

# Results of the WMT16 Metrics Shared Task

Ondřej Bojar  
Yvette Graham  
Amir Kamran  
Miloš Stanojević

WMT16, Aug 11, 2016

# Overview

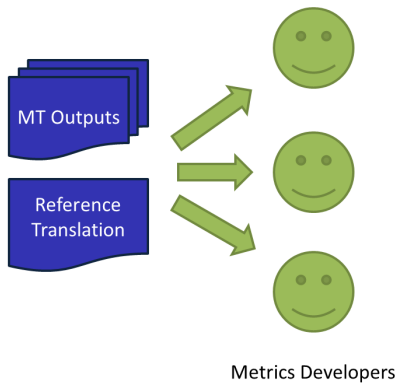
- ▶ Summary of Metrics Task.
- ▶ Updates to Metric Task in 2016.
- ▶ Results

# Metrics Task in a Nutshell

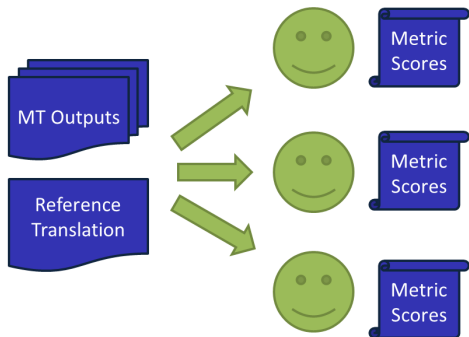
# Metrics Task in a Nutshell



# Metrics Task in a Nutshell

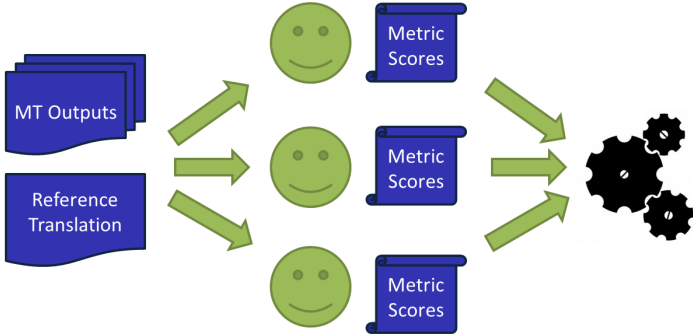


# Metrics Task in a Nutshell



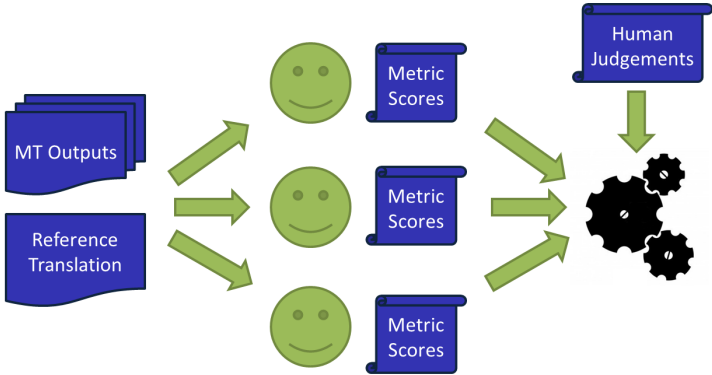
Metrics Developers

# Metrics Task in a Nutshell



Metrics Developers

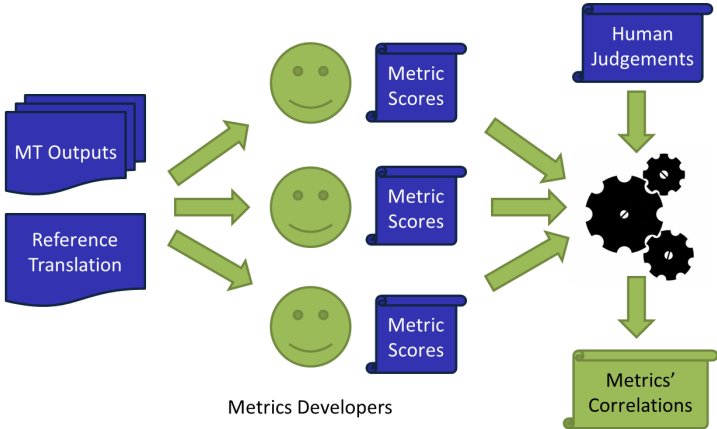
# Metrics Task in a Nutshell



Metrics Developers

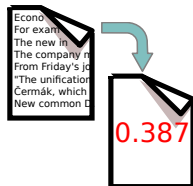


# Metrics Task in a Nutshell



# System- and Segment-Level Evaluation

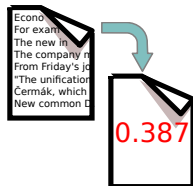
- ▶ System Level
  - ▶ Participants compute one score for the whole test set, as translated by each of the systems



# System- and Segment-Level Evaluation

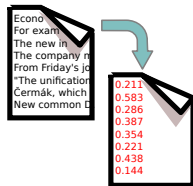
- ▶ System Level

- ▶ Participants compute one score for the whole test set, as translated by each of the systems



