

# CovScore: Evaluation of Multi-Document Abstractive Title Set Generation

Anonymous ACL submission

## Abstract

blakhsdlkajhsd This paper introduces CovScore, an automatic reference-less methodology for evaluating thematic title sets, extracted from a corpus of documents. While such extraction methods are widely used, evaluating their effectiveness remains an open question. Moreover, some existing practices heavily rely on slow and laborious human annotation procedures. Inspired by recently introduced LLM-based judge methods, we propose a novel methodology that decomposes quality into five main metrics along different aspects of evaluation. This framing simplifies and expedites the manual evaluation process and enables automatic and independent LLM-based evaluation. As a test case, we apply our approach to a corpus of Holocaust survivor testimonies, motivated both by its relevance to title set extraction and by the moral significance of this pursuit. We validate the methodology by experimenting with naturalistic and synthetic title set generation systems and compare their performance with the methodology.<sup>1</sup>

## 1 Introduction

Getting a sense of a large corpus of documents is a challenging and taxing task for humans to carry out. Much NLP work has therefore focused on creating frameworks that simplify, organize, and summarize such corpora. One common approach seeks to model the salient themes represented by the corpora using lists of descriptors (henceforth, *title sets*). Examples of such frameworks include topic modeling (Abdelrazek et al., 2023), which defines topics as distributions over words, FrameNet semantic frames (Baker et al., 1998), which grounds the documents in a predefined ontology of schematized events and states, and ATOMIC (Sap et al., 2019), which models common-sense understanding

<sup>1</sup>An implementation of our automatic methodology will be made publicly available upon publication.

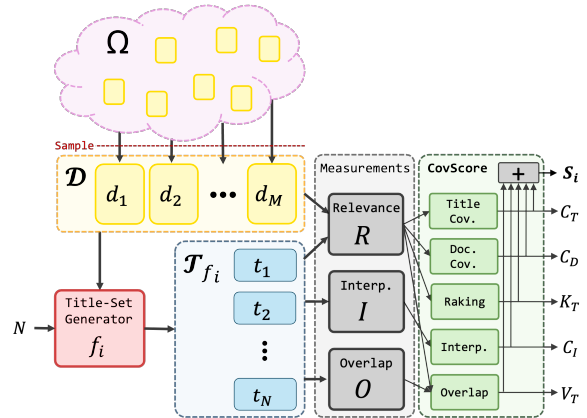


Figure 1: CovScore pipeline: A set of documents  $\mathcal{D}$  which is an accessible sample from the domain  $\Omega$ , is passed into a provided title set generation system  $f_i$ . The *Interpretability*, *Relevance* and *Overlap* signals are measured based on the set of resulting titles,  $\mathcal{T}_{f_i}$  and the document sample ( $\mathcal{D}$ ). The resulting signal is then used to compute the aspect-based scores and the aggregate score  $S_i$ .

of events using sub-event descriptions and cause-effect relations.

In recent years different solutions were proposed by the research community in an ongoing effort to automatically induce such title sets. Underlying these solutions is an attempt to automatically identify themes and their importance, followed by a title generation procedure intended to represent each such theme in a short, summarized, and easy-to-understand fashion. One well-known approach is Latent Dirichlet Allocation (LDA; Blei et al., 2003), which uses word clusters, rather than titles, as theme descriptors that are in turn used to identify themes. Other methods identify keywords, phrases, and sentences extracted from the text that may represent the themes and could serve as titles. Recently, alongside the rise of LLMs, solutions are shifting towards using generative models to implicitly identify the themes and output uniquely generated titles (e.g., Reuter et al., 2024; Garg et al.,

2021; Mishra et al., 2021). However, in contrast to the abundance of solutions, the literature lacks efficient evaluation methods, leaving the definition of what makes a “good” title set an open question.

We present a novel methodology for referenceless evaluation of title sets and show that it can be deployed manually or automatically. Acknowledging the drawbacks of using aggregate metrics (Burnell et al., 2023; Kasai et al., 2021), we report separate scores for each aspect, alongside an aggregate score.

We conduct a case study on a dataset of Holocaust survivor testimonies, collected by USC Shoah Foundation (SF).<sup>2</sup> Given the imminent passing of the last generation of Holocaust survivors, it is increasingly important that the testimonies they left be made accessible to Holocaust researchers and the public. However, due to the enormity of the collected databases (tens of thousands of testimonies), only a few of them are directly read. Our investigation will therefore support the development of stronger systems for sieving through these databases and providing a broader view of their major trends.

Other than the importance of studying these testimonies for Holocaust research, Holocaust testimonies provide an interesting test case due to the recounted common yet unique experiences. The dataset is partitioned into segments, which are manually annotated by SF with domain labels. Segments are then clustered according to their labeling to form constraint domains and are later used to generate title sets. The title sets are annotated by human annotators while measuring inter-annotator agreement (IAA) and showing correlation between human and automatic LLM-based labeling. We demonstrate the effectiveness of our methodology by experimenting with both naturalistic and synthetic title set generation systems and compare their performance by studying the intricate trade-offs existing between the different sets.

## 2 Related Work

### 2.1 Reference-Based Evaluation

Considering the general problem of free text evaluation, methods that assume an available annotated data source, most commonly rely on comparing the predicted and grounded texts. Traditionally, comparison metrics like BLEU (Papineni et al., 2002) and ROUGE (Lin, 2004) were commonly used.

<sup>2</sup><https://sfi.usc.edu/>

These methods assess the quality of the generated text by measuring N-gram overlap with the reference text. While convenient and widely used, these metrics primarily focus on surface-level similarities, often overlooking important semantic nuances, hindering the ability to truly capture the quality of the abstraction.

Newer metrics like BERTScore (Zhang et al., 2019), attempt to address this by leveraging Language Models like BERT (Devlin et al., 2018) to gauge semantic similarity. While semantic similarity methods offer some improvement over N-gram overlap, their performance can still be hampered in scenarios where context is lacking, such as when comparing titles without context. In addition, semantic similarity does not capture all aspects of interest. In the context of title set generation, it fails to gauge whether the titles themselves are interpretable or the effective size of the set. To attend to this problem, the evaluation process is often decomposed into different aspects that are measured separately (e.g., Kasai et al., 2021).

Nonetheless, the biggest hurdle for reference-based evaluation is the collection and annotation process, making such data scarce.

### 2.2 Referenceless Evaluation

Generally, referenceless evaluation metrics can be categorized into *extrinsic* and *intrinsic*. Extrinsic methods are valuable for assessing an output used as an intermediate step in a larger system (Suzuki and Fukumoto, 2014; Wu et al., 2024; Penta, 2022). However, it provides limited insight into the inherent quality of the output itself. In our work, we focus on intrinsic evaluation.

Intrinsic methods, such as Mimno et al. (2011), often exhibit weak correlation with human judgment (Stammbach et al., 2023). One such commonly used method utilizes the *intrusion* metric (Chang et al., 2009; Bhatia et al., 2018), which assesses the “coherence” of a title. This metric is used in LDA-like scenarios where a word cluster represents a title. In such cases, it is hypothesized that if a word cluster represents some induced theme, then the words it contains should be related. This property is scored based on the ability of an external validator (human or machine) to identify this intruder word (Stammbach et al., 2023). Since in this work we employ directly generated titles, this evaluation is irrelevant. However, this approach was recently adapted to the generative use case. In (Lior

158 et al., 2024), the intrusion task is used to evaluate  
159 the generated title sets by treating the whole set as  
160 a single cluster. Although this approach captures  
161 some of the aspects in a “good” title set, it lacks  
162 the direct grounding derived from our suggested  
163 aspects.

## 164 2.3 LM as a Judge

165 Another recently introduced line of work includes  
166 using “Judge” models as evaluators. At the cen-  
167 ter of this methodology is an attempt to leverage  
168 the strength of large models for automatically as-  
169 sessing the correctness of the output. Previously,  
170 the evaluation process relied on custom models  
171 specifically trained for each use-case (e.g., Bhatia  
172 et al., 2018; Gupta et al., 2014; Peyrard et al., 2017).  
173 However, training such models is difficult. Recog-  
174 nizing zero-shot and few-shot learning capabilities  
175 of LLMs (Brown et al., 2020), inspired some works  
176 (e.g., Fu et al., 2023; Huang et al., 2023; Lai et al.,  
177 2023; Kocmi and Federmann, 2023; Wang et al.,  
178 2023) to use LLMs as evaluators, instead of task-  
179 specific training.

180 Evaluating the correctness of a solution to a prob-  
181 lem is sometimes as difficult as solving the prob-  
182 lem itself. In our work, we show that reducing  
183 the evaluation into smaller measurements simpli-  
184 fies it, however further research is needed to better  
185 understand the trade-offs in such a simplification.

## 186 2.4 Manual Evaluation

187 Title set evaluation methods are often designed to  
188 allow for performing either machine or human an-  
189 notation (e.g., Chang et al., 2009; Lau et al., 2014;  
190 Nugroho et al., 2020). Coupled with the inherent  
191 difficulty of evaluating title sets, human evaluation  
192 is frequently favored (e.g., Chang et al., 2009; Lior  
193 et al., 2024). However, while flexible, it is also  
194 extremely costly and slow and therefore can only  
195 be done at a limited scale. Furthermore, common  
196 evaluation practices are hindered by subjective in-  
197 terpretations influenced by the annotators’ intrinsic  
198 biases and reading intents, the lack of context in  
199 titles, the difficulty in comparing diverse sets, and  
200 the cognitive burden of processing large amounts  
201 of information at once (Hoyle et al., 2021; Nugroho  
202 et al., 2020). We attempt to reduce the complex-  
203 ity and subjectivity of the annotation process for  
204 humans and machines alike.

## 205 3 The CovScore Methodology

206 The methodology presented in this work scores  
207 whether a given title set showcases salient themes  
208 in a corpus. To achieve this, we define a set of  
209 quantifiable evaluation aspects, defined in terms of  
210 three measurements. This framing simplifies direct  
211 evaluation processes and therefore enables a more  
212 reliable and scalable scoring system.

### 213 3.1 Formal Setting & Definitions

214 Given a *domain*  $\Omega$ , defined by the contents of a set  
215 of semantically related documents  $\{d : d \in \Omega\}$ .  
216 We further assume that only a *sample*  $\mathcal{D} \subseteq \Omega$  of  
217 the documents are available, where  $M = |\mathcal{D}|$ . The  
218 sample is fed as input to a system  $f_i(\mathcal{D}, N)$  that  
219 produces a *title set*,  $\mathcal{T}_{f_i}$ . Each *title*  $t \in \mathcal{T}_{f_i}$  is a  
220 string and  $N = |\mathcal{T}_{f_i}|$  is the size of the resulting  
221 title set. The methodology returns a score for each  
222 aspect:  $C_T, C_D, K_T, C_I, V_T$  and an aggregated  
223 overall score  $S_i$ . Each such score is referred to as  
224 an evaluation *metric* and takes values in  $[0, 1]$ .

225 Aspects are defined in terms of *measurements*  
226 that can be directly annotated, i.e., be reliably an-  
227 notated by humans. This implies that they are well-  
228 defined and not exceedingly cognitively taxing. A  
229 measurement is defined as a function of the titles  
230 and the sample and can be performed by a human  
231 or a machine. See Fig. 1 for a schematized view of  
232 the methodology.

