

# Machine Translation Evaluation

Gongbo Tang

2021-08-31

Mainly from Sara's slides.

# Why Evaluation?

- How good is a given machine translation system?
- Which one is the best system for our purpose?
- How much did we improve our system?
- How can we tune our system to become better?
- Hard problem, since many different translations acceptable  
→ semantic equivalence / similarity

# Ten Translations of a Chinese Sentence

这个机场的安全工作由以色列方面负责。

Israeli officials are responsible for airport security.

Israel is in charge of the security at this airport.

The security work for this airport is the responsibility of the Israel government.

Israeli side was in charge of the security of this airport.

Israel is responsible for the airport's security.

Israel is responsible for safety work at this airport.

Israel presides over the security of the airport.

Israel took charge of the airport security.

The safety of this airport is taken charge of by Israel.

This airport's security is the responsibility of the Israeli security officials.

(a typical example from the 2001 NIST evaluation set)

# Which translation is best? worst?

**Source** Färjetransporterna har minskat med 20,3 procent i år.

**Gloss** The-ferry-transports have decreased by 20.3 percent in year.

**Ref** Ferry transports are down by 20.3% in 2008.

---

# Which translation is best? worst?

**Source** Färjetransporterna har minskat med 20,3 procent i år.

**Gloss** The-ferry-transports have decreased by 20.3 percent in year.

**Ref** Ferry transports are down by 20.3% in 2008.

---

**Sys1** The ferry transports has reduced by 20.3% in year.

**Sys2** This year, the reduction of transports by ferry is 20,3 percent.

**Sys3** Färjetransporterna are down by 20.3% this year.

**Sys4** Ferry transports have a reduction of 20.3 percent in year.

**Sys5** Transports are down by 20.3 this year%.

# Evaluation Methods

- Subjective judgments by human evaluators
- Task-based evaluation, e.g.:
  - How much post-editing effort?
- Automatic evaluation metrics
- Test suites
- Quality estimation

# Human vs Automatic Evaluation

- Human evaluation is
  - Ultimately what we are interested in, but
  - Very time consuming
  - Not re-usable
  - Subjective
- Automatic evaluation is
  - Cheap and re-usable, but
  - Not necessarily reliable

# Human evaluation

- Adequacy/Fluency (1 to 5 scale)
- Ranking of systems (best to worst)
- Yes/no assessments (acceptable translation?)
- SSER – subjective sentence error rate (“perfect” to “absolutely wrong”)
- Usability (Good, useful, useless)
- Human post-editing time
- Error analysis



# Adequacy and Fluency

- given: machine translation output
- given: source and/or reference translation
- task: assess the quality of the machine translation output

**Adequacy:** Does the output convey the same meaning as the input sentence?

Is part of the message lost, added, or distorted?

**Fluency:** Is the output good fluent target language?  
This involves both grammatical correctness and idiomatic word choices.

# Fluency and Adequacy: Scales

<b>Adequacy</b>	
5	all meaning
4	most meaning
3	much meaning
2	little meaning
1	none

<b>Fluency</b>	
5	flawless English
4	good English
3	non-native English
2	disfluent English
1	incomprehensible

# Judge adequacy and fluency!

**Source** Färjetransporterna har minskat med 20,3 procent i år.

**Gloss** The-ferry-transports have decreased by 20.3 percent in year.

**Ref** Ferry transports are down by 20.3% in 2008.

---

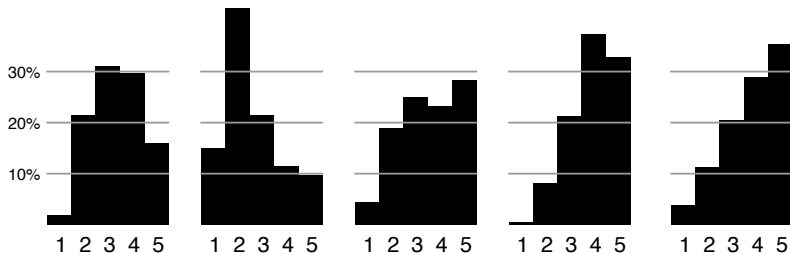
**Sys4** Ferry transports have a reduction of 20.3 percent in year.

