

# EffEval: A Comprehensive Evaluation of Efficiency for MT Evaluation Metrics

Daniil Larionov<sup>1</sup>, Jens Grünwald<sup>2</sup>, Christoph Leiter<sup>1</sup> and Steffen Eger<sup>1</sup>

<sup>1</sup> Natural Language Learning Group, Bielefeld University

<sup>2</sup> Department of Computer Science, Technical University of Darmstadt  
daniil.larionov@uni-bielefeld.de

## Main Contributions

We study 3 different aspects of computational efficiency in application to MT evaluation metrics.

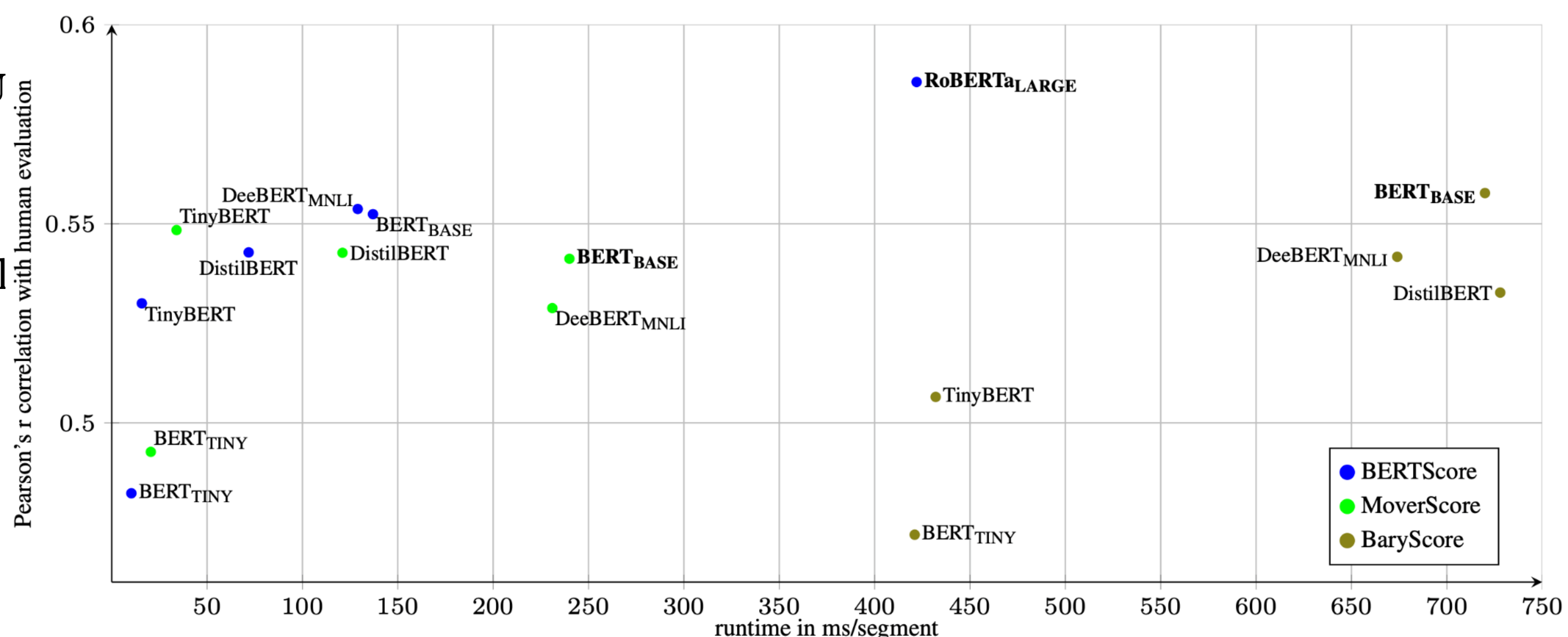
1. We replace computationally-heavy transformer models with their light-weight alternatives for metrics like **BERTScore**, **MoverScore**, **BaryScore**.
2. We switch from costly alignment techniques (Word Mover Distance - WMD) to their approximations in **MoverScore**
3. We train **COMET**-like models efficiently with adapters

## Motivation

- Inefficient metrics (like BERTScore with RoBERTa-Large) require expensive hardware to run in reasonable time, which prevents **under-resourced practitioners** from using it.
- Even with good hardware it takes too much time and energy: **71** hours for BERTScore to evaluate (30k segments × 5 language pairs × 50 MT systems)
- Metrics have loads of other applications which would benefit from an efficient option: **RL reward functions**, **mined data filtering**, **online re-ranking**

## 1. Efficient Transformers

- Evaluation datasets: WMT 15/16/21
- Measured runtime in ms/segment, on GPU and CPU, averaged across 3 runs.
- Evaluated models:
  - RoBERTa-Large - Baseline
  - BERT-Base - Smaller model
  - DistilBERT - Distilled
  - TinyBERT - Adv. Distilled
  - BERT-Tiny - BERT
  - BERT-Tiny miniatures
  - DeeBERT - Dyn. early exiting
- Best tradeoff between quality/efficiency - **TinyBERT**