### 233 3.2 Defining the Quality of a Title Set

234 Most commonly, title sets are used as a means to  
235 simplify, organize, and summarize sizable collec-  
236 tions of documents. Generally, systems achieve this  
237 goal in three steps: *identifying* recurring themes,  
238 *generating titles* that capture the essence of each  
239 theme, and determining their *importance* (Abdel-  
240 razek et al., 2023; AlSumait et al., 2009; Song et al.,  
241 2009). We use this formulation to decompose the  
242 quality of a title set into the following aspects:

#### 243 Aspect 1: Interpretability

244 Assesses whether titles in the set describe some  
245 theme in the corpus. A title describes a theme if it  
246 is largely unambiguously interpreted as that theme  
247 by the annotators. For example, within experiences  
248 of deportation during the Holocaust, a title like  
249 “sadness” can be difficult to decipher. The range of  
250 emotions present during such an experience makes  
251 it hard to understand what specific aspect of the ex-  
252 perience the title is meant to highlight and therefore

is not interpretable. Formally, we define:

$$C_I = \frac{1}{N} \sum_{t \in \mathcal{T}_{f_i}} I(t) \quad (1)$$

$I(t)$  denotes the interpretability measurement that accepts a title and outputs a score in  $[0, 1]$  for the degree of clarity and understandableness of the title to a human reader.

### Aspect 2: Coverage

Assesses whether the title set covers the sample. To quantify the coverage we define two competing metrics, both rely on the relevance measurement denoted by  $R(t, d)$ . The measurement scores the relation between the generated titles and the themes they may represent. For each title-document pair, the function returns a score in  $[0, 1]$  expressing the relevance of the title to the document, evaluating the title in context.

**Title Coverage.** indicates whether the titles in the set capture the major themes in the corpus. A major theme is a theme that recurs broadly across the corpus. Hence, a title that is related to many documents in the corpus (covers the corpus), is a title describing a major theme. For example, within experiences of deportation, the title “*Transportation to Concentration Camps*” is a major theme since it is likely to cover most, if not all, deportation experiences. To quantify this aspect, the metric computes the mean relatedness of the titles to the documents. Formally,

$$C_T = \frac{1}{N} \frac{1}{M} \sum_{t \in \mathcal{T}_{f_i}, d \in \mathcal{D}} R(t, d) \quad (2)$$

**Document Coverage.** indicates whether the title set contains titles that are not represented widely in the sample, preventing title sets from being limited to general themes and ensuring a more thorough representation of the corpus. For example, this will enable the inclusion of a more specific title like “*Transportation by a Wagon*” in the set. A quantifiable lower bound is set by the least-covered document. This means identifying the document with the lowest relevance score and its most relevant title. That score will be our measure of title set coverage. Formally,

$$C_D = \min_{d \in \mathcal{D}} \max_{t \in \mathcal{T}_{f_i}} \{R(t, d)\} \quad (3)$$

### Aspect 3: (non-)Overlap

Assesses whether the titles represent separate themes by capturing whether multiple titles overlap by the themes they induce. For example, the titles “*Transportation to Concentration Camps*” and “*Transportation by a Wagon*” may refer to the same theme. In quantifying the metric we consider both the overlap in definition and in the covered documents. Formally:

$$V_T = \frac{1}{N} \sum_{t \in \mathcal{T}_{f_i}} [1 - \max(v_{\text{def}}(t), v_{\text{cov}}(t))] \quad (4)$$

$$v_{\text{def}}(t) = \max_{t' \in \mathcal{T}_{f_i}, t \neq t'} O(t, t') \quad (5)$$

$$v_{\text{cov}}(t) = \max_{t' \in \mathcal{T}_{f_i}, t \neq t'} \sum_{d \in \mathcal{D}} R(t, d) \cdot R(t', d) \quad (6)$$

$R(t, d)$  denotes the relevance measurement defined above and  $O(t_1, t_2)$  as the overlap measurement. This measurement scores the overlap in the definition, it receives a pair of titles and outputs a score in  $[0, 1]$  for the degree that the themes expressed by the two titles overlap.

### Aspect 4: Inner-Order

Assesses whether the titles in the set are ordered by their importance. In some cases, although not all, the order of topics reflects importance, where more important topics precede less important ones in the set. For example, a title like “*Transportation to Concentration Camps*” should be ordered before “*Transportation by a Wagon*”. If the title set is well-ordered, its inner order should reflect the order of the topic’s importance. Formally,

$$K_T = \max(0, \tau(\mathcal{T}_{f_i}, \mathcal{T}')) \quad (7)$$

where  $\tau(\cdot)$  is the Kendall  $\tau$  ranking correlation coefficient (Kendall, 1948), and  $\mathcal{T}'$  is a re-ordering of  $\mathcal{T}_{f_i}$  according to the mean relevance:

$$r_t = \frac{1}{M} \sum_{d \in \mathcal{D}} R(t, d) \quad (8)$$

$R(t, d)$  denotes the relevance measurement.

## 4 Manual Evaluation

To support our claims we have created 3 manual annotation tasks where human annotators were asked to annotate the interpretability, relevance, and overlap measurements. The annotated data was then used as test data. Each measurement was carried out by 2-4 annotators with full overlap, in order to measure IAA.

### 4.1 Data

A large enough collection of sets of documents, where each such set represents a relatively constrained domain, is hard to come by. We have therefore opted to use the Holocaust Survivor Testimonies dataset collected by SF. This dataset is comprised of stories recounted by survivors based on their unique experiences and perspectives during the Holocaust. Each testimony naturally describes different experiences, but many of the themes do recur, albeit in a variety of circumstances, times, and places. We are further motivated by the recent use of this dataset in recent computational modeling work (Wagner et al., 2022, 2023).

The testimonies (see examples in Table 10) were collected as part of an oral interview in English between a survivor and an interviewer. The recordings were later transcribed into text. Since the story is told as part of an interview, the data is segmented according to the speaker sides, where most of the time survivors share their experiences while the interviewer guides the testimony with questions. Testimony lengths range from 2609 to 88105 words, with a mean length of 23536 words (Wagner et al., 2022).

In this work, we use an existing labeling of the dataset performed by SF, which identifies testimony segments that are related across survivors. The labeling system is based on a pre-defined human-generated hierarchical ontology where segments of roughly 1 minute (of audio time) were labeled with one or more ontology classes. For our purposes, we have clustered segments from multiple testimonies that share a label, to form *domains* (see §3.1). These domains represent common experiences with shared themes and therefore could be used in our experiments. A single testimony may contain multiple non-consecutive segments sharing a label. For this reason, we define a document as a concatenation of all segments in a single testimony that shares a label.

For this work, we selected 21 domains that are

relatively constrained. See Table 4 for data distributions and domain labels. The documents in each domain were then used to generate sub-titles using GPT3.5 (see Appendix B for the generation prompt). The titles as well as a random sample of 10 in-domain documents, were used as annotation data for each measurement. See final item counts in Table 1. Since contemporary LLMs rarely output uninterpretable content, for the interpretability measurement we synthetically increased the number of uninterpretable titles (negative items), to allow effective computation of IAA. As a solution, we used GPT4 to corrupt some of the generated titles. The corruption prompt as well as examples for corrupted titles can be found in Appendix B and Table 8.

### 4.2 Methodology

To annotate the data, we recruited 4 English fluent speakers with no previous expertise in Holocaust studies. The annotators were asked to perform the 3 measurements described in §3.2 and repeat the process for each domain. The annotators received guidance both frontally and through written annotation guidelines. Before each session, the annotators were asked to read all the documents in the sample that they were given (which contained 10 documents) to become familiar with the domain. The same documents were then used as references in the annotation procedure itself. Importantly, the annotators were asked to make no assumptions based on previous knowledge that did not appear in the context of the sample. During the annotation process, we followed the conclusions from Graham et al. (2013) and used Continuous Scale Rating on the scale of [0 – 100].

Finally, to gain confidence in the results, we maintained full item overlap between all the annotators and measured IAA for each measurement. We employed Krippendorff- $\alpha$  (Krippendorff, 2011) to measure the agreement. Our results indicate high levels of agreement across the different measurements. See Table 1.

## 5 LLMs as Automatic Evaluators

In this section, we examine off-the-shelf LLMs in their ability to produce title set annotations. Specifically, we test LLMs for their ability to reliably simulate human judgments on the interpretability, relevance, and overlap measurements.

Measurement	# Items	# Anno.	Agreement
Interp.	550	3	0.66
Relevance	1583	4	0.67
Overlap	464	2	0.78

Table 1: Agreement achieved on each annotation measurement, including the number of items tagged, number of annotator participants, and the resulting Krippendorff- $\alpha$  score. All items were tagged by all annotators.

## 5.1 Experimental Setup

In the following experiment we have used the annotated data collected in §4 as a test set to evaluate the performance of popular LLMs as predictors, including GPT4 (Achiam et al., 2023), GPT3.5 (Brown et al., 2020), Mixtral (Jiang et al., 2024) and LLAMA-3 (see model versions in Appendix C.1). For the last two, we used both no quantization and 4-bit quantization. For each measurement, a prompt was written (see Appendix B) for querying the model to predict a score for each data point on a scale of 1 – 100, based on the measurement definitions in §3.2.

## 5.2 Results

Table 2 reports the Spearman correlation (Spearman, 1961) between LLM’s predictions and the mean human score. The results show that LLMs can simulate the human annotations achieving high overall correlation. Even though the best model varies in each measurement, we note the GPT4’s dominance, as well as LLAMA-3 (70B) with no quantization for being a reasonable open-sourced alternative. To further substantiate this claim we report additional correlation measures in Table 6, showing the same conclusions. In addition, Table 7 presents the correlation between each annotator and the mean human score used as a test set. The high correlation further stresses the reliability of our conclusions.

## 6 Validation

To establish the validity of our methodology, we conduct a validation study. Often, validation of a new evaluation metric involves scoring the output of multiple systems and demonstrating alignment with human preferences (e.g., Papineni et al., 2002). However, as previously discussed in §2, human annotation of title sets is unreliable. Instead, we compare the methodology’s score of title sets, where we

have solid intuition as to their “true” performance. The study demonstrates that our methodology performs as expected in edge cases and effectively identifies the trade-offs between title sets. Inspired by Burnell et al. (2023), we report each aspect separately. However, acknowledging the benefits of aggregated scores, §6.3 presents a single summarized score for simpler system-level comparisons.

## 6.1 Methodology

We designed and implemented 13 title set generation systems, categorized into four groups based on their underlying generation approach. To simplify the validation process we assume each system outputs a constant number of 10 titles. Each domain is represented with a random sample of 8 documents. We use *Meta-Llama-3-8B-Instruct* as the judge model, selected for its cost-effectiveness.

We come to present the title set generation systems. Examples of outputs can be found in appendix 11.

### Baselines.

1. **Random-Letters** produces titles comprised of random sequences of English letters.
2. **Random-Words** generates titles by combining random, yet real, English words. By using actual words we expect improved results compared to Random-Letters.
3. **Domain-Name** uses the domain labels assigned by human annotators (see §4.1). Title sets are created by assigning the same domain label to every title within that set.

