

# Using MT-ComparEval

---

Roman Sudarikov, Martin Popel, Ondřej Bojar,  
Aljoscha Burchardt, Ondřej Klejch

Charles University in Prague, ÚFAL  
DFKI Berlin

Centre for Speech Technology Research, University of Edinburgh

LREC 2016 MT Eval Workshop  
Portorož, May 24

# What is MT-ComparEval?

---

- web-based tool for MT developers, who can
- check progress of a system over time or compare several MT systems
- focus on analyzing system differences
- integrate MT-ComparEval into their workflow (import translations: disk/git/REST-API)

# Try it now!

---

- <http://wmt.ufal.cz>  
all WMT 2014–2016 systems
- <http://mt-compareval.ufal.cz>  
upload and inspect your translations
- <https://github.com/choko/MT-ComparEval>  
install it (and report issues or contribute)

# wmt.ufal.cz

MT-ComparEval

Wmt.Ufal.Cz

MT-ComparEval

Fork me on GitHub

## Newstest 2016

<a href="#">Newstest 2016 en-cs</a>	
<a href="#">Newstest 2016 en-de</a>	
<a href="#">Newstest 2016 en-fi</a>	
<a href="#">Newstest 2016 en-ro</a>	
<a href="#">Newstest 2016 en-ru</a>	
<a href="#">Newstest 2016 en-tr</a>	
<a href="#">Newstest 2016 cs-en</a>	
<a href="#">Newstest 2016 de-en</a>	
<a href="#">Newstest 2016 fi-en</a>	
<a href="#">Newstest 2016 ro-en</a>	
<a href="#">Newstest 2016 ru-en</a>	
<a href="#">Newstest 2016 tr-en</a>	

Select an „experiment“

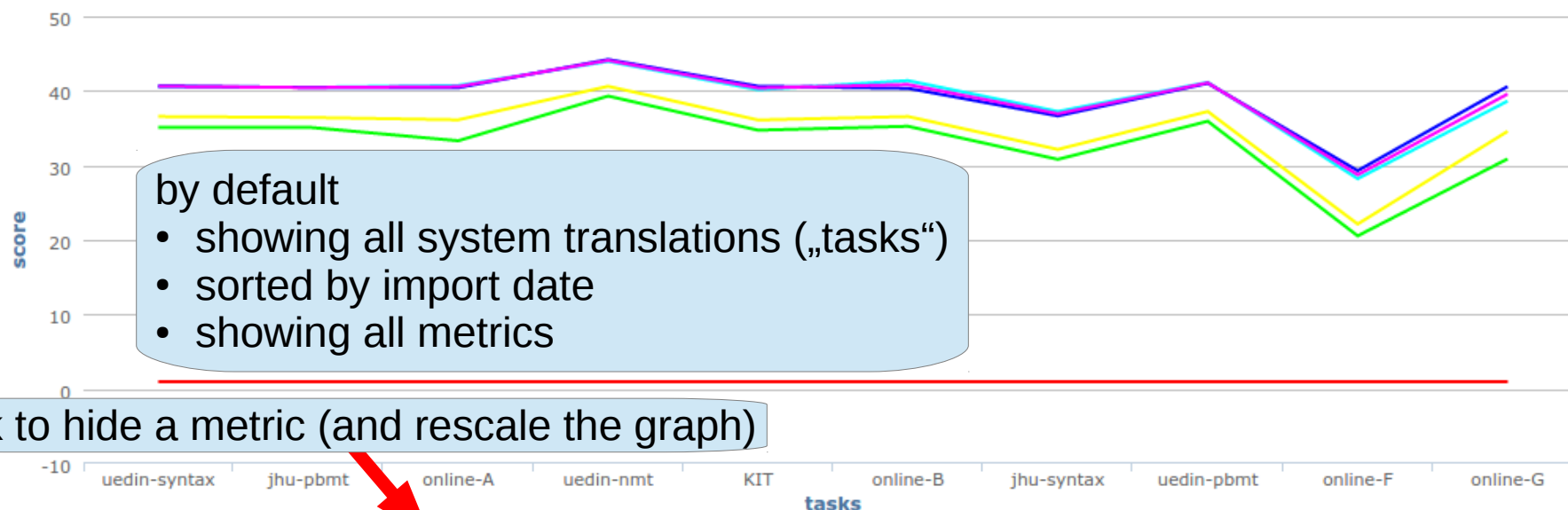
## IT Test 2016

<a href="#">IT Test 2016 en-bg</a>	
<a href="#">IT Test 2016 en-cs</a>	
<a href="#">IT Test 2016 en-de</a>	
<a href="#">IT Test 2016 en-es</a>	
<a href="#">IT Test 2016 en-eu</a>	
<a href="#">IT Test 2016 en-nl</a>	
<a href="#">IT Test 2016 en-pt</a>	

