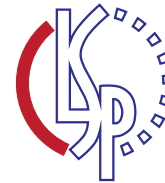


Simulated Multiple Reference Training (SMRT) Improves Low-Resource Machine Translation

Huda Khayrallah, Brian Thompson,
Matt Post & Philipp Koehn
huda@jhu.edu

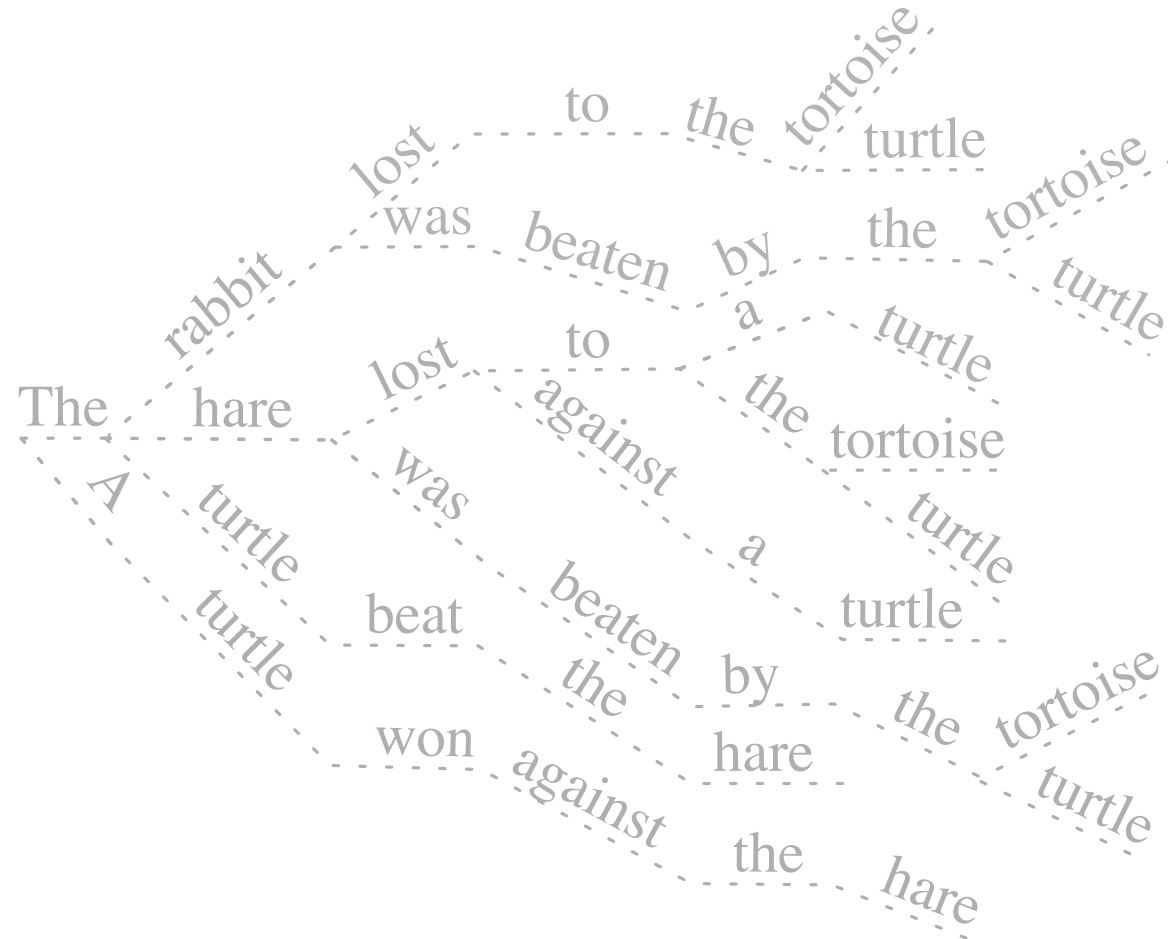


JOHNS HOPKINS
UNIVERSITY

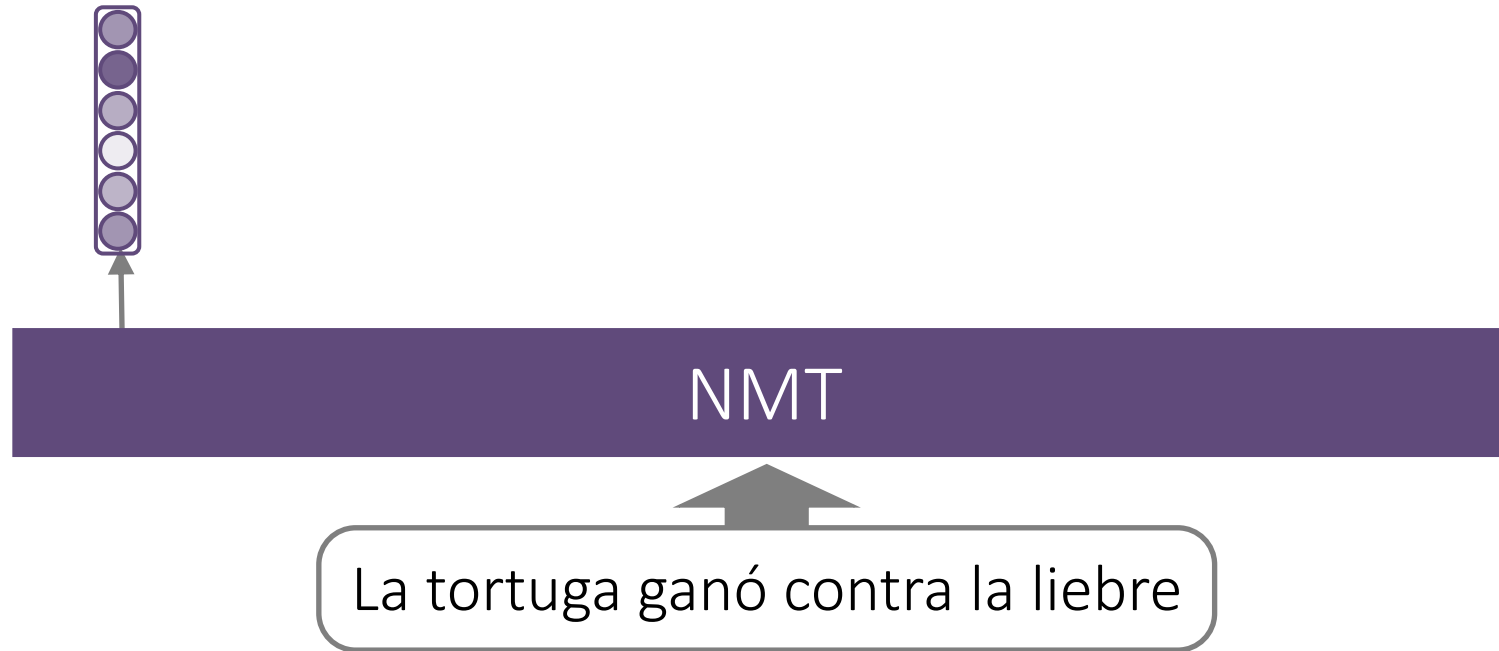


use target side *paraphrasing*
to overcome data sparsity
in *low-resource* settings

