1), BLEU (nrefs:1lcase:mixedleff:noltok:ja-mecab-0.996-IPAsmooth:explversion:2.2.1), COMET (Unbabel/wmt22-comet-da).

System	Rank	chrF	System	Rank	BLEU	System	Rank	COMET
GPT4-5shot	n/a	60.4	ONLINE-B	n/a	34.5	GPT4-5shot	n/a	83.5
ONLINE-G	n/a	59.6	GPT4-5shot	n/a	34.4	ONLINE-Y	n/a	82.5
ONLINE-A	n/a	59.4	ONLINE-G	n/a	34.0	ONLINE-B	n/a	82.3
ONLINE-B	n/a	59.4	ONLINE-A	n/a	33.8	ONLINE-W	n/a	82.2
ZengHuiMT	n/a	58.9	ONLINE-Y	n/a	33.2	ONLINE-G	n/a	82.0
ONLINE-Y	n/a	58.6	ONLINE-W	n/a	33.1	ONLINE-A	n/a	81.9
PROMT	n/a	58.4	PROMT	n/a	32.8	PROMT	n/a	80.9
ONLINE-W	n/a	58.3	Lan-BridgeMT	n/a	31.8	ONLINE-M	n/a	80.7
Lan-BridgeMT	n/a	57.4	ZengHuiMT	n/a	31.3	NLLB_MBR_BLEU	n/a	80.5
ONLINE-M	n/a	56.7	NLLB_MBR_BLEU	n/a	31.0	NLLB_Greedy	n/a	80.1
NLLB_MBR_BLEU	n/a	55.8	ONLINE-M	n/a	30.7	Lan-BridgeMT	n/a	79.9
NLLB_Greedy	n/a	55.5	NLLB_Greedy	n/a	30.3	ZengHuiMT	n/a	79.5

Table 9: Scores for the ru→en translation task: chrF (nrefs:1lcase:mixedleff:yeslnc:6lnw:0lspace:nolversion:2.2.1), BLEU (nrefs:1lcase:mixedleff:noltok:13alsmooth:explversion:2.2.1), COMET (Unbabel/wmt22-comet-da).

System	Rank	chrF	System	Rank	BLEU	System	Rank	COMET
ONLINE-B	n/a	61.9	ONLINE-B	n/a	40.4	ONLINE-G	n/a	86.6
ONLINE-A	n/a	59.0	ONLINE-A	n/a	34.8	ONLINE-W	n/a	86.6
ONLINE-G	n/a	58.9	ONLINE-G	n/a	32.9	ONLINE-B	n/a	86.2
ZengHuiMT	n/a	58.8	ONLINE-Y	n/a	32.0	GPT4-5shot	n/a	86.1
ONLINE-W	n/a	56.6	ZengHuiMT	n/a	31.6	ONLINE-Y	n/a	85.5
ONLINE-Y	n/a	56.4	ONLINE-W	n/a	31.4	ONLINE-A	n/a	85.3
GPT4-5shot	n/a	56.2	ONLINE-M	n/a	30.9	ONLINE-M	n/a	83.2
Lan-BridgeMT	n/a	55.7	Lan-BridgeMT	n/a	30.7	Lan-BridgeMT	n/a	83.1
PROMT	n/a	55.4	GPT4-5shot	n/a	30.6	NLLB_Greedy	n/a	82.9
ONLINE-M	n/a	55.1	PROMT	n/a	30.5	NLLB_MBR_BLEU	n/a	82.7
NLLB_Greedy	n/a	53.3	NLLB_MBR_BLEU	n/a	28.4	PROMT	n/a	82.3
NLLB_MBR_BLEU	n/a	53.1	NLLB_Greedy	n/a	28.2	ZengHuiMT	n/a	81.3

Table 10: Scores for the en→ru translation task: chrF (nrefs:1lcase:mixedleff:yeslnc:6lnw:0lspace:nolversion:2.2.1), BLEU (nrefs:1lcase:mixedleff:noltok:13alsmooth:explversion:2.2.1), COMET (Unbabel/wmt22-comet-da).