# Newstest 2016 de-en

Lines Bars

Tasks metric scores progress



— BREVITY-PENALTY — BLEU — BLEU-cased — PRECISION — RECALL — F-MEASURE

compare

Click to sort by BLEU

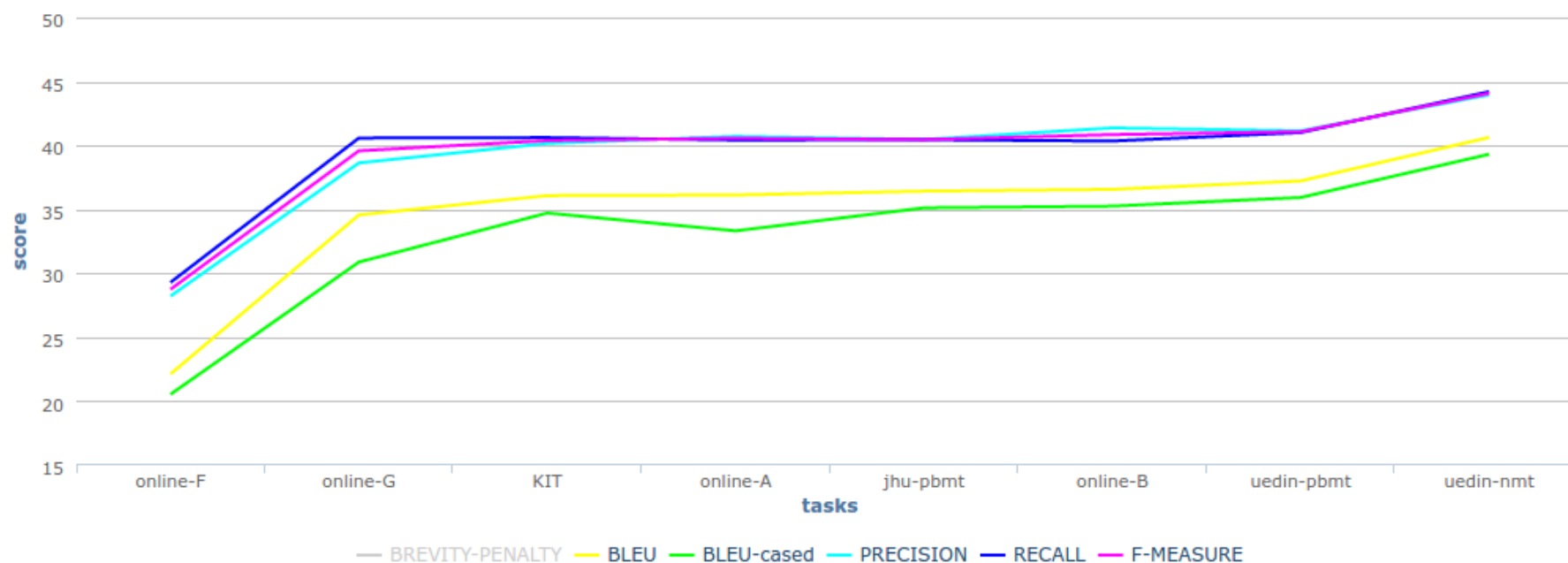
Click to hide a system

name	description	BREVITY-PENALTY	BLEU	BLEU-cased	PRECISION	RECALL	F-MEASURE	
<input type="checkbox"/> uedin-syntax		1	36.57	35.09	40.53	40.67	40.6	<a href="#">hide</a>
<input type="checkbox"/> jhu-pbmt		1	36.43	35.13	40.47	40.45	40.46	<a href="#">hide</a>
<input type="checkbox"/> online-A		0.99	36.13	33.32	40.72	40.42	40.57	<a href="#">hide</a>
<input type="checkbox"/> uedin-nmt		1	40.63	39.32	43.96	44.21	44.08	<a href="#">hide</a>
<input type="checkbox"/> KIT		1	36.09	34.72	40.16	40.62	40.39	<a href="#">hide</a>
<input type="checkbox"/> online-B		0.98	36.57	35.26	41.37	40.35	40.85	<a href="#">hide</a>

# Newstest 2016 de-en

Lines Bars

Tasks metric scores progress



compare

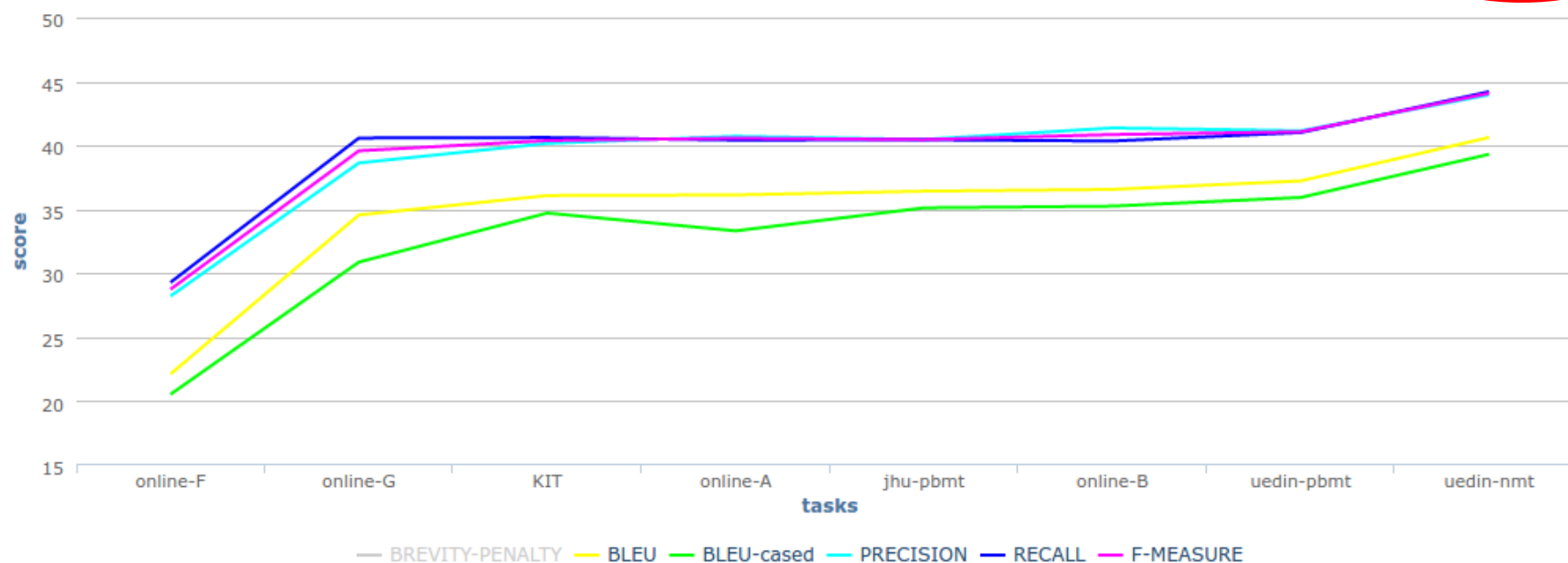
name	description	BLEU ↑	BLEU-cased	PRECISION	RECALL	F-MEASURE	
<input type="checkbox"/> online-F	1	22.13	20.55	28.22	29.29	28.75	<a href="#">hide</a>
<input type="checkbox"/> online-G	1	34.57	30.89	38.64	40.59	39.59	<a href="#">hide</a>
<input type="checkbox"/> KIT	1	36.09	34.72	40.16	40.62	40.39	<a href="#">hide</a>
<input type="checkbox"/> online-A	0.99	36.13	33.32	40.72	40.42	40.57	<a href="#">hide</a>
<input type="checkbox"/> jhu-pbmt	1	36.43	35.13	40.47	40.45	40.46	<a href="#">hide</a>
<input type="checkbox"/> online-B	0.98	36.57	35.26	41.37	40.35	40.85	<a href="#">hide</a>
<input type="checkbox"/> uedin-pbmt	1	37.23	35.94	41.15	41.02	41.08	<a href="#">hide</a>



