

# PREME: Preference-based Meeting Exploration through an Interactive Questionnaire

Anonymous ACL submission

## Abstract

The recent increase in the volume of online meetings necessitates automated tools for managing and organizing the material, especially when an attendee has missed the discussion and needs assistance in quickly exploring it. In this work, we propose a novel end-to-end framework for generating interactive questionnaires for preference-based meeting exploration. As a result, users are supplied with a list of suggested questions reflecting their preferences. Since the task is new, we introduce an automatic evaluation strategy. Namely, it measures how much the generated questions via questionnaire are answerable to ensure factual correctness and covers the source meeting for the depth of possible exploration.

## 1 Introduction

In recent years, video conferencing technology has gained substantial improvements, and thus, online meetings have become easily accessible and more prominent. Primarily due to the pandemic and work from home, the need for video calling has grown significantly. For example, the number of meeting minutes held in the Zoom application has increased by 3300% in 2021 compared to the same quarter of the previous year. Therefore, the high volume of online meetings necessitates automated tools for managing and organizing essential information for the attendees. Especially in cases when an attendee has missed an online meeting, it is critical to quickly access required information since the transcript of a 1-hour meeting averagely consists of 8000 words, and reading through is time-consuming.

Providing meeting summaries is a promising direction (Wang and Cardie, 2013; Jacquenet et al., 2019; Zhao et al., 2019; Singhal et al., 2020). However, recent works (Murray et al., 2010; Mehdad et al., 2013; Li et al., 2019; Zhu et al., 2020a) have demonstrated that approaches designed for document summarization could not effectively apply to

The following subjects were discussed in the meeting. Which subject are you more interested in?	
<input type="checkbox"/> Remote Control Cases	<input type="checkbox"/> Remote Control Design
<input type="checkbox"/> Remote Control Functions	<input type="checkbox"/> Remote Control Buttons
<input checked="" type="checkbox"/> <b>New Remote Control</b>	<input type="checkbox"/> Remote control Price
What do you want to know more about the New Remote Control?	
<input checked="" type="checkbox"/> <b>Fronts</b>	<input type="checkbox"/> Features advantages
<input type="checkbox"/> Disadvantages	<input type="checkbox"/> Think
Additional questions you might be interested in:	
- What is the new feature of the front of the remote control?	
- What are different colors of the front for the remote control?	
- What are the latest trends for a front under remote control?	
- What is the difference between front and back of the remote control?	

Figure 1: An example of exploring one of the meetings from the collection (Carletta et al., 2005) based on user preferences through an interactive questionnaire. Users may exploit the questionnaire multiple times to explore various parts of the meeting.

meetings transcripts due to the following potential reasons: (1) Standard documents are more structured compared to meetings; (2) Spoken language used in meetings is less regular than documents; and (3) The speaker role is essential. Moreover, there is little meeting data publicly available that can be used for experimentation compared to regular documents such as news or articles. In contrast with document summarization, when summarizing a meeting, different users tend different preferences on what content should be included in the summary. Recently, Zhong et al. (2021) attempted to tackle this problem by proposing a query-based multi-domain meeting summary, where a user provides a query in question form, e.g., ‘What was the discussion about the jog dial’s function when talking about changes in the current design?’ to locate the part of the transcript that related to the query and then summarize. However, when attendees have missed the meeting, they cannot formulate such questions due to no prior knowledge about the meeting. To overcome this, we aim to address the following **research challenge**: *How can attendees effectively explore a meeting content without having prior knowledge about it?*

This work is motivated by the fact that asking

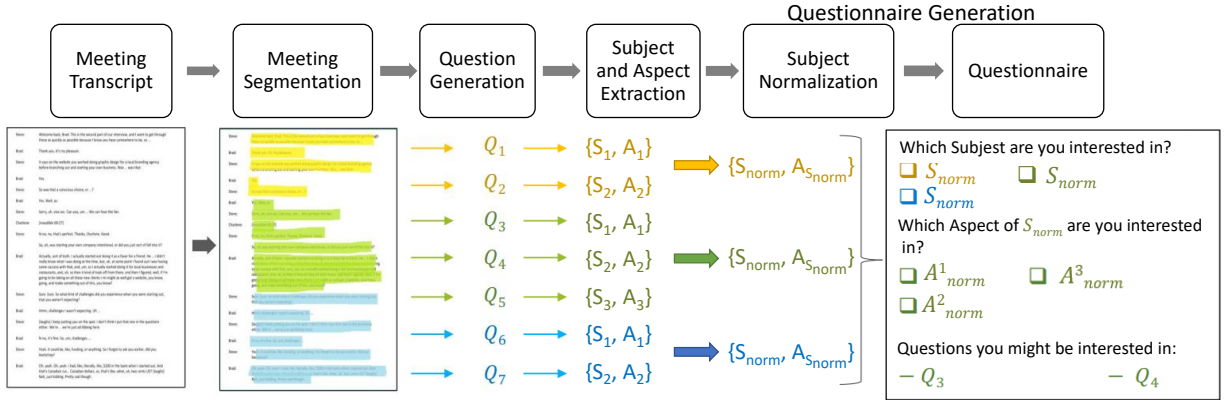


Figure 2: Overview of our framework, Preference-based Meeting Exploration through an Interactive Questionnaire (PREME), where  $Q$  is a comprehensive set of questions, and  $S_i$  and  $A_j$  are extracted pairs of subjects and aspects.

067 questions is a more efficient way for humans to  
 068 acquire information than notes in plain text (Law-  
 069 son et al., 2007, 2006). Hence, we address the  
 070 problem of preference-based meeting exploration  
 071 by automatically generating a structured interac-  
 072 tive questionnaire for a transcript that covers most  
 073 of the discussed topics and quickly walks users  
 074 through the discussed content. An example of the  
 075 desired questionnaire is shown in Fig. 1. First, the  
 076 user has the ability to express their preferences re-  
 077 garding *subjects* that have been discussed (Solbiati  
 078 et al., 2021; Huang et al., 2018; Zhang and Zhou,  
 079 2019; Sehikh et al., 2017). Next, the questionnaire  
 080 interactively suggests narrowing down their explora-  
 081 tion if possible by displaying a list of possible  
 082 related *aspects*. As a result, a list of questions re-  
 083 flecting user preferences is generated. Next, the  
 084 user can pick a question that demonstrates their  
 085 seeking needs the most and is redirected to the  
 086 meeting part containing an answer. Furthermore,  
 087 interactively asking for preferences is beneficial  
 088 for the user since the user oversees what has been  
 089 covered during the meeting they have missed.