**Sys6** Transports are down by 20.3%.

**Sys7** This year, of transports by ferry reduction is percent 20.3.

# Evaluators Disagree

- Histogram of adequacy judgments by different human evaluators



(from WMT 2006 evaluation)

# Measuring Agreement between Evaluators

- Kappa coefficient

$$K = \frac{p(A) - p(E)}{1 - p(E)}$$

- $p(A)$ : proportion of times that the evaluators agree
- $p(E)$ : proportion of time that they would agree by chance
- Example: Inter-evaluator agreement in WMT 2007 evaluation campaign

Evaluation type	$P(A)$	$P(E)$	$K$
Fluency	.400	.2	.250
Adequacy	.380	.2	.226

# Ranking Translations

- Task for evaluator: Is translation X better than translation Y?  
(choices: better, worse, equal)

- Evaluators are more consistent:

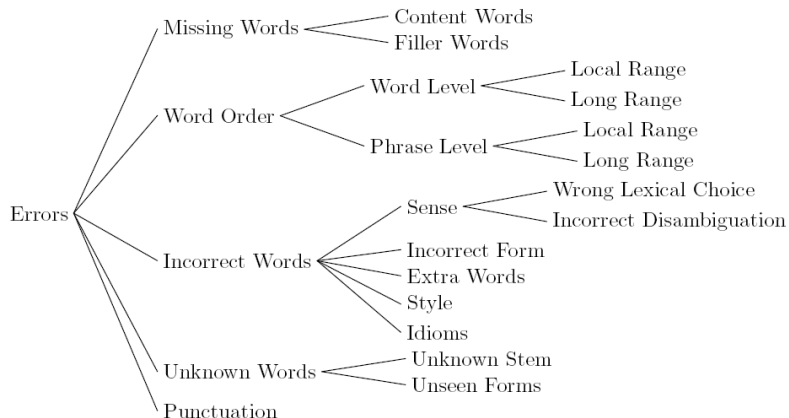
Evaluation type	$P(A)$	$P(E)$	$K$
Fluency	.400	.2	.250
Adequacy	.380	.2	.226
Sentence ranking	.582	.333	.373

# Error Analysis

- Analysis and classification of the errors from an MT system
- Many general frameworks for classification exists, e.g.
  - Flanagan, 1994
  - Vilar et al. 2006
  - Costa-jussà et al. 2012
- It is also possible to analyse specific phenomena, like compound translation, agreement, pronoun translation, ...

# Example Error Typology

Vilar et al.





# Task-Oriented Evaluation

- Machine translations is a means to an end
- Does machine translation output help accomplish a task?
- Example tasks
  - producing translations good enough for post-editing machine translation
  - information gathering from foreign language sources

# Post-Editing Machine Translation

- Measuring time spent on producing translations
  - baseline: translation from scratch (often using TMs)
  - post-editing machine translation
- Some issues:
  - time consuming
  - depends on skills of particular translators/post-editors

# Content Understanding Tests

- Given machine translation output, can monolingual target side speaker answer questions about it?
  1. basic facts: who? where? when? names, numbers, and dates
  2. actors and events: relationships, temporal and causal order
  3. nuance and author intent: emphasis and subtext
- Very hard to devise questions

# Automatic Evaluation Metrics

- Goal: computer program that computes the quality of translations
- Advantages: low cost, tunable, consistent
- Basic strategy
  - given: machine translation output
  - given: human reference translation
  - task: compute similarity between them

# Goals for Evaluation Metrics

**Low cost:** reduce time and money spent on carrying out evaluation

**Tunable:** automatically optimize system performance towards metric

**Meaningful:** score should give intuitive interpretation of translation quality

**Consistent:** repeated use of metric should give same results

**Correct:** metric must rank better systems higher

# Other Evaluation Criteria

When deploying systems, considerations go beyond quality of translations

**Speed:** we prefer faster machine translation systems

**Size:** fits into memory of available machines (e.g., handheld devices)

**Integration:** can be integrated into existing workflow

**Customization:** can be adapted to user's needs

# Metrics – overview

- Precision-based
  - BLEU, NIST, ...
- F-score-based
  - Meteor, ChrF...
- Error rates
  - WER, TER, PER, ...
- Using syntax/semantics
  - PosBleu, Meant, DepRef, ...
- Using machine learning
  - TerrorCat, Beer, CobaltF

