

Controllable Text Summarization: Unraveling Challenges, Approaches, and Prospects - A Survey

Anonymous ACL submission

Abstract

Generic text summarization approaches often fail to address the specific intent and needs of individual users. Recently, scholarly attention has turned to the development of summarization methods that are more closely tailored and controlled to align with specific objectives and user needs. Despite a growing corpus of controllable summarization research, there is no comprehensive survey available that thoroughly explores the diverse controllable attributes employed in this context, delves into the associated challenges, and investigates the existing solutions. In this survey, we formalize the Controllable Text Summarization (CTS) task, categorize controllable attributes according to their shared characteristics and objectives, and present a thorough examination of existing datasets and methods within each category. Moreover, based on our findings, we uncover limitations and research gaps, while also exploring potential solutions and future directions for CTS.

1 Introduction

Despite the significant advancements in automatic text summarization, its one-size-fits-all approach falls short in meeting the varied needs of different segments of users and application scenarios. For example, generic automatic summarization may struggle to produce easily understandable summaries of scientific documents for non-expert users or create extremely brief summaries of news stories for online feeds. Lately, a myriad of works have emerged aimed at generating more controlled (Fan et al., 2018a; Maddela et al., 2022; He et al., 2022; Zhang et al., 2023b; Pagnoni et al., 2023) and tailored text summaries that meet a wide range of user needs.

CTS task is centered around creating summaries of source documents that adhere to specific criteria. These criteria are managed through various controllable attributes (CA) or aspects like summary length (Kwon et al., 2023), writing style (Goyal

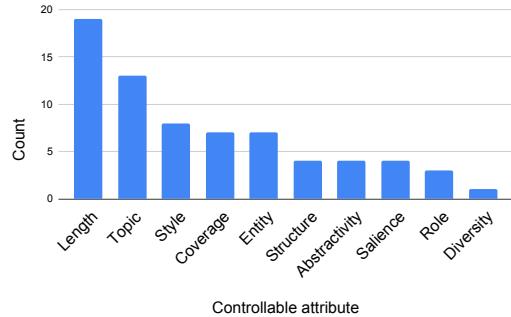


Figure 1: Number of controllable text summarization publications for various attributes.

et al., 2022), coverage of key information (Li et al., 2018; Jin et al., 2020b), content diversity (Narayan et al., 2022), and more. These criteria vary depending on the task, user needs, and specific application context. For example, length-controlled summaries (Hitomi et al., 2019) are particularly useful in situations where brevity is crucial, like in social media posts, headlines, and abstracts. In areas such as marketing, academic writing, or professional communication, a style-controlled summary (Chawla et al., 2019) is essential to ensure that the information aligns with the intended tone and messaging strategy. Similarly, topic-controlled summaries (Bahrainian et al., 2021) are commonly used in research papers, reports, and content curation, providing an emphasis on a specific topic to enhance clarity and coherence in the presented information.

There is an uneven distribution of attention within the research community towards various CAs as depicted in Fig 1. The majority of CTS works concentrate on managing length, topic, and style. This could be attributed to two main factors. First, it is comparatively simpler to develop datasets for evaluating length, topic, and style compared to aspects like structure and diversity. Second, there is a plethora of application scenarios for length or topic-oriented summaries, such as generating concise news feeds or focused legal reports.

In this survey, we collect and analyze 61 research

001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022

023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041

042
043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070

Source: (CNN)Novak **Djokovic** extended his current winning streak to 17 matches after beating Thomas Berdych 7-5, 4-6, 6-3 in the rain-interrupted **final** of the Monte Carlo Masters....After winning the Australian Open back in January, Djokovic has followed up with Masters' victories at Indian Wells and Miami. He then beat Rafa Nadal, arguably one of the greatest players on clay of all time...

Length:Long, **Coverage:**High, **Topic:**Djokovic,Final
Summary: Djokovic wins 7-5, 4-6, 6-3 after a tight match with Berdych in the Monte Carlo Masters final. Djokovic also followed up with Masters' victories at Indian Wells and Miami.

Length:Normal, **Topic:**Djokovic, **Coverage:**Normal,
Summary: It's been a sensational year for Djokovic after beating Berdych in the finals and also winning against clay expert Nadal.

Length:Short, **Coverage:**Normal, **Topic:** No Control
Summary: Djokovic wins Monte Carlo Masters after beating Berdych 7-5, 4-6, 6-3 in the finals.

Table 1: Summaries obtained by varying Controllable Attributes from MACSUM (Zhang et al., 2023b)

papers pertaining to various possible CAs. The filtration criteria for the selection of papers are described in Appendix B. Subsequently, we classify these CAs into 10 categories, grouping similar ones based on shared characteristics and objectives. Moreover, we delve into the existing datasets, evaluating their creation methods and appropriateness for the respective task in each CA category. Furthermore, we scrutinize the current CTS methodologies for each CA category, drawing comparisons between their overarching frameworks and discussing relevant limitations. Subsequently, we discuss in detail the generic and specific evaluation strategies for CAs utilized by various works. Finally, we attempt to critique the current approaches and unravel potential future research trajectories. To the best of our knowledge, it is the first comprehensive survey on CTS.

2 Task Formulation

This section introduces the Controllable Text Summarization (CTS) task by outlining its definition and offering a categorized breakdown of the different controllable attributes along with concise descriptions for each. Given a set of source documents $D = \{d_1, d_2, \dots, d_k\}$. Each document, d_i , consists of a sequence of n tokens: $\{x_{i,1}, x_{i,2}, \dots, x_{i,n}\}$. S_i is the target summary of document d_i , which comprises of a sequence of m tokens: $\{s_{i,1}, s_{i,2}, \dots, s_{i,m}\}$, where $m \ll n$. The user wants to control a set of controllable attributes C . The task can be framed as a conditional genera-

Attribute	Definition
Length	Controlling the length of the summary
Style	Controlling the readability levels, politeness, humor, and emotion
Coverage	Controlling the salient information in summary
Entity	Summary specific to pre-defined entities
Structure	Create summaries with predefined structure or order
Abstractivity	Controlling the novelty in sentence formation
Salience	Adjusting the presence of prominent information
Role	Providing role-specific summaries
Diversity	Generating semantically diverse summaries
Topic	Controlling topic-focused summary generation

Table 2: Controllable attributes definitions.

tive problem: $P(S|D, C) = \prod_i^k P(S_i|d_i, C)$

2.1 Controllable Attributes

A Controllable Attribute or Aspect (CA) refers to a user or application-driven trait of the summary designed to meet specific criteria or conditions, such as Length, Style, Role, etc. In the literature, it is evident that various authors use different terms to describe the same CAs, which exhibit similar characteristics and objectives as shown in Appendix A Table 3. Based on these observations, we group the CAs into 10 categories as listed in Table 2.

3 Related Surveys

In the literature, a multitude of surveys center around conventional text summarization methods (El-Kassas et al., 2021; Nazari and Mahdavi, 2019; Allahyari et al., 2017; Gambhir and Gupta, 2017), including task-specific surveys such as multi-document summarization (Sekine and Nobata, 2003), cross-lingual summarization (Wang et al., 2022), and dialogue-based summarization (Tugener et al., 2021). There are few surveys that concentrate on text generation techniques (Zhang et al., 2023a; Prabhumoye et al., 2020) and the causal perspective (Wang et al., 2024; Hu and Li, 2021) on the same. On the contrary, this is the first survey, that focuses on controllable summarization by offering a thorough analysis of CTS methods, challenges, and prospects.

4 Datasets

This section provides a broad overview of the CTS datasets and corresponding creation/acquisition strategies. The CTS methods are evaluated in several ways: 1) by utilizing publicly available summarization datasets, 2) by datasets derived from generic datasets, and 3) by creating human-annotated datasets.

4.1 Generic Datasets

CTS research predominantly leverages widely used news summarization datasets. Notably, about 57% of CTS studies utilize either CNN-DailyMail (Nalapati et al., 2016) or DUC (Over and Yen, 2004; Dang, 2005). Other popular datasets, including Gigaword (Napoles et al., 2012), XSum (Narayan et al., 2018), NYTimes (Sandhaus, 2008), NEWSROOM (Narayan et al., 2018), and dialogue-based SAMSUM (Gliwa et al., 2019), along with opinion-based datasets (Angelidis and Lapata, 2018; Angelidis et al., 2021), are employed for controllable summarization. However, these generic datasets lack explicit annotations and nuances to evaluate the CA-specific summarization. CTS requires specialized datasets (as detailed in Table 6) to provide evaluation opportunities for specific aspects like length, topic, style, etc.

4.2 Derived Datasets

The derived datasets are obtained by applying the aspect-specific heuristics to the widely used generic datasets. In this section, we list out few derived datasets and their creation strategies.

JAMUL. Hitomi et al. (2019) collect the length-sensitive headlines for the Japanese language. Each article consists of three headlines with varying lengths of 10, 13, and 26 characters respectively. **TS and PLS.** In order to enhance the readability of biomedical documents, Luo et al. (2022) introduce two types of summaries. The *Technical Summary* is an abstract of a peer-reviewed bio-medical research paper and the *Plain Language Summary (PLS)* is the authors submitted summary as part of the journal submission process. **Wikiasp.** In order to construct the multi-domain aspect-based summarization corpus, Hayashi et al. (2021) utilize the Wikipedia articles from 20 domains. Further, the section titles and paragraph boundaries of each article are obtained as a proxy of aspect annotation. In another study, Ahuja et al. (2022) create the **ASPECTNEWS** dataset for aspect-oriented summarization. They achieve it by utilizing articles from the CNN/DailyMail dataset and identifying documents related to ‘earthquakes’ and ‘fraud investigations’ by using the universal sentence encoder (Cer et al., 2018). Further, Mukherjee et al. (2020) collect a CA-based opinion summarization dataset consisting of tourism reviews. These are obtained from the *TripAdvisor* website and identified the relevant aspects using the unsupervised attention-

based aspect extraction technique (He et al., 2017).

4.3 Human annotated

This section provides the details of the human-annotated CTS datasets.

GranDUC. By re-annotating the DUC-2004 (Dang, 2005), Zhong et al. (2022) release a novel benchmark dataset for the granularity control. Annotators are instructed to create summaries of multiple documents with *coarse*, *medium*, and *fine* granularity levels. **Multi-LexSum.** Shen et al. (2022c) create a human-annotated corpus of 9,280 civil rights lawsuits and corresponding summaries with different degrees of granularity. The target summary length ranges from one-sentence to multi-paragraph level. **EntSUM.** (Maddela et al., 2022) is a human-annotated entity-specific controllable summarization dataset. It utilizes the articles from The New York Times Annotated Corpus (NYT) (Sandhaus, 2008) and includes annotated summaries for PERSON and ORGANISATION tags. The recent release of EntSUMV2 (Mehra et al., 2023) is the more abstractive version of EntSUM.

NEWTS. Bahrainian et al. (2022) introduce the topically focused summarization corpus by leveraging documents from CNN-DailyMail and employing crowd-sourcing to generate two distinct summaries with different thematic aspects for each document.

CSDS. Lin et al. (2021) introduces the role-oriented Chinese Customer Service Dialogue Summarization (CSDS) dataset. It is meticulously annotated, segmenting the dialogues based on their topics and summarizing each segment as a QA pair.

