BLEU (Papineni et al., 2002)



- Based on geometric mean of *n*-gram precision.
- pprox ratio of 1- to 4-grams of hypothesis confirmed by a ref. translation

SrcVom Glück der träumenden KameleConfirmedRefOn the happiness of dreaming camels1 2 3 4System A > The happiness of dreaming camels <</td>5 4 3 2System B > Dreaming of the luck that camels <</td>4 1 0 0n-grams confirmed: none, unigram, bigram, trigram, fourgram

E.g. Sys B produced 6 unigrams (4 confirmed), 4 bigrams (1 confirmed), . . .

$$\mathsf{BLEU} = \mathsf{BP} \cdot \exp\left(\frac{1}{4}\log\left(\frac{4}{6}\right) + \frac{1}{4}\log\left(\frac{1}{7}\right) + \frac{1}{4}\log\left(\frac{0}{6}\right) + \frac{1}{4}\log\left(\frac{0}{5}\right)\right)$$

BP is "brevity penalty"

 $^{1/4}$ is the "log-domain denominator" equivalent for $\sqrt[4]{\cdot}$ in geometric mean

BLEU: Avoiding Cheating



Confirmed counts "clipped" to avoid overgeneration.
"Brevity penalty" applied to avoid too short output:

$$\mathsf{BP} = \left\{ \begin{array}{ll} 1 & \text{if } c > r \\ e^{1 - r/c} & \text{if } c \leq r \end{array} \right.$$

Ref 1: The cat is on the mat .

Ref 2: There is a cat on the mat .

Candidate: The the the the the the the .

 \Rightarrow Clipping: only $\frac{3}{8}$ unigrams confirmed.

Candidate: The the .

 $\Rightarrow \frac{3}{3}$ unigrams confirmed but the output is too short. $\Rightarrow BP = e^{1-7/3} = 0.26$ strikes.

The candidate length c and "effective" ref. length r calculated over the whole test set.

References



Papineni, Salim Ward, Kishore Roukos, Todd Wei-Jing Zhu. 2002. and BLEU: Method for Automatic Evaluation of Machine Translation. In а ACL 2002, Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, pages 311-318, Philadelphia, Pennsylvania.