

## A Supplemental Material

### A.1 Proofs

In this section, we prove the Prop. 1 in the paper about viewing BERTScore (precision/recall) as a (non-optimized) Mover Distance.

As a reminder, WMD formulates as:

$$\text{WMD}(\mathbf{x}^n, \mathbf{y}^n) := \min_{F \in \mathbb{R}^{|\mathbf{x}^n| \times |\mathbf{y}^n|}} \sum_i \sum_j C_{ij} \cdot F_{ij}$$

$$\text{s.t. } \mathbf{1}^\top F \mathbf{1} = 1, \quad \mathbf{1}^\top F \mathbf{1} = 1.$$

where  $F^\top \mathbf{1} = \mathbf{f}_x^n$  and  $F \mathbf{1} = \mathbf{f}_y^n$  denote a vector of weights for each  $n$ -gram of  $\mathbf{x}^n$  and  $\mathbf{y}^n$ .

The BERTScore is defined as:

$$R_{\text{BERT}} = \frac{\sum_{y_i^1 \in \mathbf{y}^1} \text{idf}(y_i^1) \max_{x_j^1 \in \mathbf{x}^1} E(x_j^1)^\top E(y_i^1)}{\sum_{y_i^1 \in \mathbf{y}^1} \text{idf}(y_i^1)}$$

$$P_{\text{BERT}} = \frac{\sum_{x_j^1 \in \mathbf{x}^1} \text{idf}(x_j^1) \max_{y_i^1 \in \mathbf{y}^1} E(y_i^1)^\top E(x_j^1)}{\sum_{x_j^1 \in \mathbf{x}^1} \text{idf}(x_j^1)}$$

$$F_{\text{BERT}} = 2 \frac{P_{\text{BERT}} \cdot R_{\text{BERT}}}{P_{\text{BERT}} + R_{\text{BERT}}}.$$

Then,  $R_{\text{BERT}}$  can be formulated in a ‘‘quasi’’ WMD form:

$$R_{\text{BERT}}(\mathbf{x}^1, \mathbf{y}^1) := \sum_{i,j} C_{ij} \cdot F_{ij}$$

$$F_{ij} = \begin{cases} \frac{1}{M} & \text{if } x_j = \arg \max_{\hat{x}_j^1 \in \mathbf{x}^1} E(y_i^1)^\top E(\hat{x}_j^1) \\ 0 & \text{otherwise} \end{cases}$$

$$C_{ij} = \begin{cases} \frac{M}{Z} \text{idf}(y_i^1) E(x_j^1)^\top E(y_i^1) & \text{if } x_j = \arg \max_{\hat{x}_j^1 \in \mathbf{x}^1} E(y_i^1)^\top E(\hat{x}_j^1) \\ 0 & \text{otherwise} \end{cases}$$

where  $Z = \sum_{y_i^1 \in \mathbf{y}^1} \text{idf}(y_i^1)$  and  $M$  is the size of  $n$ -grams in  $\mathbf{x}^1$ . Similarly, we can have  $P_{\text{BERT}}$  in a quasi WMD form (ignored). Then,  $F_{\text{BERT}}$  can be formulated as harmonic-mean of two WMD forms of  $P_{\text{BERT}}$  and  $R_{\text{BERT}}$ .

### A.2 Routing

In this section, we study the aggregation function  $\phi$  with routing scheme, which achieved promising results on NLP tasks. Specifically, we introduce a nonparametric clustering with Kernel Density Estimation (KDE) for routing since KDE bridges a family of kernel functions with underlying empirical distributions, which often leads to computational efficiency (Zhang et al., 2018), defined as:

$$\min_{\mathbf{v}, \gamma} f(\mathbf{z}) = \sum_{i=1}^L \sum_{j=1}^T \gamma_{ij} k(d(\mathbf{v}_j - \mathbf{z}_{i,j}))$$

$$\text{s.t. } \forall i, j : \gamma_{ij} > 0, \sum_{j=1}^L \gamma_{ij} = 1.$$

where  $d(\cdot)$  is a distance function,  $\gamma_{ij}$  denotes underlying closeness between the aggregated vector  $\mathbf{v}_j$  and vector  $\mathbf{z}_i$  in the  $i$ -th layer, some instances of  $k(\cdot)$  (Wand and Jones, 1994) can be illustrated as:

$$\text{Gaussian} : k(x) \triangleq \exp\left(-\frac{x}{2}\right), \quad \text{Epanechnikov} : k(x) \triangleq \begin{cases} 1-x & x \in [0, 1) \\ 0 & x \geq 1. \end{cases}$$

One typical solution for KDE clustering to minimize  $f(\mathbf{z})$  is taking Mean Shift (Comaniciu and Meer, 2002), defined as:

$$\nabla f(\mathbf{z}) = \sum_{i,j} \gamma_{ij} k'(d(\mathbf{v}_j, \mathbf{z}_{i,j})) \frac{\partial d(\mathbf{v}_j, \mathbf{z}_{i,j})}{\partial \mathbf{v}}$$