MReD. To tackle the task of structure-controllable summarization, (Shen et al., 2022b) introduce the Meta-Review Dataset (MReD). It is created by gathering meta-reviews from the open review system and categorizing each sentence into one of nine predefined intent categories (abstract, strength, weakness, etc.). **MACSUM.** Zhang et al. (2023b) develop a human-annotated corpus to control the mix of CAs (Topic, Speaker, Length, Extractiveness, and Specificity) together. MACSUM covers source articles from CNN/DailyMail and QMSUM (Zhong et al., 2021) datasets.

5 Approaches to Controlled Summarization

Various CAs have been investigated in controllable summary generation tasks, including style (politeness, humor, formality), content (length, entities,

237 keywords), and structure. In this section, we de-
238 scribe various approaches to achieve CTS for the
239 attributes mentioned in Table 2. Additionally, we
240 list out the novel contributions and limitations for
241 each paper in the Appendix D Table 9.

242 **Length.** Earlier methods lacked length control and
243 only employed heuristics such as stopping the gen-
244 eration after a fixed number of tokens. To overcome
245 this, four different approaches to integrate length
246 as a learnable parameter are proposed.

247 **Adding length in input:** Fan et al. (2018b) propose
248 a convolutional encoder-decoder-based summariza-
249 tion system, where it quantizes summary lengths
250 into discrete bins of different size ranges. During
251 training, the input data is prepended with the
252 gold summary length represented by bin lengths.
253 Due to a fixed number of length bins, the system
254 *fails* to generate summaries of arbitrary lengths.
255 CTRLSUM (He et al., 2022) presents a generic
256 framework to generate controlled summaries using
257 keywords specific to length. Instead of controlling
258 a single attribute, Zhang et al. (2023b) allow dif-
259 ferent length attribute values (normal, short, long)
260 to be used as inputs along with the source text for
261 hard prompt tuning (Brown et al., 2020).

262 **Adding length in encoder:** Yu et al. (2021) pro-
263 pose a length context vector that is generated at
264 each decoding step derived from the positional en-
265 codings. This vector is then concatenated with the
266 decoder hidden state and encoder attention vectors.
267 The *limitation* of the system is the generation of
268 incomplete summaries for short desired lengths.
269 Liu et al. (2022b) propose a length-aware attention
270 model that adapts the source encodings based on
271 the desired length by pretraining the model. Zhang
272 et al. (2023b) add a hyperparameter for learning the
273 prefix embeddings for different attributes at each
274 layer of the encoder and decoder for soft prefix
275 tuning (Li and Liang, 2021).

276 **Adding length in decoder:** Kikuchi et al. (2016)
277 propose the first method to control length using a
278 BiLSTM encoder-decoder architecture with atten-
279 tion (Luong et al., 2015) for sentence compres-
280 sion. In each step of the decoding process, an additional
281 input for the remaining length is provided as an
282 embedding. Instead of pre-defined length ranges,
283 Liu et al. (2018) add a desired length parameter
284 at the decoding step to each convolutional block
285 of the initial layer of the convolutional encoder-
286 decoder model. Févry and Phang (2018) design an
287 unsupervised denoising auto-encoder for sentence

288 compression, where the decoder has an additional
289 input of the remaining summary length at each time
290 step. While it produces grammatically correct sum-
291 maries but nonsensical or semantically different
292 from the input. Which leads to the generation of
293 *unfaithful* summaries.

294 To handle the length constraint, Takase and
295 Okazaki (2019) propose two modifications to the
296 sinusoidal positional embeddings on the decoder
297 side: length-difference positional encoding and
298 length-ratio positional encoding. Sarkhel et al.
299 (2020) present a multi-level summarizer that mod-
300 els a multi-headed attention mechanism using a
301 series of interpretable semantic kernels to control
302 lengths, reducing the trainable parameters signif-
303 icantly. The model does *not encode* the length
304 attribute directly. Song et al. (2021) design a
305 confidence-driven generator that is trained on a
306 denoising objective with a decoder-only architec-
307 ture, where the source and summary tokens are
308 masked with position-aware beam search. Goyal
309 et al. (2022) use a mixture-of-experts model with
310 multiple transformer-based decoders for identify-
311 ing different styles or features of summaries. Kwon
312 et al. (2023) introduce the summary length predic-
313 tion task on the encoder side and this predicted
314 summary length is inserted with a length-fusion
315 positional encoding layer.

316 **Adding length in loss/reward function:** Makino
317 et al. (2019) propose a global minimum risk train-
318 ing optimization method under length constraint for
319 the neural summarization tasks which is faster and
320 generates five times fewer over-length summaries
321 on an average than others. Chan et al. (2021) use an
322 RL-based Constrained Markov Decision process
323 with a mix of attributes. Hyun et al. (2022) de-
324 vice an RL-based framework that incorporates both
325 length and quality constraints in the reward func-
326 tion to generate multiple summaries of different
327 lengths and according to the experimental results
328 present in Hyun et al. (2022), the model is computa-
329 tionally *expensive*.

330 **Style.** The generation of user-specific summaries
331 has gained significant interest, but achieving dis-
332 tinct styles has posed an enduring challenge. These
333 stylistic variations may encompass tone, readability
334 control, or the modulation of user emotions. Style
335 control aims to generate source-specific summaries
336 (Fan et al., 2018a) by utilizing the convolutional
337 encoder-decoder network.

338 Chawla et al. (2019) obtain formality-tailored sum-

339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
maries by utilizing the input-dependent reward function. The pointer-generator (See et al., 2017) network is used as the under-laying architecture and the loss function is modified with the addition of a formality-based-reward function. In another study, Jin et al. (2020a) attempt to control humor, romance, and clickbait in headlines using a multitask learning framework. By employing an inference style classifier, Cao and Wang (2021) adjust the decoder final states to obtain stylistic summaries. Moreover, obtain the lexical control by utilizing the word unit prediction that can directly constrain the output vocabulary. Similarly, Goyal et al. (2022) extend the decoder architecture to a mixture-of-experts version by using multiple decoders. The gating mechanism helps to obtain multiple summaries for a single source. However, the major limitation in this model is its *manual* gating mechanism. To control various fine-grained reading grade levels, Ribeiro et al. (2023) present three methods: instruction-prompting, reinforcement learning-based reward model, and look-ahead readability decoding approach.

Coverage. Managing the information granularity is essential to measure the semantic coverage between the source text and the summary. To regulate the granularity, Wu et al. (2021) introduce a two-stage approach, where the model incorporates a summary sketch, that encompasses user intentions and key phrases, serving as a form of weak supervision. They leverage a text-span-based conditional generation to govern the level of detail in the generated dialogue summaries. Zhong et al. (2022) proposes a multi-granular event-aware summarization method composed of four stages: event identification, unsupervised event-based summarizer pretraining, event ranking, and summary generation by adding events as hints. Extraction of events from source text may *lower* the abstractiveness. Zhang et al. (2023b) use the hard and soft-prompting strategies to control the amount of extracted text from the source in the summary. Additionally, Huang et al. (2023) utilize the natural language inference models to improve the coverage.

Entity. Entity-centric summarization concentrates on producing a summary of a document that is specific to a given target entity (Hofmann-Coyle et al., 2022). Zheng et al. (2020) extract the named entities using a pre-trained BERT (Devlin et al., 2019) based model and feed both the article and the selected entities to a bidirectional LSTM (Hochreiter

and Schmidhuber, 1997) encoder-decoder model. In another study, Liu and Chen (2021) extract the entities (speakers and non-speaker entities) from a dialogue to form a planning sequence. The entities extracted are concatenated to the source dialogue for training the conditional BART-based model. This model introduces factual *inconsistency* due to paraphrasing from a personal perspective. Maddela et al. (2022) extend the GSum (Dou et al., 2021) by feeding it either sentences or strings, which mention extracted entities as guidance. The model is an adapted version of BERTSum (Liu and Lapata, 2019), where all the sentences containing the entity string mention and its coreferent mentions are only fed. Hofmann-Coyle et al. (2022) model entity-centric extractive summarization as a sentence selection task. Building upon BERTSum (Liu and Lapata, 2019), they use a BERT (Devlin et al., 2019) based encoders to represent the sentence and target entity pair and train with a contrastive loss objective to extract sentences most relevant to the target entities.

Structure. Generic datasets lack key elements for emphasizing specific aspects in the corresponding ground truth summaries. To address this limitation and emphasize summary structure, Shen et al. (2022b) achieve structure-controllable text generation by adding a control sequence at the beginning of the input text and treating summary generation as a standalone process. However, this approach has two main *limitations*, 1) generated tokens are solely based on logits predictions without ensuring that the sequence satisfies the control signal, 2) Auto-regressive models face error propagation in generation due to self-attention, causing subsequent generations to deviate from the desired output. To overcome these challenges, the sentence beam-search (SentBS) (Shen et al., 2022a) approach produces multiple sentence options for each sentence and selects the best sentence based on both the control structure and the model’s likelihood probabilities. In a related study, Zhong and Litman (2023) utilize predicted argument role information to control the structure in legal opinion documents. Additionally, in the work of Zhang et al. (2023b), the prompt of entity chains, representing an ordered sequence of entities, is used for pre-training and fine-tuning with a planning objective to control the summary structure.

Abstractivity. It measures the degree of textual novelty between the source text and summary. See

441 et al. (2017) introduce a pointer-generator network
442 to control the source copying via *pointing* and gen-
443 erate novel sentence formations by using *generator*
444 mechanism. However, this scheme *fails* to generate
445 higher abstraction levels. Kryściński et al. (2018)
446 tackle this problem in two ways: 1) decompose
447 the decoder into a contextual network to retrieve
448 the relevant parts of the text and generate the sum-
449 mary by utilizing a pretrained model, 2) a mixed
450 RL-based objective jointly optimizes the n-gram
451 overlap with the ground truth summary. Similarly,
452 Song et al. (2020) control the copying behavior by
453 using a *mix-and-match* strategy to generate sum-
454 maries with varying n-gram copy rates. Based on
455 the *seen*, *unseen* words from the source text, the
456 system controls the copying percentage by acting
457 as a language modeling task. Moreover, methods
458 such as ControlSum (Fan et al., 2018a) allow the
459 users to explicitly specify the control attribute to fa-
460 cilitate better control. However, it does not provide
461 any supervision on *violating* the controllability. To
462 alleviate this issue, Chan et al. (2021) propose an
463 RL-based framework on the constrained Markov
464 decision process and introduced a reward to penalize
465 the violation of attribute requirement.

466 **Salience.** This attribute captures the most impor-
467 tant information in a document. In SummaRuN-
468 Ner (Nallapati et al., 2017), salience is modeled
469 as a feature in a classification objective. It uses
470 GRU-based encoders and decoders to frame sum-
471 marization as a text-to-binary sequence learning
472 task at the sentence level (Bahdanau et al., 2014)
473 (Cho et al., 2014). A binary score is assigned to
474 each sentence, indicating its membership in the
475 summary. The system performs *poorly* on out-of-
476 domain datasets. To retain key content from the
477 source, Li et al. (2018) introduce a Key Information
478 Guide Network, where keywords are identified by
479 the TextRank algorithm with a modified attention
480 mechanism that accommodates this key informa-
481 tion as an additional input. However, it focuses
482 mostly on informativeness *ignoring* coherence and
483 readability features.

