

Odds of Successful Transfer of Low-Level Concepts: A Metric for Speech-to-Speech Machine Translation

Gregory A. Sanders, Sébastien Bronsart,
Sherri Condon, Craig Schlenoff



MITRE

LREC 2008 Marrakech May 2008

NIST
National Institute of
Standards and Technology
U.S. Department of Commerce

The TRANSTAC Goal



Enable U.S. personnel speaking only English to communicate with civilian populations speaking only other languages



The TRANSTAC Program

Spoken Language Communication and Translation System for Tactical Use



- Bidirectional Speech-to-Speech Machine Translation
- Laptop or hand-held platform
- Free-form input, but in known domains
 - Medical
 - Civil Affairs
 - ✦ Examples: Sewer, Water, Electricity, Trash
 - Military operations
 - ✦ Examples: Training, Joint ops, Vehicle checkpoint
- Program sponsored by DARPA
- System performance evaluated by NIST and MITRE

What are the Low-level Concepts?



- We defined the low-level concepts to consist of the source-language content words
 - Open-class words:
 - ✦ Nouns, Verbs, Adjectives, Adverbs
 - Important prepositions and quantifiers
 - Entire verb construction (e.g., “will have been thrown”) is one concept
- Speakers choose what to give prominence via expression as a content word
- Number of such elements is determined by the speaker
 - Count is not open-ended
 - Count is not highly subjective
- Low-level concepts annotated in the source-language transcript
 - Annotated by a native speaker
 - If utterance is disfluent, count only the concepts that a fluent rendition would include

CTR, in reference annotation mode

CTR v1.3 (C:/Documents and Settings/Greg_2/Desktop/TT_July07_IA/CTR_IA_July07/491.ref)

File

Progress

0 question utterance(s) to validate yet
11 answer utterance(s) to validate yet

60%

Exchange List

- Scenario
 - Exchange #0
 - Question
 - Utterance #0 ✓
 - Utterance #1 ✓
 - Answer
 - Utterance #0 ✓
 - Exchange #1
 - Question
 - Utterance #0 ✓
 - Utterance #1 ✓
 - Answer
 - Utterance #0 ✓
 - Exchange #2
 - Question
 - Utterance #0 ✓
 - Utterance #1 ✓
 - Utterance #2 ✓
 - Answer
 - Utterance #0
 - Utterance #1

are there -- %AH what sort of problems
are there with it ?

- ⚠ what sort of(...)?
- ⚠ problems
- ⚠ with
- ⚠ it / (the waste water)

+ Add a concept

✖ Remove a concept

✖ Remove all concepts

⚠ Toggle already known

Next Not done Done

Scoring Successful Transfer



- Panel of bilingual judges who each score the MT output
 - Compare textual target-language MT output to annotated transcription of source-language utterance
 - Each low-level concept is scored:
 - ✦ Successfully transferred --- Correct
 - ✦ Deleted
 - ✦ Substituted
 - ✦ Inserted concepts are also identified by the judges
- Result stated as Odds of Successful Transfer of a low-level concept
Odds(*correct*) =
$$\frac{\text{NumCorrect}}{\text{Deletions} + \text{Substitutions} + \text{Insertions}}$$
- Progress across evaluations can be stated as an Odds Ratio

CTR in MT output scoring mode

The screenshot displays the CTR v1.3 application window. The title bar reads "CTR v1.3 (C:/Documents and Settings/Greg_2/Desktop/backup_2007_08_03_final/arabic/judge3.sys)". The interface is divided into several sections:

