

A Personalized Reinforcement Learning Summarization Service for Learning Structure from Unstructured Data

Samira Ghodratnama
Macquarie University, Australia
W.W. Grainger, USA
samira.ghodratnama@mq.edu.au
samira.ghodratnama@grainger.com

Amin Behehsti
Macquarie University, Australia
amin.behehsti@mq.edu.au

Mehrdad ZakershahraK
Macquarie University, Australia
mehrdad.zakershahrak@mq.edu.au

Abstract—The exponential growth of textual data has created a crucial need for tools that assist users in extracting meaningful insights. Traditional document summarization approaches often fail to meet individual user requirements and lack structure for efficient information processing. To address these limitations, we propose *Summation*, a hierarchical personalized concept-based summarization approach. It synthesizes documents into a concise hierarchical concept map and actively engages users by learning and adapting to their preferences. Using a Reinforcement Learning algorithm, *Summation* generates personalized summaries for unseen documents on specific topics. This framework enhances comprehension, enables effective navigation, and empowers users to extract meaningful insights from large document collections aligned with their unique requirements.

Index Terms—Document summarization, personalized summarization, hierarchical summarization, concept-based summarization.

I. INTRODUCTION

The availability of a vast amount of information on various topics has led to a phenomenon known as *information overload*, where the volume of data exceeds an individual’s capacity for effective processing within a reasonable timeframe. While this abundance of data can be valuable for analytical applications, it necessitates efficient exploration tools to harness its potential benefits without succumbing to information overload, which can strain cognitive resources. Data summaries serve as effective tools for gathering relevant information, organizing it into a coherent and manageable form, and facilitating complex question answering, insight generation, and conceptual boundary discovery [1]–[3]. Automatic document summarization has been extensively studied to address the challenges of data reduction for analysis, commercialization, management, and personalization purposes. Furthermore, users often seek information in an organized and coherent structure. However, despite the speed of document generation and the massive collections of unstructured documents, producing personalized summaries comparable to human-written ones remains challenging. Most previous work on automatic text summarization has focused on generating textual summaries rather than structured ones. These approaches typically pro-

duce a single, short, general, and flat summary that applies to all users, lacking interpretability and personalization. Moreover, they are incapable of producing more extended and detailed summaries, even if users express interest in obtaining additional information. Additionally, the lack of structure in these summaries hampers further processing, and they heavily rely on reference or gold summaries created by humans, which are subjective and costly [4], [5]. To address these limitations, we propose *Summation*, a hierarchically interactive structured summarization approach that generates personalized summaries. We emphasize the significance of the following aspects in our contribution: i) Structured summaries, ii) Personalization, iii) Interaction, and iv) The elimination of reference summaries.

Structured Summaries. Studies have demonstrated that when individuals encounter numerous documents, they seldom formulate fully-fledged summaries. Instead, they attempt to extract concepts and understand the relationships among them [6]–[8]. Consequently, structured data has become crucial in various domains. It offers a concise overview of the document collection’s contents, unveils interesting relationships, and serves as a navigational structure for further exploration of the documents. Our approach, *Summation*, provides summaries in the form of a hierarchical concept map, which caters to diverse user requirements by being interpretable, concise, and simultaneously providing an overview and detailed information.

Personalization. Existing summarization approaches typically generate a generic summary comprising a few selected sentences intended to meet the needs of all users. In contrast to such generic summaries, there is a dearth of user-centric summarization approaches that allow users to specify the desired content in the summaries [9], [10].

Interaction. Conventional summarization approaches treat a topic-related document set as input and generate a summary that captures the most salient aspects. However, research on this topic often neglects the usefulness of the approach for users, focusing primarily on the accuracy of the generated summaries. As a result, these approaches produce short (3-

6 sentences), inflexible, and flat summaries that are the same for all users. Consequently, these approaches fail to provide more extensive summaries even when users express interest in obtaining additional information.

Reference Summaries. Traditional document summarization techniques rely on reference summaries created by humans for training their systems. However, this approach is subjective and, more importantly, resource-intensive. For instance, Lin [11] reported that creating summaries for the Document Understanding Conferences (DUC) required 3,000 hours of human effort. Personalized summaries eliminate the need for such reference summaries by generating specific summary for a user instead of optimizing a summary for all users.

Our Contribution. We study the automatic creation of personalized, structured summaries, allowing the user to overview a document collection’s content without much reading quickly. The goal here is to dynamically maintain a federated summary view incrementally, resulting in a unified framework for intelligent summary generation and data discovery tools from a user-centered perspective. The unique contribution of this paper includes:

- We provide summaries in the form of a *hierarchical concept map*, labeled graphs representing concepts and relationships in a visual and concise format. Their structured nature can reveal interesting patterns in documents that users would otherwise need to discover manually. It enables providing more information than traditional approaches within the same limit size. It can be used as a navigator in the document collection. Such visualization is beneficial for decision-making systems.
- We introduce and formalize a theoretically grounded method. We propose a personalized interactive summarization approach utilizing a reinforcement learning algorithm to learn generating user-adapted results. It is the first approach to predict users’ desired structured summary to the best of our knowledge.
- We provide various evidence evaluating different aspects to prove *Summation*’s usability using human and automatic evaluation.

We divide the proposed framework into two steps. The first step is *organizer* which structure unstructured data by making a hierarchical concept map. Then *summarizer* is responsible for: i) predicting users’ preferences based on the given feedback by employing *preference learning* and ii) learning to provide personalized summaries by leveraging reinforcement learning. A general overview of the algorithm is depicted in Figure 1.

II. RELATED WORK

We categorize previous approaches into three groups including *traditional approaches*, *structured approaches*, *personalized and interactive approaches* discussed below.

