

# Generating and Evaluating Long Story Summaries with Knowledge Graphs

Anonymous ACL submission

## Abstract

Summarizing long stories is a challenging task due to their narrative complexity and the context length limits of language models. We propose a method that integrates knowledge graph retrieval with the summarization process to provide global context. We construct a knowledge graph containing entity descriptions and relations from the entire story, then retrieve relevant information from it to aid summary generation. Additionally, we propose a novel metric, KGScore, which evaluates summaries by comparing the similarity of knowledge graphs extracted from generated and reference summaries. Experimental results demonstrate that our knowledge graph retrieval method outperforms the baselines in terms of our KGScore metric and that KGScore is a reliable measure of factual consistency.

## 1 Introduction

Language models based on the transformer architecture (Vaswani et al., 2017) have been successfully trained to summarize short texts (Liu and Lapata, 2019; Zhang et al., 2020a, 2022a). However, understanding and summarizing longer documents, such as entire books, remains a challenge, largely due to the context length limit imposed by the quadratic complexity of the attention mechanism (Cao and Wang, 2023). Furthermore, stories pose unique problems as their summaries are highly abstractive in nature and require the navigation of a mix of narration and dialogue, with complex dependencies interspersed throughout the text (Kryściński et al., 2021). The difficulty is amplified by the fact that narrative texts often employ the technique of “show, don’t tell”: instead of explicit descriptions or stating of facts, the author relies on implicit information conveyed through dialogue or character actions. As a result, long stories, with their dual hurdles of extensive length and narrative intricacy, present a particularly daunting task for

summarization.

Previous research on the topic ranges from divide-and-conquer strategies that produce a summary of summaries from split-up story segments (Wu et al., 2021; Kashyap, 2022), to approaches that generate an abstractive summary of extractive samples (Hardy et al., 2022). The ability of these methods to produce factually consistent summaries is limited by the lack of a global context.

We propose the use of knowledge graphs to address this issue. Knowledge graphs represent descriptions of entities and relationships between them in a structured form and have been used to successfully improve performance on a variety of natural language generation tasks (Fan et al., 2019; Andrus et al., 2022). We frame the problem as a chapter summarization task, where the model is given the chapter text and a knowledge graph of the entire book. The model retrieves information from the knowledge graph to augment its understanding of the story and generate a chapter summary.

To generate the knowledge graph, we split the book text into small chunks and instruct a large language model to identify named entities in the text, which become the nodes of the graph. We then extract the graph edges by prompting the model to generate entity descriptions and relations. We additionally follow a series of steps to ensure that the information in the knowledge graph is accurate and relevant. During summarization, the knowledge graph edges are ranked based on their semantic similarity to a set of keywords, then retrieved and prepended to the chapter text in linearized form to be given to the summarization model.

We also propose a new metric for evaluating generated summaries, which we name KGScore. It is designed to address the limitations of existing metrics in evaluating factual consistency. The metric computes precision, recall, and F1 scores based on the cosine similarity of knowledge graph edge em-

041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081

082 beddings extracted from generated and reference  
083 summaries.

084 We evaluate the effectiveness of our approach  
085 through experiments on the BookSum Chapters  
086 dataset (Kryściński et al., 2021), which contains  
087 chapter texts and their summaries. We find that our  
088 proposed method outperforms the baseline in terms  
089 of our KGScore metric. Through an additional  
090 experiment, we also verify that KGScore is a valid  
091 measure of factuality.

## 092 2 Related Work

### 093 2.1 Long Document Summarization

094 Summarizing long documents using transformer-  
095 based language models is challenging because of  
096 the quadratic computational and memory require-  
097 ments. Existing approaches to overcoming this  
098 problem can be broadly classified into three cate-  
099 gories: divide and conquer, efficient attention, and  
100 extractive-abstractive summarization.

101 The divide-and-conquer strategy breaks down  
102 the task of summarizing a long document into  
103 smaller tasks of summarizing short sections of the  
104 document that can fit into a language model’s con-  
105 text. Summaries for each section are combined to  
106 produce the summary for the full document (Gid-  
107 iotis and Tsoumakas, 2020; Zhang et al., 2022b).  
108 This segmentation can result in reduced coherence  
109 due to a lack of global context. To address this  
110 problem, Cao and Wang (2023) introduce an exter-  
111 nal memory mechanism. Pang et al. (2023) propose  
112 a variant form of divide and conquer, where they  
113 combine a bottom-up pass using local self-attention  
114 on chunks of text with a top-down correction step  
115 to capture long-range dependencies.

116 There have also been efforts to improve the at-  
117 tention mechanism itself instead of working around  
118 its limitations, reducing the time and memory com-  
119 plexities to subquadratic levels for long sequences.  
120 This makes it feasible to fit long inputs into the  
121 model (Huang et al., 2021). On top of this, Phang  
122 et al. (2022) incorporate a pretraining step on long  
123 texts to further improve performance.

124 Extractive-abstractive summarization is a group  
125 of methods consisting of two steps: extracting rele-  
126 vant parts of the input document, then using a lan-  
127 guage model to generate an abstractive summary  
128 from the extracted snippets (Pilault et al., 2020;  
129 Zhao et al., 2020). Large language models such as  
130 OpenAI’s ChatGPT have been used for the abstrac-  
131 tive step (Lu et al., 2023).

Solutions for the more specific problem of sum-  
marizing long stories have also been explored. Wu  
et al. (2021) and Kashyap (2022) use techniques  
based on divide and conquer, while Hardy et al.  
(2022) propose an extractive-abstractive approach.  
Our method can be considered a form of divide and  
conquer; we additionally incorporate knowledge  
graphs as a way of providing global context.

### 2.2 Knowledge Graphs for Text Generation

Knowledge graphs can be used with text-generation  
tasks to supplement models with additional infor-  
mation. Here, we discuss knowledge graphs ex-  
tracted on the fly from input documents. Prior  
research, such as ASGARD (Huang et al., 2020),  
typically employs graph attention to encode the  
graph data (Zhu et al., 2021; Chen et al., 2023).  
They focus on tasks that involve synthesizing infor-  
mation from multiple documents (Fan et al., 2019),  
including the summarization of multiple news arti-  
cles (Lakshika et al., 2020).

The application of these methods has largely  
been confined to factual content, such as news arti-  
cles or academic papers, and not stories. While the  
Stanford OpenIE system (Angeli et al., 2015) is a  
popular choice for extracting relational data from  
documents to construct knowledge graphs, rule-  
based systems like this often struggle to generate  
meaningful knowledge graphs from narrative texts.  
Andrus et al. (2022) do use the OpenIE system for  
story completion and question answering tasks, but  
integrate it with GPT-3 (Brown et al., 2020). In our  
approach, we forgo the OpenIE system entirely, di-  
rectly using a large language model for knowledge  
graph construction.

### 2.3 Metrics for Summarization Evaluation

One of the most commonly used metrics for eval-  
uating summaries is ROUGE (Lin, 2004), which  
measures the overlap of n-grams between gener-  
ated and reference summaries. BERTScore (Zhang  
et al., 2020b) is another widely-used metric and  
involves computing the similarity of contextual em-  
beddings. Many of these existing metrics have been  
shown to correlate poorly with human judgments  
of quality (Novikova et al., 2017), especially for  
assessing factuality (Maynez et al., 2020).

Methods have been proposed to evaluate the fac-  
tual consistency of generated summaries (Kryscin-  
ski et al., 2020; Xie et al., 2021), including QuestE-  
val (Scialom et al., 2021), which uses question gen-  
eration and answering for this purpose. However,

these methods are often impractical to use with long documents as they involve the use of models with limited context sizes. Our proposed metric bypasses this limitation by adopting a source-free approach that compares the generated summary with the reference summary instead of the original document.

Some recent metrics make use of large language models such as ChatGPT and GPT-4 (OpenAI, 2023), guiding a model through prompts to produce evaluations (Gao et al., 2023; Liu et al., 2023). While our metric also incorporates a large language model as part of the process, we do not rely on it fully; its use is limited to the knowledge graph extraction step.

Metrics specifically targeting summaries of long documents, including long stories, have also been proposed. LongDocFACTScore (Bishop et al., 2023) is a framework that enables the extension of any preexisting metric to accommodate long documents. SNaC (Goyal et al., 2022) and BoookScore (Chang et al., 2023) are both reference-free and source-free metrics that identify errors by focusing exclusively on the content of the generated summaries. Our metric is source-free but not reference-free; it compares predicted summaries with reference summaries.