# Metrics – overview

- Precision-based
  - **BLEU**, NIST, ...
- F-score-based
  - **Meteor**, ChrF...
- Error rates
  - **WER**, **TER**, PER, ...
- Using syntax/semantics
  - PosBleu, Meant, DepRef, ...
- Using machine learning
  - TerrorCat, Beer, CobaltF



# Precision and Recall of Words

SYSTEM A: Israeli officials responsibility of airport safety

REFERENCE: Israeli officials are responsible for airport security

■ Precision  $\frac{\textit{correct}}{\textit{output-length}} = \frac{3}{6} = 50\%$

■ Recall  $\frac{\textit{correct}}{\textit{reference-length}} = \frac{3}{7} = 43\%$

■ F-measure  $\frac{\textit{precision} \times \textit{recall}}{(\textit{precision} + \textit{recall})/2} = \frac{.5 \times .43}{(.5 + .43)/2} = 46\%$



# BLEU

- N-gram overlap between machine translation output and reference translation
- Compute precision for n-grams of size 1 to 4
- Add brevity penalty (for too short translations)

$$\text{BLEU} = \min \left( 1, \frac{\text{output-length}}{\text{reference-length}} \right) \left( \prod_{i=1}^4 \text{precision}_i \right)^{\frac{1}{4}}$$

- Typically computed over the entire corpus, not single sentences

# Example

SYSTEM A: Israeli officials responsibility of airport safety  
2-GRAM MATCH 1-GRAM MATCH

REFERENCE: Israeli officials are responsible for airport security

SYSTEM B: airport security Israeli officials are responsible  
2-GRAM MATCH 4-GRAM MATCH

Metric	System A	System B
precision (1gram)	3/6	6/6
precision (2gram)	1/5	4/5
precision (3gram)	0/4	2/4
precision (4gram)	0/3	1/3
brevity penalty	6/7	6/7
BLEU	0%	52%

# Multiple Reference Translations

- To account for variability, use multiple reference translations
  - n-grams may match in any of the references
  - closest reference length used (usually)
- Example

SYSTEM:                    Israeli officials   responsibility of   airport   safety  
                                 2-GRAM MATCH            2-GRAM MATCH            1-GRAM

REFERENCES:            Israeli officials are responsible for airport security  
                                 Israel is in charge of the security at this airport  
                                 The security work for this airport is the responsibility of the Israel government  
                                 Israeli side was in charge of the security of this airport

# Notes on BLEU

- BLEU calculation:
  - NIST-BLEU
  - Moses: mteval-v\*.pl
  - Moses: multi-bleu.perl
  - Moses: multi-bleu-detok.perl
  - SacreBLEU

You should keep the same metric when you compare your models with previous models!

# METEOR: Flexible Matching

- Partial credit for matching stems

SYSTEM    Jim walk home  
REFERENCE   Joe walks home

- Partial credit for matching synonyms

SYSTEM    Jim strolls home  
REFERENCE   Joe walks home

- Use of paraphrases
- Different weights for content and function words (later versions)

# METEOR

- Both recall and precision
- Only unigrams (not higher n-grams)
- Flexible matching (Weighted P and R)
- Fluency captured by a penalty for high number of chunks

$$F_{mean} = \frac{PR}{\alpha \cdot P + (1 - \alpha) \cdot R}$$

$$Penalty = 0.5 * \gamma \cdot \left( \frac{\#chunks}{\#unigrams\_matched} \right)^\beta$$

$$Meteor = (1 - Penalty) \cdot F_{mean}$$



# METEOR: tuning

- Meteor parameters can be tuned based on human judgments

Language	$\alpha$	$\beta$	$\gamma$	$\delta$	$w_{exact}$	$w_{stem}$	$w_{syn}$	$w_{par}$
Universal	.70	1.40	.30	.70	1.00	–	–	.60
English	.85	.20	.60	.75	1.00	.60	.80	.60
French	.90	1.40	.60	.65	1.00	.20	–	.40
German	.95	1.00	.55	.55	1.00	.80	–	.20

# Word Error Rate

- Minimum number of editing steps to transform output to reference
  - match:** words match, no cost
  - substitution:** replace one word with another
  - insertion:** add word
  - deletion:** drop word
- Levenshtein distance

$$\text{WER} = \frac{\textit{substitutions} + \textit{insertions} + \textit{deletions}}{\textit{reference-length}}$$