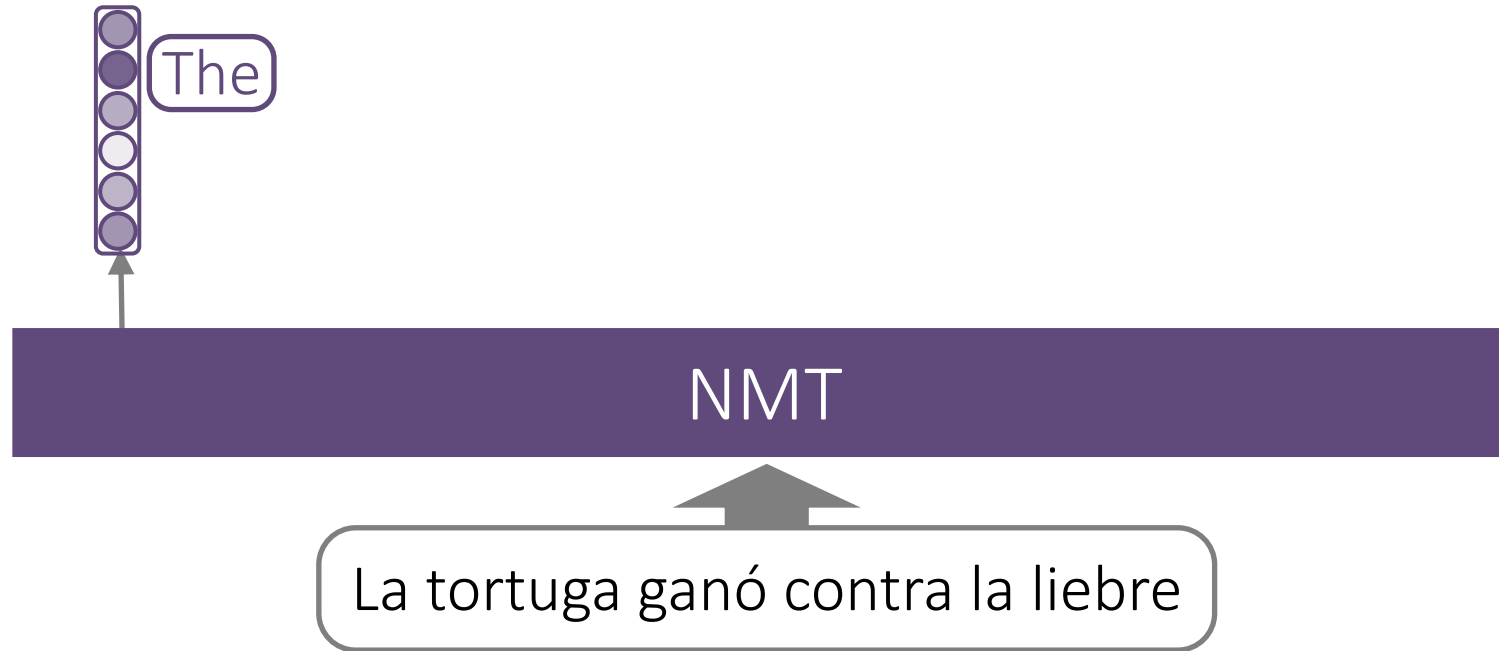
La tortuga ganó contra la liebre | The turtle beat the hare