Hence, the goal of proposed questionnaires for  
 exploration is two-fold: (1) to compactly represent  
 the discussed content; (2) to guide users to form  
 questions that express their preference regarding  
 the transcript. We require the generated question-  
 naire to satisfy the following properties:

- 090 **P1 Coverage:** coverage is the amount of the in-  
 091 formation from the source text that a question-  
 092 naire points to. The generated questionnaire  
 093 must cover the meeting as much as possible;  
 094 **P2 Answerable:** a given meeting transcript should  
 095 contain the answers to the questions generated  
 096 as a result of the questionnaire.

To address the defined challenge, we propose  
 a framework, PREME, which consists of sev-

096 eral concrete steps highlighted in Fig. 2. We  
 097 start by enchaining the method to extract meet-  
 098 ing segments (Solbiati et al., 2021). Due to the  
 099 conversational nature of the meeting, topic de-  
 100 tection from the segments is challenging (Huang  
 101 et al., 2018; Zhang and Zhou, 2019; Sehikh et al.,  
 102 2017). Thus, we indirectly extract the topics as  
 103 follows. First, we generate questions from each  
 104 segments (Brown et al., 2020). Further, we employ  
 105 a trained Conditional Random Field (CRF) model  
 106 to tag topics and aspects (Fig. 1) from generated  
 107 questions originated from each segments (Wallach,  
 108 2004). Once we got each segment’s topic list, we  
 109 proposed a strategy to normalize them to reduce  
 110 the number of options in the questionnaire.

To summarize, the main contributions are:

- 120 **C1** We propose PREME, a novel framework to en-  
 121 able meetings exploration based on user’s pref-  
 122 erences through an interactive questionnaire;  
 123 **C2** We propose a new method for subject normal-  
 124 ization which returns represent the most infor-  
 125 mative and general phrase from a set of phrases  
 126 and keywords;  
 127 **C3** We introduce a new automatic evaluation strat-  
 128 egy for measuring the effectiveness of the pro-  
 129 posed questionnaire to assess the required prop-  
 130 erties **P1** and **P2**; and  
 131 **C4** We open-source a dataset that includes 1000  
 132 questions comprehensively annotated with sub-  
 133 ject to their subjects and aspects.

## 2 Related Work 135

### 2.1 Automatic Textual Summarization 136

Automatic text summarization task has attracted  
 137 lots of attention across Natural Language Process-  
 138 ing (NLP) community recently. Many systems  
 139 are proposed to summarize documents in differ-  
 140 ent domains, including news (Rush et al., 2015;  
 141

Nallapati et al., 2017; See et al., 2017; Celikyilmaz et al., 2018; Liu and Lapata, 2019; Zhang et al., 2020), academic papers (Manakul and Gales, 2021; Huang et al., 2021) and books (Kryściński et al., 2021). Meeting summarization has also emerged as a widespread need recently. Due to the unique discourse structure of dialogues, conventional document summarization systems are facing challenges when summarizing meetings (Li et al., 2019; Zhu et al., 2020b). Thus, new models are proposed for tackling this task. Wang and Cardie (2013) employ decisions, action items in dialogues to progressively generate the summary. Oya et al. (2014) propose a template-based meeting summarization system by learning the relationship between summaries and their source meeting transcripts. Shang et al. (2018) design an unsupervised meeting summarization model with multi-sentence compression techniques. Li et al. (2019) introduce multi-modal information into meeting summarization with a hierarchical attention mechanism. Zhu et al. (2020b) propose a hierarchical meeting summarizer that can process both word-level and turn-level information of dialogs. Furthermore, it comes into sight of the community that, due to the lengthy content and distributed information, a general summary of the meetings does not necessarily satisfy what users are seeking. Thus, Query-based summarization methods become more prevailing in which the summaries are specifically and concisely generated according to user queries (Litvak and Vanetik, 2017; Nema et al., 2017; Baumel et al., 2018; Ishigaki et al., 2020; Kulkarni et al., 2020, 2021; Pasunuru et al., 2021). Recently, Zhong et al. (2021) propose a new framework of query-based summarization for meetings, in which they annotate QMSUM, a query-based multi-domain meeting dataset. Each QMSUM meetings come along with a set of queries with different levels of abstractness, i.e., general queries and specific queries. Human annotators write these queries and the summaries aligned with these queries after reading the meeting transcripts.

While query-based summarization can be a proper path to provide users with meeting information at different specificity levels, we argue that issuing such specific queries still requires a certain degree of background knowledge. In real-life scenarios, users might not be equipped with that knowledge and issue informative queries, especially when they did not attend the meeting. As a result, they can not benefit from query-based

summarization techniques to explore the meetings. Hence, we address the drawbacks of query-based summarizers by providing users with an interactive questionnaire in this work. It provides them with potential queries and allows them to explore the meetings more flexibly.

## 2.2 Evaluation of Summaries Factuality

The summaries often has called out for hallucination issues (Maynez et al., 2020). Therefore, Wang et al. (2020) propose a framework to evaluate factual consistency of summaries with the source text. Their intuition is that the summary and the source should similarly and consistently answer the factual questions about the context. Similarly, Deutsch et al. (2020) propose a Question Answering (QA)-based evaluation approach on summaries' content quality. They measure how much information is contained in a candidate summary by calculating the proportion of questions it can answer. These approaches inspired our way of thinking about automated end-to-end evaluations of generated questionnaires.

## 2.3 Question Generation and Filtering

Initial works in Question Generation task leveraged crowd-sourcing or rule-based methods to generate pre-defined question templates (Mostow and Chen, 2009; Rus et al., 2010; Lindberg et al., 2013; Fabbri et al., 2020; Mazidi and Nielsen, 2014; Labutov et al., 2015). Heilman and Smith (2010) tackled this problem in a different manner by over-generating candidate questions and then using a learning to rank framework to rank them. Ranking the questions helped filter the low-quality questions as they would rank lower.

