

INF5820/INF9820

LANGUAGE TECHNOLOGICAL APPLICATIONS

Jan Tore Lønning, Lecture 2, 31 Aug. 2016

jtl@ifi.uio.no

Machine Translation, lecture 2

2

- **Why is (machine) translation hard?**
 - Typological differences
 - Translational differences
- Evaluation in MT
 - Human evaluation of MT Quality
 - Automatic evaluation in Language Technology
 - Word precision and recall
 - BLEU

Why (machine) translation is hard.

3

Why can't we just use a dictionary?

Because:

- Languages are constructed differently (typology)
- Translation is not one-to-one

Language typology: morphology

4

□ Number of morphemes per word

- Isolating: 1,
 - Chinese, Vietnamese
- Synthetic: >1
- Polysynthetic: >>1

□ Morphemfusion:

- Agglutinative
 - putting morphemes after each other
 - Japanese, Turkish, Finnish, Sami
- Fusion
 - Russian

Washakotya'tawitsherahetkvhta'se
"He made the thing that one puts on
one's body ugly for her"
"He ruined her dress"

(Mohawk, polysynthetic, Src: Wikipedia)

(3.1) *uygarlaştıramadıklarımızdanmışsınızcasına*

uygar +laş +tır +ama +dık +lar +ımız +dan +mış +sınız +casına
civilized +BEC +CAUS +NABL +PART +PL +P1PL +ABL +PAST +2PL +AsIf

“(behaving) as if you are among those whom we could not civilize”

Turkish, agglutinative, polysynthetic J&M, Ch. 3

Language typology: Syntax

5

- Word order:
 - ▣ Subject-Verb-Object, SVO
 - ▣ SOV
 - ▣ VSO
- Prepositions vs postpositions
- Modifiers before or after:
 - ▣ Red wine vs. vin rouge
- Verb-framed vs. satellite-framed
 - ▣ Marking of direction
 - ▣ Marking of manner

Jorge swam across the river.
Jorge cruzó a nado el río.

Language typology: Markers

6

- One language may contain a marker which is lacking – or very different – in another language:
 - ▣ Tense
 - ▣ Aspect:
 - *She smiles* vs *she is smiling*
 - ▣ Case
 - ▣ Definiteness

Translational discrepancies

7

- Translation is not only about typological differences
- Even between typologically similar languages, the translation is not always one-to-one

A red oval with a black outline, containing the word "Ambiguity!" in white, bold, sans-serif font. The oval is centered horizontally and vertically on the slide.

Ambiguity!

Lexical ambiguities in SL

8

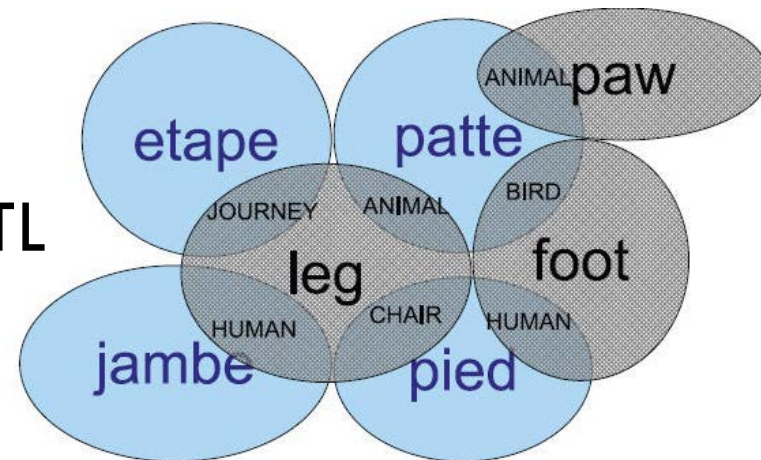
Word form	Norw: "dekket"			
POS	Noun		Verb	Adjective
Base form	"dekk"		"dekke"	
Homonymy	"dekk på båt"	"dekk på bil"		
Polysemy				
Gloss	"deck"	"tire"		

More examples		
	Norw	English
Verb/noun	løp, løper, bygg, bygget	fish, run, runs, ring
Homonymy	bygg (Noun), ball	bank, ball, bass
Polysemy	hode	head, bass (music)