**Naïve LDA-Based.** utilizes LDA (Blei et al., 2003), a widely adopted approach for extracting topics as word distributions, to generate titles. We specifically use *gensim*’s implementation of LDA (Rehurek et al., 2011). The resulting word distributions are transformed into titles using the following approaches.

1. **LDA+Prefix** titles are represented by a quoted, comma-separated list of the topic’s top  $k$  words. A prefix “*The theme defined by the following set of words:*” is then prepended to the string.
2. **LDA+GPT4** titles are generated by prompting GPT4 with the topic’s top  $k$  words (see prompt in appendix B).

Both methods use  $k \in \{1, 10, 50\}$ .

**LLM-Based.** titles are generated by prompting LLMs.

Model	Quantization	Relevance	Overlap	Interpretability
GPT 4	-	<b>0.66</b>	0.86	0.63
GPT 3.5	-	0.50	0.79	<b>0.73</b>
LLAMA 3 (8B)	None	0.45	0.85	0.54
LLAMA 3 (70B)	4-bit	0.48	0.24	0.29
LLAMA 3 (70B)	None	<u>0.62</u>	<b>0.87</b>	<u>0.66</u>
Mixtral (8x7B)	4-bit	0.43	0.83	0.65
Mixtral (8x7B)	None	0.50	0.73	0.65

Table 2: Spearman correlation between LLM and mean human annotations. The best overall model for each measurement is boldfaced and the best open-source alternative is underlined.

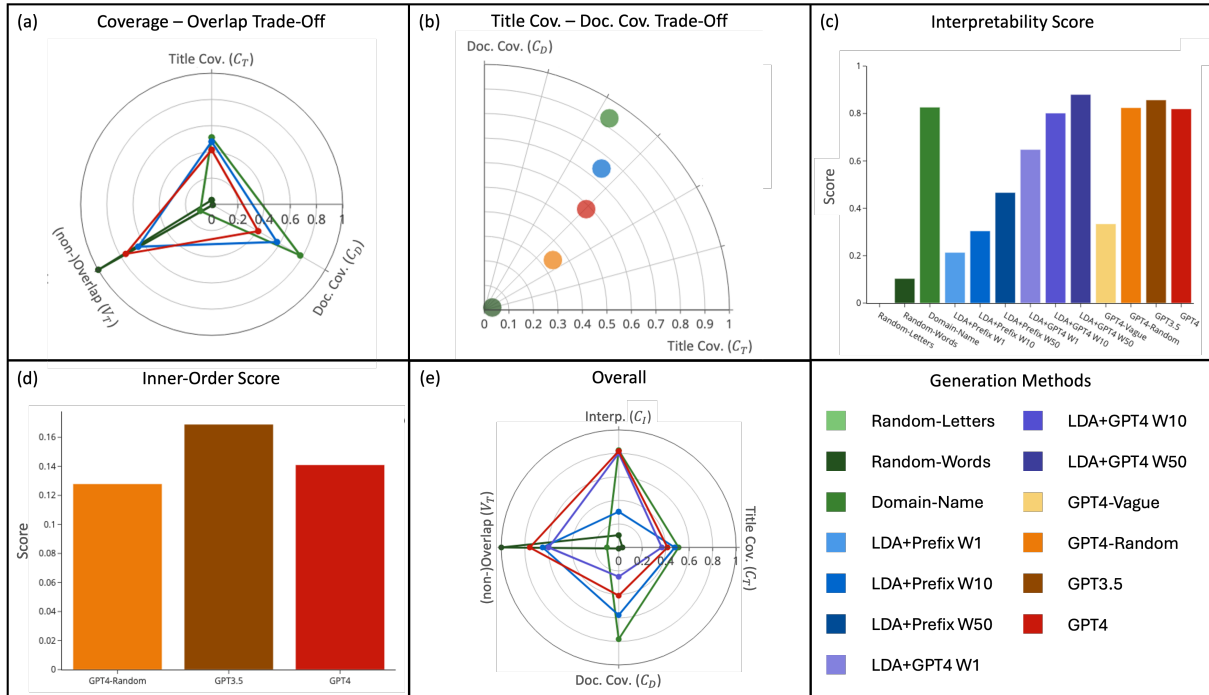


Figure 2: Validation study results; (a) shows the trade-off existing between Coverage aspects (Title Cov. and Doc. Cov.) and the non-Overlap aspect; (b) show the trade-off between Title Coverage and Document Coverage; (c) shows the Interpretability scores across systems; (d) shows the Inner-Order scores achieved by LLM based systems; (e) depicts an overall comparison of representing systems.

1. **GPT** leverages OpenAI’s GPT (Brown et al., 2020) to generate title sets. The model is used to extract common titles from a random sample of documents within a specific domain. Both, the sampling and title extraction are repeated  $N$  times, followed by a map-reduce process applied to consolidate the various extracted title sets into a single final set. We use both GPT3.5 and GPT4.
2. **GPT4-Random** samples random titles uniformly from the union of title sets from all domains, as generated by GPT4.
3. **GPT4-Vague** uses the title corruption procedure

from §4.1 to corrupt all of the titles extracted by GPT4.

## 6.2 Results

Figures 2(a)-(e) demonstrate the intricate trade-offs existing between the different generated title sets. Figure 2(a) shows the most prominent trade-off, arising between the Coverage aspect (Title Coverage and Document Coverage) and the non-Overlap aspect. Figure 2(b) presents the more subtle but nonetheless central trade-off that exists between Title Coverage and Document Coverage. Figure 2(c) shows Interpretability scores across systems.

Generation Method	Aggregated Score
Random-Letters	0.00
Random-Words	0.01
Domain-Name	0.00
LDA+Prefix W1	0.02
LDA+Prefix W10	0.13
LDA+Prefix W50	<u>0.19</u>
LDA+GPT4 W1	0.03
LDA+GPT4 W10	0.10
LDA+GPT4 W50	0.21
GPT4-Vague	0.05
GPT4-Random	0.15
GPT3.5	<b>0.22</b>
GPT4	0.19

Table 3: Aggregate scores achieved by each system. The highest scoring system is boldfaced, while the best system that does not use an LLM is underlined.

The bars in the figure are color-coded so it will be easier to distinguish between the underlying system groups. Figure 2(d) shows the Inner-Order scores achieved by LLM-based systems. Finally, Figure 2(e) depicts an overall comparison of representing systems from each generation group, considering all metrics other than Inner-Order. The results show that our methodology successfully captures the intricate trade-offs, substantiating its validity. A more thorough analysis of the results can be found in Appendix F.

### 6.3 An Aggregate Score

Along with the individual metrics, we additionally propose a single aggregate score. Using such a score could be advantageous in some scenarios such as for quick comparison between systems and as a reward function for training title set generation models. We choose the harmonic mean function to aggregate the different metrics, formally:

$$S_i(\mathcal{T}_{f_i}, \mathcal{D}) = \frac{|A_i|}{\sum_{\alpha \in A_i} \frac{1}{\alpha}} \quad (9)$$

where  $A_i$  is the set of aspect scores for the title set  $\mathcal{T}_{f_i}$ . Table 3 shows how the different systems fare on the aggregate metric. GPT3.5 outperforms all other methods, but only by a small margin.

## 7 Conclusion

We have formulated the problem of title set evaluation as a theme coverage problem and presented a

methodology for evaluating title sets by decomposing the problem into multiple quantifiable aspects. We used Holocaust survivor testimonies as a test case for studying the methodology and showed its usefulness for manual evaluation by achieving high levels of IAA. We further showed that the proposed methodology can be automated by simulating human annotations with judge models.

To validate the application of this methodology, we compared a range of systems and baselines, where the true relative order between at least a subset of them in each aspect, was clear. The study showed that our methodology successfully reflects the intricate trade-offs and relative quality of these systems, validating it as a system-level comparison metric.

Given the centrality of the task of title set extraction, and the great difficulty in evaluating the task reliably, we hope that the methodology proposed here will assist in the development of demonstrably stronger title set extraction systems.

## Limitations

The limitations of this work could be separated into data-related, and model-related limitations. First, our experiments are restricted to a single type (Holocaust survivor testimonies). However, we do not tailor our method in any way to this type, so we expect that our findings will not be directly influenced by it. Second, during the annotation process, the annotators are only presented with a small sample (10 documents) from each domain. In this work, we do not assess whether this sample sufficiently covers the entirety of the domain, which could bias the annotation process. Third, multiple parts of the same experience may be scattered throughout the testimony. To handle this problem we have defined a document as the concatenation of all of those segments. However, each such segment may have been told in a different context, which could influence the interpretation of the text. Moreover, the prior ontology labeling of the segments was done on segments of constant 1-minute length. This coarse segmentation may cause unrelated information to be included in the segment, as well as a misplacement of small but crucial segments.

Other limitations stem from the use of LLMs. First, LLMs are black box models, often trained by commercial companies that do not disclose their inner workings, limiting the replicability of the results. Second, these models are extremely ex-



