

# **IST-Unbabel 2021 Submission ÉCNICO** for the Quality Estimation Shared Task

Chrysoula Zerva, Daan Van Stigt, Ricardo Rei, Ana Farinha, Pedro Ramos, Jose de Souza, Taisiya Glushkova, Miguel Vera, Fabio Kepler, André Martins

## Highlights

- We incorporate **adapter layers** in the Openkiwi architecture [1].
- 2) We explore different types of **uncertainty**.
  - **Glass-box** features [2] extracted from the NMT models a)
  - **Aleatoric uncertainty** derived from the human annotations b)
  - **Epistemic uncertainty** of our QE model(s). **C**)
- We enhance our model with **out-of-domain** data from the Metrics shared tasks [3]

### **Tasks & Data**

- **T1** Predict sentence level quality based on direct assessment (DA) scores Score range: [-7.542, 3.178]
- Predict sentence level quality **T2.**1 based on post-edit HTER scores Score range: [0,1] Predict word level binary tags **T2.2** Labels: OK | BAD

#### DATA

- > Multi-lingual training and testing sets
- > Mix of high, medium, low resource language pairs  $\rightarrow$  Different score **distributions**
- > **Zero-shot** (blind) language pairs

We train our models using multi-lingual encoders with **adapter layers (M1)** to obtain models that generalise better. Along the same lines, we experiment with pre-training the models on the he data provided for the past **Metrics shared tasks** (M1, M2). This out-of-domain data encompasses 30 language pairs from the news domain (versus 7 in the QE dataset), including the zero-shot QE pairs.

				(	Offic	ial I	resu	lts					
	Model	Multi	En-De	En-Zh	Ro-En	Et-En	Ne-En	Si-En	Ru-En	En-Cs	En-Ja	Ps-En	Km-En
Task 1	QEMind	0.68	0.57	0.60	0.91	0.81	0.87	0.60	0.81	0.58	0.36	0.65	0.68
	HW-TSC	0.67	0.58	0.58	0.90	0.81	0.86	0.58	0.88	0.57	0.36	0.62	0.66
	IST-Unbabel	0.665	0.58	0.59	0.90	0.80	0.86	0.61	0.792	0.58	0.36	0.63	0.65
	papago (IKT)	0.66	0.57	0.567	0.901	0.76	0.85	0.60	0.79	0.57	0.33	0.64	0.66
	TUDa	0.63	0.47	0.56	0.89	0.79	0.83	0.57	0.76	0.55	0.33	0.61	0.64
	Inmon‡	0.623	3 <del>44</del>	8 <del></del> -	<u> 1922</u>		<del>11</del> 3	83 <del>44</del>		0.55	0.30	0.59	0.63
	papago (KD)	0.61	0.55	0.55	0.88	0.79	0.82	0.58	0.74	0.497	0.28	0.58	0.63
	BASELINE	0.54	0.41	0.53	0.82	0.66	0.74	0.51	0.68	0.35	0.23	0.48	0.56
Task 2.1	Model	Multi	En-De	En-Zh	Ro-En	Et-En	Ne-En	Si-En	Ru-En	En-Cs	En-Ja	Ps-En	Km-En
	HW-TSC	0.63	0.65	0.37	0.86	0.81	0.80	0.87	0.56	0.48	0.26	0.53	0.75
	IST-Unbabel	0.60	0.62	0.29	0.88	0.81	0.72	0.71	0.54	0.53	0.28	0.56	0.66
	BASELINE	0.50	0.53	0.28	0.83	0.71	0.63	0.61	0.45	0.31	0.10	0.50	0.58
Task 2.2 MT words	<u>n</u> Model	Multi	En-De	En-Zh	Ro-En	Et-En	Ne-En	Si-En	Ru-En	En-Cs	En-Ja	Ps-En	Km-En
	HW-TSC	0.53	0.51	0.35	0.66	0.60	0.67	0.84	0.45	0.38	0.25	0.45	0.63
	ST-Unbabel	0.43	0.46	0.31	0.65	0.57	0.51	0.52	0.332	0.37	0.16	0.37	0.45
	BASELINE	0.35	0.37	0.25	0.54	0.46	0.44	0.43	0.26	0.27	0.13	0.31	0.35







This work was supported by the P2020 programs MAIA (contract 045909) and Unbabel4EU (contract 042671), by the European Research Council (ERC StG DeepSPIN 758969), and by the Fundação para a Ciência e Tecnologia through contract UIDB/50008/2020.

[3] Mathur, Nitika, Johnny Wei, Markus Freitag, Qingsong Ma, and Ondřej Bojar. 2020. **Results of the wmt20 metrics shared task.** 

> **Conference on Empirical Methods in** Natural Language Processing