Traditional Approaches. A good summary should provide the maximum information about the input documents within a size limit and be fluent and natural. Different aspects for categorizing traditional multi-document summarization approaches

exist, such as the input type, the process, and the summarization goal [7], [12]. However, the main category considers the process and the output type of the summarization algorithm: *extractive* and *abstractive* approaches. The input in both cases is a set of documents, and the output is a few sentences. Abstractive summaries are generated by interpreting the main concepts of a document and then stating those contents in another format. Therefore, abstractive approaches require deep natural language processing, such as semantic representation and inference [7]. However, extractive text summarization selects some sentences from the original documents as the summary. These sentences are then concatenated into a shorter text to produce a meaningful and coherent summary [13]. Early extractive approaches focused on shallow features, employing graph structure, or extracting the semantically related words [14]. Different machine learning approaches, such as naive-Bayes, decision trees, neural networks, and deep reinforcement learning models are used for this purpose [15]–[17].

Structured Approaches. While traditional summarization approaches produce unstructured summaries, there exist few attempts on structured summaries. Structured summaries are defined by generating Wikipedia articles and biographies to extract the significant aspects of a topic using approaches such as topic modeling or an entity-aspect LDA model [18], [19]. Discovering threads of related documents is another category of structured summaries. They mostly use a machine algorithm to find the threads using a supervised approach and features such as temporal locality of stories for event recognition and time-ordering to capture dependencies [20]. A few papers have examined the relationship between summarization and hierarchies. However, the concept of hierarchy in these approaches is the relation between different elements of a document. An example is creating a hierarchy of words or phrases to organize a set of documents [21]. There is a related thread of research on identifying the hierarchical structure of the input documents and generating a summary which prioritizes the more general information according to the hierarchical structure [22]. However, the information unit is a sentence, and the hierarchy is based on time measures. Concept-based multi-document summarization is a variant of traditional summarization that produces structured summaries using concept maps. It learns to identify and merge coreferent concepts to reduce redundancy and finds an optimal summary via integer linear programming. However, it produces a single flat summary for all users [23].

Personalized and Interactive Approaches. Recently, there exist few recent attempts on personalized and interactive approaches in different NLP tasks. Unlike non-interactive systems that only present the system output to the end-user, interactive NLP algorithms ask the user to provide certain feedback forms to refine the model and generate higher-quality outcomes tailored to the user. Multiple forms of feedback also have been studied including mouse-clicks for information retrieval [24], post-edits and ratings for machine translation [24], [25], error markings for semantic parsing [26], and preferences for translation [27]. A significant category

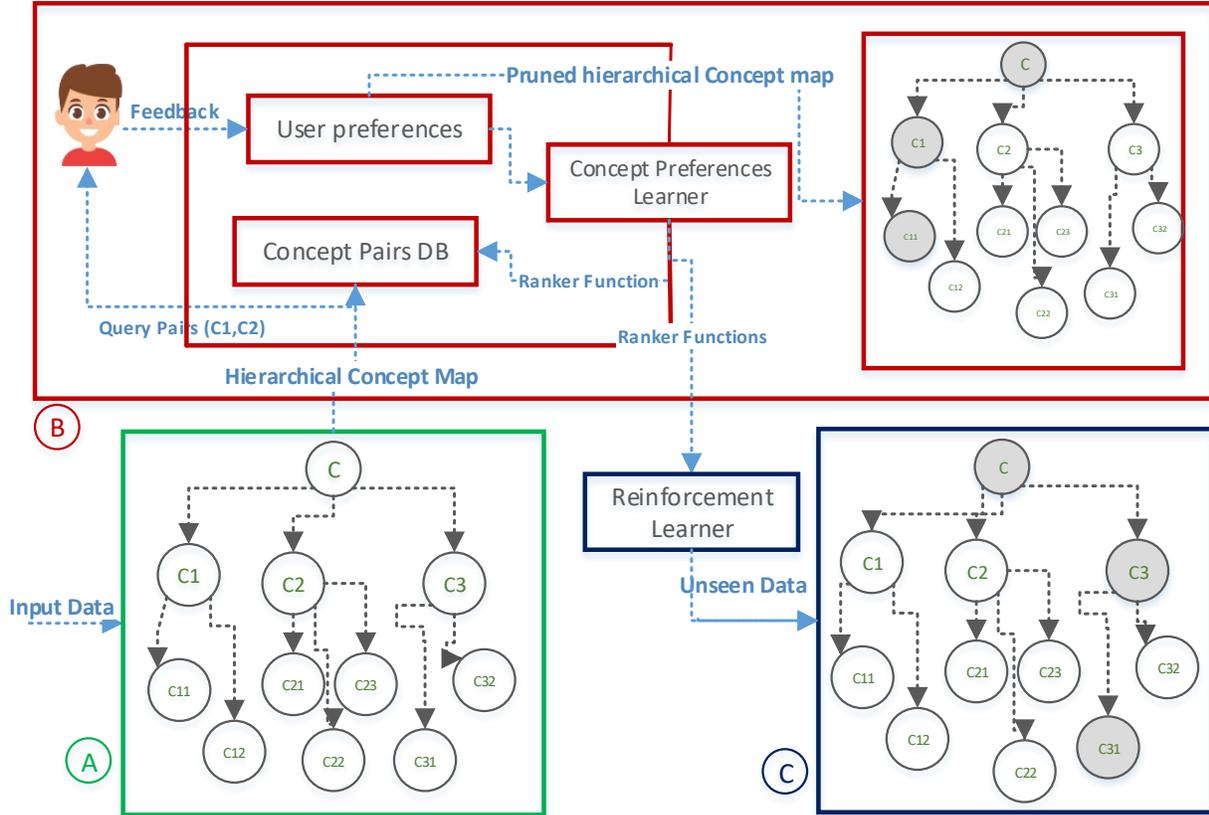


Fig. 1: An overview of *Summation*: A) The input data is converted to a hierarchical concept map (organizer). B) The summarizer is responsible for extracting and learning user preferences. C) Using reinforcement learning to predict users’ desired summary.

of interactive approaches presents the output of a given automatic summarization system to users as a draft summary, asking them to refine the results without further interaction. The refining process includes cutting, paste, and reorganize the essential elements to formulate a final summary [28], [29]. Other interactive summarization systems include the iNeATS [30] and IDS [31] systems that allow users to tune several parameters for customizing the produced summaries. Avinesh and Meyer [32] proposed the most recent interactive summarization approach that asks users to label important bigrams within candidate summaries. Their system can achieve near-optimal performance. However, labeling important bigrams is an enormous burden on the users, as users have to read through many potentially unimportant bigrams. Besides, it produces extractive summaries that are unstructured.