System	Rank	chrF	System	Rank	BLEU	System	Rank	COMET
GTCOM_Peter	n/a	69.3	ONLINE-W	n/a	47.4	ONLINE-W	n/a	87.5
ONLINE-W	n/a	69.2	GTCOM_Peter	n/a	46.4	GPT4-5shot	n/a	87.1
ONLINE-B	n/a	69.0	ONLINE-B	n/a	46.0	ONLINE-B	n/a	86.8
ZengHuiMT	n/a	68.5	ONLINE-A	n/a	45.9	GTCOM_Peter	n/a	86.3
ONLINE-A	n/a	68.3	ONLINE-Y	n/a	45.7	ONLINE-A	n/a	86.3
ONLINE-Y	n/a	68.2	ONLINE-G	n/a	44.9	ONLINE-G	n/a	86.2
GPT4-5shot	n/a	68.1	GPT4-5shot	n/a	43.9	ONLINE-Y	n/a	85.8
ONLINE-G	n/a	68.0	ZengHuiMT	n/a	43.5	Lan-BridgeMT	n/a	84.8
Lan-BridgeMT	n/a	66.2	Lan-BridgeMT	n/a	42.3	ZengHuiMT	n/a	84.4
NLLB_Greedy	n/a	62.4	NLLB_MBR_BLEU	n/a	38.1	NLLB_MBR_BLEU	n/a	84.3
NLLB_MBR_BLEU	n/a	62.4	NLLB_Greedy	n/a	37.8	NLLB_Greedy	n/a	84.2

Table 11: Scores for the uk→en translation task: chrF (nrefs:1lcase:mixedleff:yeslnc:6lnw:0lspace:nolversion:2.2.1), BLEU (nrefs:1lcase:mixedleff:noltok:13alsmooth:explversion:2.2.1), COMET (Unbabel/wmt22-comet-da).

System	Rank	chrF	System	Rank	BLEU	System	Rank	COMET
ONLINE-B	n/a	61.7	ONLINE-B	n/a	39.8	ONLINE-W	n/a	86.7
ONLINE-W	n/a	59.2	ONLINE-W	n/a	34.9	ONLINE-B	n/a	85.6
ZengHuiMT	n/a	56.4	ONLINE-A	n/a	30.3	GPT4-5shot	n/a	85.3
ONLINE-G	n/a	56.1	ONLINE-Y	n/a	29.5	ONLINE-G	n/a	85.3
ONLINE-A	n/a	55.8	ONLINE-G	n/a	28.6	ONLINE-A	n/a	83.2
ONLINE-Y	n/a	55.4	ZengHuiMT	n/a	27.8	ONLINE-Y	n/a	82.9
GTCOM_Peter	n/a	54.4	GTCOM_Peter	n/a	27.5	GTCOM_Peter	n/a	82.1
GPT4-5shot	n/a	53.0	GPT4-5shot	n/a	25.2	NLLB_Greedy	n/a	82.1
Lan-BridgeMT	n/a	51.9	NLLB_MBR_BLEU	n/a	24.9	NLLB_MBR_BLEU	n/a	81.7
NLLB_Greedy	n/a	50.8	Lan-BridgeMT	n/a	24.6	Lan-BridgeMT	n/a	80.4
NLLB_MBR_BLEU	n/a	50.5	NLLB_Greedy	n/a	24.5	ZengHuiMT	n/a	79.0

Table 12: Scores for the en→uk translation task: chrF (nrefs:1lcase:mixedleff:yeslnc:6lnw:0lspace:nolversion:2.2.1), BLEU (nrefs:1lcase:mixedleff:noltok:13alsmooth:explversion:2.2.1), COMET (Unbabel/wmt22-comet-da).

System	Rank	chrF	System	Rank	BLEU	System	Rank	COMET
HW-TSC	1	57.5	HW-TSC	1	33.6	HW-TSC	1	82.8
ONLINE-B	n/a	57.5	ONLINE-B	n/a	33.5	ONLINE-B	n/a	82.7
Yishu	n/a	57.4	Yishu	n/a	33.4	Yishu	n/a	82.7
ZengHuiMT	n/a	54.6	ONLINE-A	n/a	28.3	GPT4-5shot	n/a	81.6
ONLINE-G	n/a	53.9	Lan-BridgeMT	n/a	27.3	Lan-BridgeMT	n/a	81.2
ONLINE-A	n/a	53.4	IOL_Research	2	27.2	ONLINE-G	n/a	80.9
GPT4-5shot	n/a	53.1	ZengHuiMT	n/a	27.0	ONLINE-Y	n/a	80.6
Lan-BridgeMT	n/a	53.1	GPT4-5shot	n/a	26.8	ONLINE-A	n/a	80.3
ONLINE-W	n/a	52.5	ONLINE-G	n/a	26.6	ZengHuiMT	n/a	79.6
IOL_Research	2	52.4	ONLINE-W	n/a	26.4	ONLINE-W	n/a	79.3
ONLINE-Y	n/a	52.3	ONLINE-Y	n/a	25.0	IOL_Research	2	79.2
ONLINE-M	n/a	49.7	ONLINE-M	n/a	23.5	ONLINE-M	n/a	77.7
ANVITA	3	47.1	ANVITA	3	21.8	NLLB_MBR_BLEU	n/a	76.8
NLLB_Greedy	n/a	46.1	NLLB_Greedy	n/a	20.5	ANVITA	3	76.6
NLLB_MBR_BLEU	n/a	45.8	NLLB_MBR_BLEU	n/a	19.8	NLLB_Greedy	n/a	76.4

Table 13: Scores for the zh→en translation task: chrF (nrefs:1lcase:mixedleff:yeslnc:6lnw:0lspace:nolversion:2.2.1), BLEU (nrefs:1lcase:mixedleff:noltok:13alsmooth:explversion:2.2.1), COMET (Unbabel/wmt22-comet-da).