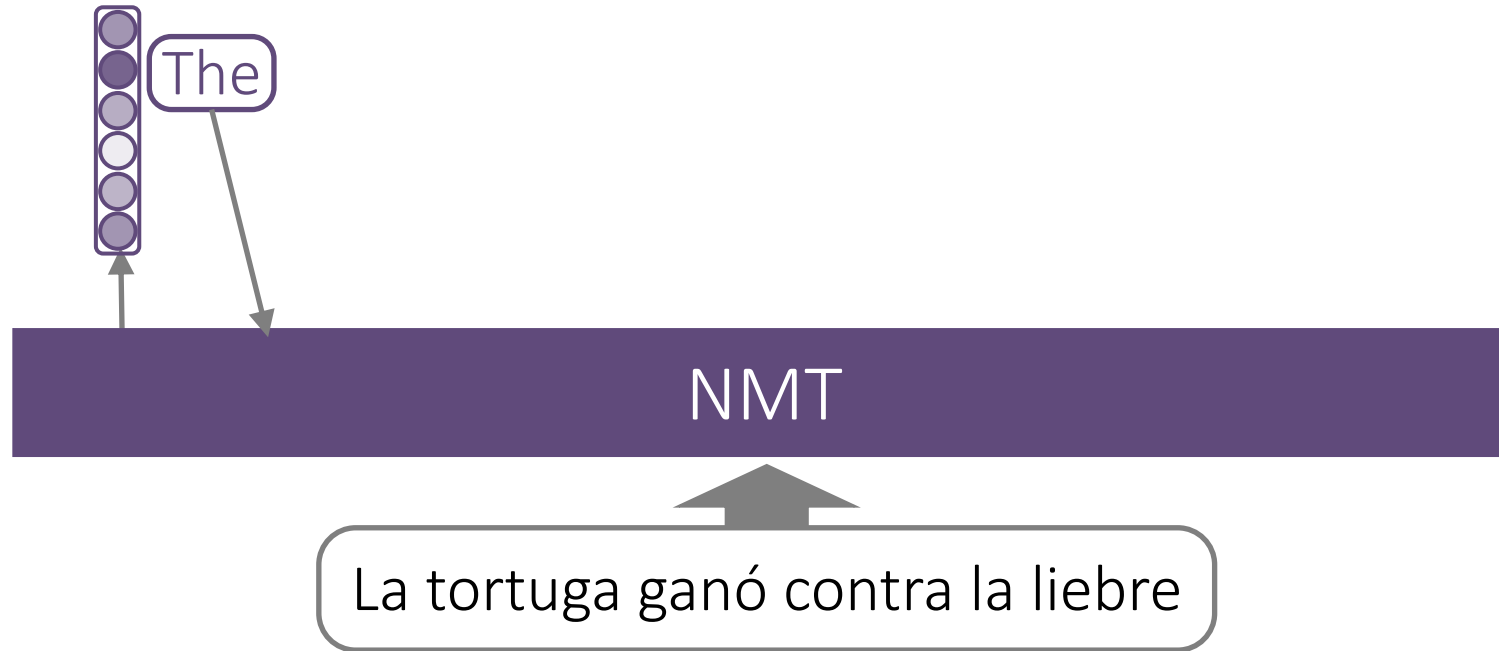
Machine Translation



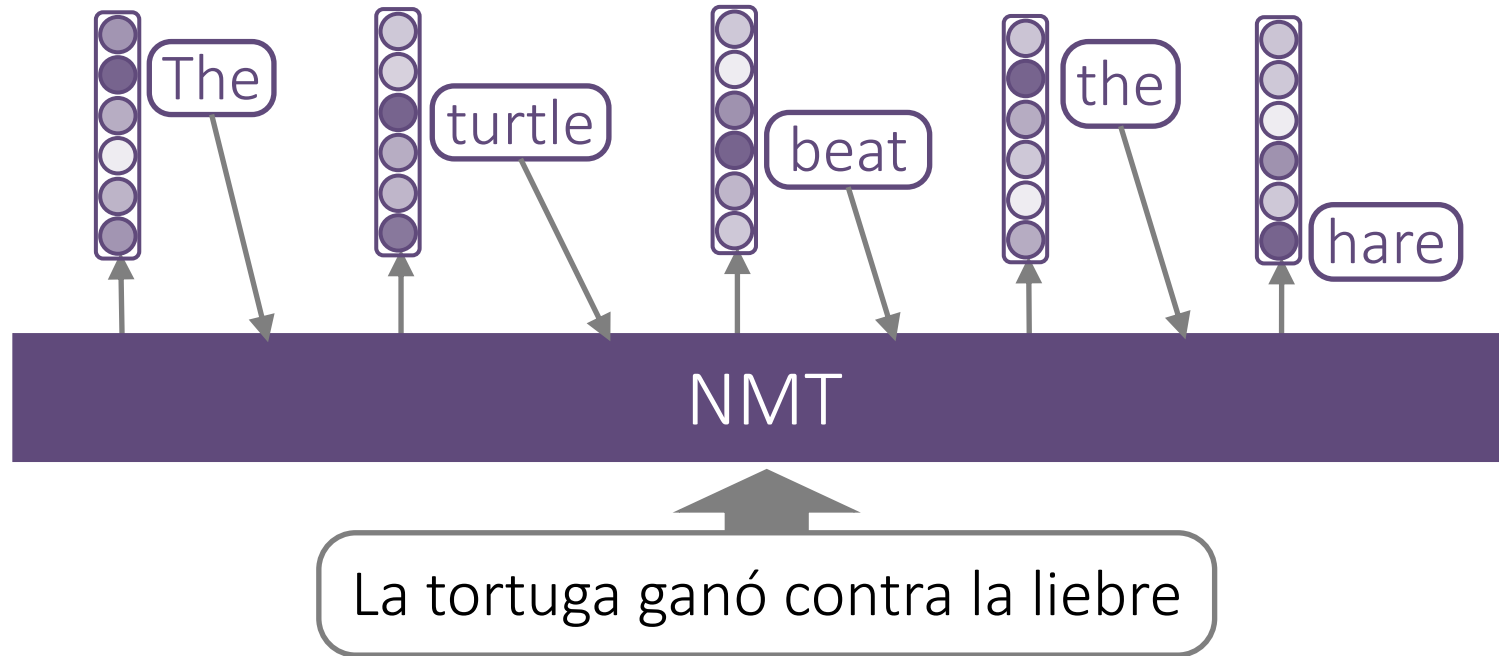
Machine Translation



Machine Translation

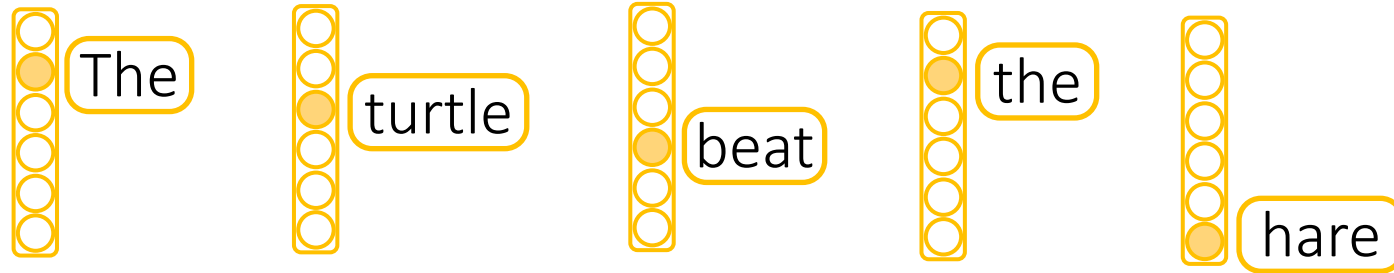


Machine Translation



La tortuga ganó contra la liebre | The turtle beat the hare

La tortuga ganó contra la liebre | The turtle beat the hare



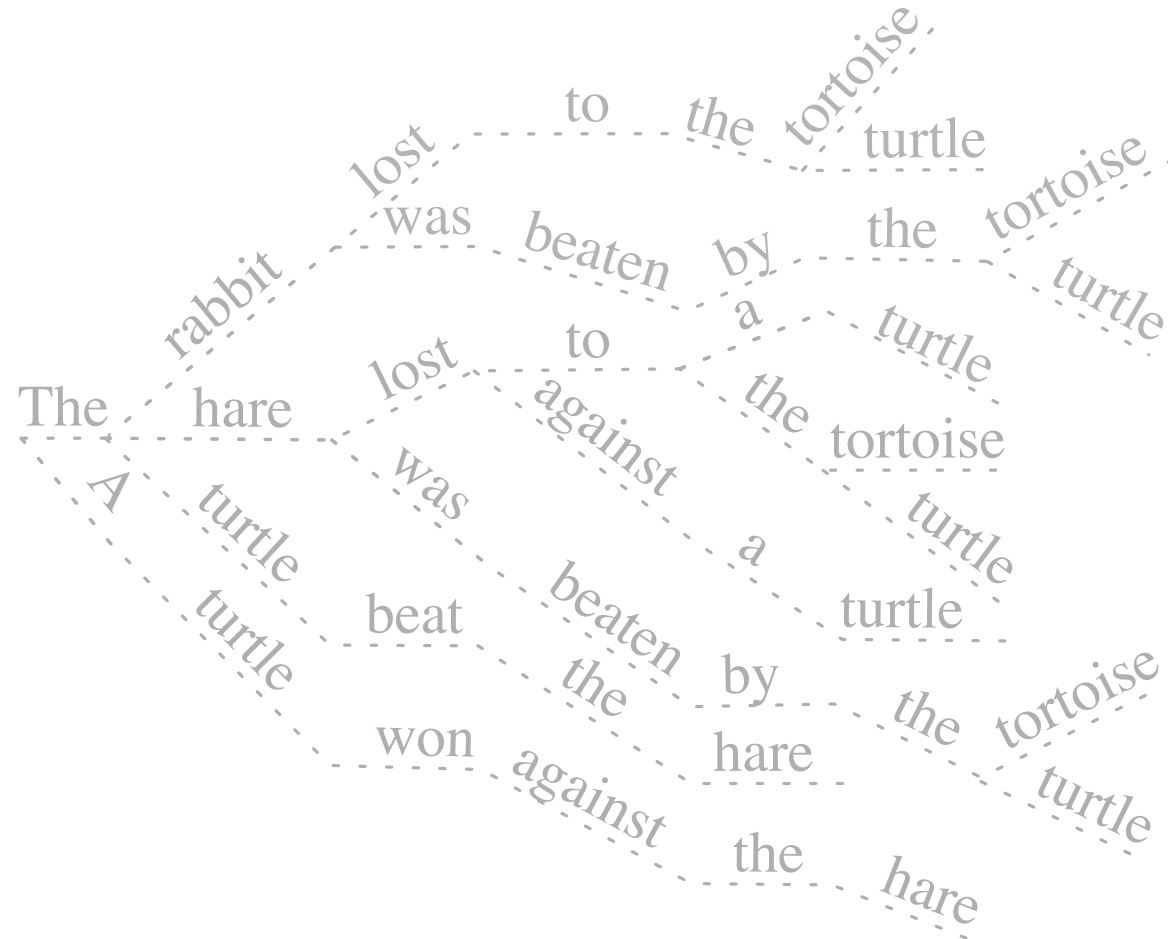
NLL Objective

$$-\sum_{v \in \mathcal{V}} \left[\underbrace{\mathbb{1}\{y_i = v\}}_{\text{Gold Target}} \times \log \underbrace{p_{\text{MT}}(y_i = v \mid x; y_{j < i})}_{\text{MT Model output}} \right]$$

