# In-Context-Learning

## Experiments on Manual Template Engineering

Ann Weitz, Xiao Yang, Annabelle Runge, Lea Krumbach,
Emanuel Nierstenhöfer & Sandy Rodrigues

June 12, 2024

Outline

# WHAT IS MANUAL TEMPLATE ENGINEERING?

# What is Manual Template Engineering?

- **Definition:**

  - Creating specific input formats (templates) for AI models.

  - Guides the responses of AI.

- **Purpose:**

  - Improves the accuracy and quality of AI responses.

# Advantages

- **Improved Accuracy:** Enhances the precision of responses by reducing ambiguity.

- **Consistency:** Ensures uniformity in responses across different instances.

- **Efficiency:** Reduces the need for extensive post-processing or corrections.

- **User Satisfaction:** Leads to more relevant and satisfactory interactions for end-users.

# 2. PROMPTING TECHNIQUES

**What prompting techniques do you know?**

(or: How would you structure your prompts?)

# Tips for structurizing your prompts

- Be specific but avoid unnecessary details
- Use Keywords
  - "Write"
  - "Classify"
  - "Summarize"
  - "Translate"
  - "Order"
- Experiment with different prompts
- Context Setting
- Separate instruction and context (e.g., '', """"')
- Articulate the desired output format

# Prompting techniques

1. Zero-Shot Prompting

2. Few-Shot Prompting

3. Chain-of-Thought Prompting

4. Generate Knowledge Prompting

5. Tree of Thoughts

# Zero-Shot Prompting: Definition and Advantages

- **Definition:** AI models can perform tasks without specific training.

- **Advantages:**
  - **Versatility:** Handles various tasks without task-specific training.
  - **Efficiency:** Saves time and resources by not needing task-specific data.
  - **Adaptability:** Quickly adjusts to new tasks with minimal modifications.

# Zero-Shot Prompting: Example and Applications

- **Example:**
    - **Prompt:**
      "'Classify the text into neutral, negative or positive.
      Text: I think the vacation is okay.
      Sentiment: "'
    - **Response:** Neutral

- **Applications:**
    - Language translation
    - Text summarization
    - Question answering
    - Content generation

# Few-Shot Prompting: Definition and Advantages

- **Definition:** AI models learn from a few examples to perform tasks.

- **Advantages:**
  - **Flexibility:** Adapts to various tasks with minimal examples.
  - **Scalability:** Scales efficiently with a small dataset.
  - **Accuracy:** Maintains high performance with limited data.

# Few-Shot Prompting: Example and Applications

- **Example:**
  - **Prompt:**
    "'This is awesome! // Negative
    This is bad! // Positive
    Wow that movie was rad! // Positive
    What a horrible show¡"//
  - **Response:** Negative

- **Applications:**
  - Content summarization
  - Sentiment analysis
  - Text classification
  - Document categorization

# Arithmetic Tasks



**Standard Prompting**

**Model Input**

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

**Model Output**

A: The answer is 27. ✖

Wei et al. 2022

# Chain-of-Thought Prompting

**Aim**: enable complex reasoning capabilities through intermediate reasoning steps; generate a chain of thought

**Why**?
$\rightarrow$ insight into reasoning path of LM (facilitates debugging)
$\rightarrow$ useful for math word problems, commonsense reasoning, and symbolic manipulation

# Chain-of-Thought Prompting



**Standard Prompting**

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The answer is 27. ✗

**Chain-of-Thought Prompting**

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had 23 - 20 = 3. They bought 6 more apples, so they have 3 + 6 = 9. The answer is 9. ✓

Wei et al. 2022

# Limitations

- High scale models $\rightarrow$ high performance
- Effectiveness of LM reliant on complexity of problem
- Uncertainty whether LM is actually "reasoning"
- Costly to serve in real-world applications

# Zero-shot CoT Prompting

Add "Let's think step-by-step" to the original prompt

# Zero-shot CoT Prompting

Add "Let's think step-by-step" to the original prompt



## (a) Few-shot

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?
A: The answer is 11.

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?
A:

(Output) The answer is 8. ✗

## (b) Few-shot-CoT

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?
A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?
A:

(Output) The juggler can juggle 16 balls. Half of the balls are golf balls. So there are 16 / 2 = 8 golf balls. Half of the golf balls are blue. So there are 8 / 2 = 4 blue golf balls. The answer is 4. ✓

## (c) Zero-shot

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?
A: The answer (arabic numerals) is

(Output) 8 ✗

## (d) Zero-shot-CoT (Ours)

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?
A: Let's think step by step.