SQUASH (Krishna and Iyyer, 2019) is one of the recent works in which authors used question generation methods to convert a document into a hierarchy of question-answer pairs with the focus on questions' granularity level. They employed a neural encoder-decoder model trained on three reading comprehension data sets, i.e., SQuAD (Rajpurkar et al., 2016), QuAC (Choi et al., 2018), and CoQ (Reddy et al., 2019) to generate the questions, and further, they filtered out the unanswerable questions using some heuristics and question answering models. While question generation using question answering data sets seems a general approach, this method does not work well on meeting-related questions generated due to many reasons, including: (1) Different structure of meetings com-



pared to documents; (2) There is not many question-answering datasets available from meetings; (3) Sometimes, the answer to questions generated from meetings could be very long, making it hard to fit the context in neural models. In our work, we introduce an automatic method that can generate questions regarding the meeting to overcome the high price of collecting with annotators as in (Zhong et al., 2021).

## 2.4 Questionnaire Organization

Obtaining users preferences has always shown to be a challenging task (Jiang et al., 2008; Rokach and Kisilevich, 2012; Anava et al., 2015; Christakopoulou et al., 2016; Sepliarskaia et al., 2018). The task becomes more challenging when we aim to minimize the number of interactions with users to get to know their preferences. For example, in (Sepliarskaia et al., 2018), the authors reformulate this task as an optimization problem. They propose a static questionnaire by choosing a minimal and diverse set of questions to solve the cold start problem in recommender systems. Similarly, in Liu et al. (2019) proposed a dynamic questionnaire generation method for search of clinical trials. They perform real-time dynamic question generation to select criteria at a time by maximizing its relevance score that reflects its potential to rule out ineligible trials from primary search results. Quiz-style question generation has also been explored recently by Lelkes et al. (2021). The authors have formulated the problem as two sequence to sequence tasks, including the question-answer generation step and incorrect answer generation step. We argue that while the former step seems relevant to our work, it could not be adapted to meeting transcripts since their proposed dataset has been trained on factual question answering data sets and cannot be used for meeting purposes. All in all, we can conclude that creating questionnaires are still under exploration in different domain. Hence, our effort in organizing a questionnaire, especially for meetings, is timely and useful for future research in NLP area.

## 3 Proposed Framework: PREME

This section explains the proposed novel methodology to explore meetings based on users' preferences through an interactive questionnaire, called PREME. An overview of our methodology is shown in Fig. 2 in which we first apply a topic segmentation method (Solbiati et al., 2021) on

meeting transcript to retrieve segments with different topics from the meeting (Section 3.1). Then, we generate a set of all possible questions from each segment (Section 3.2). Further, we extract the most informative part of the questions, i.e., the subject and aspect of each question (Section 3.3). In the last step, we map the normalized subjects and aspects with generated questions and form the questionnaire (Section 3.4).

### 3.1 Meeting Segmentation

A meeting transcript can be extremely long and contain discussions of various topics. Therefore, our goal is to divide the meeting text into a sequence of topically coherent chunks. Thus, we adopted an unsupervised topic segmentation method based on the contextualized presentation of meeting (Solbiati et al., 2021). In this topic segmentation method, the authors compute the BERT embeddings for every utterance of the meeting transcript. Further, they curated blocks of utterances and performed a block-wise max-pooling operation to generate contextualized embedding for each block. Then, the semantic similarity between two adjacent blocks is captured, and a change in the topic is detected if two adjacent blocks show similarity below a certain threshold. This approach has several advantages, including: (1) It is unsupervised; (2) Since we are putting barriers in between the meeting text, we are just converting the meeting into smaller pieces, and we are not losing any part of the meeting.

### 3.2 Question Generation

For question generation from a segment, we leveraged the powerful GPT-3 model (Brown et al., 2020).<sup>1</sup> An impressive capability of the GPT-3 is to generate very realistic results from few training samples or even no training sample (few-shot and zero-shot learning). The variety of the generated content can be controlled using a temperature hyperparameter. For question extraction in each segment, the API is called with different temperature values between [0-1] with a 0.05 margin, where the value closer to 1 means more diversified questions. We then repeat the process ten times for each specific temperature. A list of questions is extracted based on random initialization in each API call, meaning different results are achieved even with the same hyperparameters. We extracted

<sup>1</sup>GPT-3 is a large autoregressive Transformer-based language model developed by OpenAI, with 175 billion parameters. The model is accessible through API calls<sup>2</sup>.

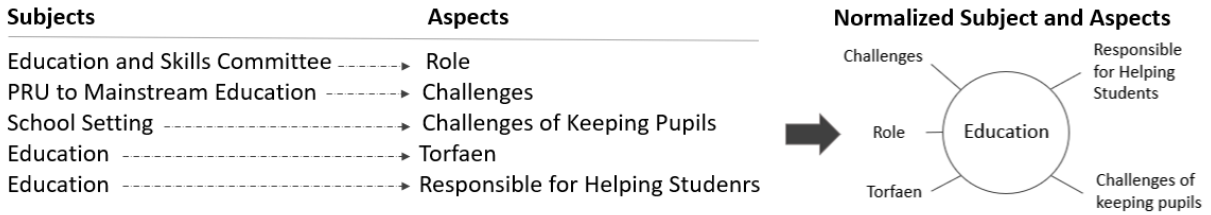


Figure 3: An example of how extracted subjects and aspects from a given segment are normalized.

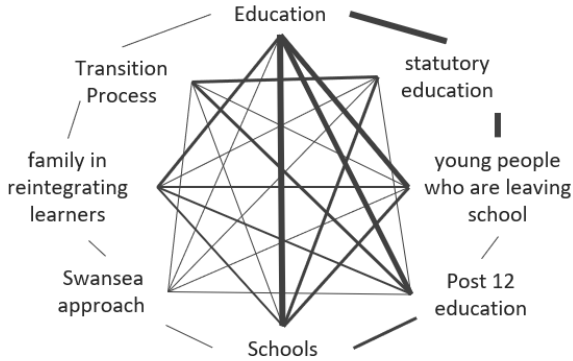


Figure 4: An example of subject-network built for one extracted segments from (Janin et al., 2003). Here, the edge weights is related to the semantic similarity between each nodes and edges with higher weights are shown with higher width. In this network, the node “Education” gained the highest PageRank value.

