

A Linguistically Motivated Test Suite to Semi-Automatically Evaluate German—English Machine Translation Output

Vivien Macketanz, Eleftherios Avramidis, Aljoscha Burchardt, He Wang, Renlong Ai, Shushen Manakhimova, Ursula Strohriegel, Sebastian Möller, Hans Uszkoreit

¹German Research Center for Artificial Intelligence (DFKI)
Speech and Language Technology Lab, Berlin

MT-TestSuite

- around 10,000 test items to evaluate German <> English MT outputs
- 13/14 categories per language direction, divided into more than 100 linguistic phenomena each, with min. 20 test items per phenomenon
- set of rules (fixed strings and regex) for the semi-automatic evaluation
- half of the test items and rules available at <https://github.com/DFKI-NLP/mt-testsuite>

Test Suite Creation Process

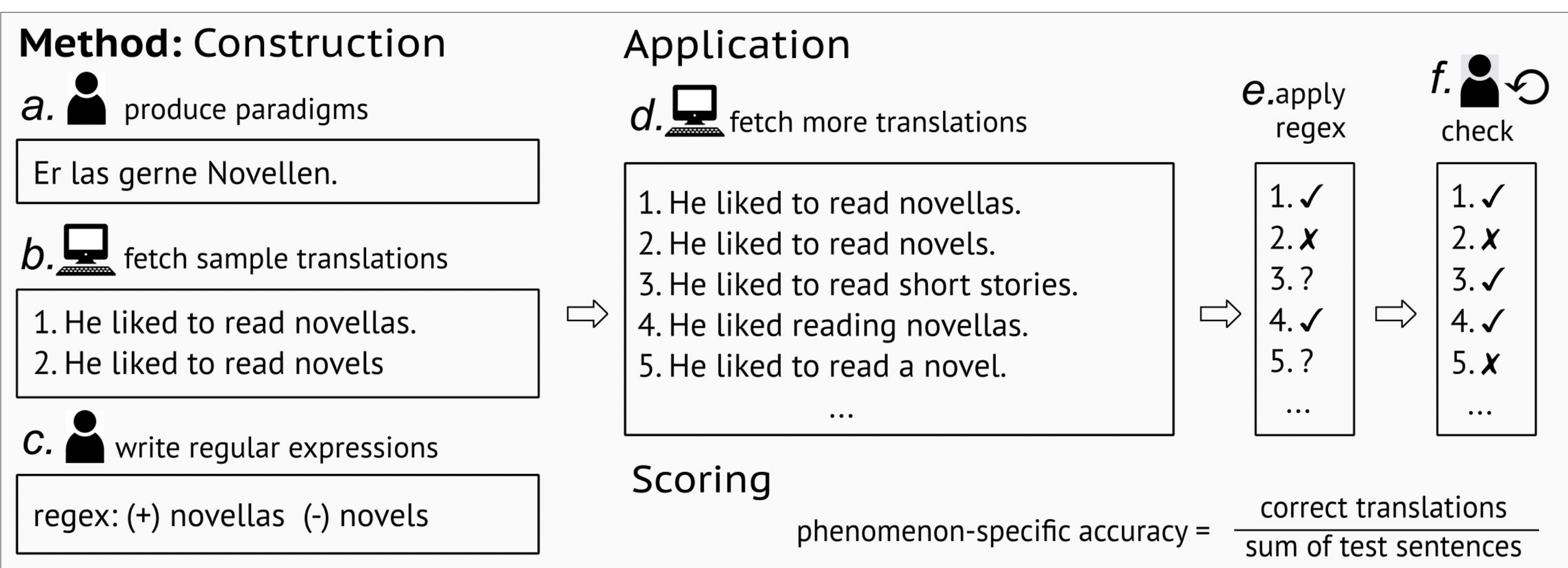


Figure 1: Example of the preparation of the test suite for one test item.

Test Suite Examples

	Collocation	False Friends	Modal negated preterite
Test item	Simon trinkt am liebsten <i>lieblichen</i> Wein.	Dieser Autor schreibt hauptsächlich <i>Novellen</i> .	Ihr <i>durftet nicht</i> lesen.
WMT 18	Simon prefers to drink <i>adorable</i> wine.	This author writes mainly <i>Novellen</i> .	You <i>are not allowed to read</i> .
WMT 19	Simon loves to drink <i>lovely</i> wine.	This author writes mainly <i>novellas</i>.	You <i>are not allowed to read</i> .
WMT 20	Simon prefers to drink <i>lovely</i> wine.	This author writes mainly <i>novellas</i>.	You <i>don't read</i> .
WMT 21	Simon prefers to drink <i>sweet</i> wine.	This author mainly writes <i>novels</i> .	You <i>were not allowed to read</i>.