484 Deutsch and Roth (2023) model salience in terms
485 of noun phrases using QA signals where the gen-
486 eration of the summary is conditioned on these
487 identified phrases. This approach is *not applicable*
488 to languages for which question generation and
489 question answering models are not available. In
490 long document CLS tasks, summarization systems
491 often *fail* to respond to user queries. To resolve this

492 issue, Pagnoni et al. (2023) propose a pre-training
493 approach that involves two tasks of salient infor-
494 mation identification from sentences having the
495 highest self ROUGE score and a question genera-
496 tion system to generate questions whose answers
497 are the salient sentences.

498 **Role.** Role-oriented dialogue summarization gen-
499 erates summaries for different roles/agents present
500 in a dialogue (e.g. doctor and patient) (Liang et al.,
501 2022). Lin et al. (2021) propose the CSDS dataset
502 (see Section 4.3) and benchmarked a variety of
503 existing state-of-the-art summarization models for
504 the task of generating agent and user surveys. They
505 find that agent summaries generated by the exist-
506 ing methods *lack* key information needed to be ex-
507 tracted from dialogues of the other role. To bridge
508 this gap, Lin et al. (2022) build a role-aware sum-
509 marization model for two users (agent and user)
510 present in the dataset. They use two separate de-
511 coders for generating the user and agent summaries
512 by utilizing user and agent masks. A role atten-
513 tion mechanism is introduced to each decoder so
514 that it can leverage the overall context by attend-
515 ing to the hidden states of the other role. Liang
516 et al. (2022) use a role-aware centrality scoring
517 model that computes role-aware centrality scores
518 for each utterance, which measures the relevance
519 between the utterance and the role prompts (signaling
520 whether the summary is for the user or agent).
521 This is then used to reweight the attention scores
522 for each utterance, which is then used by the de-
523 coder to generate the summary.

524 **Diversity.** Traditional decoding strategies, like
525 beam search, excel at generating single summaries
526 but often *struggle* to produce diverse ones. Tech-
527 niques such as top-k and nucleus sampling are ef-
528 fective in generating diverse outputs but may sac-
529 rifice faithfulness. In response to these challenges,
530 Narayan et al. (2022) introduce compositional sam-
531 pling, a decoding method to obtain diverse sum-
532 maries. This method initiates by planning a seman-
533 tic composition (Narayan et al., 2021) of the target
534 in the form of entity chains, and then leverages
535 beam search to generate diverse summaries.

536 **Topic.** Long documents often cover multiple topics,
537 and a generic summary might not fully encompass
538 the diverse scope. Krishna and Srinivasan (2018)
539 train a topic-conditioned pointer-generator network
540 (See et al., 2017) by concatenating one hot encod-
541 ing representation of the topic with the embedding
542 of each token in the input document. However,

news categories are used as the predefined topics, that *limits* the generalization to other tasks. To handle diverse topics, Tan et al. (2020) utilize external knowledge sources like Wikipedia and ConceptNet to create a weakly supervised summarization framework compatible with any encoder-decoder architecture. Suhara et al. (2020) propose an unsupervised method, where aspect-specific opinions are extracted from a set of reviews by a pre-trained opinion extractor, and the summary of the opinion is generated by a generator model trained to reconstruct the reviews from the opinions. Similarly, given a set of reviews for a product (e.g. Hotels), Amplayo et al. (2021) train a Multiple Instance Learning (MIL) model, to extract the predictions for aspect (like cleanliness) codes at the document, sentence, token level (Mukherjee et al., 2020). These predicted aspects transform the input such that relevant sentences and keywords along with aspect tokens are fed into the pre-trained T5 (Raffel et al., 2020) model.

Hsu and Tan (2021) introduce the task of generating decision-supportive summaries. The focus is on predicting future Yelp ratings from the set of reviews using a Longformer-based (Beltagy et al., 2020) regression model. They propose an iterative algorithm that selects the sentences of the summary from a set of representative sentences. Mukherjee et al. (2022) extend topic-focused summarization for multimodal documents by creating a joint image-text context vector.

6 Evaluation Strategies

This section catalogs and briefly describes the variety of automatic and human evaluation metrics that are being used to evaluate the summaries generated by the different methods studied in this paper.

6.1 Automatic Evaluation

The automatic evaluation metrics can be categorized based on how they are defined. We categorize the metrics into n-gram-based, language-model-based, and aspect-specific.

N-gram based evaluation metrics like ROUGE (Lin, 2004), BLEU (Papineni et al., 2002) are based on matching n-grams from candidate summaries to a set of reference summaries. ROUGE is the most widely used metric in CTS literature. **Language-model based** metrics are computed using pre-trained language models (PLM) like BERT (Devlin et al., 2019) or BART (Lewis

et al., 2019). One class of approach computes the distance between the PLM embeddings of the reference and the generated summary Another way is based on computing the log probability of the generated text conditioned on input text as demonstrated in BARTScore (Yuan et al., 2021). **Summarization specific** metrics including ROUGE-WE (Ng and Abrecht, 2015), S³ (Peyrard et al., 2017), Sentence Mover’s Similarity (SMS) (Clark et al., 2019), SummQA (Scialom et al., 2019), BLANC (Vasilyev et al., 2020), and SUPERT (Gao et al., 2020), (*Lite*)³Pyramid (Zhang and Bansal, 2021) are prominent for controllable summary evaluation. **Aspect specific** metrics do not fall cleanly into either of the above-mentioned categories. These metrics focused on evaluating specific controllable aspects such as Flesh Reading Ease (Flesch, 1948), Gunning Fog Index, and Coleman Liau Index for readability, control correlation, and error rate (Zhang et al., 2023b) for topic, abstractivity and role attributes. Appendix C Table 7 describes more details about the automatic evaluation metrics.

6.2 Human Evaluation

Human evaluation is an indicator of the robustness and effectiveness of different summarization systems on specific aspects that cannot be directly captured by automatic evaluation metrics. These aspects include generic properties of a summary such as truthfulness (Song et al., 2020; Hyun et al., 2022), relevance (Goyal et al., 2022; He et al., 2022; Shen et al., 2022b), fluency (Narayan et al., 2022; Suhara et al., 2020), and readability (Cao and Wang, 2021; Kryściński et al., 2018) or specific properties such as completeness (Yu et al., 2021; Liu et al., 2022a) for length-controlled summaries, coverage (Mukherjee et al., 2020, 2022) for the entity, and topic-controlled summary generation. Broadly two kinds of scoring mechanisms are used for human evaluation: binary and rank-based. The rank-based scores usually range from 1 to 5. Despite these widely adapted mechanisms, human evaluation of summarization is challenging due to ambiguity and subjectivity. Aspects like coherence and fluency help mitigate ambiguity, but remain subjective to individual annotators. Accurately defining annotation descriptions is crucial, yet achieving a standardized approach across annotators remains difficult (Iskender et al., 2021; Ito et al., 2023). The details about different human evaluation metrics are detailed in Appendix C Table 8.

7 Challenges and Future Prospects

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
692

Generic vs specialized benchmarks. We observe that more than 75% of CTS works either utilize or alter the generic news summarization datasets to evaluate the controllable summarization. As shown in Table 2, out of the 10 categories, we could find CA-specific datasets for only seven categories. We envisage that conducting evaluations with specialized datasets that align closely with real-world application scenarios or user requirements will help better in assessing the practical utility, robustness, and performance of CTS. It is evident from our survey of CTS systems that evaluations are often confined to specific domains, like news, possibly due to the abundance of available datasets in that domain. However, this narrow focus limits the evaluation of the CTS model’s robustness.

Standardization of metrics. The goal of the CTS task is to produce CA-specific summaries, warranting a metric tailored to capture the nuances of this particular attribute.

We observe, that comparing models for a specific CA-based CTS task is challenging due to the use of varying metrics, leading each study to redo evaluations for a fair comparison with prior work. Standardizing CA-specific evaluation metrics could offer a valuable solution.

Explainability. For effectively controlling user or application-specific attributes, it is imperative to leverage the understanding of the decision-making process within CTS systems. Also, this comprehension is essential for users or stakeholders, enabling them to discern how the system generates summaries from source text. This holds particular significance in applications where human decision-making or interpretation plays a pivotal role, such as in legal, medical, or financial domains. The existing CTS efforts lack proper emphasis on the explainability aspects, which can be readily addressed through the incorporation of suitable explanation methodologies (Abnar and Zuidema, 2020; Sundararajan et al., 2017; Lundberg and Lee, 2017).

Multi-lingual, multi-model, and code-mixed CTS. The existing literature on CTS predominantly focuses on works in English, with only one study addressing the topic in a Japanese context. We could not find any studies and datasets related to multilingual and code-mixed CTS approaches. Moreover, the task of controllable summarization in multi-modal and multi-document settings remains largely unexplored, presenting unique challenges for models to address and offering avenues

for intriguing research problems.

Multi-CA control. Even though, few of the works perform multi-attribute controllable summarization (Goyal et al., 2022; He et al., 2022; Zhang et al., 2023b), we observe that existing works predominantly investigate combinations of length and entity attributes (see Appendix F Table 5). As a future research direction, it’s essential to design models that consider other important combinations of control attributes, such as length, style, and saliency. Furthermore, creating standardized multi-CA benchmarks is crucial to facilitate the evaluations.

Reproducibility. In the detailed analysis outlined in Table 10, we note that 35% of research studies do not share their code publicly. Furthermore, 25% of the papers did not carry out any human evaluation, and among the remaining studies, 79% did not conduct Inter Annotator Agreement (IAA) assessments. The lack of reproducibility (Ito et al., 2023; Gao et al., 2023; Iskender et al., 2021) measures hinders the scientific community’s ability to validate and build upon existing work. On the other hand, the human study component should be a must for a text summarization evaluation scheme, otherwise, we are potentially overlooking essential aspects of real-world applicability.

Standing on the shoulders of LLMs. The rise and success of large language models (LLMs) have opened up unparalleled possibilities for leveraging their capabilities across diverse stages of the Natural Language Processing (NLP) pipeline. In the context of CTS, LLMs can be fine-tuned to grasp context-specific nuances about CAs without the need for a dedicated training set. Additionally, when it comes to evaluating CTS models, LLMs can serve as effective substitutes for human experts or judges (similar to Liu et al. (2023)), offering an efficient method for assessing performance.

8 Conclusions

We present a comprehensive survey on controllable text summarization (CTS) by offering a detailed analysis, from formalizing various controllable attributes, classifying them based on shared characteristics, and delving into existing datasets, proposed models, associated limitations, and evaluation strategies. Moreover, we discuss the challenges and prospects, making it a helpful guide for researchers interested in CTS. We plan to keep the GitHub repository regularly updated with the latest CTS works.

743 9 Limitations

744 Although we attempt to conduct a rigorous analysis
745 of existing literature on controllable summarization,
746 some works might have been possibly left out due
747 to variations in search keywords. Furthermore, due
748 to limited space, our survey primarily concentrates
749 solely on the high-level aspects of the approaches,
750 omitting a very fine-grained experimental compari-
751 son. Finally, our exploration of multilingual works
752 was limited as we encountered challenges in find-
753 ing them, likely influenced by the relatively low
754 attention from the research community. We aim to
755 further investigate the potential reasons behind the
756 challenges associated with multilingual CTS tasks.

757 10 Ethics statement

758 To uphold transparency and accountability, the
759 papers utilized in this survey are detailed in Ap-
760 pendix E Table 10. We have provided a com-
761 prehensive set of papers, accompanied by our qualita-
762 tive classification and annotations, enabling public
763 scrutiny and examination. Moreover, to alleviate
764 qualitative bias, each paper underwent review by
765 at least three different individuals independently,
766 aiming to minimize misclassification. We adhere
767 to the same methodology to validate the presence
768 of diverse observations in each paper. By incorpo-
769 rating these ethical considerations, we affirm our
770 dedication to conducting research in an ethical and
771 accountable manner.