Lexical choice in transfer

9

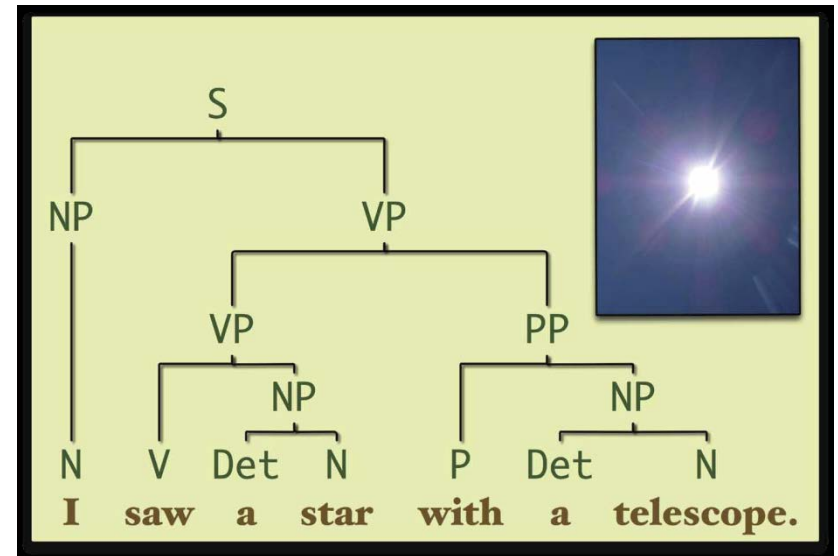
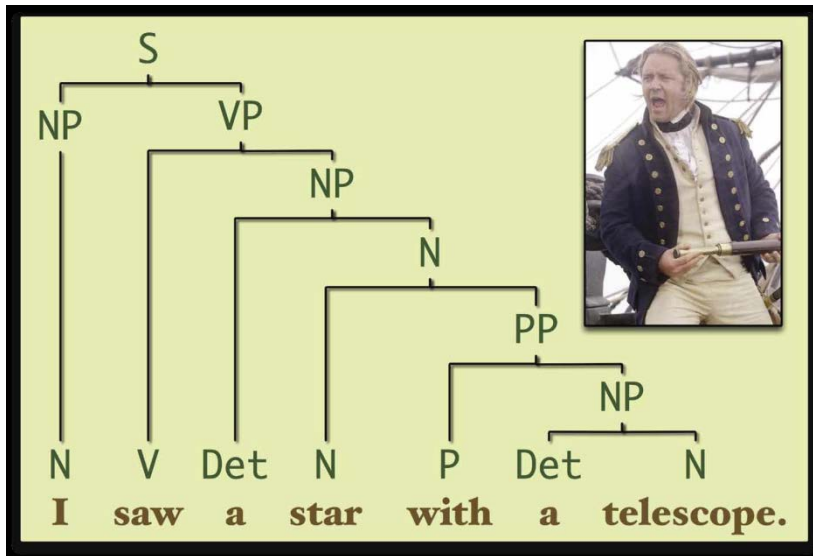
- The TL may make more distinctions than SL
 - ▣ No: *tak*, Eng: *ceiling/roof*
 - ▣ Eng: *grandmother*,
No: *farmor/mormor*
- Context dependent choice in TL
 - ▣ Strong tea, powerful government
 - ▣ *Dekke på bordet* → *set the table*
 - ▣ *Dekke bordet* → *set/cover the table*
- Languages may draw different distinctions
 - ▣ *Morgen* – *morning*, *legg* – *leg*



Syntactic ambiguities in SL

10

□ Global ambiguities



□ Local ambiguities:

- De kontrollerte bilene → They controlled the cars
- De kontrollerte bilene er i orden → The controlled cars are OK

Structural mismatch

11

- Thematic divergence/argument switching
 - E: I like *Mary*.
 - S: *Mary* me gusta.
- Head switching:
 - E: *Kim* likes to swim.
 - G: *Kim* schwimmt gern.
- More divergence:
 - N: *Han* heter *Paul*.
 - E: *His* name is *Paul*.
 - F: *Il* s'appell *Paul*.
- Idiomatic expressions



Beyond sentence meaning

12

- Tracking the referent,
No: **den/det** **han/hun**
- Metaphors, idioms

- Change,
- Rhime, rythm
- Deliberate ambiguity, humor
- ...