- ▶ Segment Level

- ▶ Participants compute one score for each sentence of each system's translation



# Nine Years of Metrics Task

	'07	'08	'09	'10	'11	'12	'13	'14	'15	'16
Participating Teams	-	6	8	14	9	8	12	12	11	<b>9</b>
Evaluated Metrics	11	16	38	26	21	12	16	23	46	<b>16</b>
Baseline Metrics					2		5	6	7	<b>9</b>
System-level										
Spearman Rank Corr.	●	●	●	●	●	●	●	○		
Pearson Corr. Coeff.							○	●	●	●
Segment-level										
Ratio of Concordant Pairs	●									
Kendall's $\tau$				●	●	●	*	*	*	*
Pearson Corr. Coeff.										●

● main and ○ secondary score reported for the system-level evaluation.

●, \* and ☆ are slightly different variants regarding ties.

- ▶ Stable number of participating teams.
- ▶ A growing set of “baseline metrics”.
- ▶ Stable but gradually improving evaluation methods.

# Updates to Metrics Task in 2016

- ▶ More Domains
  - ▶ News, IT, Medical.
- ▶ Two Golden Truths in News Task
  - ▶ Relative Ranking, Direct Assessment.
- ▶ Third golden truth in Medical Domain.
- ▶ Confidence for Sys-level Computed Differently.
  - ▶ Participants needed to score 10K systems.
- ▶ More languages (18 pairs):
  - ▶ Basque, Bulgarian, Czech, Dutch, Finnish, German, Polish, Portuguese, Romanian, Russian, Spanish, and Turkish
  - ▶ Paired with English in one or both directions.

# Metrics Task Madness

Track	Test set	Systems				English into													
		News Task	Tuning Task	IT Task	HimL Year 1	Hybrid	cs	de	ro	fi	ru	tr	bg	es	eu	nl	pl	pt	
RRsysNews	newstest2016	✓	✓			✓	•	•	•	•	•	•	•						
RRsysIT	it-test2016			✓		✓	•	•						•	•	•	•		•
DAsysNews	newstest2016	✓	✓			✓	•	•	•	•	•	•							
RRsegNews	newstest2016	✓	✓			✓	•	•	•	•	•	•							
DAssegNews	newstest2016	✓				✓	•	•	•	•	•	•							
HUMEsseg	himl2015				✓								•						•

“✓”: sets of underlying MT systems

“•”: language pairs covered in the evaluation

“.” language pairs planned but abandoned

# Metrics Task Madness

Track	Test set	Systems					English into												
		News Task	Tuning Task	IT Task	HimL	Year 1 Hybrid	cs	de	ro	fi	ru	tr	bg	es	eu	nl	pl	pt	
RRsysNews	newstest2016	✓	✓			✓	•	•	•	•	•	•							
RRsysIT	it-test2016			✓		✓	•	•	•	•	•	•	•	•	•	•	•	•	•
DAsysNews	newstest2016	✓	✓			✓	•	•	•	•	•	•	•	•	•	•	•	•	•
RRsegNews	newstest2016	✓	✓				•	•	•	•	•	•	•	•	•	•	•	•	•
DAssegNews	newstest2016	✓					•	•	•	•	•	•	•	•	•	•	•	•	•
HUMEseg	himl2015				✓		•	•	•	•	•	•	•	•	•	•	•	•	•

“✓”: sets of underlying MT systems

“•”: language pairs covered in the evaluation

“.” language pairs planned but abandoned

For participants, this was cut down to the standard:

Econo  
For exam  
The new in  
The company n  
From Friday's jo  
"The unificatio  
Cermák, which  
New common D



0.387

Sys-level

and seg-level

Econo  
For exam  
The new in  
The company n  
From Friday's jo  
"The unificatio  
Cermák, which  
New common D



0.211  
0.583  
0.286  
0.387  
0.354  
0.221  
0.438  
0.144

scoring.

# Metrics Task Domains

- ▶ WMT16 News Task
  - ▶ Systems and language pairs from the main translation task.
  - ▶ Truth: Primarily RR, DA into English and Russian.
- ▶ WMT16 IT Task
  - ▶ IT domain.
  - ▶ Only out of English.
  - ▶ Interesting target languages: (Czech, German,) Bulgarian, Spanish, Basque, Dutch, Portuguese.
  - ▶ Truth: Only RR
- ▶ HimL Medical Texts
  - ▶ Just one system per target language.
  - ▶ (So only seg-level evaluation.)
  - ▶ Truth: A new semantics-based metric.



# Golden Truths

- ▶ Relative Ranking (RR)
  - ▶ 5-way relative comparison.
    - ▶ Interpreted as 10 pairwise comparison.
    - ▶ Identical outputs deduplicated.
    - ▶ Finally converted to a score using TrueSkill.
- ▶ Direct Assessment (DA)
  - ▶ Absolute adequacy judgement over individual sentences.
    - ▶ Judgements from each worker standardized.
    - ▶ Multiple judgements of a candidate averaged.
    - ▶ Finally averaged over all sentences of a system.
  - ▶ Fluency optionally to resolve ties.
  - ▶ Provided by Turkers (only English and Russian).
    - ▶ Planned but not done with Researchers.
- ▶ HUME
  - ▶ A composite score of manual judgements of meaning preservation.
  - ▶ Used only in the “medical” track.