# Newstest 2016 de-en

Lines Bars

Tasks metric scores progress



compare

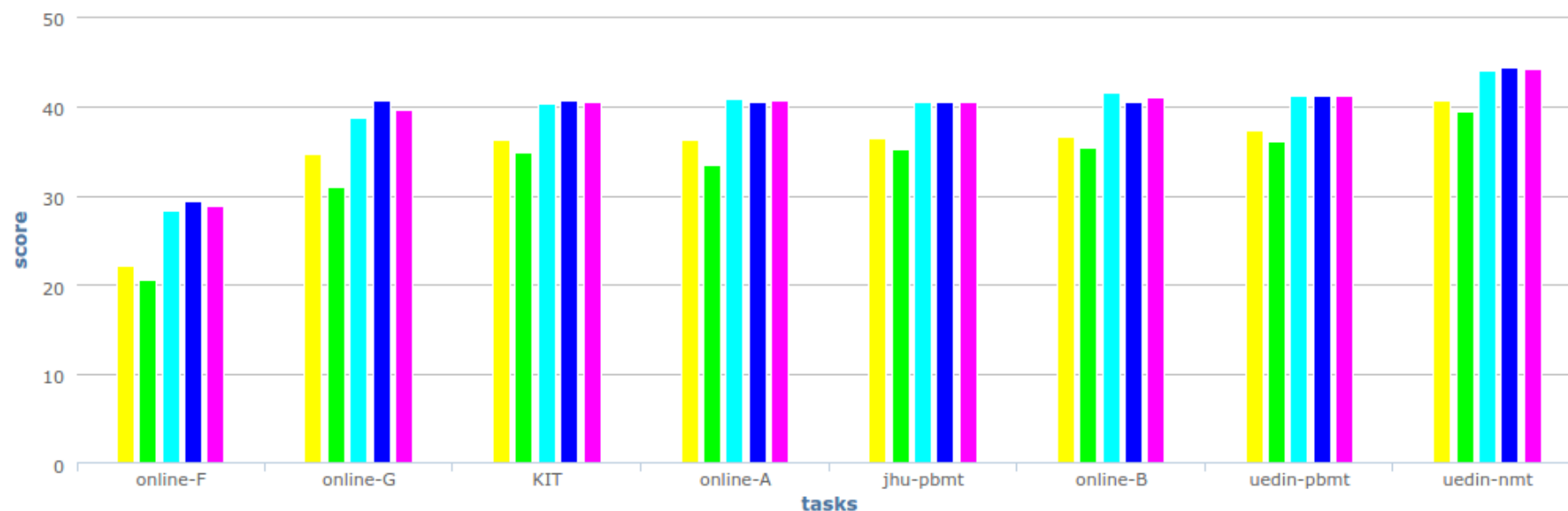
name	description	BLEU ↑	BLEU-cased	PRECISION	RECALL	F-MEASURE	
<input type="checkbox"/> online-F	1	22.13	20.55	28.22	29.29	28.75	<a href="#">hide</a>
<input type="checkbox"/> online-G	1	34.57	30.89	38.64	40.59	39.59	<a href="#">hide</a>
<input type="checkbox"/> KIT	1	36.09	34.72	40.16	40.62	40.39	<a href="#">hide</a>
<input type="checkbox"/> online-A	0.99	36.13	33.32	40.72	40.42	40.57	<a href="#">hide</a>
<input type="checkbox"/> jhu-pbmt	1	36.43	35.13	40.47	40.45	40.46	<a href="#">hide</a>
<input type="checkbox"/> online-B	0.98	36.57	35.26	41.37	40.35	40.85	<a href="#">hide</a>
<input type="checkbox"/> uedin-pbmt	1	37.23	35.94	41.15	41.02	41.08	<a href="#">hide</a>

# Newstest 2016 de-en

Lines

Bars

Tasks metric scores progress

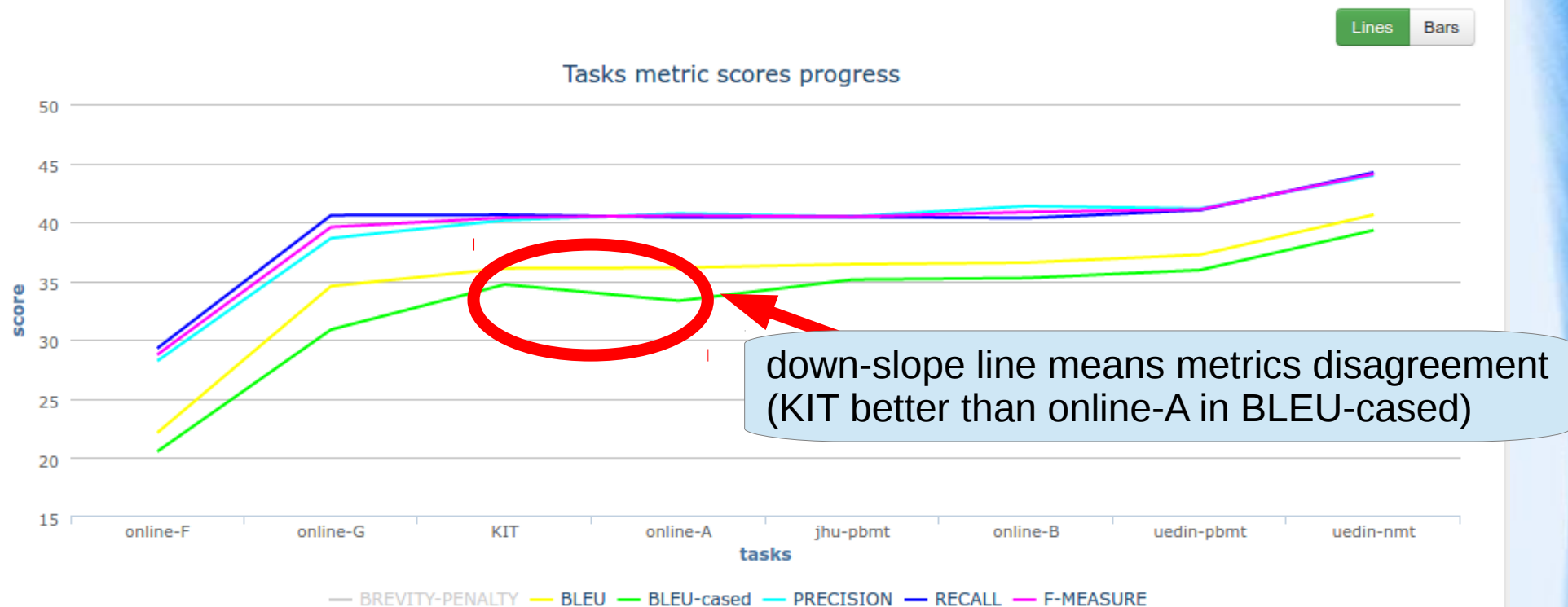


☐ BREVITY-PENALTY
 ☒ BLEU
 ☒ BLEU-cased
 ☒ PRECISION
 ☒ RECALL
 ☒ F-MEASURE

name	description	BLEU ↑	BLEU-cased	PRECISION	RECALL	F-MEASURE	
<input type="checkbox"/> online-F	1	22.13	20.55	28.22	29.29	28.75	<a href="#">hide</a>
<input type="checkbox"/> online-G	1	34.57	30.89	38.64	40.59	39.59	<a href="#">hide</a>
<input type="checkbox"/> KIT	1	36.09	34.72	40.16	40.62	40.39	<a href="#">hide</a>
<input type="checkbox"/> online-A	0.99	36.13	33.32	40.72	40.42	40.57	<a href="#">hide</a>
<input type="checkbox"/> jhu-pbmt	1	36.43	35.13	40.47	40.45	40.46	<a href="#">hide</a>
<input type="checkbox"/> online-B	0.98	36.57	35.26	41.37	40.35	40.85	<a href="#">hide</a>
<input type="checkbox"/> uedin-pbmt	1	37.23	35.94	41.15	41.02	41.08	<a href="#">hide</a>