Machine Translation, lecture 2

13

- Why is (machine) translation hard?
 - ▣ Typological differences
 - ▣ Translational differences
- **Evaluation in MT**
 - ▣ **Human evaluation of MT Quality**
 - ▣ Automatic evaluation in Language Technology
 - ▣ Word precision and recall
 - ▣ BLEU

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.

NIST evaluation task 2001, from Koehn: SMT

Translation quality – Human eval.

15

- Given output of MT system + either
 1. Source text + reference translation (bilingual evaluator)
 2. Source text only (bilingual evaluator)
 3. Reference translation only (monolingual evaluator)
 4. Nothing (output only) (only fluency)
- Rate the translations (one sentence a time)
- Across several dimensions, typically
 - ▣ Adequacy: Does the output convey the same as the original/reference translation?
 - ▣ Fluency: Is this good target language?
 - ▣ and maybe several other dimensions

Judge Sentence

You have already judged 14 of 3064 sentences, taking 86.4 seconds per sentence.

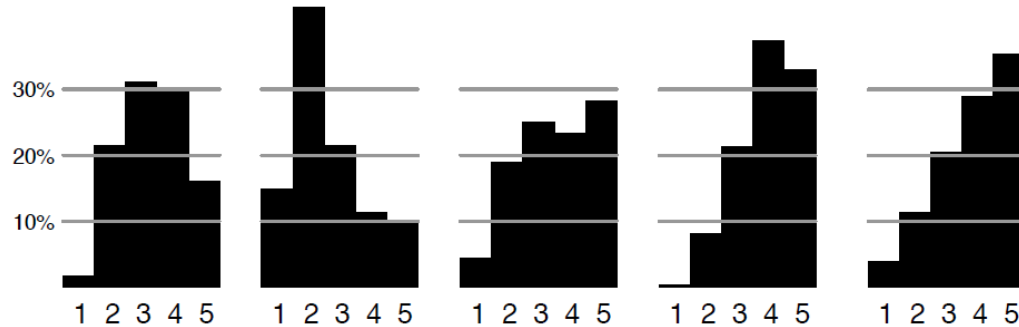
Source: les deux pays constituent plutôt un laboratoire nécessaire au fonctionnement interne de l'ue .

Reference: rather , the two countries form a laboratory needed for the internal working of the eu .