# Effects of DA vs. RR for Metrics Task

## Benefits:

- ▶ More principled golden truth.
- ▶ Possibly more reliable, *assuming enough judgements*.

## Negative aspects:

- ▶ Sampling for sys-level and seg-level is different.
- ▶ Perhaps impossible for seg-level out of English:
  - ▶ Too few Turker annotations.
  - ▶ ~~Too few researchers.~~ (Repeated judgements work as well.)

This year, only English and Russian news systems have DA judgements.

# Participants

<b>Metric</b>	<b>Participant</b>
BEER	ILLIC – UvA (Stanojević and Sima'an, 2015)
CHARACTER	RWTH Aachen University (Wang et al., 2016)
CHRF1,2,3	Humboldt University of Berlin (Popović, 2016)
WORDF1,2,3	Humboldt University of Berlin (Popović, 2016)
DEPCHECK	Charles University, no corresponding paper
DPMFCOMB- -WITHOUT-RED	Chinese Academy of Sciences and Dublin City University (Yu et al., 2015)
MPEDA	Jiangxi Normal University (Zhang et al., 2016)
UOW.REVAL	University of Wolverhampton (Gupta et al., 2015)
UPF-COBALT	Universitat Pompeu Fabra (Fomicheva et al., 2016)
COBALTF	Universitat Pompeu Fabra (Fomicheva et al., 2016)
METRICSF	Universitat Pompeu Fabra (Fomicheva et al., 2016)
DTED	University of St Andrews, (McCaffery and Nederhof, 2016)

# Standard Presentation of the Results

	cs-en		de-en		fi-en		ro-en		ru-en		tr-en	
Human	RR	DA	RR	DA	RR	DA	RR	DA	RR	DA	RR	DA
Systems	6	6	10	10	9	9	7	7	10	10	8	8
MPEDA	<b>.996</b>	<b>.993</b>	<b>.956</b>	.937	<b>.967</b>	<b>.976</b>	<b>.938</b>	<b>.932</b>	<b>.986</b>	.929	<b>.972</b>	<b>.982</b>
UGW.REVAL	<b>.993</b>	<b>.986</b>	<b>.949</b>	<b>.985</b>	<b>.958</b>	<b>.970</b>	<b>.919</b>	<b>.957</b>	<b>.990</b>	<b>.976</b>	<b>.977</b>	<b>.958</b>
BEER	<b>.996</b>	<b>.990</b>	.949	.879	<b>.964</b>	<b>.972</b>	<b>.908</b>	.852	<b>.986</b>	.901	<b>.981</b>	<b>.982</b>
CHRF1	.993	.986	.934	.868	<b>.974</b>	<b>.980</b>	.903	.865	<b>.984</b>	.898	<b>.973</b>	.961
CHRF2	<b>.992</b>	.989	.952	.893	<b>.957</b>	<b>.967</b>	.913	.886	<b>.985</b>	.918	.937	.933
CHRF3	.991	.989	<b>.958</b>	.902	.946	.958	.915	<b>.892</b>	<b>.981</b>	.923	.918	.917
CHARACTER	<b>.997</b>	<b>.995</b>	<b>.985</b>	.929	.921	.927	<b>.970</b>	<b>.883</b>	.955	<b>.930</b>	.799	.827
MTEVALNIST	.988	.978	.887	.801	.924	.929	.834	.807	.966	.854	<b>.952</b>	<b>.938</b>
MTEVALBLEU	<b>.992</b>	<b>.989</b>	.905	.808	.858	.864	.899	.840	.962	.837	.899	.895
MOSESCDER	<b>.995</b>	<b>.988</b>	.927	.827	.846	.860	.925	.800	.968	.855	.836	.826
MOSESTER	<b>.983</b>	<b>.969</b>	.926	.834	.852	.846	.900	.793	.962	.847	.805	.788
WORDF2	<b>.991</b>	<b>.985</b>	.897	.786	.790	.806	.905	.815	.955	.831	.807	.787
WORDF3	<b>.991</b>	<b>.985</b>	.898	.787	.786	.803	.909	.818	.955	.833	.803	.786
WORDF1	<b>.992</b>	<b>.984</b>	.894	.780	.796	.808	.890	.804	.954	.825	.806	.776
MOSESPER	.981	.970	.843	.730	.770	.767	.791	.748	<b>.974</b>	.887	<b>.947</b>	.940
MOSESBLEU	<b>.991</b>	<b>.983</b>	.880	.757	.752	.759	.878	.793	.950	.817	.765	.739
MOSESWER	<b>.982</b>	<b>.967</b>	.926	.822	.773	.768	.895	.762	.958	.837	.680	.651

newstest2016

- ▶ Bold in RR indicates “official winners”.
- ▶ Some setups fairly non-discerning, here e.g. csen:
  - ▶ All but CHRF1, CHRF3, MTEVALNIST and MOSESPER tie at top.