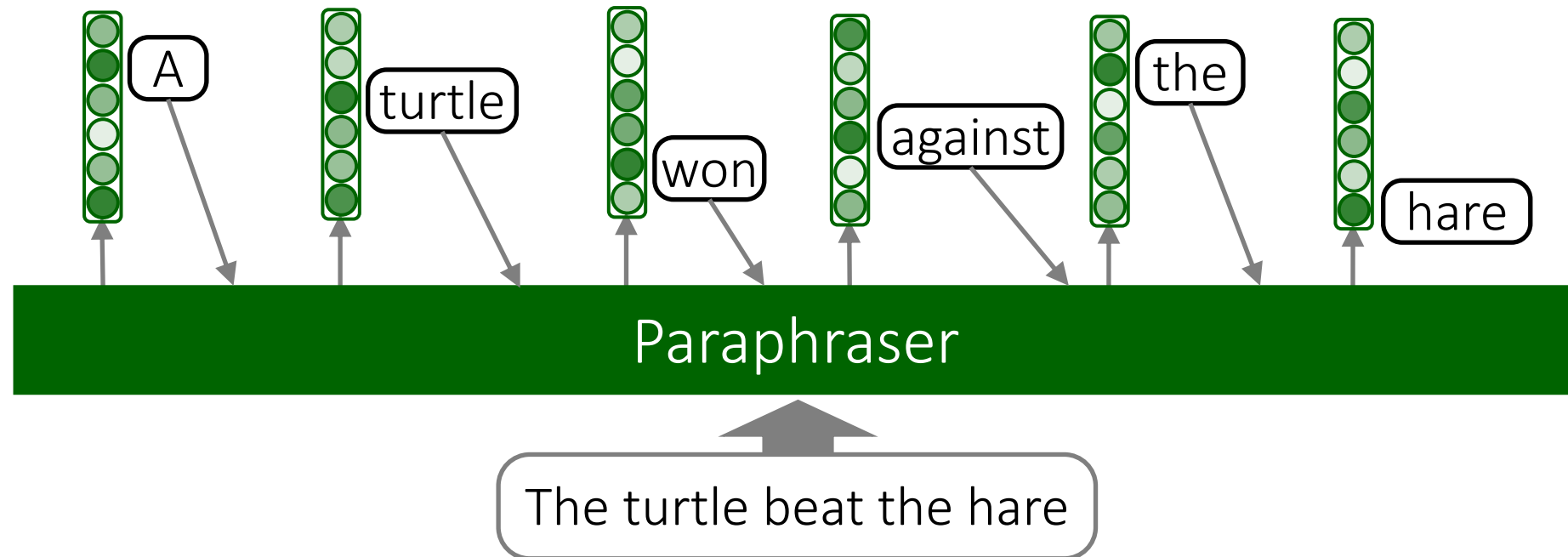
Cross Entropy( , )

Gold Target MT Model
output

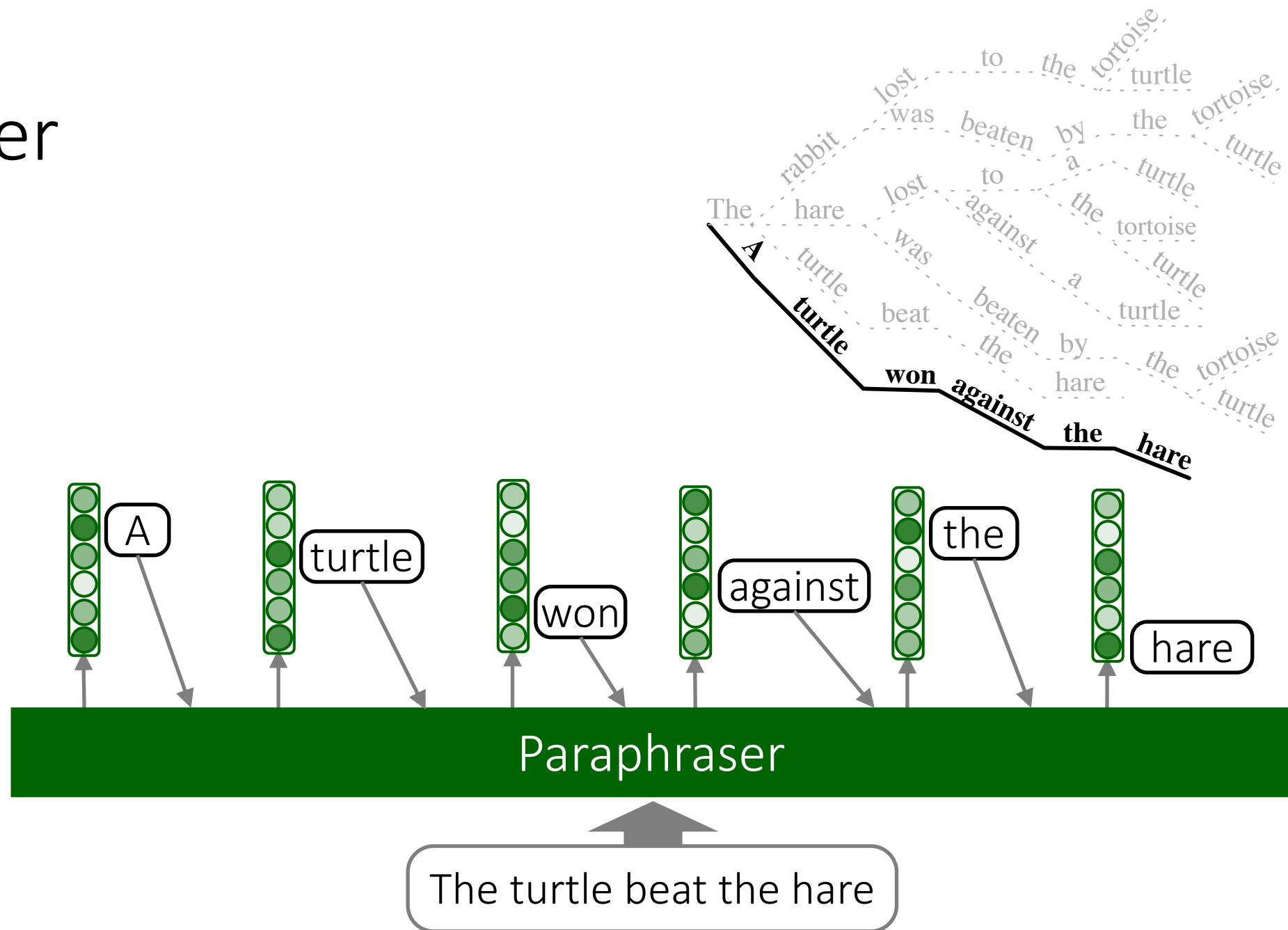
La tortuga ganó contra la liebre | The turtle beat the hare