772 References

773 Samira Abnar and Willem Zuidema. 2020. Quantify-
774 ing attention flow in transformers. In *Proceedings*
775 of the 58th Annual Meeting of the Association for
776 Computational Linguistics, pages 4190–4197, On-
777 line. Association for Computational Linguistics.

778 Ojas Ahuja, Jiacheng Xu, Akshay Gupta, Kevin
779 Horecka, and Greg Durrett. 2022. ASPECTNEWS:
780 Aspect-oriented summarization of news documents.
781 In *Proceedings of the 60th Annual Meeting of the*
782 *Association for Computational Linguistics (Volume*
783 *1: Long Papers)*, pages 6494–6506, Dublin, Ireland.
784 Association for Computational Linguistics.

785 Mehdi Allahyari, Seyedamin Pouriyeh, Mehdi Assefi,
786 Saeid Safaei, Elizabeth D Trippe, Juan B Gutier-
787 rez, and Krys Kochut. 2017. Text summariza-
788 tion techniques: a brief survey. *arXiv preprint*
789 *arXiv:1707.02268*.

790 Reinald Kim Amplayo, Stefanos Angelidis, and Mirella
791 Lapata. 2021. Aspect-controllable opinion summa-
792 rization. In *Proceedings of the 2021 Conference on*

793 *Empirical Methods in Natural Language Processing*,
794 pages 6578–6593, Online and Punta Cana, Domini-
795 can Republic. Association for Computational Lin-
796 guistics.

797 Stefanos Angelidis, Reinald Kim Amplayo, Yoshihiko
798 Suhara, Xiaolan Wang, and Mirella Lapata. 2021.
799 Extractive opinion summarization in quantized trans-
800 former spaces. *Transactions of the Association for*
801 *Computational Linguistics*, 9:277–293.

802 Stefanos Angelidis and Mirella Lapata. 2018. Sum-
803 marizing opinions: Aspect extraction meets senti-
804 ment prediction and they are both weakly supervised.
805 In *Proceedings of the 2018 Conference on Empiri-
806 cal Methods in Natural Language Processing*, pages
807 3675–3686, Brussels, Belgium. Association for Com-
808 putational Linguistics.

809 Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Ben-
810 gio. 2014. Neural machine translation by jointly
811 learning to align and translate. *arXiv preprint*
812 *arXiv:1409.0473*.

813 Seyed Ali Bahrainian, Sheridan Feucht, and Carsten
814 Eickhoff. 2022. NEWTS: A corpus for news topic-
815 focused summarization. In *Findings of the Associa-
816 tion for Computational Linguistics: ACL 2022*, pages
817 493–503, Dublin, Ireland. Association for Compu-
818 tational Linguistics.

819 Seyed Ali Bahrainian, George Zerveas, Fabio Crestani,
820 and Carsten Eickhoff. 2021. Cats: Customizable
821 abstractive topic-based summarization. *ACM Trans-
822 actions on Information Systems (TOIS)*, 40(1):1–24.

823 Iz Beltagy, Matthew E. Peters, and Arman Cohan.
824 2020. Longformer: The long-document transformer.
825 *arXiv:2004.05150*.

826 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
827 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
828 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
829 Askell, Sandhini Agarwal, Ariel Herbert-Voss,
830 Gretchen Krueger, Tom Henighan, Rewon Child,
831 Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens
832 Winter, Chris Hesse, Mark Chen, Eric Sigler, Ma-
teusz Litwin, Scott Gray, Benjamin Chess, Jack
833 Clark, Christopher Berner, Sam McCandlish, Alec
834 Radford, Ilya Sutskever, and Dario Amodei. 2020.
835 Language models are few-shot learners. In *Ad-
836 vances in Neural Information Processing Systems*,
837 volume 33, pages 1877–1901. Curran Associates,
838 Inc.

839 Shuyang Cao and Lu Wang. 2021. Inference time style
840 control for summarization. In *Proceedings of the*
841 *2021 Conference of the North American Chapter of*
842 *the Association for Computational Linguistics: Hu-*
843 *man Language Technologies*, pages 5942–5953, On-
844 line. Association for Computational Linguistics.

845 Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua,
846 Nicole Limtiaco, Rhomni St John, Noah Constant,
847 Mario Guajardo-Cespedes, Steve Yuan, Chris Tar,

849	et al. 2018. Universal sentence encoder.	<i>arXiv preprint arXiv:1803.11175</i> .	907	
850			908	
851	Hou Pong Chan, Lu Wang, and Irwin King.	2021. Controllable summarization with constrained markov decision process.	<i>Transactions of the Association for Computational Linguistics</i> , 9:1213–1232.	909
852			910	
853	Kushal Chawla, Balaji Vasan Srinivasan, and Niyati Chhaya.	2019. Generating formality-tuned summaries using input-dependent rewards.	In <i>Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)</i> , pages 833–842, Hong Kong, China. Association for Computational Linguistics.	911
854				
855	Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio.	2014. Learning phrase representations using RNN encoder–decoder for statistical machine translation.	In <i>Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 1724–1734, Doha, Qatar. Association for Computational Linguistics.	912
856				
857	Elizabeth Clark, Asli Celikyilmaz, and Noah A. Smith.	2019. Sentence mover’s similarity: Automatic evaluation for multi-sentence texts.	In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 2748–2760, Florence, Italy. Association for Computational Linguistics.	913
858				
859	Hoa Trang Dang.	2005. Overview of duc 2005.	In <i>Proceedings of the document understanding conference</i> , volume 2005, pages 1–12. Citeseer.	914
860				
861	Daniel Deutsch and Dan Roth.	2023. Incorporating question answering-based signals into abstractive summarization via salient span selection.	In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 575–588.	915
862				
863	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova.	2019. BERT: Pre-training of deep bidirectional transformers for language understanding.	In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.	916
864				
865	Zi-Yi Dou, Pengfei Liu, Hiroaki Hayashi, Zhengbao Jiang, and Graham Neubig.	2021. GSUM: A general framework for guided neural abstractive summarization.	In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 4830–4842, Online. Association for Computational Linguistics.	917
866				
867	Wafaa S El-Kassas, Cherif R Salama, Ahmed A Rafea, and Hoda K Mohamed.	2021. Automatic text summarization: A comprehensive survey.	<i>Expert systems with applications</i> , 165:113679.	918
868				
869				
870				
871	Angela Fan, David Grangier, and Michael Auli.	2018a. Controllable abstractive summarization.	In <i>Proceedings of the 2nd Workshop on Neural Machine Translation and Generation</i> , pages 45–54, Melbourne, Australia. Association for Computational Linguistics.	919
872				
873	Angela Fan, David Grangier, and Michael Auli.	2018b. Controllable abstractive summarization.	<i>ACL 2018</i> , page 45.	920
874				
875	Thibault Févry and Jason Phang.	2018. Unsupervised sentence compression using denoising auto-encoders.	In <i>Proceedings of the 22nd Conference on Computational Natural Language Learning</i> , pages 413–422, Brussels, Belgium. Association for Computational Linguistics.	921
876				
877	Rudolf Flesch.	1948. A new readability yardstick.	<i>Journal of Applied Psychology</i> , 32:221–233.	922
878				
879	Mahak Gambhir and Vishal Gupta.	2017. Recent automatic text summarization techniques: a survey.	<i>Artificial Intelligence Review</i> , 47:1–66.	923
880				
881	Mingqi Gao, Jie Ruan, and Xiaojun Wan.	2023. A reproduction study of the human evaluation of role-oriented dialogue summarization models.	In <i>Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems</i> , Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.	924
882				
883	Yang Gao, Wei Zhao, and Steffen Eger.	2020. SUPERT: Towards new frontiers in unsupervised evaluation metrics for multi-document summarization.	In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 1347–1354, Online. Association for Computational Linguistics.	925
884				
885	Bogdan Gliwa, Iwona Mochol, Maciej Biesek, and Aleksander Wawer.	2019. SAMSum corpus: A human-annotated dialogue dataset for abstractive summarization.	In <i>Proceedings of the 2nd Workshop on New Frontiers in Summarization</i> , pages 70–79, Hong Kong, China. Association for Computational Linguistics.	926
886				
887	Tomas Goldsack, Zheheng Luo, Qianqian Xie, Carolina Scarton, Matthew Shardlow, Sophia Ananiadou, and Chenghua Lin.	2023. Overview of the biolaysumm 2023 shared task on lay summarization of biomedical research articles.	In <i>The 22nd Workshop on Biomedical Natural Language Processing and BioNLP Shared Tasks</i> , pages 468–477, Toronto, Canada. Association for Computational Linguistics.	927
888				
889	Tanya Goyal, Nazneen Rajani, Wenhao Liu, and Wojciech Kryscinski.	2022. HydraSum: Disentangling style features in text summarization with multi-decoder models.	In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 464–479, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	928
890				
891				
892				
893				
894				
895				
896				
897				
898				
899				
900				
901				
902				
903				
904				
905				
906				

