AUTOPROMPT: Eliciting Knowledge from Language Models with Automatically Generated Prompts

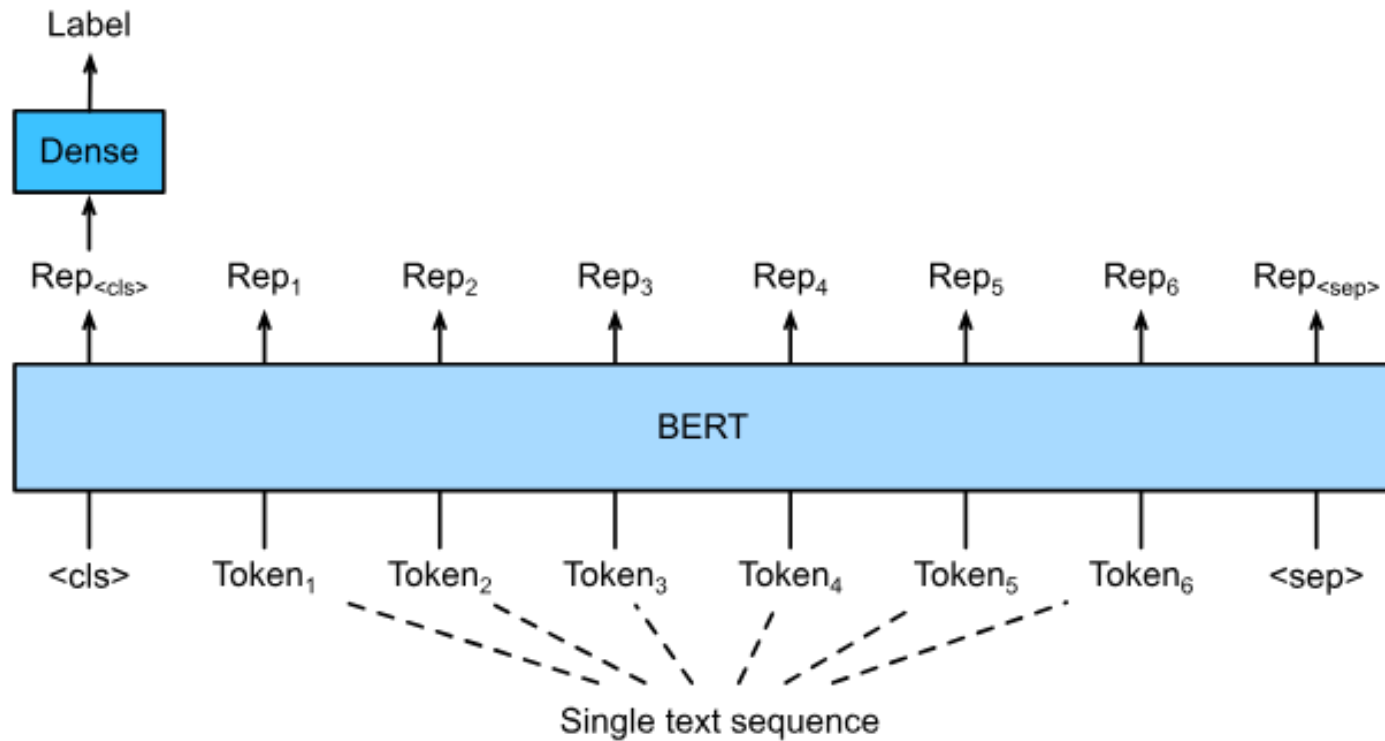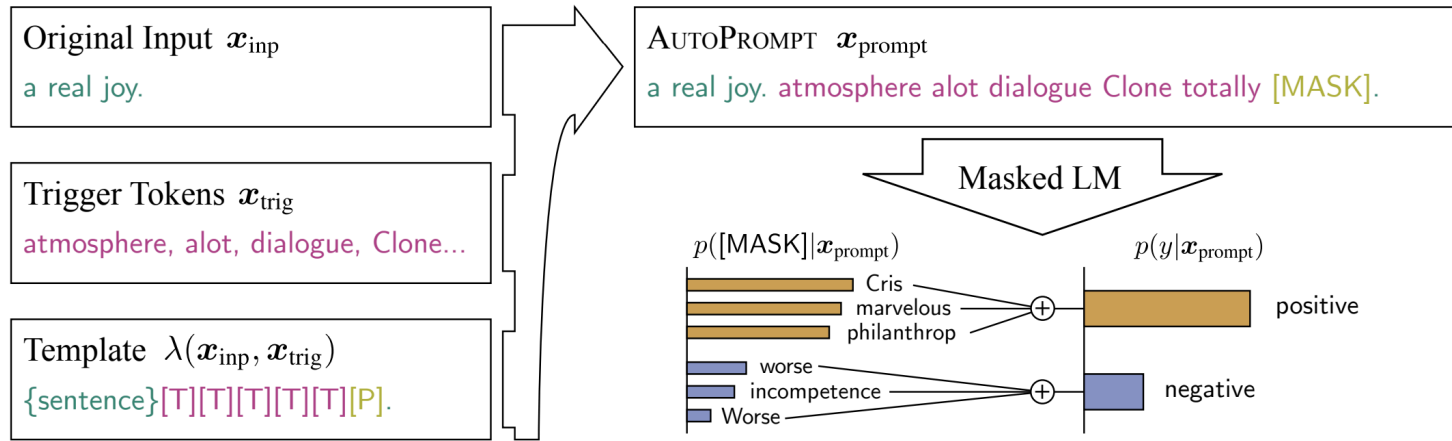Jiayu Yao