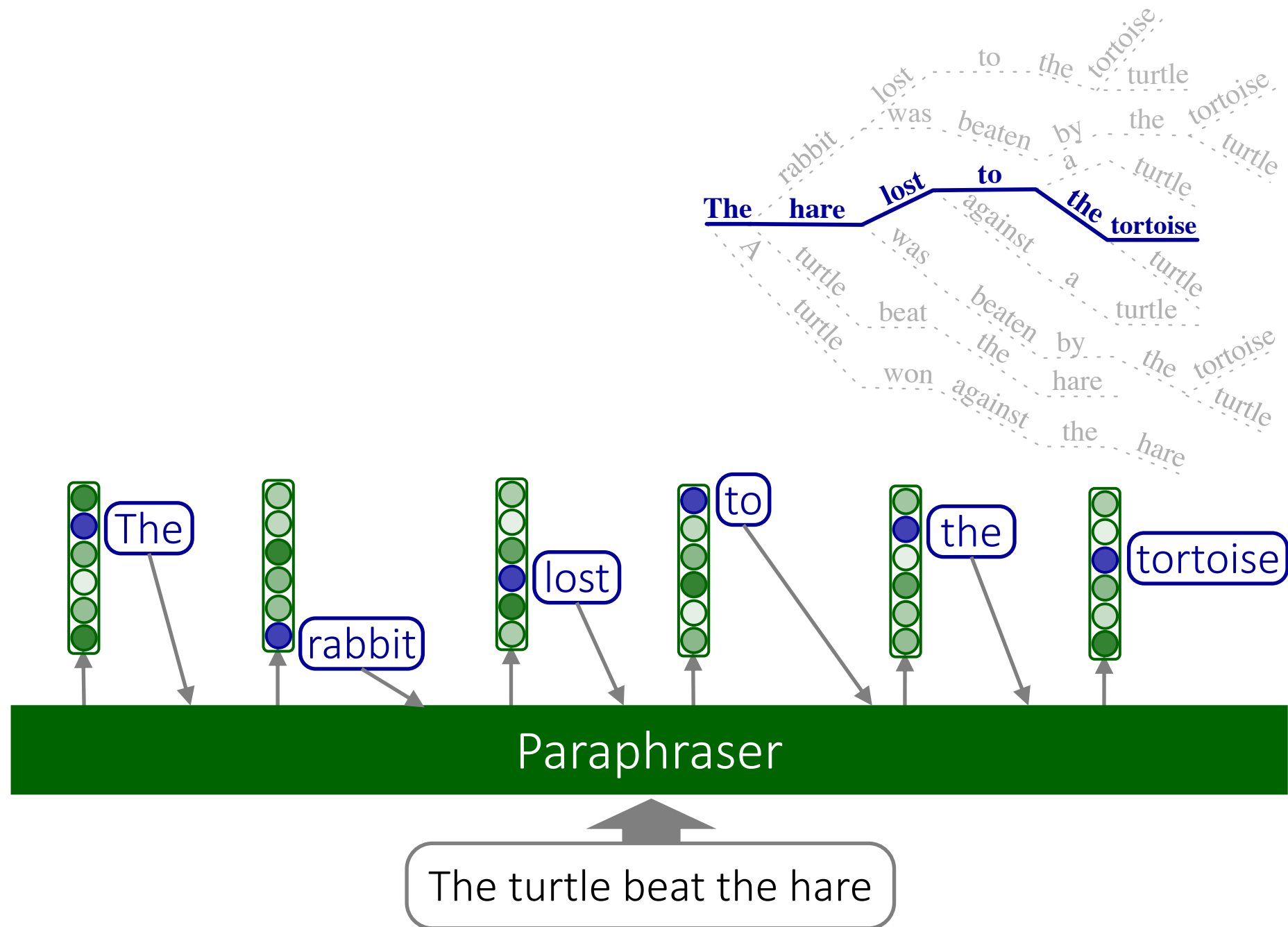
Paraphraser



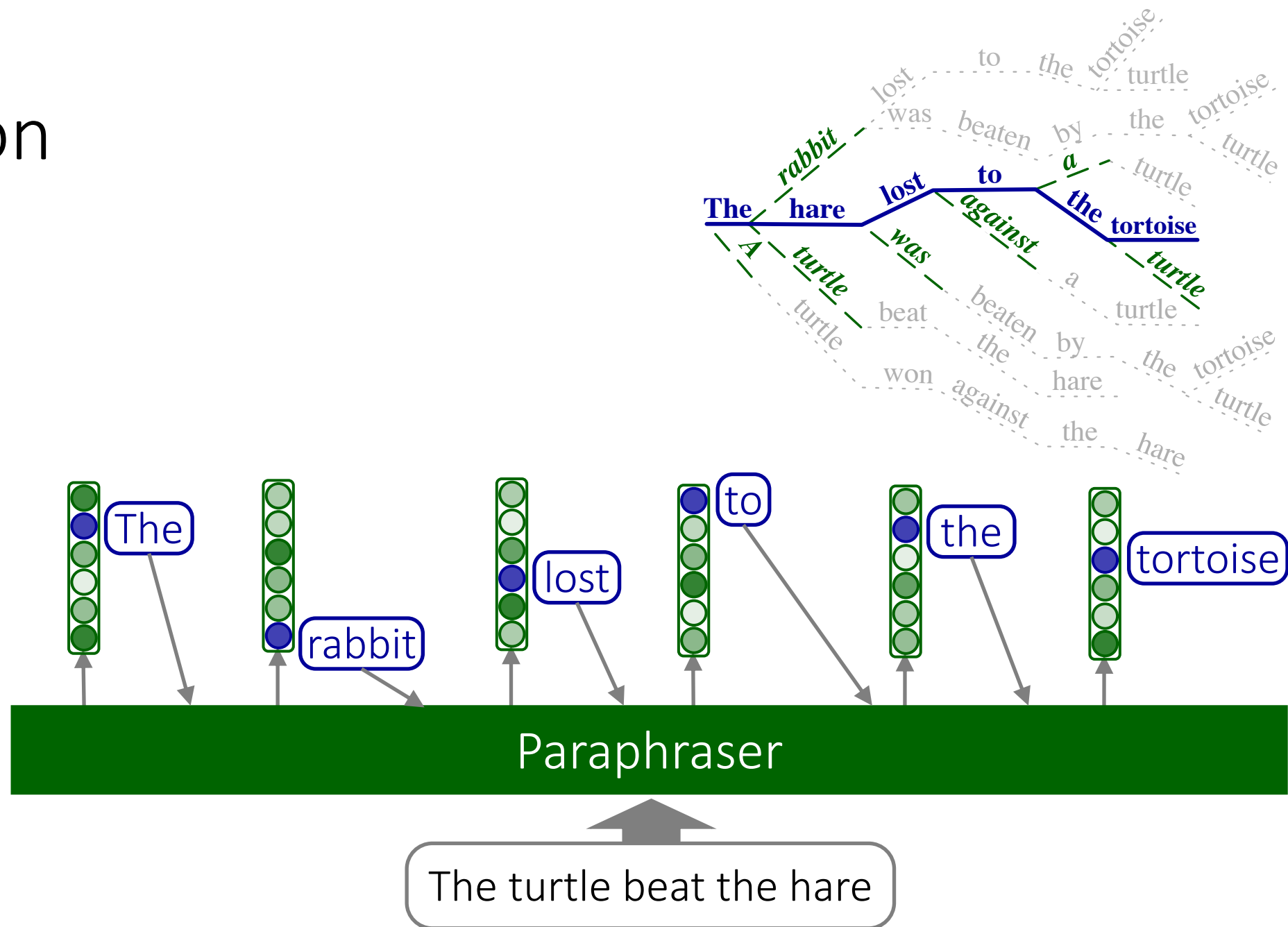
Paraphraser



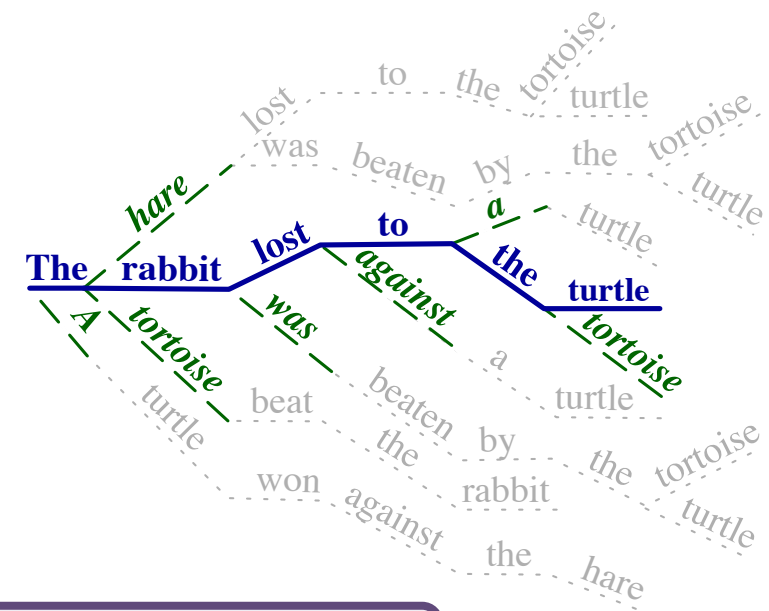
Sampling



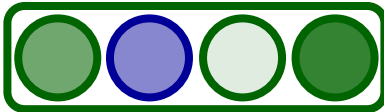

Distribution



SMRT Objective

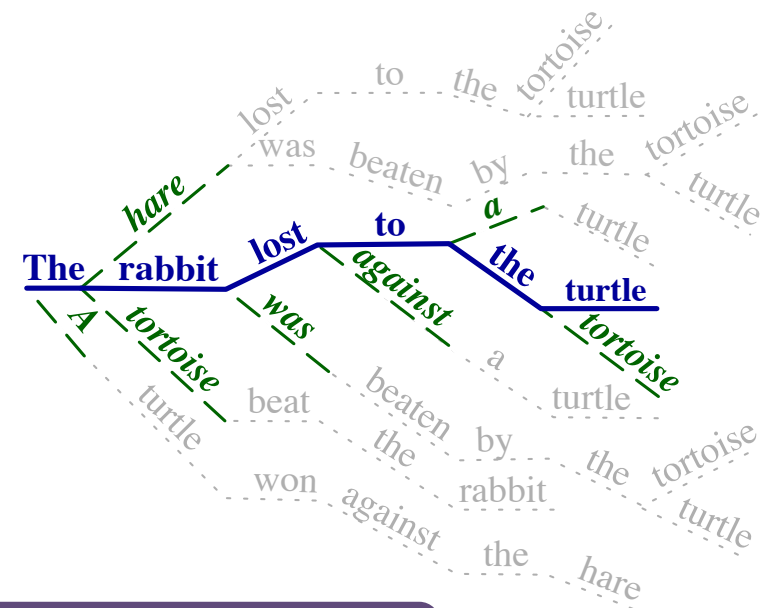


$$-\sum_{v \in \mathcal{V}} \left[\underbrace{p_{\text{para}}(y'_i = v \mid y; y'_{j < i})}_{\text{Paraphraser Output}} \times \log \underbrace{p_{\text{MT}}(y'_i = v \mid x; y'_{j < i})}_{\text{MT Model output}} \right]$$

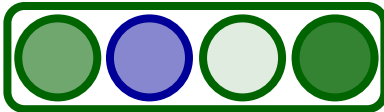

Cross Entropy( , )

Paraphraser Output MT Model output

SMRT Objective