961	Hiroaki Hayashi, Prashant Budania, Peng Wang, Chris Ackerson, Raj Neervannan, and Graham Neubig.	Dongmin Hyun, Xiting Wang, Chayoung Park, Xing Xie, and Hwanjo Yu.	1017
962	2021. <i>WikiAsp: A dataset for multi-domain aspect-based summarization</i> . <i>Transactions of the Association for Computational Linguistics</i> , 9:211–225.	<i>Generating multiple-length summaries via reinforcement learning for unsupervised sentence summarization</i> .	1018
963		In <i>Findings of the Association for Computational Linguistics: EMNLP 2022</i> , pages 2939–2951, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	1019
964			1020
965			1021
966	Junxian He, Wojciech Kryscinski, Bryan McCann, Nazneen Rajani, and Caiming Xiong. 2022. <i>CTRL-sum: Towards generic controllable text summarization</i> .		1022
967	In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 5879–5915, Abu Dhabi, United Arab Emirates.		1023
968	Association for Computational Linguistics.		1024
969			
970	Ruidan He, Wee Sun Lee, Hwee Tou Ng, and Daniel Dahlmeier. 2017. <i>An unsupervised neural attention model for aspect extraction</i> .	Neslihan Iskender, Tim Polzehl, and Sebastian Möller.	1025
971	In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 388–397, Vancouver, Canada. Association for Computational Linguistics.	2021. <i>Reliability of human evaluation for text summarization: Lessons learned and challenges ahead</i> .	1026
972		In <i>Proceedings of the Workshop on Human Evaluation of NLP Systems (HumEval)</i> , Online. Association for Computational Linguistics.	1027
973			1028
974			1029
975			1030
976	Ruidan He, Wee Sun Lee, Hwee Tou Ng, and Daniel Dahlmeier. 2017. <i>An unsupervised neural attention model for aspect extraction</i> .	Takumi Ito, Qixiang Fang, Pablo Mosteiro, Albert Gatt, and Kees van Deemter.	1031
977	In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 388–397, Vancouver, Canada. Association for Computational Linguistics.	2023. <i>Challenges in reproducing human evaluation results for role-oriented dialogue summarization</i> .	1032
978		In <i>Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems</i> , Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.	1033
979			1034
980	Yuta Hitomi, Yuya Taguchi, Hideaki Tamori, Ko Kikuta, Jiro Nishitoba, Naoaki Okazaki, Kentaro Inui, and Manabu Okumura. 2019. <i>A large-scale multi-length headline corpus for analyzing length-constrained headline generation model evaluation</i> .		1035
981	In <i>Proceedings of the 12th International Conference on Natural Language Generation</i> , pages 333–343, Tokyo, Japan.		1036
982	Association for Computational Linguistics.		
983			
984	Sepp Hochreiter and Jürgen Schmidhuber. 1997. <i>Long short-term memory</i> .	Di Jin, Zhijing Jin, Joey Tianyi Zhou, Lisa Orii, and Peter Szolovits.	1037
985	<i>Neural computation</i> , 9(8):1735–1780.	2020a. <i>Hooks in the headline: Learning to generate headlines with controlled styles</i> .	1038
986		In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 5082–5093, Online. Association for Computational Linguistics.	1039
987			1040
988	Ella Hofmann-Coyle, Mayank Kulkarni, Lingjue Xie, Mounica Maddela, and Daniel Preotiuc-Pietro. 2022. <i>Extractive entity-centric summarization as sentence selection using bi-encoders</i> .		1041
989	In <i>Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)</i> , pages 326–333, Online only. Association for Computational Linguistics.		1042
990			1043
991			
992	Chao-Chun Hsu and Chenhao Tan. 2021. <i>Decision-focused summarization</i> .	Hanqi Jin, Tianming Wang, and Xiaojun Wan.	1044
993	In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 117–132, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	2020b. <i>Semsum: Semantic dependency guided neural abstractive summarization</i> .	1045
994		In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 34, pages 8026–8033.	1046
995			1047
996			1048
997			
998			
999			
1000			
1001	Zhitong Hu and Li Erran Li. 2021. <i>A causal lens for controllable text generation</i> .	Yuta Kikuchi, Graham Neubig, Ryohei Sasano, Hiroya Takamura, and Manabu Okumura.	1049
1002	<i>Advances in Neural Information Processing Systems</i> , 34:24941–24955.	2016. <i>Controlling output length in neural encoder-decoders</i> .	1050
1003		In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 1328–1338.	1051
1004			1052
1005			1053
1006			1054
1007	Kung-Hsiang Huang, Sifffi Singh, Xiaofei Ma, Wei Xiao, Feng Nan, Nicholas Dingwall, William Yang Wang, and Kathleen McKeown. 2023. <i>SWING: Balancing coverage and faithfulness for dialogue summarization</i> .	Kundan Krishna, Aniket Murhekar, Saumitra Sharma, and Balaji Vasan Srinivasan.	1055
1008	In <i>Findings of the Association for Computational Linguistics: EACL 2023</i> , Dubrovnik, Croatia.	2018. <i>Vocabulary tailored summary generation</i> .	1056
1009	Association for Computational Linguistics.	In <i>Proceedings of the 27th International Conference on Computational Linguistics</i> , pages 795–805.	1057
1010			1058
1011			1059
1012			
1013			
1014			
1015			
1016			
1068	Wojciech Kryściński, Romain Paulus, Caiming Xiong, and Richard Socher.	2018. <i>Improving abstraction in text summarization</i> .	1069
1069	In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 1808–1817, Brussels, Belgium.		1070
1070	Association for Computational Linguistics.		1071
1071			1072
1072			1073
1073			

1074	Jingun Kwon, Hidetaka Kamigaito, and Manabu Okumura. 2023. Abstractive document summarization with summary-length prediction . In <i>Findings of the Association for Computational Linguistics: EACL 2023</i> , pages 606–612.	1132
1075		1133
1076		1134
1077		1135
1078		1136
1079	Mike Lewis, Yinhai Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension . In <i>Annual Meeting of the Association for Computational Linguistics</i> .	1137
1080		1138
1081		1139
1082		1140
1083		1141
1084		1142
1085		1143
1086	Chenliang Li, Weiran Xu, Si Li, and Sheng Gao. 2018. Guiding generation for abstractive text summarization based on key information guide network . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)</i> , pages 55–60, New Orleans, Louisiana. Association for Computational Linguistics.	1144
1087		1145
1088		1146
1089		1147
1090		1148
1091		1149
1092		1150
1093		1151
1094		1152
1095	Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation . 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 4582–4597, Online. Association for Computational Linguistics.	1153
1096		1154
1097		1155
1098		1156
1099		1157
1100		1158
1101		1159
1102		1160
1103	Xinnian Liang, Chao Bian, Shuangzhi Wu, and Zhoujun Li. 2022. Towards modeling role-aware centrality for dialogue summarization . In <i>Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)</i> , pages 43–50, Online only. Association for Computational Linguistics.	1161
1104		1162
1105		1163
1106		1164
1107		1165
1108		1166
1109		1167
1110		1168
1111		1169
1112	Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries . In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain. Association for Computational Linguistics.	1170
1113		1171
1114		1172
1115		1173
1116	Haitao Lin, Liqun Ma, Junnan Zhu, Lu Xiang, Yu Zhou, Jiajun Zhang, and Chengqing Zong. 2021. CSDS: A fine-grained Chinese dataset for customer service dialogue summarization . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 4436–4451, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	1174
1117		1175
1118		1176
1119		1177
1120		1178
1121		1179
1122		1180
1123		1181
1124	Haitao Lin, Junnan Zhu, Lu Xiang, Yu Zhou, Jiajun Zhang, and Chengqing Zong. 2022. Other roles matter! enhancing role-oriented dialogue summarization via role interactions . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2545–2558, Dublin, Ireland. Association for Computational Linguistics.	1182
1125		1183
1126		1184
1127		1185
1128		1186
1129		1187
1130		1188
1131		1189
520	Puyuan Liu, Xiang Zhang, and Lili Mou. 2022a. A character-level length-control algorithm for non-autoregressive sentence summarization . <i>Advances in Neural Information Processing Systems</i> , 35:29101–29112.	1190
521		1191
522	Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. G-eval: NLG evaluation using gpt-4 with better human alignment , may 2023. <i>arXiv preprint arXiv:2303.16634</i> .	1192
523		1193
524	Yang Liu and Mirella Lapata. 2019. Text summarization with pretrained encoders . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3730–3740, Hong Kong, China. Association for Computational Linguistics.	1194
525		1195
526	Yizhu Liu, Qi Jia, and Kenny Zhu. 2022b. Length control in abstractive summarization by pretraining information selection . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 6885–6895, Dublin, Ireland. Association for Computational Linguistics.	1196
527		1197
528	Yizhu Liu, Zhiyi Luo, and Kenny Zhu. 2018. Controlling length in abstractive summarization using a convolutional neural network . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 4110–4119.	1198
529		1199
530	Zhengyuan Liu and Nancy F. Chen. 2021. Controllable neural dialogue summarization with personal named entity planning . In <i>Conference on Empirical Methods in Natural Language Processing</i> .	1200
531		1201
532	Scott M Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions . <i>Advances in neural information processing systems</i> , 30.	1202
533		1203
534	Zheheng Luo, Qianqian Xie, and Sophia Ananiadou. 2022. Readability controllable biomedical document summarization . In <i>Findings of the Association for Computational Linguistics: EMNLP 2022</i> , pages 4667–4680, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	1204
535		1205
536	Minh-Thang Luong, Hieu Pham, and Christopher D Manning. 2015. Effective approaches to attention-based neural machine translation . In <i>Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing</i> , pages 1412–1421.	1206
537		1207
538	Mounica Maddela, Mayank Kulkarni, and Daniel Preotiuc-Pietro. 2022. EntSUM: A data set for entity-centric extractive summarization . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3355–3366, Dublin, Ireland. Association for Computational Linguistics.	1208
539		1209
540	Takuya Makino, Tomoya Iwakura, Hiroya Takamura, and Manabu Okumura. 2019. Global optimization	1210
541		1211

1187	under length constraint for neural text summarization. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 1039–1048, Florence, Italy. Association for Computational Linguistics.	Shashi Narayan, Yao Zhao, Joshua Maynez, Gonçalo Simões, Vitaly Nikolaev, and Ryan McDonald. 2021. Planning with learned entity prompts for abstractive summarization. <i>Transactions of the Association for Computational Linguistics</i> , 9:1475–1492.	1245
1188			1246
1189			1247
1190			1248
1191			1249
1192	Dhruv Mehra, Lingjue Xie, Ella Hofmann-Coyle, Mayank Kulkarni, and Daniel Preotiuc-Pietro. 2023. EntSUMv2: Dataset, models and evaluation for more abstractive entity-centric summarization. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 5538–5547, Singapore. Association for Computational Linguistics.	Narges Nazari and MA Mahdavi. 2019. A survey on automatic text summarization. <i>Journal of AI and Data Mining</i> , 7(1):121–135.	1250
1193			1251
1194			1252
1195			1253
1196			1254
1197			1255
1198			1256
1199			1257
1200	Rajdeep Mukherjee, Hari Chandana Peruri, Uppada Vishnu, Pawan Goyal, Sourangshu Bhattacharya, and Niloy Ganguly. 2020. Read what you need: Controllable aspect-based opinion summarization of tourist reviews. In <i>Proceedings of the 43rd international ACM SIGIR conference on research and development in information retrieval</i> , pages 1825–1828.	Jun-Ping Ng and Viktoria Abrecht. 2015. Better summarization evaluation with word embeddings for ROUGE. In <i>Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing</i> , pages 1925–1930, Lisbon, Portugal. Association for Computational Linguistics.	1258
1201			
1202			
1203			
1204			
1205			
1206			
1207	Sourajit Mukherjee, Anubhav Jangra, Sriparna Saha, and Adam Jatowt. 2022. Topic-aware multimodal summarization. In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2022</i> , pages 387–398, Online only. Association for Computational Linguistics.	Paul Over and James Yen. 2004. An introduction to duc-2004. <i>National Institute of Standards and Technology</i> .	1259
1208			1260
1209			1261
1210			
1211			
1212			
1213	Ramesh Nallapati, Feifei Zhai, and Bowen Zhou. 2017. Summarunner: A recurrent neural network based sequence model for extractive summarization of documents. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 31.	Artidoro Pagnoni, Alex Fabbri, Wojciech Kryscinski, and Chien-Sheng Wu. 2023. Socratic pretraining: Question-driven pretraining for controllable summarization. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 12737–12755, Toronto, Canada. Association for Computational Linguistics.	1262
1214			1263
1215			1264
1216			1265
1217			1266
1218	Ramesh Nallapati, Bowen Zhou, Cicero dos Santos, Çağlar Gülcöhre, and Bing Xiang. 2016. Abstractive text summarization using sequence-to-sequence RNNs and beyond. In <i>Proceedings of the 20th SIGNLL Conference on Computational Natural Language Learning</i> , Berlin, Germany. Association for Computational Linguistics.	Kishore Papineni, Salim Roukos, Todd Ward, and Wei Jing 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, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.	1267
1219			1268
1220			
1221			
1222			
1223			
1224			
1225	Courtney Napoles, Matthew R Gormley, and Benjamin Van Durme. 2012. Annotated gigaword. In <i>Proceedings of the joint workshop on automatic knowledge base construction and web-scale knowledge extraction (AKBC-WEKEX)</i> , pages 95–100.	Maxime Peyrard, Teresa Botschen, and Iryna Gurevych. 2017. Learning to score system summaries for better content selection evaluation. In <i>Proceedings of the Workshop on New Frontiers in Summarization</i> , pages 74–84, Copenhagen, Denmark. Association for Computational Linguistics.	1269
1226			1270
1227			1271
1228			1272
1229			1273
1230	Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. Don’t give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 1797–1807, Brussels, Belgium. Association for Computational Linguistics.	Shrimai Prabhumoye, Alan W Black, and Ruslan Salakhutdinov. 2020. Exploring controllable text generation techniques. In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 1–14, Barcelona, Spain (Online). International Committee on Computational Linguistics.	1274
1231			1275
1232			
1233			
1234			
1235			
1236			
1237	Shashi Narayan, Gonçalo Simões, Yao Zhao, Joshua Maynez, Dipanjan Das, Michael Collins, and Mirella Lapata. 2022. A well-composed text is half done! composition sampling for diverse conditional generation. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1319–1339, Dublin, Ireland. Association for Computational Linguistics.	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>Journal of Machine Learning Research</i> , 21(140):1–67.	1276
1238			1277
1239			1278
1240			1279
1241			1280
1242			1281
1243			
1244			