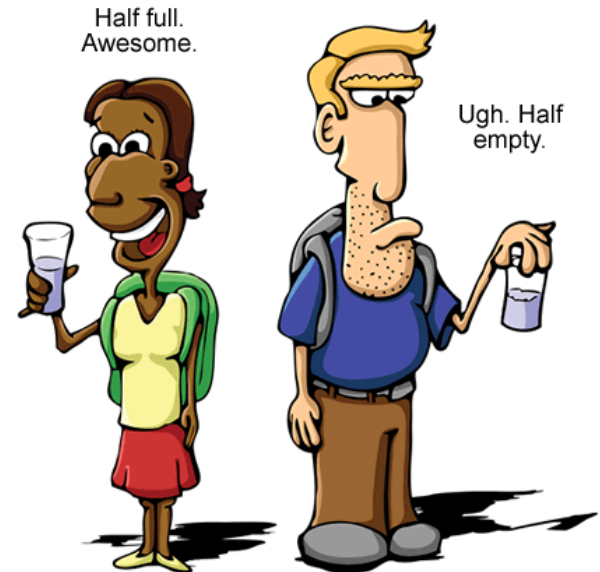
Translation	Adequacy	Fluency
both countries are rather a necessary laboratory the internal operation of the eu .	<input type="radio"/> 1 <input type="radio"/> 2 <input type="radio"/> 3 <input type="radio"/> 4 <input checked="" type="radio"/> 5	<input type="radio"/> 1 <input type="radio"/> 2 <input type="radio"/> 3 <input type="radio"/> 4 <input checked="" type="radio"/> 5
both countries are a necessary laboratory at internal functioning of the eu .	<input type="radio"/> 1 <input type="radio"/> 2 <input checked="" type="radio"/> 3 <input type="radio"/> 4 <input type="radio"/> 5	<input type="radio"/> 1 <input type="radio"/> 2 <input checked="" type="radio"/> 3 <input type="radio"/> 4 <input type="radio"/> 5
the two countries are rather a laboratory necessary for the internal workings of the eu .	<input type="radio"/> 1 <input type="radio"/> 2 <input type="radio"/> 3 <input checked="" type="radio"/> 4 <input type="radio"/> 5	<input type="radio"/> 1 <input type="radio"/> 2 <input type="radio"/> 3 <input checked="" type="radio"/> 4 <input type="radio"/> 5
the two countries are rather a laboratory for the internal workings of the eu .	<input type="radio"/> 1 <input type="radio"/> 2 <input checked="" type="radio"/> 3 <input type="radio"/> 4 <input type="radio"/> 5	<input type="radio"/> 1 <input type="radio"/> 2 <input type="radio"/> 3 <input type="radio"/> 4 <input checked="" type="radio"/> 5
the two countries are rather a necessary laboratory internal workings of the eu .	<input type="radio"/> 1 <input type="radio"/> 2 <input checked="" type="radio"/> 3 <input type="radio"/> 4 <input type="radio"/> 5	<input type="radio"/> 1 <input type="radio"/> 2 <input checked="" type="radio"/> 3 <input type="radio"/> 4 <input type="radio"/> 5
Annotator: Philipp Koehn Task: WMT06 French-English	<input type="button" value="Annotate"/>	
Instructions	5= All Meaning 4= Most Meaning 3= Much Meaning 2= Little Meaning 1= None	5= Flawless English 4= Good English 3= Non-native English 2= Disfluent English 1= Incomprehensible

Challenges in human TQ eval.

17



- What's in a number?
 - ▣ People use the scales differently
 - ▣ Normalize?
- More reliable alternative:
 - ▣ Evaluate several systems at once
 - ▣ Which translation is better?



Machine Translation, lecture 2

18

- Why is (machine) translation hard?
 - ▣ Typological differences
 - ▣ Translational differences
- Evaluation in MT
 - ▣ Human evaluation of MT Quality
 - ▣ **Automatic evaluation in Language Technology**
 - ▣ Word precision and recall
 - ▣ BLEU

Evaluation in language technology

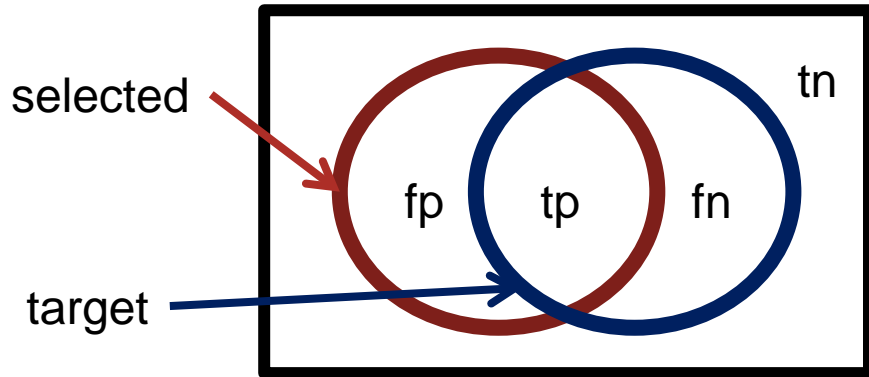
19

- Example 1: Tagging
 - Task: Assign part of speech tags to words in text
 - The/**DT** grand/**JJ** jury/**NN** commented/**VBD** ...
 - Gold standard: A hand-annotated corpus
 - Run your tagger on the gold standard
 - Compare the results with the gold standard
 - Accuracy: $\frac{\#(\text{correct tags})}{\#\text{words}}$
- Experimental set up:
 - Split an annotated corpus in two parts:
 - Training
 - Testing (=gold standard) not used in training