System	Rank	chrF	System	Rank	BLEU	System	Rank	COMET
HW-TSC	1	53.8	HW-TSC	1	58.6	ONLINE-B	n/a	88.1
Yishu	n/a	53.0	ONLINE-A	n/a	58.5	Yishu	n/a	88.1
ONLINE-B	n/a	52.9	Yishu	n/a	57.6	HW-TSC	1	87.3
ONLINE-A	n/a	52.8	ONLINE-B	n/a	57.5	GPT4-5shot	n/a	87.1
IOL_Research	2	51.9	IOL_Research	2	56.9	ONLINE-W	n/a	86.8
ONLINE-M	n/a	50.6	ONLINE-M	n/a	54.9	Lan-BridgeMT	n/a	86.6
ONLINE-Y	n/a	49.8	ONLINE-Y	n/a	54.2	ONLINE-Y	n/a	86.5
ONLINE-G	n/a	49.4	ONLINE-G	n/a	54.1	ONLINE-A	n/a	86.2
ONLINE-W	n/a	47.3	ZengHuiMT	n/a	52.9	IOL_Research	2	85.3
ZengHuiMT	n/a	47.0	ONLINE-W	n/a	52.1	ZengHuiMT	n/a	84.3
Lan-BridgeMT	n/a	46.8	Lan-BridgeMT	n/a	50.2	ONLINE-M	n/a	84.2
GPT4-5shot	n/a	46.5	GPT4-5shot	n/a	49.6	ONLINE-G	n/a	83.8
ANVITA	3	36.9	ANVITA	3	38.9	NLLB_Greedy	n/a	75.7
NLLB_Greedy	n/a	26.3	NLLB_Greedy	n/a	27.4	ANVITA	3	75.6
NLLB_MBR_BLEU	n/a	21.1	NLLB_MBR_BLEU	n/a	19.1	NLLB_MBR_BLEU	n/a	71.5

Table 14: Scores for the en→zh translation task: chrF (nrefs:1lcase:mixedleff:yeslnc:6lnw:0lspace:nolversion:2.2.1), BLEU (nrefs:1lcase:mixedleff:noltok:zhlsmooth:explversion:2.2.1), COMET (Unbabel/wmt22-comet-da).

	GPT4-5sh.	ONLINE-W.	GTCOM_Pe.	ONLINE-B.	ONLINE-A.	CUNI-Tra.	ONLINE-G.	MUNI-NLP.	ONLINE-Y.	NLLB_Gre.	NLLB_MBR.	Lan-Brid.
CUNI-GA	0.1	<u>1.5</u>	<u>2.0</u>	<u>2.1</u>	<u>2.7</u>	<u>2.9</u>	<u>3.2</u>	<u>3.9</u>	<u>4.4</u>	<u>4.6</u>	<u>4.6</u>	<u>4.9</u>
GPT4-5shot	0.0	<u>1.4</u>	<u>1.9</u>	<u>2.0</u>	<u>2.6</u>	<u>2.8</u>	<u>3.1</u>	<u>3.8</u>	<u>4.3</u>	<u>4.5</u>	<u>4.5</u>	<u>4.8</u>
ONLINE-W		0.0	<u>0.5</u>	<u>0.6</u>	<u>1.2</u>	<u>1.4</u>	<u>1.7</u>	<u>2.4</u>	<u>2.9</u>	<u>3.1</u>	<u>3.1</u>	<u>3.4</u>
GTCOM_Peter			0.0	0.1	<u>0.7</u>	<u>0.9</u>	<u>1.2</u>	<u>1.9</u>	<u>2.4</u>	<u>2.6</u>	<u>2.6</u>	<u>2.9</u>
ONLINE-B				0.0	<u>0.6</u>	<u>0.8</u>	<u>1.1</u>	<u>1.8</u>	<u>2.3</u>	<u>2.5</u>	<u>2.5</u>	<u>2.8</u>
ONLINE-A					0.0	<u>0.2</u>	<u>0.5</u>	<u>1.2</u>	<u>1.7</u>	<u>1.9</u>	<u>1.9</u>	<u>2.2</u>
CUNI-Transformer						0.0	<u>0.3</u>	<u>1.0</u>	<u>1.5</u>	<u>1.7</u>	<u>1.7</u>	<u>2.0</u>
ONLINE-G							0.0	<u>0.7</u>	<u>1.2</u>	<u>1.4</u>	<u>1.4</u>	<u>1.7</u>
MUNI-NLP								0.0	<u>0.5</u>	<u>0.7</u>	<u>0.7</u>	<u>1.0</u>
ONLINE-Y									0.0	0.2	0.2	<u>0.5</u>
NLLB_Greedy										0.0	0.0	0.3
NLLB_MBR_BLEU											0.0	0.3

Table 15: Statistical significance testing of the COMET score difference for each system pair for the es→uk.