1301	<i>Conference on Empirical Methods in Natural Language Processing</i> , pages 379–389.	1358
1302		1359
1303	Evan Sandhaus. 2008. The new york times annotated corpus . <i>Linguistic Data Consortium, Philadelphia</i> , 6(12):e26752.	1360
1304		1361
1305		1362
1306	Ritesh Sarkhel, Moniba Keymanesh, Arnab Nandi, and Srinivasan Parthasarathy. 2020. Interpretable multi-headed attention for abstractive summarization at controllable lengths . In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 6871–6882, Barcelona, Spain (Online). International Committee on Computational Linguistics.	1363
1307		1364
1308		1365
1309		1366
1310		1367
1311		1368
1312		1369
1313	Thomas Scialom, Sylvain Lamprier, Benjamin Piwowarski, and Jacopo Staiano. 2019. Answers unite! unsupervised metrics for reinforced summarization models . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3246–3256, Hong Kong, China. Association for Computational Linguistics.	1370
1314		1371
1315		1372
1316		1373
1317		
1318		
1319		
1320		
1321		
1322	Abigail See, Peter J. Liu, and Christopher D. Manning. 2017. Get to the point: Summarization with pointer-generator networks . In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1073–1083, Vancouver, Canada. Association for Computational Linguistics.	1374
1323		1375
1324		1376
1325		1377
1326		1378
1327		1379
1328		1380
1329	Satoshi Sekine and Chikashi Nobata. 2003. A survey for multi-document summarization . In <i>Proceedings of the HLT-NAACL 03 Text Summarization Workshop</i> , pages 65–72.	1381
1330		
1331		
1332		
1333	Chenhui Shen, Liying Cheng, Lidong Bing, Yang You, and Luo Si. 2022a. SentBS: Sentence-level beam search for controllable summarization . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 10256–10265, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	1382
1334		1383
1335		1384
1336		1385
1337		1386
1338		1387
1339		1388
1340	Chenhui Shen, Liying Cheng, Ran Zhou, Lidong Bing, Yang You, and Luo Si. 2022b. MReD: A meta-review dataset for structure-controllable text generation . In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 2521–2535, Dublin, Ireland. Association for Computational Linguistics.	1389
1341		1390
1342		1391
1343		1392
1344		1393
1345		1394
1346		1395
1347	Zejiang Shen, Kyle Lo, Lauren Yu, Nathan Dahlberg, Margo Schlanger, and Doug Downey. 2022c. MultiLexsum: Real-world summaries of civil rights lawsuits at multiple granularities . <i>Advances in Neural Information Processing Systems</i> , 35:13158–13173.	1396
1348		1397
1349		1398
1350		1399
1351	Kaiqiang Song, Bingqing Wang, Zhe Feng, and Fei Liu. 2021. A new approach to overgenerating and scoring abstractive summaries . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1392–1404, Online. Association for Computational Linguistics.	1400
1352		1401
1353		
1354		
1355		
1356		
1357		
	Kaiqiang Song, Bingqing Wang, Zhe Feng, Ren Liu, and Fei Liu. 2020. Controlling the amount of verbatim copying in abstractive summarization . In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 34, pages 8902–8909.	1402
		1403
		1404
		1405
		1406
	Yoshihiko Suhara, Xiaolan Wang, Stefanos Angelidis, and Wang-Chiew Tan. 2020. OpinionDigest: A simple framework for opinion summarization . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 5789–5798, Online. Association for Computational Linguistics.	1407
		1408
		1409
		1410
	Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic attribution for deep networks . In <i>International conference on machine learning</i> , pages 3319–3328. PMLR.	1411
		1412
	Sho Takase and Naoaki Okazaki. 2019. Positional encoding to control output sequence length . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 3999–4004, Minneapolis, Minnesota. Association for Computational Linguistics.	1413
		1414
	Bowen Tan, Lianhui Qin, Eric Xing, and Zhiting Hu. 2020. Summarizing text on any aspects: A knowledge-informed weakly-supervised approach . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 6301–6309, Online. Association for Computational Linguistics.	1415
		1416
	Don Tuggener, Margot Mieskes, Jan Deriu, and Mark Cieliebak. 2021. Are we summarizing the right way? a survey of dialogue summarization data sets . In <i>Proceedings of the Third Workshop on New Frontiers in Summarization</i> , pages 107–118, Online and in Dominican Republic. Association for Computational Linguistics.	1417
		1418
	Oleg Vasilyev, Vedant Dharnidharka, and John Bohannon. 2020. Fill in the BLANC: Human-free quality estimation of document summaries . In <i>Proceedings of the First Workshop on Evaluation and Comparison of NLP Systems</i> , pages 11–20, Online. Association for Computational Linguistics.	1419
		1420
	Jiaan Wang, Fandong Meng, Duo Zheng, Yunlong Liang, Zhixu Li, Jianfeng Qu, and Jie Zhou. 2022. A survey on cross-lingual summarization . <i>Transactions of the Association for Computational Linguistics</i> , 10:1304–1323.	1421
		1422
	Junli Wang, Chenyang Zhang, Dongyu Zhang, Haibo Tong, Chungang Yan, and Changjun Jiang. 2024. A recent survey on controllable text generation: a causal perspective . <i>Fundamental Research</i> .	1423
		1424
	Chien-Sheng Wu, Linqing Liu, Wenhao Liu, Pontus Stenetorp, and Caiming Xiong. 2021. Controllable	1425
		1426

- 1413 abstractive dialogue summarization with sketch su-
1414 pervision. In *Findings of the Association for Com-*
1415 *putational Linguistics: ACL-IJCNLP 2021*, pages
1416 5108–5122, Online. Association for Computational
1417 Linguistics.
- 1418 Zhongyi Yu, Zhenghao Wu, Hao Zheng, Zhe XuanYuan,
1419 Jefferson Fong, and Weifeng Su. 2021. **LenAtten:**
1420 **An effective length controlling unit for text summa-**
1421 **rization.** In *Findings of the Association for Com-*
1422 *putational Linguistics: ACL-IJCNLP 2021*, pages
1423 363–370, Online. Association for Computational Lin-
1424 guistics.
- 1425 Weizhe Yuan, Graham Neubig, and Pengfei Liu. 2021.
1426 **Bartscore: Evaluating generated text as text genera-**
1427 **tion.** In *Advances in Neural Information Processing*
1428 *Systems*, volume 34, pages 27263–27277. Curran As-
1429 sociates, Inc.
- 1430 Hanqing Zhang, Haolin Song, Shaoyu Li, Ming Zhou,
1431 and Dawei Song. 2023a. **A survey of controllable**
1432 **text generation using transformer-based pre-trained**
1433 **language models.** *ACM Computing Surveys*, 56(3):1–
1434 37.
- 1435 Shiyue Zhang and Mohit Bansal. 2021. **Finding a bal-**
1436 **anced degree of automation for summary evaluation.**
1437 In *Proceedings of the 2021 Conference on Empiri-*
1438 *cal Methods in Natural Language Processing*, pages
1439 6617–6632, Online and Punta Cana, Dominican Re-
1440 public. Association for Computational Linguistics.
- 1441 Yusen Zhang, Yang Liu, Ziyi Yang, Yuwei Fang, Yulong
1442 Chen, Dragomir Radev, Chenguang Zhu, Michael
1443 Zeng, and Rui Zhang. 2023b. **Macsum: Control-**
1444 **lable summarization with mixed attributes.** *Transac-*
1445 *tions of the Association for Computational Linguis-*
1446 *tics*, 11:787–803.
- 1447 Changmeng Zheng, Yi Cai, Guanjie Zhang, and Qing Li.
1448 2020. **Controllable abstractive sentence summariza-**
1449 **tion with guiding entities.** In *Proceedings of the 28th*
1450 *International Conference on Computational Linguis-*
1451 *tics*, pages 5668–5678, Barcelona, Spain (Online).
1452 International Committee on Computational Linguis-
1453 tics.
- 1454 Ming Zhong, Yang Liu, Suyu Ge, Yuning Mao, Yizhu
1455 Jiao, Xingxing Zhang, Yichong Xu, Chenguang Zhu,
1456 Michael Zeng, and Jiawei Han. 2022. **Unsupervised**
1457 **multi-granularity summarization.** In *Findings of the*
1458 *Association for Computational Linguistics: EMNLP*
1459 *2022*, pages 4980–4995, Abu Dhabi, United Arab
1460 Emirates. Association for Computational Linguistics.
- 1461 Ming Zhong, Da Yin, Tao Yu, Ahmad Zaidi, Mutethia
1462 Mutuma, Rahul Jha, Ahmed Hassan Awadallah, Asli
1463 Celikyilmaz, Yang Liu, Xipeng Qiu, and Dragomir
1464 Radev. 2021. **QMSum: A new benchmark for query-**
1465 **based multi-domain meeting summarization.** In *Pro-*
1466 *ceedings of the 2021 Conference of the North Amer-*
1467 *ican Chapter of the Association for Computational*
1468 *Linguistics: Human Language Technologies*, pages
1469 5905–5921, Online. Association for Computational
1470 Linguistics.
- Yang Zhong and Diane Litman. 2023. **Strong–**
structure controllable legal opinion summary gen-
eration. *arXiv preprint arXiv:2309.17280*.
- 1471
1472
1473

A Controllable Attributes Classification

Different terms have been used to describe the same CAs, which exhibit similar characteristics and objectives (such as “Salience: Key information”, and “Coverage: Granularity”). Additionally, numerous attributes can be encapsulated by a representative class; for instance, “Style” may serve as a class encompassing Tone, Readability, Humor, Romance, and similar aspects, facilitating their classification within the same category.

Class	Attribute
Style	Tone, Readability, Humor, Romance, Clickbait
Coverage	Coverage, Granularity
Entity	Entity, Keyword
Topic	Topic, Aspect, Decision of interest, Opinion based on user interest
Abstractivity	Abstractiveness, Extractiveness, Novelty
Salience	Salience, Key information

Table 3: Merging of attributes into representative classes.

B Survey papers selection criteria

We have collected a total of 105 papers pertaining to the controllable text summarization task (see Table 4). Among these, six papers are pertinent to CTS, albeit they have not undergone peer review. Additionally, 23 papers touch upon the summarization aspect to some extent, although they may not be directly aligned with controllable summarization. Furthermore, we have excluded 15 papers as they primarily discuss controllable text generation or focus on enhancing the summarization task without specifically controlling any CTS attributes. Post to applying the above three filters we left out 61 peer-reviewed and relevant papers to CTS.