## 3 Methods

### 3.1 Chapter Summarization Task

The task is formulated as follows. Given the full text of a book  $\mathcal{B} = \{\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_n\}$  consisting of  $n$  chapters, and a chapter index  $k$  ( $1 \leq k \leq n$ ), the model  $\mathbb{M}$  must learn to generate the chapter summary  $\mathcal{Y}_k$  corresponding to the chapter text  $\mathcal{C}_k$ :

$$\mathbb{M} : (\mathcal{B}, k) \rightarrow \mathcal{Y}_k \quad (1 \leq k \leq n) \quad (1)$$

The model may utilize information from all parts of the book as needed. However, we assume that the full book text is too long to be given to the model as input in its original form, while individual chapters are not. This task can be considered the first half of a divide-and-conquer approach, where the chapters are segments of the book that can fit into the model’s context size. The second half of the process would be to combine the generated summaries for each chapter into a summary for the full book, but we do not focus on that part here.

### 3.2 Summarization with Knowledge Graph Retrieval

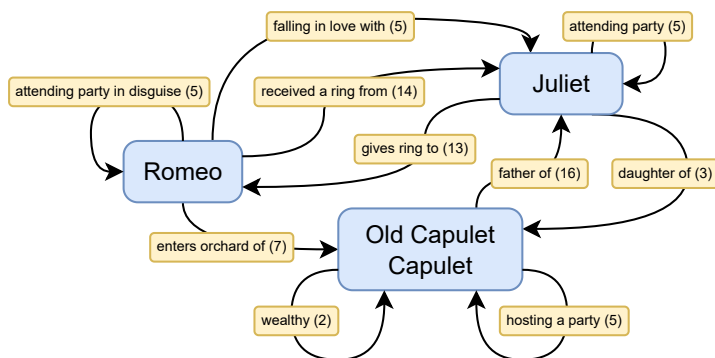
In our method, we generate a knowledge graph containing information from the entire book text. Each node in the graph represents a named entity in the story, such as a character, organization, or location, and each directed edge represents a <subject, predicate, object> triple (e.g., <Romeo, is in love with, Juliet>), where the source node and target node correspond to the subject and object, respectively. As an exception, in self-loops (i.e., edges with the same source and target node), the object is ignored and the edge represents a <subject, predicate> pair (an entity description or an action with no object; e.g., <Romeo, is in love>). Multiple edges with different predicates can exist between a single (subject, object) pair. Figure 1a shows an example of a generated knowledge graph.

During training and inference, we retrieve information from the knowledge graph and provide it to the summarization model along with the full chapter text. This additional information acts as global context that is lacking when only the chapter text is provided. For example, if a character that was introduced in a previous chapter reappears in the current chapter, it could be difficult for the model to determine the identity of the character by only examining the current chapter. Information from previous chapters would be helpful global context that helps the model “remember” who the character was, and any details about the character in the generated summary would be more likely to be correct.

### 3.3 Book Knowledge Graph Generation

**Knowledge Extraction** We use a large language model to extract the knowledge graph nodes and edges from the book text. We split the text at paragraph boundaries so that each segment, when inserted into the prompt, fits into the model’s context size. The prompt begins with an instruction to identify named entities and knowledge graph edges. This is followed by an example containing a story excerpt and lists of corresponding entities and edges. Finally, the book segment of interest is given as the task for the model. The complete prompt can be found in Appendix E. For our experiments, we use OpenAI’s gpt-3.5-turbo-0613 model.

**Names Graph** Before building the knowledge graph, we parse the named entities from the model



(a) Knowledge graph.

```

<subject>Romeo
  <predicate>attending party in disguise
    <object>Juliet
  <predicate>falling in love with
    <predicate>received a ring from
      <object>Capulet
    <predicate>enters orchard of
      <subject>Juliet
    <predicate>attending party
      <object>Romeo
    <predicate>gives ring to
      <object>Capulet
    <predicate>daughter of
      <subject>Capulet
    <predicate>hosting a party
    <predicate>wealthy
      <object>Juliet
    <predicate>father of

```

(b) Linearized knowledge graph.

Figure 1: (a) Part of a knowledge graph generated from William Shakespeare’s *Romeo and Juliet*, with three named entities. “Old Capulet/Capulet” is a single entity with two names. The numbers in parentheses following the predicates are chapter numbers. (b) The same graph in linearized form.

278 responses to generate a *names graph*. The model returns  
 279 a list of names (aliases or name variations) for  
 280 each entity that appears in a given story segment.  
 281 The purpose of the names graph is to consolidate  
 282 this information and keep track of names that refer  
 283 to the same entity over the span of the entire  
 284 book. For example, names *A* and *B* could be identified  
 285 as aliases of an entity in one section of the  
 286 book, while the same is done for names *B* and *C* in  
 287 a different section; the names graph will indicate  
 288 that names *A*, *B*, and *C* are all names of the same  
 289 entity. A node is created for each distinct name,  
 290 and undirected edges are created between nodes  
 291 that represent names of the same entity. By repeating  
 292 this process, we obtain a graph of all identified  
 293 names in the book. If two nodes are connected (*i.e.*,  
 294 there exists a path of edges between them), their  
 295 names refer to the same entity.

296 **Knowledge Graph Initialization** We initialize the  
 297 knowledge graph by giving each name in the names  
 298 graph its own node in the knowledge graph, regardless  
 299 of whether they belong to the same entity. We  
 300 parse the knowledge graph edges from the model  
 301 responses and perform a processing step to ensure  
 302 that both the subject and object are named entities.  
 303 If there is no object, we repeat the subject as the  
 304 object to create a self-loop. Each processed edge  
 305 is added to the knowledge graph as a directed edge  
 306 from the subject node to the object node.

307 **Node Merging** At this point, the knowledge graph  
 308 may contain multiple nodes representing the same

entity under different names. In this step, we merge  
 309 these nodes into one node to obtain a graph with  
 310 one node per entity, with each node having a list  
 311 of names for its entity. For each edge in the names  
 312 graph (connecting two names that refer to a single  
 313 entity), we identify the nodes in the knowledge  
 314 graph that contain the names that the edge connects.  
 315 If the names belong to two different nodes, we  
 316 merge the nodes by combining their name lists and  
 317 transferring the edges of one node to the other.  
 318

319 One problem is that the names graph occasionally  
 320 contains incorrect connections (*i.e.*, edges between  
 321 names of distinct entities). We employ two  
 322 heuristics to prevent the merging of entity nodes  
 323 when this occurs. First, if the two nodes have an  
 324 edge between them (representing a  $\langle \text{subject}, \text{predicate}, \text{object} \rangle$   
 325 triple), we do not merge the nodes. It is unlikely  
 326 that a triple would have the same entity as the  
 327 subject and object while also referring to it using  
 328 different names. Second, we check the degree (the  
 329 number of edges entering or leaving the node) of  
 330 each of the two nodes, excluding any self-loops  
 331 from the count. We do not merge the nodes if  
 332 both edge counts exceed a maximum threshold  
 333 value. The reasoning is that if both nodes have a  
 334 high degree, they are both important entities and  
 335 are likely to be different. In our experiments, we  
 336 set the threshold to 3.

337 **Node Removal** The final step is to remove nodes  
 338 in the knowledge graph with a degree that is less  
 339 than a minimum threshold value (excluding self-loops).  
 340 This ensures that only the most relevant



341 entities remain in the graph. This can also eliminate  
 342 erroneously identified entities. The removal of a  
 343 node and its edges can affect the edge counts of  
 344 other nodes, so we repeat the process of counting  
 345 and removing until no more nodes are removed.  
 346 We use a minimum degree of 2 in our experiments.

### 347 3.4 Knowledge Graph Edge Retrieval

348 **Edge Ranking** Instead of providing the entire  
 349 knowledge graph to the model with the chapter  
 350 text, we retrieve a subset of the graph that would  
 351 be the most helpful for generating a summary. We  
 352 experimentally select a set of keywords and corre-  
 353 sponding weights to score and rank the knowledge  
 354 graph edges. Each keyword is designed to focus  
 355 on an important narrative aspect, such as charac-  
 356 ter relationships (with the “relation” keyword) and  
 357 events (“happen”). The full set of keywords and  
 358 weights can be found in Appendix A.

