

Learning to Plan and Generate Text with Citations

Anonymous ACL submission

Abstract

The increasing demand for the deployment of LLMs in information-seeking scenarios has spurred efforts in creating verifiable systems, which generate responses to queries along with supporting evidence. In this paper, we explore the *attribution* capabilities of plan-based models which have been recently shown to improve the faithfulness, grounding, and controllability of generated text. We conceptualize plans as a sequence of questions which serve as *blueprints* of the generated content and its organization. We propose two attribution models that utilize different variants of blueprints, an *abstractive* model where questions are generated from scratch, and an *extractive* model where questions are copied from the input. Experiments on long-form question-answering show that planning consistently improves attribution quality. Moreover, the citations generated by blueprint models are more accurate compared to those obtained from LLM-based pipelines lacking a planning component.

1 Introduction

Large language models (LLMs) have demonstrated remarkable abilities to engage in creative conversations (Thoppilan et al., 2022; OpenAI, 2023), summarize information from contextual cues (Goyal et al., 2022; Zhang et al., 2023), and deliver zero-shot performance on a wide range of previously unseen predictive and generative tasks (Brown et al., 2020; Chowdhery et al., 2022). They are also becoming increasingly useful in information-seeking scenarios, ranging from answering simple questions (Roberts et al., 2020; Rae et al., 2021) to generating responses to search-like queries (Nakano et al., 2021; Menick et al., 2022).

The increasing demand for the deployment of LLMs in information-seeking scenarios has further spurred efforts in creating verifiable systems, which generate responses to queries along with supporting evidence. The evidence can take the form of

a URL pointing to a short segment of text which supports an answer (Bohnet et al., 2022), an attribution report with evidence snippets (Gao et al., 2022), quotes cited verbatim from pages retrieved from a search engine (Menick et al., 2022), and references to passages extracted while browsing (Nakano et al., 2021; Gao et al., 2023). In fact, this last type of evidence has been recently adopted in the form of in-line citations by commercial search engines such as BingChat¹ and perplexity.ai².

Regardless of how the evidence is presented, recent approaches tend to rely on a retrieval system (e.g., a commercial search engine) to obtain passages relevant to a query, while an LLM conditions on them to generate a response (Menick et al., 2022; Nakano et al., 2021; Bohnet et al., 2022). Other work generates an answer to the input query first and subsequently retrieves relevant evidence in a post-processing step (Bohnet et al., 2022). Alternatively, the retrieved evidence can be used to further revise the generated response rendering it more consistent with the evidence (Gao et al., 2022).

Despite recent efforts, it remains an open question how to best develop models with a *built-in* mechanism for attribution to external evidence. A related question is whether said mechanism contributes to generating more factually faithful output. Large-scale evaluation studies paint a worrying picture. Liu et al. (2023b) find that long-form responses from existing search engines frequently contain unsupported statements or inaccurate citations, while Bohnet et al. (2022) show that model performance on attribution varies greatly (between 46% and 71%) across different architectures for the simpler question answering task.

In this paper, we explore the attribution capabilities of plan-based models which have been shown to be less prone to hallucinations and more controllable (Moryossef et al., 2019; Puduppully