	CUNI-GA	ONLINE-B.	GPT4-5sh.	ONLINE-A.	CUNI-Doc.	GTCom_Pe.	ONLINE-M.	Lan-Brid.	CUNI-Tra	NLLB_Gre.	ONLINE-Y.	NLLB_MBR.	ONLINE-G.	ZengHuiM.
ONLINE-W	1.0	1.9	2.4	3.4	3.5	4.1	4.4	4.5	4.6	4.7	4.8	4.9	5.9	6.4
CUNI-GA	0.0	0.9	1.4	2.4	2.5	3.1	3.4	3.5	3.6	3.7	3.8	3.9	4.9	5.4
ONLINE-B		0.0	0.5	1.5	1.6	2.2	2.5	2.6	2.7	2.8	2.9	3.0	4.0	4.5
GPT4-5shot			0.0	1.0	1.1	1.7	2.0	2.1	2.2	2.3	2.4	2.5	3.5	4.0
ONLINE-A				0.0	0.1	0.7	1.0	1.1	1.2	1.3	1.4	1.5	2.5	3.0
CUNI-DocTransformer					0.0	0.6	0.9	1.0	1.1	1.2	1.3	1.4	2.4	2.9
GTCom_Peter						0.0	0.3	0.4	0.5	0.6	0.7	0.8	1.8	2.3
ONLINE-M							0.0	0.1	0.2	0.3	0.4	0.5	1.5	2.0
Lan-BridgeMT								0.0	0.1	0.2	0.3	0.4	1.4	1.9
CUNI-Transformer									0.0	0.1	0.2	0.3	1.3	1.8
NLLB_Greedy										0.0	0.1	0.2	1.2	1.7
ONLINE-Y											0.0	0.1	1.1	1.6
NLLB_MBR_BLEU												0.0	1.0	1.5
ONLINE-G													0.0	0.5

Table 16: Statistical significance testing of the COMET score difference for each system pair for the en→cs.

	ONLINE-W.	ONLINE-B.	ONLINE-A.	ONLINE-Y.	ONLINE-M.	ONLINE-G.	GTCOM_Pe.	NLLB_MBR.	ZengHuiM.	Lan-Brid.	NLLB_Gre.	AIRC.
GPT4-5shot	0.3	0.7	0.8	1.4	1.5	1.7	3.6	4.9	5.2	5.4	6.4	7.6
ONLINE-W	0.0	0.4	0.5	1.1	1.2	1.4	3.3	4.6	4.9	5.1	6.1	7.3
ONLINE-B		0.0	0.1	0.7	0.8	1.0	2.9	4.2	4.5	4.7	5.7	6.9
ONLINE-A			0.0	0.6	0.7	0.9	2.8	4.1	4.4	4.6	5.6	6.8
ONLINE-Y				0.0	0.1	0.3	2.2	3.5	3.8	4.0	5.0	6.2
ONLINE-M					0.0	0.2	2.1	3.4	3.7	3.9	4.9	6.1
ONLINE-G						0.0	1.9	3.2	3.5	3.7	4.7	5.9
GTCOM_Peter							0.0	1.3	1.6	1.8	2.8	4.0
NLLB_MBR_BLEU								0.0	0.3	0.5	1.5	2.7
ZengHuiMT									0.0	0.2	1.2	2.4
Lan-BridgeMT										0.0	1.0	2.2
NLLB_Greedy											0.0	1.2

Table 17: Statistical significance testing of the COMET score difference for each system pair for the de→en.

	GPT4-5sh.	ONLINE-B.	ONLINE-Y.	ONLINE-A.	ONLINE-G.	ONLINE-M.	Lan-Brid.	ZengHuiM.	NLLB_MBR.	NLLB_Gre.	AIRC.
ONLINE-W 0.0	0.5	0.7	1.4	1.8	3.0	3.8	5.1	6.1	7.5	7.6	12.6
GPT4-5shot	0.0	0.2	0.9	1.3	2.5	3.3	4.6	5.6	7.0	7.1	12.1
ONLINE-B		0.0	0.7	1.1	2.3	3.1	4.4	5.4	6.8	6.9	11.9
ONLINE-Y			0.0	0.4	1.6	2.4	3.7	4.7	6.1	6.2	11.2
ONLINE-A				0.0	1.2	2.0	3.3	4.3	5.7	5.8	10.8
ONLINE-G					0.0	0.8	2.1	3.1	4.5	4.6	9.6
ONLINE-M						0.0	1.3	2.3	3.7	3.8	8.8
Lan-BridgeMT							0.0	1.0	2.4	2.5	7.5
ZengHuiMT								0.0	1.4	1.5	6.5
NLLB_MBR_BLEU									0.0	0.1	5.1
NLLB_Greedy										0.0	5.0

Table 18: Statistical significance testing of the COMET score difference for each system pair for the en→de.