III. THE PROPOSED APPROACH (SUMMATION)

The ultimate goal of summarization is to provide a concise, understandable, and interpretable summary tailored to the users’ needs. However, making such a summary is challenging due to massive document collection, the speed of generated documents, and the unstructured format. In this regard, *Summation* aims to make structured summaries to facilitate further processes to make it concise and easily understandable while engaging users to create their personalized summaries. This

novel framework has two components: *organizer* and the *summarizer*. First, we discuss the problem definition, and then each component is explained.

Problem Definition. The input is a set of documents $D = \{D_1, D_2, \dots, D_N\}$ and each document consists of a sequence of sentences $S = [s_1, s_2, \dots, s_n]$. Each sentence s_i is a set of concepts $\{c_1, c_2, \dots, c_k\}$, where a concept can be a word (unigram) or a sequence of words. The output is a personalized hierarchical concept map. This novel framework has two components, an organizer and a summarizer, explained in Sec. III-A and III-B, respectively.

A. Adding Structure to Unstructured Data

The first step is to structure unstructured information by making a hierarchical concept map. A concept map is a graph with directed edges, where nodes indicate concepts and edges indicate relations. Both concepts and relations are sequences of related words representing a semantic unit. Consequently, the first step in creating a concept map is to identify all concepts and relations. Here, we propose hierarchical clustering to form the hierarchical concept map.

1) *Concept and Relation Extraction.*: Concepts come in different syntactic types, including nouns, proper nouns, more complex noun phrases, and verb phrases that describe activities [23]. For this purpose, we used open information

extraction (OIE) [33] through which the entities and relations are obtained directly from the text. OIE finds binary propositions from a set of documents in the form of (con_1, R, con_2) , which are equivalent to the desired concepts and relations. For example, the output for the sentence, ‘cancer treatment is underpinned by the Pharmaceutical Benefits Scheme’, is: *Cancer treatment* $\xrightarrow{\text{is underpinned}}$ *by the Pharmaceutical Benefits Scheme*

Balancing precision and recall in extracting concepts is a challenging task. A high precision causes to define all identified spans as mentions of concepts. Therefore, some constructions are usually missed, which leads to lowering the recall. On the other hand, a high recall is necessary since missed concepts can never be in summary. Obtaining a higher recall may extract too many mentions, including false positives. Generalizability is also essential. The reason is that extracting a particular syntactic structure might generate only correct mentions, causing too broad mentions. Ideally, a proper method applies to many text types. To avoid meaningless and long concepts, we processed the OIE results such that concepts with less than one noun token or more than five tokens are omitted. The original nouns also replace pronouns. If an argument is a conjunction indicating conj-dependency in the parse tree, we split them.

2) *Concept Map Construction.*: Among various extracted concepts and relations, multiple expressions can refer to the same concept while not using precisely the same words; that is, they can also use synonyms or paraphrases. However, distinguishing similar concepts to group them is challenging and subjective. For example, adding a modifier can completely change the meaning of a concept based on the purpose of summarization. Consequently, grouping them may lead to propositions that are not stated in the document. Therefore, we need to group every subset that contains mentions of a single, unique concept. Scalability is another critical issue. For example, pairwise comparisons of concepts cause a quadratic run-time complexity applicable only to limited-sized document sets. The same challenges exist for relation grouping. However, we first grouped all mentions by the concepts’ pairs, and then performed relation grouping. Therefore, this task’s scope and relevance are much smaller than when concepts are used. Therefore, in practise, comparison-based quadratic approaches are feasible. Moreover, as the final goal is to create a defined size summary, the summary size significantly affects the level of details in grouping concepts. This is because the distinction between different mentions of a concept might not be required, as it is a subjective task. Ideally, the decision to merge must be made based on the final summary map’s propositions to define the necessary concept granularity.

We further propose hierarchical conceptual clustering using k-means with word embedding vectors to tackle this problem, as it spans a semantic space. Therefore, word embedding clusters give a higher semantic space, grouping semantically similar word classes under the Euclidean metric constraint defined below. Before defining the proposed hierarchical con-

ceptual clustering, we review word embedding schemes used in the proposed model.

Word Embedding. Word embedding is a learnt representation of text such that the same meaning words have similar representations. Different techniques can be used to learn a word embedding from the text. Word2Vec [34] is an example of a statistical model for learning a word embedding representation from a text corpus, utilising different architectures. As such, we used skip-gram and bag of character n-grams in our experiments. The skip-gram model uses the current word for predicting the surrounding words by increasing the weights of nearby context words more than other words using a neural network model. One drawback of skip-gram is its inability to detect rare words. In another model, authors define an embedding method by representing each word as the sum of the vector representations of its character n-grams, known as ‘bag of character n-grams’ [35]. If the training corpus is small, character n-grams will outperform the skip-gram (of words) approach.¹