615	pensive to use, either as services or by running	David M Blei, Andrew Y Ng, and Michael I Jordan.	665
616	them locally on multiple high-end GPUs. Since	2003. Latent dirichlet allocation. <i>Journal of machine</i>	666
617	our method requires employing such models, the	<i>Learning research</i> , 3(Jan):993–1022.	667
618	high cost may pose a limitation in some contexts.	Tom Brown, Benjamin Mann, Nick Ryder, Melanie	668
619	However, we expect this cost to rapidly decline in	Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind	669
620	the near future.	Neelakantan, Pranav Shyam, Girish Sastry, Amanda	670
621	<b>Ethics Statement</b>	Askell, et al. 2020. Language models are few-shot	671
622	In this work, we have been aided by in-house an-	learners. <i>Advances in neural information processing</i>	672
623	notators, who were employed by the university and	<i>systems</i> , 33:1877–1901.	673
624	given instructions beforehand. We abided by the	Ryan Burnell, Wout Schellaert, John Burden, Tomer D	674
625	instructions provided by the SF. We note that the	Ullman, Fernando Martinez-Plumed, Joshua B	675
626	witnesses identified themselves by name, and so the	Tenenbaum, Danaja Rutar, Lucy G Cheke, Jascha	676
627	testimonies are open and not anonymous by design.	Sohl-Dickstein, Melanie Mitchell, et al. 2023. Re-	677
628	We intend to release our scripts, but those will not	think reporting of evaluation results in ai. <i>Science</i> ,	678
629	include any of the data received from the archives;	380(6641):136–138.	679
630	the data and trained models used in this work will	Jonathan Chang, Sean Gerrish, Chong Wang, Jordan	680
631	not be given to a third party without the consent of	Boyd-Graber, and David Blei. 2009. Reading tea	681
632	the relevant archives. The testimonies can be made	leaves: How humans interpret topic models. <i>Ad-</i>	682
633	accessible for browsing and research by requesting	<i>vances in neural information processing systems</i> , 22.	683
634	permission from the SF archive. Holocaust tes-	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	684
635	timonies are by nature, sensitive material. Users	Kristina Toutanova. 2018. Bert: Pre-training of deep	685
636	should exercise caution when applying LLMs for	bidirectional transformers for language understand-	686
637	Holocaust testimonies, to avoid incorrect represen-	ing. <i>arXiv preprint arXiv:1810.04805</i> .	687
638	tation of the told stories.	David Freedman, Robert Pisani, and Roger Purves.	688
639	<b>References</b>	2007. Statistics (international student edition).	689
640	Aly Abdelrazek, Yomna Eid, Eman Gawish, Walaa Med-	<i>Pisani, R. Purves, 4th edn. WW Norton &amp; Company,</i>	690
641	hat, and Ahmed Hassan. 2023. Topic modeling al-	<i>New York.</i>	691
642	gorithms and applications: A survey. <i>Information</i>	Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei	692
643	<i>Systems</i> , 112:102131.	Liu. 2023. Gptscore: Evaluate as you desire. <i>arXiv</i>	693
644	Josh Achiam, Steven Adler, Sandhini Agarwal, Lama	<i>preprint arXiv:2302.04166</i> .	694
645	Ahmad, Ilge Akkaya, Florencia Leoni Aleman,	Krishna Garg, Jishnu Ray Chowdhury, and Cornelia	695
646	Diogo Almeida, Janko Altmenschmidt, Sam Altman,	Caragea. 2021. Keyphrase generation beyond the	696
647	Shyamal Anadkat, et al. 2023. Gpt-4 technical report.	boundaries of title and abstract. <i>arXiv preprint</i>	697
648	<i>arXiv preprint arXiv:2303.08774</i> .	<i>arXiv:2112.06776</i> .	698
649	Loulwah AlSumait, Daniel Barbará, James Gentle, and	Yvette Graham, Timothy Baldwin, Alistair Moffat, and	699
650	Carlotta Domeniconi. 2009. Topic significance rank-	Justin Zobel. 2013. Continuous measurement scales	700
651	ing of lda generative models. In <i>Machine Learn-</i>	in human evaluation of machine translation. In <i>Pro-</i>	701
652	<i>ing and Knowledge Discovery in Databases: Euro-</i>	<i>ceedings of the 7th Linguistic Annotation Workshop</i>	702
653	<i>pean Conference, ECML PKDD 2009, Bled, Slovenia,</i>	<i>and Interoperability with Discourse</i> , pages 33–41.	703
654	<i>September 7-11, 2009, Proceedings, Part I 20</i> , pages	Pooja Gupta, Nisheeth Joshi, and Iti Mathur. 2014.	704
655	67–82. Springer.	Quality estimation of machine translation outputs	705
656	Collin F Baker, Charles J Fillmore, and John B Lowe.	through stemming. <i>arXiv preprint arXiv:1407.2694</i> .	706
657	1998. The berkeley framenet project. In <i>COLING</i>	Alexander Hoyle, Pranav Goel, Andrew Hian-Cheong,	707
658	<i>1998 Volume 1: The 17th International Conference</i>	Denis Peskov, Jordan Boyd-Graber, and Philip	708
659	<i>on Computational Linguistics</i> .	Resnik. 2021. Is automated topic model evaluation	709
660	Shraey Bhatia, Jey Han Lau, and Timothy Baldwin.	broken? the incoherence of coherence. <i>Advances</i>	710
661	2018. Topic intrusion for automatic topic model eval-	<i>in neural information processing systems</i> , 34:2018–	711
662	uation. In <i>Proceedings of the 2018 Conference on</i>	2033.	712
663	<i>Empirical Methods in Natural Language Processing</i> ,	Fan Huang, Haewoon Kwak, and Jisun An. 2023. Is	713
664	pages 844–849.	chatgpt better than human annotators? potential and	714
		limitations of chatgpt in explaining implicit hate	715
		speech. In <i>Companion proceedings of the ACM web</i>	716
		<i>conference 2023</i> , pages 294–297.	717

718	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, et al. 2024. Mixtral of experts. <i>arXiv preprint arXiv:2401.04088</i> .	772	Kishore Papineni, Salim Roukos, Todd Ward, and Weijing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In <i>Proceedings of the 40th annual meeting of the Association for Computational Linguistics</i> , pages 311–318.	776
723	Wenxiang Jiao, Wenxuan Wang, Jen-tse Huang, Xing Wang, Shuming Shi, and Zhaopeng Tu. 2023. Is chatgpt a good translator? yes with gpt-4 as the engine. <i>arXiv preprint arXiv:2301.08745</i> .	777	Antonio Penta. 2022. Enhance topics analysis based on keywords properties. <i>arXiv preprint arXiv:2203.04786</i> .	779
727	Jungo Kasai, Keisuke Sakaguchi, Lavinia Dunagan, Jacob Morrison, Ronan Le Bras, Yejin Choi, and Noah A Smith. 2021. Transparent human evaluation for image captioning. <i>arXiv preprint arXiv:2111.08940</i> .	780	Maxime Peyrard, Teresa Botschen, and Iryna Gurevych. 2017. Learning to score system summaries for better content selection evaluation. In <i>Proceedings of the Workshop on New Frontiers in Summarization</i> , pages 74–84.	784
732	Maurice George Kendall. 1948. Rank correlation methods. <i>Oxford University Press</i> .	785	Radim Rehurek, Petr Sojka, et al. 2011. Gensim—statistical semantics in python. Retrieved from <i>gensim.org</i> .	787
734	Tom Kocmi and Christian Federmann. 2023. Large language models are state-of-the-art evaluators of translation quality. <i>arXiv preprint arXiv:2302.14520</i> .	788	Arik Reuter, Anton Thielmann, Christoph Weisser, Sebastian Fischer, and Benjamin Säfken. 2024. Gp-topic: Dynamic and interactive topic representations. <i>arXiv preprint arXiv:2403.03628</i> .	791
737	Klaus Krippendorff. 2011. Computing krippendorff’s alpha-reliability. In <i>Departmental Papers (ASC)</i> , page 43.	792	Maarten Sap, Ronan Le Bras, Emily Allaway, Chandra Bhagavatula, Nicholas Lourie, Hannah Rashkin, Brendan Roof, Noah A Smith, and Yejin Choi. 2019. Atomic: An atlas of machine commonsense for if-then reasoning. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 33, pages 3027–3035.	798
740	Huiyuan Lai, Antonio Toral, and Malvina Nissim. 2023. Multidimensional evaluation for text style transfer using chatgpt. <i>arXiv preprint arXiv:2304.13462</i> .	799	Yangqiu Song, Shimei Pan, Shixia Liu, Michelle X Zhou, and Weihong Qian. 2009. Topic and keyword re-ranking for lda-based topic modeling. In <i>Proceedings of the 18th ACM conference on Information and knowledge management</i> , pages 1757–1760.	803
743	Jey Han Lau, David Newman, and Timothy Baldwin. 2014. Machine reading tea leaves: Automatically evaluating topic coherence and topic model quality. In <i>Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 530–539.	800	Charles Spearman. 1961. The proof and measurement of association between two things. <i>Am. J. Psychol.</i> , 15:72.	806
746	Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In <i>Text summarization branches out</i> , pages 74–81.	801	Dominik Stammbach, Vilem Zouhar, Alexander Hoyle, Mrinmaya Sachan, and Elliott Ash. 2023. Re-visiting automated topic model evaluation with large language models. <i>arXiv preprint arXiv:2305.12152</i> .	810
749	Gili Lior, Yoav Goldberg, and Gabriel Stanovsky. 2024. Leveraging collection-wide similarities for unsupervised document structure extraction. <i>arXiv preprint arXiv:2402.13906</i> .	802	Yoshimi Suzuki and Fumiyo Fukumoto. 2014. Detection of topic and its extrinsic evaluation through multi-document summarization. In <i>Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 241–246.	816
752	David Mimno, Hanna Wallach, Edmund Talley, Miriam Leenders, and Andrew McCallum. 2011. Optimizing semantic coherence in topic models. In <i>Proceedings of the 2011 conference on empirical methods in natural language processing</i> , pages 262–272.	804	Eitan Wagner, Renana Keydar, and Omri Abend. 2023. Event-location tracking in narratives: A case study on holocaust testimonies. In <i>The 2023 Conference on Empirical Methods in Natural Language Processing</i> .	820
755	Prakhar Mishra, Chaitali Diwan, Srinath Srinivasa, and Gopalakrishnan Srinivasaraghavan. 2021. Automatic title generation for text with pre-trained transformer language model. In <i>2021 IEEE 15th International Conference on Semantic Computing (ICSC)</i> , pages 17–24. IEEE.	805	Eitan Wagner, Renana Keydar, Amit Pinchevski, and Omri Abend. 2022. Topical segmentation of spoken narratives: A test case on holocaust survivor testimonies. <i>arXiv preprint arXiv:2210.13783</i> .	824
756	Robertus Nugroho, Cecile Paris, Surya Nepal, Jian Yang, and Weiliang Zhao. 2020. A survey of recent methods on deriving topics from twitter: algorithm to evaluation. <i>Knowledge and information systems</i> , 62:2485–2519.	806		

- 825 Jiaan Wang, Yunlong Liang, Fandong Meng, Zengkui  
826 Sun, Haoxiang Shi, Zhixu Li, Jinan Xu, Jianfeng  
827 Qu, and Jie Zhou. 2023. Is chatgpt a good nlg  
828 evaluator? a preliminary study. *arXiv preprint*  
829 *arXiv:2303.04048*.
- 830 Xiaobao Wu, Thong Nguyen, and Anh Tuan Luu. 2024.  
831 A survey on neural topic models: Methods, applica-  
832 tions, and challenges. *Artificial Intelligence Review*,  
833 57(2):1–30.
- 834 Xianjun Yang, Yan Li, Xinlu Zhang, Haifeng Chen, and  
835 Wei Cheng. 2023. Exploring the limits of chatgpt  
836 for query or aspect-based text summarization. *arXiv*  
837 *preprint arXiv:2302.08081*.
- 838 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q  
839 Weinberger, and Yoav Artzi. 2019. Bertscore: Eval-  
840 uating text generation with bert. *arXiv preprint*  
841 *arXiv:1904.09675*.

## A Data Distributions

Overall, in the data processing stage, we have extracted 572 different domains, where each domain contains 1-999 documents with an average of 105 documents and an overall mean document length of 86 sentences. For our purposes, we have selected a subset of 21 domains. Table 4 depicts the labels given to each one of these domains by SF, the number of documents it contains, and the mean length of a document in the domain. Figure 3 shows an overall distribution of all the available domains. Table 5 includes examples of testimony segments and their corresponding ontology labels assigned by SF.

Experience	# Documents	Ave. length (sentences)
Deportation To Concentration Camps	308	41.5
Family Interactions	900	124.9
Living Conditions	815	101.1
Forced Marches	345	51.6
Jewish Religious Observances	700	83.7
Anti-Jewish Regulations	597	49.9
Antisemitism	672	55.0
Armed Forces	541	70.5
Food and Drink	449	61.9
Forced Labor	530	162.8
Hiding	450	118.7
Housing Conditions	356	57.3
Immigration	633	113.2
Jewish Holidays	503	62.2
Kapos	138	64.4
Liberation	567	36.3
Military Activities	551	71.3
Post-Liberation Recovery	398	42.6
Sanitary and Hygienic Conditions	178	39.4
Soldiers	621	64.6
Transportation Routes	347	40.8

Table 4: Domain data distributions. Each domain is labeled by USC’s annotators. Each document is a concatenation of all segments in a testimony that were labeled as belonging to this experience.

