



# MT-ComparEval: Graphical evaluation interface for Machine Translation development

Ondřej Klejch, Eleftherios Avramidis, Aljoscha Burchardt, Martin Popel

klejch@ufal.mff.cuni.cz, {eleftherios.avramidis,aljoscha.burchardt}@dfki.de, popel@ufal.mff.cuni.cz

Charles University in Prague, Faculty of Mathematics and Physics, Institute of Formal and Applied Linguistics  
German Research Center for Artificial Intelligence (DFKI), Language Technology Lab



● Try it now (all WMT 2014-2015 results): <http://wmt.ufal.cz>

○ Install it (and report issues or contribute): <https://github.com/choko/MT-ComparEval>

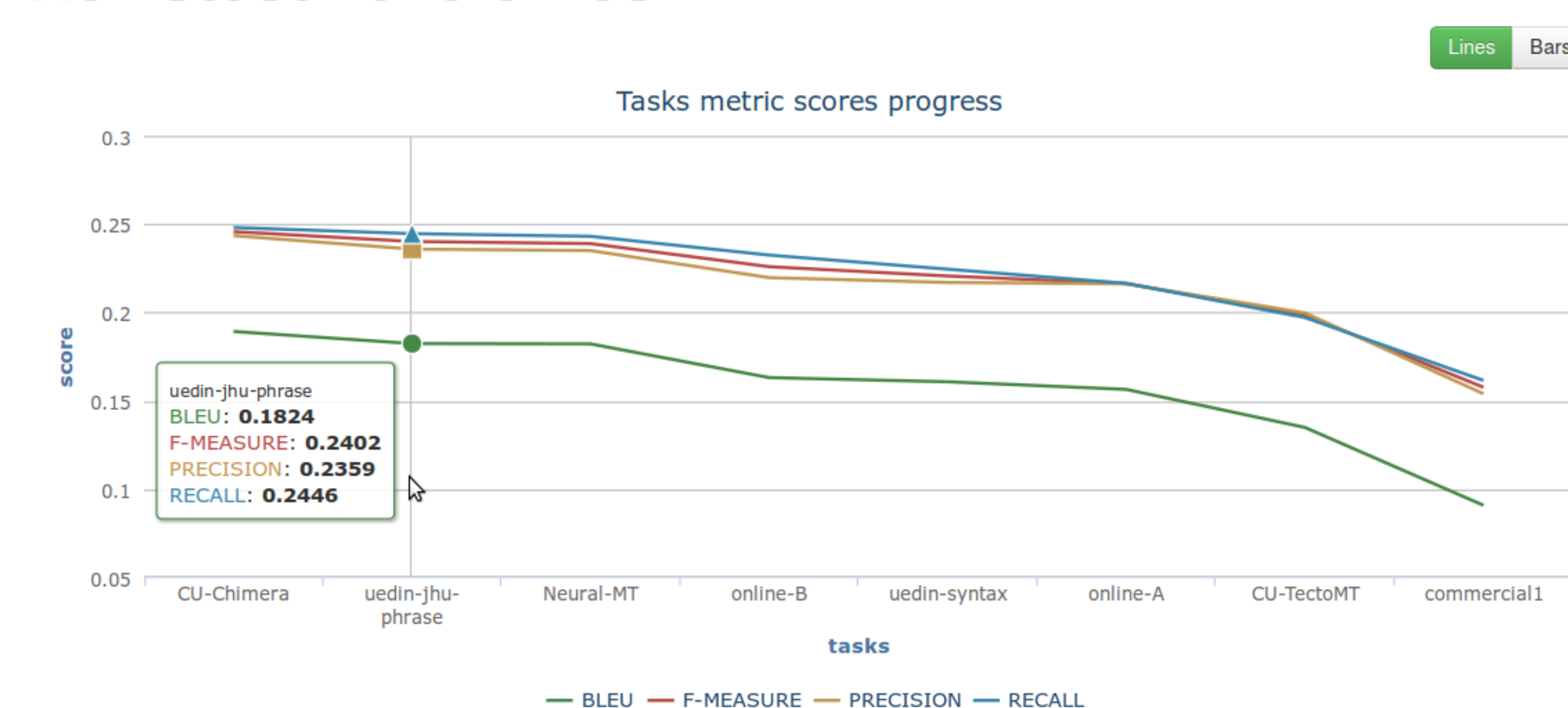
- web-based tool for MT developers
- check progress of a system over time or compare several MT systems
- focus on analyzing system differences
- can be integrated into your workflow (import translations from disk/git/...)

① Select an "Experiment", e.g. English-Czech WMT15

② See BLEU, F-measure, ... Select two "Tasks", i.e. system translations

Experiments	
Newstest 2015 en-cs	
Newstest 2015 en-cs tuning task	
Newstest 2015 en-de	
Newstest 2015 en-fi	
Newstest 2015 en-ru	
Newstest 2015 cs-en	

Newstest 2015 en-cs



name	description	BLEU ↓	F-MEASURE	PRECISION	RECALL	
<input type="checkbox"/> CU-Chimera		0.1893	0.2458	0.2435	0.2481	hide
<input type="checkbox"/> uedin-jhu-phrase		0.1824	0.2402	0.2359	0.2446	hide
<input type="checkbox"/> Neural-MT		0.1822	0.239	0.2351	0.2431	hide

③ Sentences pane with various diffs highlighted

Sentences sorted by diffBLEU

No more excuse for missing examples!

	BLEU	BLEU-cis	F-MEASURE	F-MEASURE-cis	PRECISION	PRECISION-cis	RECALL	RECALL-cis
Neural-MT	0.1173	0.1173	0.1864	0.1864	0.1864	0.1864	0.1864	0.1864
CU-Chimera	0.6989	0.6989	0.7025	0.7025	0.7025	0.7025	0.7025	0.7025
Diff	-0.5816	-0.5816	-0.5161	-0.5161	-0.5161	-0.5161	-0.5161	-0.5161

- confirmed n-gram = occurs in the reference (light yellow and blue highlight)
- improving n-gram = confirmed n-gram occurring in only one of the systems (dark yellow and blue highlight)
- worsening n-gram = unconfirmed, occurring in only one of the systems (red highlighting)
- diff (LCS underlined in green)

- See the sentences with biggest improvement/worsening
- See the 1-grams...4-grams with biggest improvement/worsening
- Hints for improving the systems

④ Statistics pane (paired bootstrap resampling, ...)



⑤ Confirmed n-grams pane

⑥ Unconfirmed n-grams

Click on any n-gram to see all its occurrences

Source	Reference	Pilot 0.00	Pilot 1.00
Wie öffne ich ein Dokument in Libreoffice ?	How do I open a document in LibreOffice ?	As I open a document in LibreOffice ?	How do I open a document in Libreoffice ?
Wie verschicke ich eine Datei über Skype ?	How do I send a file using Skype ?	As I always send a file on Skype ?	How do I send a file above Skype ?