- Progress:** Shows "5 question utterance(s) to validate yet" and "10 answer utterance(s) to validate yet" with a progress bar at 98%.
- Exchange List:** A tree view on the left showing a list of exchanges. The current exchange is "Exchange #1", which contains a "Question" (Utterance #0) and an "Answer" (Utterance #0, #1, #2).
- Main Text Area:** Displays the source text "what kind of job would your company be able to do right now ?" and the target text "شئو نوع الشغل اللي إنت تروح تقدر تسوي هسا".
- Scoring Legend:** A list of words with their corresponding scores: "what kind(...)? <-> شئو نوع" (Correct), "job <-> الشغل" (Correct), "would be able to do / can do <-> تقدر تسوي" (Correct), "your <-> إنت" (Correct), "company" (Deleted), "right now <-> هسا" (Correct), and "تروح" (Inserted).
- Buttons:** Includes "Correct", "Substituted", "Inserted", "Deleted", and "Remove".
- Adequacy of translation:** Radio buttons for "Completely adequate", "Tending towards adequate" (selected), "Tending towards inadequate", and "Inadequate".
- Navigation:** "Next" button and status indicators for "Done" and "Not done yet".

Judgments of Semantic Adequacy



- We asked our bilingual judges to also give a single judgment of semantic adequacy for each utterance on a four-point scale
 - Completely adequate
 - Tending towards adequate
 - Tending towards inadequate
 - Inadequate
- Judges assigned this utterance-level score immediately after scoring the low-level concepts in the utterance
- We consider these judgments to be our *benchmark* score
 - We compare our other metrics to it

CTR in MT output scoring mode

The screenshot shows the CTR v1.3 application window. The title bar reads "CTR v1.3 (C:/Documents and Settings/Greg_2/Desktop/backup_2007_08_03_final/arabic/judge3.sys)". The interface is divided into several sections:

- Progress:** A green progress bar indicates 98% completion. Text above it says "5 question utterance(s) to validate yet" and "10 answer utterance(s) to validate yet".
- Exchange List:** A tree view on the left shows the structure of the data. It includes "Exchange #1" with "Question" and "Answer" sub-items, and "Exchange #0" which is currently selected and highlighted in yellow.
- Main Display Area:** The top part shows the English question: "what kind of job would your company be able to do right now ?". The bottom part shows the Arabic translation: "شئو نوع الشغل اللي إنت تروح تقدر تسوي هسا".
- Scoring Legend:** A list of words from the English sentence is shown with their corresponding Arabic equivalents and a color-coded score:
 - C (Correct): what kind(...)?, الشئو نوع
 - C (Correct): job, الشغل
 - C (Correct): would be able to do / can do, تقدر تسوي
 - C (Correct): your, إنت
 - D (Deleted): company
 - C (Correct): right now, هسا
 - I (Inserted): تروح
- Control Panel:** Buttons for "Correct", "Substituted", "Inserted", and "Deleted" are present. A "Remove" button is also available. Below these are radio buttons for "Adequacy of translation": "Completely adequate", "Tending towards adequate" (selected), "Tending towards inadequate", and "Inadequate".
- Navigation:** A "Next" button is at the bottom left. A status bar at the bottom right shows "Done" (checked) and "Not done yet" (unchecked).

Training the Judges for Semantic Adequacy



- We explained the intended use and purpose of the system
 - Asked judges to assign scores that reflect how well the translations would serve that purpose
- We gave the judges a substantial set of exemplars for each of the four possible scores
 - The exemplars were taken from a previous eval, and were utterances on which the (different) set of judges from that eval had a high level of agreement
- We had the judges discuss several example translations as a group
 - Made sure each judge was offering appropriate reasons for their choice of score --- made sure they understand the task
- For Arabic, we told the judges to favor translations into Iraqi dialect, not the standard written language (MSA or Fus'ha)

Converting Odds to Probability of Correct Transfer



$\text{NumCorrect} / (\text{Deletions} + \text{Substitutions} + \text{Insertions})$