¹<https://bing.com/new>

²<https://perplexity.ai>

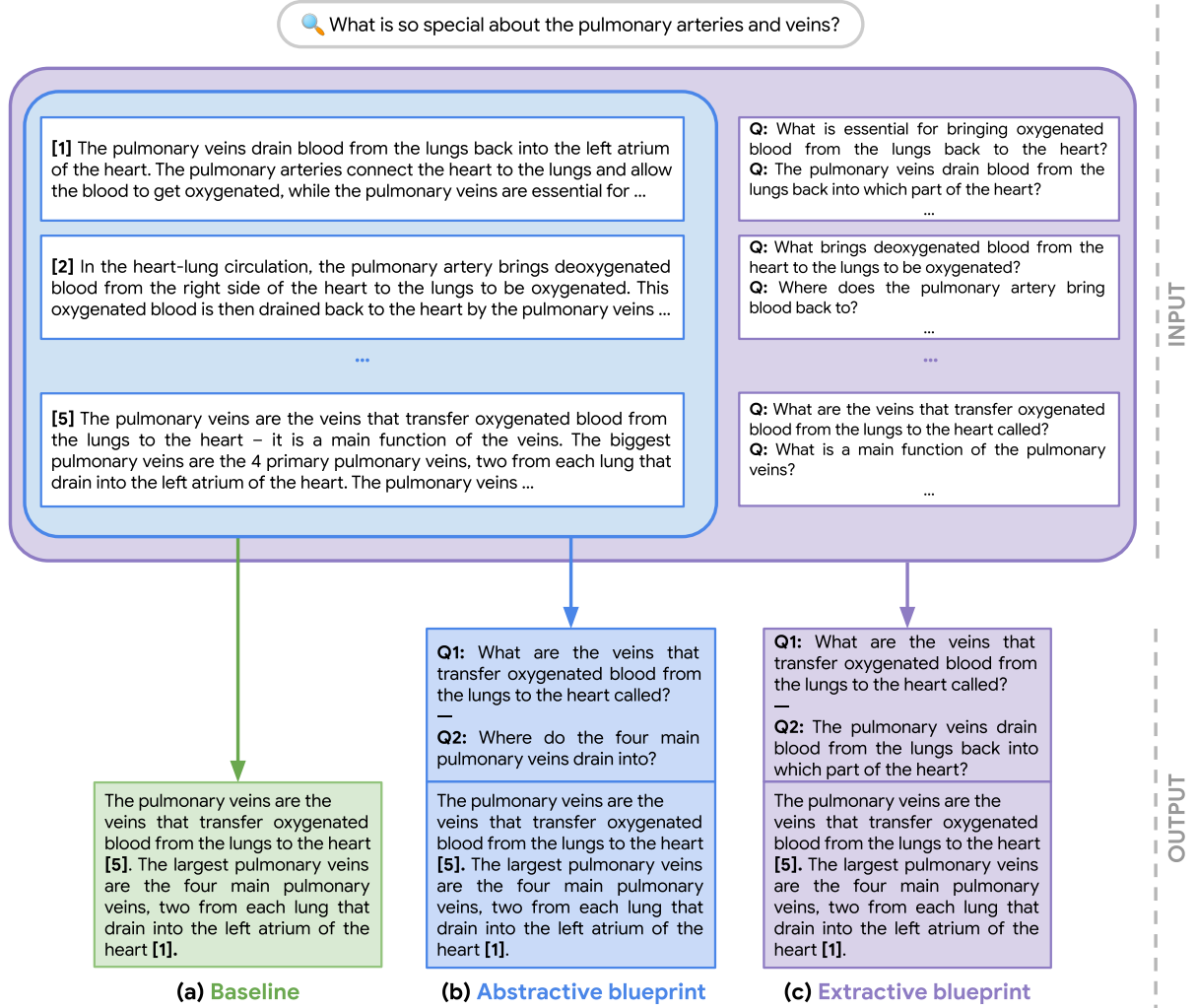


Figure 1: Query (top), followed by most relevant (abridged) passages, and summaries (bottom) with in-line citations. Summary (a) is the output of a vanilla sequence-to-sequence model trained to generate long answers with citations. Summaries (b) and (c) are the output of models with abstractive and extractive plans, respectively. Citations for plan-based models can have different formats (e.g., references to the question plan; see Section 4.2).

et al., 2019; Narayan et al., 2021, 2023; Huot et al., 2023b,a). We focus on long-form question answering (Fan et al., 2019; Vig et al., 2022; Xu and Lapata, 2022) which aims to generate summaries from a set of passages that answer a specific query. We simulate how a search engine might synthesize passages of high relevance to a user query by assuming access to a retriever, and some way of verifying the output, i.e., by citing sources (see Figure 1). Our models operate on retrieved passages and learn to plan and generate summaries with attribution. On account of being more expressive, plan-based models allow us to formalize different forms of attribution, e.g., plans can be verified via citations to passages, while summaries can be verified through citations to the plan, passages, or both.

Our models conceptualize text plans as a se-

quence of questions operating as *blueprints* for generation, determining what to say and in which order (Narayan et al., 2023; Liu et al., 2023a; Huot et al., 2023b). Questions as a planning mechanism are ideally suited to attribution, since they provide a natural link between retrieved passages and their summaries. We define two models which differ on whether the question-based plan is *generated* (see Figure 1(b)) or *copied* from input passages (Figure 1(c)) and explore whether explicit planning has any bearing on citation quality. Our contributions can be summarized as follows:

- We develop automatic methods to annotate training data with plans *and* citations, and fine-tune several Transformer models (Vaswani et al., 2017) to generate attributed text.³

³Our data and code are available from [xxx.yyy.zzz](#).

- Experimental results on the AQuAMuSe dataset (Kulkarni et al., 2020) demonstrate that plans consistently improve attribution quality. Furthermore, summary quality improves with an extractive blueprint model.
- Out-of-domain experiments on the ALCE benchmark (Gao et al., 2023) show that, once acquired, attribution is a robust skill (across information-seeking tasks). In terms of attribution quality, our models are competitive with (and sometimes better than) pipelines that heavily rely on large language models.

2 Related Work

LLMs have made significant advances in recent years in generating high quality natural language responses (Brown et al., 2020; Chowdhery et al., 2022). Despite impressive capabilities, they also demonstrate high levels of opacity: it is unclear where the information they generate comes from, which can undermine their trustworthiness (Cheng et al., 2022).

A variety of techniques aim to remedy this by building language models that provide references to support generated text. The majority of existing work has focused on models which *learn* to generate citations. For example, Nakano et al. (2021) and Menick et al. (2022) use reinforcement learning from human preferences to train language models to answer questions and provide supporting evidence (in the form of retrieved snippets). Bohnet et al. (2022) generate both an answer and a URL to a web page from which a paragraph is subsequently selected as support for the answer. They also experiment with attribution as a post-processing step where no learning is involved and use retrieval to select sources supporting their model’s output.

Along similar lines, Gao et al. (2022) propose a two-stage approach where generated text is first post-edited to be made attributable (to web content) and then revised accordingly. In addition to modeling, efforts have been also directed to the creation of evaluation protocols and benchmarks for assessing whether a statement is supported by provided evidence (Rashkin et al., 2021; Bohnet et al., 2022; Liu et al., 2023b).

Our work focuses on plan-based models which have been shown to be less opaque and more controllable (Moryossef et al., 2019; Puduppully et al., 2019; Narayan et al., 2021, 2023; Huot et al., 2023b). We draw inspiration from Narayan

et al. (2023) who formulate plans as a sequence of question-answer pairs. However, their models vary from ours in several respects: (a) they do not perform attribution or use a retrieval step to find relevant passages; (b) in our information seeking setting, we express plans as a sequence of questions only, whose answers can be found in (cited) passages; (c) we introduce extractive blueprint models wherein attribution comes for free (the questions express the content of input passages).

3 Problem Formulation

We follow a formulation of query-focused summarization common in the literature (Xu and Lapata 2020; Vig et al. 2022; *inter alia*). Let q denote an information-seeking request (e.g., “What is so special about the pulmonary arteries and veins?” in Figure 1) and $\mathcal{P} = \{p_1, \dots, p_n\}$ be a set of passages most relevant to q (see top left in Figure 1). We adopt a two-step approach where k most relevant passages are first retrieved and possibly reranked based on the query alone, and then fed to a secondary model which synthesizes the passages into a final summary $\mathcal{S} = \{s_1, \dots, s_m\}$ consisting of sentences s_i . An *attributed* summary further includes supporting evidence. For the sake of simplicity, we assume this is expressed by embedded in-line citations (see bottom left in Figure 1). Each s_i has a set of citations $\mathcal{C}_i = \{c_{i,1}, \dots, c_{i,k}\}$ where $c_{i,j}$ is a passage ID signifying that sentence i is citing said passage. For the summary (a) in Figure 1, $\mathcal{C}_1 = \{[1]\}$, and $\mathcal{C}_2 = \{[2]\}$.