359 For a book knowledge graph  $G$  and the set of  
 360 nodes  $N_k$  of entities mentioned in the current chap-  
 361 ter text  $\mathcal{C}_k$ , we obtain the induced subgraph  $G[N_k]$   
 362 that only contains the nodes in  $N_k$  and edges be-  
 363 tween them. We further judge that future informa-  
 364 tion is mostly unnecessary for generating a chapter  
 365 summary and remove edges from  $G[N_k]$  that were  
 366 extracted from later in the story than the current  
 367 chapter. We define the remaining set of edges as  
 368  $E_k$ .

369 We use Sentence-BERT (Reimers and Gurevych,  
 370 2019) to compute similarity scores between embed-  
 371 dings for each keyword in the set of keywords  $Q$   
 372 and the predicate portion of each  $\langle$ subject, predi-  
 373 cate, object $\rangle$  edge in  $E_k$ . Let  $s_{ij} = \cos\_sim(p_i, q_j)$   
 374 be the cosine similarity score between the embed-  
 375 ding of the  $i$ -th predicate  $p_i$  ( $1 \leq i \leq |E_k|$ ) and  
 376 the embedding of the  $j$ -th keyword  $q_j$  ( $1 \leq j \leq |Q|$ ).  
 377 Then the normalized similarity score  $\tilde{s}_{ij}$  is com-  
 378 puted as:

$$379 \quad \tilde{s}_{ij} = \frac{s_{ij} - \mu_j}{\sigma_j} \quad (2)$$

380 where  $\mu_j$  and  $\sigma_j$  are the mean and standard devia-  
 381 tion, respectively, of the scores for the  $j$ -th keyword  
 382 across all edges. This normalization is required to  
 383 fairly weight the scores, as some keywords may  
 384 have generally higher or lower scores than others.  
 385 The weighted score  $W_{ij}$  for each edge and keyword  
 386 is then:

$$387 \quad W_{ij} = \tilde{s}_{ij} \cdot w_j \quad (3)$$

388 where  $w_j$  is the weight associated with the  $j$ -th  
 389 keyword. The final aggregated score  $S_i$  for the  $i$ -th

edge is the sum of its weighted scores across all  
 keywords:

$$390 \quad S_i = \sum_j W_{ij} \quad (4) \quad 392$$

This score is used to rank the edges in  $E_k$ . 393

394 **Edge Linearization** We provide the retrieved  
 395 knowledge graph edges to the summarization  
 396 model as a linearized string of text prepended be-  
 397 fore the chapter text. Starting from the highest-  
 398 ranked edge (*i.e.*, the edge with the highest score),  
 399 we gather edges in  $E_k$  until the context length limit  
 400 is reached, taking into account the length of the  
 401 chapter text. We then arrange the gathered edges by  
 402 their subject entities so that edges with a common  
 403 subject are grouped together. In each group, we cat-  
 404 egorize the edges by their objects. Additionally, we  
 405 sort the subjects and objects by the total number of  
 406 appearances in the chapter text in decreasing order.  
 407 This puts entities important to the chapter near the  
 408 front of the linearized text. Self-loops, which are  
 409 treated as edges with no object, are placed before  
 410 other edges with the same subject.

411 The linearization format depends on the level  
 412 of access that is available for the summarization  
 413 model. For models that are able to be finetuned  
 414 for the task, we use a format that includes three  
 415 new special tokens:  $\langle$ subject $\rangle$ ,  $\langle$ object $\rangle$ , and  
 416  $\langle$ predicate $\rangle$ . The model is expected to learn the  
 417 meaning of these tokens during the finetuning pro-  
 418 cess. Each group of edges sharing a subject is  
 419 preceded by the  $\langle$ subject $\rangle$  token and the name of  
 420 the subject. If the entity has multiple names, the  
 421 name that appears the most frequently in the chap-  
 422 ter text is used. Inside each group, every new object  
 423 is marked by the  $\langle$ object $\rangle$  token and the object  
 424 name (again, the most commonly used name). This  
 425 part is omitted for edges with no object. Finally, a  
 426  $\langle$ predicate $\rangle$  token and the predicate text are ap-  
 427 pended for each predicate. An example is shown in  
 428 Figure 1b. The input to the summarization model  
 429 is the string of linearized edges, followed by an ad-  
 430 ditional  $\langle$ chapter $\rangle$  special token and the chapter  
 431 text.

432 If the summarization model is one that must  
 433 be used without finetuning (*e.g.*, a model that can  
 434 only be accessed through an inference API), the  
 435 linearization format simply consists of each edge  
 436 on its own line, with the subject, predicate, and  
 437 object separated by a single space and semicolon  
 438 (*e.g.*, Juliet; gives ring to; Romeo); no ad-  
 439 ditional special tokens are used. The subjects and

objects for each edge are always specified, even when they are the same as those in the previous edge. The intent is to present the edge information in a format that the model can understand without further training.

### 3.5 Knowledge Graph Similarity Metric (KGScore)

Existing metrics such as ROUGE and BERTScore may not be good measures of the quality of abstractive summaries (Novikova et al., 2017), especially when it comes to factuality and faithfulness to the original text (Maynez et al., 2020). As we believe the main role of the knowledge graphs is to provide global context and enhance factual consistency, a metric that explicitly checks for adherence to established facts would be helpful for evaluating the effect of the knowledge graphs. For stories, many of these facts have to do with character descriptions, actions taken by characters, and relationships between them, all of which can be represented in a knowledge graph.

We propose a novel metric that we call KGScore, which measures the quality of a summary by computing the similarity between two knowledge graphs. Let  $G_g$  and  $G_r$  represent knowledge graphs extracted from a generated and reference summary, respectively. They contain edges  $E_g$  and  $E_r$ . For an edge  $e_g \in E_g$  with subject  $s_{e_g}$  and object  $o_{e_g}$ , let  $E_{r,e_g}$  be the subset of edges in  $E_r$  that have the same subject  $s_{e_r}$  and object  $o_{e_r}$  as  $e_g$ :

$$E_{r,e_g} = \{e_r \in E_r \mid s_{e_r} = s_{e_g} \wedge o_{e_r} = o_{e_g}\} \quad (5)$$

$E'_g$  is the subset of edges in  $E_g$  for which  $E_{r,e_g}$  is not empty:

$$E'_g = \{e_g \in E_g \mid E_{r,e_g} \neq \emptyset\} \quad (6)$$

Then the precision KGScore  $P_{KG}$  is defined as follows:

$$P_{KG} = \frac{1}{|E_g|} \sum_{e_g \in E'_g} \max_{e_r \in E_{r,e_g}} \cos\_sim(p_{e_g}, p_{e_r}) \quad (7)$$

where  $\cos\_sim(p_{e_g}, p_{e_r})$  is the cosine similarity between the embeddings of the predicates of  $e_g$  and  $e_r$ . The precision score is roughly equivalent to a measure of the proportion of the information in the generated summary graph that also exists in the reference summary graph. For the recall score

$R_{KG}$ , the direction is reversed:

$$R_{KG} = \frac{1}{|E_r|} \sum_{e_r \in E'_r} \max_{e_g \in E_{g,e_r}} \cos\_sim(p_{e_g}, p_{e_r}) \quad (8)$$

where  $E_{g,e_r}$  and  $E'_r$  are defined in the same way as  $E_{r,e_g}$  and  $E'_g$ , but with the roles of the two graphs swapped. Finally, the F1 score  $F_{KG}$  is the harmonic mean of the precision and recall scores:

$$F_{KG} = \frac{2 \cdot P_{KG} \cdot R_{KG}}{P_{KG} + R_{KG}} \quad (9)$$

To generate the two knowledge graphs  $G_g$  and  $G_r$ , we follow a process similar to the one for book knowledge graph generation, with some modifications. Instead of using a single prompt and model to identify both named entities and knowledge graph edges, we split the process into two steps.

First, we use spaCy (Montani et al., 2023) to find named entities in the reference summary (we use version 3.7.3 of the en-core-web-trf pipeline). This is faster than the previous approach of using a large language model, but it is limited in that it cannot identify aliases or name variations. We consider this an acceptable compromise, as summaries are generally short and it is unlikely that multiple names are used for a single character. Consequently, we skip the steps of generating a names graph and merging nodes that represent the same entity.