## B LLM Prompts

Throughout our work we have used the following prompts when employing LLMs:

### Relevance Score

#### System Prompt:

```
You are a helpful Holocaust researcher assistant. You will
perform the following instructions as best as you can.
You will be presented with a topic and a text. Rate on a
scale of 1 to {max-rate} whether the topic describes a part
of the text ("1" = does not describe, "{mid-rate}" = somewhat
describes, "{max-rate}" = describes well).
```

```
Provide reasoning for the rate in one sentence only.
```

```
Please output the response in the following JSON format:
```

```
{
```

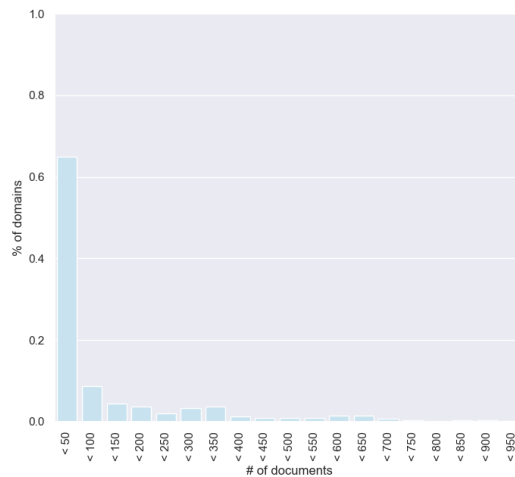


Figure 3: Size distribution of domains in terms of number of documents. Note that most domains contain less than 50 documents.

```
"rate": <rate> 863
"reasoning": <reasoning> 864
} 865
```

User Prompt: 866

```
Topic: "{topic}", 867
Text: ""{document}"" 868
```

**non-Overlap Score** 869

System Prompt: 870

You are a helpful Holocaust researcher assistant. You will 871  
perform the following instructions as best as you can. You 872  
will be presented with two topics: topic1 and topic2. Rate 873  
on a scale of 1 to {max-rate} whether topic1 have the same 874  
meaning as topic2 ("0" = different meaning, "{mid-rate}" = 875  
somewhat similar meaning, "{max-rate}" = same meaning). Pro- 876  
vide reasoning for the rate in one sentence only. 877

Please output the response in the following JSON format: 878

```
{ 879
"rate": <rate> 880
"reasoning": <reasoning> 881
} 882
```

User Prompt: 883

```
topic1: "{topic1}", 884
topic2: "{topic2}" 885
```

**Interpretability Score** 886

System Prompt: 887

You are a helpful Holocaust researcher assistant. You will 888  
perform the following instructions as best as you can. You 889

890 will be presented with a title representing a topic. Rate on  
891 a scale of 1 to {max-rate} whether the topic represented by  
892 the title is interpretable to humans ("0" = not interpretable,  
893 "{mid-rate}" = somewhat interpretable, "{mid-rate}" = easily  
894 interpretable). Provide reasoning for the rate in one sen-  
895 tence only.

896 Please output the response in the following JSON format:

```
897 {  
898   "rate": <rate>  
899   "reasoning": <reasoning>  
900 }
```

901 User Prompt:

```
902   topic1: "{topic1}",  
903   topic2: "{topic2}"
```

### 904 **Title Corruption**

905 Following is a title, that represents a theme. Corrupt the  
906 title such that the theme could not be easily understood by a  
907 human reader. The title must be short and readable. You may  
908 make the title vague, metaphorical, or designed to pique cu-  
909 riosity without directly revealing the topic

910  
911 Title: {title}

912 New Title:

### 913 **LDA Word-Cluster Conversion to Titles**

914 Following is a list of words extracted with an LDA model, rep-  
915 resenting an LDA cluster. Please give a title to the topic  
916 this cluster represents

917  
918 Cluster words: [{"", ".join(words)}]

919 Title:

### 920 **LLM-based title set Generation**

#### 921 Single title set Generation

922 You are a Holocaust researcher. You will perform the follow-  
923 ing instructions as best as you can. You will be displayed  
924 multiple texts. Please make a list of {NUM-TOPICS} unique top-  
925 ics that are common for all of the following texts. Make sure  
926 that the topics are general in their description, relevant to  
927 the texts, distinct, comprehensive, specific, interpretable,  
928 and short.

929 Desired format:

```
930  
931 1. <topic1>  
932 2. <topic2>  
933 3. <topic3>  
934 ...
```

935  
936 Text 1: <text1>

Text 2: <text2> 937  
 Text 3: <text3> 938  
 Text 4: <text4> 939  
 ... 940  
 Text <N>: <textN> 941  
 942

Sets Aggregation 943

You will be presented with a set of topic titles. Please 944  
 choose {NUM-TOPICS} distinct titles that best describe the 945  
 set. Make sure that the topics are distinct, comprehensive, 946  
 specific, interpretable, and short. 947  
 948

Desired format: 949

- 1. <topic1> 950
- 2. <topic2> 951
- 3. <topic3> 952
- ... 953

- 1. <topic1> 956
- 2. <topic2> 957
- 3. <topic3> 958
- ... 959
- <N>. <topicN> 960  
 961

**C Models and Computations** 962

**C.1 LLM Model Versions** 963

Since off-the-shelf LLM are updated by the day, we report the exact model versions used in this work in 964  
 Table 5. 965

**C.2 Computational Cost** 966

During the experimentation stage of our work, we employed different LLM models. To run the models 967  
 we have used both the University’s GPU infrastructure (mainly used 3 GPUs with memory of 48GB each) 968  
 and AWS Cloud services (EC2, AWS Bedrock). We report the model versions in §C.1. The different 969  
 properties (e.g. number of parameters) of these models can be found online based on the version, if 970  
 published by developers. Overall we estimate the computational cost of about 2 weeks of GPU run time. 971

Developer	Model Family	Version
OpenAI	GPT	gpt-4-0125-preview, gpt-3.5-turbo-0125
Meta	LLAMA	Meta-Llama-3-8B-Instruct, Meta-Llama-3-70B-Instruct
Mistral	Mistral	Mixtral 8x7B

Table 5: LLM model versions used in this work, grouped by model family

972  
973  
974  
975  
976

## D Additional Results

### D.1 Judge Model Evaluation

To further support our claim that LLMs can be used as judge models for measurement annotation, Table 6, depicts additional correlation measures, and Table 7 shows the correlation between human annotations and the mean human score used as the test set.

Model	Quant.	Relevance			Overlap			Interpretability		
		Pear.	Spear.	Kend.	Pear.	Spear.	Kend.	Pear.	Spear.	Kend.
GPT 4	-	<b>0.70</b>	<b>0.66</b>	<b>0.52</b>	0.89	0.86	0.74	0.72	0.63	0.47
GPT 3.5	-	0.53	0.50	0.40	0.82	0.79	0.68	<b>0.76</b>	<b>0.73</b>	<b>0.59</b>
LLAMA 3 (8B)	None	0.44	0.45	0.37	0.88	0.85	0.76	0.67	0.54	0.43
LLAMA 3 (70B)	4-bit	0.39	0.48	0.40	0.31	0.24	0.22	0.16	0.29	0.22
LLAMA 3 (70B)	None	<u>0.64</u>	<u>0.62</u>	<u>0.49</u>	<b>0.90</b>	<b>0.87</b>	<b>0.77</b>	0.67	<u>0.66</u>	<u>0.51</u>
Mixtral (8x7B)	4-bit	0.44	0.43	0.34	0.88	0.83	0.72	0.73	0.65	0.50
Mixtral (8x7B)	None	0.53	0.50	0.39	0.79	0.73	0.65	<u>0.74</u>	0.65	<u>0.51</u>

Table 6: An extension of table 2. Showing Pearson (Freedman et al., 2007), Spearman (Spearman, 1961) and Kendall (Kendall, 1948) correlation between LLM and mean human annotations. The best overall model for each measurement is boldfaced and the best open-source alternative is underlined.

	Relevance			Overlap			Interpretability		
	Pear.	Spear.	Kend.	Pear.	Spear.	Kend.	Pear.	Spear.	Kend.
Annotator 1	<b>0.93</b>	0.67	<u>0.58</u>	<u>0.95</u>	<u>0.93</u>	<u>0.89</u>	<b>0.92</b>	<b>0.91</b>	<b>0.8</b>
Annotator 2	<u>0.85</u>	<b>0.95</b>	<b>0.89</b>	-	-	-	-	-	-
Annotator 3	0.90	<u>0.66</u>	<u>0.58</u>	<b>0.95</b>	<b>0.96</b>	<b>0.91</b>	<u>0.86</u>	<u>0.77</u>	<u>0.66</u>
Annotator 4	0.92	0.71	0.62	-	-	-	0.91	0.83	0.71

Table 7: Correlation of each annotator with the mean human annotation used as the test set. The annotators with max./min. correlation for each metric is boldfaced/underlined respectively.

Original Title	Corrupted Title
“Fear of being shot by Germans”	“Trepidation Under Teutonic Projectiles”
“Inhumane conditions in the concentration camps”	“Unkind States at Encampment Zones”
“Disbelief”	“Dissonant Credence”
“Encounter with Russian soldiers”	“Conflux with Rus Algid Militants”
“Russian liberation”	“Slavic Unshackling”
“Discovery of bodies and evidence of mass killings”	“Unearthed Enigmas: Corporeal Clusters & Mortality Indices”
“Food”	“Nourishment Alchemization Elements”
“Hospitals and medical treatment”	“Healing Havens and Remedial Maneuvers”
“Red Cross”	“Crimson Intersection”
“Bombings and attacks”	“Explosive Events and Assaults Unclear”

Table 8: Examples of title corruptions generated using GPT4.

## E Title Set Generation Systems

Examples of generated title sets for each generation system are shown in Table 11.

977  
978



## Coverage – Overlap Trade-Off

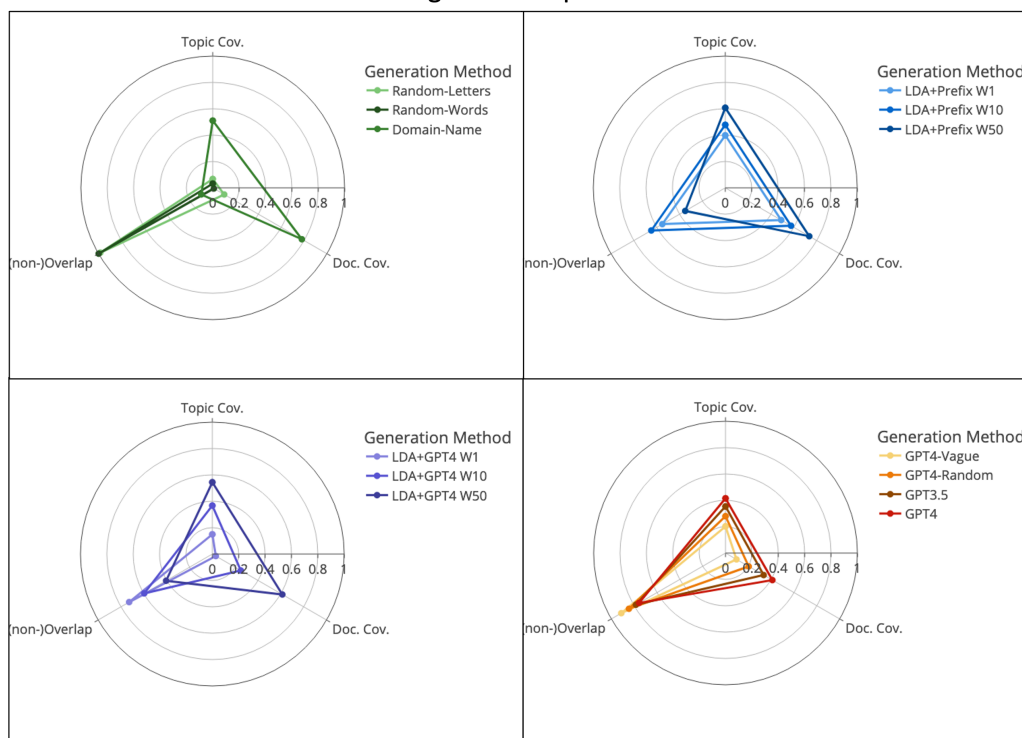


Figure 4: Coverage - Overlap trade-off for all systems, grouped by generation approach.