In the following we first explain how to train a sequence-to-sequence model to generate citations, and then describe the proposed plan-based models.

3.1 Sequence-to-Sequence Model

Datasets for query-focused summarization such as AQuAMuSe (Kulkarni et al., 2020) or ASQA (Stelmakh et al., 2022) typically consist of queries Q , passages \mathcal{P} , and target summaries \mathcal{S} . A vanilla sequence-to-sequence model, generates summary S given passages P and query q , hence modeling the conditional probability $p(\mathcal{S}|P, q)$. In the case of attributed generation, the model is trained on summaries with citations assumed to be part of the summary’s token vocabulary. In Figure 1, the query, passages, and summary (a) with in-line references to passages [1] and [5] would constitute a training instance for this model.

In order to generate text with attribution, we

must somehow annotate target summaries with citations. In experiments, we obtain silver-standard citations, by measuring the entailment between passage $p_j \in \mathcal{P}$ as the premise and sentence $s_i \in \mathcal{S}$ as the hypothesis. We annotate S_i with $c_{i,j}$ for the passage p_j with the highest entailment score (see Section 4.1 for more details).

3.2 Blueprint Models with Attribution

Narayan et al. (2023) formalize three types of blueprint models, which are purportedly suited to different generation tasks (e.g., long vs short output/input). We repurpose their end-to-end variant as it is a straightforward extension of the standard approach described above. Our model encodes query q and passages \mathcal{P} , and then generates, $\mathcal{B}; \mathcal{S}$, the concatenation of blueprint \mathcal{B} and output summary \mathcal{S} in one decoding step. We define \mathcal{B} as a set of questions $B = \{b_1, \dots, b_k\}$ that are answered in \mathcal{S} and act as a proxy for content selection and planning. This definition is more flexible (compared to question-answer pairs), as we do not require answers to be extracted from passages, and are not limited to a specific style of answers or questions such as those manifested in SQuAD (Rajpurkar et al., 2018). Questions can be more general while answers can be abstractive, and represent a variety of syntactic constructs beyond noun phrases (e.g., verbs, clauses, even sentences).

We assume the model has access to training data consisting of queries, passages, blueprints, and summaries (See Figure 1). Blueprints are not generally available, we explain how we obtain these automatically in the next section and define two blueprint models: an abstractive model (Figure 1(b)) and an extractive model (Figure 1(c)); these differ on whether the question-based plan is generated or copied from the input.

Abstractive Model We augment the training data with blueprints by generating questions for each target summary. We adopt an overgenerate-and-rank strategy, generating multiple questions to be filtered based on their overlap with the summary.

A challenge with blueprint models is that the output may be answering a mix of general and specific questions. To avoid generating only one style of question, we employ two datasets originally developed for question answering, namely SQuAD (Rajpurkar et al., 2018) and Natural Questions (Kwiatkowski et al., 2019b). QA datasets typically consist of question/paragraph/answer triples

but can be repurposed for question generation; to produce general-purpose questions, we assume the paragraph is the input and the question is the output, while for more specific questions, we concatenate the paragraph with the sentence where the answer span is found and the model’s aim is to generate a question for the specific answer. We fine-tune T5-3B (Raffel et al., 2020) to obtain these two flavors of question generation models (see Appendix A for training details).

For each question generation model, we sample (using beam search) 10 questions per sentence and another 10 per summary so as to have a diverse question set. For each summary sentence, we then select the question with the highest lexical overlap from the set of questions specific to that sentence *and* the set of summary-level questions. Blueprint questions are sorted according to the order of appearance of their corresponding summary sentences. We add a special character in the blueprint (– in Figure 1) to delineate the questions corresponding to each sentence. We provide an example of how questions are generated and subsequently filtered in Appendix C.

Extractive Model In the extractive model, the questions align more closely with the input passages. Specifically, we explicitly verbalize the questions each passage might answer, i.e., $\mathcal{P}_{ext} = \{(p_1, Q_1), \dots, (p_n, Q_n)\}$, where p_i are retrieved passages and Q_i corresponding questions. Blueprint \mathcal{B} then consists of a subset of questions whose answers are found in the input. The extractive model is functionally equivalent to the abstractive one, the only difference being that it encodes passages \mathcal{P}_{ext} (instead of \mathcal{P}). This model also outputs $\mathcal{B}; \mathcal{S}$, however, the decoder learns to *copy* questions from the input in order to construct \mathcal{B} . The extractive model has two advantages. Firstly, it reduces the risk of hallucinations and irrelevant information since the questions are directly copied from the input. Secondly, attribution comes for free, as we know which passage gave rise to which questions (see Figure 1). An obvious drawback is computational efficiency, as questions must be generated for every input passage.

We create a large number of questions for each passage using the same generation models described above. We then greedily select five questions Q_i for each p_i such that they have highest lexical overlap with the passage and minimal overlap with each other. Subsequently, we construct the ex-

tractive blueprints under the assumption that their questions ought to be aligned to the input passages *and* output summary. We first discard questions that cannot be answered in the target summary using an answerability classifier, a T5 model trained on a repurposed version of SQuAD v2 (Rajpurkar et al., 2018) that predicts whether a question can be answered by a passage (see Appendix B for training details). The remaining questions are compared against the target summary and for each sentence we select a question based on lexical overlap. As in the abstractive model, a special character is added to mark which blueprint questions correspond to which sentences (see Appendix C for an example).