For the next part of the process, we opt for a locally run model instead of an OpenAI model and pair it with the Guidance library<sup>1</sup> to constrain the model output to text that can be parsed into valid knowledge graph edges. This is to increase accuracy and lower costs at the expense of longer generation times. More specifically, we use Mixtral 8x7B (Jiang et al., 2024), a mixture-of-experts model with 13 billion active parameters. The prompt includes the named entities from the reference summary, as well as three few-shot examples of extracting edges from summaries, formatted as a multi-turn conversation. The full prompt is in Appendix F. We use the same entities for both the reference and generated summaries to maximize the overlap between entities in the two sets of edges, which is important for computing reliable KGScore values. We omit the original final step of removing nodes with few edges because the number of edges is generally small for a summary.

<sup>1</sup><https://github.com/guidance-ai/guidance>

528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
  
539  
540  
  
541  
542  
543  
544  
545  
546  
547  
548  
  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574

## 4 Experiments

### 4.1 Dataset

For our experiments, we choose the BookSum Chapters dataset (Kryściński et al., 2021), which contains chapter texts and their summaries from over 200 English-language books, including novels, plays, and short stories. We filter the dataset to improve its quality and better align it with our purposes, and we are left with 7255 examples in the training set and 1155 examples in the validation set; details are in Appendix B.

### 4.2 Effectiveness of Knowledge Graph Retrieval

To verify the effectiveness of our knowledge graph retrieval method, we finetune two LongT5 (Guo et al., 2022) models for summarization on the filtered BookSum Chapters dataset: a baseline model trained with only the chapter texts as input (LongT5-No-KG) and a model trained using our proposed method (LongT5-KG). Finetuning details and parameters are included in Appendix C.

We additionally apply our method to OpenAI’s gpt-4-1106-preview model with no finetuning, using the simple edge linearization format as described in Section 3.4. We employ Chain of Density (Adams et al., 2023), an iterative prompting technique that produces entity-dense summaries, to maximize the effect of the entity information contained in the knowledge graph edges. We label the baseline results without our knowledge graph retrieval method as GPT-4-No-KG and the results from our method as GPT-4-KG.

The results are summarized in Table 1. Due to resource constraints, we randomly select 100 chapters from the test set of the BookSum Chapters dataset and report evaluation results on this smaller subset. Results for three sets of metrics are included: ROUGE, BERTScore, and our proposed KGScore. Details on the evaluation procedure are in Appendix D. For the LongT5 models, the model using our knowledge graph retrieval method (LongT5-KG) achieves higher scores across all metrics compared to its baseline counterpart (LongT5-No-KG). For the GPT-4 results, our method outperforms the baseline in terms of KGScore while receiving slightly lower scores for ROUGE and BERTScore.

### 4.3 Validity of KGScore

We perform an additional experiment to verify the hypothesis that our proposed KGScore metric is a valid measure of factual consistency. We filter the training set of the BookSum Chapters dataset for summaries whose word counts fall within the range of 300 to 450 words, then gather pairs of summaries of the same book chapter. We select the 100 most similar summary pairs and use them as the baseline dataset for our experiment. Although all of these summaries are human written, we randomly select one summary from each pair as a “prediction” (“generated”) summary and the other as a “reference” summary for the purpose of calculating evaluation metrics.

Next, we create a modified version of the dataset by identifying the named entities in each prediction summary using spaCy (Montani et al., 2023) and shuffling their locations, ensuring that entities only get swapped with other entities of the same type (e.g., person or location). Much of the factual information included in this new summary is likely to be inaccurate. While the baseline dataset consists of pairs of summaries containing similar information (as they are summaries of the same chapter), the altered prediction summaries in the entity-shuffled dataset factually deviate from the reference summaries. Therefore, a reliable factuality metric should produce a significantly lower score for the entity-shuffled dataset in comparison to the baseline dataset.

The results of evaluating the two datasets on ROUGE, BERTScore, and KGScore are shown in Table 2. The decrease in metric values from the baseline to the entity-shuffled dataset is substantially greater in KGScore compared to ROUGE and BERTScore, which exhibit relatively small reductions.

## 5 Analysis

Combining and examining the results of the experiment in Section 4.2, where our knowledge graph retrieval method attains better KGScore results than the baselines, and the results of the entity shuffle experiment in Section 4.3, which show that KGScore is significantly more sensitive to variations in factual accuracy than ROUGE and BERTScore, we claim that our method successfully improves factual consistency in summaries as intended. This improvement may not always be detectable through traditional metrics, as evident in the GPT-4 re-

575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624

Model	ROUGE-1	ROUGE-2	ROUGE-L	BERTScore F1	$P_{KG}$	$R_{KG}$	$F_{KG}$
LongT5-No-KG	28.80	5.48	13.96	53.82	22.54	15.69	17.73
LongT5-KG	<b>30.07</b>	<b>5.91</b>	<b>14.51</b>	<b>54.58</b>	<b>23.07</b>	<b>16.59</b>	<b>18.34</b>
GPT-4-No-KG	<b>25.06</b>	<b>3.76</b>	<b>14.03</b>	<b>56.28</b>	23.26	18.67	20.00
GPT-4-KG	24.14	3.60	13.76	55.96	<b>25.57</b>	<b>20.17</b>	<b>21.75</b>

Table 1: Evaluation results on a subset of the BookSum Chapters test set.  $P_{KG}$ ,  $R_{KG}$ , and  $F_{KG}$  are average KGScore values as defined in Section 3.5. The best scores among each model category (LongT5 and GPT-4) are in bold.

	ROUGE-1	ROUGE-2	ROUGE-L	BERTScore F1	$P_{KG}$	$R_{KG}$	$F_{KG}$
Baseline	50.45	13.15	23.51	63.62	26.03	24.67	24.93
Shuffled entities	50.34	11.79	21.59	60.44	16.10	14.33	14.78
Change (%)	-0.2	-10.3	-8.1	-5.0	<b>-38.1</b>	<b>-41.9</b>	<b>-40.7</b>

Table 2: Results of the entity shuffle experiment.  $P_{KG}$ ,  $R_{KG}$ , and  $F_{KG}$  are average KGScore values as defined in Section 3.5. Changes in metric values after the entity shuffling are shown, with KGScore results in bold.

sults, where the application of our method leads to slightly worse ROUGE and BERTScore values than the baseline.

## 5.1 KGScore Trends

In Tables 1 and 2, the precision KGScore ( $P_{KG}$ ) is higher than the recall KGScore ( $R_{KG}$ ) in all cases. This is because the named entities are identified from the reference summary and used for extracting knowledge graph edges in both the reference and generated summaries, as described in Section 3.5. All entities contained in the edges can be found in the reference texts, while some are missing in the generated texts. This imbalance could be removed by gathering named entities from both summaries, but this could potentially introduce incorrectly hallucinated entities from the generated summaries.

Another related observation is that the KGScore values are low overall, ranging in the 10s and 20s out of a theoretical maximum of 100 (%). This could also be the symptom of an entity-matching problem, as a single entity may sometimes appear under different names in the prediction and reference summaries. Adding a step in the metric computation process to identify these cases could help, but if overly eager deductions are made about which names refer to the same entity, a new problem could arise where entities that should be kept separate are merged into one.

## 5.2 Qualitative Evaluation

To verify the interpretation that our knowledge graph retrieval method improves factual consistency among entities, we qualitatively evaluate a

small sample of generated chapter summaries. Examples of these summaries are in Appendix G.

For the finetuned LongT5 models, we find that the summaries generated by both the LongT5-No-KG model and the LongT5-KG model are of low quality. They are similar to extractive summaries, repeating large sections of story text verbatim, and they often contain errors. This basic deficiency in summarization performance makes it difficult to determine the effects of using our method.

In comparison, we perceive the GPT-4 summaries to be higher-quality, which is corroborated by the higher KGScore values in Table 1. It is interesting to note that the ROUGE-1 and ROUGE-2 scores are much lower than those of the LongT5 models; this could be another indication that ROUGE does not align well with human judgments.

Although GPT-4-KG receives better KGScore results than GPT-4-No-KG, we find it challenging to identify specific examples of retrieved knowledge graph edges improving summary quality. It could be that the degree of improvement in this case is too subtle for humans to readily notice.

## 6 Conclusion

We present a method of summarizing long stories that employs knowledge graph retrieval to provide global context. We also propose a novel metric, KGScore, which evaluates summaries based on knowledge graph similarity. Experimental results indicate that our approach may enhance factual consistency and that our KGScore metric is an effective measure of it.



## 690 Limitations

691 Besides the improvement in KGScore, we have not  
692 presented additional evidence that our knowledge  
693 graph retrieval method enhances summary quality,  
694 such as a large-scale human evaluation.