Conceptual Hierarchical Clustering. Given word (concept) embeddings learnt from a corpus, $\{v_{w_1}, v_{w_2}, \dots, v_{w_T}\}$, we propose a novel recursive clustering algorithm to form a hierarchical concept map, H . This variable denotes a set of concept maps organised into a hierarchy that incrementally maintains hierarchical summaries from the most general node (root) to the most specific summary (leaves). Within this structure, any non-leaf summary generalises the content of its children nodes. Hierarchical summarization has two critical strengths in the context of large-scale summarization. First, the initial information under review is small and grows upon users’ request, so as not to overwhelm them. Second, the parent-to-child links facilitate user navigation and drilling down for more details on interesting topics. The hierarchical conceptual clustering minimizes the objective function Eq. 1 over all k clusters as $C = \{c_1, c_2, \dots, c_k\}$.

$$J = \sum_{k=1}^K \sum_{t=1}^{|T|} |v_{w_t} - c_k|^2 + \alpha \min_{c \in C} \text{size}(c), \quad (1)$$

where c_k is the randomly selected centre k – th cluster, and T is the number of word vectors. The second term is the evenness of the clusters, added to avoid clusters with small sizes. α tunes the evenness factor, which was defined by employing a grid search over a development set. We also implemented hierarchical clustering top-down at each time, optimising Eq. 1. After defining the clusters, we must find the concept that best represents every concept at the lower levels to ensure hierarchical abstraction. A concise label is the desired label for each node; however, shortening mentions can introduce propositions that are not asserted by a text. For example, the concept labelled ‘students’ can change in meaning where the emphasis is on a few students or some students. To this end, a centre of a cluster at each level of the hierarchy was defined as a label. The inverse distance to

¹We used fastText for word embedding: <https://fasttext.cc/docs/en/support.html>

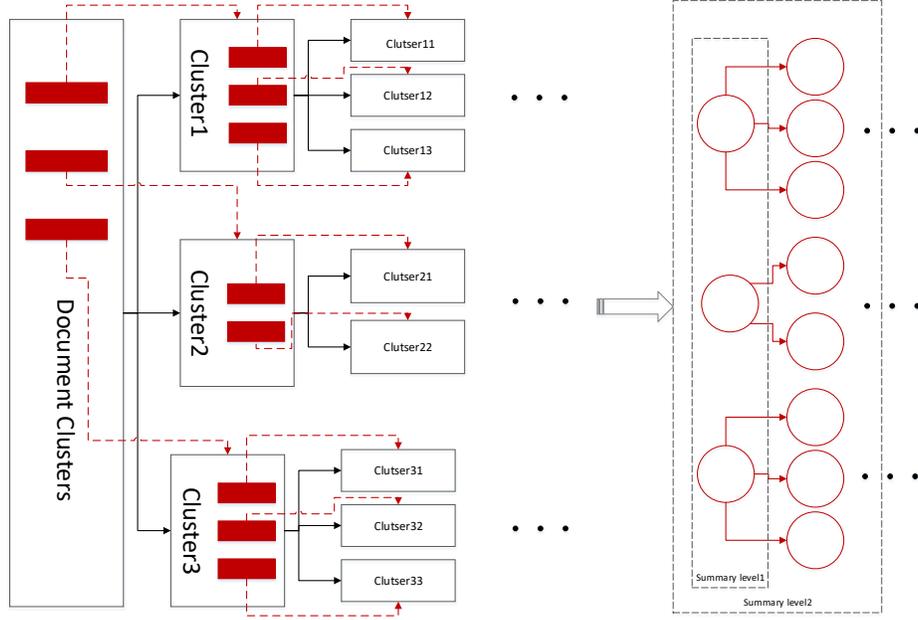


Fig. 2: The hierarchical conceptual clustering architecture.

the cluster centres is the membership degree or the similarity to each label. The cluster distance for a word w_t is defined as $d_{v_{w_t}}$. Consequently, the membership of each word w_t in cluster c_k to its label is the inverse distance defined in Eq. 2.

$$m_{v_{w_t}} = \frac{1}{d_{v_{w_t}}} = \frac{1}{|c_k - v_{w_t}|^2} \quad \forall w_t \in c_k \quad (2)$$

We then fine-tuned K within the 5–50 range based on the dataset size and chose the cluster number according to gap statistic value [36]. The output H can be directly used as a new dataset for other actions, such as browsing, querying, data mining process, or any other procedures requiring a reduced but structured version of data. The hierarchical clustering can also be pruned at each level to represent a summarised concept map for different purposes or users. Therefore, H is fed to the summariser for pruning to generate a personalized summary. Moreover, by using preference-based learning and RL, we learn users' preferences in making personalized summaries for unseen topic-related documents, discussed in Sec. III-B.

B. Summarizer

The hierarchical concept map produced in the previous step is given to the summariser to make the desired summaries for users based on their given preferences. Therefore, the summariser consists of two phases—(i) predicting user preferences and (ii) generating the desired summary.

1) *Predicting User Preference.*: The first step towards creating personalized summaries is to understand users' interests. It can be extracted implicitly based on users' profiles, browsing history, likes or dislikes, or retweeting in social media [37]. When this information is not available, interaction with users is an alternative to retrieve user's perspectives. The user

feedback can be in any form, such as mouse-click or post-edits, as explained in Section II. Preference-based interactive approaches are another form of feedback that puts a lower cognitive burden on human subjects [38]. For instance, asking users to select one concept among “cancer treatment” and “cancer symptoms” is more straightforward than asking for giving a score to each of these concepts. Therefore, in this paper, to reduce users' cognitive load, queries are in the form of concept preference. Preference learning is a classification method that learns to rank instances based on the observed preference information. It trains based on a set of pairwise preferred items and obtaining the total ranking of objects [39].