339 five questions on average per segment in each call.  
 340 Finally, a union across all runs is used to form our  
 341 question pool.

### 3.3 Subject and Aspect Extraction

343 Every of the generated questions has a *subject* that  
 344 it refers to, i.e., the principal matter that attendees  
 345 have discussed. In addition, some questions might  
 346 cover more details about a given subject, and they  
 347 might particularly ask about a certain related *aspect*.  
 348 We aim to extract the primary subjects from any  
 349 question and the detailed aspect if it is mentioned.  
 350 Table 1 shows a few examples of annotated questions  
 351 with regard of their *subjects* and *aspects*. For  
 352 instance, in the question “*What is the arrow symbol  
 353 on the remote control for?*”, “remote control”  
 354 is mentioned as a subject. There are some specific  
 355 aspects of the subject, i.e., the “arrow symbol”.

356 In order to extract the subjects and aspects from  
 357 the questions, we use CRF (Wallach, 2004). We  
 358 examined SOTA keyword extraction and contex-  
 359 tualized neural embedding-based topic extraction  
 360 models; however, the CRF model seems to work  
 361 the best among them. To train the CRF model, we  
 362 were required to have annotated questions with sub-  
 363 jects and aspects labels. We designed an annotation  
 364 study using a crowd-sourcing platform, where we

Table 1: Examples of annotated questions with their subjects and aspects for a product meeting from (Carletta et al., 2005). **Subjects** are highlighted in red and **Aspects** are highlighted in green.

Q1	What is the <b>arrow symbol</b> on the <b>remote control</b> for?
Q2	What are the <b>main frustrations people have</b> with the <b>remote control</b> ?
Q3	How will the <b>logo and color scheme</b> be incorporated into the <b>product</b> ?
Q4	What are <b>pros and cons</b> of having a <b>remote with a large number of buttons</b> ?
Q5	What is the <b>most difficult part of the project</b> <b>from the industrial engineer’s point of view</b> ?

asked well-trained annotators to label 1000 questions with their subject and aspects. Each question has been assigned to two annotators, and we report the agreement rate between annotators in Section 4. Further, we employ the trained CRF model to extract subjects and aspects from the questions.

### 3.4 Questionnaire Generation

Given a meeting transcript, for each of its segment  $T$  which was initially supposed to coherently point out one subject, we generate  $Q_T$ , a set of generated questions from  $T$ . Further, We create a set  $S_{Q_T}$  by extracting the subjects from each question in  $Q_T$ . Therefore, for the segment  $T$ , we have  $|Q_T|$  number of subjects. Extracted subjects from a question set with the same origin segment must be normalized so that one comprehensive, general, and informative subject presents a segment. The more the selected subject representative covers other concepts in  $S_{Q_T}$ , the better normalization we employed. This subject normalization reduces the number of subjects shown to the user at the first step of the questionnaire and will decrease the user’s effort, and it is aligned with our goal, i.e., figuring out users’ preferences by asking them the minimum number of questions. In other words, our goal is to select a single subject  $S_{norm}$  from  $S_{Q_T}$  which represents  $S_{Q_T}$  in the most informative way. To do so, we define the notion of the subject

network as follows.

**Definition 3.1.** Given a segment  $T$ , a set of generated questions  $Q_T$ , and extracted subjects  $S_{Q_T}$ , a subject-network for  $G(S_{Q_T})$  is denoted as  $G(S_{Q_T}) = (\mathbb{V}, \mathbb{E}, w)$ . It is a weighted undirected graph, where  $\mathbb{V} = \{s_i \in S_{Q_T}\}$ , and  $\mathbb{E} = \{e_{s_i, s_j} : \forall s_i, s_j \in \mathbb{V}\}$ . The function  $w : \mathbb{E} \rightarrow [0, 1]$  is the cosine similarity between the semantic relatedness of the contextualized embedding vectors of two incident subjects of an edge  $e_{s_i, s_j}$ , i.e.,  $v_{s_i}$  and  $v_{s_j}$ .

In Def. 3.1, we propose a subject-network for the question set  $Q_T$  where subjects are connected, and edge weights represent the semantic similarity between the two subjects. We hypothesize that the node with highest similarity and connection to others is the most central one. In other words, since it has great similarity to other subjects, there is a high probability that it points to a more generic concept and that covers the other subjects. Hence, the node  $S_{norm}$  should have high centrality attribute to represent the main subject of segment  $S$ . We employed PageRank (Haveliwala, 2003) value to find the most important and informative node in this network. Similarly, PageRank has shown to have a high correlation with the most important nodes and has been used in tackling different tasks such as quantifying term’s specificity or ranking problems in different information retrieval tasks (Arabzadeh et al., 2020, 2019; Kurland and Lee, 2010). We measure the PageRank score of each node and select the node with the highest PageRank value as the representative subject  $S_{norm}$  of the subject set  $S_{Q_T}$  for segment  $T$ . In other words, we represent each segment  $T$  by subject  $S_{norm}$  where  $PageRank(S_{norm}) > PageRank(s_i)$  for every  $s_i \in \mathbb{V}$ .

Fig. 4 displays a subject-network generated from extracted subjects from one of the meetings’ segments in the QMSUM dataset. subjects such as “Education”, “Schools,” “Young people who are leaving school” are included in this subject set and represented by nodes in this subject-network. Further, we connect every pair of nodes in this graph, and the edge weight is directly related to their semantic similarity. As presented in Fig. 4, some nodes have higher edge weights which their connected lines are shown with greater width. We measure page rank in this weighted network. Here “Education” got the highest PageRank value in this subject-network. Hence, we present these subjects

by one subject, i.e., “Education”. “Education” can be a promising representative for these subjects as it covers more specific concepts such as “schools”, “statutory education,” and “post 12 education.”