4 AQuAMuSe Experiments

4.1 Experimental Setting

Data The bulk of our experiments use AQuAMuSe (Kulkarni et al., 2020), a long-form question answering dataset, which simulates how a search engine might synthesize documents of high relevance to a user query. It consists of queries taken from Natural Questions (Kwiatkowski et al., 2019a), passages from Common Crawl, and multi-sentence summaries from Wikipedia. We use the splits released in Kulkarni et al. (2020) which contain 6,599, 714, and 849 examples for training, development, and testing, respectively. The average query/response length is 9.2/107.3 words, and the number of input passages is 10.

All our models operate on the same passages which are reranked using T5-R (Huebscher et al., 2022), a T5 11B-based encoder trained with a classification loss on the MS MARCO dataset (Nguyen et al., 2016). T5-R reorders passages $p \in \mathcal{P}$ based on their relevance with query q . Target summaries in AQuAMuSe (Kulkarni et al., 2020) were annotated with in-line citations (see summaries in Figure 1) using an entailment classifier (T5-11B; Raffel et al. 2020) fine-tuned on the ANLI dataset (Nie et al., 2020); neutral and contradicting pairs were classified as non-entailing (Honovich et al., 2022a). We consider citations part of the summary’s token vocabulary, but exclude them when evaluating summary quality.

Comparison Models We compare several models, examining how attribution and planning interact: a sequence-to-sequence baseline which does not have a plan or perform attribution (–Blueprint –Attribution); a model which generates citations

without a blueprint (–Blueprint +Attribution; Section 3.1); its mirror image, i.e., a model with a blueprint (abstractive or extractive) but no attribution (+Blueprint_{A|E} –Attribution); and a blueprint model with attribution (+Blueprint_{A|E} +Attribution). In all cases, we fine-tune a LongT5 model (Guo et al. 2022, 3B parameters).

Evaluation Metrics We evaluate the quality of the summaries and for models that output citations whether these are correct. We measure the summary’s *relevance* using ROUGE-L⁴ (Lin, 2004) which computes the longest common subsequence between output summaries and reference text. We perform this comparison after removing any citations present in the response.

We quantify the extent to which generated summaries are *faithful* to their input using textual entailment (Maynez et al., 2020; Falke et al., 2019; Narayan et al., 2022; Honovich et al., 2022b). Our classifier is trained on the Adversarial NLI dataset (Nie et al., 2020); for each sentence in the summary, we check whether it is entailed by the input passages (and report the average across all sentences to obtain an overall score). Again, we remove citations from the output to measure faithfulness.

In addition, we evaluate *citation quality* by checking whether the citations mention the *right* passages. We extend AutoAIS (Bohnet et al., 2022), a recently proposed metric for automatically measuring the attribution of short answers (in a QA setting). For long-form responses, we take sentences with in-line citations and check whether these are entailed by all the passages they cite (Amplayo et al., 2023). We then report the proportion of entailed sentences. We use the same ANLI classifier (Nie et al., 2020) described above.

For blueprint models, we would additionally expect the plan to convey information relayed by the input passages. There is no point in introducing an intermediate representation if it does not facilitate grounding to the input. We measure *blueprint quality*, by leveraging an answerability classifier (see Appendix B for details) that checks whether a question is answered by a passage. We report the proportion of questions in the blueprint that are answered by the passages it references. Questions in the blueprint cite passages *implicitly* through their respective summary sentences. Recall that by construction there is a one-to-one correspondence between blueprint questions and summary

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

Retrieval-based Models		ROUGE-L	ANLI
–Blueprint	–Attribution	63.80	87.20
+Blueprint _A	–Attribution	64.17	87.88
+Blueprint _E	–Attribution	72.25	87.92
–Blueprint	+Attribution	63.72	86.73
+Blueprint _A	+Attribution	63.49	88.05
+Blueprint _E	+Attribution	72.98	88.01

Table 1: AQuAMuSe results: summary quality.

Retrieval-based Models		Answerability	AutoAIS
–Blueprint	+Attribution	—	72.64
+Blueprint _A	+Attribution	92.27	74.16
+Blueprint _E	+Attribution	97.97	74.35

Table 2: AQuAMuSe results: attribution quality.

sentences (which in turn cite input passages).

4.2 Results

The extractive blueprint performs best across evaluation metrics. Table 1 presents our results on AQuAMuSe when evaluating summary relevance and faithfulness. In the first block we compare models that do not generate citations. As can be seen plan-based models deliver better summaries in terms of ROUGE-L, with slightly enhanced faithfulness (see ANLI column), with the extractive blueprint model taking the lead. Our blueprint models are overall better compared to Narayan et al. (2023), achieving a new state of the art in the AQuAMuSe dataset. However, as mentioned earlier, the two approaches are not directly comparable as they differ in terms of the input they expect and blueprint definition. In the second block of Table 1 we compare models with an attribution mechanism. In general, faithfulness increases for plan-based models and relevance improves by 10 points for the extractive blueprint model.

In Table 2 we evaluate attribution quality for models which output citations. We observe gains in AutoAIS for our blueprint models over performing attribution without a plan. Attribution quality is slightly higher for the extractive model, which also creates more grounded blueprints, 97.97% of the questions are answered by the passages it references. We show example output in Appendix H.

The extractive blueprint generates the most abstractive summaries. We investigate the hypothesis that blueprint models might perform better

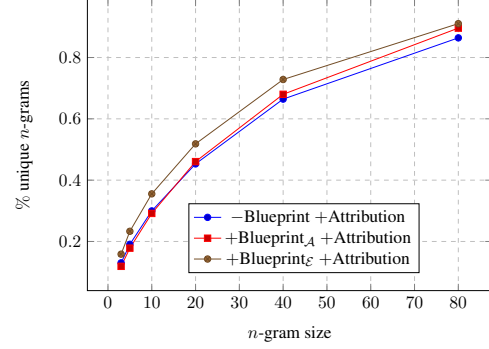


Figure 2: Unique n -grams in generated summary.

simply because it is easier for them to copy sentences or passages from the input. In Figure 2 we evaluate the proportion of unique n -grams in the generated summary (compared to the input passages) for different n -gram sizes (where $n = \{3, 5, 10, 20, 40, 80\}$) for the +Attribution models in Table 1. As can be seen, the extractive model generates the most abstractive responses, while the –Blueprint baseline copies longer chunks of text from the input. This suggests that blueprints help models consolidate information from multiple sources, while still being faithful to the input.