H is the hierarchical concept map, where at the i -th level of the hierarchy there exist m_i nodes defining a label l . $L = \{l_{11}, \dots, l_{nm_i}\}$ is the set of all labels, where l_{i1} indicates the first node at i -th level of the hierarchy and n is the number of levels, and L_i indicates the labels at i -th level. We queried users with a set of pairwise concepts at the same levels, $\{p(l_{i1}, l_{i2}), p(l_{i2}, l_{i3}), \dots, p(l_{im_i-1}, l_{im_i})\}$, where $p(l_{i1}, l_{i2})$ is defined in Eq. 3.

$$p(l_{i1}, l_{i2}) = \begin{cases} 1, & \text{if } l_{i1} > l_{i2} \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

where $>$ indicates the preference of l_{i1} over l_{i2} . Preference learning aims to predict the overall ranking of concepts, which requires transforms concepts into real numbers, called utility function. The utility function U such that $l_i > l_j \rightarrow U(l_i) > U(l_j)$, where U is a function $U : C \rightarrow \mathbb{R}$. In this problem, the ground-truth utility function (U) measures each concept's importance based on users' attitudes, defined as a regression learning problem. According to U , we defined the

ranking function, R , measuring the importance of each concept towards other concepts based on users’ attitude. This is defined in Eq. 4.

$$R(l_i) = \sum \mathbb{1}\{U(l_i) > U(l_j)\}, \forall l_i, l_j \in L \quad (4)$$

where $\mathbb{1}$ is the indicator function. The Bradley–Terry model [40], [41] is a probability model widely used in preference learning. Given a pair of individuals l_i and l_j drawn from some population, the model estimates the probability that the pairwise comparison $l_i > l_j$ is true. Having n observed preference items, the model approximates the ranking function R by computing the maximum likelihood estimate in Eq. 5.

$$J_x(w) = \sum_{i \in n} [p(l_i, l_j) \log F(l_i, l_j; w) + p(l_j, l_i) \log F(l_j, l_i; w)] \quad (5)$$

where $F(l)$ is the logistic function defined in Eq. 6.

$$F(l_i, l_j; w) = \frac{1}{1 + \exp[U^*(l_j; w) - U^*(l_i; w)]} \quad (6)$$

Here, U^* is the approximation of U parameterised by w , which can be learnt using different function approximation techniques. In our problem, a linear regression model was designed for this purpose, defined as $U(l; w) = w^T \phi(l)$, where $\phi(l)$ is the representation feature vector of the concept l . For any $l_i, l_j \in L$, the ranker prefers l_i over l_j if $w^T \phi(l_i) > w^T \phi(l_j)$.

By maximizing the $J_x(w)$ in Eq. 5, $w^* = \operatorname{argmax}_w J_x(w)$, the resulting w^* using stochastic gradient ascent optimisation will be used to estimate U^* , and consequently the approximated ranking function $R^* : C \rightarrow \mathbb{R}$. Thus, *Summation* learns a ranking over concepts and uses the ranking to generate personalized summaries.

2) *Generating Personalized Summaries.*: The summarization task is to transform the input (a cluster of documents) d to the best summary among all possible summaries, called $Y(d)$, for the learnt preference ranking function. This problem can be defined as a sequential decision-making problem, starting from the root, sequentially selecting concepts and adding them to a draft summary. Therefore, it can be defined as an MDP problem.

An MDP is a tuple (S, A, R, T) , where S is the set of states, A is the set of actions, $R(s, a)$ is the reward for performing an action (a) in a state (s), and T is the set of terminal states. In our problem, a state is a draft summary, and A includes two types of action—either adding a new concept to the current draft summary or terminating the construction process if it reaches users’ limit size. The reward function R returns an evaluation score in one of the termination states or 0 in other states.

A policy $\pi(s, a) : S \times A \rightarrow R$ in an MDP defines the selection of actions in state s . The goal of RL algorithms is to learn a policy that maximises the accumulated reward. The learnt policy trained on specific users’ interests is used on unseen data at the test time (in this problem to generate summaries in new and related topic documents).

Algorithm 1 Summation

Input: Document cluster d
Output: Summary (H) and optimal policy

- 1: **procedure** SUMMATION
- 2: **Organiser:**
- 3: $Concepts\ and\ Relations \leftarrow$ Concept and relations extraction (d)
- 4: $H \leftarrow$ Hierarchical conceptual clustering ($Concepts\ and\ Relations$)
- 5: **Summarizer (User preference learner (iteratively)):**
- 6: $User\ preferences \leftarrow$ Query pairs (user)
- 7: $Ranker\ function \leftarrow$ Preference learner ($User\ preferences$)
- 8: **Summarizer (RL learner):**
- 9: $Optimal\ policy \leftarrow$ Policy learner ($Ranker\ function$)
- 10: **return** Summary (H) and optimal policy

We defined the reward as the summation of all concepts’ importance included in the summary. A policy π defines the strategy to add concepts to the draft summary to build a user’s desired summary. We defined π as the probability of choosing a summary of y among all possible summaries within the limit size using different hierarchy paths, $Y(d)$, denoted as $\pi(y)$. The expected reward of performing policy π , where $R(y)$ is the reward for selecting summary y , is defined in Eq. 7.

$$R^{RL}(\pi|d) = \mathbb{E}_{y \in Y(d)} R(y) = \sum_{y \in Y(d)} \pi(y) R(y) \quad (7)$$

The goal of MDP is to find the optimal policy π^* that has the highest expected reward. Therefore, the optimal policy, π^* , is the function that finds the desired summary for a given input based on user feedback (Eq. 8).

$$\pi^* = \operatorname{argmax} R^{RL}(\pi|d) = \operatorname{argmax} \sum_{y \in Y(d)} \pi(y) R(y) \quad (8)$$

We also used the linear temporal difference algorithm to obtain π^* . The process is explained in Algorithm 1.

IV. EVALUATION

In this section, we present the experimental setup for assessing our summarization model’s performance. We discuss the datasets, give implementation details, and explain how system output was evaluated.

A. Datasets and Evaluation

We evaluated *Summation* using three commonly employed benchmark datasets from the Document Understanding Conferences (DUC) ². Each dataset contains a set of document clusters accompanied by several human-generated summaries used for training and evaluation. Details are explained in Table I

²Produced by the National Institute Standards and Technology (<https://duc.nist.gov/>)

TABLE I: Dataset description: indicating the number of documents, number of document clusters, and the average number of sentences in each document.