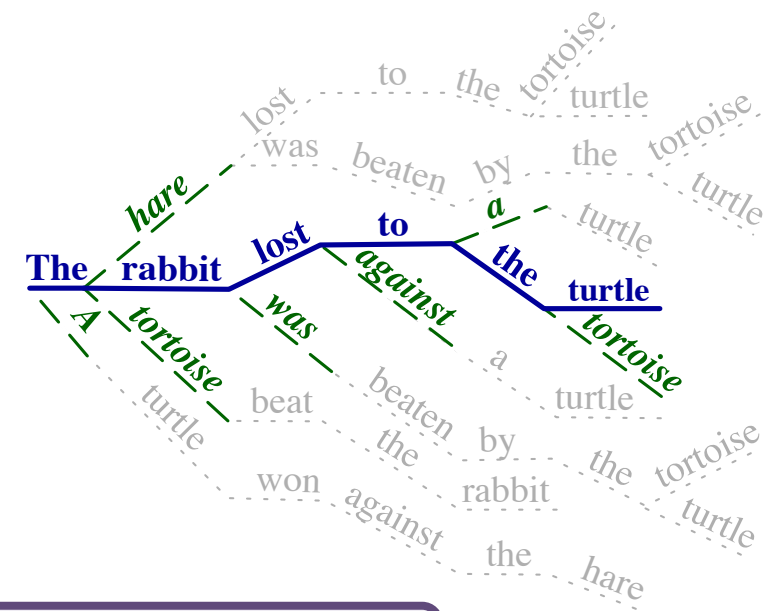


$$-\sum_{v \in \mathcal{V}} \left[\underbrace{p_{\text{para}}(y'_i = v \mid y; y'_{j < i})}_{\text{Paraphraser Output (teacher)}} \times \log \underbrace{p_{\text{MT}}(y'_i = v \mid x; y'_{j < i})}_{\text{MT Model output (student)}} \right]$$

Cross Entropy( , )

Paraphraser Output (teacher) MT Model Output (student)

SMRT Objective



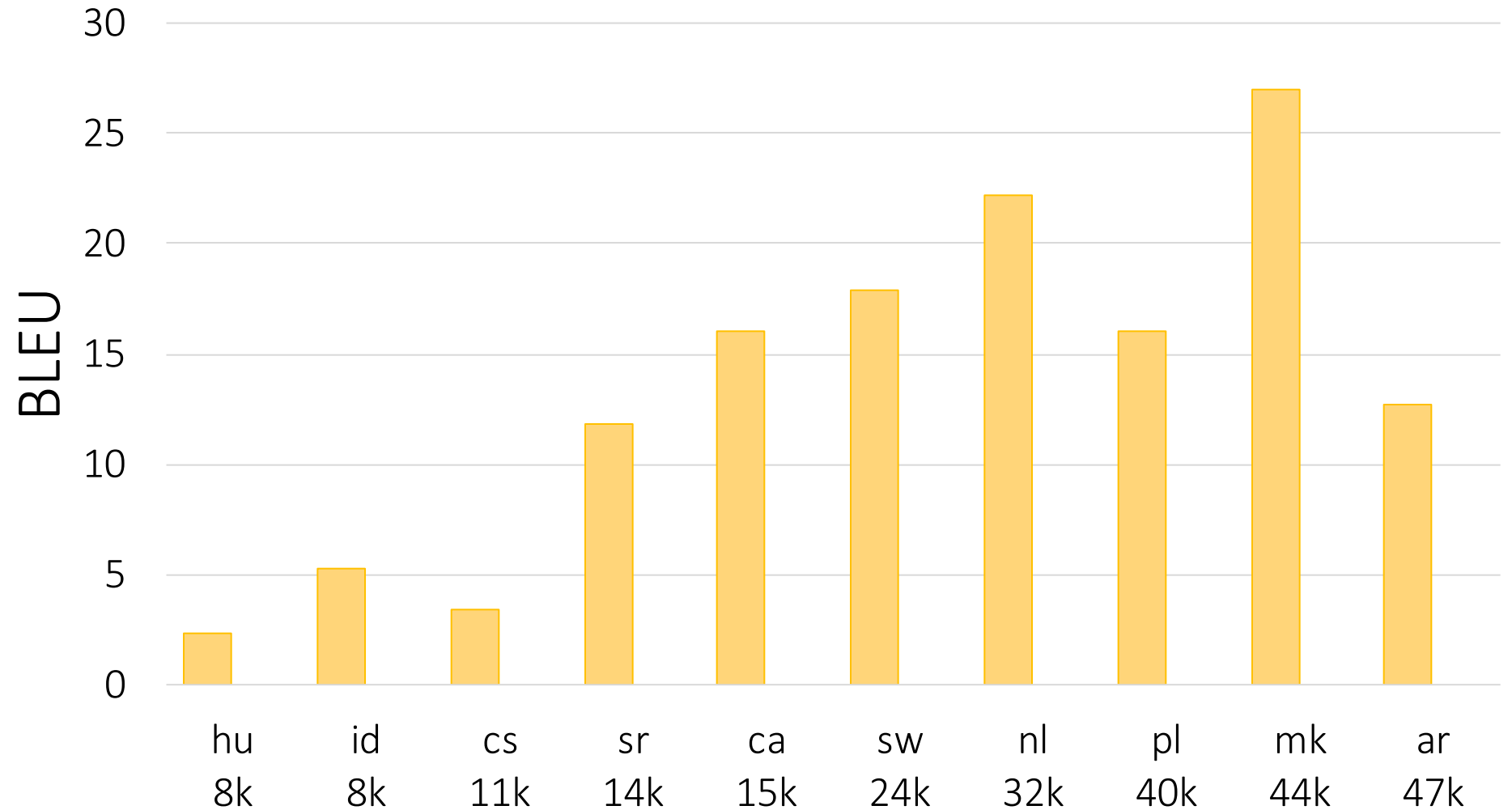
$$- \sum_{v \in \mathcal{V}} \left[\underbrace{p_{\text{para}}(y'_i = v \mid y; y'_{j < i})}_{\text{Paraphraser Output}} \times \log \underbrace{p_{\text{MT}}(y'_i = v \mid x; y'_{j < i})}_{\text{MT Model output}} \right]$$