## F Validation Results

This section shows a thorough analysis of the results presented in §6, further expanding on the trade-offs arising between the different title sets due to the different generation approaches. Figures 4, 5, 6, 7 show the same trade-offs presented in Figure 2 but extended to include all tested systems.

**Coverage-Overlap Trade-Off** Throughout the study, intricate trade-offs emerged between individual aspects. The most prominent trade-off arises between the Coverage aspect (Title Cov. and Doc. Cov.) and the non-Overlap aspect. While it is easy to generate title sets that achieve high Coverage or non-Overlap scores, excelling in both is challenging.

To check whether our methodology reflects this trade-off, Figure 2(a) compares 4 generation systems, one from each group of methods. The first two methods are extreme cases of high non-overlap/low coverage and low non-overlap/high coverage, respectively.

Since Random-Words generates titles randomly, its title sets should not contain titles that cover the documents nor are overlapping. Domain-Name utilizes the domain names assigned by the annotators which were intended to describe the entire domain and therefore most of the documents should be covered by its sets. As a middle ground, we also examine LDA-Prefix W10 and GPT4. These two systems represent naturalistic systems and therefore are expected to reflect better balance. The figure demonstrates that our methodology successfully captures the coverage-overlap trade-off. Random-Word and Domain-Name tend toward high non-overlap/low coverage and low non-overlap/high coverage respectively, and LDA-Prefix W10 and GPT4 are more balanced between all 3 aspects where the first is more coverage oriented, indicating higher-level and less diverse titles while the latter is more non-Overlap oriented indicating a more specific and diverse set.

Examining the more elaborated Figure 4, we note that simpler methods (either a small number of words in the output or older versions) achieve lower coverage scores than more complex ones, where the coverage levels improve from system to system. However, this improvement is often achieved at the expense of the non-overlapping of the titles. This is most visible in the case of LDA-based methods, where the best coverage-achieving methods rely on 50 words in each topic cluster, however, they score

## Overall

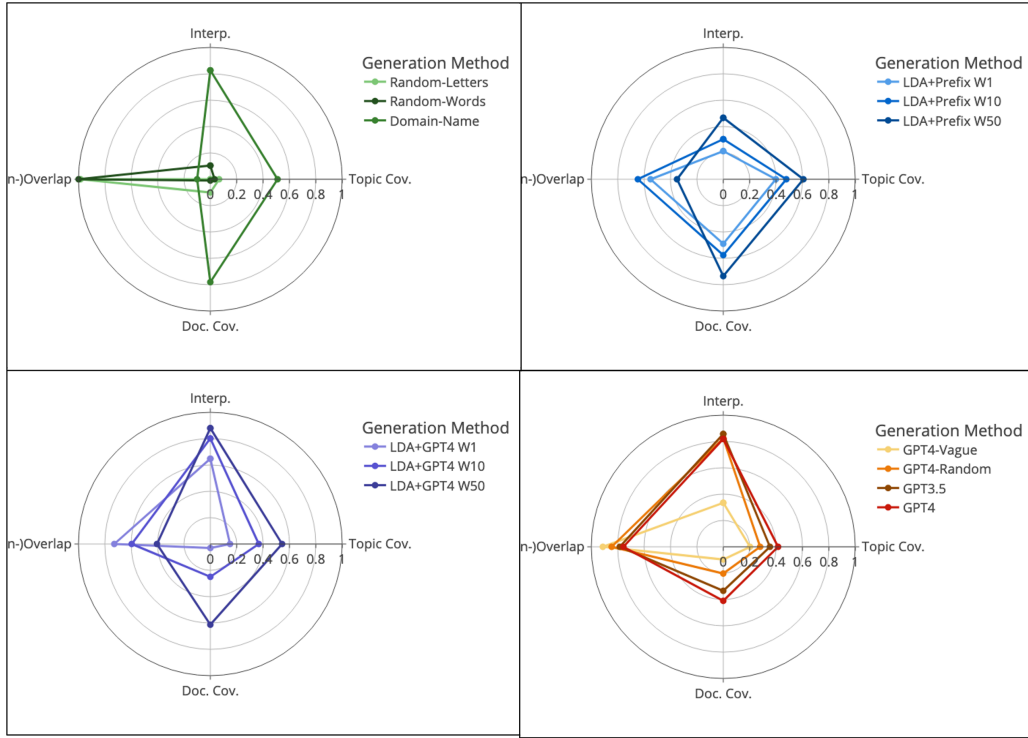


Figure 5: Overall comparison of all aspects other than Inner-Order, for all systems. Grouped by generation approach.

1005 much lower on the non-Overlap aspect than the base version of 1 word per topic cluster. This indicates  
 1006 that a larger number of words in the topic cluster helps in defining the topic they represent. Alongside, the  
 1007 topics become more general, causing overlapping. These results align with the increasing mean number  
 1008 of overlapping words between LDA clusters as the number of words in the cluster increases (see Table 9)

$k$	Mean Word Overlap
1	0.40
10	0.55
50	0.60

Table 9: Mean number of exact word overlap between pairs of LDA top  $k$  words clusters for varying number of words in a cluster. The table shows that the overlap between clusters increases as the number of words in the cluster increases.

1009 **Title Coverage and Document Coverage Tradeoff.** A more subtle but central trade-off exists between  
 1010 Title Coverage and Document Coverage metrics. Figure 2(b) depicts this trade-off. Here too, the  
 1011 methodology successfully gives a low score to Random-Words which generate titles that do not represent  
 1012 any real theme and therefore should not cover any document in the sample. Alongside, the methodology  
 1013 scores highly on the title sets generated by Domain-Name which renders high-level titles that should be  
 1014 relevant to most documents in the sample. Results further indicate that GPT4-Random achieves higher  
 1015 scores than Random-Word. Since GPT4-Random generates Holocaust-related titles, this demonstrates the  
 1016 methodology’s ability to capture fine-grained quality differences.

1017 Examining the more elaborated Figure 6, we notice that the methodology also captures the trade-off that  
 1018 arises between systems that generate higher-level and non-diverse titles (Domain-Name and LDA-based)  
 1019 to LLM-based systems which generate more specific and diverse titles. Indeed, the firsts achieve higher  
 1020 overall coverage scores at the expense of leaning towards Document Coverage over Title Coverage, while

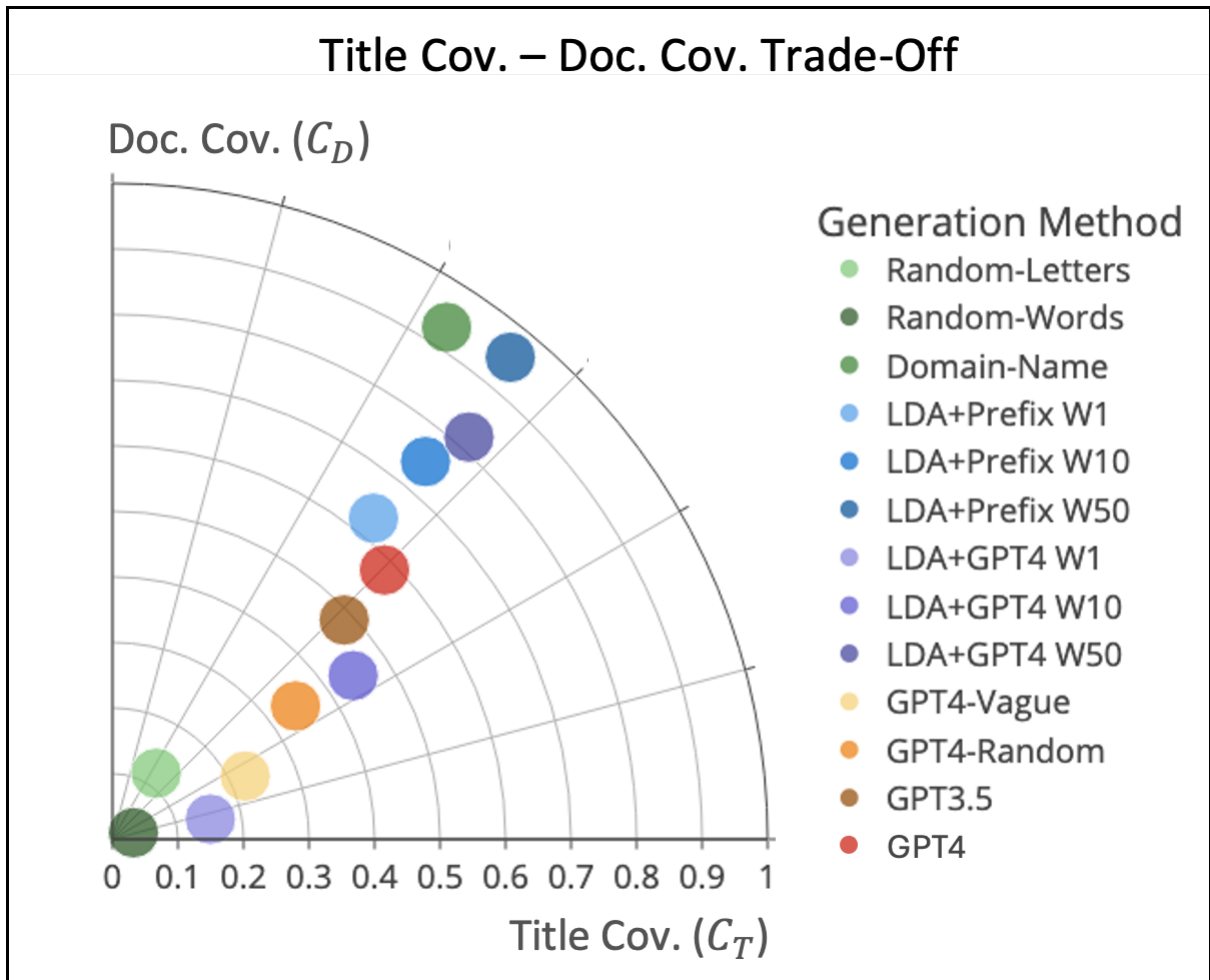


Figure 6: Title Coverage - Document Coverage trade-off for all systems.