(Output) There are 16 balls in total. Half of the balls are golf balls. That means that there are 8 golf balls. Half of the golf balls are blue. That means that there are 4 blue golf balls. ✓
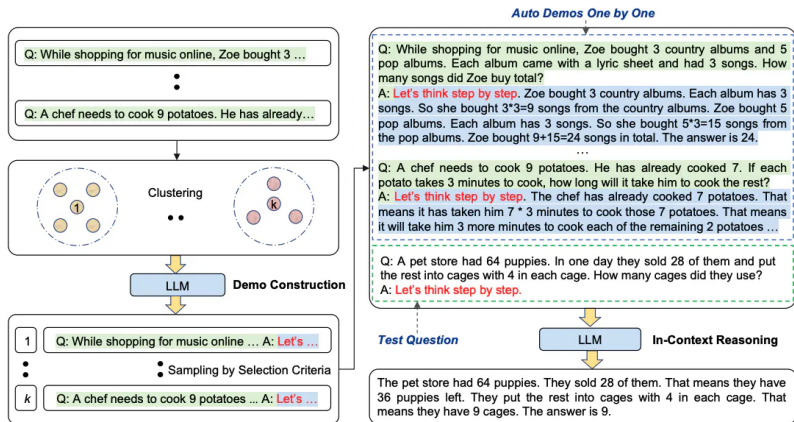
Kojima et al. 2022

# Automatic CoT

Eliminate manual efforts by leveraging LLMs with "Let's think step by step" prompt to generate reasoning chains for demonstrations one by one

Two stages:
(1) **Question clustering**: partition questions of a given dataset into clusters
(2) **Demonstration sampling**: select a representative question from each cluster and generate its reasoning chain using zero-shot-CoT
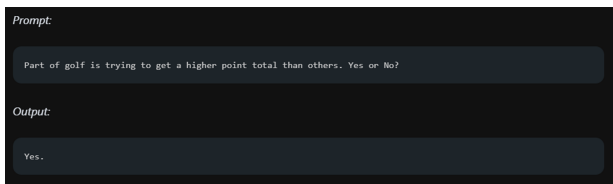
# Automatic CoT



Zhang et al. 2022

# Generate Knowledge Prompting

**Popular idea:** incorporate knowledge to help the model make more accurate ideas

# Generate Knowledge Prompting

**Popular idea**: incorporate knowledge to help the model make more accurate ideas

$\rightarrow$ Can the model also be used to generate knowledge before making a prediction? (Liu et al. 2021)

# Generate Knowledge Prompting: Knowledge

### Prompt

```
Input: A rock is the same size as a pebble.
Knowledge: A pebble is a clast of rock with a particle size of 4 to 64 millimetres based on the Udden-Wentworth
scale of sedimentology. Pebbles are generally considered larger than granules (2 to 4 millimetres diameter) and
smaller than cobbles (64 to 256 millimetres diameter).
Input: Part of golf is trying to get a higher point total than others.
Knowledge:
```

# Generate Knowledge Prompting: Knowledge

Prompt

```
Input: A rock is the same size as a pebble.
Knowledge: A pebble is a clast of rock with a particle size of 4 to 64 millimetres based on the Udden-Wentworth
scale of sedimentology. Pebbles are generally considered larger than granules (2 to 4 millimetres diameter) and
smaller than cobbles (64 to 256 millimetres diameter).
Input: Part of golf is trying to get a higher point total than others.
Knowledge:
```

Knowledge 1

```
The objective of golf is to play a set of holes in the least number of strokes. A round of golf typically
consists of 18 holes. Each hole is played once in the round on a standard golf course. Each stroke is counted as
one point, and the total number of strokes is used to determine the winner of the game.
```

# Generate Knowledge Prompting: Knowledge

Prompt

```
Input: A rock is the same size as a pebble.
Knowledge: A pebble is a clast of rock with a particle size of 4 to 64 millimetres based on the Udden-Wentworth
scale of sedimentology. Pebbles are generally considered larger than granules (2 to 4 millimetres diameter) and
smaller than cobbles (64 to 256 millimetres diameter).
Input: Part of golf is trying to get a higher point total than others.
Knowledge:
```

Knowledge 1

```
The objective of golf is to play a set of holes in the least number of strokes. A round of golf typically
consists of 18 holes. Each hole is played once in the round on a standard golf course. Each stroke is counted as
one point, and the total number of strokes is used to determine the winner of the game.
```

Knowledge 2

```
Golf is a precision club-and-ball sport in which competing players (or golfers) use many types of clubs to hit
balls into a series of holes on a course using the fewest number of strokes. The goal is to complete the course
with the lowest score, which is calculated by adding up the total number of strokes taken on each hole. The
player with the lowest score wins the game.
```

Liu et al. 2021

# Generate Knowledge Prompting: Prediction

$\rightarrow$ Integrate knowledge and get a prediction

Example based on Knowledge 1:

```
Question: Part of golf is trying to get a higher point total than others. Yes or No?
Knowledge: The objective of golf is to play a set of holes in the least number of strokes. A round of golf
typically consists of 18 holes. Each hole is played once in the round on a standard golf course. Each stroke is
counted as one point, and the total number of strokes is used to determine the winner of the game.
Explain and Answer:

No, the objective of golf is not to get a higher point total than others. Rather, the objective is to play a set
of holes in the least number of strokes. The total number of strokes is used to determine the winner of the
game, not the total number of points.
```