$$- \sum_{v \in \mathcal{V}} \left[\underbrace{\mathbb{1}\{y_i = v\}}_{\text{Gold Target}} \times \log \underbrace{p_{\text{MT}}(y_i = v \mid x; y_{j < i})}_{\text{MT Model output}} \right]$$

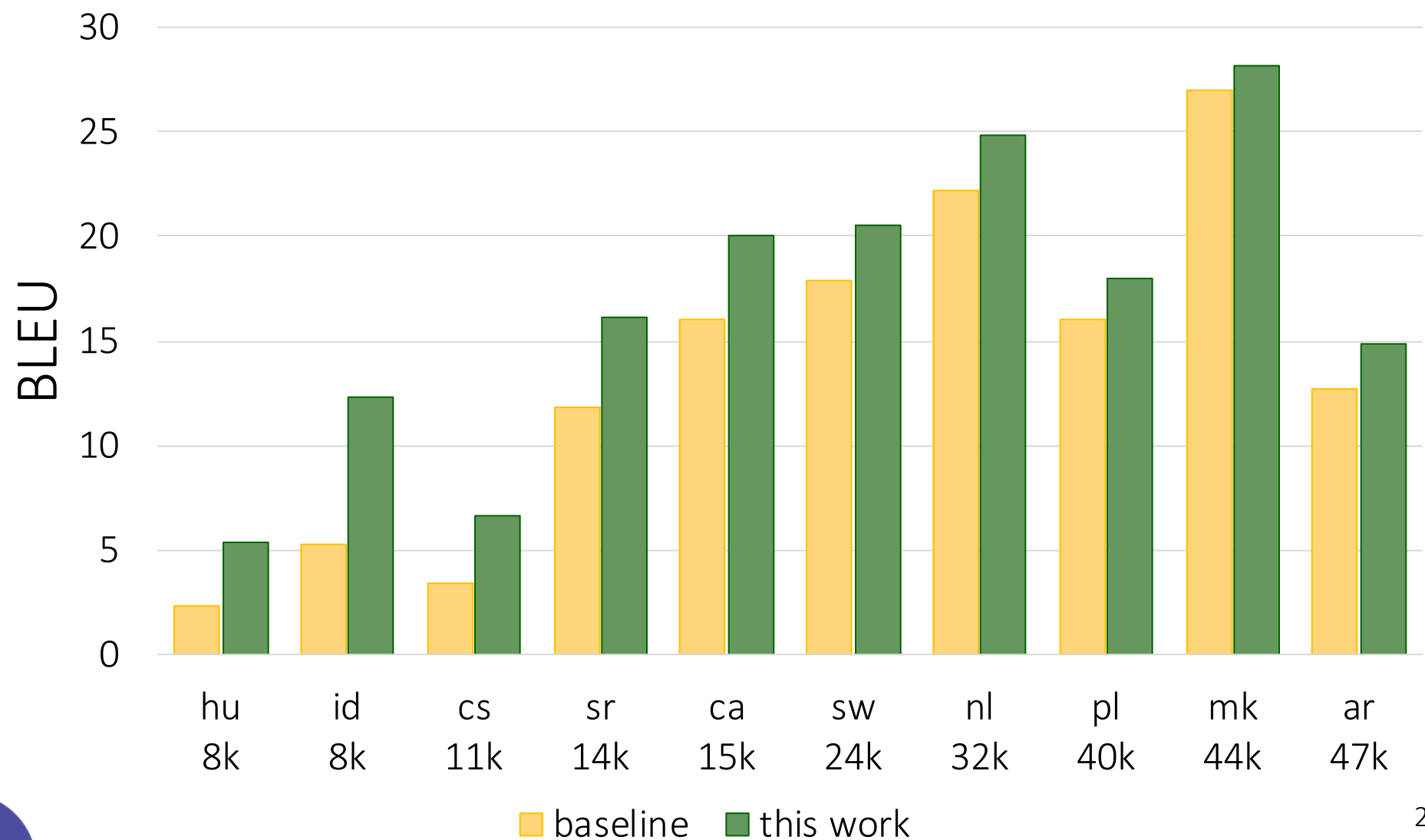
Experimental Details

- transformer model in fairseq (Ott et al., 2019)
- Global Voices corpora (Tiedemann, 2012)
 - (+ MATERIAL corpora in paper)
- Use SMRT w/ 50% probability, NLL otherwise
- English Paraphraser trained on ParaBank2 (Hu et al., 2019)
- 4k SentencePiece vocab (Kudo & Richardson 2018)
- Code, Global Voices Data splits & paraphraser released:
data.statmt.org/SMRT