Table 1: Output examples from 2018 to 2021. *phenomenon*, **correct**

Test Suite Application Process

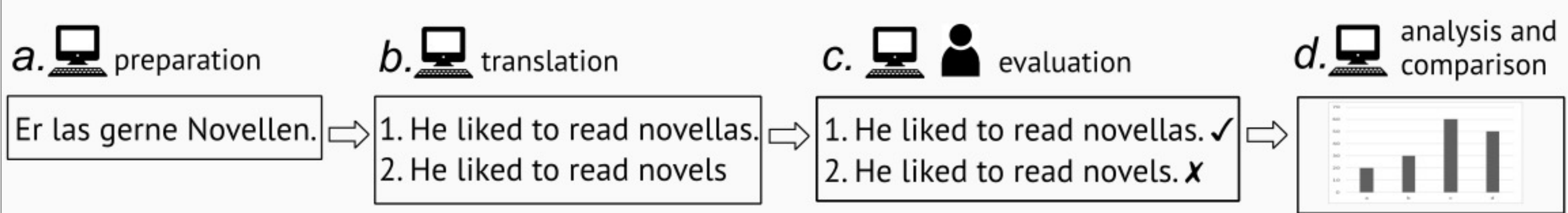


Figure 2: Example of the application of the test suite for one test item.

Application Example: WMT Test Suite Track

- the MT-TestSuite was submitted to the WMT 2018-2021 (De-En) and 2021 (En-De)
- translations from systems of the news translation task were evaluated with the test suite

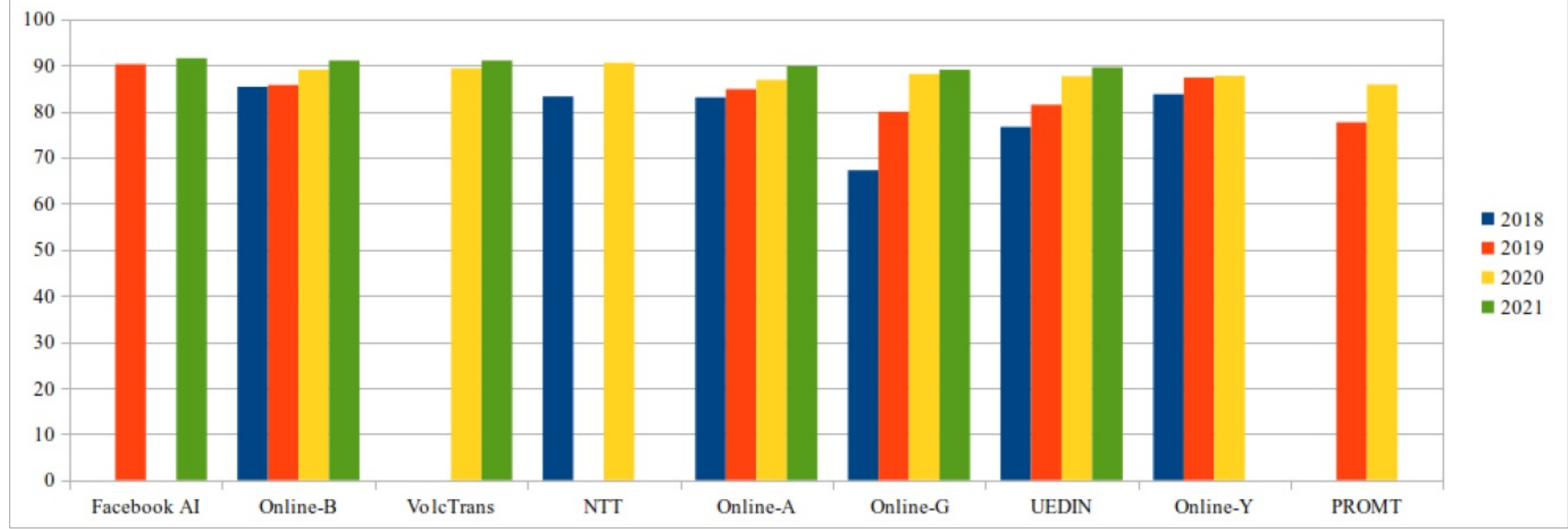


Figure 3: System improvements (accuracy macro-average on the test suite from 2018 to 2021 for German to English, including the systems that appear at least once in the past two years).

Further Application Examples

- WMT metrics track
- Domain-specific test suite
- Quality estimation
- Portuguese—English test suite

Future Work

- automatic test item creator for test set expansion
- new publicly available evaluation tool
- include more language pairs, currently we are working on extending to English—Russian