- Because we count insertions as errors, our odds calculation is not quite canonical $P(\text{correct}) / (1 - P(\text{correct}))$
- As $P(\text{correct})$ approaches 1.0, $\text{Odds}(\text{correct})$ approaches ∞
 - Typical automated MT metrics behave mathematically more like $P(\text{correct})$ than like $\text{Odds}(\text{correct})$
 - Correlation with automated MT metrics calls for a statistic that behaves like $P(\text{correct})$, but with insertions taken into account
- Adjusted Probability Correct

$$\text{AdjP}(\text{correct}) = 1 - (1 / (\text{Odds}(\text{correct}) + 1))$$

Other metrics are also important



- Concepts vary in importance -- some concepts are crucial
 - Utt: There are new IEDs along the road from here to Fallujah.
 - MT: There are no IEDs along the road from here to Fallujah.
- Low-level concept transfer metric gives all concepts equal weight
 - Utterance-level human judgments of semantic adequacy weigh the crucial errors appropriately
- Low-level concept transfer metric does not consider fluency
 - Even badly fractured syntax may be given a pass
 - Many automated MT metrics (e.g., BLEU, METEOR) do effectively consider fluency, as do utterance-level human judgments

Other Metrics We Calculated



- Source-language ASR was scored with Word Error Rate
- MT was scored with several commonly used metrics
 - BLEU
 - METEOR
 - TER
 - HTER --- only completed for translations into English

Discussion of Results



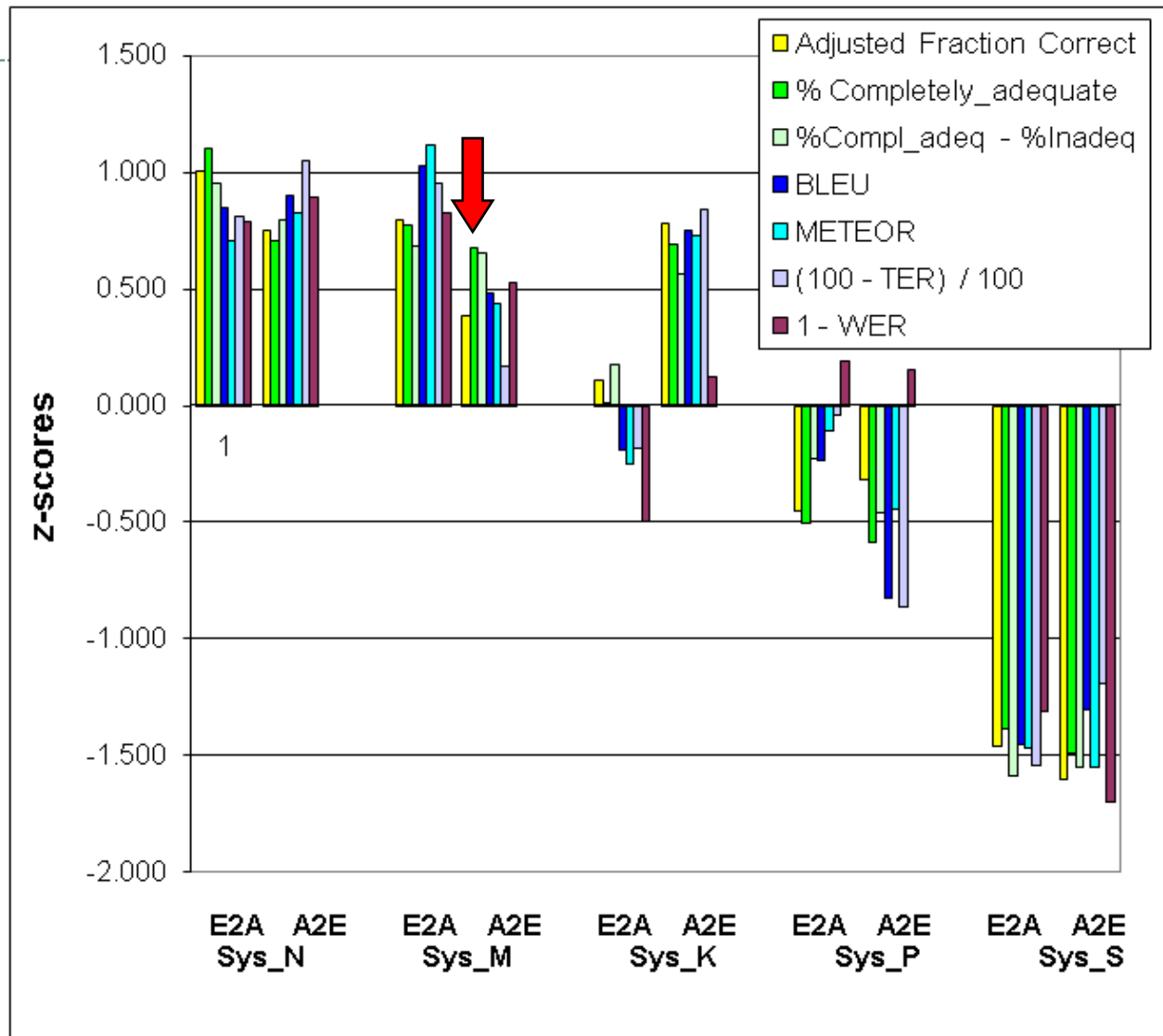
- Between January 2007 and July 2007 systems made large improvements in this metric
 - For English to Iraqi Arabic the median value over the five systems improved to 4.32 from 1.55 (an odds ratio of 2.79)
 - For Iraqi Arabic to English the median value improved to 3.15 from 2.46 (an odds ratio of 1.28)
- Scores on AdjP(*correct*) strongly correlated to the utterance-level judgments of semantic adequacy
 - Pooling all data for each system, Pearson correlation over the five systems
 - ✦ $R = 0.997$ for English to Iraqi Arabic
 - ✦ $R = 0.978$ for Iraqi Arabic to English
 - ✦ $R = 0.997$ for English to SurpriseLanguage
 - ✦ $R = 0.960$ for SurpriseLanguage to English