# Example

	0	1	2	3	4	5	6
Israeli	1	0	1	2	3	4	5
officials	2	1	0	1	2	3	4
are	3	2	1	1	2	3	4
responsible	4	3	2	2	2	3	4
for	5	4	3	3	3	3	4
airport	6	5	4	4	4	3	4
security	7	6	5	5	5	4	4

	0	1	2	3	4	5	6
Israeli	1	1	2	2	3	4	5
officials	2	2	2	3	2	3	4
are	3	3	3	3	3	2	3
responsible	4	4	4	4	4	3	2
for	5	5	5	5	5	4	3
airport	6	5	6	6	6	5	4
security	7	6	5	6	7	6	5

Metric	System A	System B
word error rate (WER)	57%	71%

## Other error rates

- PER – position-independent word error rate
  - Does not consider the order of words
- TER – translation edit rate
  - Adds the operation SHIFT – the movement of a contiguous sequence of words an arbitrary distance
- SER – sentence error rate
  - The percentage of sentences that are identical to reference sentences

# Metrics using syntax/semantics

- Posbleu, Bleu calculated on part-of-speech
- ULC – Overlap of:
  - shallow parsing
  - dependency and constituent parsing
  - named entities
  - semantic roles
  - discourse representation structures
- Using dependency structures
- Meant, semantic roles
- Considerations:
  - parsers/taggers do not perform well on misformed MT output
  - parsers/tagger not available for all languages

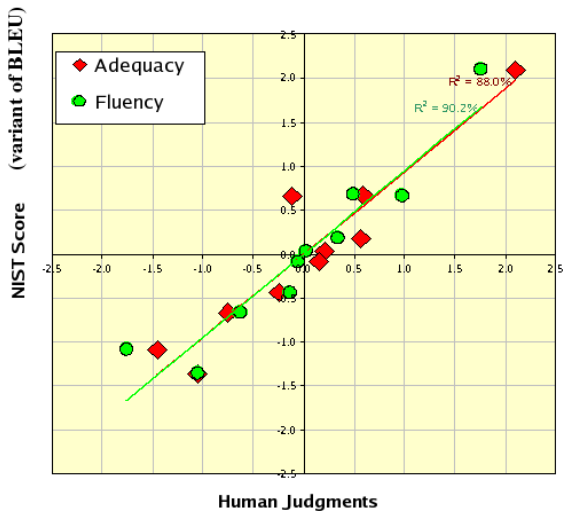
# Critique of Automatic Metrics

- Ignore relevance of words  
(names and core concepts more important than determiners and punctuation)
- Operate on local level  
(do not consider overall grammaticality of the sentence or sentence meaning)
- Scores are meaningless  
(scores very test-set specific, absolute value not informative)
- Human translators score low on BLEU  
(possibly because of higher variability, different word choices)

# Evaluation of Evaluation Metrics

- Automatic metrics are low cost, tunable, consistent
  - But are they correct?
- Yes, if they correlate with human judgement

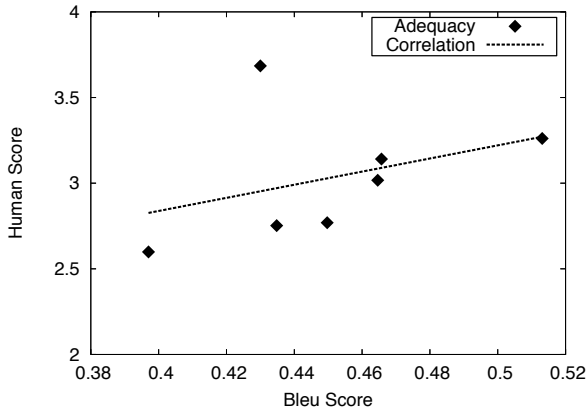
# Correlation with Human Judgement





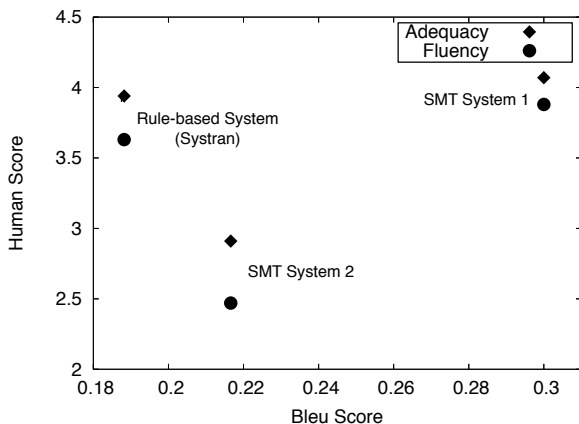
# Evidence of Shortcomings of Automatic Metrics