Blueprint plans allow us to define different citation formats. Our experiments have so far adopted in-line citations as a common format for all models. We now explore alternatives which are unique to our plan-based models and examine whether these have any bearing on performance. Specifically, we (a) augment in-line passage citations with references to blueprint questions; (b) place in-line citations in the blueprint (after each question) rather than in the output summary; and (c) have implicit citations for extractive models only; in this case, we do not generate any citations but attribution can be inferred on the basis of the questions being aligned to summary sentences and copied from the input passages. Figure 3 shows examples and definitions of these formats.

Table 3 reports results on plan quality (checking whether the questions can be answered by the cited passages) and citation correctness (AutoAIS). We see that formats assuming an explicit connection between passages and questions improve answerability, however, at the expense of AutoAIS. Overall, in-line citations strike the best balance between answerability and faithfulness. We further examine the quality of blueprint plans in Appendix D.

Human evaluation corroborates automatic attribution results. In addition to automatic eval-

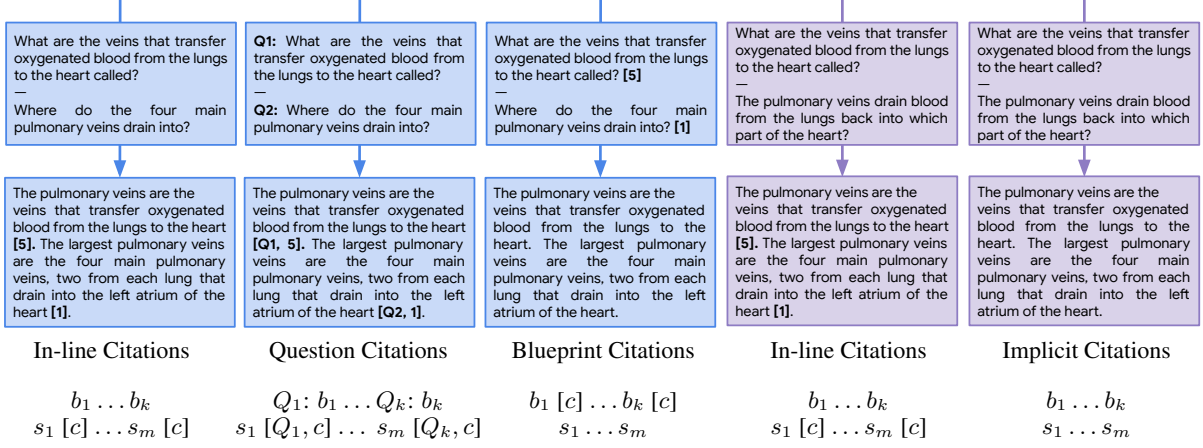


Figure 3: Blueprints (top), corresponding summaries (bottom), and different citation formats: b_i is a blueprint question, s_i is a summary sentence and c , Q are citations. Abstractive/extractive blueprints are colored in blue/purple.

+Blueprint _A +Attribution	Answer	AutoAIS
In-line Citations	92.27	74.16
Question Citations	92.95	60.86
Blueprint Citations	93.67	61.49
+Blueprint _E +Attribution	Answer	AutoAIS
In-line Citations	97.97	74.35
Implicit Citations	98.20	63.79

Table 3: Results for different citation formats.

	AQuAMuSe	FSupp	PSupp	NSupp	Contra
–Blueprint +Attribution	46.91*	5.21	46.91*	0.69	
+Blueprint _A +Attribution	49.41*	5.62	43.64	1.33	
+Blueprint _E +Attribution	58.28	4.10	36.90	0.71	
	ELI5 & ASQA	FSupp	PSupp	NSupp	Contra
Vicuna-13B	43.97*	1.88	54.02*	0.13	
ChatGPT (Vanilla)	76.51*	1.95	20.21*	1.33	
+Blueprint _A +Attribution	88.11	3.50	7.51	0.88	

Table 4: Proportion of sentences with citations Fully Supported by the referenced passages, Partially Supported, Not Supported and in Contradiction. Systems marked with * are significantly different the model in **bold** ($p < 0.01$, using bootstrap resampling).

uation, we conducted a human study to assess the attribution quality of our models. Participants were asked (see Appendix E for details on our instructions and study) to verify whether the citations in the generated summaries are correct, classifying each sentence as fully supported (FSupp) by the referenced passages, partially supported (PSupp), unsupported (NSupp), or contradictory (Contra). We elicited judgments on 50 instances (randomly sampled from the AQUAMuse test set) for the systems shown in Table 4 (first block). We find that the extractive system is perceived as the most accurate among attribution models achieving the highest proportion of correct citations.

5 ALCE Experiments

5.1 Experimental Setting

Data We further examine whether the attribution capabilities of our models transfer to other domains, by performing (zero-shot) experiments on a recently collated benchmark called ALCE (Gao et al., 2023). ALCE contains evaluation sets from ASQA (Stelmakh et al., 2022), which focuses on ambiguous queries and long-form answers, ELI5

(Fan et al., 2019), which focuses on how/why/what questions that require explanations, and QAMPARI (Rubin et al., 2022) where the answers are lists of entities. These datasets were selected such that the answers are factual, long-form, and require information from multiple sources. In experiments, we use the passages provided by the ALCE benchmark with no further reranking. Details on ALCE and representative examples are in Appendix G.