695 The LongT5 model chosen for finetuning is a  
696 relatively old model with limited performance, as  
697 described in Section 5.2. More recent models, such  
698 as LLaMA 2 (Touvron et al., 2023), could be better  
699 suited for future experiments.

700 It can be argued that the observed score improve-  
701 ments using our method are not very significant. To  
702 achieve greater improvements, methods to more ef-  
703 fectively provide information to the summarization  
704 model through the knowledge graph edges could be  
705 explored. For example, providing temporal infor-  
706 mation along with the edges, either in the form of  
707 chapter numbers or relative positions in the story,  
708 could be beneficial. The model can then gain a  
709 sense of time by placing events and details in re-  
710 lation to each other and the current chapter. This  
711 can be especially helpful if significant changes oc-  
712 cur to an entity over the course of the story (e.g.,  
713 a death or a relocation). Another potential point  
714 of enhancement could lie in the way the retrieved  
715 edges are provided to the model. For instance, in-  
716 stead of simply prepending a linearized string to  
717 the context, a graph neural network could be used  
718 to process the knowledge graph information.

## 719 Ethics Statement

720 We do not anticipate any ethical issues. Our work  
721 deals with producing summaries of existing sto-  
722 ries, so malicious applications would be limited,  
723 although the generated summaries could be inaccur-  
724 ate or reflect harmful content in the source text.

725 All tools used are open source, including the  
726 scripts to download the BookSum Chapters dataset,  
727 which are released under the BSD 3-Clause Li-  
728 cense. Considering the nature of the dataset, which  
729 consists of book texts and their summaries, its con-  
730 tent was not inspected for offensive material or per-  
731 sonally identifiable information. Where specified,  
732 usage of tools and datasets was done in accordance  
733 with their intended use.

## 734 References

735 Griffin Adams, Alex Fabbri, Faisal Ladhak, Eric  
736 Lehman, and Noémie Elhadad. 2023. From Sparse to  
737 Dense: GPT-4 Summarization with Chain of Density

Prompting. In *Proceedings of the 4th New Frontiers  
in Summarization Workshop*, pages 68–74. 738  
739

Berkeley R. Andrus, Yeganeh Nasiri, Shilong Cui, Ben-  
jamin Cullen, and Nancy Fulda. 2022. Enhanced  
Story Comprehension for Large Language Mod-  
els through Dynamic Document-Based Knowledge  
Graphs. *Proceedings of the AAAI Conference on  
Artificial Intelligence*, 36(10):10436–10444. 740  
741  
742  
743  
744  
745

Gabor Angeli, Melvin Jose Johnson Premkumar, and  
Christopher D. Manning. 2015. Leveraging Lin-  
guistic Structure For Open Domain Information Ex-  
traction. In *Proceedings of the 53rd Annual Meet-  
ing of the Association for Computational Linguistics  
and the 7th International Joint Conference on Natu-  
ral Language Processing (Volume 1: Long Papers)*,  
pages 344–354. 746  
747  
748  
749  
750  
751  
752  
753

Jennifer A. Bishop, Qianqian Xie, and Sophia Anani-  
adou. 2023. LongDocFACTScore: Evaluating the  
Factuality of Long Document Abstractive Summari-  
sation. *arXiv preprint arXiv:2309.12455*. 754  
755  
756  
757

Tom Brown, Benjamin Mann, Nick Ryder, Melanie  
Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind  
Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
Askell, Sandhini Agarwal, Ariel Herbert-Voss,  
Gretchen Krueger, Tom Henighan, Rewon Child,  
Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens  
Winter, Chris Hesse, Mark Chen, Eric Sigler, Ma-  
teusz Litwin, Scott Gray, Benjamin Chess, Jack  
Clark, Christopher Berner, Sam McCandlish, Alec  
Radford, Ilya Sutskever, and Dario Amodei. 2020.  
Language Models are Few-Shot Learners. In *Ad-  
vances in Neural Information Processing Systems*,  
volume 33, pages 1877–1901. 758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770

Shuyang Cao and Lu Wang. 2023. AWESOME: GPU  
Memory-constrained Long Document Summariza-  
tion using Memory Mechanism and Global Salient  
Content. *arXiv preprint arXiv:2305.14806*. 771  
772  
773  
774

Yapei Chang, Kyle Lo, Tanya Goyal, and Mohit Iyyer.  
2023. BoookScore: A systematic exploration  
of book-length summarization in the era of LLMs.  
*arXiv preprint arXiv:2310.00785*. 775  
776  
777  
778

Tong Chen, Xuwei Wang, Tianwei Yue, Xiaoyu Bai,  
Cindy X. Le, and Wenping Wang. 2023. Enhancing  
Abstractive Summarization with Extracted Knowl-  
edge Graphs and Multi-Source Transformers. *Ap-  
plied Sciences*, 13(13):7753. 779  
780  
781  
782  
783

Angela Fan, Claire Gardent, Chloé Braud, and An-  
toine Bordes. 2019. Using Local Knowledge Graph  
Construction to Scale Seq2Seq Models to Multi-  
Document Inputs. In *Proceedings of the 2019 Confer-  
ence on Empirical Methods in Natural Language Pro-  
cessing and the 9th International Joint Conference  
on Natural Language Processing (EMNLP-IJCNLP)*,  
pages 4186–4196. 784  
785  
786  
787  
788  
789  
790  
791

Mingqi Gao, Jie Ruan, Renliang Sun, Xunjian Yin, Ship-  
ing Yang, and Xiaojun Wan. 2023. Human-like Sum-  
marization Evaluation with ChatGPT. *arXiv preprint  
arXiv:2304.02554*. 792  
793  
794  
795

796	Alexios Gidiotis and Grigorios Tsoumakas. 2020. A Divide-and-Conquer Approach to the Summarization of Long Documents. <i>IEEE/ACM Transactions on Audio, Speech, and Language Processing</i> , 28:3029–3040.	<i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 9332–9346.	853 854 855
801	Tanya Goyal, Junyi Jessy Li, and Greg Durrett. 2022. SNaC: Coherence Error Detection for Narrative Summarization. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 444–463.	Wojciech Kryściński, Nazneen Rajani, Divyansh Agarwal, Caiming Xiong, and Dragomir Radev. 2021. BookSum: A Collection of Datasets for Long-form Narrative Summarization. <i>arXiv preprint arXiv:2105.08209</i> .	856 857 858 859 860
806	Mandy Guo, Joshua Ainslie, David Uthus, Santiago Ontanon, Jianmo Ni, Yun-Hsuan Sung, and Yinfei Yang. 2022. LongT5: Efficient Text-To-Text Transformer for Long Sequences. In <i>Findings of the Association for Computational Linguistics: NAACL 2022</i> , pages 724–736.	M. V. P. T. Lakshika, H. A. Caldera, and W. V. Welgama. 2020. Abstractive Web News Summarization Using Knowledge Graphs. In <i>2020 20th International Conference on Advances in ICT for Emerging Regions (ICTer)</i> , pages 300–301.	861 862 863 864 865
812	Hardy Hardy, Miguel Ballesteros, Faisal Ladhak, Muhammad Khalifa, Vittorio Castelli, and Kathleen McKeown. 2022. Novel Chapter Abstractive Summarization using Spinal Tree Aware Sub-Sentential Content Selection. <i>arXiv preprint arXiv:2211.04903</i> .	Chin-Yew Lin. 2004. ROUGE: A Package for Automatic Evaluation of Summaries. In <i>Text Summarization Branches Out</i> , pages 74–81.	866 867 868
817	Luyang Huang, Shuyang Cao, Nikolaus Parulian, Heng Ji, and Lu Wang. 2021. Efficient Attentions for Long Document Summarization. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1419–1436.	Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. G-Eval: NLG Evaluation using GPT-4 with Better Human Alignment. <i>arXiv preprint arXiv:2303.16634</i> .	869 870 871 872
823	Luyang Huang, Lingfei Wu, and Lu Wang. 2020. Knowledge Graph-Augmented Abstractive Summarization with Semantic-Driven Cloze Reward. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 5094–5107.	Yang Liu and Mirella Lapata. 2019. Text Summarization with Pretrained Encoders. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3730–3740.	873 874 875 876 877 878
829	Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, L�elio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Th�eophile Gervet, Thibaut Lavril, Thomas Wang, Timoth�e Lacroix, and William El Sayed. 2024. Mixture of Experts. <i>arXiv preprint arXiv:2401.04088</i> .	Guang Lu, Sylvia B. Larcher, and Tu Tran. 2023. Hybrid Long Document Summarization using C2F-FAR and ChatGPT: A Practical Study. <i>arXiv preprint arXiv:2306.01169</i> .	879 880 881 882
830	Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, L�elio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Th�eophile Gervet, Thibaut Lavril, Thomas Wang, Timoth�e Lacroix, and William El Sayed. 2024. Mixture of Experts. <i>arXiv preprint arXiv:2401.04088</i> .	Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. On Faithfulness and Factuality in Abstractive Summarization. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 1906–1919.	883 884 885 886 887
840	Perna Kashyap. 2022. COLING 2022 Shared Task: LED Finteuning and Recursive Summary Generation for Automatic Summarization of Chapters from Novels. In <i>Proceedings of The Workshop on Automatic Summarization for Creative Writing</i> , pages 19–23.	Ines Montani, Matthew Honnibal, Matthew Honnibal, Adriane Boyd, Sofie Van Landeghem, and Henning Peters. 2023. spaCy: Industrial-strength Natural Language Processing in Python.	888 889 890 891
845	Diederik P. Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In <i>3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings</i> .	Jekaterina Novikova, Ondr�ej Du�sek, Amanda Cercas Curry, and Verena Rieser. 2017. Why We Need New Evaluation Metrics for NLG. In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 2241–2252.	892 893 894 895 896
850	Wojciech Kryscinski, Bryan McCann, Caiming Xiong, and Richard Socher. 2020. Evaluating the Factual Consistency of Abstractive Text Summarization. In	OpenAI. 2023. GPT-4 Technical Report. <i>arXiv preprint arXiv:2303.08774</i> .	897 898
852	Wojciech Kryscinski, Bryan McCann, Caiming Xiong, and Richard Socher. 2020. Evaluating the Factual Consistency of Abstractive Text Summarization. In	Bo Pang, Erik Nijkamp, Wojciech Kryscinski, Silvio Savarese, Yingbo Zhou, and Caiming Xiong. 2023. Long Document Summarization with Top-down and Bottom-up Inference. In <i>Findings of the Association for Computational Linguistics: EACL 2023</i> , pages 1267–1284.	899 900 901 902 903 904