# Generate Knowledge Prompting: Prediction

$\rightarrow$ Integrate knowledge and get a prediction

Example based on Knowledge 2:

```
Question: Part of golf is trying to get a higher point total than others. Yes or No?
Knowledge: Golf is a precision club-and-ball sport in which competing players (or golfers) use many types of
clubs to hit balls into a series of holes on a course using the fewest number of strokes. The goal is to
complete the course with the lowest score, which is calculated by adding up the total number of strokes taken on
each hole. The player with the lowest score wins the game.
Explain and Answer:

Yes, part of golf is trying to get a higher point total than others. Each player tries to complete the course
with the lowest score, which is calculated by adding up the total number of strokes taken on each hole. The
player with the lowest score wins the game.
```
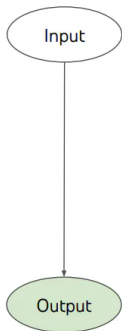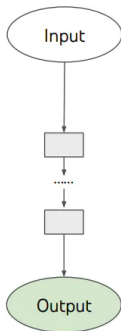
# Tree of Thoughts

Main idea:

1. Generate multiple reasoning paths for one question/problem
2. Evaluate reasoning paths in accordance with criteria
3. Expand and prune
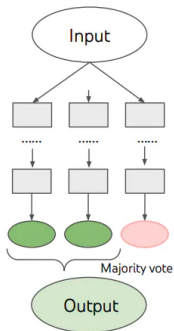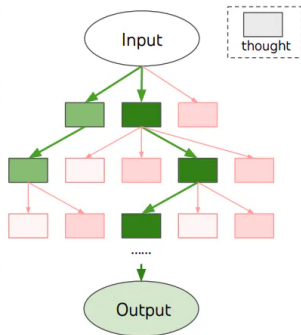4. Choose final path based on the highest score

# Tree of Thoughts



(a) Input-Output Prompting (IO)

(c) Chain of Thought Prompting (CoT)

(c) Self Consistency with CoT (CoT-SC)

(d) Tree of Thoughts (ToT)

Yao et al. 2024

# 3. LANGUAGE MODELS AS KNOWLEDGE BASES

## 3.1. The LAMA probe

**LA**nguage **M**odel **A**nalysis probe (Petroni et al., 2019)

$\rightarrow$ How much knowledge is present in pretrained Language Models?
$\rightarrow$ Can pretrained LLMs outperform state-of-the-art NLP methods in receiving knowledge?
$\rightarrow$ How does the performance of LLMs differ for different kinds of knowledge (relational, common sense, factual)?

**"Knowledge:"**

- ▶ (subject, relation, object)
- ▶ (question, answer)

## 3.1. The LAMA probe

**Procedure**

- ▶ manually convert "knowledge" (from existing knowledge sources) into cloze-statements
- ▶ Example: (Einstein, born_in, Ulm) → "Einstein was born in [MASK]"
- ▶ ask models to predict the masked token/missing object ([MASK])

Assumption: LLM "knows" a fact, if it can predict a single object [MASK] or answer [MASK] token.

# 3.1. The LAMA probe: Considerations

## 1. Manually defined templates

| Relation | Query | Answer | Generation |
|----------|-------|--------|------------|
| P19 | Francesco Bartolomeo Conti was born in ___. | Florence | Rome [-1.8], **Florence [-1.8]**, Naples [-1.9], Milan [-2.4], Bologna [-2.5] |
| P20 | Adolphe Adam died in ___. | Paris | **Paris [-0.5]**, London [-3.5], Vienna [-3.6], Berlin [-3.8], Brussels [-4.0] |
| P279 | English bulldog is a subclass of ___. | dog | dogs [-0.3], breeds [-2.2], **dog [-2.4]**, cattle [-4.3], sheep [-4.5] |
| P37 | The official language of Mauritius is ___. | English | **English [-0.6]**, French [-0.9], Arabic [-6.2], Tamil [-6.7], Malayalam [-7.0] |

## 2. Single token prediction
$\rightarrow$ only single token objects as prediction targets

# 3.1. The LAMA probe: Considerations

### 3. Object slot predictions
$\rightarrow$ only query object slots *not* subject or relation slots
$\rightarrow$ relations can be expressed with many different wordings: what would be the correct pattern for a relation?

### 4. Intersection of Vocabularies
$\rightarrow$ intersection of the vocabulary all models were trained on - about 21k tokens
$\rightarrow$ every model can only rank tokens of that vocabulary
$\rightarrow$ the larger the vocabulary, the harder to rank correct token

# 3.2. Knowledge sources

- ▶ Google-RE
- ▶ T-REx
- ▶ ConceptNet
- ▶ SQuAD

## 3.2. Knowledge sources

### 1. **Google-RE**

```
{"pred": "/people/person/date_of_birth", "sub": "/m/0j240kx", "obj": "1944",
"evidences": [{"url": "http://en.wikipedia.org/wiki/Gao_Yu_(journalist)",
"snippet": "Gao Yu (born 1944) is a Chinese journalist.",
"considered_sentences": ["Gao Yu (born 1944) is a Chinese journalist ."]}],
[], "sub_label": "Gao Yu", [], "obj_label": "1944", [],
"masked_sentences": ["Gao Yu (born [MASK]) is a Chinese journalist ."]}
```