Comparison Models Gao et al. (2023) use ALCE to examine the citation abilities of LLMs. They report experiments with two base LLMs, namely ChatGPT (gpt-3.5-turbo-0301) and LLaMA (Touvron et al., 2023), both of which were prepended with detailed instructions and a few demonstrations for in-context learning. Their ChatGPT pipelines include (a) Vanilla, an end-to-end model where the input is the query and top-5 passages; (b) Summarization condenses the passages first with a separate LLM and thus can process more of them; (c) Snippet extracts snippets

from passages using an LLM; (d) Inline Search, where instead of providing passages in the input, the LLM calls search in-line whenever needed; and (e) Closed Book, where no passages are provided and attribution happens post-hoc (i.e., for each sentence, find the best passage to cite). LLaMA-based pipelines are versions of the Vanilla approach described above with either (f) LLaMA-13B, or (g) Vicuna-13B (Chiang et al., 2023).⁵ We compare our abstractive model (+Blueprint_A+Attribution) against these pipelines. We assume that a question generation model is not available at transfer time to obtain questions for input passages (and create blueprints for the extractive model).

Evaluation Metrics We follow Gao et al. (2023) and evaluate system output in terms of *correctness* and *attribution*. Correctness is measured differently for each dataset; for ASQA, they calculate the recall of correct short answers (Stelmakh et al., 2022); for ELI5, they use InstructGPT (text-davinci-003; Ouyang et al., 2022) to generate sub-claims from the target response, and an NLI model (TRUE; Honovich et al., 2022a) to predict whether the response entails each sub-claim; for QAMPARI, they calculate precision and recall by checking the exact match to the gold-standard list of answers, considering recall to be 100% if the prediction includes at least five correct answers (Rubin et al., 2022); here, we only use recall to measure correctness since precision is calculated over a list of predicted answers and our models output *one* long answer (given the AQuAMuSe training). Attribution is evaluated using citation recall, which determines if the output is entirely supported by cited passages, and citation precision, which identifies any irrelevant citations. We report F1 for metrics that have both precision and recall.

5.2 Results

The abstractive blueprint competes with LLMs in attribution. Our results are summarized in Table 5. On average, we observe that our model is better than competing systems in terms of attribution. We suspect this is due to the fact that it has been explicitly fine-tuned with outputs that have citations, which in turn improves attribution, rendering the model robust across datasets. In terms of correctness, our model performs on average better

⁵We did not include systems that additionally use response reranking and interaction strategies as these are orthogonal improvements and do not consistently improve performance.

		ASQA		ELI5		QAMPARI		Average	
		C	A	C	A	C	A	C	A
ChatGPT	Vanilla (5-psg)	40.4	73.0	12.0	50.5	20.8	20.7	24.4	48.1
	Summarization (10-psg)	43.3	65.2	12.3	49.8	22.3	24.6	26.0	46.5
	Snippet (10-psg)	41.4	61.1	14.3	47.5	22.9	23.9	26.2	44.2
	Inline Search	32.4	58.2	18.6	44.6	18.7	14.9	23.2	39.3
	Closed Book	38.3	26.7	18.6	15.5	24.7	10.0	27.2	17.4
	LLaMA-13B	26.9	12.6	3.9	3.9	9.4	6.9	13.4	7.8
Vicuna-13B		31.9	50.6	10.0	17.4	14.9	12.9	18.9	27.0
LongT5 3B (10-psg) +Blueprint _A +Attribution		33.8	77.8	5.2	60.9	12.9	10.8	17.3	49.8

Table 5: Results on ALCE benchmark: LLM pipelines and proposed abstractive blueprint model.

than LLaMA-13B but worse than the other systems, which is not surprising as it is applied to new datasets without access to in-domain training. In contrast, comparison LLM pipelines have been exposed to few-shot demonstrations.

Humans find Blueprint citations more accurate compared to those produced by LLMs. To further validate the results in Table 5, we conducted human evaluation comparing ChatGPT Vanilla, Vicuna-13B, and our abstractive blueprint model. Using the same instructions from Section 4.2, we collected judgments over 100 randomly sampled instances (50 from ELI5 and 50 from ASQA).⁶ As shown in Table 4, our model is perceived to have the most accurate citations, significantly outperforming both Vicuna-13B and ChatGPT Vanilla.

6 Conclusion

In this paper we focused on retrieval augmented generation with citations. We explored the attribution capabilities of plan-based models and proposed different variants depending on how the plans are created (i.e., abstractively or extractively). Our experiments revealed several interesting findings. Summary quality improves for blueprint models when these learn to generate citations, both in terms of relevance and faithfulness. Conversely, this is not the case for models without an intermediate planning stage, in fact attribution slightly hurts performance (see second line in Table 1). The attribution mechanism is also important; formulating the attribution as in-line citations seems beneficial as well as a tight alignment between the input and the blueprint questions. Finally, attribution is a transferable skill for blueprint models, across datasets and information-seeking tasks.

⁶For Gao’s (2023) systems, we use the output from https://github.com/princeton-nlp/ALCE/tree/main/human_eval. Appendix F shows results for each dataset.

Limitations

In this work, we have focused exclusively on encoder-decoder models. Future work could explore fine-tuning (on attribution annotations) a wider range of models, including decoder-only ones. In the future we would further like to explore more efficient ways of indexing passages with questions, as well as using the questions as an additional training signal for improving information retrieval. Finally, an important consideration with generative models is the problem of misinformation. While the work we present here makes a step towards improving the faithfulness and factual consistency of summarization systems, it is important to note that our models still make mistakes, in particular in out-of-domain scenarios.