Next, the extracted aspects from each question set should be mapped to their representative subject. First, we remove the redundant and repetitive aspects and subjects by removing those who have highly similar n-grams. In addition, There might be several subjects existing in  $S_{Q_T}$  which all point out to  $S_{norm}$ , and they might be semantically very similar. Thus, in this step, we must be concerned not to lose any aspect because of subject normalization. We aim to map every aspect from  $S_{norm}$  and every  $s_i$  in  $S_{Q_T}$  which is highly similar to  $S_{norm}$  to maximize the potential of questions we might want to show at the end of the questionnaire. For instance, in Fig 3 we display a few extracted subjects and aspects from one segment. If we only consider “education” and its related aspect, we will lose many aspects that users might be interested in, and as a result, the questionnaire coverage will drop. On the other hand, if we merge the highly similar representative subjects with, e.g., “school setting” and “Education and Skills Committee,” we will have a broader host of questions to suggest to users. Therefore, we will filter out dissimilar subjects from  $S_{Q_T}$  to  $S_{norm}$  and map extracted aspects from filtered  $S_{Q_T}$  to  $S_{norm}$  as it is shown in Fig. 3. As a result, if “education” is the subject of interest for a user, they have the opportunity to select which aspects of education they are more interested in, such as “Role” of education or “challenges” of education. Finally, we will show users the questions in which the selected aspects and normalized subjects have appeared.

## 4 Evaluation Methodology

### 4.1 Dataset

For experiments, we use the QMSUM dataset (Zhong et al., 2021), which includes 232 different type of meetings: product (Carletta et al., 2005), academic (Janin et al., 2003), and committee<sup>34</sup>. The dataset consists of 162, 35 and 35 meetings for training, validation and testing purposes respectively. Each meeting comes with a set of general and specific questions; the general ones are out of the scope of this work since they

<sup>3</sup><https://www.ourcommons.ca/Committees/en/Home>

<sup>4</sup><https://record.assembly.wales/>



Table 2: Annotators agreement on annotated questions with respect to subjects and aspects using Krippendorff’s score (Krippendorff, 2011)

	Subject	Aspect
Hard [Exact Match]	0.459	0.415
Soft [At least one term matched]	0.490	0.485

refer to very broad concepts, e.g., “*summarize the whole meeting.*” Further evaluations are conducted on the QMSUM test set.

## 4.2 Evaluating Framework Components

The proposed framework consists of several steps (Fig. 2). The used *meeting segmentation* (Solbiati et al., 2021) is the best in class model<sup>5</sup>. Hence, we refer to original paper for evaluation results.

**Evaluating Question Generation** We evaluate the quality of our generated questions by measuring the fraction of generated questions by human annotators in QMSUM that we covered in PREME. We assume the specific queries in the QMSUM dataset enjoy relatively high quality because annotators issued them after comprehensively reading the transcript (gold standard questions). Hence, Fig. 5 reports the similarity between most similar questions generated by PREME and the gold questions by three different similarity metrics i.e., Sentence-BERT similarity (Reimers and Gurevych, 2019), Rouge F-1 score (Lin, 2004), and BLEU score (Papineni et al., 2002). We assume a questions from QMSUM is covered if there is atleast a question generated by PREME that has similarity is higher than a certain threshold  $t \in [1, 0.9, \dots, 0.1, 0]$ . We report the percentage of ‘Covered/Not Covered’ questions based on different similarity matching thresholds. Based on Fig. 5 we conclude while we cover a relatively fair number of specific questions, there is still room for improvement. However, we should note that the questions in QMSUM are very limited, and initially, they were not supposed to cover all possible questions that one could raise from the meeting. Additionally, we observe that questions in QMSUM, which are issued by humans, include more abstractive questions while our generated questions inclined toward more factual ones.

**Evaluating Subject and Aspect Extraction** To assess the quality of the collected dataset, we mea-

<sup>5</sup>The topic segmentation method has been evaluated on two of the three meetings used in this paper and has shown to outperform baseline (Hearst, 1997; Beeferman et al., 1999; Badjatiya et al., 2018)

Table 3: CRF performance on extracting subjects and aspects of questions using 10-fold cross validation

	Precision	Recall	F1-Score
Subject	0.64	0.69	0.67
Aspect	0.89	0.80	0.84
N/A	0.63	0.73	0.68

Table 4: Test set statistics and PREME Performance: Average number of generated questions and Coverage.

	#Meetings	Average # Turns	Average # Questions	Coverage (%)
Academic	9	893	1257	83.07%
Committee	6	214	1105	64.04%
Product	20	569	724	86.25%
All	35	591	927	81.62%

sure Krippendorff’s alpha agreement between annotators (Krippendorff, 2011) for extracted *subject* and *aspect* of the 1000 questions generated from the training set. Tab. 2 shows annotators have agreement  $\sim 0.4$ , which is “Moderate”, and it is a good agreement for such a challenging task. Since different annotators might selected different section of the text, Tab. 2 reports both *hard* and *soft* agreements. we trained the CRF model using *crfsuite* library and evaluated it by 10-fold cross-validation. Given each term in the questions, the model predicts whether the term is considered the subject, aspect, or not applicable for labeling (N/A). Tab. 3 shows the result of the CRF model evaluation in terms of precision, recall, and F1 scores. We notice that the model shows better performance on detecting aspects compared to the subject.

## 4.3 Evaluating Questionnaires

To the best of our knowledge, we are first to propose a preference-based questionnaire as a way for meeting exploration; thus, no particular gold standard benchmark or evaluation metrics. We introduce a new evaluation strategy that satisfies the desired properties on coverage (P1) and the existence of answers in the transcript (P2). Since we require users to express their preference, it makes it challenging to simulate ‘*enough imaginative context*’ among annotators. The proposed automatic metrics give good insights if our framework is ready to be tested through a user study in the future.

For our experimentation, we utilize the model SOTA called Locator in (Zhong et al., 2021) in which, given the query, it can extract the relevant spans from the meeting. The Locator employs a hierarchical ranking-based model structure based on CNN (Kim, 2014) and Transformers (Vaswani

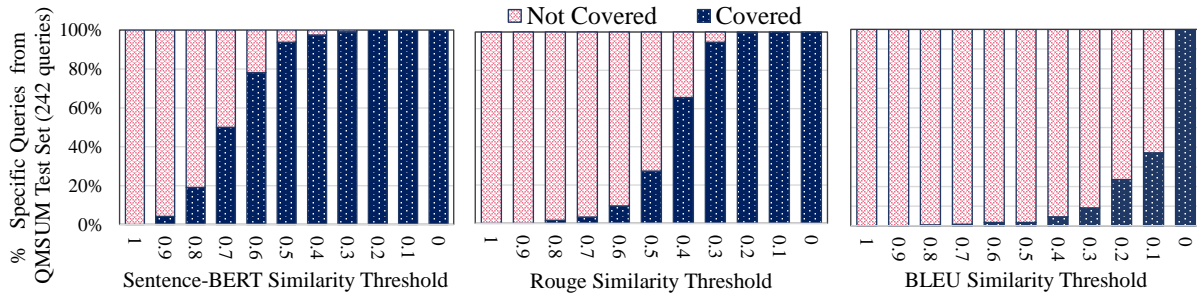


Figure 5: Coverage of specific queries in QMSUM test set among our generated questions considering different similarity metrics and threshold as coverage definition.

et al., 2017) architecture. The Locator embeds each utterance of the meeting and feeds it to a CNN network by capturing the local features, and utilize Transformer layers to obtain contextualized turn-level representations. In addition, the speaker’s embedding is also concatenated to the features list. Finally, the model uses MLP to score each turn, and the turns with the highest scores are considered the relevant spans for each question.