C Evaluation Approaches

We have listed the automatic and human evaluation methodologies along with their respective metric details in Table 7 and Table 8. The automatic evaluation metrics are categorized into three groups: embedding-based, n-gram-based, and miscellaneous. Additionally, we present a compilation

Criteria	Number of papers
arXiv version	6
Not relevant	23
Enhancement	15
Relevant	61
Total	105

Table 4: Survey papers filtration criteria.

of papers organized by aspects, each associated with the relevant metrics, along with concise descriptions. As for human evaluation, we specify the corresponding metrics and provide definitions based on the attributes under consideration.

D Model Descriptions

As outlined in Table 9, we augment novel contributions, utilized dataset, and the corresponding limitation for each paper, all aligned with the respective controllable attribute.

E Survey papers checklist explanation

To underscore the comprehensiveness of our survey, as mentioned in Table 10, we include 23 features for each paper. For easier understanding, we briefly describe each feature in the master table below.

- *Paper*: Citation of the paper.
- *Year*: Year of the publication.
- *Venue*: Paper published conference or journal.
- *Controllable attribute*: Controllable attribute(s) concentrates in the paper.
- *Controlling more than one aspect*: Whether the paper handles more than one controllable aspect or not?
- *Model type*: Type of the model used in the paper such as encoder-decoder, encoder, or decoder architecture.
- *Training strategy*: Training approaches employed to perform CTS task.
- *Approach*: Type of the training approach employed to perform CTS task.
- *Code access*: Whether the code is publicly accessible or not?
- *Code link*: Address of the public repository.
- *Dataset*: Dataset utilized in the paper.

	Controllable Attributes						
	Length	Entity	Style	Abstractivity	Coverage	Saliency	Topic
Fan et al. (2018a)	✓	✓	✓	✗	✗	✗	✗
Zhang et al. (2023b)	✓	✗	✗	✗	✓	✗	✓
Chan et al. (2021)	✓	✓	✗	✓	✗	✗	✗
See et al. (2017)	✗	✗	✗	✓	✓	✗	✗
Pagnoni et al. (2023)	✗	✓	✗	✗	✗	✓	✗
He et al. (2022)	✓	✓	✗	✗	✗	✗	✗
Nallapati et al. (2017)	✗	✗	✗	✓	✗	✓	✗

Table 5: Support for multiple controllable attributes across various models.

Dataset	Controllable attribute(s)	Human-annotated	Size	Domain	Dataset URL
Multi-LexSum (Shen et al., 2022c)	Coverage	Yes	9280	Legal	https://tinyurl.com/22ksfase
GranDUC (Zhong et al., 2022)	Coverage	Yes	50	News	https://tinyurl.com/2x72ubrw
TS and PLS (Luo et al., 2022)	Style	No	28124	Biomedical	https://tinyurl.com/yck3v9px
MACSUM (Zhang et al., 2023b)	Length, Coverage, Topic	Yes	9686	News, meetings	https://tinyurl.com/3d2dsc7u
NEWTS (Bahrainian et al., 2022)	Topic	Yes	6000	News	https://tinyurl.com/36hzk3ew
WikiAsp (Hayashi et al., 2021)	Topic	No	320272	Encyclopedia	https://tinyurl.com/3u45hfbn
ASPECTNEWS (Ahuja et al., 2022)	Topic	No	2000	News	https://tinyurl.com/bdzxs8ej
Tourism ASPECTS (Mukherjee et al., 2020)	Topic	No	7000	Reviews	https://tinyurl.com/ypjhhrxv
EntSUM (Maddela et al., 2022)	Entity	Yes	2788	News	https://tinyurl.com/2pz9vzyw
JAMUL (Hitomi et al., 2019)	Length	No	1932398	News	https://tinyurl.com/3s3ecua9
CSDS (Lin et al., 2021)	Role	Yes	10700	Dialogues	https://tinyurl.com/adk7zc7u
MReD (Shen et al., 2022b)	Structure	Yes	7089	Meta reviews	https://tinyurl.com/4nn87fd6

Table 6: List of controllable summarization datasets.

- *Source*: Source of the dataset used in the paper.
- *Nature of the data*: Dataset creation/acquisition strategy.
- *Data release*: Public availability of the dataset.
- *Domain*: The corresponding domain of the dataset.
- *Data link*: Public repository link to the dataset.
- *Metric name*: Name of the metric used in the paper.
- *Proposed new metric*: Names of the proposed new automatic evaluation metrics.
- *Human evaluation*: Human evaluation performed or not?
- *Metric names*: Name of the metrics used to perform human evaluation.
- *IAA*: Whether Inter Annotator Agreement assessment performed or not?
- *Limitation*: Any limitations of the paper mentioned or not?
- *Reproducibility*: Rate the reproducibility of the paper.

From the master table, we have represented our observations in Figures 2, 3, 4, 5, 6.

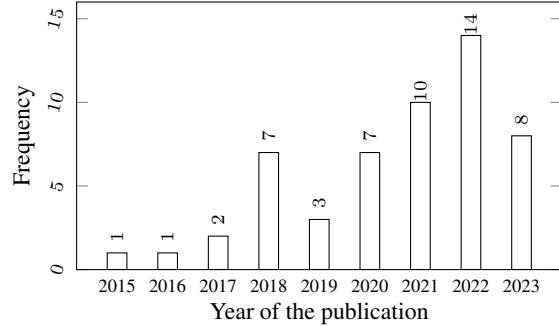


Figure 2: Year-wise papers published in CTS to handle various controllable attributes.

F Multi CA control

We list out the papers that attempt to control more than one controllable attribute in Table 5. We note that current research primarily explores combinations of length and entity attributes compared to other combinations.

Automatic Evaluation				
Type of metric	Attribute	Papers	Metrics	Description
Embedding-based (Language Model)	General	Lin et al. (2022), Liang et al. (2022) Song et al. (2020), Shen et al. (2022a), Cao and Wang (2021) , Deutsch and Roth (2023), Chan et al. (2021), Pagnoni et al. (2023), Lin et al. (2022), Liang et al. (2022), Narayan et al. (2022), Zhong and Litman (2023), Ribeiro et al. (2023), Shen et al. (2022c), Maddela et al. (2022), Lin et al. (2021)	MoverScore BERTScore	Computed using pretrained language models, either by computing similarity scores between reference and generated text embeddings or through likelihood computation of the generated text.
		Huang et al. (2023) Zheng et al. (2020)	BartScore Bert-Reo	
		Luo et al. (2022)	Masked Noun Phrase-based Text Complexity, Ranked NP Based Text Complexity, Masked Random Token-Based Text Complexity	
	Readability			
Ngram Based	General	Lin et al. (2022), (Liang et al., 2022), Jin et al. (2020a), Narayan et al. (2022) All except* Zhang et al. (2023b), Goldsack et al. (2023), Cao and Wang (2021), Hsu and Tan (2021), Hofmann-Coyle et al. (2022)	BLEU	These metrics are based on matching ngram tokens between reference and generated summaries
		Jin et al. (2020a), Sarkhel et al. (2020) Jin et al. (2020b)	ROUGE	
			METEOR	
			Word Mover's Distance	
	Miscellaneous			
		Length	Absolute Length, Compression Ratio, Length Variance, Var, Bin Percentage	
		Entity	QA-F1 Entity Planning, Entity Specificity	Non-normative metrics proposed by authors to evaluate specific controlled aspect
		Topic, Speaker, Length, Extractiveness, Specificity	Control Correlation, Control Error Rate	
		Abstractiveness, Degree of Specificity	Abstractiveness, Degree of specificity	
	Readability	Goyal et al. (2022)	Dale-Chall	
		Goyal et al. (2022), Cao and Wang (2021)	Flesch Reading Ease, Gunning Fog Index, Coleman Liau Index	

Table 7: Automatic evaluation metrics for controllable summarization, “General” refers to all controllable attributes.

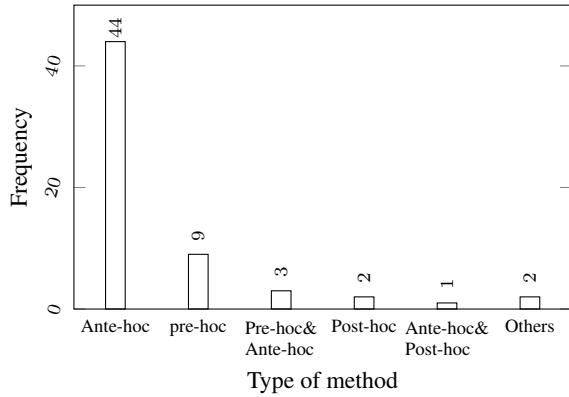


Figure 3: Various training approaches utilized to perform CTS tasks.

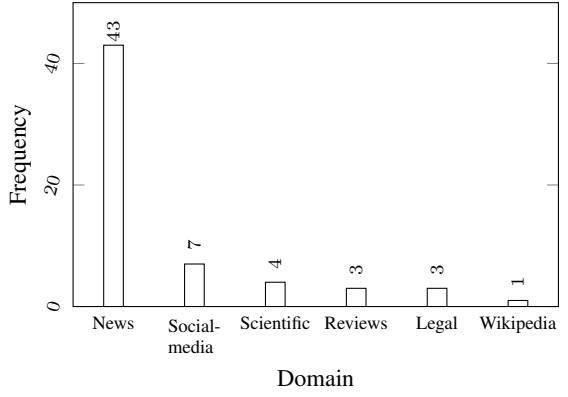


Figure 4: Domains utilized in CTS; most of the existing CTS tasks build on news domain data due to ease in accessibility.