	ONLINE-A.	GPT4-5sh.	GTCOM_Pe.	ONLINE-G.	ZengHuiM.	ONLINE-Y.	UvA-LTL.	NLLB_MBR.	NLLB_Gre.	Samsung_.	Lan-Brid.
ONLINE-B 0.0	2.9	3.0	3.2	4.3	4.3	5.0	5.2	7.0	7.1	7.3	7.5
ONLINE-A	0.0	0.1	0.3	1.4	1.4	2.1	2.3	4.1	4.2	4.4	4.6
GPT4-5shot		0.0	0.2	1.3	1.3	2.0	2.2	4.0	4.1	4.3	4.5
GTCOM_Peter			0.0	1.1	1.1	1.8	2.0	3.8	3.9	4.1	4.3
ONLINE-G				0.0	0.0	0.7	0.9	2.7	2.8	3.0	3.2
ZengHuiMT					0.0	0.7	0.9	2.7	2.8	3.0	3.2
ONLINE-Y						0.0	0.2	2.0	2.1	2.3	2.5
UvA-LTL							0.0	1.8	1.9	2.1	2.3
NLLB_MBR_BLEU								0.0	0.1	0.3	0.5
NLLB_Greedy									0.0	0.2	0.4
Samsung_Research_Philippines										0.0	0.2

Table 19: Statistical significance testing of the COMET score difference for each system pair for the he→en.

	ONLINE-A.	GPT4-5sh.	GTCOM_Pe.	ONLINE-Y.	UvA-LTL.	Samsung_.	Lan-Brid.	NLLB_Gre.	ZengHuiM.	NLLB_MBR.	ONLINE-G.
ONLINE-B	0.7	1.5	1.7	1.7	2.2	2.7	3.4	3.5	3.7	3.9	4.2
ONLINE-A	0.0	0.8	1.0	1.0	1.5	2.0	2.7	2.8	3.0	3.2	3.5
GPT4-5shot		0.0	0.2	0.2	0.7	1.2	1.9	2.0	2.2	2.4	2.7
GTCOM_Peter			0.0	0.0	0.5	1.0	1.7	1.8	2.0	2.2	2.5
ONLINE-Y				0.0	0.5	1.0	1.7	1.8	2.0	2.2	2.5
UvA-LTL					0.0	0.5	1.2	1.3	1.5	1.7	2.0
Samsung_Research_Philippines						0.0	0.7	0.8	1.0	1.2	1.5
Lan-BridgeMT							0.0	0.1	0.3	0.5	0.8
NLLB_Greedy								0.0	0.2	0.4	0.7
ZengHuiMT									0.0	0.2	0.5
NLLB_MBR_BLEU										0.0	0.3

Table 20: Statistical significance testing of the COMET score difference for each system pair for the en→he.

	GPT4-5sh.	ONLINE-W.	NAIST-NI.	ONLINE-Y.	ONLINE-B.	ONLINE-A.	GTCOM_Pe.	ANVITA.	Lan-Brid.	ZengHuiM.	ONLINE-G.	ONLINE-M.	KYB.	NLLB_MBR.	AIRC.	NLLB_Gre.
SKIM	0.6	1.7	2.1	2.4	2.5	3.0	3.8	4.5	4.7	4.8	6.2	6.5	7.4	8.8	9.5	9.7
GPT4-5shot	0.0	1.1	1.5	1.8	1.9	2.4	3.2	3.9	4.1	4.2	5.6	5.9	6.8	8.2	8.9	9.1
ONLINE-W		0.0	0.4	0.7	0.8	1.3	2.1	2.8	3.0	3.1	4.5	4.8	5.7	7.1	7.8	8.0
NAIST-NICT			0.0	0.3	0.4	0.9	1.7	2.4	2.6	2.7	4.1	4.4	5.3	6.7	7.4	7.6
ONLINE-Y				0.0	0.1	0.6	1.4	2.1	2.3	2.4	3.8	4.1	5.0	6.4	7.1	7.3
ONLINE-B					0.0	0.5	1.3	2.0	2.2	2.3	3.7	4.0	4.9	6.3	7.0	7.2
ONLINE-A						0.0	0.8	1.5	1.7	1.8	3.2	3.5	4.4	5.8	6.5	6.7
GTCOM_Peter							0.0	0.7	0.9	1.0	2.4	2.7	3.6	5.0	5.7	5.9
ANVITA								0.0	0.2	0.3	1.7	2.0	2.9	4.3	5.0	5.2
Lan-BridgeMT									0.0	0.1	1.5	1.8	2.7	4.1	4.8	5.0
ZengHuiMT										0.0	1.4	1.7	2.6	4.0	4.7	4.9
ONLINE-G											0.0	0.3	1.2	2.6	3.3	3.5
ONLINE-M												0.0	0.9	2.3	3.0	3.2
KYB													0.0	1.4	2.1	2.3
NLLB_MBR_BLEU														0.0	0.7	0.9
AIRC															0.0	0.2