# Motivation

Determine whether knowledge are learned during language model training process

PROMPTING TO ELICIT KNOWLEDGE FROM LM

# Method

# Method



- Find Trigger Prompt
- Find Label Token Set

# Method

$$p(y|\boldsymbol{x}_{\text{prompt}}) = \sum_{w \in \mathcal{V}_y} p([\text{MASK}] = w | \boldsymbol{x}_{\text{prompt}})$$

$$\mathcal{V}_{\text{cand}} = \underset{w \in \mathcal{V}}{\text{top-}k} \left[ \boldsymbol{w}_{\text{in}}^{T} \nabla \log p(y|\boldsymbol{x}_{\text{prompt}}) \right]$$

# Label Projection

Replace $y$ with $\boldsymbol{w}_{out}$ (token output embedding)

$$s(y, w) \;=\; p(y|\boldsymbol{w}_{\text{out}}).$$

$$\mathcal{V}_y = \underset{w \in \mathcal{V}}{\text{top-}k}\left[s(y, w)\right]$$

# Experiment

| Model | Dev | Test |
|---|---|---|
| BiLSTM | - | $82.8^{\dagger}$ |
| BiLSTM + ELMo | - | $89.3^{\dagger}$ |
| BERT (linear probing) | 85.2 | 83.4 |
| BERT (finetuned) | - | $93.5^{\dagger}$ |
| RoBERTa (linear probing) | 87.9 | 88.8 |
| RoBERTa (finetuned) | - | $96.7^{\dagger}$ |
| BERT (manual) | 63.2 | 63.2 |
| BERT (AUTOPROMPT) | 80.9 | 82.3 |
| RoBERTa (manual) | 85.3 | 85.2 |
| RoBERTa (AUTOPROMPT) | 91.2 | 91.4 |

Table 1: **Sentiment Analysis** performance on the SST-2 test set of supervised classifiers (top) and fill-in-the-blank MLMs (bottom). Scores marked with $\dagger$ are from the GLUE leaderboard: http://gluebenchmark.com/leaderboard.

# Discussion

- LLM Evaluation
  - Conventional Evaluation
  - Co-evaluation
- Prompt Guidence

# Thanks