$\rightarrow$ about 60k facts from Wikipedia
$\rightarrow$ 3 relations (place_of_birth, date_of_birth, place_of_death)
$\rightarrow$ manually defined templates

### 2. **T-REx**

$\rightarrow$ 41 relations with about 1000 facts per relation from Wikidata
$\rightarrow$ manually defined templates

# 3.3. Language Models vs. Baselines

**Language Models:**

- ▶ fairseq-fconv (Fs)
- ▶ Transformer-XL (Txl)
- ▶ ELMo base (Eb)
- ▶ ELMo 5.5 (E5B)
- ▶ BERT base (Bb)
- ▶ BERT large (Bl)

## 3.3. Language Models vs. Baselines

Exercise 2 (see colab)

## 3.3. Language Models vs. Baselines

**Metrics for LLMs: Ranking and mean precision at k (P@k)**

▶ model generates output-prediction layer (logits) for possible objects which are just unnormalized numbers (eg. [2.3, -0.5, 4.6])

▶ softmax function is applied to those logits which converts the raw scores into probabilities that sum up to 1 (eg. [0.878,0.045,0.077])

▶ those probabilities are ranked in descending order (first position = highest probability)

▶ k is the number of predictions we consider

▶ if ground truth object is among these top k predictions, it's counted as a correct prediction

▶ calculate mean precision by dividing the correct predicted objects by all predicted objects

# 3.3. Language Models vs. Baselines

**Baselines:**
= existing methods/systems commonly used for *relation knowledge extraction*

- ▶ Freq (Freq)
- ▶ Relation Extraction with naiive entity linking ($RE_n$)
- ▶ Relation Extraction with oracle entity linking ($RE_o$)
- ▶ DrQA

# 3.5. Results

**Results with <span style="color:purple">p@1:</span>**

| Corpus | Relation | Statistics | | Baselines | | KB | | LM | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | #Facts | #Rel | Freq | DrQA | $RE_n$ | $RE_o$ | Fs | Txl | Eb | E5B | Bb | Bl |
| Google-RE | birth-place | 2937 | 1 | 4.6 | - | 3.5 | 13.8 | 4.4 | 2.7 | 5.5 | 7.5 | 14.9 | **16.1** |
| | birth-date | 1825 | 1 | 1.9 | - | 0.0 | **1.9** | 0.3 | 1.1 | 0.1 | 0.1 | 1.5 | 1.4 |
| | death-place | 765 | 1 | 6.8 | - | 0.1 | 7.2 | 3.0 | 0.9 | 0.3 | 1.3 | 13.1 | **14.0** |
| | Total | 5527 | 3 | 4.4 | - | 1.2 | 7.6 | 2.6 | 1.6 | 2.0 | 3.0 | 9.8 | **10.5** |
| T-REx | 1-1 | 937 | 2 | 1.78 | - | 0.6 | 10.0 | 17.0 | 36.5 | 10.1 | 13.1 | 68.0 | **74.5** |
| | $N$-1 | 20006 | 23 | 23.85 | - | 5.4 | **33.8** | 6.1 | 18.0 | 3.6 | 6.5 | 32.4 | 34.2 |
| | $N$-$M$ | 13096 | 16 | 21.95 | - | 7.7 | **36.7** | 12.0 | 16.5 | 5.7 | 7.4 | 24.7 | 24.3 |
| | Total | 34039 | 41 | 22.03 | - | 6.1 | **33.8** | 8.9 | 18.3 | 4.7 | 7.1 | 31.1 | 32.3 |
| ConceptNet | Total | 11458 | 16 | 4.8 | - | - | - | 3.6 | 5.7 | 6.1 | 6.2 | 15.6 | **19.2** |
| SQuAD | Total | 305 | - | - | 37.5 | - | - | 3.6 | 3.9 | 1.6 | 4.3 | 14.1 | 17.4 |

Table 2: Mean precision at one (P@1) for a frequency baseline (Freq), DrQA, a relation extraction with naïve entity linking ($RE_n$), oracle entity linking ($RE_o$), fairseq-fconv (Fs), Transformer-XL large (Txl), ELMo original (Eb), ELMo 5.5B (E5B), BERT-base (Bb) and BERT-large (Bl) across the set of evaluation corpora.