References

- Reinald Kim Amplayo, Peter J Liu, Yao Zhao, and Shashi Narayan. 2023. [SMART: Sentences as basic units for text evaluation](#). In *The Eleventh International Conference on Learning Representations*.
- Bernd Bohnet, Vinh Q Tran, Pat Verga, Roei Aharoni, Daniel Andor, Livio Baldini Soares, Jacob Eisenstein, Kuzman Ganchev, Jonathan Herzig, Kai Hui, et al. 2022. Attributed question answering: Evaluation and modeling for attributed large language models. *arXiv preprint arXiv:2212.08037*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Ruijia Cheng, Alison Smith-Renner, Ke Zhang, Joel Tetreault, and Alejandro Jaimes-Larrarte. 2022. [Mapping the design space of human-AI interaction in text summarization](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 431–455, Seattle, United States. Association for Computational Linguistics.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. [Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality](#).
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*.
- Tobias Falke, Leonardo F. R. Ribeiro, Prasetya Ajie Utama, Ido Dagan, and Iryna Gurevych. 2019. [Ranking generated summaries by correctness: An interesting but challenging application for natural language inference](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2214–2220, Florence, Italy. Association for Computational Linguistics.
- Angela Fan, Yacine Jernite, Ethan Perez, David Grangier, Jason Weston, and Michael Auli. 2019. [ELI5: Long form question answering](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3558–3567, Florence, Italy. Association for Computational Linguistics.
- Luyu Gao, Zhuyun Dai, Panupong Pasupat, Anthony Chen, Arun Tejasvi Chaganty, Yicheng Fan, Vincent Y Zhao, Ni Lao, Hongrae Lee, Da-Cheng Juan, et al. 2022. Attributed text generation via post-hoc research and revision. *arXiv preprint arXiv:2210.08726*.
- Shuai Gao. 2022. [Système de traduction automatique neuronale français-mongol \(historique, mise en place et évaluations\) \(French-Mongolian neural machine translation system \(history, implementation, and evaluations\) machine translation \(hereafter abbreviated MT\) is currently undergoing rapid development, during which less-resourced languages nevertheless seem to be less developed\)](#). In *Actes de la 29e Conférence sur le Traitement Automatique des Langues Naturelles. Volume 2 : 24e Rencontres Etudiants Chercheurs en Informatique pour le TAL (RECITAL)*, pages 97–110, Avignon, France. ATALA.
- Tianyu Gao, Howard Yen, Jiatong Yu, and Danqi Chen. 2023. Enabling large language models to generate text with citations. *arXiv preprint arXiv:2305.14627*.
- Tanya Goyal, Junyi Jessy Li, and Greg Durrett. 2022. [News summarization and evaluation in the era of gpt-3](#).
- Mandy Guo, Joshua Ainslie, David Uthus, Santiago Ontanon, Jianmo Ni, Yun-Hsuan Sung, and Yinfei Yang. 2022. [LongT5: Efficient text-to-text transformer for long sequences](#). In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 724–736, Seattle, United States. Association for Computational Linguistics.

Or Honovich, Roei Aharoni, Jonathan Herzig, Hagai Taitelbaum, Doron Kukliansy, Vered Cohen, Thomas Scialom, Idan Szpektor, Avinatan Hassidim, and Yossi Matias. 2022a. TRUE: Re-evaluating factual consistency evaluation . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 3905–3920, Seattle, United States. Association for Computational Linguistics.	778
Or Honovich, Roei Aharoni, Jonathan Herzig, Hagai Taitelbaum, Doron Kukliansy, Vered Cohen, Thomas Scialom, Idan Szpektor, Avinatan Hassidim, and Yossi Matias. 2022b. TRUE: Re-evaluating factual consistency evaluation . In <i>Proceedings of the Second DialDoc Workshop on Document-grounded Dialogue and Conversational Question Answering</i> , pages 161–175, Dublin, Ireland. Association for Computational Linguistics.	779
Michelle Chen Huebscher, Christian Buck, Massimiliano Ciaramita, and Sascha Rothe. 2022. Zero-shot retrieval with search agents and hybrid environments. <i>arXiv preprint arXiv:2209.15469</i> .	780
Fantine Huot, Joshua Maynez, Chris Alberti, Reinald Kim Amplayo, Priyanka Agrawal, Constanza Fierro, Shashi Narayan, and Mirella Lapata. 2023a. μplan: Summarizing using a content plan as cross-lingual bridge . <i>arXiv preprint arXiv:2305.14205</i> .	781
Fantine Huot, Joshua Maynez, Shashi Narayan, Reinald Kim Amplayo, Kuzman Ganchev, Annie Priyadarshini Louis, Anders Sandholm, Dipanjan Das, and Mirella Lapata. 2023b. Text-blueprint: An interactive platform for plan-based conditional generation . In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations</i> , pages 105–116, Dubrovnik, Croatia. Association for Computational Linguistics.	782
Sayali Kulkarni, Sheide Chammas, Wan Zhu, Fei Sha, and Eugene Ie. 2020. Aquamuse: Automatically generating datasets for query-based multi-document summarization. <i>arXiv preprint arXiv:2010.12694</i> .	783
Tom Kwiakowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. 2019a. Natural questions: a benchmark for question answering research. <i>Transactions of the Association for Computational Linguistics</i> , 7:453–466.	784
Tom Kwiakowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Matthew Kelcey, Jacob Devlin, Kenton Lee, Kristina N. Toutanova, Llion Jones, Ming-Wei Chang, Andrew Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019b. Natural questions: a benchmark for question answering research. <i>Transactions of the Association of Computational Linguistics</i> .	785
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.	786
Danyang Liu, Mirella Lapata, and Frank Keller. 2023a. Visual storytelling with question-answer plans . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 5800–5813, Singapore. Association for Computational Linguistics.	787
Nelson F. Liu, Tianyi Zhang, and Percy Liang. 2023b. Evaluating verifiability in generative search engines .	788
Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. On faithfulness and factuality in abstractive summarization . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 1906–1919, Online. Association for Computational Linguistics.	789
Jacob Menick, Maja Trebacz, Vladimir Mikulik, John Aslanides, Francis Song, Martin Chadwick, Mia Glaese, Susannah Young, Lucy Campbell-Gillingham, Geoffrey Irving, et al. 2022. Teaching language models to support answers with verified quotes. <i>arXiv preprint arXiv:2203.11147</i> .	790
Sewon Min, Julian Michael, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2020. AmbigQA: Answering ambiguous open-domain questions . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 5783–5797, Online. Association for Computational Linguistics.	791
Amit Moryossef, Yoav Goldberg, and Ido Dagan. 2019. Step-by-step: Separating planning from realization in neural data-to-text generation . 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 2267–2277, Minneapolis, Minnesota. Association for Computational Linguistics.	792
Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, Xu Jiang, Karl Cobbe, Tyna Eloundou, Gretchen Krueger, Kevin Button, Matthew Knight, Benjamin Chess, and John Schulman. 2021. WebGPT: Browser-assisted question-answering with human feedback . <i>CoRR</i> , abs/2112.09332.	793
Shashi Narayan, Joshua Maynez, Reinald Kim Amplayo, Kuzman Ganchev, Annie Louis, Fantine Huot, Anders Sandholm, Dipanjan Das, and Mirella Lapata. 2023. Conditional Generation with a Question-Answering Blueprint . <i>Transactions of the Association for Computational Linguistics</i> , 11:974–996.	794
Shashi Narayan, Gonalo Simões, Yao Zhao, Joshua Maynez, Dipanjan Das, Michael Collins, and	795