905	Jason Phang, Yao Zhao, and Peter J. Liu. 2022.	State-of-the-Art Natural Language Processing. In	964
906	Investigating Efficiently Extending Transformers	<i>Proceedings of the 2020 Conference on Empirical</i>	965
907	for Long Input Summarization. <i>arXiv preprint</i>	<i>Methods in Natural Language Processing: System</i>	966
908	<i>arXiv:2208.04347</i> .	<i>Demonstrations</i> , pages 38–45.	967
909	Jonathan Pilault, Raymond Li, Sandeep Subramanian,	Jeff Wu, Long Ouyang, Daniel M. Ziegler, Nisan Sti-	968
910	and Chris Pal. 2020. On Extractive and Abstractive	ennon, Ryan Lowe, Jan Leike, and Paul Christiano.	969
911	Neural Document Summarization with Transformer	2021. Recursively Summarizing Books with Human	970
912	Language Models. In <i>Proceedings of the 2020 Con-</i>	Feedback. <i>arXiv preprint arXiv:2109.10862</i> .	971
913	<i>ference on Empirical Methods in Natural Language</i>		
914	<i>Processing (EMNLP)</i> , pages 9308–9319.		
915	Nils Reimers and Iryna Gurevych. 2019. Sentence-	Yuexiang Xie, Fei Sun, Yang Deng, Yaliang Li, and	972
916	BERT: Sentence Embeddings using Siamese BERT-	Bolin Ding. 2021. Factual Consistency Evaluation	973
917	Networks. In <i>Proceedings of the 2019 Conference on</i>	for Text Summarization via Counterfactual Estima-	974
918	<i>Empirical Methods in Natural Language Processing</i>	tion. In <i>Findings of the Association for Computa-</i>	975
919	<i>and the 9th International Joint Conference on Natu-</i>	<i>tional Linguistics: EMNLP 2021</i> , pages 100–110.	976
920	<i>ral Language Processing (EMNLP-IJCNLP)</i> , pages		
921	3982–3992.	Jingqing Zhang, Yao Zhao, Mohammad Saleh, and	977
922	Thomas Scialom, Paul-Alexis Dray, Sylvain Lamprier,	Peter J. Liu. 2020a. PEGASUS: Pre-training with	978
923	Benjamin Piwowarski, Jacopo Staiano, Alex Wang,	extracted gap-sentences for abstractive summariza-	979
924	and Patrick Gallinari. 2021. QuestEval: Summariza-	tion. In <i>Proceedings of the 37th International Confer-</i>	980
925	tion Asks for Fact-based Evaluation. In <i>Proceedings</i>	<i>ence on Machine Learning</i> , volume 119 of <i>ICML'20</i> ,	981
926	<i>of the 2021 Conference on Empirical Methods in</i>	pages 11328–11339.	982
927	<i>Natural Language Processing</i> , pages 6594–6604.		
928	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	Mengli Zhang, Gang Zhou, Wanting Yu, Ningbo Huang,	983
929	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	and Wenfen Liu. 2022a. A Comprehensive Survey	984
930	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	of Abstractive Text Summarization Based on Deep	985
931	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton	Learning. <i>Computational Intelligence and Neuro-</i>	986
932	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,	<i>science</i> , 2022:e7132226.	987
933	Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,		
934	Cynthia Gao, Vedanuj Goswami, Naman Goyal, Antho-	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q.	988
935	ny Hartshorn, Saghar Hosseini, Rui Hou, Hakan	Weinberger, and Yoav Artzi. 2020b. BERTScore:	989
936	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,	Evaluating Text Generation with BERT. <i>arXiv</i>	990
937	Isabel Kloumann, Artem Korenev, Punit Singh Koura,	<i>preprint arXiv:1904.09675</i> .	991
938	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-		
939	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-	Yusen Zhang, Ansong Ni, Ziming Mao, Chen Henry Wu,	992
940	tinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-	Chenguang Zhu, Budhaditya Deb, Ahmed Awadallah,	993
941	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-	Dragomir Radev, and Rui Zhang. 2022b. Summ <sup>N</sup> :	994
942	stein, Rashi Rungta, Kalyan Saladi, Alan Schel-	A Multi-Stage Summarization Framework for Long	995
943	ten, Ruan Silva, Eric Michael Smith, Ranjan Sub-	Input Dialogues and Documents. In <i>Proceedings</i>	996
944	ramanian, Xiaoqing Ellen Tan, Binh Tang, Ross	<i>of the 60th Annual Meeting of the Association for</i>	997
945	Taylor, Adina Williams, Jian Xiang Kuan, Puxin	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	998
946	Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, An-	pages 1592–1604.	999
947	gela Fan, Melanie Kambadur, Sharan Narang, Aure-		
948	lien Rodriguez, Robert Stojnic, Sergey Edunov, and	Yao Zhao, Mohammad Saleh, and Peter J. Liu.	1000
949	Thomas Scialom. 2023. Llama 2: Open Founda-	2020. SEAL: Segment-wise Extractive-Abstractive	1001
950	tion and Fine-Tuned Chat Models. <i>arXiv preprint</i>	Long-form Text Summarization. <i>arXiv preprint</i>	1002
951	<i>arXiv:2307.09288</i> .	<i>arXiv:2006.10213</i> .	1003
952	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	Chenguang Zhu, William Hinthorn, Ruochen Xu,	1004
953	Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz	Qingkai Zeng, Michael Zeng, Xuedong Huang, and	1005
954	Kaiser, and Illia Polosukhin. 2017. Attention is All	Meng Jiang. 2021. Enhancing Factual Consistency	1006
955	you Need. In <i>Advances in Neural Information Pro-</i>	of Abstractive Summarization. In <i>Proceedings of</i>	1007
956	<i>cessing Systems</i> , volume 30.	<i>the 2021 Conference of the North American Chap-</i>	1008
957	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien	<i>ter of the Association for Computational Linguistics:</i>	1009
958	Chaumond, Clement Delangue, Anthony Moi, Pier-	<i>Human Language Technologies</i> , pages 718–733.	1010
959	ric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz,		
960	Joe Davison, Sam Shleifer, Patrick von Platen, Clara	<b>A Keywords for Edge Retrieval</b>	1011
961	Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven	Table 3 lists the keywords and weights for knowl-	1012
962	Le Scao, Sylvain Gugger, Mariama Drame, Quentin	edge graph edge retrieval. The weight for each	1013
963	Lhoest, and Alexander Rush. 2020. Transformers:	keyword is empirically chosen based on the per-	1014
		ceived effectiveness at retrieving useful edges.	1015

Keyword	Weight
relation	30
happen	15
conflict	10
desire	10
emotion	10
role	10
think	5
location	5
personality	5

Table 3: Keywords and weights for knowledge graph edge retrieval.