Dataset	Doc-Num	Cluster-Num	Sentence
DUC1	30	308	378
DUC2	59	567	271
DUC4	50	500	265

Automatic Evaluation. We evaluate the quality of summaries using $ROUGE_N$ measure [11]³ defined as:

The three variants of ROUGE (ROUGE-1, ROUGE-2, and ROUGE-L) are used. We used the limited length ROUGE recall-only evaluation (75 words) to avoid being biased.

Human Evaluation. For this purpose, we hired fifteen Amazon Mechanical Turk (AMT)⁴ workers to attend tasks without any specific prior background required. Then five document clusters are randomly selected from the DUC datasets. Each evaluator was presented with three documents to avoid any subjects’ bias and was given two minutes to read each article. To make sure human subjects understood the study’s objective, we asked workers to complete a qualification task first. They were required to write a summary of their understanding. We manually removed spam from our results.

B. Results and Analysis

Summation was evaluated from different evaluation aspects, first from the organiser’s output, and then concerning the hierarchical concept map (H), which can be served individually to users as the structured summarised data. Next, we evaluated H using both human and automatic evaluation techniques to answer the following questions:

- Do users prefer hierarchical concept maps to explore new and complex topics?
- How much do users learn from a hierarchical concept map?
- How coherent is the produced hierarchical concept map?
- How informative are summaries in the form of a hierarchical concept map?

Personalized summaries generated on test data were also evaluated from various perspectives to analyse the effect of RL and preference learning, including:

- The impact of different features in approximating the proposed preference learning.
- The role of the query budget in retrieving pairwise preferences.
- The performance of RL algorithm and the information coverage in terms of ROUGE.
- Users’ perspectives on learned summaries based on their given feedback.

Hierarchical Concept Map Evaluation. To answer the questions in Sec. IV, we performed three experiments. First, within the same limit size as the reference summaries, we

³We run ROUGE 1.5.5: <http://www.berouge.com/Pages/default.aspx> with parameters -n 2 -m -u -c 95 -r 1000 -f A -p 0.5 -t 0

⁴<https://www.mturk.com/>

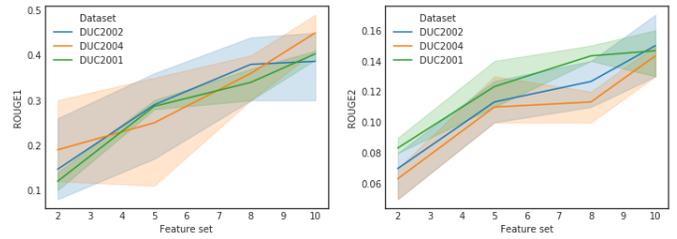


Fig. 3: Evaluating different feature set for estimating the ranker function.

compared the summaries produced by three models—using ExDos, which is a traditional approach; using a traditional hierarchical approach [22]; and using a structured summarization approach [23] on selected documents (with ROUGE-1 and ROUGE-2 scores based on the reference summaries). The average ROUGE-1 for Summation was 0.65 and ROUGE-2 was 0.48. The structured approach [23] showed similar performance with ROUGE-1 and ROUGE-2 at 0.65 and 0.45, respectively. Meanwhile, traditional hierarchical approaches [22] produced a ROUGE-1 of 0.27 and ROUGE-2 of 0.18. In the same task, the percentage of covered unigrams and bigrams based on documents were also compared. Both Summation and the structured approach covered approximately 4% unigrams and 2% bigrams, but dropped below 1% in both cases when testing the hierarchical approaches. In the third experiment, all competitors’ outputs were rated based on three measures, including usability in exploring new topics, level of informativeness, and coherency. Summation’s rate for the first and second criteria was 96% and 94%, respectively. However, it was 34% for coherency. We removed all concepts with low similarity to their parents based on a different threshold at each level. After repeating the same experiment, and rate of coherency increased to 76%.

Feature Analysis. Before evaluating the effect of conceptual preference, it is important to explain the ground-truth concept ranker function (U) and the approximate function (U^*), indicating the importance of concepts. To estimate the approximate function (U^*), we defined a linear model $U^*(c) = W^T \phi(c)$, where ϕ are the features. To this end, a set of features (whose importance was validated in ExDos) was used, including surface-level and linguistic-level features. Surface-level features include frequency-based features (TF-IDF, RIDF, gain and word co-occurrence), word-based features (upper-case words and signature words), similarity-based features (Word2Vec and Jaccard measure) and named entities. Linguistic features are generated using semantic graphs and include the average weights of connected edges, the merge status of concepts as a binary feature, the number of concepts merged with a concept, and the number of concepts connected to the concept. We defined different combinations of features with different sizes, $\{2, 5, 8, 10\}$, starting from the most critical one. Then, we repeated the experiments for 10 cluster documents. We used the concepts included in the reference summary as preferences, and then evaluated the concept coverage in

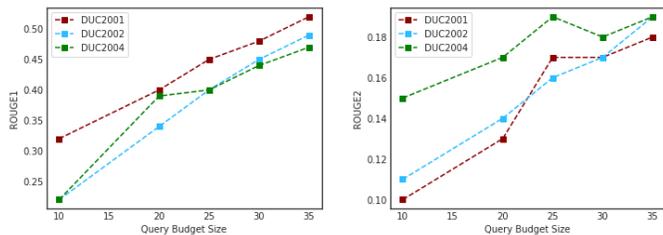


Fig. 4: The effect of query budget in ROUGE1 and ROUGE2 score.

TABLE II: Comparing *Summation* with benchmark datasets.