# Newstest 2016 de-en

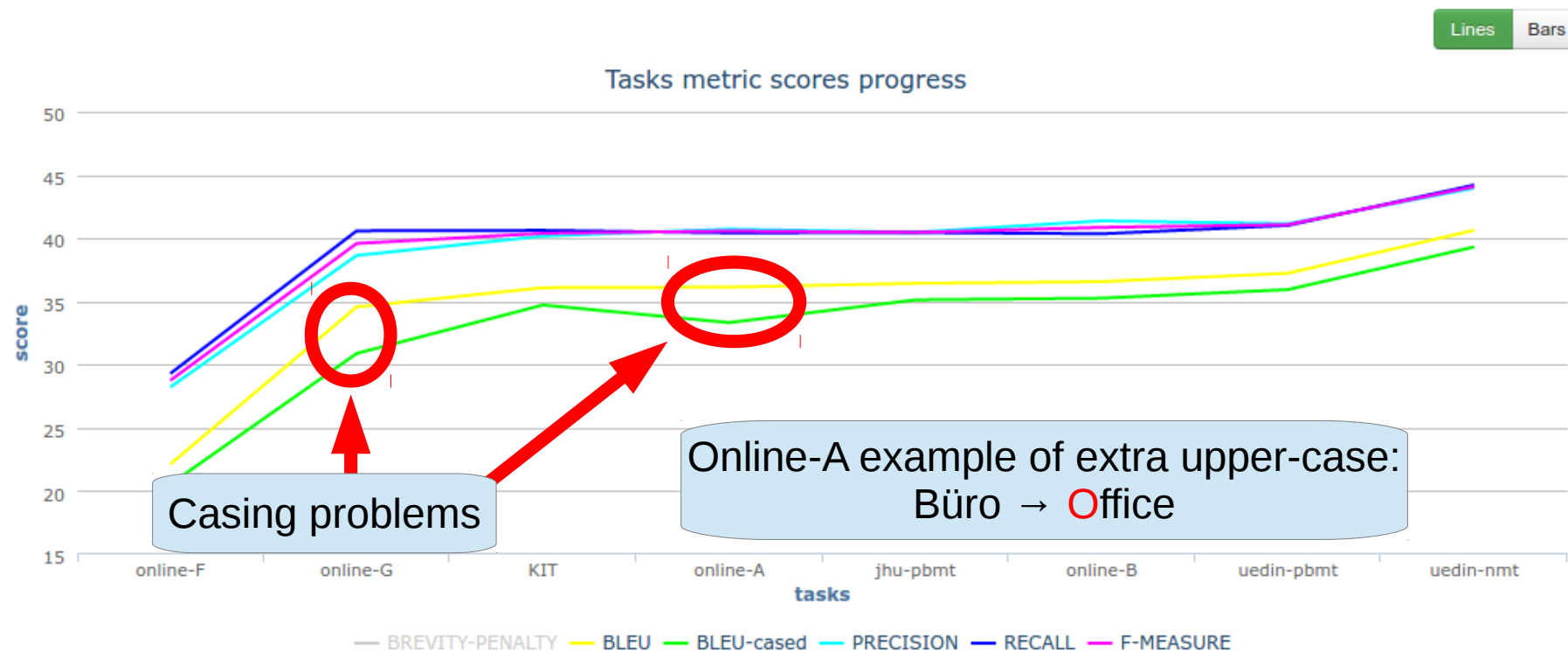


compare

BLEU (case insensitive) vs. BLEU-cased

name	description	BLEU ↑	BLEU-cased	PRECISION	RECALL	F-MEASURE	
<input type="checkbox"/> online-F	1	22.13	20.55	28.22	29.29	28.75	<a href="#">hide</a>
<input type="checkbox"/> online-G	1	34.57	30.89	38.64	40.59	39.59	<a href="#">hide</a>
<input type="checkbox"/> KIT	1	36.09	34.72	40.16	40.62	40.39	<a href="#">hide</a>
<input type="checkbox"/> online-A	0.99	36.13	33.32	40.72	40.42	40.57	<a href="#">hide</a>
<input type="checkbox"/> jhu-pbmt	1	36.43	35.13	40.47	40.45	40.46	<a href="#">hide</a>
<input type="checkbox"/> online-B	0.98	36.57	35.26	41.37	40.35	40.85	<a href="#">hide</a>
<input type="checkbox"/> uedin-pbmt	1	37.23	35.94	41.15	41.02	41.08	<a href="#">hide</a>