## D Evaluation

For the finetuned LongT5 models, we evaluate the checkpoints with the lowest validation losses during training, with 4 beams and no\_repeat\_ngram\_size set to 3. We use version 0.1.2 of the rouge-score Python package<sup>4</sup> and version 0.3.13 of the bert-score package<sup>5</sup> for ROUGE and BERTScore, respectively, with microsoft/deberta-xlarge-mnli<sup>6</sup> as the BERTScore model.

## B BookSum Chapters Dataset Filtering

We remove books that are found in both the training set and either the validation or test sets (book IDs 1130, 1526, 1783, and 1798). We also exclude books that are not narrative texts (61, 1232, 1404, 3207, 3420, 3755, 4320, 7370, 11224, 13434, and 34901) and books that are collections of multiple stories (221, 416, 610, 1429, and 15859). For some summaries, the script provided to build the BookSum dataset either fails to download the summary (because it is no longer available on the web) or downloads an incomplete version of it (empty or just a few words); we do not use those. Finally, in order to accommodate the context length limit of our model, we set chapter length and summary length limits of 16384 and 1024 tokens, respectively, and only use text pairs that fall under both limits.

## C LongT5 Finetuning

We use the Hugging Face Transformers library (Wolf et al., 2020) and begin finetuning from the pretrained long-t5-tglobal-base checkpoint<sup>2</sup>, which has 250 million parameters. For computing similarity scores with SentenceBERT, we use the all-MiniLM-L6-v2 model<sup>3</sup>. We train on 8 NVIDIA A100 GPUs for approximately 70 hours total for the two models, using the following parameters: Adam optimizer (Kingma and Ba, 2015), cosine learning rate scheduling with 2e-4 maximum learning rate and 1 epoch linear warm-up, batch size 128, 48 epochs. The hyperparameters were selected manually without tuning due to resource constraints.

<sup>2</sup><https://huggingface.co/google/long-t5-tglobal-base>

<sup>3</sup><https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2>

<sup>4</sup><https://pypi.org/project/rouge-score>

<sup>5</sup><https://pypi.org/project/bert-score>

<sup>6</sup><https://huggingface.co/microsoft/deberta-xlarge-mnli>



## E Prompt for Book Knowledge Graph Generation 1059

Read part of a story, then identify named entities and generate knowledge graph edges. 1060

[Begin story excerpt] 1061

“Christmas won’t be Christmas without any presents,” grumbled Jo. “It’s so dreadful to be poor!” sighed Meg, looking out the window at the snow-covered streets of Concord. “I don’t think it’s fair for some girls to have plenty of pretty things, and other girls nothing at all,” added little Amy, with an injured sniff. “We’ve got Father and Mother, and each other,” said Beth contentedly from her corner. The four young faces brightened at the cheerful words, but darkened again as Jo said sadly, “We haven’t got Father, and shall not have him for a long time.” She didn’t say “perhaps never,” but each silently added it, thinking of Father far away, where the fighting was. 1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069

As young readers like to know ‘how people look’, we will take this moment to give them a little sketch of the four sisters. Margaret March, the eldest of the four, was sixteen, and very pretty, with large eyes, plenty of soft brown hair, a sweet mouth, and white hands. Fifteen-year-old Jo March was very tall, thin, and brown, and never seemed to know what to do with her long limbs. Elizabeth, or Beth, as everyone called her, was a rosy, smooth-haired, bright-eyed girl of thirteen, with a shy manner, a timid voice, and a peaceful expression which was seldom disturbed. Amy, the youngest, was a regular snow maiden, with blue eyes, and yellow hair curling on her shoulders. 1070  
1071  
1072  
1073  
1074  
1075  
1076

The clock struck six and, having swept up the hearth, Beth put a pair of slippers down to warm. Somehow the sight of the old shoes had a good effect upon the girls, for Mother was coming, and everyone brightened to welcome her. Jo sat up to hold the slippers nearer to the blaze. “They are quite worn out. Marmee must have a new pair.” “I thought I’d get her some with my dollar,” said Beth. “No, I shall!” cried Amy. “I’ll tell you what we’ll do,” said Beth, “let’s each get her something for Christmas, and not get anything for ourselves.” “Let Marmee think we are getting things for ourselves, and then surprise her. We must go shopping tomorrow afternoon,” said Jo, marching up and down. 1077  
1078  
1079  
1080  
1081  
1082  
1083

“Glad to find you so merry, my girls,” said a cheery voice at the door, and the girls turned to welcome a tall, motherly lady. She was not elegantly dressed, but the girls thought the gray cloak and unfashionable bonnet covered the most splendid mother in the world. As they gathered about the table, Mrs. March said, with a particularly happy face, “I’ve got a treat for you after supper.” A quick, bright smile went round like a streak of sunshine. Beth clapped her hands, and Jo tossed up her napkin, crying, “A letter! A letter! Three cheers for Father!” “Yes, a nice long letter. He is well, and he sends all sorts of loving wishes for Christmas, and an especial message to you girls,” said Mrs. March, patting her pocket as if she had got a treasure there. “I think it was so splendid in Father to go as chaplain when he was too old to be drafted, and not strong enough for a soldier,” said Meg warmly, proud of her father’s work with the Union Army. 1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092

[End story excerpt] 1093

Named entities (include all aliases and name variations): 1094  
1095

Jo / Jo March 1096

Meg / Margaret / Margaret March 1097

Amy 1098

Beth / Elizabeth 1099

March sisters 1100

Mrs. March / Marmee / Mother 1101

Father 1102

Concord 1103

Union Army 1104  
1105

Knowledge graph edges (select up to 15 most important, `subject(s); predicate; [object(s)]` format, named entities only, predicate: five words max): 1106  
1107

1. Jo, Meg, Amy, Beth; in; March sisters 1108

2. March sisters; daughters of; Mrs. March, Father 1109

- 1110 3. Mrs. March; mother of; March sisters
- 1111 4. Father; father of; March sisters
- 1112 5. March sisters, Mrs. March; living in; Concord
- 1113 6. Father; away fighting in war
- 1114 7. Father; chaplain in; Union Army
- 1115 8. Meg; sixteen years old
- 1116 9. Jo; fifteen years old
- 1117 10. Beth; thirteen years old
- 1118 11. Beth; shy
- 1119 12. Amy; youngest among; March sisters
- 1120 13. March sisters; complained about not getting presents
- 1121 14. March sisters; decided to buy presents for; Mrs. March
- 1122 15. Mrs. March; brought home a letter from; Father

1123  
 1124 [Begin story excerpt]  
 1125 {story excerpt}  
 1126 [End story excerpt]

## 1127 **F Prompt for Summary Knowledge Graph Generation**

1128 Your task is to create a list of knowledge graph edges from a chapter summary. Here are the rules you  
 1129 must follow:

- 1130 - A knowledge graph edge is in the format “<subject(s)>; [None] or <object(s)>; <predicate>”.
- 1131 - Subjects and objects must be named entities.
- 1132 - If there are multiple subjects or objects, separate them with commas.
- 1133 - Predicates describe a relationship or action between the subjects and objects.
- 1134 - Predicates should not contain names.
- 1135 - Keep the predicates short (four words max).
- 1136 - If the knowledge graph edge is a description or action of a single entity with no object, use “[None]” in  
 1137 place of the object(s).
- 1138 - Only include information explicitly provided in the summary.
- 1139 - Find a diverse set of edges, never repeating similar edges.
- 1140 - Order the edges by importance rather than appearance in the summary, with the most important edges first.