Post-edited output vs. statistical systems (NIST 2005)



# Evidence of Shortcomings of Automatic Metrics

## Rule-based vs. statistical systems



# Metric Research

- Active development of new metrics
  - syntactic similarity
  - semantic equivalence or entailment
  - metrics targeted at reordering
  - trainable metrics
  - neural network-based:
    - DREEM (Distributed Representations Evaluation Metrics)
    - BERTscore
  - etc.
  
- Evaluation campaigns that rank metrics (using Pearson's correlation coefficient)

# Correlations of metrics with human ranking

Metric	de-en	en-de
BLEU	.88	.76
WER	.93	.82
PER	.84	.73
ChrF1	.93	.87
ChrF3	.96	.90
Beer	.95	.88

(System level, WMT 2017)

# Correlations of metrics with human ranking

Metric	de-en	en-de
BLEU	.90	.79
METEOR	.96	.88
TER	.83	.85
WER	.67	.83
TERRORCAT	.96	.95
DEPREF-ALIGN	.97	–

(System level, WMT 2013)

# Correlations of metrics with human ranking

Metric	de-en	en-de
SentBLEU	.27	.45
ChrF1	.37	.45
ChrF3	.35	.46
Cobalt-F	.42	.59

(Segment level, WMT 2017)

# Correlations of metrics with human ranking

Metric	de-en	en-de
BLEU	.23	.18
METEOR	.26	.24
TERRORCAT	.25	.21
DEPREF-ALIGN	.26	–

(Segment level, WMT 2013)

# Automatic Metrics: Conclusions

- Automatic metrics essential tool for system development
- Not fully suited to rank systems of different types
- Reasonable results on system level evaluation, but not on sentence level
- Evaluation metrics still open challenge



# Test suites / Challenge sets

- Create a test set targeting a specific phenomena you want to evaluate
- Translate it using MT systems of interest
- Evaluate how well your specific phenomena is translated
  - Automatically
  - By humans
- Was used to some extent for rule-based MT
- Has had a recent revival (e.g. WMT 2018)
- Test suites can be reused (but may require human scoring)

# Challenge Set: An example

category	English	German (correct)	German (contrastive)
NP agreement	[...] <b>of the American Congress</b>	[...] <b>des</b> amerikanischen <b>Kongresses</b>	* [...] <b>der</b> amerikanischen <b>Kongresses</b>
subject-verb agr.	[...] that the <b>plan will</b> be approved	[...], dass der <b>Plan</b> verabschiedet <b>wird</b>	* [...], dass der <b>Plan</b> verabschiedet <b>werden</b>
separable verb particle	he is <b>resting</b>	er <b>ruht</b> sich <b>aus</b>	* er <b>ruht</b> sich <b>an</b>
polarity	the timing [...] is <b>uncertain</b>	das Timing [...] ist <b>unsicher</b>	das Timing [...] ist <b>sicher</b>
transliteration	Mr. <b>Ensign's</b> office	Senator <b>Ensigns</b> Büro	Senator <b>Enisgns</b> Büro

Examples from Lingeval97 data set (Sennrich, 2017)

# Quality Estimation

- For standard automatic metrics, a reference translation is needed
- In a realistic translation scenario, we do not have reference translations
- It is very useful for a translator who is presented MT output to know:
  - Is it good enough as it is
  - Can it be easily edited
  - Can it be edited with some effort
  - Is it completely useless
- This task is called quality estimation

# Quality Estimation – Details

- Automatic evaluation without a reference
- Typically modelled as a machine learning task
- Using features such as:
  - How long is the sentence?
  - What is the length difference between the source and target?
  - How common are the words and n-grams in the source sentence?
  - How ambiguous are the words in the source sentence?
  - How many punctuation marks are there in the sentence?
- Train on judgments of fluency/adequacy, post-editing effort, or post-editing time

# Hypothesis Testing

- Situation
  - system A has score  $x$  on a test set
  - system B has score  $y$  on the same test set
  - $x > y$
- Is system A really better than system B?
- In other words:  
Is the difference in score **statistically significant**?