Comparing all the Metrics

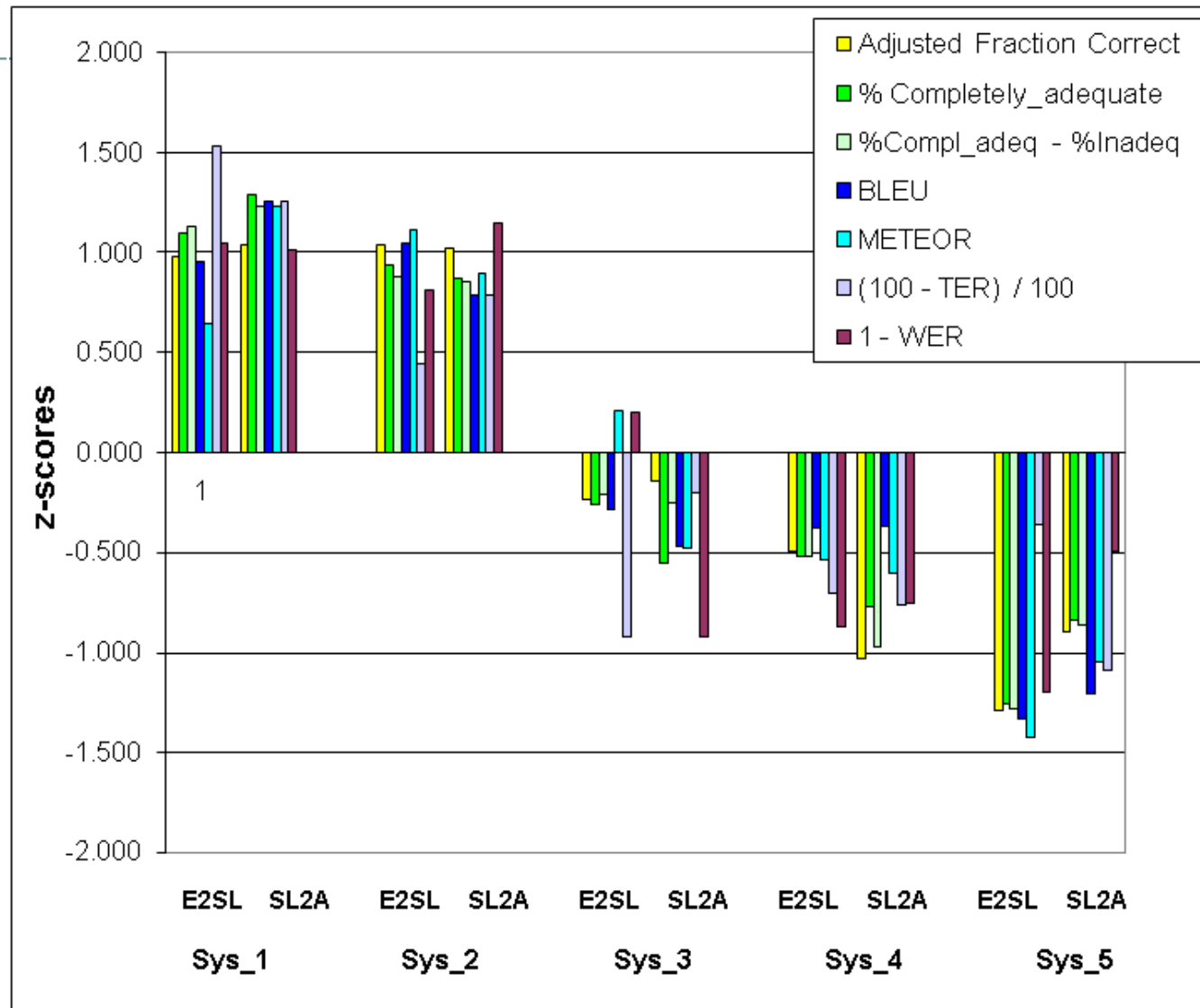


- For each language pair, separately, and each direction (to/from English) separately, we calculated mean and standard deviation, then converted all values to standard normal z statistics
- Result shown in the following synoptic overview graphs

Synoptic Overview for Arabic



Synoptic Overview for SurpriseLanguage



HTER for Iraqi Arabic to English



- HTER based on a human post-editing the MT output as necessary so that it has the correct meaning (fix the semantic errors)
 - HTER is a measure of the minimum number of edits necessary
- Key wrinkle in TER and HTER: a block move counts as one edit
 - Moving a string of any number of words by any distance
- Looking at HTER for each of the nine scenarios, for each of the four strongest systems (thus $4 \times 9 = 36$ data points)
 - Pearson correlation of HTER with *AdjP(correct)* is $R = 0.905$
 - Pearson correlation SemAdeq with *AdjP(correct)* is $R = -0.833$
- Omitting the hardest and easiest scenario to eliminate outlier effects (thus, $4 \times 7 = 28$ data points)
 - Pearson correlation of HTER with *AdjP(correct)* is $R = 0.849$
 - Pearson correlation SemAdeq with *AdjP(correct)* is $R = -0.790$

Inter-judge Agreement on Semantic Adequacy



- We had six judges for Arabic, and five for the surprise language
- Values of Cohen's kappa for pairwise inter-judge agreement, over the Arabic judges:
 - Exact match pairwise kappa range 0.178 to 0.435 (median 0.294)
 - ✦ Very low values -- not good
 - If we count the disagreements by just one level as being matches, then the pairwise kappa range is 0.508 to 0.805 (with median 0.611)
 - ✦ We regard this as an acceptable level of agreement
- For odds of successful transfer, there was fairly close agreement between the mean and median values over our set of judges
- Considering all this, we suggest that a reasonably large set of judges is necessary, as outlier judges are likely

Conclusions



- Odds of successful transfer of a low-level concept appears to be a relatively useful quantitative metric for information transfer
 - Strong correlation to human judgments of semantic adequacy
 - Strong correlation to the most common automated MT metrics, such as BLEU and METEOR
- The metric is labor-intensive
 - More useful for summative evaluation
- Training the judges carefully is important
 - Important to provide guidelines, with several examples of what counts as the same and what counts as different. Tricky tricky issues arise.
- Using a panel of **several** bilingual judges appears important
 - There were notably forgiving and harsh judges (outliers).
- Getting some judges to mark insertions is difficult; this can bias results.

For Further Info



- Over time, various TRANSTAC papers, presentations, guidelines documents, and so forth, will appear in the web pages for the NIST Speech Group

<http://www.nist.gov/speech>