Table 21: Statistical significance testing of the COMET score difference for each system pair for the ja→en.

	ONLINE-W.	ONLINE-Y.	GPT4-5sh.	SKIM.	NAIST-NI.	ZengHuiM.	ONLINE-A.	Lan-Brid.	ONLINE-M.	ANVITA.	KYB.	AIRC.	ONLINE-G.	NLLB_Gre.	NLLB_MBR.
ONLINE-B	0.7	0.9	1.2	1.6	2.0	2.9	3.0	3.7	4.9	5.5	7.4	7.5	7.8	8.9	10.5
ONLINE-W	0.0	0.2	0.5	0.9	1.3	2.2	2.3	3.0	4.2	4.8	6.7	6.8	7.1	8.2	9.8
ONLINE-Y		0.0	0.3	0.7	1.1	2.0	2.1	2.8	4.0	4.6	6.5	6.6	6.9	8.0	9.6
GPT4-5shot			0.0	0.4	0.8	1.7	1.8	2.5	3.7	4.3	6.2	6.3	6.6	7.7	9.3
SKIM				0.0	0.4	1.3	1.4	2.1	3.3	3.9	5.8	5.9	6.2	7.3	8.9
NAIST-NICT					0.0	0.9	1.0	1.7	2.9	3.5	5.4	5.5	5.8	6.9	8.5
ZengHuiMT						0.0	0.1	0.8	2.0	2.6	4.5	4.6	4.9	6.0	7.6
ONLINE-A							0.0	0.7	1.9	2.5	4.4	4.5	4.8	5.9	7.5
Lan-BridgeMT								0.0	1.2	1.8	3.7	3.8	4.1	5.2	6.8
ONLINE-M									0.0	0.6	2.5	2.6	2.9	4.0	5.6
ANVITA										0.0	1.9	2.0	2.3	3.4	5.0
KYB											0.0	0.1	0.4	1.5	3.1
AIRC												0.0	0.3	1.4	3.0
ONLINE-G													0.0	1.1	2.7
NLLB_Greedy														0.0	1.6

Table 22: Statistical significance testing of the COMET score difference for each system pair for the en→ja.

	ONLINE-Y.	ONLINE-B.	ONLINE-W.	ONLINE-G.	ONLINE-A.	PROMT.	ONLINE-M.	NLLB_MBR.	NLLB_Gre.	Lan-Brid.	ZengHuiM.
GPT4-5shot	1.0	1.2	1.3	1.5	1.6	2.6	2.8	3.0	3.4	3.6	4.0
ONLINE-Y	0.0	0.2	0.3	0.5	0.6	1.6	1.8	2.0	2.4	2.6	3.0
ONLINE-B		0.0	0.1	0.3	0.4	1.4	1.6	1.8	2.2	2.4	2.8
ONLINE-W			0.0	0.2	0.3	1.3	1.5	1.7	2.1	2.3	2.7
ONLINE-G				0.0	0.1	1.3	1.3	1.5	1.9	2.1	2.5
ONLINE-A					0.0	1.1	1.2	1.4	1.8	2.0	2.4
PROMT						1.0	0.2	0.4	0.8	1.0	1.4
ONLINE-M						0.0	0.0	0.2	0.6	0.8	1.2
NLLB_MBR_BLEU								0.0	0.4	0.6	1.0
NLLB_Greedy									0.0	0.2	0.6
Lan-BridgeMT										0.0	0.4

Table 23: Statistical significance testing of the COMET score difference for each system pair for the ru→en.