1141  
 1142 Summary:

1143 The trio continues their arduous journey toward Canyon B, 300 km east of their home colony. Tensions  
 1144 rise as Janek and AC-293 clash over the best route to take. The Narrator, who has been silent for the  
 1145 past few days, finally speaks up and suggests that they take a shortcut through the mountains. Janek  
 1146 and AC-293 are skeptical, but they agree to try it. The Narrator leads them through a narrow pass, and  
 1147 they soon find themselves in an open valley. The Martian landscape is breathtaking, and the trio sets  
 1148 up camp for the night. The next morning, they continue their journey through the valley. Suddenly,  
 1149 they hear a loud noise and see a cloud of dust rising in the distance. They realize that a sandstorm  
 1150 is approaching, and they must find shelter quickly. They spot a cave in the distance and rush toward  
 1151 it. In the haste, Janek trips and falls, injuring his leg. The Narrator and AC-293 help him to his  
 1152 feet, and they make it to the cave just in time. They wait out the storm and emerge the next morning  
 1153 to find that their rover has been buried in sand. They dig it out and continue their journey toward Canyon B.  
 1154

- 1155 - Narrator, Janek, AC-293; Canyon B; continues journey toward
- 1156 - Janek; AC-293; argues with
- 1157 - AC-293; Janek; argues with
- 1158 - Narrator; [None]; suggests shortcut through mountains
- 1159 - Narrator, Janek, AC-293; [None]; camps in valley

- Narrator, Janek, AC-293; [None]; notices approaching sandstorm	1160
- Narrator, Janek, AC-293; [None]; rushes toward cave	1161
- Janek; [None]; falls and injures leg	1162
- Narrator, AC-293; Janek; helps to feet	1163
- Narrator, Janek, AC-293; [None]; waits out storm	1164
- Narrator, Janek, AC-293; [None]; digs out rover	1165
- Janek, AC-293; [None]; skeptical of shortcut	1166
- Janek, AC-293; [None]; agrees to shortcut	1167
- Narrator; [None]; leads through narrow pass	1168
- Narrator, Janek, AC-293; [None]; makes it to cave	1169
- Canyon B; [None]; east of home colony	1170

Summary:

This chapter focuses on the dazzling Grand Acorn Gala at Furrington Grove. Sir Reginald, the host of the gala, is a wealthy hedgehog who has a reputation for being a bit of a snob. He is also a collector of rare artifacts, and he has invited the guests to bring their own treasures to the gala. Xander, a young fox, and his father, Yorick, are among the guests. Yorick is a famous explorer, and he has brought a rare artifact to the gala, a golden leaf. Xander is eager to show off his father's treasure, but when he approaches Sir Reginald, the host dismisses the golden leaf as common and uninteresting. Undeterred, Xander explores the gala and encounters Penelope, a wise owl, who shares a secret about the golden leaf: it was once part of a magical tree that grew in the forest. Xander is intrigued by the story, and he decides to investigate further. He discovers that the tree was destroyed by a terrible storm, and the golden leaf was the only thing that survived. Xander is determined to find the tree and restore it to its former glory.

- Sir Reginald; Grand Acorn Gala; host of	1184
- Sir Reginald; [None]; artifact collector	1185
- Xander; [None]; young fox	1186
- Yorick; Xander; father of	1187
- Yorick; [None]; brings golden leaf	1188
- Sir Reginald; [None]; dismisses golden leaf	1189
- Penelope; Xander; shares secret with	1190
- Xander; [None]; investigates golden leaf	1191
- Xander; [None]; wants to restore tree	1192
- Yorick; [None]; famous explorer	1193
- Xander; Sir Reginald; approaches	1194
- Xander; Penelope; encounters	1195
- Penelope; [None]; wise owl	1196
- Sir Reginald; [None]; wealthy	1197
- Sir Reginald; [None]; snob	1198
- Grand Acorn Gala; Furrington Grove; takes place at	1199

Summary:

In this chapter, Amelia, a seasoned investigator with the Justice & Integrity Taskforce, dives into a high-stakes art theft case at the Metropolitan Museum of Art in New York City. The missing painting leads her to Marco Santos, an enigmatic artist tied to the Paris Art Collective. Crossing the Atlantic, Amelia unravels a web of clues in the underground European art scene, exposing the ruthless European Art Syndicate behind the theft. As Amelia delves deeper, she unearths Marco's troubled past, connecting him to his controversial exhibitions at the Louvre and the Tate Modern. An unexpected alliance forms between them as they race to unveil the truth behind the theft. They are not the sole seekers of the stolen masterpiece, however. Dr. Harold Blackwood, CEO of Blackwood Enterprises, lurks in the shadows, manipulating the chaos for his personal gain. With a cunning mind and powerful connections, Blackwood poses a hidden threat. His motives remain veiled in mystery as he holds secret meetings in Paris, London,

1212 and New York. Amelia and Marco navigate a treacherous path, not only against the European Art  
1213 Syndicate but also against the calculated machinations of Dr. Harold Blackwood, who seeks to outsmart  
1214 them and claim the stolen artwork for himself.

- 1215
- 1216 - Amelia; Metropolitan Museum of Art; investigates theft at
- 1217 - European Art Syndicate; [None]; behind theft
- 1218 - Amelia; European Art Syndicate; exposes
- 1219 - Amelia; Marco Santos; forms alliance with
- 1220 - Marco Santos; Amelia; forms alliance with
- 1221 - Amelia, Marco Santos; [None]; races for truth
- 1222 - Dr. Harold Blackwood; Amelia, Marco Santos; tries to outsmart
- 1223 - Dr. Harold Blackwood; [None]; seeks stolen artwork
- 1224 - Amelia; Justice & Integrity Taskforce; investigator at
- 1225 - Marco Santos; [None]; enigmatic artist
- 1226 - Marco Santos; Paris Art Collective; has ties to
- 1227 - Amelia; Marco Santos; discovers past of
- 1228 - Marco Santos; Louvre, Tate Modern; held controversial exhibition at
- 1229 - European Art Syndicate; [None]; ruthless
- 1230 - Dr. Harold Blackwood; Blackwood Enterprises; CEO of
- 1231 - Dr. Harold Blackwood; [None]; cunning
- 1232 - Dr. Harold Blackwood; Paris, London, New York; holds meetings in
- 1233 - Metropolitan Museum of Art; New York City; located in

1234

1235 Summary:  
1236 {summary }

## 1237 G Examples of Generated Summaries

1238 The following are summaries generated from Chapter 8 of *Tess of the d'Urbervilles* by Thomas Hardy.  
1239 KGScore F1 values ( $F_{KG}$ ) are also provided. The reference summary is omitted due to copyright.

### 1240 **LongT5-No-KG** ( $F_{KG} = 35.54$ ):

1241 Having mounted beside her, Alec d'Urberville drove rapidly along the crest of the first hill, chatting  
1242 compliments to Tess as they went, the cart with her box being left far behind. Rising still, an immense  
1243 landscape stretched around them on every side; behind, the green valley of her birth, before, a gray country  
1244 of which she knew nothing except from her first brief visit to Trantridge.

### 1245 **LongT5-KG** ( $F_{KG} = 35.51$ ):

1246 D'Urberville and Tess are on their way to Trantridge. She is afraid that she will be thrown out of the  
1247 carriage if she touches his arm. She tries to hold on to him, but he insists that she never touch his arm  
1248 again. They reach the summit of another declivity, and he asks her to put one kiss on those holmberry  
1249 lips, or even on that warmed cheek, and she will stop. He says that he will break both their necks, but  
1250 she refuses to move. He tells her that she is mighty sensitive for a cottage girl, and that she should be  
1251 made sorry for the kiss he has given her. She says that she would not have come if he had known. He  
1252 swore at her, and then he hems her in between the gig and the hedge. He laughs heartily at the sight of his  
1253 misdemeanour, and says he would never do it any more against his will now.

### 1254 **GPT-4-No-KG** ( $F_{KG} = 26.82$ ):

1255 Tess, in white muslin, dreads Alec's driving and his horse Tib's temperament. Her defiance peaks when  
1256 she wipes away Alec's forced kiss with her handkerchief. Despite her involuntary hold on him during  
1257 their descent near Melbury Down and Wingreen, she asserts her resolve, opting to walk to Trantridge, her  
1258 strategic silence marking her commitment to her family's welfare.



**GPT-4-KG** ( $F_{KG} = 31.32$ ):

The emotional stakes are high as Tess's big tear and flushed face betray her distress during the tumultuous ride with d'Urberville, whose injured tone and passionate outbursts oscillate between anger and feigned concern. Tess's strategic silence and defiance, coupled with her evasion of his kiss, underscore her determination to resist his control. D'Urberville's laugh at her resistance and his role as a passionate companion add complexity to their dynamic. As they near Trantridge, Tess's resolve is clear, despite the emotional toll of the journey.

1259  
1260  
1261  
1262  
1263  
1264  
1265