To measure the coverage (to satisfy **P1**), we adopt the newly proposed QA-style of evaluation (Deutsch et al., 2020; Wang et al., 2020). We define the coverage as the fraction of the meeting that is covered by the located relevant spans. We believe that that is a promising indicator of questionnaire informativeness. We run our experiments on the QMSUM test set. Tab. 4 shows the details of this test set. We over generate the questions and after removing the duplicates, on average, the questionnaire has 1257 unique questions from Academic meetings, 1105 questions from Committee meetings, and 724 questions from Product meetings. Further, Tab. 4 reports the percentage of utterances covered in each meeting. On average, our proposed questionnaire can cover 81% of the meeting. We also compared the coverage on the three types of meetings, i.e., the product vs. education and academic. As shown in this Table, While our generated questionnaire covered Committee meetings the least (64%), the Product and Academic meetings show higher coverage (over 80%).

Further, we evaluate how much the generated questions in PREME are answerable (to satisfy **P2**). Inspired by (Krishna and Iyer, 2019), we run a pretrained QA model (Sanh et al., 2019) over generated questions and report the confidence score for each QA pair in Fig. 6. We use DistilBERT fine-tuned on SQUAD (Rajpurkar et al., 2016) dataset<sup>6</sup>. We observe that more than 73% of generated questions from PREME on meetings in test set

<sup>6</sup><https://huggingface.co/distilbert-base-cased-distilled-squad>

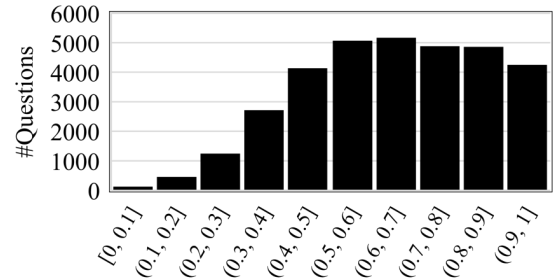


Figure 6: Histogram of Confidence Scores of Question-Answering (Sanh et al., 2019) model on generated questions from PREME.

of QMSUM shows confidence score higher than 0.5 and more than 42% of questions shows confidence score greater than 0.7. The results confirm that a promising portion of generated questions from PREME are answerable.

## 5 Conclusions and Future Work

Due to the increasing amount of meeting transcripts, there is a need for automatic tools for interactive preference-driven exploration that allows to quickly examine a meeting even if it was not attended or has been forgotten. We have proposed an end-to-end framework, called PREME, that allows automatically build a questionnaire that will enable users to explore the most of discussed subjects and their aspects if desired. As a result, users are supplied with high-quality questions about the meetings that express their information needs, and answers can be found in the transcript. We have proposed an automatic end-to-end evaluation strategy to demonstrate the desired properties (**P1** and **P2**) of the generated questionnaires since simulating actual preferences is challenging with hired annotators. The future works should include a user study that would enable real user interactions with generated questionnaires.

We publicly release the collected dataset of annotated questions concerning its subjects and aspects, the code for questionnaires generation, and our evaluation procedure to carry forward the proposed state-of-the-art for the newly formulated problem.



636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691

## References

Oren Anava, Shahar Golan, Nadav Golbandi, Zohar Karnin, Ronny Lempel, Oleg Rokhlenko, and Oren Somekh. 2015. Budget-constrained item cold-start handling in collaborative filtering recommenders via optimal design. In *Proceedings of the 24th international conference on world wide web*, pages 45–54.

Negar Arabzadeh, Fattane Zarrinkalam, Jelena Jovanovic, Feras Al-Obeidat, and Ebrahim Bagheri. 2020. Neural embedding-based specificity metrics for pre-retrieval query performance prediction. *Information Processing & Management*, 57(4):102248.

Negar Arabzadeh, Fattaneh Zarrinkalam, Jelena Jovanovic, and Ebrahim Bagheri. 2019. Geometric estimation of specificity within embedding spaces. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*, pages 2109–2112.

Pinkesh Badjatiya, Litton J Kurisinkel, Manish Gupta, and Vasudeva Varma. 2018. Attention-based neural text segmentation. In *European Conference on Information Retrieval*, pages 180–193. Springer.

Tal Baumel, Matan Eyal, and Michael Elhadad. 2018. Query focused abstractive summarization: Incorporating query relevance, multi-document coverage, and summary length constraints into seq2seq models. *arXiv preprint arXiv:1801.07704*.

Doug Beeferman, Adam Berger, and John Lafferty. 1999. Statistical models for text segmentation. *Machine learning*, 34(1):177–210.

Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*.

Jean Carletta, Simone Ashby, Sebastien Bourban, Mike Flynn, Mael Guillemot, Thomas Hain, Jaroslav Kadlec, Vasilis Karaiskos, Wessel Kraaij, Melissa Kronenthal, et al. 2005. The ami meeting corpus: A pre-announcement. In *International workshop on machine learning for multimodal interaction*, pages 28–39. Springer.

Asli Celikyilmaz, Antoine Bosselut, Xiaodong He, and Yejin Choi. 2018. Deep communicating agents for abstractive summarization. *arXiv preprint arXiv:1803.10357*.

Eunsol Choi, He He, Mohit Iyyer, Mark Yatskar, Wentau Yih, Yejin Choi, Percy Liang, and Luke Zettlemoyer. 2018. Quac: Question answering in context. *arXiv preprint arXiv:1808.07036*.