	ONLINE-W.	ONLINE-B.	GPT4-5sh.	ONLINE-Y.	ONLINE-A.	ONLINE-M.	Lan-Brid.	NLLB_Gre.	NLLB_MBR.	PROMT.	ZengHuM.
ONLINE-G	0.0	0.4	0.5	1.1	1.3	3.4	3.5	3.7	3.9	4.3	5.3
ONLINE-W	0.0	0.4	0.5	1.1	1.3	3.4	3.5	3.7	3.9	4.3	5.3
ONLINE-B		0.0	0.1	0.7	0.9	3.0	3.1	3.3	3.5	3.9	4.9
GPT4-5shot			0.0	0.6	0.8	2.9	3.0	3.2	3.4	3.8	4.8
ONLINE-Y				0.0	0.2	2.3	2.4	2.6	2.8	3.2	4.2
ONLINE-A				0.0	0.0	2.1	2.2	2.4	2.6	3.0	4.0
ONLINE-M						0.0	0.1	0.3	0.5	0.9	1.9
Lan-BridgeMT							0.0	0.2	0.4	0.8	1.8
NLLB_Greedy								0.0	0.2	0.6	1.6
NLLB_MBR_BLEU									0.0	0.4	1.4
PROMT										0.0	1.0

Table 24: Statistical significance testing of the COMET score difference for each system pair for the en→ru.

	GPT4-5sh.	ONLINE-B.	GTCOM_Pe.	ONLINE-A.	ONLINE-G.	ONLINE-Y.	Lan-Brid.	ZengHuiM.	NLLB_MBR.	NLLB_Gre.
ONLINE-W	0.4	0.7	1.2	1.2	1.3	1.7	2.7	3.1	3.2	3.3
GPT4-5shot	0.0	0.3	0.8	0.8	0.9	1.3	2.3	2.7	2.8	2.9
ONLINE-B		0.0	0.5	0.5	0.6	1.0	2.0	2.4	2.5	2.6
GTCOM_Peter			0.0	0.0	0.1	0.5	1.5	1.9	2.0	2.1
ONLINE-A				0.0	0.1	0.5	1.5	1.9	2.0	2.1
ONLINE-G					0.0	0.4	1.4	1.8	1.9	2.0
ONLINE-Y						0.0	1.0	1.4	1.5	1.6
Lan-BridgeMT							0.0	0.4	0.5	0.6
ZengHuiMT								0.0	0.1	0.2
NLLB_MBR_BLEU									0.0	0.1

Table 25: Statistical significance testing of the COMET score difference for each system pair for the uk→en.

	ONLINE-B.	GPT4-5sh.	ONLINE-G.	ONLINE-A.	ONLINE-Y.	GTCOM_Pe.	NLLB_Gre.	NLLB_MBR.	Lan-Brid.	ZengHuiM.
ONLINE-W	1.1	1.4	1.4	3.5	3.8	4.6	4.6	5.0	6.3	7.7
ONLINE-B	0.0	0.3	0.3	2.4	2.7	3.5	3.5	3.9	5.2	6.6
GPT4-5shot		0.0	0.0	2.1	2.4	3.2	3.2	3.6	4.9	6.3
ONLINE-G			0.0	2.1	2.4	3.2	3.2	3.6	4.9	6.3
ONLINE-A				0.0	0.3	1.1	1.1	1.5	2.8	4.2
ONLINE-Y					0.0	0.8	0.8	1.2	2.5	3.9
GTCOM_Peter						0.0	0.0	0.4	1.7	3.1
NLLB_Greedy							0.0	0.4	1.7	3.1
NLLB_MBR_BLEU								0.0	1.3	2.7
Lan-BridgeMT									0.0	1.4

Table 26: Statistical significance testing of the COMET score difference for each system pair for the en→uk.

	ONLINE-B.	Yishu.	GPT4-5sh.	Lan-Brid.	ONLINE-G.	ONLINE-Y.	ONLINE-A.	ZengHuiM.	ONLINE-W.	IOL_Rese.	ONLINE-M.	NLLB_MBR.	ANVITA.	NLLB_Gre.
HW-TSC	0.1	0.1	<u>1.2</u>	<u>1.6</u>	<u>1.9</u>	<u>2.2</u>	<u>2.5</u>	<u>3.2</u>	<u>3.5</u>	<u>3.6</u>	<u>5.1</u>	<u>6.0</u>	<u>6.2</u>	<u>6.4</u>
ONLINE-B	0.0	0.0	<u>1.1</u>	<u>1.5</u>	<u>1.8</u>	<u>2.1</u>	<u>2.4</u>	<u>3.1</u>	<u>3.4</u>	<u>3.5</u>	<u>5.0</u>	<u>5.9</u>	<u>6.1</u>	<u>6.3</u>
Yishu	0.0	0.0	<u>1.1</u>	<u>1.5</u>	<u>1.8</u>	<u>2.1</u>	<u>2.4</u>	<u>3.1</u>	<u>3.4</u>	<u>3.5</u>	<u>5.0</u>	<u>5.9</u>	<u>6.1</u>	<u>6.3</u>
GPT4-5shot			0.0	0.4	0.7	1.0	1.3	2.0	2.3	2.4	3.9	4.8	5.0	5.2
Lan-BridgeMT				0.0	0.3	<u>0.6</u>	<u>0.9</u>	<u>1.6</u>	<u>1.9</u>	<u>2.0</u>	<u>3.5</u>	<u>4.4</u>	<u>4.6</u>	<u>4.8</u>
ONLINE-G				0.0	0.3	0.0	0.6	1.3	1.6	1.7	3.2	4.1	4.3	4.5
ONLINE-Y					0.0	0.0	0.3	<u>1.0</u>	<u>1.3</u>	<u>1.4</u>	<u>2.9</u>	<u>3.8</u>	<u>4.0</u>	<u>4.2</u>
ONLINE-A							0.0	0.7	1.0	1.1	2.6	3.5	3.7	3.9
ZengHuiMT								0.0	0.3	0.4	<u>1.9</u>	<u>2.8</u>	<u>3.0</u>	<u>3.2</u>
ONLINE-W									0.0	0.1	<u>1.6</u>	<u>2.5</u>	<u>2.7</u>	<u>2.9</u>
IOL_Research										0.0	<u>1.5</u>	<u>2.4</u>	<u>2.6</u>	<u>2.8</u>
ONLINE-M											0.0	<u>0.9</u>	<u>1.1</u>	<u>1.3</u>
NLLB_MBR_BLEU												0.0	0.2	0.4
ANVITA													0.0	0.2

