# A zero-resource approach for CS text generation and filtering for automatic speech recognition

**Olga Iakovenko** (UoS), M. Umar Farooq (UoS), Jie Chi (UoE), Elaf Islam (UoS), Brian Lu (JHU)

Mentor: Hexin Liu

# Problem

---

Little textual CS data available for training satisfactory LMs for CS ASR systems.

Automatic CS text generation does not take into account some of the prior linguistic knowledge.
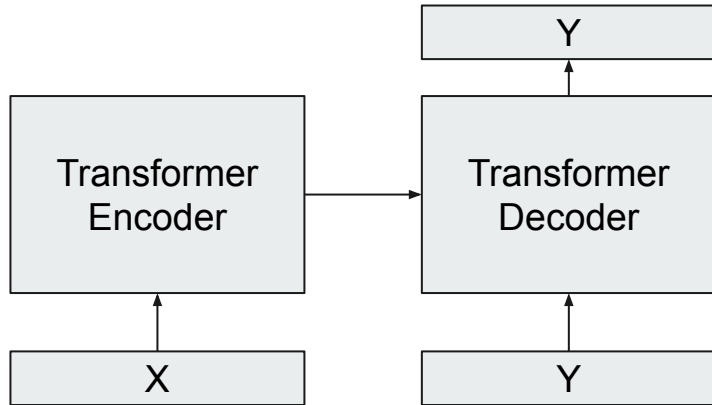
okay kay 让我拿出我的 calculator

# Solution

———

I. Generate arbitrary CS data using constrained beam search decoding within an NMT system trained with parallel data.

II. Automatically select the examples from the data which seem as the most natural CS examples
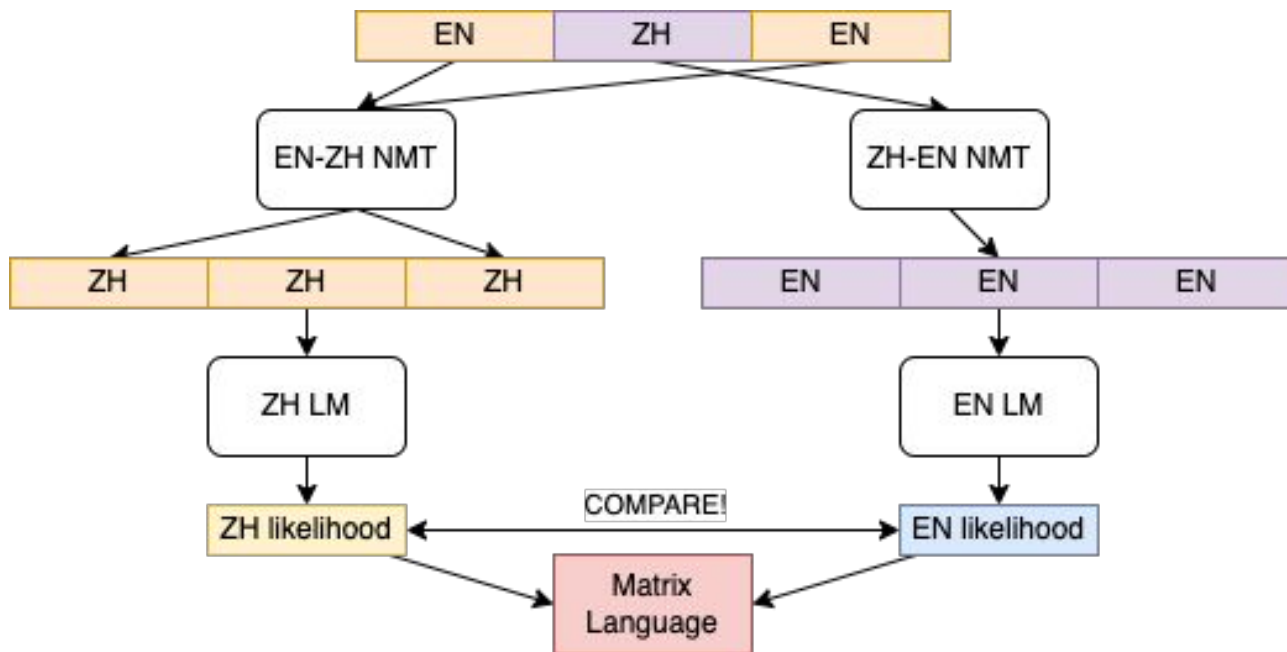
# I. Arbitrary CS text generation

