BiVert: Bidirectional Vocabulary Evaluation using Relations for Machine Translation

Carinne Cherf, Yuval Pinter LREC-Coling 2024



Initial Research

Searching for translation mistakes through homonyms.

What are homonyms? What were we looking for? almost, generally, roughly, what are other relatively, about, around, ballpark figure, circa, closely, words for trunk very close? comparatively park bark bank ock 👁 👁 SAW D

🚺 Thesaurus.plus

Evaluation Methods Today

Reference Based

• BLEU, ROUGE



• BertScore



• MAUVE, MoverScore





3

Evaluation Methods Today

Quality Estimation

CometKiwi











Back-translation

 $l_1: {
m EN}$ $l_2: {
m RU}$ s: Public debt will balloon, creating financial challenges around the world. $\int {
m DT}$ t: Государственный долг будет расти, создавая финансовые проблемы по всему миру. $\int {
m BT}$

s': Public debt will grow, creating financial problems around the world.

Figure 1: Example of a direct translation from English to Russian using the system we wish to evaluate, and its back-translation using a state-of-the-art translation system suitable for BIVERT.

Word Pair Relations

s: Public debt will balloon, creating financial challenges around the world.

Figure 2: An example of final words alignment using the linear sum assignment problem algorithm.

Categories:

- 1. Same
- 2. Extra

6. Derivation: happy - happiness

5. Inflection: eat - ate (same lemma)

3. Missing

b. Derivation: happy - ha7. Sense

4. Stop Words

Word Pair Alignment

The students read quietly in the library

They quietly studied books in the library

	The	students	read
They	1- 0.3507	1- 0.4420	1- 0.4078
studied	1- 0.3738	1- 0.2949	1- 0.3170
books	1- 0.3866	1- 0.3187	1- 0.3178

Categories:

- 1. Same
- 2. Extra
- 3. Missing
- 4. Stop Words

- 5. Inflection: eat ate (same lemma)
- 6. Derivation: happy happiness
- 7. Sense

Word Pair Alignment



Option 1:



Option 2:

Step 1:Token PairingAlignment algorithm between
embeddings of individual tokens_in \leftrightarrow _un

Step 2: Concatenate following tokens of starting tokens to form full word pairs.

Sense Relation - BabelNet



Sense Relation



Figure 4.3: Fragment of a semantic graph between the two words *challenge* and *problem*. The hatched grey edges connect roots to their senses, and the red edges represent hypernym relations between the nodes contents. Shortest path highlighted.

Sense Relation

$$w(a,b) = \frac{w(a \rightarrow_r b) + w(b \rightarrow_{r^{-1}} a)}{2d}$$
$$w(x \rightarrow_r y) = \max_r -\frac{\max_r - \min_r}{n_r(X)}$$

where \rightarrow_r is a relation of type r and r^{-1} is its inverse; d is the depth of the deeper of the two nodes; \max_r and \min_r are the maximum and minimum weights possible for a relation of type r; and $n_r(X)$ is the number of relations of type r leaving node X. The final graph score S(a,b) from root a to b is given by the normalized sum of the edge weights along the path between them:

$$S(a,b) = 2 \times \left(0.5 - \frac{1}{\sum_{e \in P(a \rightsquigarrow b)} w(e)}\right)$$
(By Michael Sussna,1993)

Final Score

1. Same	= 0
2. Extra	= $1/len(s)$
3. Missing	= $1/len(s)$
4. Stop Words	= $1/len(s)$
5. Inflection	= cosine similarity
6. Derivation	= cosine similarity
7. Sense	= $S(a,b) = 2 \times \left(0.5 - \frac{1}{\sum_{e \in P(a \rightsquigarrow b)} w(e)}\right)$

 $score_t = \sum_{i=1}^{n} f_t(a, b).$

Where t is a category in set: {Same, Inflection, Extra, Derivation, $(\overline{a,b}) \in t$ Missing, Sense, Stop Words}

$$bivert_similarity(s, s') = \sum_{t \in Types} \alpha_t \times score_t$$

Experiments and Results

- WMT Metrics Task 2021 (train), 2022 (test)
- State-of-the-art MT system: Marian MT
- Gradient Boosting Regression