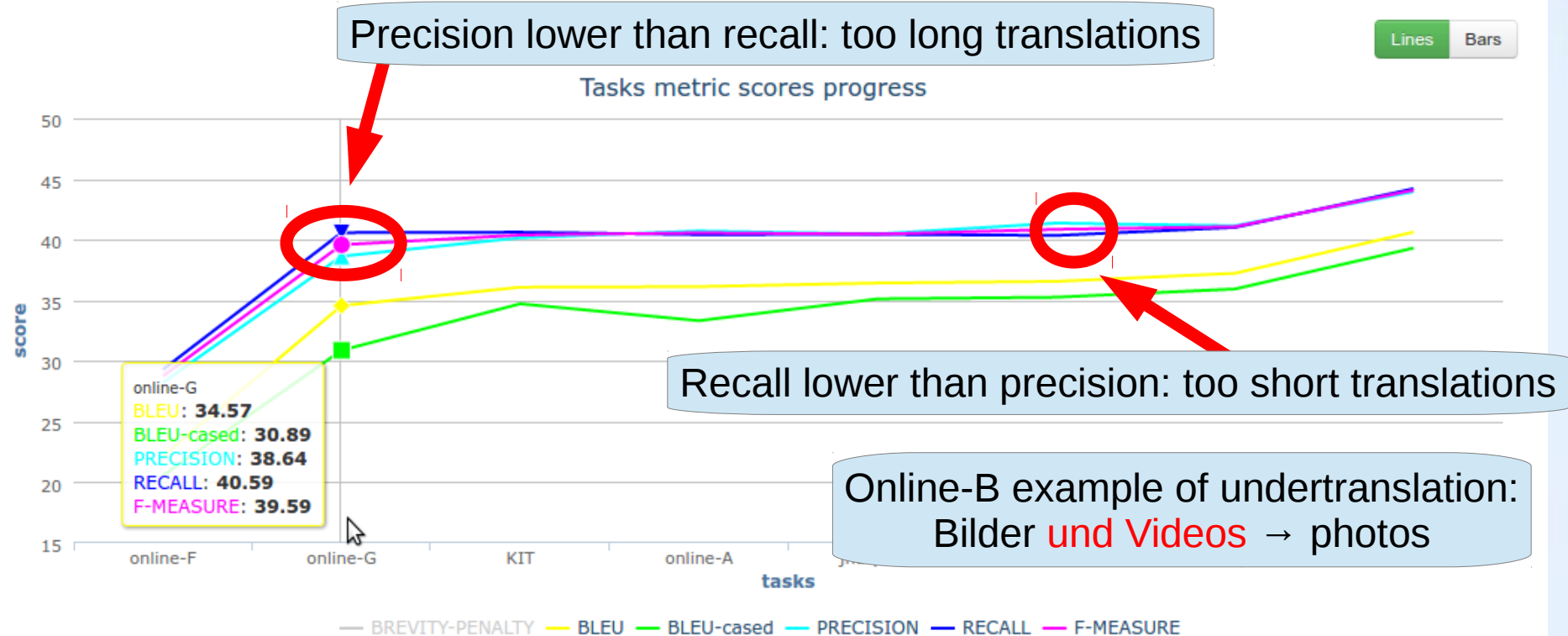
# Newstest 2016 de-en



compare

name	description	BLEU ↑	BLEU-cased	PRECISION	RECALL	F-MEASURE	
<input type="checkbox"/> online-F	1	22.13	20.55	28.22	29.29	28.75	<a href="#">hide</a>
<input type="checkbox"/> online-G	1	34.57	30.89	38.64	40.59	39.59	<a href="#">hide</a>
<input type="checkbox"/> KIT	1	36.09	34.72	40.16	40.62	40.39	<a href="#">hide</a>
<input type="checkbox"/> online-A	0.99	36.13	33.32	40.72	40.42	40.57	<a href="#">hide</a>
<input type="checkbox"/> jhu-pbmt	1	36.43	35.13	40.47	40.45	40.46	<a href="#">hide</a>
<input type="checkbox"/> online-B	0.98	36.57	35.26	41.37	40.35	40.85	<a href="#">hide</a>
<input type="checkbox"/> uedin-pbmt	1	37.23	35.94	41.15	41.02	41.08	<a href="#">hide</a>

# Newstest 2016 de-en



compare

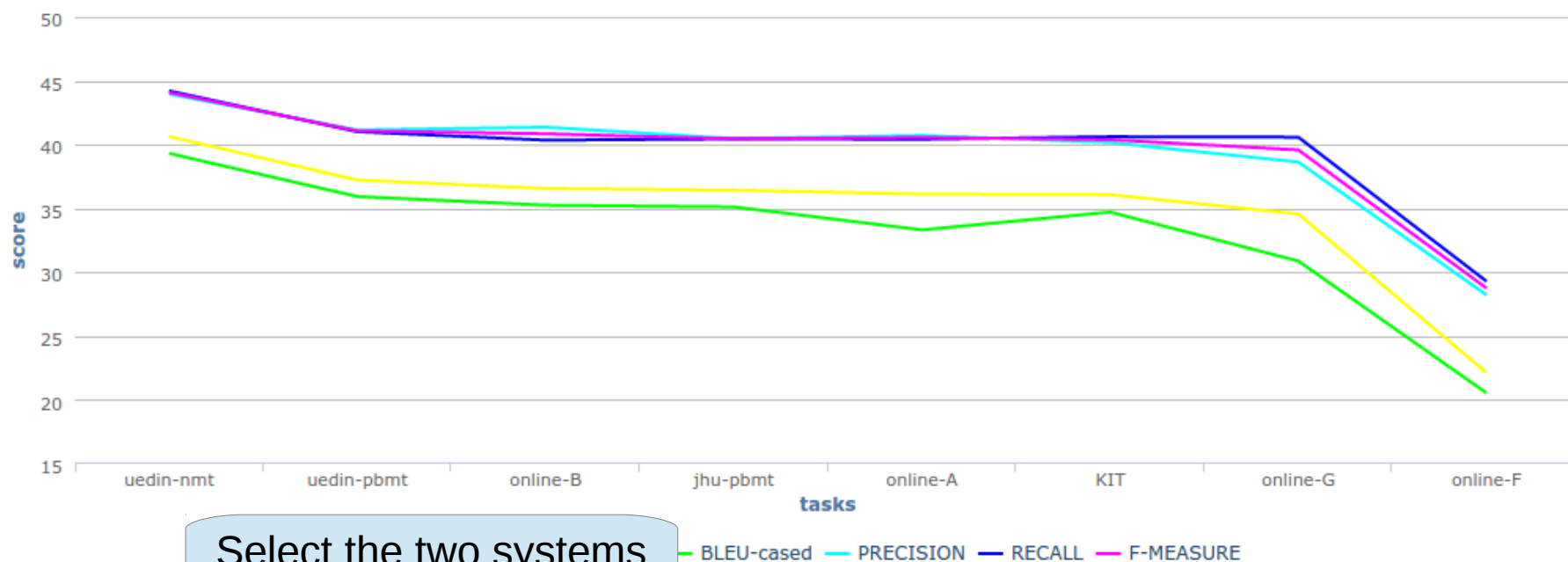
name	description	BLEU ↑	BLEU-cased	PRECISION	RECALL	F-MEASURE	
<input type="checkbox"/> online-F	1	22.13	20.55	28.22	29.29	28.75	<a href="#">hide</a>
<input type="checkbox"/> online-G	1	34.57	30.89	38.64	40.59	39.59	<a href="#">hide</a>
<input type="checkbox"/> KIT	1	36.09	34.72	40.16	40.62	40.39	<a href="#">hide</a>
<input type="checkbox"/> online-A	0.99	36.13	33.32	40.72	40.42	40.57	<a href="#">hide</a>
<input type="checkbox"/> jhu-pbmt	1	36.43	35.13	40.47	40.45	40.46	<a href="#">hide</a>
<input type="checkbox"/> online-B	0.98	36.57	35.26	41.37	40.35	40.85	<a href="#">hide</a>
<input type="checkbox"/> uedin-pbmt	1	37.23	35.94	41.15	41.02	41.08	<a href="#">hide</a>

# Newstest 2016 de-en

How to see the differences  
(and find example sentences)?

Lines Bars

Tasks metric scores progress



Select the two systems  
and click "compare".

compare

name	description	BLEU ↓	BLEU-cased	PRECISION	RECALL	F-MEASURE	
<input type="checkbox"/> uedin-nmt	1	40.63	39.32	43.96	44.21	44.08	hide
<input type="checkbox"/> uedin-pbmt	1	37.23	35.94	41.15	41.02	41.08	hide
<input checked="" type="checkbox"/> online-B	0.98	36.57	35.26	41.37	40.35	40.85	hide
<input checked="" type="checkbox"/> jhu-pbmt	1	36.43	35.13	40.47	40.45	40.46	hide
<input type="checkbox"/> online-A	0.99	36.13	33.32	40.72	40.42	40.57	hide
<input type="checkbox"/> KIT	1	36.09	34.72	40.16	40.62	40.39	hide
<input type="checkbox"/> online-G	1	34.57	30.89	38.64	40.59	39.59	hide

Sentences

Statistics

Confirmed n-grams

Unconfirmed n-grams

You can quickly compare with another system

## Sentences

### Options

#### N-grams highlighting options

- ☒ Highlight confirmed n-grams
- ☒ Highlight improving n-grams
- ☒ Highlight worsening n-grams

#### Diff highlighting options

- ☐ Show diff with reference
- ☒ Show diff for online-B
- ☐ Show diff for jhu-pbmt
- ☐ Show diff with each other

#### Sentences visibility options

- ☒ Show source
- ☒ Show reference
- ☒ Show online-B
- ☒ Show jhu-pbmt
- ☒ Show sentence level metrics