---



Transformer encoder-decoder translation

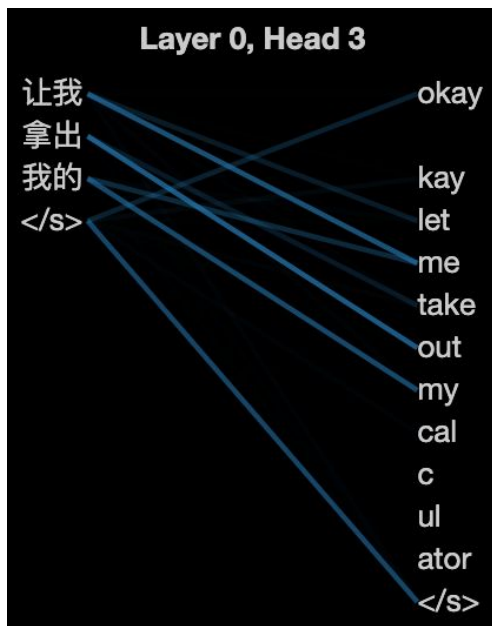Constrained beam search: token masking for creating CS points

# II. Synthetic data filtering

———

*Idea 1*: There is a main (Matrix) language that defines the grammar
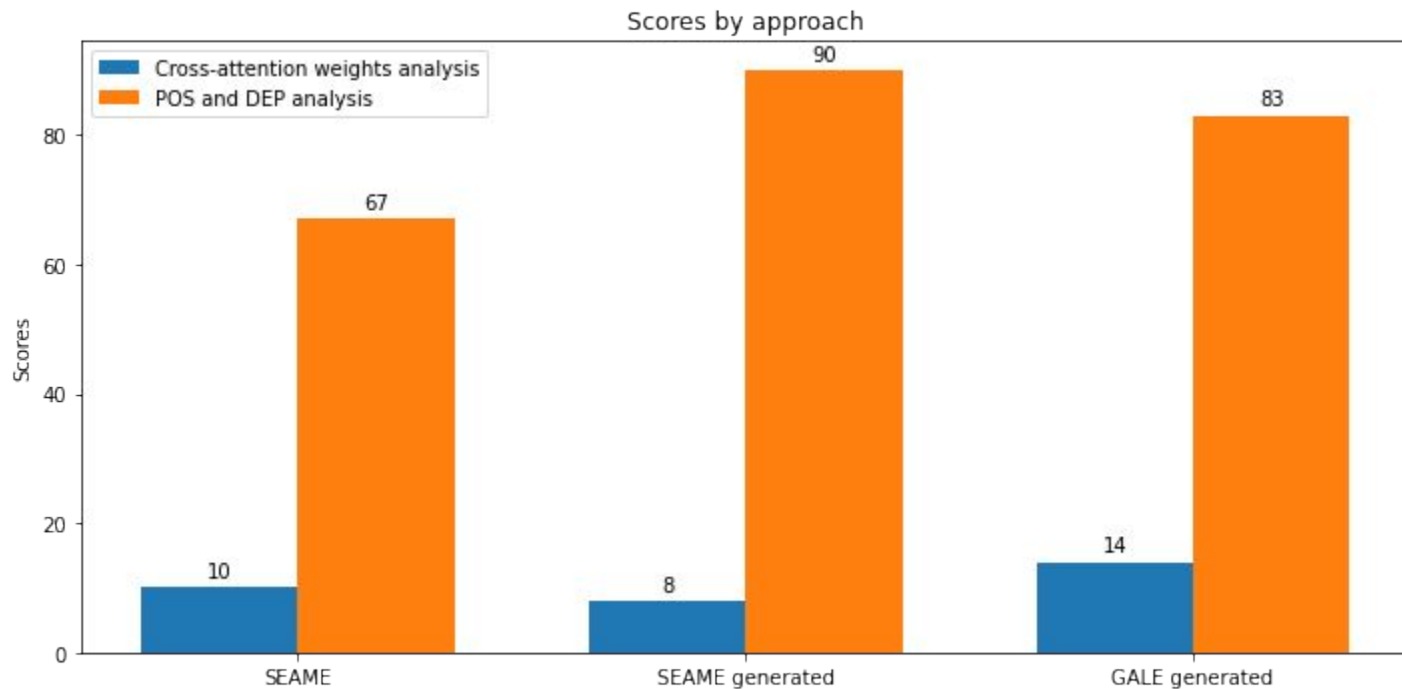
# II. Synthetic data filtering