Firstly,  $\mathbf{v}_j^{\tau+1}$  can be updated while  $c_{ij}^{\tau+1}$  is fixed:

$$\mathbf{v}_j^{\tau+1} = \frac{\sum_i \gamma_{ij}^{\tau} k'(d(\mathbf{v}_j^{\tau}, \mathbf{z}_{i,j})) \mathbf{z}_{i,j}}{\sum_{i,j} k'(d(\mathbf{v}_j^{\tau}, \mathbf{z}_{i,j}))}$$

Intuitively,  $\mathbf{v}_j$  can be explained as a final aggregated vector from  $L$  contextualized layers. Then, we adopt SGD to update  $\gamma_{ij}^{\tau+1}$ :

$$\gamma_{ij}^{\tau+1} = \gamma_{ij}^{\tau} + \alpha \cdot k(d(\mathbf{v}_j^{\tau}, \mathbf{z}_{i,j}))$$

where  $\alpha$  is the hyper parameter to control step size. The routing process is summarized in Algorithm 1.

---

**Algorithm 1** Aggregation by Routing

---

```

1: procedure ROUTING( $\mathbf{z}_{ij}, \ell$ )
2: Initialize  $\forall i, j : \gamma_{ij} = 0$ 
3: while true do
4:   foreach representation  $i$  and  $j$  in layer  $\ell$  and  $\ell + 1$  do  $\gamma_{ij} \leftarrow \text{softmax}(\gamma_{ij})$ 
5:   foreach representation  $j$  in layer  $\ell + 1$  do
6:      $\mathbf{v}_j \leftarrow \sum_i \gamma_{ij} k'(d(\mathbf{v}_j, \mathbf{z}_i)) \mathbf{z}_i / \sum_i k'(d(\mathbf{v}_j, \mathbf{z}_i))$ 
7:   foreach representation  $i$  and  $j$  in layer  $\ell$  and  $\ell + 1$  do  $\gamma_{ij} \leftarrow \gamma_{ij} + \alpha \cdot k(d(\mathbf{v}_j, \mathbf{z}_i))$ 
8:   loss  $\leftarrow \log(\sum_{i,j} \gamma_{ij} k(d(\mathbf{v}_j, \mathbf{z}_i)))$ 
9:   if  $|\text{loss} - \text{preloss}| < \epsilon$  then
10:    break
11:  else
12:    preloss  $\leftarrow$  loss
13: return  $\mathbf{v}_j$ 

```

---

**Best Layer and Layer-wise Consolidation** Table 6 compares our Word Mover based metric combining BERT representations on different layers with stronger BERT representations consolidated from these layers ( $p$ -means and routing). We often see that which layer has best performance is task-dependent, and our WMD-based metrics with  $p$ -means or routing schema come close to the oracle performance obtained from best layers.

Metrics	Direct Assessment						
	cs-en	de-en	fi-en	lv-en	ru-en	tr-en	zh-en
WMD-1 + BERT + LAYER 8	.6361	.6755	.8134	.7033	.7273	.7233	.7175
WMD-1 + BERT + LAYER 9	.6510	.6865	.8240	.7107	.7291	<b>.7357</b>	<b>.7195</b>
WMD-1 + BERT + LAYER 10	.6605	<b>.6948</b>	<b>.8231</b>	.7158	<b>.7363</b>	.7317	.7168
WMD-1 + BERT + LAYER 11	<b>.6695</b>	.6845	.8192	.7048	.7315	.7276	.7058
WMD-1 + BERT + LAYER 12	.6677	.6825	.8194	<b>.7188</b>	.7326	.7291	.7064
WMD-1 + BERT + ROUTING	.6618	.6897	.8225	.7122	.7334	.7301	.7182
WMD-1 + BERT + PMEANS	.6623	.6873	.8234	.7139	.7350	.7339	.7192

Table 6: Absolute Pearson correlations with segment-level human judgments on WMT17 to-English translations.

**Experiments** Table 7, 8 and 9 show correlations between metrics (all baseline metrics and Word Mover based metrics) and human judgments on machine translation, text summarization and dialogue response generation tasks.

Setting	Metrics	Direct Assessment							
		cs-en	de-en	fi-en	lv-en	ru-en	tr-en	zh-en	Average
Existing Metrics	BLEND	0.594	0.571	0.733	0.594	0.622	0.671	0.661	0.635
	RUSE	0.624	0.644	0.750	0.697	0.673	0.716	0.691	0.685
	SENTBLEU	0.435	0.432	0.571	0.393	0.484	0.538	0.512	0.481
	CHRF++	0.523	0.534	0.678	0.520	0.588	0.614	0.593	0.579
	METEOR++	0.552	0.538	0.720	0.563	0.627	0.626	0.646	0.610
	BERTSCORE-F1	0.670	0.686	0.820	0.710	0.729	0.714	0.704	0.719
Word-Mover	WMD-1 + W2V	0.392	0.463	0.558	0.463	0.456	0.485	0.481	0.471
	WMD-1 + BERT + ROUTING	0.658	0.689	0.823	0.712	0.733	0.730	0.718	0.723
	WMD-1 + BERT + MNLI + ROUTING	0.665	0.705	0.834	0.744	0.735	0.752	0.736	0.739
	WMD-2 + BERT + MNLI + ROUTING	0.676	0.706	0.831	0.743	0.734	0.755	0.732	0.740
	WMD-1 + BERT + PMEANS	0.662	0.687	0.823	0.714	0.735	0.734	0.719	0.725
	WMD-1 + BERT + MNLI + PMEANS	0.670	0.708	0.835	0.746	0.738	0.762	0.744	0.743
	WMD-2 + BERT + MNLI + PMEANS	0.679	0.710	0.832	0.745	0.736	0.763	0.740	0.743