# News RR Winners Across Languages

Metric	# Wins	Language Pairs
BEER	11	cse, encs, ende, enfi, enro, enru, entr, fi, roen, ruen, tren
UoW.ReVal	6	cse, deen, fi, roen, ruen, tren
chrF2	6	cse, encs, enro, entr, fi, ruen
chrF1	5	encs, enro, fi, ruen, tren
chrF3	4	deen, enfi, entr, ruen
mosesCDER	4	cse, enfi, enru, entr
CharacTer	3	cse, deen, roen
mosesBLEU	3	cse, encs, enfi
mosesPER	3	enro, ruen, tren
mtevalBLEU	3	cse, encs, enro
wordF1	3	cse, encs, enro
wordF2	3	cse, encs, enro
mosesTER	2	cse, encs
mtevalNIST	2	encs, tren
wordF3	2	cse, entr
mosesWER	1	cse

# Graphical Presentation of Significant Wins

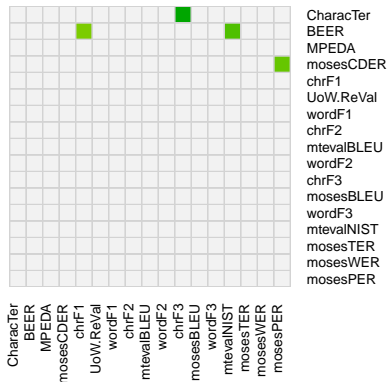
- ▶ Williams (1959) test of significant improvement in Pearson correlation.
  - ▶ Green cell indicates that the metric in the row has significantly better correlation than the metric in the column.

So for Czech-English RR, we have:

# Graphical Presentation of Significant Wins

- ▶ Williams (1959) test of significant improvement in Pearson correlation.
  - ▶ Green cell indicates that the metric in the row has significantly better correlation than the metric in the column.

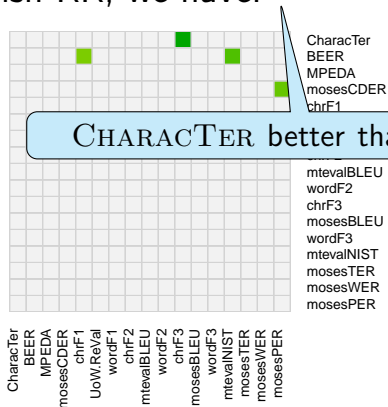
So for Czech-English RR, we have:



# Graphical Presentation of Significant Wins

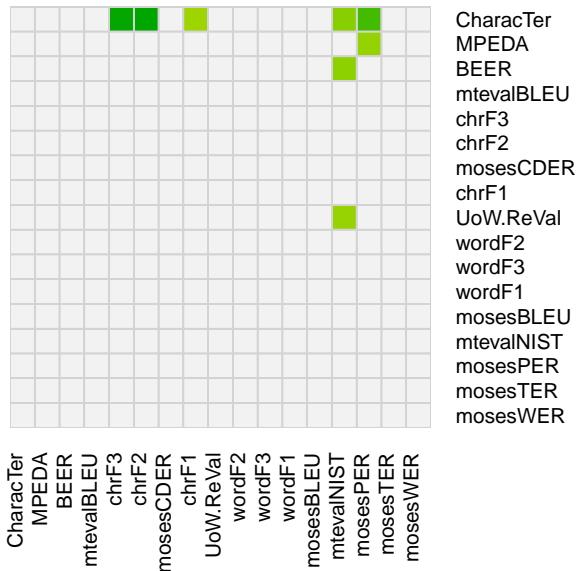
- ▶ Williams (1959) test of significant improvement in Pearson correlation.
  - ▶ Green cell indicates that the metric in the row has significantly better correlation than the metric in the column.

So for Czech-English RR, we have:

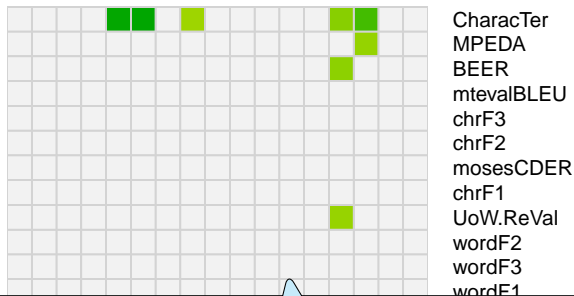




# Czech-English Direct Assessments



# Czech-English Direct Assessments

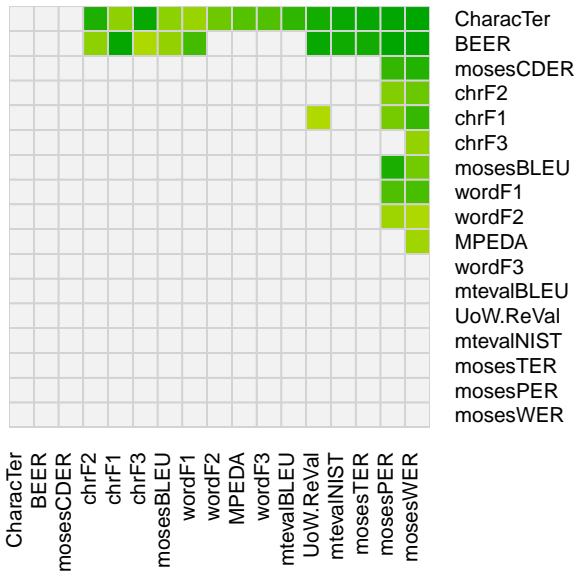


With just 6 systems,  
correlations do not differ reliably.

CharacTer  
MPEDA  
BEER  
mtevalBLEU  
chrF3  
chrF2  
mosesCDER  
chrF1  
UoW.ReVal  
wordF2  
wordF3  
wordF1  
mosesBLEU  
mtevalNIST  
mosesPER  
mosesTER  
mosesWER

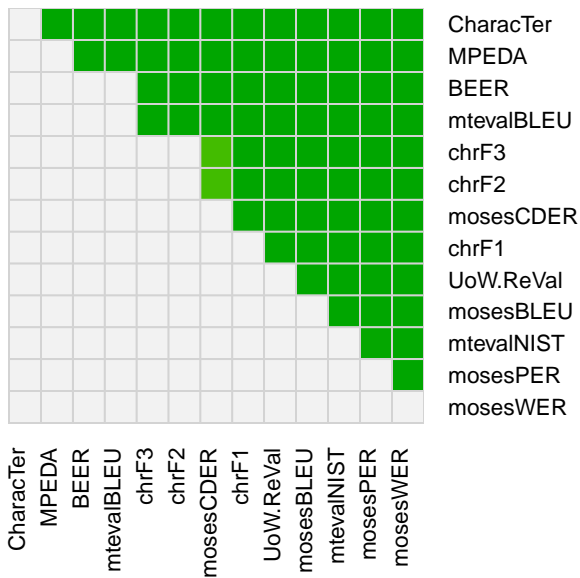
mosesWER

# Czech-English RR with Tuning Systems

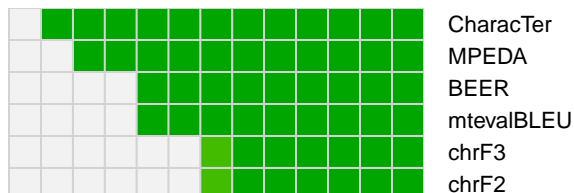




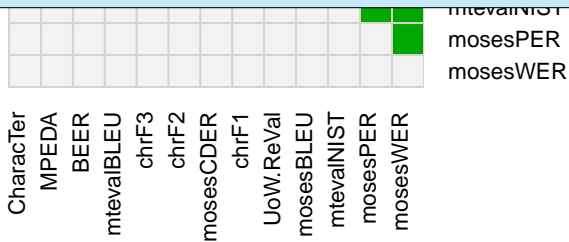
# Czech-English DA with Hybrids



# Czech-English DA with Hybrids



10,000 synthesized systems  
allow to find almost  
total ordering.



CharacTer  
MPEDA  
BEER  
mtevalBLEU  
chrF3  
chrF2  
mosesCDER  
chrF1  
UoW.ReVal  
mosesBLEU  
mtevalNIST  
mosesPER  
mosesWER

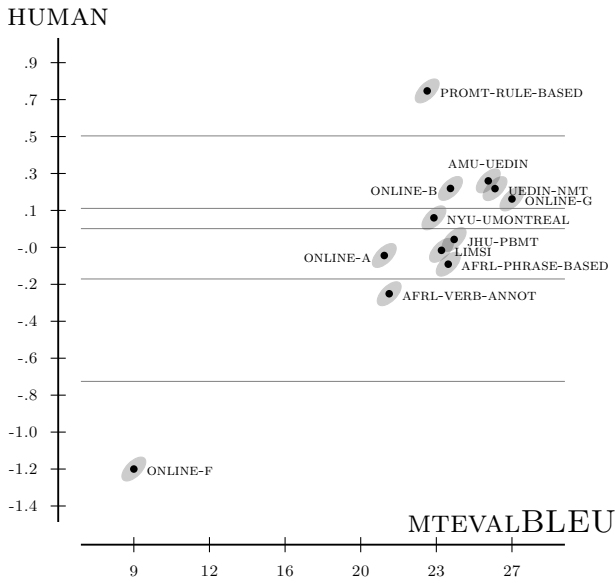
# “Hybrids” = Hybrid Super-Sampling

- ▶ 10,000 “new systems” constructed by mixing sentences.
- ▶ Puts extra burden on task participants.
  - ▶ Need to score 10k “system” outputs, full test set each.
  - ▶ 200MB–1.1GB bzipped input file per language pair.
- ▶ Allows to distinguish sys-level metrics much better.
- ▶ Applicable to both RR and DA.
  - ▶ Done with DA only now because RR human judgements of individual sentences have to be carried over to these 10k systems.