	Train	Predict
English-German	19,501	19,725
English-Russian	$12,\!000$	19,725
Chinese-English	$16,\!124$	$28,\!124$

	Extra	Missing	Stopword	Inflection	Derivation	Sense
English-German	0.121	0.134	0.188	0.101	0.092	0.360
English-Russian	0.112	0.164	0.196	0.087	0.063	0.375
Chinese-English	0.172	0.203	0.126	0.000	0.000	0.497

Experiments and Results

	System	Human	BiVert	2
	bleu_bestmbr	9.615	0.614	Languaga pair
	bleurt_bestmbr	9.555	0.609	Language pan
	comet_bestmbr	9.567	0.627	Human Translation Included
	JDExploreAcademy	9.581	0.620	indinan indindición indiada
	Lan-Bridge	9.435	0.608	DEDTO
	M2M100_1.2D-D4 Opling A	0.072	0.598	BERIScore
$English \rightarrow Cormon$	Online B	0.585	0.626	Cross OF
⊔ingiisii → German	Online-G	9,510	0.605	CIOSS-QL
	Online-W	9.684	0.624	COMETKiwi
	Online-Y	9.480	0.615	
	OpenNMT	9.329	0.603	MS-COMET-QE-22
	PROMT	9.297	0.596	
	QUARTZ_TuneReranking	9.462	0.622	UniTE-src
	refB	9.634	0.614	MATESE-OE
	bleu_bestmbr	9.715	0.296	
	comet_bestmbr	9.677	0.286	COMET-QE
	e franslation	9.417	0.295	KC BEDTScore
	IDEvelope A codemu	9.470	0.292	NG-DENT SCOLE
	J DExploreAcademy	9.079	0.279	HWTSC-TLM
	M2M100_1_2B_B4	0.208	0.230	
$English \rightarrow Russian$	Online-A	9.561	0.212	HWTSC-Teacher-Sim
English / Hussian	Online-B	9.701	0.310	DUV
	Online-G	9.687	0.333	BiVert
	Online-W	9.789	0.320	
	Online-Y	9.608	0.263	
	PROMT	9.548	0.296	Table 5.4. System lavel I
	QUARTZ_TuneReranking	9.375	0.310	Table 5.4. System-level I
	SRPOL	9.434	0.305	RIVERT scores compared
	AISP-SJTU	9.682	0.443	Diventi scores, compared
	$bleu_bestmbr$	9.701	0.435	tion Included" refers to one
	bleurt_bestmbr	9.749	0.452	tion monuted refers to one
	$comet_bestmbr$	9.714	0.459	excluded from the total con
	HuaweiTSC	9.692	0.443	
	JDExploreAcademy	9.718	0.446	Table 5.2 Highest reference
	Lan-Bridge	9.753	0.460	rabie 5.2. mgnest reference
Chinese \rightarrow English	LanguageX	9 727	0.434	

Language pair	eng-deu	eng-deu	eng-rus	zho-eng	zho-eng
Human Translation Included	yes	no	no	yes	no
BERTScore	0.338	0.428	0.811	0.843	0.924
Cross-QE	0.643	0.661	0.806	0.817	0.870
COMETKiwi	0.592	0.674	0.763	0.795	0.866
MS-COMET-QE-22	0.417	0.539	0.672	0.799	0.897
UniTE-src	0.509	0.509	0.779	0.791	0.874
MATESE-QE	0.363	0.337	0.637	0.741	0.767
COMET-QE	0.480	0.502	0.468	0.544	0.569
KG-BERTScore	0.369	0.400	0.612	0.617	0.743
HWTSC-TLM	0.311	0.428	0.597	0.368	0.460
HWTSC-Teacher-Sim	0.290	0.385	0.675	0.294	0.356
BiVert	0.694	0.703	0.657	0.376	0.239

Table 5.4: System-level Pearson correlation between human scores and BIVERT scores, compared to other evaluation metrics. "Human Translation Included" refers to one other system, refB which may be included or excluded from the total correlation calculation. See system-level scores in Table 5.2. Highest reference-free scores are **bolded**.

Future Work

- Idioms and phrases
- Language categories
- Switching between system roles