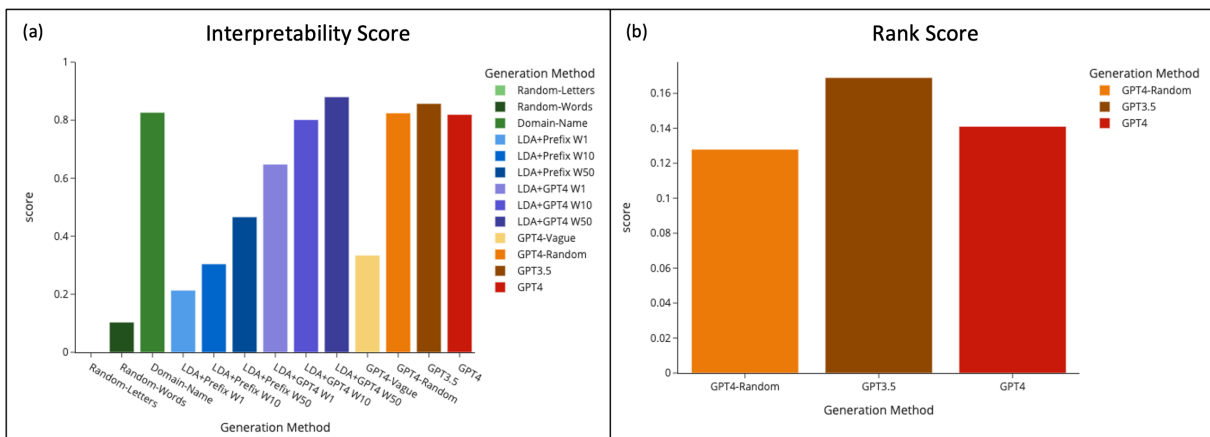


Figure 7: Interpretability and Inner-Order scores for all systems (that participate).

1021 the latter better balances between the two but achieves a lower overall coverage score.

1022 **Interpretability Trade-Off.** Figure 2(c), 7(a) shows Interpretability across systems. The bars in the  
1023 figure are color-coded so it will be easier to distinguish between the different system groups. Examining  
1024 the results we first notice that the methodology successfully captures the low interpretability built into  
1025 Random-Letters and Random-Words, while human-generated titles (Domain-Name) and systems that  
1026 employ LLMs (excluding GPT4-Vague) achieve the highest scores. In the case of GPT4-Vague, the system  
1027 was specifically designed to output uninterpretable titles, which aligns with its low score. Furthermore,  
1028 LLM-based methods achieve comparable scores to humans, aligning with recent claims that LLMs achieve  
1029 high fluency (Yang et al., 2023; Lai et al., 2023; Jiao et al., 2023). Additionally, we note that systems  
1030 based only on LDA (LDA-Prefix) are ranked in the low to mid-score ranges. This aligns with the main  
1031 drawback of LDA-based topics which are difficult to interpret. Finally, comparing the LDA-Based method  
1032 to LLM-based methods we can see that the methodology successfully captures an improvement in the  
1033 interpretability score of LDA-Based systems when increasing the number of words, while the score of  
1034 LLM-based systems remains steady. This phenomenon is attributed to the fact that increasing the number  
1035 of words in an LDA cluster adds substantial useful information, whereas changing the LLM version  
1036 doesn't necessarily enhance its ability to generate high-quality titles.

1037 **Inner-Order Performance.** Figure 2(d), 7(b) shows the inner-order scores achieved by LLM-based  
1038 systems. While LDA-based methods inherently neglect inner ordering, when designing the LLM-based  
1039 methods we did not specify any ordering instruction in the generation prompt (see Appendix B). In this  
1040 comparison, we choose to only include systems that were under our control, and for this reason, we choose  
1041 to only include LLM-based systems. The results show that our methodology successfully captures the  
1042 lack of ordering instruction by not significantly surpassing the random baseline. We note however that  
1043 this result may be easily improved by better prompt engineering.

1044 **Overall Comparison.** Figure 2(e) depicts an overall comparison of representing systems from each  
1045 generation group, considering all aspects other than the Inner-Order aspect. We notice that both the  
1046 LDA-based and LLM-based systems, which correspond to applicable systems, achieve high scores on  
1047 all aspects compared to the baseline methods. However, it is also hard to tell which model outperforms.  
1048 Examining the separate metrics, we notice the intricate trade-offs between the systems. While LLM-based  
1049 methods tend to distribute evenly across aspects, LDA-based methods tend towards higher-level titles,  
1050 which correspond to high coverage at the expense of non-Overlap and Interpretability. These conclusions  
1051 are further stressed in the full system comparison depicted in Figure 5.

## 1052 G Annotation Guidelines

1053 The following includes the annotation guidelines provided for each measurement annotations. Before  
1054 passing the guidelines to the annotators, a short frontal meeting was conducted where we introduced our  
1055 research and the specific goals of the annotation session. We introduced the data (Holocaust Testimonies)  
1056 and discussed its subtleties and sensitivities. Finally, the guidelines and examples were presented  
1057 and discussed. During the meeting, we have answered any questions raised by the annotators. Each  
1058 measurement received its own annotation guidelines and was conducted independently: first relevance,  
1059 then overlap, and finally interpretability.

### 1060 Relevance

1061 Following is a collection of passages extracted from Holocaust Testi-  
1062 monies. Please read thoroughly each one of the documents. When you  
1063 finish, you will be shown a passage from the collection along with a  
1064 set of titles, each title represents a theme. For each passage-title  
1065 pair, please indicate how relevant is the title to the given passage  
1066 (0 - not relevant at all, 100 - very relevant).

## Overlap

Attached are the files required to tag the Overlap task. The files include:

- A text file containing a collection of passages for annotation (the same passages you have already seen). It is worth opening the file in "Word" for ease of reading.
- An Excel file containing pairs of titles under the same domain in which you will have to fill in the overlap scores.

The file contains 4 columns: "domain": the label given to the domain by SF; "topic 1", "topic 2": Titles relevant to the domain and that are to be scored; "score": the appropriate score in your opinion from 0 to 100 according to the definition below; "reasoning": your explanation for the score in a short sentence.

### Task definition:

- Open the text file and read all the passages (you should already be familiar with these passages)
- Open the Excel file. For each pair of titles, give a score between 0 and 100 for the degree to which the themes defined by the two titles overlap, in the context of the passages (0 = no overlap at all, 50 = there is a partial overlap, 100 = there is a complete overlap / the titles have the same meaning).

## Interretability

Attached are Excel files containing titles and a text file containing experiences from Holocaust Testimonies. The experiences are the same experiences from previous tasks, but please go through them and read them again. The Excel file contains the titles for labeling.

Task definition: For each title, give a score of 0-100 for the degree to which the title is understandable (75-100 = the theme is understandable, 50-75 = the theme is partially understandable, 25-50 = the theme is poorly understandable, 0-25 = it is not possible to understand what is the intended theme). An understandable title is a title that the theme it induces can be easily understood from the title's text, in the context of the documents. If the theme is clear but not related to the documents you have seen, please give a score regardless of the documents and make a note in the "notes" column. In addition, you must give a one-sentence explanation of the score. The explanation should be noted in the "explanation" column.

### Highlights:

- Do you know which parts of the story the title refers to?
- Can you find an example in the text that links to the title?
- It should be noted that one title may include several topics that are not clearly related (in the context of the documents) such that it may not be clear which theme the title describes overall.
- Some titles describe features of the theme but do not give a clear and understandable name to the theme. Points should be deducted for this.

1117  
1118  
1119

- Pay attention to the wording, points must be deducted for titles that are not clearly worded.
- Points must be deducted in case there is unnecessary information.

Table 10: Examples of segments extracted from the testimonies and the corresponding ontology labels assigned by SF. Speakers are denoted as either “INT” for the interviewer or the first letter of the first and last name of the survivor. Note that multiple labels are possible for the same segment.

Labels	Segment
“Deportation to Concentration Camps”, “Jewish Prayers”	“before. INT: When they left– when– when they told you to get out of your home, where did they– SK: We were– my mother was baking cookies. INT: Yes? SK: We should have for the trip. And they come in, the Gendarmes, but from our same village. We know them. They said, listen, Günczler [NON-ENGLISH], you have to pack your package. You can bring only– I know the exact details, all. And you have to come up here, in front of the house, five in a row. And I’ll come back in 20 minutes, or whatever, and you have to be ready. So my mother put us the clothes on and the food for the kids, whatever we could. And we– we were waiting there. And they took us for the night to this big [NON-ENGLISH], has a big shul. And there we sit in there. But this is there. I shouldn’t repeat it. INT: No, no, it’s OK. SK: I will talk about it. Or if you want to start, and then I’ll tell you. INT: No, no, no. Just tell me. SK: Now? OK. So when– so that night, we sit in the shul, everybody and their luggage, and the men saying”
“Deportation to Concentration Camps”, “Forced Marches”	“it was all organized by the transport [? Leitung, ?] you know? Everything was seemingly made by our own people. INT: Did you see any Germans? RS: No, no. I didn’t. INT: What did you see? How long did the journey take, the walk? RS: Well, it was about four kilometers. INT: Did you arrive at day? What time of day did you arrive? RS: It was night. It was night. INT: Were you marching in the dark? RS: Yes. INT: Were any orders given to you? RS: No, no. INT: Was anybody hit or any punishments given? RS: No. I couldn’t see anything. There were Czech gendarmes around, and some SS men. But they didn’t touch anybody. INT: What nationality”
“Living conditions”, “Protected houses (Budapest)”	“didn’t get along very well. We never did get along very well with her. And all her things were there. And we used all her thing. And we didn’t have our own sheets, and our own pillow cases, and our own beddings. But we– all of us moved, like three– three or four of us moved into a small room, where she stayed with my– In the meantime, my sister actually left, too. She was– she was hiding somewhere. We didn’t know where. At one point she disappeared, and my father and I took off the stars, and were looking for her all day long. That was in summer– must have been July or August. We’re looking for her all– all day long, and then it turned out that she went with– to yoga teacher. At that time when nobody in Budapest even”

Table 11: Examples of title sets generated by each system for the domain “Antisemitism”.

<b>Generation Method</b>	<b>Title Rank</b>	<b>Generated Title</b>
Random-Letters	1	DTrHXGOEuctmGDuQd
	2	tHTbUhnToumKgtEedNlkRo
	3	zCPYogMzYgObhMZYiDNexdyZ
	4	lluAvbK
	5	KkhtVdgzUcAD
	6	qQDlywcXWxvzEhtRjid
	7	JsdcvRfzjTIAYq
	8	ZTPazuWwfFTwnZKoINUU
	9	PloDhuTCp
	10	EZXckfQkRmxGhcS
Random-Words	1	brachtmema diatomin
	2	garfish obscuring asterisks
	3	select serjeantry vavasories
	4	fathers raylet integrate
	5	restrengthen hoplonemertine
	6	perfectible spondylexarthrosis obtrusiveness
	7	conventionalism
	8	hotter incoalescence
	9	demulce
	10	underpainting extending circumrotate
Domain-Name	1	antisemitism
	2	antisemitism
	3	antisemitism
	4	antisemitism
	5	antisemitism
	6	antisemitism
	7	antisemitism
	8	antisemitism
	9	antisemitism
	10	antisemitism
LDA+Prefix W1	1	The theme defined by the following set of words: “int”.
	2	The theme defined by the following set of words: “int”.
	3	The theme defined by the following set of words: “know”.
	4	The theme defined by the following set of words: “jewish”.
	5	The theme defined by the following set of words: “int”.
	6	The theme defined by the following set of words: “int”.
	7	The theme defined by the following set of words: “int”.
	8	The theme defined by the following set of words: “jewish”.
	9	The theme defined by the following set of words: “int”.
	10	The theme defined by the following set of words: “int”.