Winners according to DA hybrids:

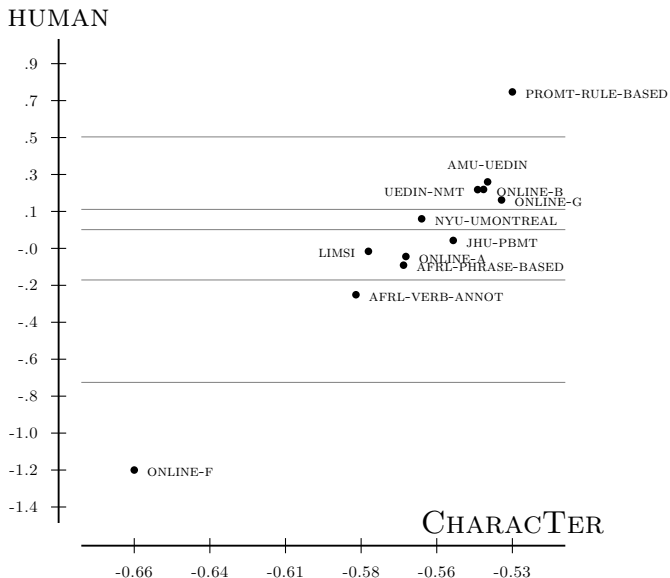
Metric	# Wins	Language Pairs
UoW.ReVal	3	deen, roen, ruen
CharacTer	2	csen, enru
MPEDA	2	fien, tren
BEER	1	tren

# Ex. English-Russian RR BLEU

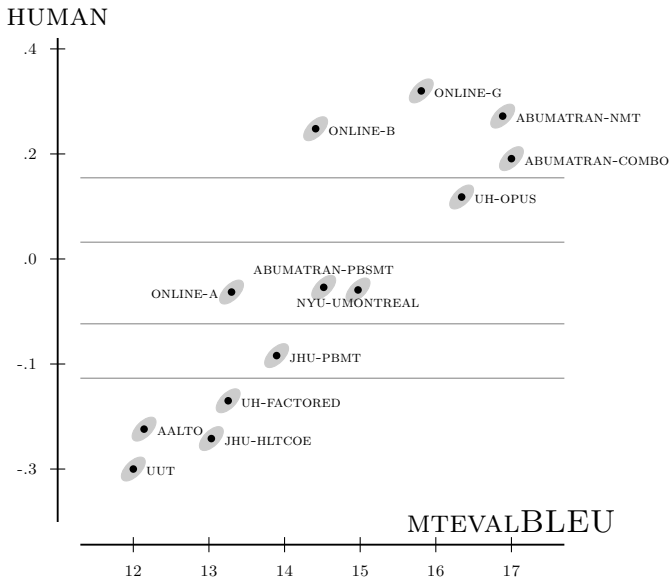




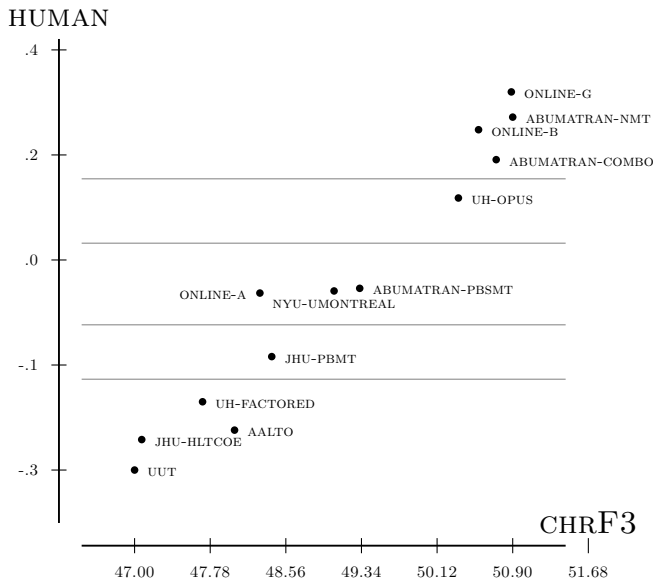
# Ex. English-Russian RR CharacTer



# Ex. English-Finnish RR BLEU



# Ex. English-Finnish RR chrF3



# Sys-Level Metrics on IT Task

- ▶ To test metrics in domain-specific setting.
- ▶ Unfortunately, often too few participating systems.