Source	Konkrete Zahlen nannte sie nicht und verwies auf Gespräche am Mittwoch in Berlin .					
Reference	She gave no specific figures and referred to talks in Berlin on Wednesday .					
online-B	Concrete numbers they did not name , citing talks on Wednesday in Berlin .					
jhu-pbmt	She gave no specific figures and referred to talks in Berlin on Wednesday .					
	BREVITY-PENALTY	BLEU	BLEU-cased	PRECISION	RECALL	F-MEASURE
online-B	1	8.89	8.89	14.56	14.56	14.56
jhu-pbmt	1	100	100	100	100	100
Diff	0.0000	-91.1100	-91.1100	-85.4400	-85.4400	-85.4400

Source	Das Quartiersbüro ist geöffnet Montag , Mittwoch und Freitag von 9 bis 14 Uhr , Dienstag und Donnerstag von 16 bis 18 Uhr und samstags von 10 Uhr bis 12 Uhr .					
--------	--	--	--	--	--	--



online-B



jhu-pbmt

BLEU



Sentences

Statistics

Confirmed n-grams

Unconfirmed n-grams

## Sentences

the selected metric

4 panes

## Options

## N-grams highlighting options

- ☒ Highlight confirmed n-grams
- ☒ Highlight improving n-grams
- ☒ Highlight worsening n-grams

## Diff highlighting options

- ☐ Show diff with reference
- ☒ Show diff for online-B
- ☐ Show diff for jhu-pbmt

## Sentences visibility options

- ☒ Show source
- ☒ Show reference
- ☒ Show online-B

Sentences sorted according to the difference in the selected metric. Top: sentences where JHU-PBMT outperforms Online-B the most. Very useful also for regression testing (new version vs. baseline).

Reverse ordering  
Top: OnlineB wins

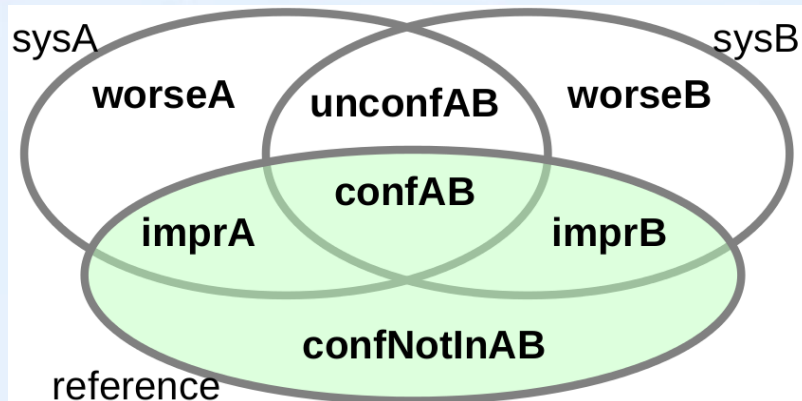
Source	Konkrete Zahlen nannte sie nicht und verwies auf Gespräche am Mittwoch in Berlin .					
Reference	She gave no specific figures and referred to talks in Berlin on Wednesday .					
online-B	Concrete numbers they did not name , citing talks on Wednesday in Berlin .					
jhu-pbmt	She gave no specific figures and referred to talks in Berlin on Wednesday .					
	BREVITY-PENALTY	BLEU	BLEU-cased	PRECISION	RECALL	F-MEASURE
online-B	1	8.89	8.89	14.56	14.56	14.56
jhu-pbmt	1	100	100	100	100	100
Diff	0.0000	-91.1100	-91.1100	-85.4400	-85.4400	-85.4400

Source	Das Quartiersbüro ist geöffnet Montag , Mittwoch und Freitag von 9 bis 14 Uhr , Dienstag und Donnerstag von 16 bis 18 Uhr und samstags von 10 Uhr bis 12 Uhr .
--------	--



# Color highlighting

Source	Wie hoch sind die Kosten ?
Reference	What are the costs ?
online-B	How high are the costs ?
jhu-pbmt	What are the costs ?



- **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)

# Sentence pane tricks

Source	Sie schicken uns Ihre Bilder und Videos oder eine SMS an 61124.
online-B	Send us your photos or SMS to 61124th
jhu-pbmt	Send us your pictures and videos , or an SMS to 61124.
Source	Fischer ist Spezialist für Herzschrittmacher und Defibrillatoren .
online-B	Fischer is a specialist for pacemakers and defibrillators .
jhu-pbmt	Fischer is a specialist in cardiac pacemakers and defibrillators .
Source	Konrad wollte einen Container am Stephansplatz als Büro aufstellen .
online-B	Konrad wanted to set up a container at Stephansplatz offices .
jhu-pbmt	Konrad wanted to set up a container on Stephansplatz as an office .
Source	Wie hoch sind die Kosten ?
online-B	How high are the costs ?
jhu-pbmt	What are the costs ?

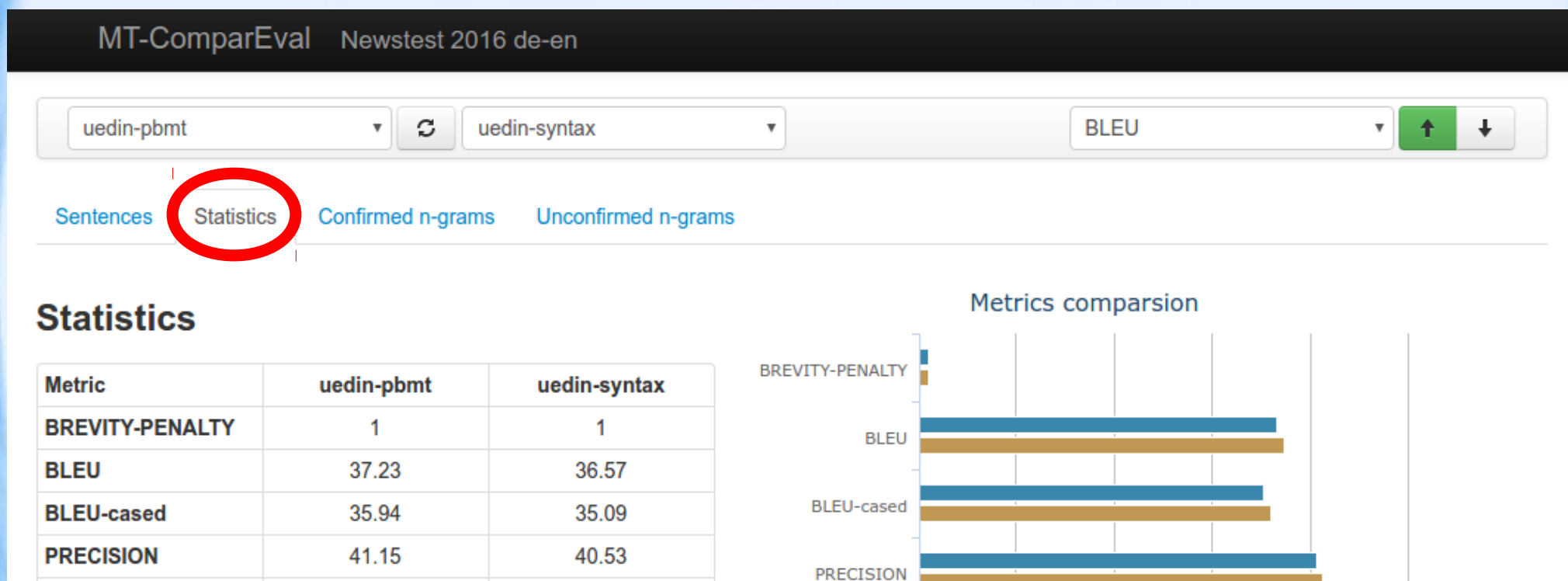
## Sentences visibility options

- ☒ Show source
- ☒ Show reference
- ☒ Show online-B
- ☒ Show jhu-pbmt
- ☒ Show sentence level metrics

Looking for nice example sentences?