Model	ROUGE1	ROUGE2	ROUGEL
Traditional Structured [23]	0.346	0.090	0.251
Traditional Hierarchical [22]	0.211	0.013	0.149
<i>Summation</i>	0.731	0.651	0.681

a concept map compared to the reference summaries using ROUGE-1 and ROUGE-2. The results reported in Fig. 3 show that the model’s performance improved after adding more features.

Summary Evaluation. To avoid subjectivity in the evaluation process, we used the reference summaries as feedback. The mentioned concepts that exist in reference summaries receive the maximum score by the ranked function. We compared the summaries produced by three models, including the traditional approach (ExDos), a range of hierarchical approaches [22], and a structured summarization approach [23], each tested on randomly selected documents from three datasets using ROUGE-1, ROUGE-2 and ROUGE-L scores based on the references summaries. The average results reported in Table II show the supremacy of *Summation* in selecting specific contents.

Query Budget Size. We also measure the effectiveness of the users’ query budget size in the process. The pairwise preferences are defined based on the reference summaries, defining in a dictionary format. We selected the query size among the selection of $\{10, 15, 20, 25, 30, 35\}$, demonstrating the user’s number of feedback. The results are reported in Figure 4. As expected, by increasing the number of feedback, the ROUGE score increases significantly. However, the difference rate decreases through the process.

Human Analysis. Since the goal of *Summation* is to help users make their desired summary, we conducted two human experiments to evaluate the model. In the first experiments, to assess the possibility of finding their desired information, they were asked to answer a given question about each topic. Their level of confidence in answering questions and their answers were recorded. An evaluator assessed their accuracy in answering questions. Among the fifteen workers, 86.67% were completely confident in their answers. However, 57% answered completely accurately. In another task, after querying users for feedback, we ask them to select some concepts as the summary for the test data. Then the outputs were also shown to users, and they all approved their satisfaction. Besides, an

evaluator manually compared them and reported more than 80% correlation between outputs.

V. CONCLUSION AND FUTURE WORK

Extensive information in various formats is producing from single or multiple simultaneous sources in different systems and applications. For instance, data can be structured, such as data in SQL databases, unstructured stored in NoSQL systems, semi-structured like web server logs, or streaming data from a sensor. We propose a summarization approach based on a hierarchical concept map to tackle the variety and volume of big generated data. We trained our approach using document collections as input and employed users’ feedback to generate desired summaries for users, which can be extended to other data types. Many future directions are possible. First, capturing users’ interests is a significant challenge in providing practical personalized information. The reason is that users are reluctant to specify their preferences as entering lists of interests may be a tedious and time-consuming process. Therefore, techniques that extract implicit information about users’ preferences are the next step for making useful personalized summaries. Another potential direction is to use human feedback records to provide personalized summaries on new domains using transfer learning. Moreover, we aim to use fuzzy clustering to make a hierarchical concept map.

ACKNOWLEDGEMENT

We acknowledge the Centre for Applied Artificial Intelligence at Macquarie University, Sydney, Australia, for funding this research.

REFERENCES

- [1] F. Schiliro, A. Beheshti, S. Ghodratnama, F. Amouzgar, B. Benatallah, J. Yang, Q. Z. Sheng, F. Casati, and H. R. Motahari-Nezhad, “icop: Iot-enabled policing processes,” in *Service-Oriented Computing–ICSOC 2018 Workshops: ADMS, ASOCA, ISYyCC, CloTS, DDBS, and NLS4IoT, Hangzhou, China, November 12–15, 2018, Revised Selected Papers 16*. Springer, 2019, pp. 447–452.
- [2] F. Amouzgar, A. Beheshti, S. Ghodratnama, B. Benatallah, J. Yang, and Q. Z. Sheng, “isheets: A spreadsheet-based machine learning development platform for data-driven process analytics,” in *Service-Oriented Computing–ICSOC 2018 Workshops: ADMS, ASOCA, ISYyCC, CloTS, DDBS, and NLS4IoT, Hangzhou, China, November 12–15, 2018, Revised Selected Papers 16*. Springer, 2019, pp. 453–457.
- [3] A. Beheshti, F. Schiliro, S. Ghodratnama, F. Amouzgar, B. Benatallah, J. Yang, Q. Z. Sheng, F. Casati, and H. R. Motahari-Nezhad, “iprocess: Enabling iot platforms in data-driven knowledge-intensive processes,” in *Business Process Management Forum: BPM Forum 2018, Sydney, NSW, Australia, September 9-14, 2018, Proceedings 16*. Springer, 2018, pp. 108–126.
- [4] S. Ghodratnama, A. Beheshti, M. Zakershahrak, and F. Sobhanmanesh, “Intelligent narrative summaries: from indicative to informative summarization,” *Big Data Research*, vol. 26, p. 100257, 2021.
- [5] U. Khanna, S. Ghodratnama, A. Beheshti *et al.*, “Transformer-based models for long document summarisation in financial domain,” in *Proceedings of the 4th Financial Narrative Processing Workshop@ LREC2022*, 2022, pp. 73–78.
- [6] A. Beheshti, S. Ghodratnama, M. Elahi, and H. Farhood, *Social Data Analytics*. CRC Press, 2022.
- [7] V. Gupta and G. S. Lehal, “A survey of text summarization extractive techniques,” *Journal of emerging technologies in web intelligence*, vol. 2, no. 3, pp. 258–268, 2010.