---
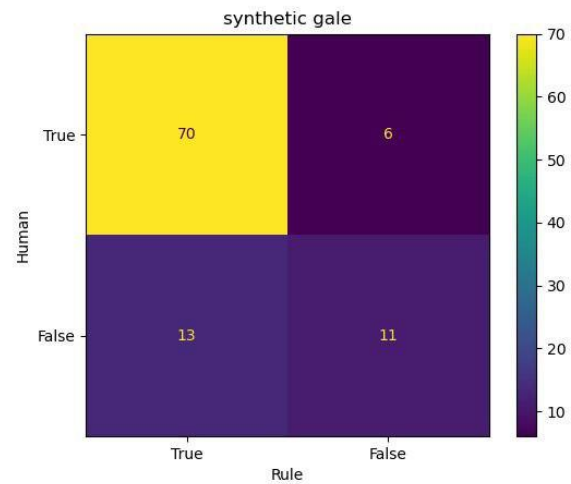
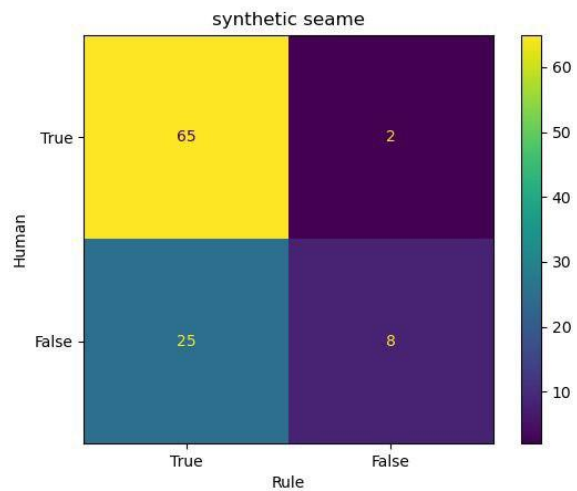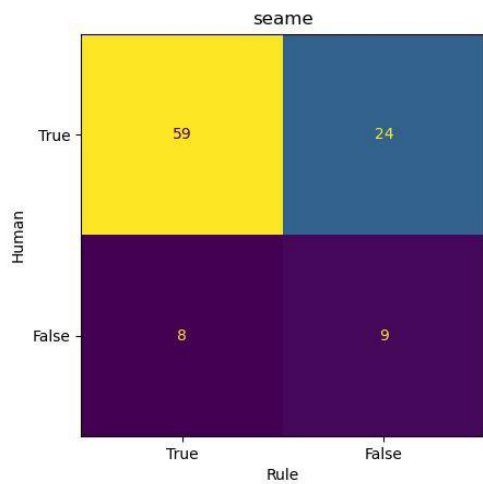*Idea 2*: Some words & word categories tend to switch more



Cross-attention analysis of incomplete translations

POS tag & dependency analysis: function words do not appear in isolation!

# Results



Scores by approach

# Results

‒ ‒ ‒

**Thank you!**

Demo of the project:
https://t.me/cs_assessment
_bot

@CS_ASSESSMENT_BOT