- Look for short sentences (for slides).
- Search for blue-highlighted content words with no equivalent in the other system's translation.
- Or use Hjerson-omissions metric (off by default).
- Show just the reference and search for blue.
- More sentences are loaded as you scroll down.

# Statistics pane



uedin-pbmt

uedin-syntax

BLEU

Sentences

Statistics

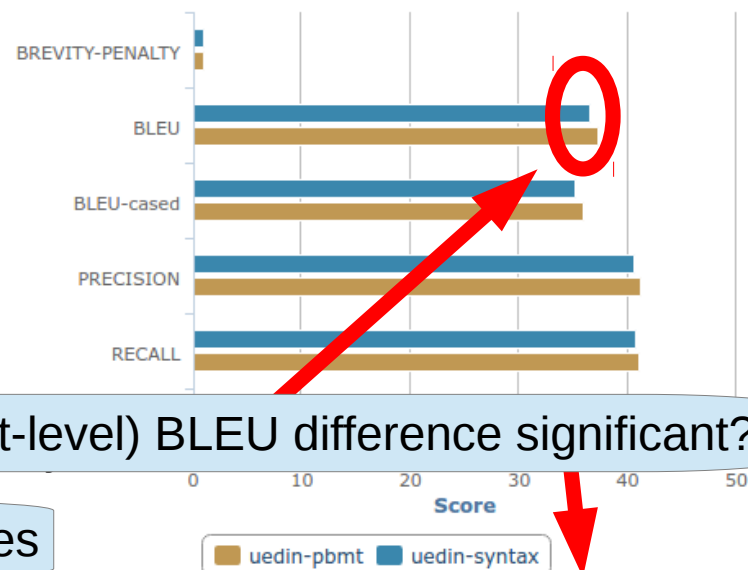
Confirmed n-grams

Unconfirmed n-grams

## Statistics

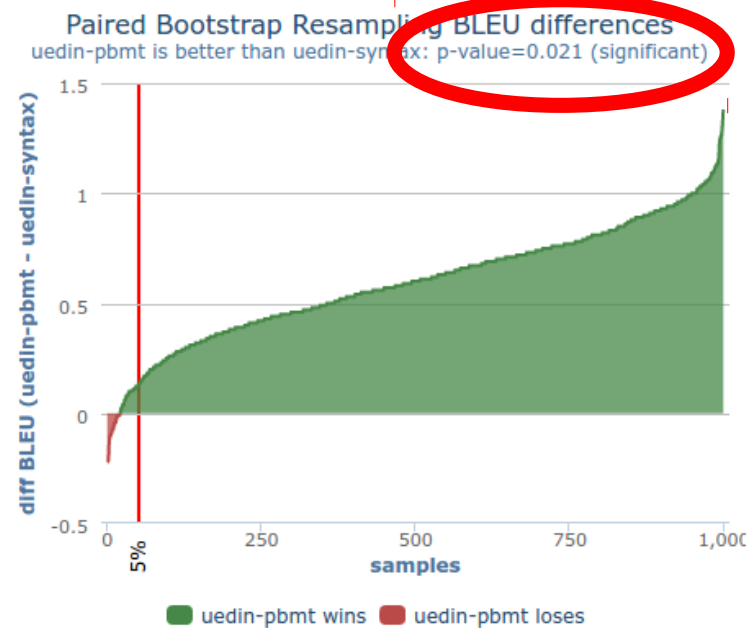
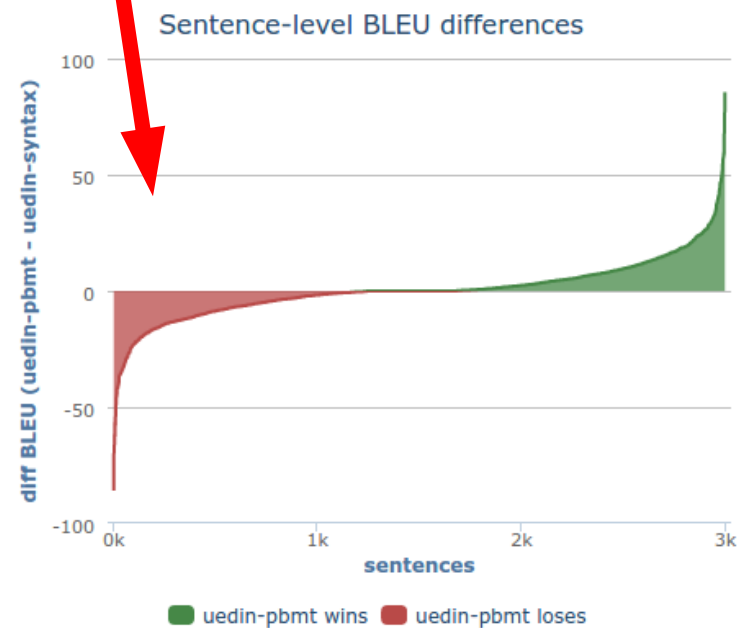
Metric	uedin-pbmt	uedin-syntax
BREVITY-PENALTY	1	1
BLEU	37.23	36.57
BLEU-cased	35.94	35.09
PRECISION	41.15	40.53
RECALL	41.02	40.67
F-MEASURE	41.08	40.6

## Metrics comparison



Is the (document-level) BLEU difference significant?