Human Evaluation			
Aspect(s)	Papers	Metrics	Short description
Abstractivity, Length, Style, Topic, Coverage Role, Entity	Sarkhel et al. (2020), Song et al. (2020) Liu et al. (2022b), Cao and Wang (2021) Ampayo et al. (2021), Wu et al. (2021), Jin et al. (2020b), Kwon et al. (2023), Lin et al. (2022), Liu and Chen (2021), Zheng et al. (2020), Bahraianian et al. (2021), Suhara et al. (2020), Lin et al. (2021)	Informativeness	Has the summary covered key content of the input text?
Structure, Length, Entity, Salience, Topic	Kwon et al. (2023)	Conciseness/Granularity	Is the key information presented in a crisp way?
Structure, Style, Topic Length, Entity, Abstractivity, Coverage, Role, Diversity	Goyal et al. (2022), Tan et al. (2020) Yu et al. (2021), Shen et al. (2022b), Song et al. (2020), Liu et al. (2022b), Févry and Phang (2018), Zheng et al. (2020), Shen et al. (2022a), Cao and Wang (2021), Ampayo et al. (2021), Hyun et al. (2022), Chan et al. (2021), Jin et al. (2020b), Liu et al. (2022a), Lin et al. (2022), Zhong et al. (2022), Jin et al. (2020a), Lin et al. (2021)	Fluency/Grammaticality	Are the sentences in a summary grammatically correct?
Role, Topic, Diversity	Narayan et al. (2022), Suhara et al. (2020), Lin et al. (2022), Lin et al. (2021), Mukherjee et al. (2020)	Non-redundancy/Diversity	Is the summary conveying diverse information?
Topic	Krishna et al. (2018)	Contextual Appropriateness	Is the substituted word more readable in the summary?
Style, Diversity	Goyal et al. (2022), Narayan et al. (2022), Chan et al. (2021), Jin et al. (2020b), Zhong et al. (2022), Huang et al. (2023)	Faithfulness/Factuality	Does the summary present factually correct content with respect to the source?
Style, Topic, Entity, Structure	Goyal et al. (2022), Zhong and Litman (2023)	Coherence	Is the summary composed of correlated sentences?
Style, Structure, Length, Entity, Abstractivity, Coverage, Topic, Diversity	Goyal et al. (2022), He et al. (2022), Shen et al. (2022b), Chan et al. (2021), Krishna and Srinivasan (2018), Shen et al. (2022a), Cao and Wang (2021), Zhong et al. (2022), Jin et al. (2020a), Kryściński et al. (2018), Luo et al. (2022), Huang et al. (2023)	Relevance	Does the summary contain relevant information regarding the user provided attribute (topic/entity)?
Abstractivity, Length Length, Entity	Song et al. (2020), Hyun et al. (2022), Févry and Phang (2018), Huang et al. (2023) He et al. (2022), Yu et al. (2021)	Truthfulness/Fidelity Accuracy/Correctness	Has the summary successfully preserved the meaning of the original text? Is the information in the summary accurate?
Style	Kryściński et al. (2018), Ribeiro et al. (2023), Cao and Wang (2021)	Readability	Is the text inside the summary readable?
Length	Yu et al. (2021), Liu et al. (2022a)	Completeness	Does the summary contain incomplete text?
Topic, Structure	Zhong and Litman (2023), Mukherjee et al. (2020), Mukherjee et al. (2022)	Coverage	Does the summary include all the topics or aspects defined in the source?

Table 8: Human evaluation metrics for controllable text summarization.

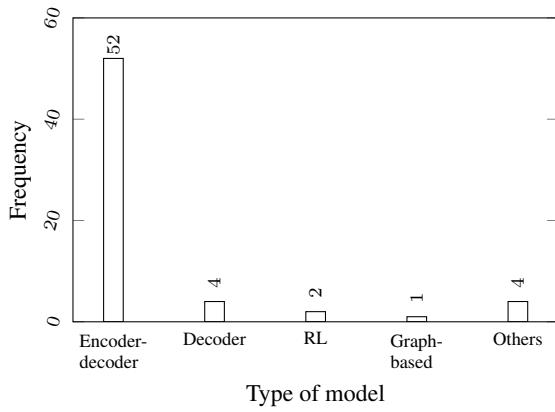


Figure 5: Type of models used in CTS; the majority of the models fall under standard sequence-to-sequence architecture.

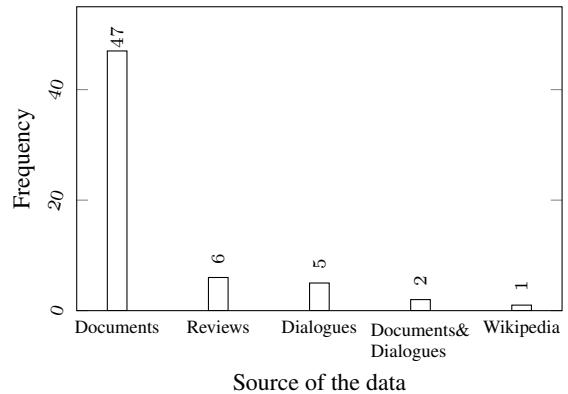


Figure 6: Source of the datasets used for the CTS task. The majority of the data samples are of ‘document’ type.

Aspect	Paper	Novel contribution	Dataset(s)	Limitations
Structure	Shen et al. (2022b)	Prepend structure prompt to the input	MRed	Subsequent generations deviate from the desired output
	Shen et al. (2022a)	Sentence-beam approach	MRed	Decoding methods significantly impact performance
	Zhong and Litman (2023)	Utilize predicted-role argument to control the structure	CanLII	Computationally expensive
Abstractivity	See et al. (2017)	Pointer-generator network	CNNDM	Failed to achieve higher abstraction and ineffective in core text selection
	Kryściński et al. (2018)	Decouples the decoder into a contextual network and mixed RL objective to encourage abstraction	CNNDM	Less readable summaries
	Song et al. (2020)	Mix-and-match strategy to generate summaries with various degree of copying levels	Gigaword, NEWSROOM	Poor performance in cross-domain settings
	Chan et al. (2021)	RL-based framework on constrained markov decision process to penalize the violation of control requirement	CNNDM, NEWSROOM	Poor performance for highly abstractive targets
Diversity	Narayan et al. (2022)	Compositional sampling decoding method	CNNDM, XSum	Generates unfaithful summaries for highly abstractive targets
	Fan et al. (2018b)	Convolutional encoder-decoder to generate stylistic summaries by adding the source prompt to the input	CNNDM, DUC2004	Repetitive and longer summaries
	Chawla et al. (2019)	RL-based method to generate formally-tailored summaries	CNNDM, Webis-TLDR-17	Poor performance in informal summary settings
	Jin et al. (2020a)	Multi-task learning framework with style-dependent layer normalization and style-guided encoder attention	NYT, CNN, Humor, Romance, Clickbait corpus	Poor performance on English Gigaword dataset
	Cao and Wang (2021)	Novel decoding methods: decode state adjustment, word unit prediction based	Hyperpartisan News detection dataset	-
	Goyal et al. (2022)	Mixture of experts strategy	CNNDM, XSum, NEWSROOM	-
	Luo et al. (2022)	Readability control of bio-medical documents	LS, PLS	Fail to handle fine-grained readability control
	Ribeiro et al. (2023)	Fine-grained readability control	CNNDM	Style insights may not generalize beyond English newswire datasets
	Wu et al. (2021)	A two-stage control generation strategy	SAMSUM	-
Coverage	Zhong et al. (2022)	Unsupervised framework to multi-granularity summary generation	Multi-NEWS, arXiv, DUC2004	Events extraction from source may effect the abstractivity
	Huang et al. (2023)	Utilize the NLI models to improve the coverage	DIALOGSUM, SAMSUM	Partially addressing the factuality problem
Role	Lin et al. (2022)	Decoders for user and agent summaries and attention divergence loss for the same topic	CSDS, MC	-
	Liang et al. (2022)	Role aware centrality scores to reweight context representations for decoding	CSDS, MC	-
Entity	Zheng et al. (2020)	Controllable neural network with guiding entities	Gigaword, DUC 2004	Performance poorer than SOTA models
	Liu and Chen (2021)	Graph convolutional network based coreference fusion layer and entity conditioned Summary Generation	SAMSUM	Paraphrasing introduces factual inconsistencies in person-specific summaries
	Hofmann-Coyle et al. (2022)	Model as a sentence selection task using transformer based biencoder with a cosine similarity based loss and adapting contrastive loss	EntSUM	-
Salience	Nallapati et al. (2017)	Summarization as a sentence selection task with salience as a feature using sequence-to-sequence model	CNNDM	Poor performance on out-of-domain datasets
	Li et al. (2018)	Key information guided network with modified attention	CNNDM	Coverage mechanism not implemented
	Deutsch and Roth (2023)	Model salience in terms of noun phrases by incorporating QA signals	CNNDM, DUC-2004	Performance relies on question generation and answering models
	Pagnoni et al. (2023)	Unsupervised pretraining involving salient sentence selection	QMSum, SQuALITY	Computationally expensive
Length	Kikuchi et al. (2016)	Remaining words provided as additional input to decoder	Gigaword	Poor performance on DUC-2004
	Fan et al. (2018b)	Convolutional encoder-decoder, summary length grouping into bins and the source document prepend with length bin's value	CNNDM	Fails to generate summaries of arbitrary lengths
	Liu et al. (2018)	Remaining number of tokens replaced by characters at the decoder	CNNDM, DMQA	Fails to generalize to new control aspects at test time
	Févry and Phang (2018)	Unsupervised denoising auto-encoder for the task of sentence compression and the decoder provided with an additional input of the remaining summary length at each time step	Gigaword	Unfaithful summary generation in some cases
	Makino et al. (2019)	Global minimum risk training optimization method under length constraint	CNNDM, Mainichi	Fails to control length
	Sarkhel et al. (2020)	Multi-level summarizer with a multi-headed attention mechanism using a series of timestep independent semantic kernels	MSR Narratives and Thinking-Machines	Fail to encode desired length
	Takase and Okazaki (2019)	Extension to the sinusoidal positional embeddings to preserve the length constraint with length-difference positional encoding and length-ratio positional encoding	JAMUS corpus (Japanese)	Poor performance when desired target length is unseen
	Yu et al. (2021)	Concatenate the length context vector with the decoder hidden state and other attention vectors	CNNDM	Incomplete shorter summary generation
	Song et al. (2021)	Confidence driven generator trained on a denoising objective with a decoder only architecture with masked source and summary tokens	Gigaword, NEWSROOM	Poor performance on large datasets
	Chan et al. (2021)	Used a reinforcement learning based Constrained Markov Decision Process to control length along with constraints on a mix of attributes such as abstractivity and covered entity	CNNDM, NEWSROOM DUC-2002	Length control only at word level
	Liu et al. (2022a)	Dynamic programming algorithm based on the Connectionist Temporal Classification model	Gigaword, DUC2004	Poor performance compared to autoregressive models
	Goyal et al. (2022)	Mixture-of-expert model with multiple decoders	CNNDM, XSum, NEWSROOM	No insights about style diversity in non-English and non-newswire datasets
	He et al. (2022)	A generic framework using keywords	CNNDM, arXiv, BIGPATENT	High reliance on the quality of extracted keywords
	Liu et al. (2022b)	Length aware attention model adapting the source encodings	CNNDM, XSum	Performance directly proportional to the summary length
	Zhong et al. (2022)	Events identification with unsupervised summary generation	GranuDUC, MultiNews, DUC2004, arXiv	Fails to capture abstractness due to event extraction
	Hyun et al. (2022)	RL based framework incorporating both the length and quality constraints in the reward function	DUC2004	Computationally expensive
	Kwon et al. (2023)	Summary length prediction task on the encoder side and encoded this information inserting a length-fusion positional encoding layer	CNNDM, NYT, WikiHow	Performance decreases with increase in summary length variance
	(Zhang et al., 2023b)	Hard prompt tuning and soft prefix tuning	CNNDM, QMSum	Low specificity in long generated summaries
Topic	Krishna and Srinivasan (2018)	RNN based attention model to generate multiple topic conditioned summaries	CNNDM	News categories provide predefined topics, limiting generalization to other tasks.
	Tan et al. (2020)	Extends topic based summarization to arbitrary topics, integrating external knowledge from ConceptNet and Wikipedia	CNNDM, MA News, All the News	-
	Suhara et al. (2020)	Framework for opinion summarization	HOTEL, Yelp	-
	Amplayo et al. (2021)	Multi-Instance Learning and a document preprocessing mechanism	SPACE, OPOSUM+	Incapable of handling unseen aspects
	Mukherjee et al. (2020)	Iterative sentence extraction algorithm	YELP	Poor performance in absence of attributes
	Mukherjee et al. (2022)	Topic-aware multimodal summarization system	MSMO	Output quality relies on data size

Table 9: CTS models descriptions and corresponding limitations.

Table 10: Master survey table (* marked fields are filled to best of our understanding based on the available information).