Common evaluation measures in LT

20



$$\square \text{ Recall} = \frac{tp}{tp + fn}$$

$$\square \text{ Precision} = \frac{tp}{tp + fp}$$

$$\square \text{ F-score} = \frac{1}{\alpha \frac{1}{P} + (1-\alpha) \frac{1}{R}}$$

$$\square F_1 = \frac{1}{0.5 \frac{1}{P} + (1-0.5) \frac{1}{R}} = \frac{2PR}{R+P}$$

		Actual (gold)	
		target	Not target
System perform	selected	tp: True positive	fp: False positive
	Not selected	fn: False negative	tn: True negative

Some remarks

21

- Precision and recall:
 - ▣ Comes from Information Retrieval (IR)
 - ▣ Have become (too?) popular in language technology
- Useful when:
 - ▣ There is more than one target/correct answer
 - ▣ The targets are known
 - ▣ The true negatives are many, uninteresting or unknown
 - ▣ The targets are not ranked
- Statistical significance tests are more easily available for accuracy than for P, R, F

Machine Translation, lecture 2

22

- Why is (machine) translation hard?
 - ▣ Typological differences
 - ▣ Translational differences
- Evaluation in MT
 - ▣ Human evaluation of MT Quality
 - ▣ Automatic evaluation in Language Technology
 - ▣ **Word precision and recall**
 - ▣ BLEU

Adapting P, R, F to MT-eval

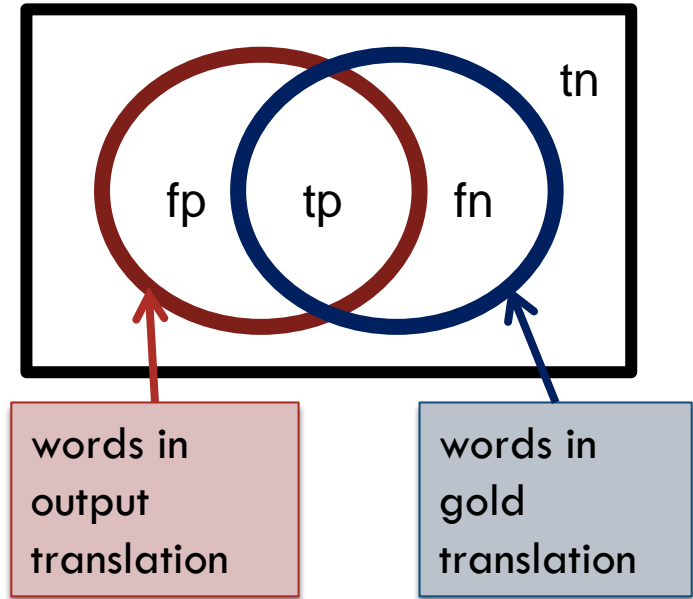
23

□ Precision = $\frac{\textit{correct}}{\textit{output.length}}$

□ Recall = $\frac{\textit{correct}}{\textit{ref.length}}$

□ $F_1 =$

$$\frac{2}{\frac{1}{R} + \frac{1}{P}} = \frac{2}{\frac{\textit{ref.length}}{\textit{correct}} + \frac{\textit{output.length}}{\textit{correct}}} = \frac{2\textit{correct}}{\textit{output.length} + \textit{ref.length}}$$



Precision and Recall of Words

SYSTEM A: Israeli officials responsibility of airport safety
REFERENCE: Israeli officials are responsible for airport security

- Precision

$$\frac{\text{correct}}{\text{output-length}} = \frac{3}{6} = 50\%$$

- Recall

$$\frac{\text{correct}}{\text{reference-length}} = \frac{3}{7} = 43\%$$

- F-measure

$$\frac{\text{precision} \times \text{recall}}{(\text{precision} + \text{recall})/2} = \frac{.5 \times .43}{(.5 + .43)/2} = 46\%$$

Precision and Recall



Metric	System A	System B
precision	50%	100%
recall	43%	100%
f-measure	46%	100%

$\frac{6}{7} \approx 0.86$

$\frac{12}{13} \approx 0.92$

flaw: no penalty for reordering

Position-independent error rate

26

- Similar measure to (word) recall+precision
- Reports mistakes – not correctness
- We skip the details - formula