Konstantina Christakopoulou, Filip Radlinski, and Katja Hofmann. 2016. Towards conversational recommender systems. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 815–824.

Daniel Deutsch, Tania Bedrax-Weiss, and Dan Roth. 2020. Towards question-answering as an automatic metric for evaluating the content quality of a summary. *CoRR*, abs/2010.00490. 692  
693  
694  
695

Alexander R. Fabbri, Patrick Ng, Zhiguo Wang, Ramesh Nallapati, and Bing Xiang. 2020. Template-based question generation from retrieved sentences for improved unsupervised question answering. *CoRR*, abs/2004.11892. 696  
697  
698  
699  
700

Taher H Haveliwala. 2003. Topic-sensitive pagerank: A context-sensitive ranking algorithm for web search. *IEEE transactions on knowledge and data engineering*, 15(4):784–796. 701  
702  
703  
704

Marti A Hearst. 1997. Text tiling: Segmenting text into multi-paragraph subtopic passages. *Computational linguistics*, 23(1):33–64. 705  
706  
707

Michael Heilman and Noah A Smith. 2010. Good question! statistical ranking for question generation. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 609–617. 708  
709  
710  
711  
712  
713

Luyang Huang, Shuyang Cao, Nikolaus Parulian, Heng Ji, and Lu Wang. 2021. Efficient attentions for long document summarization. *arXiv preprint arXiv:2104.02112*. 714  
715  
716  
717

Tai-Chia Huang, Chia-Hsuan Hsieh, and Hei-Chia Wang. 2018. Automatic meeting summarization and topic detection system. *Data Technologies and Applications*. 718  
719  
720  
721

Tatsuya Ishigaki, Hen-Hsen Huang, Hiroya Takamura, Hsin-Hsi Chen, and Manabu Okumura. 2020. Neural query-biased abstractive summarization using copying mechanism. In *European Conference on Information Retrieval*, pages 174–181. Springer. 722  
723  
724  
725  
726

François Jacquenet, Marc Bernard, and Christine Largeron. 2019. Meeting summarization, a challenge for deep learning. In *International Workshop Conference on Artificial Neural Networks*, pages 644–655. Springer. 727  
728  
729  
730  
731

Adam Janin, Don Baron, Jane Edwards, Dan Ellis, David Gelbart, Nelson Morgan, Barbara Peskin, Thilo Pfau, Elizabeth Shriberg, Andreas Stolcke, et al. 2003. The icsi meeting corpus. In *2003 IEEE International Conference on Acoustics, Speech, and Signal Processing, 2003. Proceedings.(ICASSP'03).*, volume 1, pages I–I. IEEE. 732  
733  
734  
735  
736  
737  
738

Bin Jiang, Jian Pei, Xuemin Lin, David W Cheung, and Jiawei Han. 2008. Mining preferences from superior and inferior examples. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 390–398. 739  
740  
741  
742  
743

Yoon Kim. 2014. Convolutional neural networks for sentence classification. *CoRR*, abs/1408.5882. 744  
745

746	Klaus Krippendorff. 2011. Computing krippendorff’s alpha-reliability.	801
747		802
748	Kalpesh Krishna and Mohit Iyyer. 2019. Generating question-answer hierarchies. <i>arXiv preprint arXiv:1906.02622</i> .	803
749		804
750		805
751	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> .	806
752		807
753		808
754		809
755		810
756	Sayali Kulkarni, Sheide Chammas, Wan Zhu, Fei Sha, and Eugene Ie. 2020. Aquamuse: Automatically generating datasets for query-based multi-document summarization. <i>arXiv preprint arXiv:2010.12694</i> .	811
757		812
758		813
759		814
760	Sayali Kulkarni, Sheide Chammas, Wan Zhu, Fei Sha, and Eugene Ie. 2021. Comsum and sibert: A dataset and neural model for query-based multi-document summarization. In <i>International Conference on Document Analysis and Recognition</i> , pages 84–98. Springer.	815
761		816
762		817
763		818
764		819
765		820
766	Oren Kurland and Lillian Lee. 2010. Pagerank without hyperlinks: Structural reranking using links induced by language models. <i>ACM Transactions on Information Systems (TOIS)</i> , 28(4):1–38.	821
767		822
768		823
769		824
770	Igor Labutov, Sumit Basu, and Lucy Vanderwende. 2015. Deep questions without deep understanding. In <i>Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 889–898.	825
771		826
772		827
773		828
774		829
775		830
776		831
777	Timothy J. Lawson, James H. Bodle, Melissa A. Houlette, and Richard R. Haubner. 2006. Guiding questions enhance student learning from educational videos. <i>Teaching of Psychology</i> , 33(1):31–33.	832
778		833
779		834
780		835
781	Timothy J Lawson, James H Bodle, and Tracy A McDonough. 2007. Techniques for increasing student learning from educational videos: Notes versus guiding questions. <i>Teaching of Psychology</i> , 34(2):90–93.	836
782		837
783		838
784		839
785		840
786	Adam D Lelkes, Vinh Q Tran, and Cong Yu. 2021. Quiz-style question generation for news stories. In <i>Proceedings of the Web Conference 2021</i> , pages 2501–2511.	841
787		842
788		843
789		844
790	Manling Li, Lingyu Zhang, Heng Ji, and Richard J. Radke. 2019. <a href="#">Keep meeting summaries on topic: Abstractive multi-modal meeting summarization</a> . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 2190–2196, Florence, Italy. Association for Computational Linguistics.	845
791		846
792		847
793		848
794		849
795		850
796		851
797	Chin-Yew Lin. 2004. <a href="#">ROUGE: A package for automatic evaluation of summaries</a> . In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain. Association for Computational Linguistics.	852
798		853
799		854
800		855
		856
	David Lindberg, Fred Popowich, John Nesbit, and Phil Winne. 2013. Generating natural language questions to support learning on-line. In <i>Proceedings of the 14th European Workshop on Natural Language Generation</i> , pages 105–114.	
	Marina Litvak and Natalia Vanetik. 2017. <a href="#">Query-based summarization using MDL principle</a> . In <i>Proceedings of the MultiLing 2017 Workshop on Summarization and Summary Evaluation Across Source Types and Genres</i> , pages 22–31, Valencia, Spain. Association for Computational Linguistics.	
	Cong Liu, Chi Yuan, Alex M Butler, Richard D Carvajal, Ziran Ryan Li, Casey N Ta, and Chunhua Weng. 2019. Dquest: dynamic questionnaire for search of clinical trials. <i>Journal of the American Medical Informatics Association</i> , 26(11):1333–1343.	
	Yang Liu and Mirella Lapata. 2019. Text summarization with pretrained encoders. <i>arXiv preprint arXiv:1908.08345</i> .	
	Potsawee Manakul and Mark Gales. 2021. Long-span summarization via local attention and content selection. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 6026–6041.	
	Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. <a href="#">On faithfulness and factuality in abstractive summarization</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 1906–1919, Online. Association for Computational Linguistics.	
	Karen Mazidi and Rodney D. Nielsen. 2014. <a href="#">Linguistic considerations in automatic question generation</a> . In <i>Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 321–326, Baltimore, Maryland. Association for Computational Linguistics.	
	Yashar Mehdad, Giuseppe Carenini, Frank Tompa, and Raymond Ng. 2013. Abstractive meeting summarization with entailment and fusion. In <i>Proceedings of the 14th European Workshop on Natural Language Generation</i> , pages 136–146.	
	Jack Mostow and Wei Chen. 2009. Generating instruction automatically for the reading strategy of self-questioning. In <i>AIED</i> , pages 465–472.	
	Gabriel Murray, Giuseppe Carenini, and Raymond Ng. 2010. Generating and validating abstracts of meeting conversations: a user study. In <i>Proceedings of the 6th International Natural Language Generation Conference</i> .	
	Ramesh Nallapati, Feifei Zhai, and Bowen Zhou. 2017. Summarunner: A recurrent neural network based sequence model for extractive summarization of documents. In <i>Thirty-First AAAI Conference on Artificial Intelligence</i> .	