Distribution of sentence-level BLEU differences



# Confirmed n-grams pane

MT-ComparEval Newstest 2016 en-cs

uedin-nmt ↕ NYU-UMontreal ↕

Sentences Statistics **Confirmed n-grams** Unconfirmed n-grams

## n-grams confirmed by the reference

1-gram

“se” confirmed in Uedin-nmt 588 times  
“se” confirmed in NYU-UMontreal 504 times

uedin-nmt wins		NYU-UMontreal wins	uedin-n
, 3536 - 3425 = 111		000 30 - 9 = 21	, "
se 588 - 504 = 84		které 101 - 87 = 14	" řekl
v 660 - 590 = 70		Ale 40 - 26 = 14	."
o 214 - 164 = 50		řekl 100 - 87 = 13	, že
na 561 - 513 = 48		nebo 74 - 61 = 13	a



# Unconfirmed n-grams pane

MT-ComparEval Newstest 2016 en-cs

uedin-nmt NYU-UMontreal

Sentences Statistics Confirmed n-grams **Unconfirmed n-grams**

## n-grams unconfirmed by the reference

1-gram

uedin-nmt loses	NYU-UMontreal loses	uedin-n
, 1156 - 894 = 262	pro	
<b>se 419 - 288 = 131</b>	ve	
ne 109 - 7 = 102	rovnez	
; 80 - 8 = 72	bylo	& #

Thus, Uedin-nmt uses "se" more often than NYU-UMontreal: 84 times more confirmed, 131 times more unconfirmed.



uedin-nmt



NYU-UMontreal

BLEU



Sentences

Statistics

Confirmed n-grams

Unconfirmed n-grams

## n-grams unconfirmed by the reference

Encoding problems

## 1-gram

uedin-nmt loses

,	1156 - 894 = 262
se	419 - 288 = 131
ne	109 - 7 = 102
;	80 - 8 = 72
&	71 - 0 = 71
#	69 - 0 = 69
160	69 - 0 = 69
o	165 - 98 = 67
v	336 - 284 = 52
že	241 - 189 = 52

NYU-UMontreal loses

pro	147 - 95 = 52
ve	128 - 91 = 37
rovněž	39 - 5 = 34
bylo	78 - 45 = 33
jejich	63 - 34 = 29
.	96 - 71 = 25
byla	74 - 51 = 23
být	66 - 43 = 23
řekl	115 - 95 = 20
rok	24 - 4 = 20

## 2-gram

uedin-nmt loses

, ne	101 - 2 = 99
ne,	101 - 3 = 98
160 ;	69 - 0 = 69
& #	69 - 0 = 69
# 160	69 - 0 = 69
, ze	242 - 189 = 53
, a	76 - 29 = 47
" řekl	42 - 0 = 42
, který	117 - 95 = 22
000 &	22 - 0 = 22

NYU-UMontreal loses

",	210 - 31 = 179
, řekl	131 - 25 = 106
".	82 - 48 = 34
, řekla	20 - 1 = 19
, které	91 - 76 = 15
, protože	22 - 10 = 12
, říká	23 - 11 = 12
řikal ,	16 - 4 = 12
řikají ,	15 - 4 = 11
předtím ,	12 - 1 = 11

uedin-nmt



NYU-UMontreal

BLEU

[Sentences](#)[Statistics](#)[Confirmed n-grams](#)[Unconfirmed n-grams](#)

## n-grams unconfirmed by the reference

## 1-gram

## uedin-nmt loses

,	1156 - 894 = 262
se	419 - 288 = 131
ne	109 - 7 = 102
;	80 - 8 = 72
&	71 - 0 = 71
#	69 - 0 = 69
160	69 - 0 = 69
o	165 - 98 = 67
v	336 - 284 = 52
že	241 - 189 = 52

## NYU-UMontreal loses

pro	147 - 95 = 52
ve	128 - 91 = 37
rovněž	39 - 5 = 34
bylo	78 - 45 = 33
jejich	63 - 34 = 29
.	96 - 71 = 25
byla	74 - 51 = 23
být	66 - 43 = 23
řekl	115 - 95 = 20
rok	24 - 4 = 20

## 2-gram

## uedin-nmt loses

, ne	101 - 2 = 99
ne ,	101 - 3 = 98
160 ;	69 - 0 = 69
& #	
# 160	
, že	
, a	76 - 29 = 47
" řekl	42 - 0 = 42
, který	117 - 95 = 22
000 &	22 - 0 = 22

## NYU-UMontreal loses

", ,	210 - 31 = 179
, řekl	131 - 25 = 106
".	82 - 48 = 34
	20 - 1 = 19
	91 - 76 = 15
	22 - 10 = 12
, říká	23 - 11 = 12
řikal ,	16 - 4 = 12
řikají ,	15 - 4 = 11
předtím ,	12 - 1 = 11

98 vs. 0  
What is this?  
Let's click on the ngram.

## 3-gram

## uedin-nmt loses

, ne ,	100 - 1 = 99
ne , ne	98 - 0 = 98
& # 160	69 - 0 = 69
# 160 ;	69 - 0 = 69
, " řekl	42 - 0 = 42
000 & #	22 - 0 = 22
, že se	49 - 28 = 21
, a to	25 - 4 = 21
; 000 &	21 - 0 = 21
160 ; 000	21 - 0 = 21

## NYU-UMontreal loses

", řekl	100 - 1 = 99
", řekla	13 - 0 = 13
řikal , že	16 - 4 = 12
řikají , že	14 - 3 = 11
v České republice	12 - 1 = 11
", říká	11 - 0 = 11
říkala , že	12 - 2 = 10
, která byla	13 - 4 = 9
", prohlásil	9 - 0 = 9
říká , že	22 - 14 = 8

## uedin-nmt loses

, ne , ne	98 - 0 = 98
ne , ne ,	98 - 0 = 98
& # 160 ;	69 - 0 = 69
000 & # 160	22 - 0 = 22
; 000 & #	21 - 0 = 21
# 160 ; 000	21 - 0 = 21
160 ; 000 &	21 - 0 = 21
; & # 160	18 - 0 = 18
# 160 ; &	18 - 0 = 18
160 ; & #	18 - 0 = 18

## 4-gram

## NYU-UMontreal loses

", řekl .	7 - 0 = 7
, že je to	18 - 14 = 4
, řekl Navrátil .	4 - 0 = 4
", řekl Wenger	4 - 0 = 4
, řekl Wenger .	4 - 0 = 4
", řekl Navrátil	4 - 0 = 4
, že to byla	4 - 0 = 4
v Československé armády by	4 - 0 = 4
", řekla .	4 - 0 = 4

We see the Sentence pane with a filter for sentences containing n-gram “ne, ne,” which is unconfirmed only in Uedin-nmt (ie. worsening).

You are displaying sentences with worsening n-gram , **ne** , **ne**. [Show all sentences](#)

Source	A two - out walk to right fielder J . D . Martinez brought up Victor Martinez , who singled up the middle for the first run of the game .
Reference	Dva odchody pro pravého polaře J . D . Martineze vynesly <b>Victora Martineze</b> , <b>který</b> jako první oběhl všechny mety .
uedin-nmt	<div> <div>Ne</div> <div>, ne ,</div> <div>ne , ne ,</div> <div>ne , ne ,</div> <div>ne , ne ,</div> <div>ne , ne , ne , ne ,</div> </div>
NYU-UMontreal	Martinez vychoval <b>Victora Martineze</b> , <b>který</b> zpíval uprostřed prvního běhu hry .

There is just one such translation. It contains 100 tokens “ne” (no).

Uedin-nmt is the overall winner of WMT16 (including en → cs).

Apparently, there are still ways how to improve it.

There is just one such translation. It contains 100 tokens “ne” (no).

Uedin-nmt is the overall winner of WMT16 (including en → cs).

Apparently, there are still ways how to improve it.

There is just one such translation. It contains 100 tokens “ne” (no).

Uedin-nmt is the overall winner of WMT16 (including en → cs).

Apparently, there are still ways how to improve it.