Table 7: Absolute Pearson correlations with segment-level human judgments on WMT17 to-English translations.

Setting	Metrics	TAC-2008				TAC-2009			
		Responsiveness		Pyramid		Responsiveness		Pyramid	
		$r$	$\rho$	$r$	$\rho$	$r$	$\rho$	$r$	$\rho$
Existing Metrics	$S_{full}^3$	0.696	0.558	0.753	0.652	0.731	0.552	0.838	0.724
	$S_{best}^3$	0.715	0.595	0.754	0.652	0.738	0.595	0.842	0.731
	TF*IDF-1	0.176	0.224	0.183	0.237	0.187	0.222	0.242	0.284
	TF*IDF-2	0.047	0.154	0.049	0.182	0.047	0.167	0.097	0.233
	ROUGE-1	0.703	0.578	0.747	0.632	0.704	0.565	0.808	0.692
	ROUGE-2	0.695	0.572	0.718	0.635	0.727	0.583	0.803	0.694
	ROUGE-1-WE	0.571	0.450	0.579	0.458	0.586	0.437	0.653	0.516
	ROUGE-2-WE	0.566	0.397	0.556	0.388	0.607	0.413	0.671	0.481
	ROUGE-L	0.681	0.520	0.702	0.568	0.730	0.563	0.779	0.652
	FRAME-1	0.658	0.508	0.686	0.529	0.678	0.527	0.762	0.628
	FRAME-2	0.676	0.519	0.691	0.556	0.715	0.555	0.781	0.648
	BERTSCORE-F1	0.724	0.594	0.750	0.649	0.739	0.580	0.823	0.703
Word-Mover	WMD-1 + W2V	0.669	0.559	0.665	0.611	0.698	0.520	0.740	0.647
	WMD-1 + BERT + ROUTING	0.729	0.601	0.763	0.675	0.740	0.580	0.831	0.700
	WMD-1 + BERT + MNLI + ROUTING	0.734	0.609	0.768	0.686	0.747	0.589	0.837	0.711
	WMD-2 + BERT + MNLI + ROUTING	0.731	0.593	0.755	0.666	0.753	0.583	0.827	0.698
	WMD-1 + BERT + PMEANS	0.729	0.595	0.755	0.660	0.742	0.581	0.825	0.690
	WMD-1 + BERT + MNLI + PMEANS	0.736	0.604	0.760	0.672	0.754	0.594	0.831	0.701
	WMD-2 + BERT + MNLI + PMEANS	0.734	0.601	0.752	0.663	0.753	0.586	0.825	0.694

Table 8: Correlation of automatic metrics with summary-level human judgments for TAC-2008 and TAC-2009.

Setting	Metrics	BAGEL			SFHOTEL		
		Inf	Nat	Qual	Inf	Nat	Qual
Existing Metrics	BLEU-1	0.225	0.141	0.113	0.107	0.175	0.069
	BLEU-2	0.211	0.152	0.115	0.097	0.174	0.071
	BLEU-3	0.191	0.150	0.109	0.089	0.161	0.070
	BLEU-4	0.175	0.141	0.101	0.084	0.104	0.056
	ROUGE-L	0.202	0.134	0.111	0.092	0.147	0.062
	NIST	0.207	0.089	0.056	0.072	0.125	0.061
	CIDER	0.205	0.162	0.119	0.095	0.155	0.052
	METEOR	0.251	0.127	0.116	0.111	0.148	0.082
	BERTSCORE-F1	0.267	0.210	0.178	0.163	0.193	0.118
Word-WMD	WMD-1 + W2V	0.222	0.079	0.123	0.074	0.095	0.021
	WMD-1 + BERT + ROUTING	0.294	0.209	0.156	0.208	0.256	0.178
	WMD-1 + BERT + MNLI + ROUTING	0.278	0.180	0.144	0.211	0.252	0.175
	WMD-2 + BERT + MNLI + ROUTING	0.279	0.182	0.147	0.204	0.252	0.172
	WMD-1 + BERT + PMEANS	0.298	0.212	0.163	0.203	0.261	0.182
	WMD-1 + BERT + MNLI + PMEANS	0.285	0.195	0.158	0.207	0.270	0.183
	WMD-2 + BERT + MNLI + PMEANS	0.284	0.194	0.156	0.204	0.270	0.182

Table 9: Spearman correlation with utterance-level human judgments for BAGEL and SFHOTEL datasets.