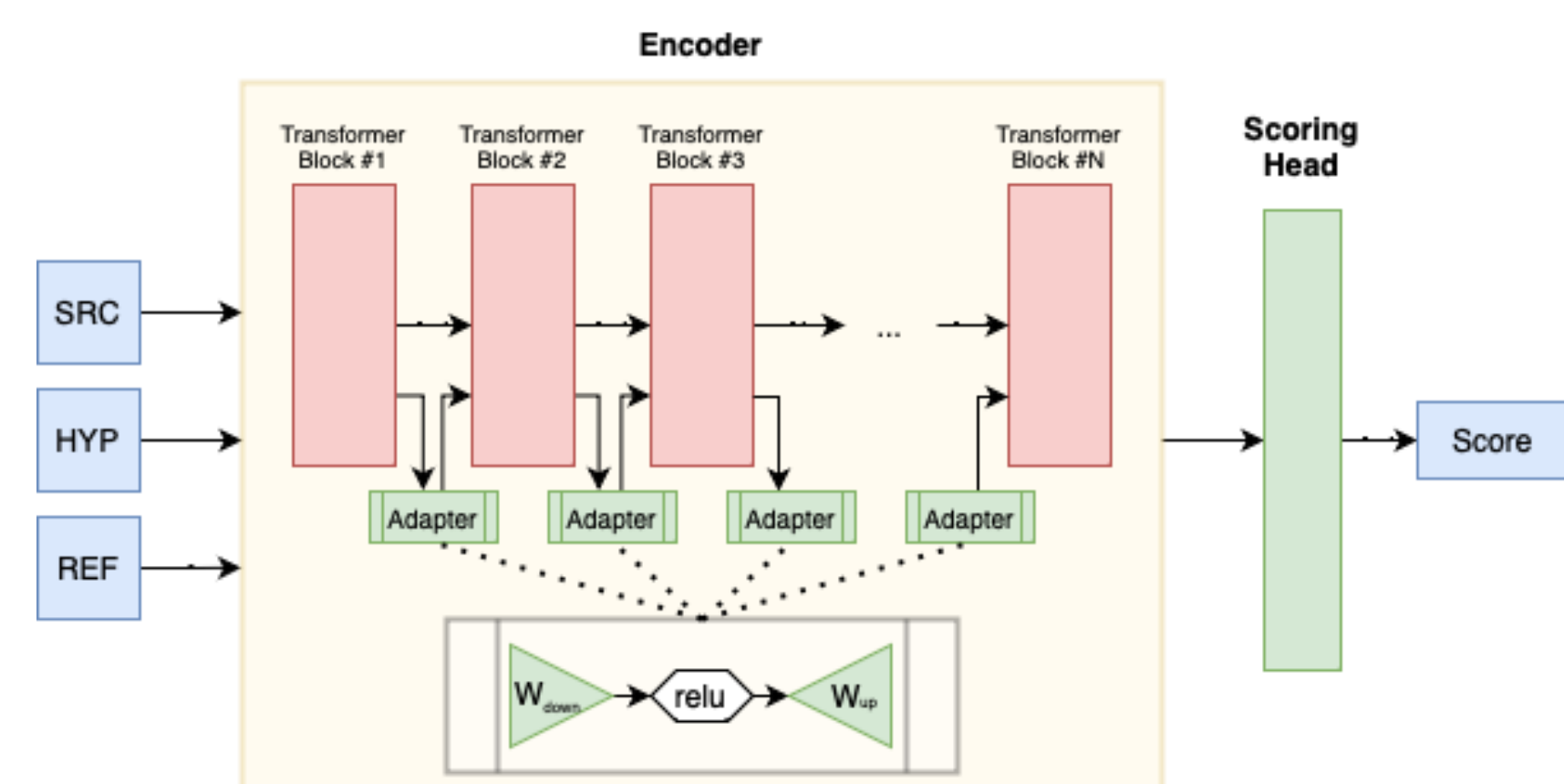
## 2. Approximations of Word Mover Distance

- MoverScore uses **WMD** as measure of distance between two texts.
- WMD calculates minimum cumulative distance that words from one text needs to travel to match the other text. Complexity is  **$O(p^3 \log p)$** .
- Word Centroid Distance (**WCD**) and Relaxed Word Mover Distance (**RWMD**) are two fast approximations of WMD.

Step	WMD	WCD	RWMD
get BERT embeddings	285.499	287.915	291.122
calculate distance matrix	0.829	0.005	0.782
calculate distance	5.602	0.616	0.449

## 3. COMET + Adapters = ❤️?

- COMET is a trained MT evaluation metric, based on XLM-RoBERTa-Large. It is quite large and takes days to train fully. **An ideal case for adapters!**
- We tested various adapter configurations: Pfeiffer et al, Houlsby et al, Parallel adapter, Compacter, (IA)<sup>3</sup>. Measured **Mem.** in **MB/Token**, **Fwd.** and **Bwd.** speed in **Tokens/Second** and **Kendall-tau** correlation on WMT21
- Along with large COMET we also tested smaller COMETINHO.
- Results:
  - Adapters **decrease** GPU RAM usage for training and **improve** backward pass speed
  - Models with adapter can **outperform** the ones without them!
  - Simpler adapters works better



Config	Mem.↓	Fwd.↑	Bwd.↑	$\tau$ ↑
pfeiffer	4.88	5123	<b>4808</b>	0.273
parallel	4.97	5128	4525	<b>0.289</b>
houlsby	4.87	4607	4036	0.273
compacter	<b>4.80</b>	3649	3049	0.269
(IA) <sup>3</sup>	5.76	<b>5195</b>	4712	0.268
no adapters	7.32	<u>6247</u>	2238	0.275
reference	-	-	-	<b>0.290</b>

Config	Mem.↓	Fwd.↑	Bwd.↑	$\tau$ ↑
pfeiffer	0.770	25499	25774	0.252
parallel	<b>0.741</b>	26109	<b>26113</b>	<b>0.252</b>
houlsby	0.769	23746	21678	0.252
compacter	0.776	18382	15671	0.243
(IA) <sup>3</sup>	0.997	<b>27075</b>	24804	0.248
no adapters	1.012	<u>31836</u>	18941	0.243
reference	-	-	-	<b>0.241</b>



Github