# Core Concepts

- Null hypothesis
  - assumption that there is no real difference
- P-Levels (statistical significance)
  - related to probability that there is a true difference
  - p-level  $p < 0.01$  = more than 99% chance that difference is real
  - typically used: p-level 0.05 or 0.01
- Confidence Intervals
  - given that the measured score is  $x$
  - what is the true score (on an infinite size test set)?
  - interval  $[x - d, x + d]$  contains true score with, e.g., 95% probability
  - whether the intervals overlap → different or not

# Pairwise Comparison

Score itself is not interesting/meaningful, what else I can do?

- Example
  - Given a test set of 100 sentences
  - System A better on 60 sentence
  - System B better on 40 sentences
- Is system A really better?

# Sign Test

- Using binomial distribution
  - system A better with probability  $p_A$
  - system B better with probability  $p_B (= 1 - p_A)$
  - probability of system A better on  $k$  sentences out of a sample of  $n$  sentences

$$\binom{n}{k} p_A^k p_B^{n-k} = \frac{n!}{k!(n-k)!} p_A^k p_B^{n-k}$$

- Null hypothesis:  $p_A = p_B = 0.5$

$$\binom{n}{k} p^k (1-p)^{n-k} = \binom{n}{k} 0.5^n = \frac{n!}{k!(n-k)!} 0.5^n$$



# Examples

$n$	$p \leq 0.01$		$p \leq 0.05$	
5	-	-	-	-
10	$k = 10$	$\frac{k}{n} = 1.00$	$k \geq 9$	$\frac{k}{n} \geq 0.90$
20	$k \geq 17$	$\frac{k}{n} \geq 0.85$	$k \geq 15$	$\frac{k}{n} \geq 0.75$
50	$k \geq 35$	$\frac{k}{n} \geq 0.70$	$k \geq 33$	$\frac{k}{n} \geq 0.66$
100	$k \geq 64$	$\frac{k}{n} \geq 0.64$	$k \geq 61$	$\frac{k}{n} \geq 0.61$

Given  $n$  sentences  
system has to be better in at least  $k$  sentences  
to achieve statistical significance at specified p-level

# Bootstrap Resampling

- Described methods require scores at sentence level
- But: common metrics such as BLEU are computed for whole corpus
- Sampling
  - 1 test set of 2000 sentences, sampled from large collection
  - 2 compute the BLEU score for this set
  - 3 repeat step 1–2 for 1000 times
  - 4 ignore 25 highest and 25 lowest obtained BLEU scores  
→ 95% confidence interval
- Bootstrap resampling
  - Sample 1000 sentences from the same 2000 sentences, with replacement

# WMT Shared Tasks

Link: <http://www.statmt.org/wmt21/>

- Evaluation Tasks
  - Quality Estimation
  - Evaluation metrics
- Translation Tasks
  - News Translation
  - Similar languages translation
  - Biomedical translation
  - European low resource multilingual translation
  - Large-scale multilingual translation
  - Trilingual MT
  - Unsupervised very low-resource MT
  - Terminology
- Other Tasks
  - Automatic post-editing

# Summary

- MT evaluation is hard
- Human evaluation is expensive
- Automatic evaluation is cheap, but not always fair
- What is typically used in MT research:
  - Bleu! (sometimes with significance testing)
  - Maybe another/several other metrics (typically Meteor, TER, ChrF)
  - Maybe significance testing of metric improvements
  - Maybe some human judgments
    - Ranking of systems
    - Targeted analysis of specific phenomenon
- → Be careful when you argue about MT quality!

# Outlook

- Thursday:
  - Assignment 1: MT evaluation
  - Work in pairs
  - 11-12 oral examination
- Next week:
  - SMT, 2 lectures
  - Assignment 2: Moses and Phrase-based SMT

# Reminder

- Course registration:
  - send me your choice this week, 5LN711 (7.5 hp) or 5LN718 (5hp)
  - please spread this information to your classmates who might missed
- Snowy cluster registration:
  - remember to register an account to access GPU machines
  - refer to Lecture 1 for details