Petroni et al. (2019)

# 3.5. Results

**Results with p@1:**

| Corpus | Relation | Statistics | | Baselines | | KB | | LM | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | #Facts | #Rel | Freq | DrQA | $RE_n$ | $RE_o$ | Fs | Txl | Eb | E5B | Bb | Bl |
| Google-RE | birth-place | 2937 | 1 | 4.6 | - | 3.5 | 13.8 | 4.4 | 2.7 | 5.5 | 7.5 | 14.9 | **16.1** |
| | birth-date | 1825 | 1 | 1.9 | - | 0.0 | **1.9** | 0.3 | 1.1 | 0.1 | 0.1 | 1.5 | 1.4 |
| | death-place | 765 | 1 | 6.8 | - | 0.1 | 7.2 | 3.0 | 0.9 | 0.3 | 1.3 | 13.1 | **14.0** |
| | Total | 5527 | 3 | 4.4 | - | 1.2 | 7.6 | 2.6 | 1.6 | 2.0 | 3.0 | 9.8 | **10.5** |

- used "standard" template for each relation
- Suprising: $RE_o$ baseline has seen at least one sentence per fact
- But: BERT prob. has sentence in training data (trained on Wikipedia)

# 3.5. Results

**Results with p@1:**

| Corpus | Relation | Statistics | | Baselines | | KB | | LM | | | | | |
|--------|----------|------------|------|-----------|------|------|--------|-----|------|-----|-----|------|------|
| | | #Facts | #Rel | Freq | DrQA | $RE_n$ | $RE_o$ | Fs | Txl | Eb | E5B | Bb | Bl |
| T-REx | 1-1 | 937 | 2 | 1.78 | - | 0.6 | 10.0 | 17.0 | 36.5 | 10.1 | 13.1 | 68.0 | **74.5** |
| | N-1 | 20006 | 23 | 23.85 | - | 5.4 | **33.8** | 6.1 | 18.0 | 3.6 | 6.5 | 32.4 | 34.2 |
| | N-M | 13096 | 16 | 21.95 | - | 7.7 | **36.7** | 12.0 | 16.5 | 5.7 | 7.4 | 24.7 | 24.3 |
| | Total | 34039 | 41 | 22.03 | - | 6.1 | **33.8** | 8.9 | 18.3 | 4.7 | 7.1 | 31.1 | 32.3 |

- ▶ BERT way better than $RE_o$ for 1-1 relations (eg. *capital_of*)
- ▶ results N-1 BERT_large $\approx$ results $RE_o$
- ▶ $RE_o$ unrivaled for N-M relations
- ▶ general results BERT_large $\approx$ general results $RE_o$

# 3.5. Results

**Conclusion:**

- ▶ could be complicated to improve the performance of RE by providing additional data
- ▶ RE performs similar to BERT_large in general and doesn't need complicated pipelines
- ▶ LMs could become an useful alternative for traditionally extracted knowledge bases
- ▶ in the future: with LLMs that are trained on even more data, they might be able to replace knowledge bases

# 4. SMALL LANGUAGE MODELS ARE ALSO FEW-SHOT LEARNERS

# 4.1 General idea

- ▶ Paper by Schick and Schütze (LMU) published in June 2021
- ▶ GPT-3 achieves great results on SuperGLUE tasks by priming
- ▶ Two problems:
    - ▶ GPT-3 is a LLM and has a large carbon footprint
    - ▶ Examples are limited to a few due to size of the context window
- ▶ Solution: Use Pattern-Exploiting Training (PET)

Schick and Schütze 2020

# 4.2 Pattern-Exploiting Training (PET)

- ▶ PET combines the idea of reformulating tasks as cloze questions with regular gradient-based finetuning
- ▶ PET additionally requires unlabeled data, unlabeled data is much easier to obtain than labeled examples for many real-world applications.
- ▶ Crucially, PET only works when the answers to be predicted by the LM correspond to a single token in its vocabulary; this is a severe limitation as many tasks cannot easily be worded that way.
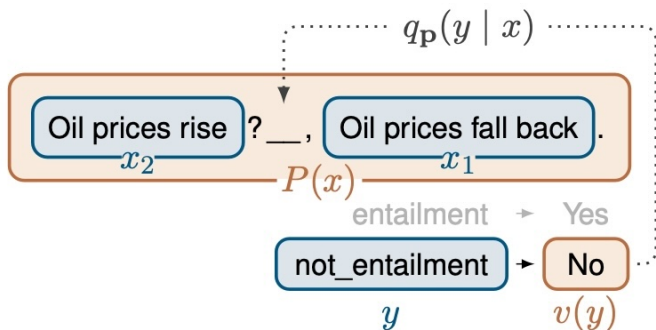
## 4.3 Pattern-Verbalizer Pairs

Each PVP $p = (P, v)$ consists of:

- A pattern $P : X \rightarrow T^*$ maps inputs to cloze questions containing a single mask. ($T^*$: set of all token sequences)
- A verbalizer $v : Y \rightarrow T$ maps each output to a single token representing its task-specific meaning in the pattern. ($T$: vocabulary)

## 4.3 Pattern-Verbalizer Pairs



Application of a PVP $p = (P, v)$ for recognizing textual entailment:

- ▶ An input $x = (x_1, x_2)$ is converted into a cloze question $P(x)$.
- ▶ $q_p(y|x)$ for each $y$ is derived from the probability of $v(y)$ being a plausible choice for the masked position.