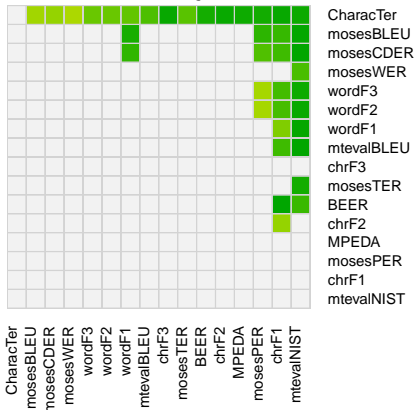
	en-bg	en-cs	en-de	en-es	en-eu	en-nl	en-pt
<b>Human</b>	RR	RR	RR	RR	RR	RR	RR
<b>Systems</b>	2	5	10	4	2	4	4
CHARACTER	<b>1.000</b>	<b>0.901</b>	<b>0.930</b>	<b>0.963</b>	<b>1.000</b>	<b>0.927</b>	0.976
CHRF3	<b>1.000</b>	<b>0.831</b>	0.700	<b>0.938</b>	<b>1.000</b>	<b>0.961</b>	0.990
CHRF2	<b>1.000</b>	0.837	0.672	0.933	<b>1.000</b>	<b>0.959</b>	0.986
BEER	<b>1.000</b>	<b>0.744</b>	0.621	0.931	<b>1.000</b>	<b>0.983</b>	<b>0.989</b>
CHRF1	<b>1.000</b>	<b>0.845</b>	0.588	0.915	<b>1.000</b>	<b>0.951</b>	0.967
MTEVALNIST	<b>1.000</b>	<b>0.905</b>	0.524	0.926	<b>1.000</b>	0.722	<b>0.993</b>
MPEDA	<b>1.000</b>	0.620	0.599	<b>0.951</b>	<b>1.000</b>	0.856	<b>0.989</b>
MOSESTER	<b>1.000</b>	<b>0.616</b>	0.628	<b>0.908</b>	<b>1.000</b>	0.835	<b>0.994</b>
MTEVALBLEU	<b>1.000</b>	<b>0.750</b>	0.621	<b>0.976</b>	<b>1.000</b>	0.596	<b>0.997</b>
MOSESWER	<b>1.000</b>	<b>0.009</b>	0.656	0.916	<b>1.000</b>	<b>0.903</b>	0.991
MOSECDER	<b>1.000</b>	0.181	0.652	<b>0.932</b>	<b>1.000</b>	<b>0.914</b>	<b>0.997</b>
WORDF1	<b>1.000</b>	0.240	0.644	<b>0.959</b>	<b>1.000</b>	<b>0.911</b>	<b>0.997</b>
WORDF2	<b>1.000</b>	0.266	0.652	<b>0.965</b>	<b>1.000</b>	<b>0.900</b>	<b>0.997</b>
WORDF3	<b>1.000</b>	0.274	0.655	<b>0.966</b>	<b>1.000</b>	<b>0.897</b>	<b>0.996</b>
MOSEBLEU	<b>1.000</b>	0.296	0.650	<b>0.974</b>	<b>1.000</b>	0.886	<b>0.992</b>
MOSEPER	<b>1.000</b>	0.307	0.548	0.911	<b>1.000</b>	<b>0.938</b>	<b>0.998</b>

ittest2016

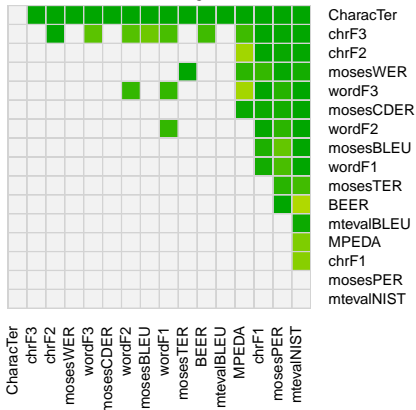
... so only English-German tells us something.

# English-German RR News vs. IT

News  
15 Systems



IT  
10 Systems



▶ CHARACTER wins in both domains.

# Segment-Level News Task Evaluation

## Relative Ranking

- ⊕ Genuine comparisons
- ⊖ 5-way comparison hard?
- ⊖ Non-standard Kendall's  $\tau$
- ⊖ Conf. estimation unclear

## Direct Assessment

- ⊖ Only 1 candidate shown
- ⊕ Principled Pearson
- ⊖ Distinct sampling needed

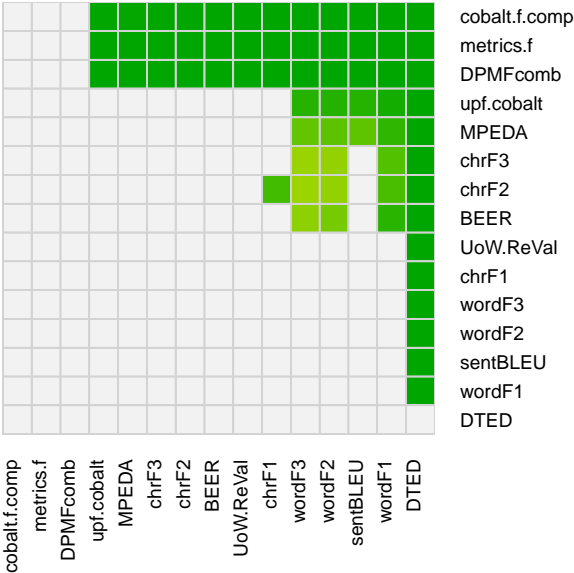
# Segment-Level News Task Results

- ▶ DA and RR correlate at .85–.99 (.92 avg across langs).
- ▶ Top RR metric always among DA winners.
- ▶ RR Winners