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

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

Levenshtein distance used in

- spell-checking
- OCR
- Translation memory

Example

		Israeli	officials	responsibility	of	airport	safety
	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

		airport	security	Israeli	officials	are	responsible
	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%

Machine Translation, lecture 2

29

- Why is (machine) translation hard?
 - ▣ Typological differences
 - ▣ Translational differences
- Evaluation in MT
 - ▣ Human evaluation of MT Quality
 - ▣ Automatic evaluation in Language Technology
 - ▣ Word precision and recall
 - ▣ **BLEU**

BLEU

30

- A Bilingual Evaluation Understudy Score
- Main ideas:
 - Use several reference translations
 - Count precision of n-grams:
 - For each n-gram in output:
does it occur in at least one reference?
 - Don't count recall but use a penalty for brevity
 - Why not recall?

BLEU

31

$$P_n = \frac{\sum_{C \in \{Candidates\}} \sum_{n-gram \in C} Count_{clip}(n-gram, C, C.refs)}{\sum_{C \in \{Candidates\}} \sum_{n-gram \in C} Count(n-gram, C)}$$

- Candidates:
 - ▣ the set of sentences output by trans. system
- Count(n -gram, C):
 - ▣ the number of times n -gram occurs in C
- Count_{clip}(n -gram, C , $C.refs$):
 - ▣ the number of times the n .gram occurs in both
 - C and
 - the reference translation for the same sentence where n .gram occurs most frequent

- **Technicality:**
 - ▣ If the same n-gram has several occurrences in a candidate translation sentence, it should not be counted more times than the number of occurrences in the reference sentence with the largest number of occurrences of the same n-gram.

Example, p_3

33

- Hyp, C:
 - ▣ One of the girls gave one of the boys one of the boys.
- C-Refs:
 - ▣ A girl gave a boy one of the toy cars
 - ▣ One of the girls gave a boy one of the cars.

#

Example, p_3

34

- Hyp, C:
 - One of the girls gave one of the boys one of the boys.
- C-Refs:
 - A girl gave a boy one of the toy cars
 - One of the girls gave a boy one of the cars.
- $\text{Count_clip}(\text{one of the}, C, C\text{-refs})=2$

one of the	of the girls	the girls gave	girls gave one
2 (3)	1	1	0 (1)

gave one of	of the boys	the boys one	boys one of
0 (1)	0 (2)	0 (1)	0 (1)

- $P_3 = 4/11$

BLEU

35

- How to combine the n-gram precisions?

$$p_1 \times p_2 \times \cdots \times p_n = \prod_{i=1}^n p_i$$

- Remember

$$\ln\left(\prod_{i=1}^n p_i\right) = \ln(p_1 \times p_2 \times \cdots \times p_n) = \ln(p_1) + \ln(p_2) + \cdots + \ln(p_n) = \sum_{i=1}^n \ln p_i$$

- One can add weights, typically $a_i = 1/n$

$$\ln(p_1^{a_1} \times p_2^{a_2} \times \cdots \times p_n^{a_n}) = a_1 \ln(p_1) + a_2 \ln(p_2) + \cdots + a_n \ln(p_n)$$

- How long n-grams?

- ▣ Max 4-grams seems to work best

Brevity penalty

36

- c is the length of the candidates
- r is the length of the reference translations:
 - ▣ for each C choose the R most similar in length

- Penalty applies if $c < r$:

- ▣ $BP = 1$ if $c \geq r$
- ▣ $BP = e^{(1-r/c)}$ otherwise

$$c = \sum_{C \in \text{Candidates}} \text{length}(C)$$

$$r = \sum_{C \in \text{Candidates}} \text{length}(R.\text{sim}.C)$$

- $BLEU = BP \cdot \exp \sum_{i=1}^n w_n \ln p_i$

- $\ln BLEU = \min(1 - \frac{r}{c}, 0) + \sum_{i=1}^n w_n \ln p_i$

This is correct
Error in K:SMT