Table 27: Statistical significance testing of the COMET score difference for each system pair for the zh→en.

	HW-TSC.	GPT4-5sh.	ONLINE-W.	Lan-Brid.	ONLINE-Y.	ONLINE-A.	IOL_Rese.	ZengHuiM.	ONLINE-M.	ONLINE-G.	NLLB_Gre.	ANVITA.	NLLB_MBR.
ONLINE-B	0.0	<u>0.8</u>	<u>1.0</u>	<u>1.3</u>	<u>1.5</u>	<u>1.6</u>	<u>1.9</u>	<u>2.8</u>	<u>3.8</u>	<u>3.9</u>	<u>4.3</u>	<u>12.4</u>	<u>16.6</u>
Yishu	0.0	<u>0.8</u>	<u>1.0</u>	<u>1.3</u>	<u>1.5</u>	<u>1.6</u>	<u>1.9</u>	<u>2.8</u>	<u>3.8</u>	<u>3.9</u>	<u>4.3</u>	<u>12.4</u>	<u>16.6</u>
HW-TSC		0.0	0.2	<u>0.5</u>	<u>0.7</u>	<u>0.8</u>	<u>1.1</u>	<u>2.0</u>	<u>3.0</u>	<u>3.1</u>	<u>3.5</u>	<u>11.6</u>	<u>15.8</u>
GPT4-5shot			0.0	0.3	<u>0.5</u>	<u>0.6</u>	<u>0.9</u>	<u>1.8</u>	<u>2.8</u>	<u>2.9</u>	<u>3.3</u>	<u>11.4</u>	<u>15.6</u>
ONLINE-W				0.0	<u>0.2</u>	<u>0.3</u>	<u>0.6</u>	<u>1.5</u>	<u>2.5</u>	<u>2.6</u>	<u>3.0</u>	<u>11.1</u>	<u>15.3</u>
Lan-BridgeMT					0.0	0.1	<u>0.4</u>	<u>1.3</u>	<u>2.3</u>	<u>2.4</u>	<u>2.8</u>	<u>10.9</u>	<u>15.1</u>
ONLINE-Y						0.0	<u>0.3</u>	<u>1.2</u>	<u>2.2</u>	<u>2.3</u>	<u>2.7</u>	<u>10.8</u>	<u>15.0</u>
ONLINE-A							0.0	<u>0.9</u>	<u>1.9</u>	<u>2.0</u>	<u>2.4</u>	<u>10.5</u>	<u>14.7</u>
IOL_Research								0.0	<u>1.0</u>	<u>1.1</u>	<u>1.5</u>	<u>9.6</u>	<u>13.8</u>
ZengHuiMT									0.0	0.1	<u>0.5</u>	<u>8.6</u>	<u>12.8</u>
ONLINE-M										0.0	<u>0.4</u>	<u>8.5</u>	<u>12.7</u>
ONLINE-G											0.0	<u>8.1</u>	<u>12.3</u>
NLLB_Greedy												0.0	<u>4.2</u>
ANVITA												0.0	<u>4.1</u>

Table 28: Statistical significance testing of the COMET score difference for each system pair for the en→zh.

References

- Philipp Koehn. 2004. [Statistical significance tests for machine translation evaluation](#). In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 388–395, Barcelona, Spain. Association for Computational Linguistics.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. [Bleu: a method for automatic evaluation of machine translation](#). In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Maja Popović. 2015. [chrF: character n-gram F-score for automatic MT evaluation](#). In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.
- Matt Post. 2018. [A call for clarity in reporting BLEU scores](#). In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.
- Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. [COMET: A neural framework for MT evaluation](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.