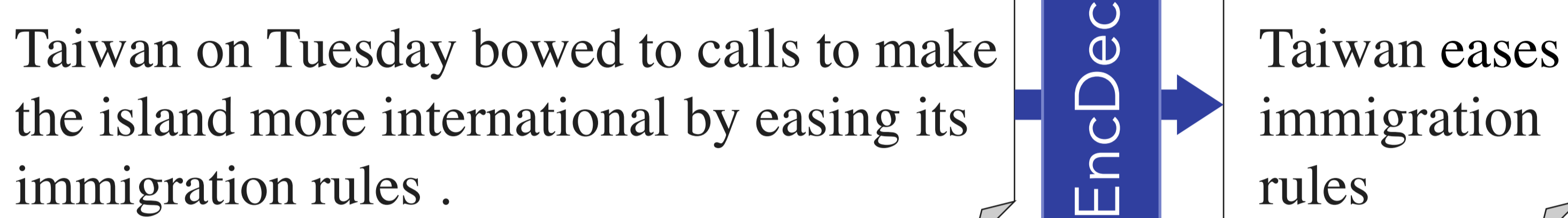


Summary

- Attention matrix is **not an optimal choice** for interpreting the output of Encoder-Decoder-based models (e.g. low alignment acc.)
- Proposed method, Unsupervised Alignment Module (UAM), models **token-wise alignment** between the source and target
- UAM provides **better interpretability of Encoder-Decoder-based models** through alignments

Task: Headline Generation

Input X : 1st Sentence of an Article Output Y : Headline



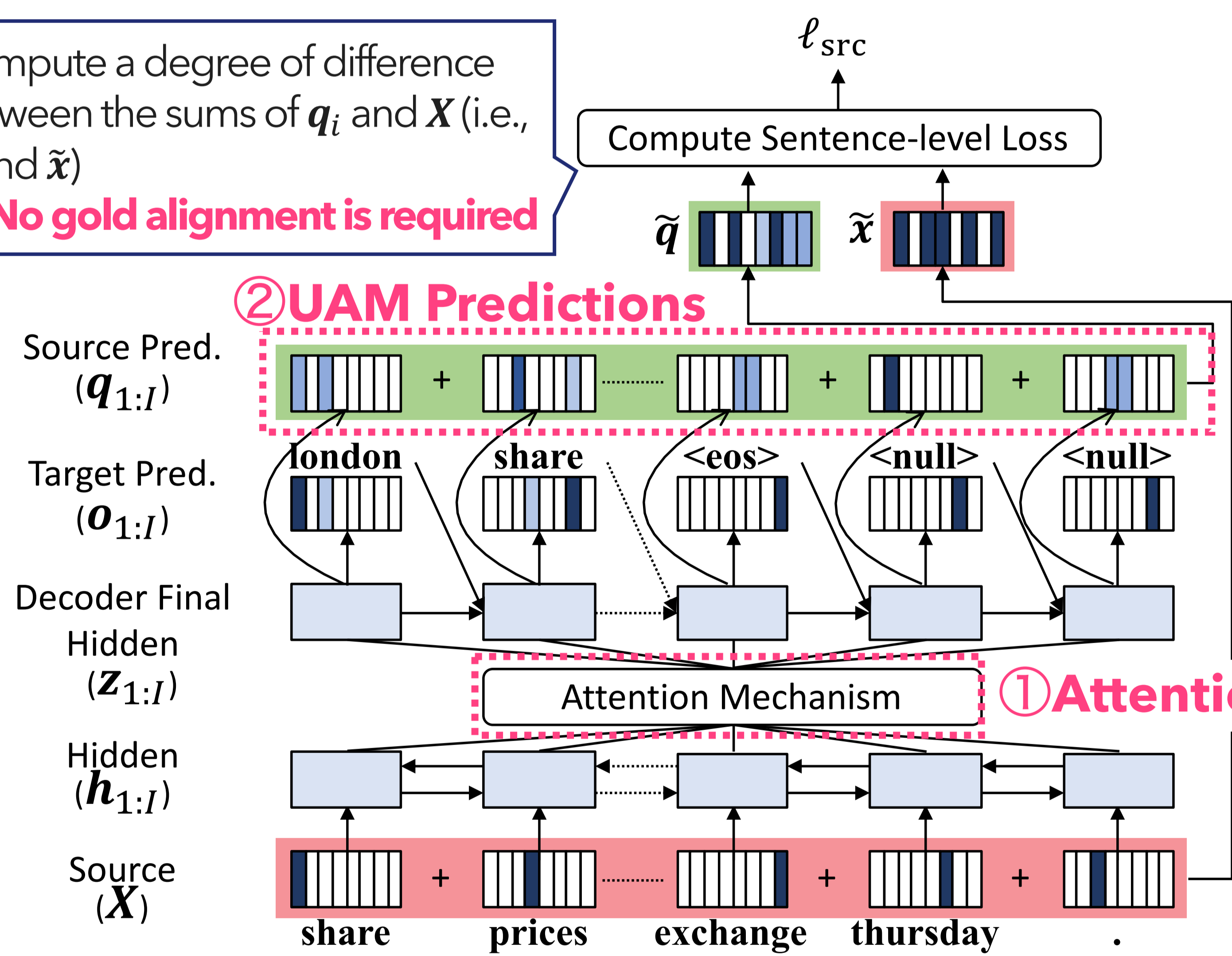
4M sentence pairs are available from Gigaword corpus