- [8] A. Beheshti, B. Benatallah, H. R. Motahari-Nezhad, S. Ghodratnama, and F. Amouzgar, "Bp-sparql: A query language for summarizing and analyzing big process data," in *Process Querying Methods*. Springer, 2021, pp. 21–48.
- [9] S. Ghodratnama, "Towards personalized and human-in-the-loop document summarization," *arXiv preprint arXiv:2108.09443*, 2021.
- [10] S. Ghodratnama, M. Zakershaharak, and A. Beheshti, "Summary2vec: learning semantic representation of summaries for healthcare analytics," in *2021 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2021, pp. 1–8.
- [11] C.-Y. Lin, "Rouge: A package for automatic evaluation of summaries," *Text Summarization Branches Out*, 2004.
- [12] S. Ghodratnama, M. Zakershaharak, and F. Sobhanmanesh, "Am i rare? an intelligent summarization approach for identifying hidden anomalies," in *Service-Oriented Computing–ICSOC 2020 Workshops: AIOps, CFTIC, STRAPS, AI-PA, AI-IOTS, and Satellite Events, Dubai, United Arab Emirates, December 14–17, 2020, Proceedings*. Springer, 2021, pp. 309–323.
- [13] P. Mehta and P. Majumder, "Effective aggregation of various summarization techniques," *Information Processing & Management*, vol. 54, no. 2, pp. 145–158, 2018.
- [14] H. P. Edmundson, "New methods in automatic extracting," *Journal of the Association for Computing Machinery*, vol. 16, no. 2, pp. 264–285, 1969.
- [15] S. Ghodratnama, A. Beheshti, M. Zakershaharak, and F. Sobhanmanesh, "Extractive document summarization based on dynamic feature space mapping," *IEEE Access*, vol. 8, pp. 139 084–139 095, 2020.
- [16] Y. Wu and B. Hu, "Learning to extract coherent summary via deep reinforcement learning," in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [17] S. Narayan, S. B. Cohen, and M. Lapata, "Ranking sentences for extractive summarization with reinforcement learning," in *Association for Computational Linguistics*, 2018.
- [18] X. Liu, Z. Nie, N. Yu, and J.-R. Wen, "Biosnowball: automated population of wikis," in *ACM SIGKDD international conference on Knowledge discovery and data mining*, 2010.
- [19] S. Ghodratnama and R. Boostani, "An efficient strategy to handle complex datasets having multimodal distribution," in *ISCS 2014: Interdisciplinary Symposium on Complex Systems*. Springer, 2015, pp. 153–163.
- [20] R. Nallapati, A. Feng, F. Peng, and J. Allan, "Event threading within news topics," in *The thirteenth ACM international conference on Information and knowledge management*, 2004.
- [21] A. Haghighi and L. Vanderwende, "Exploring content models for multi-document summarization," in *Annual Conference of the Association for Computational Linguistics*, 2009.
- [22] J. Christensen, S. Soderland, G. Bansal *et al.*, "Hierarchical summarization: Scaling up multi-document summarization," in *ACL*, 2014, pp. 902–912.
- [23] T. Falke, "Automatic structured text summarization with concept maps," Ph.D. dissertation, Technische Universität, 2019.
- [24] A. Borisov, M. Wardenaar, I. Markov, and M. de Rijke, "A click sequence model for web search," in *ACM SIGIR Conference on Research & Development in IR*, 2018, pp. 45–54.
- [25] J. Kreutzer, S. Khadivi, E. Matusov, and S. Riezler, "Can neural machine translation be improved with user feedback?" *arXiv preprint arXiv:1804.05958*, 2018.
- [26] C. Lawrence and S. Riezler, "Counterfactual learning from human proofreading feedback for semantic parsing," *arXiv preprint arXiv:1811.12239*, 2018.
- [27] "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [28] C. Orasan and L. Hasler, "Computer-aided summarisation-what the user really wants," in *LREC*, 2006, pp. 1548–1551.
- [29] M. Narita, K. Kurokawa, and T. Utsuro, "A web-bojanowski2017enrichingased english abstract writing tool using a tagged ej parallel corpus," in *LREC*, 2002.
- [30] A. Leuski, C.-Y. Lin, and E. Hovy, "ineats: interactive multi-document summarization," in *41th Annual Meeting of the Association for Computational Linguistics*, 2003, pp. 125–128.
- [31] S. Jones, S. Lundy, and G. W. Paynter, "Interactive document summarisation using automatically extracted keyphrases," in *International Conference on System Sciences*. IEEE, 2002, pp. 1160–1169.
- [32] P. Avinesh and C. M. Meyer, "Joint optimization of user-desired content in multi-document summaries by learning from user feedback," in *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2017, pp. 1353–1363.
- [33] O. Etzioni, M. Banko, S. Soderland, and D. S. Weld, "Open information extraction from the web," *Communications of the ACM*, vol. 51, no. 12, pp. 68–74, 2008.
- [34] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," *NIPS*, vol. 26, pp. 3111–3119, 2013.
- [35] P. Bojanowski, E. Grave, A. Joulin, and T. Mikolov, "Enriching word vectors with subword information," *Transactions of the Association for Computational Linguistics*, vol. 5, pp. 135–146, 2017.
- [36] R. Tibshirani, G. Walther, and T. Hastie, "Estimating the number of clusters in a data set via the gap statistic," *Journal of the Royal Statistical Society*, vol. 63, no. 2, pp. 411–423, 2001.
- [37] A. Alhindi, U. Kruschwitz, C. Fox, and M.-D. Albakour, "Profile-based summarisation for web site navigation," *ACM Transactions on Information Systems (TOIS)*, vol. 33, no. 1, pp. 1–39, 2015.
- [38] M. Zopf, "Estimating summary quality with pairwise preferences," in *North American Chapter of the Association for Computational Linguistics, Volume 1*, 2018, pp. 1687–1696.
- [39] J. Fürnkranz and E. Hüllermeier, "Preference learning and ranking by pairwise comparison," in *Preference learning*. Springer, 2010, pp. 65–82.
- [40] M. Szummer and E. Yilmaz, "Semi-supervised learning to rank with preference regularization," in *International conference on Information and knowledge management*, 2011, pp. 269–278.
- [41] M. Zakershaharak and S. Ghodratnama, "Are we on the same page? hierarchical explanation generation for planning tasks in human-robot teaming using reinforcement learning," *arXiv preprint arXiv:2012.11792*, 2020.