**iPET**: Iterative variant of PET for improved learning through iterations

**Process:**

- ▶ **Initial Training:** Train an ensemble of MLMs using PET
- ▶ **Generate New Training Set:** For each model $M_i$:
  - ▶ Select a random subset of other models
  - ▶ Generate a new training set $T_i$
  - ▶ Assign labels to unlabeled examples based on the subset's most confident predictions
- ▶ **Retrain Models:** Retrain each $M_i$ on $T_i$
- ▶ **Iterate:** Repeat the process, increasing the size of $T_i$ by a constant factor in each iteration

**Benefits:**

- ▶ Enhanced Learning: Models learn from different patterns and data points
- ▶ Progressive Improvement: Gradual increase in training data size leads to better model performance

# 4.4 GLUE and SuperGLUE

**GLUE**

- Multi-task benchmark platform for Natural Language Understanding (NLU) tasks
- Consists of 9 tasks
    - CoLa: Corpus of Linguistic Acceptability
    - QQP: Quora Question Pairs
- Performance of LM's surpassed level of non-expert humans quickly

Wang et al. (2019b)

# 4.4 GLUE and SuperGLUE

**GLUE Leaderboard**

# 4.4 GLUE and SuperGLUE

**SuperGLUE**

- ▶ New and improved benchmark with more difficult and more diverse tasks, total of 8
- ▶ Retained the two hardest tasks of GLUE: Winograd Schema Challenge and Recognizing Textual Entailment
- ▶ New tasks include CommitmentBank, Words in Context and Reading Comprehension with Commonsense Reasoning

Wang et al. (2019a)