Proposed Method: UAM

Key Idea: Predict Source-side Tokens

$$G_1(\theta, \gamma) = \frac{1}{|\mathcal{D}|} \sum_{(X, Y) \in \mathcal{D}} \left(\underbrace{\ell_{\text{trg}}(Y', X, \theta)}_{\text{EncDec Loss}} + \underbrace{\ell_{\text{src}}(\tilde{x}, X, Y', \gamma, \theta)}_{\text{UAM Loss}} \right)$$

Compute a degree of difference between the sums of q_i and X (i.e., \tilde{q} and \tilde{x})
→ No gold alignment is required

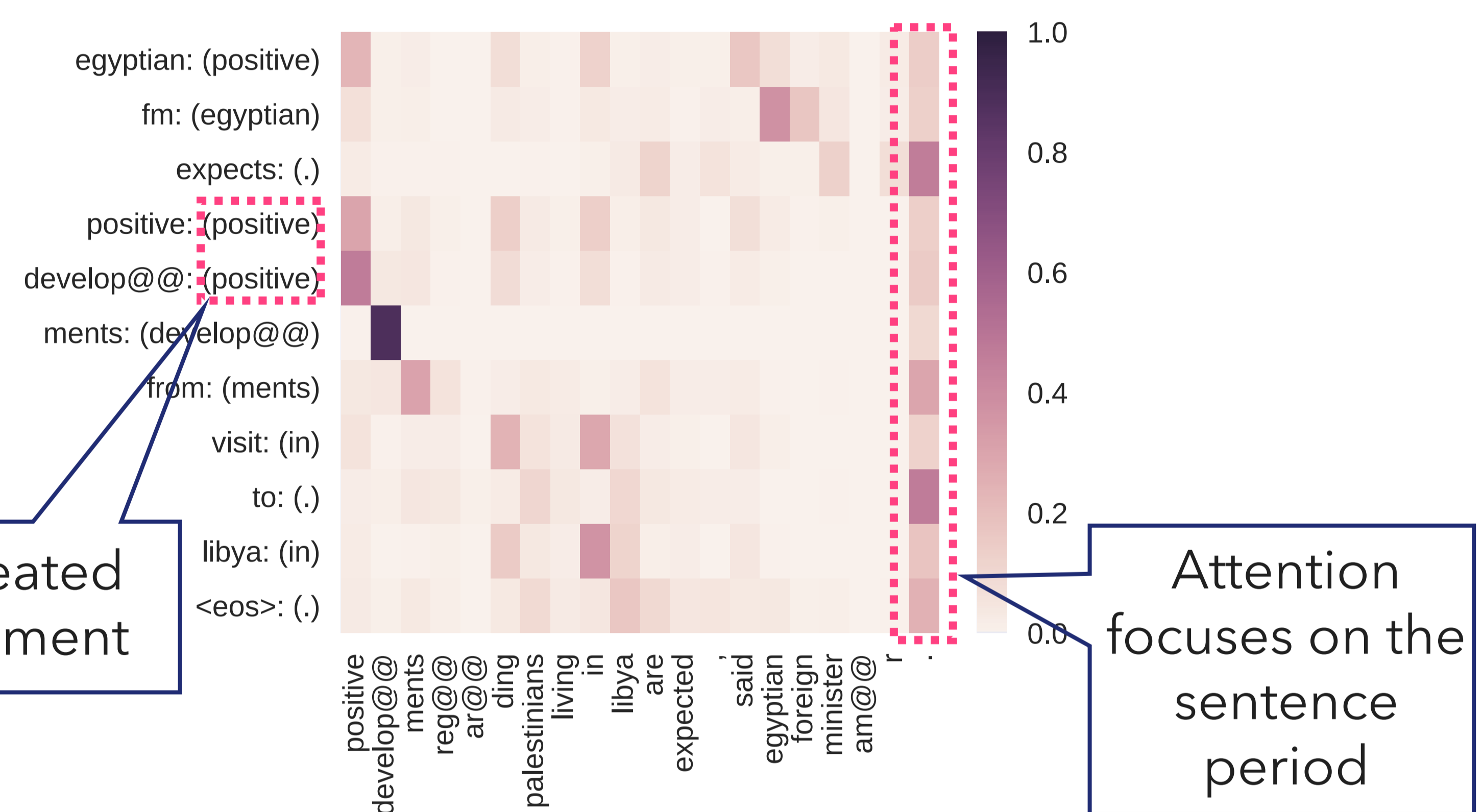


Analysis: Attention vs. UAM

Qualitative Analysis

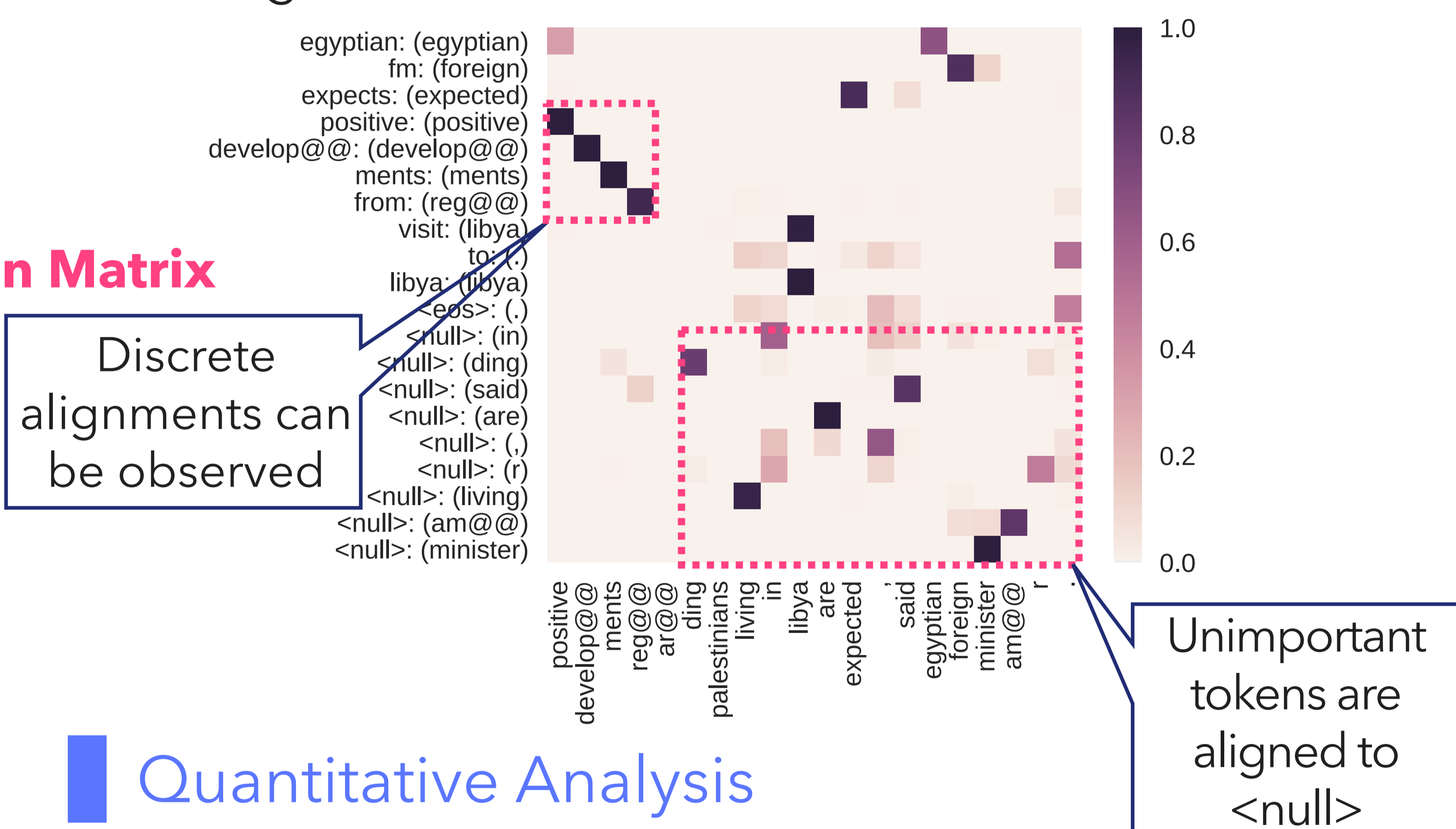
① Attention Matrix

- ☹️ Repeated alignments to same source-side tokens
- ☹️ Attention values are distributed to several tokens



② UAM Predictions

- 😊 Repeated alignments are rarely observed
- 😊 Alignments are more discrete than attention



Results (ROUGE Scores)

Table: ROUGE F1 Evaluation Results

Model	Test (Ours)			Test (Zhou)		
	RG-1	RG-2	RG-L	RG-1	RG-2	RG-L
Baseline (EncDec)	46.80	24.48	43.74	46.79	24.75	43.62
Baseline +UAM	46.91	24.86	43.87	46.89	24.93	43.68

Quantitative Analysis

Table: Alignment Accuracy

Model	Test (Ours)	Test (Zhou)
Baseline (EncDec)	8.60	5.97
Baseline+UAM	52.52	50.91

UAM improves ROUGE scores
→ <null> alignment handles selection of unimportant information?

UAM alignments are significantly more accurate than that of the attention matrix
→ better interpretation of the model