Metric	# Wins	Language Pairs
BEER	4	encs, ende, enro, entr
DPMFcomb	3	cseu, fienu, ruen
metrics-f	3	deeu, roenu, tren
chrF2	1	enru
chrF3	1	enfi

- ▶ Williams' test for DA reveals more top-performing metrics:
  - ▶ cobalt-f (deeu, ruen), MPEDA (enru)

# Ex. Russian-English DA Significance

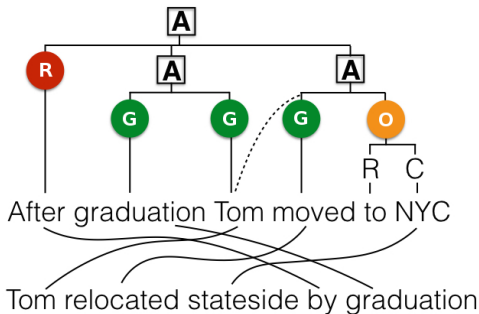




# Semantic Golden Truth (HUME)

HUME (Birch et al., 2016) uses two-stage annotation:

1. Semantic annotation (structure) of *source*.
2. Correctness assessment of corresponding parts of candidate.



- ▶ Final sentence-level score aggregated over source components.

# HUME in Metrics Task

- ▶ A first probe.
- ▶ One test set:
  - ▶ Medical texts from Cochrane and NHS24
  - ▶ Translated by year 1 MT systems of the EU project HimL.
  - ▶ Source English annotated once.
  - ▶ Targets: Czech, German, Romanian, Polish
  - ▶ ~340 sentences
- ▶ Used only in segment-level evaluation.

# Results of Semantic Evaluation

Direction	en-cs	en-de	en-ro	en-pl
<i>n</i>	339	330	349	345
CHRF3	<b>.544</b>	<b>.480</b>	<b>.639</b>	.413
CHRF2	.537	<b>.479</b>	.634	<b>.417</b>
BEER	.516	<b>.480</b>	.620	<b>.435</b>
CHRF1	.506	<b>.467</b>	.611	<b>.427</b>
MPEDA	.468	<b>.478</b>	.595	<b>.425</b>
WORDF3	.413	.425	.587	.383
WORDF2	.408	.424	.583	.383
WORDF1	.392	.415	.569	.381
SENTBLEU	.349	.377	.550	.328

- ▶ Bold again indicates metrics not significantly outperformed by any other (Williams, 1959).
- ▶ CHRF3 and other character-level metrics clearly win.
- ▶ SENTBLEU by far the worst.

# Summary

- ▶ The **2017 golden truth** will follow the main translation task.
- ▶ Whether DA or RR, we will **use hybrids** for sys-level.
- ▶ Domain-specific evaluation of metrics needs **enough systems** to participate (or plan seg-level evaluation).
- ▶ Top metrics consider again **character sequences** and are **trained**.
- ▶ Even “semantics” seems well captured by **character-level metrics**.

# References

- Alexandra Birch, Barry Haddow, Ondřej Bojar, and Omri Abend. 2016. Hume: Human ucca-based evaluation of machine translation. *arXiv preprint arXiv:1607.00030* .
- Marina Fomicheva, Núria Bel, Lucia Specia, Iria da Cunha, and Anton Malinovskiy. 2016. CobaltF: A Fluent Metric for MT Evaluation. In *Proceedings of the First Conference on Machine Translation*. Association for Computational Linguistics, Berlin, Germany.
- Rohit Gupta, Constantin Orăsan, and Josef van Genabith. 2015. ReVal: A Simple and Effective Machine Translation Evaluation Metric Based on Recurrent Neural Networks. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Lisbon, Portugal.
- Martin McCaffery and Mark-Jan Nederhof. 2016. DTED: Evaluation of Machine Translation Structure Using Dependency Parsing and Tree Edit Distance. In *Proceedings of the First Conference on Machine Translation*. Association for Computational Linguistics, Berlin, Germany.
- Maja Popović. 2016. chrF deconstructed: beta parameters and n-gram weights. In *Proceedings of the First Conference on Machine Translation*. Association for Computational Linguistics, Berlin, Germany.
- Miloš Stanojević and Khalil Sima'an. 2015. BEER 1.1: ILLC UvA submission to metrics and tuning task. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*. Association for Computational Linguistics, Lisboa, Portugal.
- Weiyue Wang, Jan-Thorsten Peter, Hendrik Rosendahl, and Hermann Ney. 2016. CharacTer: Translation Edit Rate on Character Level. In *Proceedings of the First Conference on Machine Translation*. Association for Computational Linguistics, Berlin, Germany.
- Evan James Williams. 1959. *Regression analysis*, volume 14. Wiley New York.
- Hui Yu, Qingsong Ma, Xiaofeng Wu, and Qun Liu. 2015. CASICT-DCU Participation in WMT2015 Metrics Task. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*. Association for Computational Linguistics, Lisboa, Portugal.