# 4.4 GLUE and SuperGLUE
## SuperGLUE Leaderboard



Leaderboard Version: **2.0**

| Rank | Name | Model | URL | Score | BoolQ | CB | COPA | MultiRC | ReCoRD | RTE | WiC | WSC | AX-b | AX-g |
|------|------|-------|-----|-------|-------|-----|------|---------|--------|-----|-----|-----|------|------|
| 1 | JDExplore d-team | Vega v2 | ☑ | 91.3 | 90.5 | 98.6/99.2 | 99.4 | 88.2/62.4 | 94.4/93.9 | 96.0 | 77.4 | 98.6 | -0.4 | 100.0/50.0 |
| 2 | Liam Fedus | ST-MoE-32B | ☑ | 91.2 | 92.4 | 96.9/98.0 | 99.2 | 89.6/65.8 | 95.1/94.4 | 93.5 | 77.7 | 96.6 | 72.3 | 96.1/94.1 |
| 3 | Microsoft Alexander v-team | Turing NLR v5 | ☑ | 90.9 | 92.0 | 95.9/97.6 | 98.2 | 88.4/63.0 | 96.4/95.9 | 94.1 | 77.1 | 97.3 | 67.8 | 93.3/95.5 |
| 4 | ERNIE Team - Baidu | ERNIE 3.0 | ☑ | 90.6 | 91.0 | 98.6/99.2 | 97.4 | 88.6/63.2 | 94.7/94.2 | 92.6 | 77.4 | 97.3 | 68.6 | 92.7/94.7 |
| 5 | Yi Tay | PaLM 540B | ☑ | 90.4 | 91.9 | 94.4/96.0 | 99.0 | 88.7/63.6 | 94.2/93.3 | 94.1 | 77.4 | 95.9 | 72.9 | 95.5/90.4 |
| 6 | Zirui Wang | T5 + UDG, Single Model (Google Brain) | ☑ | 90.4 | 91.4 | 95.8/97.6 | 98.0 | 88.3/63.0 | 94.2/93.5 | 93.0 | 77.9 | 96.6 | 69.1 | 92.7/91.9 |
| 7 | DeBERTa Team - Microsoft | DeBERTa / TuringNLRv4 | ☑ | 90.3 | 90.4 | 95.7/97.6 | 98.4 | 88.2/63.7 | 94.5/94.1 | 93.2 | 77.5 | 95.9 | 66.7 | 93.3/93.8 |
| | SuperGLUE Human Baselines | SuperGLUE Human Baselines | | 89.8 | 89.0 | 95.8/98.9 | 100.0 | 81.8/51.9 | 91.7/91.3 | 93.6 | 80.0 | 100.0 | 76.6 | 99.3/99.7 |
| 8 | T5 Team - Google | T5 | ☑ | 89.3 | 91.2 | 93.9/96.8 | 94.8 | 88.1/63.3 | 94.1/93.4 | 92.5 | 76.9 | 93.8 | 65.6 | 92.7/91.9 |
| 10 | SPoT Team - Google | Frozen T5 1.1 + SPoT | ☑ | 89.2 | 91.1 | 95.8/97.6 | 91.6 | 87.9/61.9 | 93.3/92.4 | 92.9 | 75.8 | 93.8 | 66.9 | 83.1/82.5 |
| 11 | Huawei Noah's Ark Lab | NEZHA-Plus | ☑ | 86.7 | 87.8 | 94.4/96.0 | 93.6 | 84.6/55.1 | 90.1/89.6 | 89.1 | 74.6 | 93.2 | 58.0 | 87.1/74.4 |
| 12 | Alibaba PAI&ICBU | PAI Albert | ☑ | 86.1 | 88.1 | 92.4/96.4 | 91.8 | 84.6/54.7 | 89.0/88.3 | 88.8 | 74.1 | 93.2 | 73.6 | 98.3/99.2 |
| 13 | Infosys : DAWN : AI Research | RoBERTa-iCETS | | 85.7 | | | 91.2 | 86.4/58.2 | 89.9/89.3 | 89.9 | 72.9 | 89.0 | 81.8 | 88.8/81.5 |
| 14 | Tencent Jarvis Lab | RoBERTa (ensemble) | | 85.9 | 88.2 | 92.5/95.6 | 90.8 | 84.4/63.4 | 91.5/91.0 | 87.9 | 74.1 | 91.8 | 57.6 | 89.3/75.6 |
| 15 | Zhuiyi Technology | RoBERTa-mtl-adv | | 85.7 | 87.1 | 92.4/95.6 | 91.2 | 85.1/54.3 | 91.7/91.3 | 88.1 | 72.1 | 91.8 | 58.5 | 91.0/78.1 |
| 16 | Facebook AI | RoBERTa | ☑ | 84.6 | 87.1 | 90.5/95.2 | 90.6 | 84.4/52.5 | 90.6/90.0 | 88.2 | 69.9 | 89.0 | 57.9 | 91.0/78.1 |
| 17 | Anuar Sharafudinov | AiLabs Team. Transformers | | 82.6 | 88.1 | 91.6/94.8 | 86.8 | 85.1/54.7 | 82.8/79.8 | 88.9 | 74.1 | 78.8 | 100.0 | 100.0/100.0 |
| 18 | Ying Luo | FSL++(ALBERT)-Few-Shot(32 Examples) | | 77.7 | 81.1 | 87.8/92.0 | 87.0 | 77.3/38.4 | 81.9/81.1 | 75.1 | 60.5 | 86.4 | 33.9 | 94.4/63.5 |
| 19 | Rathin Bector | Text to Text PETL | | 76.5 | 76.3 | 86.9/92.4 | 80.2 | 80.4/44.8 | 82.2/81.3 | 78.1 | 67.5 | 74.0 | 38.1 | 97.2/83.7 |
| 20 | CASIA | INSTALL(ALBERT)-few-shot | ☑ | 76.6 | 78.4 | 85.9/92.0 | 85.6 | 75.9/35.1 | 84.3/83.5 | 74.9 | 60.9 | 84.9 | -0.4 | 100.0/50.0 |
| 21 | Rakesh Radhakrishnan Menon | ADAPET (ALBERT) - few-shot | | 76.0 | 80.0 | 82.3/92.0 | 85.4 | 76.2/35.7 | 86.1/85.5 | 75.0 | 53.5 | 85.6 | -0.4 | 100.0/50.0 |
| 22 | Timo Schick | iPET (ALBERT) - Few-Shot (32 Examples) | ☑ | 75.4 | 81.2 | 79.9/88.8 | 90.8 | 74.1/31.7 | 85.9/85.4 | 70.8 | 49.3 | 88.4 | 38.2 | 97.8/57.9 |
| 23 | Adrian de Wynter | Bort (Alexa AI) | ☑ | 74.1 | 83.7 | 81.9/86.4 | 89.6 | 83.7/54.1 | 49.8/49.0 | 81.2 | 70.1 | 65.8 | 48.0 | 96.1/61.5 |
| 24 | IBM Research AI | BERT-mtl | | 73.5 | 84.8 | 89.6/94.0 | 73.8 | 73.2/30.5 | 74.6/74.0 | 84.1 | 66.2 | 61.0 | 29.6 | 97.8/57.3 |
| 25 | Ben Mann | GPT-3 few-shot - OpenAI | ☑ | 71.8 | 76.4 | 52.0/75.6 | 92.0 | 75.4/30.5 | 91.1/90.2 | 69.0 | 49.4 | 80.1 | 21.1 | 90.4/55.3 |
| 26 | SuperGLUE Baselines | BERT++ | ☑ | 71.5 | 79.0 | 84.8/90.4 | 73.8 | 70.0/24.1 | 72.0/71.3 | 79.0 | 69.6 | 64.4 | 38.0 | 99.4/51.4 |

