

# Knowledge and Story Comprehension

Eduardo Badillo\*

Geronimo Walker†

Svein Gonzalez‡

March 25, 2024

## Abstract

It has been demonstrated that pre-trained language models have a good performance on numeric, and common-sense reasoning tasks with little or no fine-tuning needed. This performance is further enhanced when more context or knowledge of the input is provided to the model at the moment of inference; for it enables the model to leverage the knowledge it already contains from its pretraining and use it more effectively on downstream tasks. In this paper we study the impact of incorporating additional knowledge, with different structures and retrieval techniques, on tasks that require both reasoning and comprehension (ClozeStory Dataset) to determine a correct ending of a story. Our code is available at Github.

## 1 Introduction

There are two main ways in which task performance can be enhanced in pre-trained language models, by incorporating knowledge and by fine-tuning. The latter has shown that larger and larger models with fine-tuning diminish the need to integrate auxiliary knowledge (Khashabi et al., 2020; Lourie et al., 2021). Yet, incorporating knowledge on top of large-scale models (Liu et al. 2022) has shown a marked improvement on performance when no fine-tuning is involved. Suggesting that the models are able to exploit more effectively their pre-training knowledge when additional task-relevant information is provided.

Obtaining and encoding relevant, contextual information isn't trivial. There might not be enough quality data for the task, and the way it is retrieved from the external knowledge structure and paired with the relevant keywords or entities in each sample must also be thought of. The methods we will be using are knowledge graphs (Aglionbi et al. 2022) and knowledge prompting (Liu et al. 2022). We will add the relevant knowledge contained in each on top of pre-trained models, without fine-tuning to test whether their inclusion amounts to an increased performance on the story comprehension task. These approaches have been used on common-sense, numeric reasoning tasks where the query or question is configured in a syllogistic manner. But its usefulness on comprehension datasets, where samples have no

straightforward outcomes, has yet to be tested. As baseline we will be using the models without knowledge structures and no fine-tuning, and during the experimentation phase we will train the models with and without knowledge to test whether their inclusion increases performance.

## 2 Dataset and Model

Each sample in the Story Cloze dataset comprises four distinct sentences alongside a correct ending sentence and an incorrect ending sentence. The dataset is divided into two training sets, one for 2016 (45,495 stories) and one for 2017 (52,664 stories). These datasets did not include an incorrect ending like the validation set does. To obtain incorrect endings for each sample, in order to fulfill our binary classification approach, we prompted a gpt model through an open ai api to generate them for us. We asked the gpt 3.5 turbo model to generate appropriate incorrect endings, providing several examples of how we did it manually. This generative process is also how we obtained relevant knowledge during the knowledge prompt experiments (section 3).

To process the dataset for training and inference our procedure involved concatenating each set of four input sentences into a cohesive paragraph, followed by duplicating each paragraph and appending both the correct and incorrect endings. This resulted in

---

\*UC Berkeley, Email: eduardo.badillo@berkeley.edu

†UC Berkeley, Email: geronimowalker@ischool.berkeley.edu

‡UC Berkeley, Email: sggonzal@ischool.berkeley.edu

two stories—one with a valid conclusion, and the other an invalid conclusion. We conducted our training with Bert for Sequence Classification and Bert for Multiple Choice, both equipped with a standard 12-layer architecture and a classification head.

### 3 Prompt-based Knowledge

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maece-nas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

### 4 Triples

Triples are dependency representations of sentences. For example, take the sentence “David noticed he had put on a lot of weight recently”, using Spacy, we can generate triples, most commonly Subject-Verb-Object (SVOs), to represent the sentence with the most relevant context and information. Drawing inspiration from Aglinoby and Teugel (2022), where graph representations of possible answers to questions are used to determine the highest ranking answer, our model will use extracted triples from the sentences comprising a story and compute a cosine similarity between the triples and two possible endings. Extracted triples are initially captured by a simple ap-

proach, speech tagging identification.

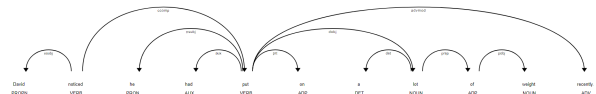


Figure 1: Spacy Dependency Example

Using SPACY, we parse through all our sentences capturing words that are labeled “ROOT” (for verb), “nsubj” (for subject), or “dobj” / “pobj” (for object) into our SVO triples. It is important to note that “ROOT” can be a successful tag for verbs, but isn’t a guarantee. As noted in Figure 1, certain words such as “put” are stronger central nodes than others since they share more dependencies. This will be important to improve upon to optimally capture significant SVOs.

After generating SVOs, one for each story sentence, we use a bert model to generate cosine similarities for our different endings. We do so by tokenizing our SVOs then running them through a bert model to obtain a mean pooled output for each triple. The story endings are treated the same. Next we compute a cosine similarity between all the story sentences and the two possible endings (Figure 2). The initial attempt to use triples hasn’t yielded remarkable results, only presenting a 50 percent accuracy.

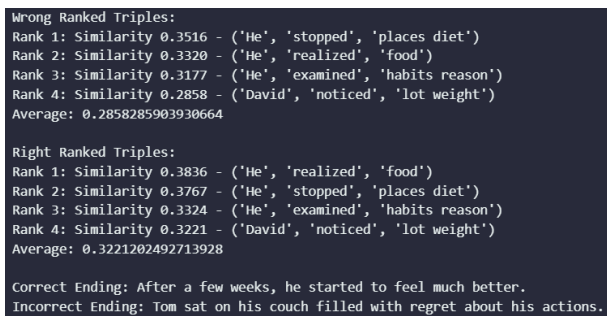


Figure 2: Consine Similarity

However, the next models will investigate triples with better selection and ordering of the triples. Currently, all triples are weighed evenly when the cosines are averaged (Figure 2). The next model iteration will note the ending triples of the story more using a BertSeqtoSeq model.