857	Preksha Nema, Mitesh M. Khapra, Anirban Laha, and	Abigail See, Peter J. Liu, and Christopher D. Manning.	911
858	Balaraman Ravindran. 2017. <a href="#">Diversity driven attention model for query-based abstractive summarization</a> .	2017. <a href="#">Get to the point: Summarization with pointer-generator networks</a> . <i>CoRR</i> , abs/1704.04368.	912
859	In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1063–1072, Vancouver, Canada. Association for Computational Linguistics.		913
860			
861		Imran Sehikh, Dominique Fohr, and Irina Illina. 2017.	914
862		Topic segmentation in asr transcripts using bidirectional rnns for change detection. In <i>2017 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)</i> , pages 512–518. IEEE.	915
863			916
864	Tatsuro Oya, Yashar Mehdad, Giuseppe Carenini, and Raymond Ng. 2014. <a href="#">A template-based abstractive meeting summarization: Leveraging summary and source text relationships</a> .	Anna Sepliarskaia, Julia Kiseleva, Filip Radlinski, and Maarten de Rijke. 2018. Preference elicitation as an optimization problem. In <i>Proceedings of the 12th ACM Conference on Recommender Systems</i> , pages 172–180.	917
865			918
866			919
867			920
868			921
869			922
870			923
871			924
872	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.	Guokan Shang, Wensi Ding, Zekun Zhang, Antoine Tixier, Polykarpos Meladianos, Michalis Vazirgiannis, and Jean-Pierre Lorré. 2018. <a href="#">Unsupervised abstractive meeting summarization with multi-sentence compression and budgeted submodular maximization</a> . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 664–674, Melbourne, Australia. Association for Computational Linguistics.	925
873			926
874			927
875			928
876			929
877			930
878			931
879			932
880			933
881			934
882			935
883			936
884			937
885			938
886			939
887			940
888			941
889			942
890			943
891			944
892			945
893			946
894			947
895			948
896			949
897			950
898			951
899			952
900			953
901			954
902			955
903			956
904			957
905			958
906			959
907			960
908			961
909			962
910			963
			964
			965



- 966 gap-sentences for abstractive summarization. In *International Conference on Machine Learning*, pages  
967 11328–11339. PMLR.  
968
- 969 Leilan Zhang and Qiang Zhou. 2019. Topic segmenta-  
970 tion for dialogue stream. In *2019 Asia-Pacific Sig-  
971 nal and Information Processing Association Annual  
972 Summit and Conference (APSIPA ASC)*, pages 1036–  
973 1043. IEEE.
- 974 Zhou Zhao, Haojie Pan, Changjie Fan, Yan Liu, Lin-  
975 lin Li, Min Yang, and Deng Cai. 2019. Abstrac-  
976 tive meeting summarization via hierarchical adap-  
977 tive segmental network learning. In *The World Wide  
978 Web Conference*, pages 3455–3461.
- 979 Ming Zhong, Da Yin, Tao Yu, Ahmad Zaidi, Mutethia  
980 Mutuma, Rahul Jha, Ahmed Hassan Awadallah, Asli  
981 Celikyilmaz, Yang Liu, Xipeng Qiu, and Dragomir  
982 Radev. 2021. [QMSum: A new benchmark for query-  
983 based multi-domain meeting summarization](#). In *Pro-  
984 ceedings of the 2021 Conference of the North Amer-  
985 ican Chapter of the Association for Computational  
986 Linguistics: Human Language Technologies*, pages  
987 5905–5921, Online. Association for Computational  
988 Linguistics.
- 989 Chenguang Zhu, Ruochen Xu, Michael Zeng, and Xue-  
990 dong Huang. 2020a. [A hierarchical network for ab-  
991 stractive meeting summarization with cross-domain  
992 pretraining](#). In *Proceedings of the 2020 Conference  
993 on Empirical Methods in Natural Language Pro-  
994 cessing: Findings, EMNLP 2020, Online Event, 16-  
995 20 November 2020*, pages 194–203. Association for  
996 Computational Linguistics.
- 997 Chenguang Zhu, Ruochen Xu, Michael Zeng, and Xue-  
998 dong Huang. 2020b. [A hierarchical network for ab-  
999 stractive meeting summarization with cross-domain  
1000 pretraining](#). In *Findings of the Association for Com-  
1001 putational Linguistics: EMNLP 2020*, pages 194–  
1002 203, Online. Association for Computational Linguis-  
1003 tics.