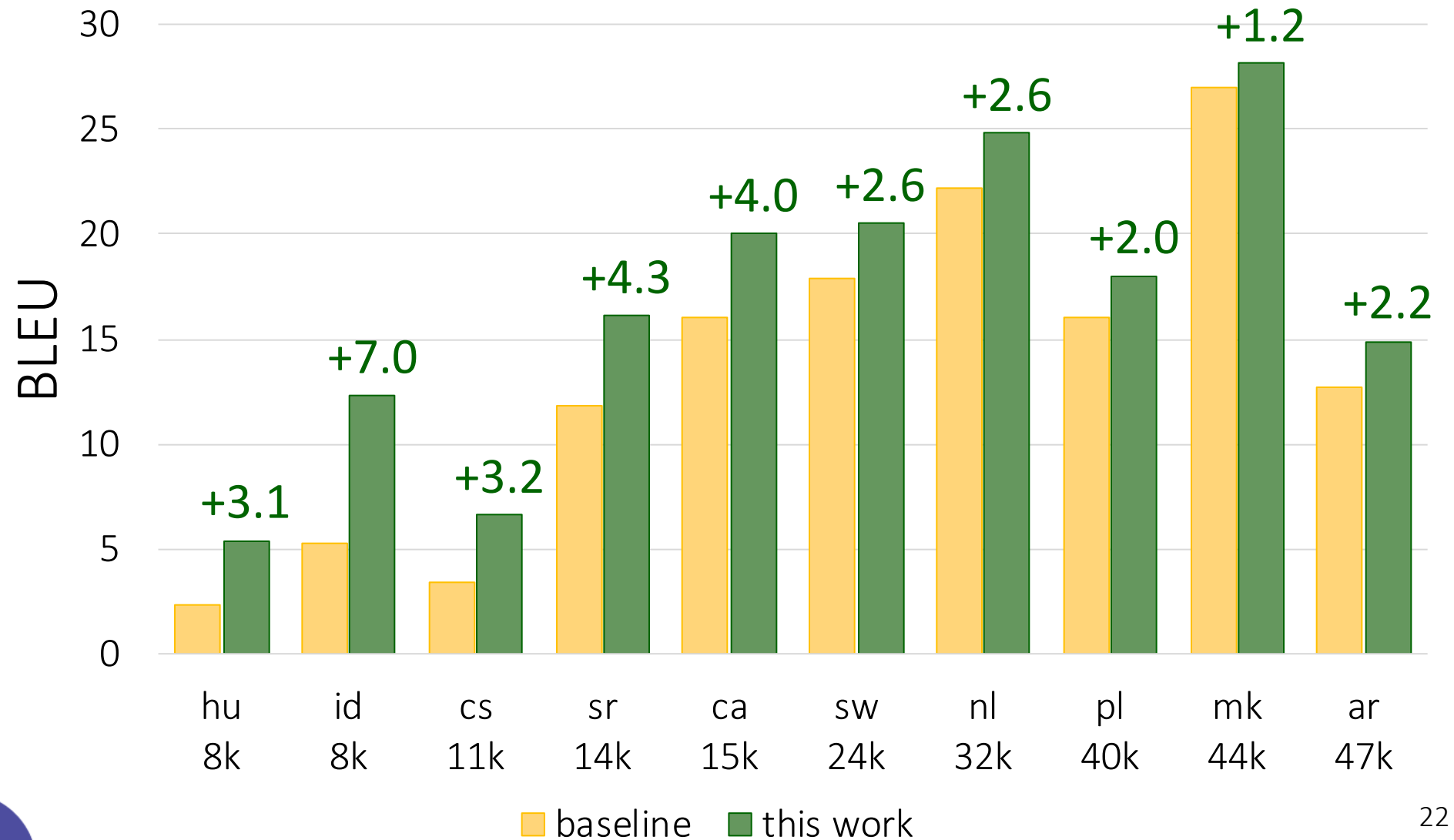
Results



Results



Results



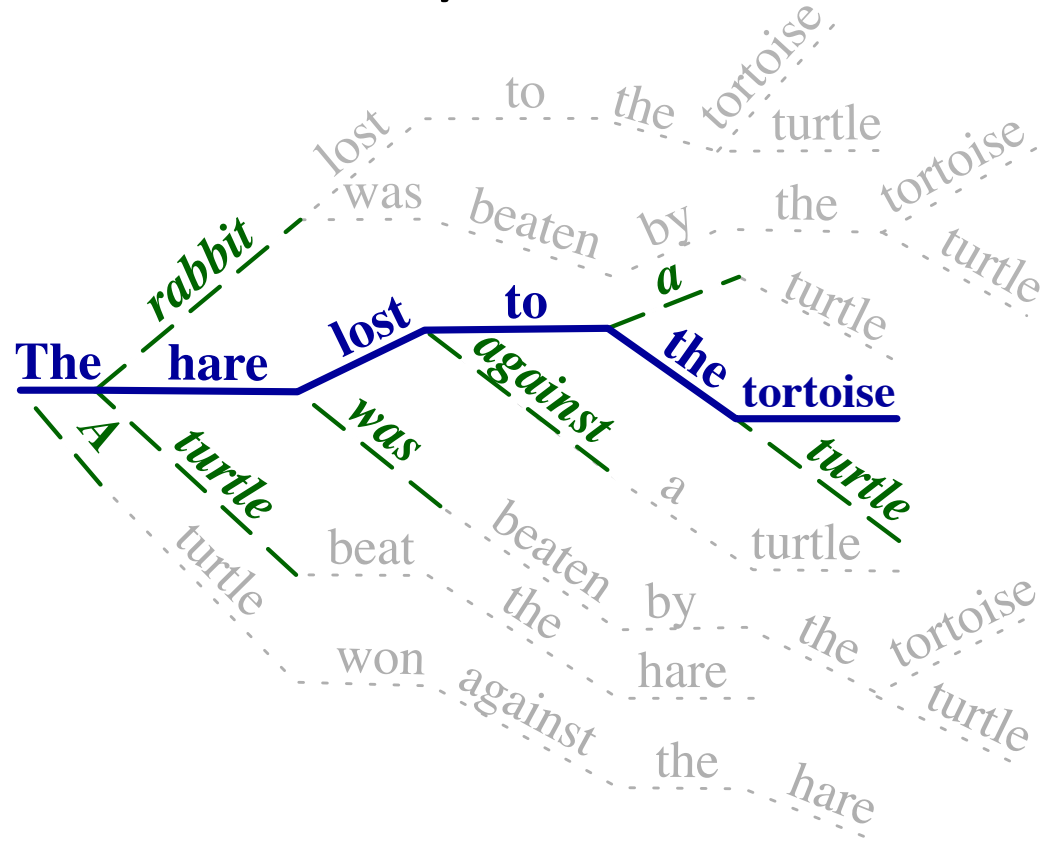
But wait, there's more! (in the paper)

- Comparison to back-translation
 - Both work, SMRT is better than BT in very low resource
 - Can combine for larger improvement
- Data ablation
 - Larger improvements in lower resource settings
- Method ablation
 - Both sampling and the distribution in the loss are helpful
- Sequence-Level Paraphrastic Augmentation
 - It works; SMRT is better

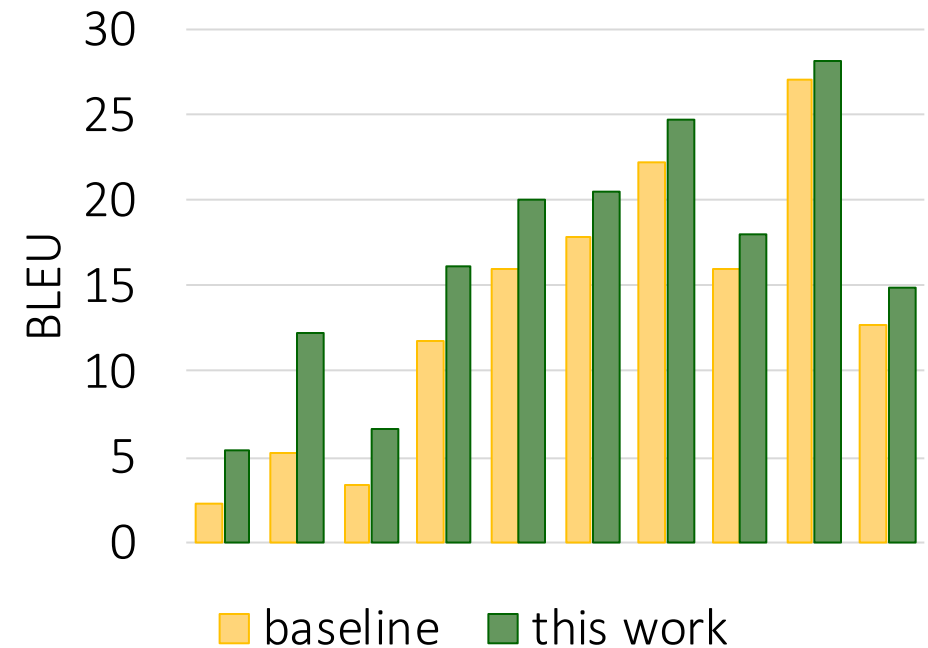
But wait, there's more! (other papers)

- Thompson & Post have a new multilingual paraphraser (Prism) that works for 39 languages and is a great MT metric (@EMNLP)
- Khayrallah & Sedoc apply SMRT + Prism to Chatbots (@EMNLP findings)
- This work will be published @EMNLP

Summary



Questions? Hiring?
Huda Khayrallah
huda@jhu.edu



Code, Global Voices Data splits & paraphraser
released: data.statmt.org/smrt