Generation Method	Title Rank	Generated Title
LDA+Prefix W10	1	The theme defined by the following set of words: “int”, “school”, “jewish”, “would”, “us”, “know”, “one”, “remember”, “went”, “time”.
	2	The theme defined by the following set of words: “int”, “know”, “school”, “jewish”, “time”, “jews”, “jew”, “one”, “went”, “seconds”.
	3	The theme defined by the following set of words: “know”, “int”, “one”, “school”, “jewish”, “remember”, “would”, “time”, “pauses”, “seconds”.
	4	The theme defined by the following set of words: “jewish”, “know”, “int”, “used”, “jews”, “like”, “people”, “school”, “would”, “go”.
	5	The theme defined by the following set of words: “int”, “jewish”, “know”, “like”, “jews”, “people”, “went”, “said”, “yes”, “remember”.
	6	The theme defined by the following set of words: “int”, “know”, “would”, “school”, “remember”, “jewish”, “one”, “like”, “seconds”, “pauses”.
	7	The theme defined by the following set of words: “int”, “going”, “would”, “one”, “bg”, “english”, “non”, “put”, “went”, “jew”.
	8	The theme defined by the following set of words: “jewish”, “int”, “know”, “one”, “school”, “seconds”, “pauses”, “jews”, “well”, “would”.
	9	The theme defined by the following set of words: “int”, “know”, “seconds”, “pauses”, “jews”, “people”, “jewish”, “came”, “would”, “see”.
	10	The theme defined by the following set of words: “int”, “know”, “school”, “go”, “jewish”, “went”, “people”, “us”, “came”, “one”.
LDA+Prefix W50	1	The theme defined by the following set of words: “int”, “school”, “jewish”, “would”, “us”, “know”, “one”, “remember”, “went”, “time”, “yes”, “go”, “came”, “well”, “jews”, “children”, “said”, “like”, “even”, “get”, “first”, “home”, “pauses”, “think”, “seconds”, “people”, “say”, “jew”, “could”, “got”, “non”, “going”, “much”, “back”, “parents”, “never”, “day”, “come”, “polish”, “started”, “called”, “town”, “high”, “always”, “used”, “lot”, “knew”, “father”, “boys”, “german”.
	2	The theme defined by the following set of words: “int”, “know”, “school”, “jewish”, “time”, “jews”, “jew”, “one”, “went”, “seconds”, “pauses”, “yeah”, “go”, “children”, “came”, “remember”, “first”, “said”, “yes”, “would”, “going”, “us”, “well”, “father”, “say”, “people”, “like”, “antisemitism”, “ml”, “non”, “hitler”, “war”, “told”, “parents”, “english”, “years”, “little”, “mother”, “polish”, “anti”, “think”, “german”, “mean”, “friends”, “used”, “mb”, “house”, “thing”, “old”, “started”.
	3	The theme defined by the following set of words: “know”, “int”, “one”, “school”, “jewish”, “remember”, “would”, “time”, “pauses”, “seconds”, “jews”, “go”, “went”, “little”, “like”, “jew”, “really”, “hl”, “laughs”, “father”, “first”, “said”, “came”, “got”, “non”, “child”, “well”, “mean”, “think”, “say”, “took”, “want”, “could”, “kind”, “course”, “teacher”, “quite”, “things”, “started”, “us”, “even”, “thing”, “english”, “yes”, “knew”, “come”, “grade”, “boy”, “house”, “high”.
	4	The theme defined by the following set of words: “jewish”, “know”, “int”, “used”, “jews”, “like”, “people”, “school”, “would”, “go”, “non”, “went”, “us”, “jew”, “one”, “remember”, “polish”, “time”, “english”, “war”, “said”, “yeah”, “got”, “came”, “lot”, “seconds”, “pauses”, “antisemitism”, “see”, “poland”, “say”, “even”, “children”, “come”, “always”, “could”, “sb”, “back”, “mother”, “well”, “good”, “going”, “little”, “many”, “get”, “called”, “think”, “way”, “took”, “home”.



Generation Method	Title Rank	Generated Title	
LDA+Prefix W50	5	The theme defined by the following set of words: “int”, “jewish”, “know”, “like”, “jews”, “people”, “went”, “said”, “yes”, “remember”, “mother”, “came”, “us”, “would”, “go”, “jk”, “father”, “well”, “school”, “could”, “fs”, “polish”, “time”, “one”, “non”, “little”, “seconds”, “pauses”, “english”, “think”, “name”, “get”, “yeah”, “used”, “see”, “lot”, “yiddish”, “two”, “war”, “lived”, “never”, “something”, “really”, “home”, “years”, “oh”, “tell”, “say”, “told”, “german”.	
	6	The theme defined by the following set of words: “int”, “know”, “would”, “school”, “remember”, “jewish”, “one”, “like”, “seconds”, “pauses”, “said”, “go”, “well”, “people”, “came”, “went”, “time”, “yes”, “jews”, “used”, “think”, “us”, “going”, “jew”, “mother”, “always”, “father”, “things”, “children”, “say”, “got”, “come”, “oh”, “could”, “little”, “much”, “day”, “first”, “really”, “back”, “knew”, “home”, “name”, “course”, “see”, “also”, “get”, “two”, “started”, “never”.	
	7	The theme defined by the following set of words: “int”, “going”, “would”, “one”, “bg”, “english”, “non”, “put”, “went”, “jew”, “tape”, “hiding”, “well”, “little”, “police”, “day”, “pauses”, “take”, “hit”, “seconds”, “course”, “go”, “two”, “thrown”, “discuss”, “ways”, “rocks”, “among”, “got”, “ok”, “number”, “next”, “time”, “way”, “think”, “poland”, “know”, “polish”, “boy”, “bad”, “couple”, “guns”, “kids”, “father”, “killed”, “laughs”, “three”, “say”, “us”, “jk”.	
	8	The theme defined by the following set of words: “jewish”, “int”, “know”, “one”, “school”, “seconds”, “pauses”, “jews”, “well”, “would”, “like”, “said”, “people”, “antisemitism”, “us”, “non”, “time”, “mother”, “think”, “went”, “go”, “used”, “kids”, “lived”, “yes”, “things”, “little”, “friends”, “say”, “er”, “name”, “even”, “years”, “german”, “children”, “family”, “father”, “polish”, “always”, “english”, “came”, “hl”, “way”, “home”, “called”, “poland”, “lot”, “felt”, “quite”, “got”.	
	9	The theme defined by the following set of words: “int”, “know”, “seconds”, “pauses”, “jews”, “people”, “jewish”, “came”, “would”, “see”, “one”, “well”, “time”, “went”, “said”, “polish”, “like”, “go”, “us”, “say”, “war”, “remember”, “could”, “school”, “non”, “yes”, “many”, “back”, “years”, “english”, “right”, “always”, “going”, “something”, “good”, “poland”, “first”, “think”, “get”, “started”, “name”, “father”, “yeah”, “antisemitism”, “told”, “called”, “things”, “wanted”, “took”, “little”.	
	10	The theme defined by the following set of words: “int”, “know”, “school”, “go”, “jewish”, “went”, “people”, “us”, “came”, “one”, “jews”, “remember”, “would”, “like”, “said”, “time”, “father”, “going”, “well”, “used”, “back”, “yes”, “could”, “really”, “pauses”, “seconds”, “little”, “home”, “mother”, “non”, “never”, “children”, “say”, “see”, “friends”, “always”, “antisemitism”, “get”, “war”, “got”, “house”, “poland”, “oh”, “mean”, “even”, “polish”, “family”, “many”, “take”, “years”.	
	LDA+GPT4 W1	1	Title: Programming and Development
		2	Programming Languages
		3	Understanding or Knowledge Acquisition
		4	Title: Judaism
5		Numerical Data & Analysis	
6		Programming and Software Development	

<b>Generation Method</b>	<b>Title Rank</b>	<b>Generated Title</b>
LDA+GPT4 W1	7	Title: Mathematics or Numerical Analysis
	8	Title: Jewish Culture or Religion
	9	Title: Programming and Software Development
	10	“Given the extremely limited information (only the word “int”), it’s challenging to assign a precise and meaningful title. The word “int” is commonly used in programming languages to represent “integer,” a data type used to store whole numbers. Therefore, without more context or additional words from the cluster, a title could broadly refer to programming or computer science concepts related to numerical data types. Thus, a potential title could be: Title: “Numeric Data Types in Programming”
LDA+GPT4 W10	1	Jewish Education and Experiences
	2	Jewish Education and Community
	3	Title: Educational Reflections in Jewish Contexts
	4	Jewish Cultural and Educational Experiences
	5	Jewish Community and Historical Events
	6	“School Experiences and Memories in Jewish Communities”
	7	Title: Language and Identity Discussions
	8	Jewish Education and Cultural Discussions
	9	Discussion on Jewish Historical Events and Perspectives
	10	Jewish Community and Education Experiences
LDA+GPT4 W50	1	Title: Memories of Jewish Life and Education
	2	Jewish Life and Antisemitism in Historical Context
	3	Jewish Education and Childhood Memories
	4	Jewish Experience and Antisemitism in Poland
	5	Jewish Life and Heritage Through Personal Narratives
	6	Title: “Personal Memories and Cultural Reflections”
	7	Title: Conversational Narratives and Conflicts
	8	Jewish Identity and Experiences Amid Antisemitism
	9	Jewish Experience and Antisemitism in Poland
	10	Jewish Family and Antisemitism in Poland
GPT3.5	1	Anti-Semitism in Schools and Education
	2	Segregation and Discrimination in Jewish Communities
	3	Jewish Identity and Pride
	4	Fear and Violence against Jews
	5	Holocaust Survivors’ Experiences
	6	Prejudice and Stereotypes
	7	Impact of Nazi Ideology on Education
	8	Discrimination in Schools
	9	Jewish Family and Community Support
	10	Anti-Semitic Propaganda and Indoctrination
GPT4	1	Antisemitism
	2	Jewish Education
	3	Jewish Community Life
	4	Personal Experiences of Discrimination

<b>Generation Method</b>	<b>Title Rank</b>	<b>Generated Title</b>
GPT4	5	Impact of Nazi Policies
	6	Jewish-Gentile Relations
	7	School Experiences
	8	Family Dynamics
	9	Resistance and Survival Strategies
	10	Post-War Experiences
GPT4-Vague	1	Anisdeitsm
	2	Hebraic Pedagogy Enigmas
	3	Judaic Communal Existence
	4	Experiential Encodings of Differential Treatment
	5	Policy Influence of N-Axis Entities
	6	JewGent Nexus Dynamics
	7	Educational Episodes
	8	Kinetic Household Constructs
	9	Defiance and Endurance Tactics
	10	Ex-Combat Aftermaths
GPT4-Random	1	Survival Strategies
	2	Encounters with Local Populations
	3	Smuggling and Black Market
	4	Violence and Persecution
	5	Daily Routine
	6	Immigration and Resettlement
	7	Ghettoization
	8	Post-War Migration
	9	Curfews
	10	Forced Labor