# 4.4 GLUE and SuperGLUE

Exercise 3, see Colab

# 4.5 Results

| | Model | Params (M) | BoolQ Acc. | CB Acc. / F1 | COPA Acc. | RTE Acc. | WiC Acc. | WSC Acc. | MultiRC EM / F1a | ReCoRD Acc. / F1 | Avg – |
|---|---|---|---|---|---|---|---|---|---|---|---|
| dev | GPT-3 Small | 125 | 43.1 | 42.9 / 26.1 | 67.0 | 52.3 | 49.8 | 58.7 | 6.1 / 45.0 | 69.8 / 70.7 | 50.1 |
| | GPT-3 Med | 350 | 60.6 | 58.9 / 40.4 | 64.0 | 48.4 | 55.0 | 60.6 | 11.8 / 55.9 | 77.2 / 77.9 | 56.2 |
| | GPT-3 Large | 760 | 62.0 | 53.6 / 32.6 | 72.0 | 46.9 | 53.0 | 54.8 | 16.8 / 64.2 | 81.3 / 82.1 | 56.8 |
| | GPT-3 XL | 1,300 | 64.1 | 69.6 / 48.3 | 77.0 | 50.9 | 53.0 | 49.0 | 20.8 / 65.4 | 83.1 / 84.0 | 60.0 |
| | GPT-3 2.7B | 2,700 | 70.3 | 67.9 / 45.7 | 83.0 | 56.3 | 51.6 | 62.5 | 24.7 / 69.5 | 86.6 / 87.5 | 64.3 |
| | GPT-3 6.7B | 6,700 | 70.0 | 60.7 / 44.6 | 83.0 | 49.5 | 53.1 | 67.3 | 23.8 / 66.4 | 87.9 / 88.8 | 63.6 |
| | GPT-3 13B | 13,000 | 70.2 | 66.1 / 46.0 | 86.0 | 60.6 | 51.1 | 75.0 | 25.0 / 69.3 | 88.9 / 89.8 | 66.9 |
| | GPT-3 | 175,000 | 77.5 | 82.1 / 57.2 | 92.0 | 72.9 | **55.3** | 75.0 | 32.5 / 74.8 | **89.0 / 90.1** | 73.2 |
| | PET | 223 | 79.4 | 85.1 / 59.4 | **95.0** | 69.8 | 52.4 | **80.1** | **37.9 / 77.3** | 86.0 / 86.5 | 74.1 |
| | iPET | 223 | **80.6** | **92.9 / 92.4** | **95.0** | **74.0** | 52.2 | **80.1** | 33.0 / 74.0 | 86.0 / 86.5 | **76.8** |
| test | GPT-3 | 175,000 | 76.4 | 75.6 / 52.0 | **92.0** | 69.0 | 49.4 | 80.1 | 30.5 / 75.4 | **90.2 / 91.1** | 71.8 |
| | PET | 223 | 79.1 | 87.2 / 60.2 | 90.8 | 67.2 | **50.7** | **88.4** | **36.4 / 76.6** | 85.4 / 85.9 | 74.0 |
| | iPET | 223 | **81.2** | **88.8 / 79.9** | 90.8 | **70.8** | 49.3 | **88.4** | 31.7 / 74.1 | 85.4 / 85.9 | **75.4** |
| | SotA | 11,000 | *91.2* | *93.9 / 96.8* | *94.8* | *92.5* | *76.9* | *93.8* | *88.1 / 63.3* | *94.1 / 93.4* | *89.3* |

▶ Better than Chat GPT-3 on most of the tasks, but not SOTA

# 4.6 Analysis of the results

**What can influence the performance?**

# 4.6 Analysis of the Results

**What can influence the performance?**

- ▶ Patterns/Templates
- ▶ Labeled and unlabeled data usage
- ▶ Model type
- ▶ Training examples

# 5. SUMMARY

# Summary

- Very intuitively and easy to understand

- Performance can vary greatly depending on multiple factors

- **But:** It can be very time and cost intensive.

- **Solution:** Automated Template Learning (next week :))

# Questions and Discussion

# Thank You for Your attention!

# Literature I

Kojima, T., Gu, S. S., Reid, M., Matsuo, Y., and Iwasawa, Y. (2022). Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213.

Liu, J., Liu, A., Lu, X., Welleck, S., West, P., Bras, R. L., Choi, Y., and Hajishirzi, H. (2021). Generated knowledge prompting for commonsense reasoning. *arXiv preprint arXiv:2110.08387*.

Petroni, F., Rocktäschel, T., Lewis, P., Bakhtin, A., Wu, Y., Miller, A. H., and Riedel, S. (2019). Language models as knowledge bases? *arXiv preprint arXiv:1909.01066*.

Schick, T. and Schütze, H. (2020). It's not just size that matters: Small language models are also few-shot learners. *arXiv preprint arXiv:2009.07118*.

# Literature II

Wang, A., Pruksachatkun, Y., Nangia, N., Singh, A., Michael, J., Hill, F., Levy, O., and Bowman, S. R. (2019a). SuperGLUE: A stickier benchmark for general-purpose language understanding systems. *arXiv preprint 1905.00537*.

Wang, A., Singh, A., Michael, J., Hill, F., Levy, O., and Bowman, S. R. (2019b). GLUE: A multi-task benchmark and analysis platform for natural language understanding. In the Proceedings of ICLR.

Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., Le, Q. V., Zhou, D., et al. (2022). Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.

Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T., Cao, Y., and Narasimhan, K. (2024). Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36.

# Literature III

Zhang, Z., Zhang, A., Li, M., and Smola, A. (2022). Automatic chain of thought prompting in large language models. *arXiv preprint arXiv:2210.03493*.