834	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.	890
835		891
836		892
837		
838		893
839		894
840	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.	895
841		896
842		897
843		
844		898
845	Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016. Ms marco: A human generated machine reading comprehension dataset .	899
846		900
847		901
848		902
849	Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. 2020. Adversarial NLI: A new benchmark for natural language understanding . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4885–4901, Online. Association for Computational Linguistics.	903
850		904
851		905
852		906
853		907
854		908
855		909
856	OpenAI. 2023. Gpt-4 technical report .	
857	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Gray, et al. 2022. Training language models to follow instructions with human feedback. In <i>Advances in Neural Information Processing Systems</i> .	910
858		911
859		912
860		913
861		914
862		915
863	Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Dmytro Okhonko, Samuel Broscheit, Gautier Izacard, Patrick S. H. Lewis, Barlas Oguz, Edouard Grave, Wen-tau Yih, and Sebastian Riedel. 2021. The web is your oyster - knowledge-intensive NLP against a very large web corpus . <i>CoRR</i> , abs/2112.09924.	916
864		
865		917
866		918
867		919
868		920
869		921
870	Ratish Puduppully, Li Dong, and Mirella Lapata. 2019. Data-to-text generation with content selection and planning. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , pages 6908–6915.	922
871		923
872		924
873		925
874	Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susanah Young, et al. 2021. Scaling language models: Methods, analysis & insights from training gopher. <i>arXiv preprint arXiv:2112.11446</i> .	926
875		927
876		928
877		929
878		930
879		931
880	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.	932
881		933
882		934
883		935
884		936
885		937
886	Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know what you don’t know: Unanswerable questions for SQuAD . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 784–789, Melbourne, Australia. Association for Computational Linguistics.	938
887		939
888		940
889		941
		942
		943
		944
		945
		946
		947

H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 5998–6008. Curran Associates, Inc.

Jesse Vig, Alexander Fabbri, Wojciech Kryscinski, Chien-Sheng Wu, and Wenhao Liu. 2022. [Exploring neural models for query-focused summarization](#). In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 1455–1468, Seattle, United States. Association for Computational Linguistics.

Yumo Xu and Mirella Lapata. 2020. [Coarse-to-fine query focused multi-document summarization](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3632–3645, Online. Association for Computational Linguistics.

Yumo Xu and Mirella Lapata. 2022. [Document summarization with latent queries](#). *Transactions of the Association for Computational Linguistics*, 10:623–638.

Tianyi Zhang, Faisal Ladhak, Esin Durmus, Percy Liang, Kathleen McKeown, and Tatsunori B. Hashimoto. 2023. [Benchmarking large language models for news summarization](#).

A Question Generation Implementation

Two question generation models were used to obtain questions for input passages and blueprints: one model was trained on SQuAD v2 (Rajpurkar et al., 2018) and another one on Natural Questions (Kwiatkowski et al., 2019b). Both question generators were built on top of a T5 3B (Raffel et al., 2020) base model, and expect either of the following input formats:

1. General-Purpose QG, where the model generates a question that can be answered anywhere in the passage:

```
Generate      question      »>
[PASSAGE]
```

2. Sentence-Specific QG, where the model generates a question specific to a sentence in the passage:

```
Generate      question      »>
[PASSAGE] » [SENTENCE]
```

where [PASSAGE] and [SENTENCE] are the passage and the sentence where the answer to the question (to-be generated) is found.

B Answerability Classifier Implementation

Our answerability classifier repurposes the SQuAD v2 (Rajpurkar et al., 2018) dataset that includes unanswerable questions. Specifically, a T5 11B model is fine-tuned to accept as input a question and a passage in the following format:

```
question:      [question]  context:
[passage]
```

and to generate Yes as output in cases where [question] can be answered given [passage] as context and No otherwise.

C Blueprint Selection Examples

C.1 Abstractive Blueprints

Table 6 shows a small-scale example of how questions are selected in abstractive blueprints. Given the summary in the example, we generate three SQuAD-style questions (marked with subscript Sq) and three NQ-style questions (marked with subscript N) for each summary sentence. In addition, we generate three SQuAD and three NQ summary-level questions, giving us a total of 24

questions. Note that in our experiments, we use 10 SQuAD/NQ questions instead of three. For each sentence, we then select the question with the highest lexical overlap from the union of the sentence-specific and summary-level questions. In our example, the question for S_1 is selected from the summary-level set. Finally, we collate these questions together with diacritic – to create the blueprint.

C.2 Extractive Blueprints

Table 7 shows a small-scale example of how questions are selected in extractive blueprints. For each of the five input passages (10 in our experiments), we generate the same amount of questions using the method described above. That is, for each passage we generate sentence-level and summary-level questions (similar to what is shown in Table 6) using both the SQuAD and the NQ question generation models. We then greedily select five questions for each passage, such that they have the highest lexical overlap with the passage and minimal overlap with each other. Once we have the passage questions, we are ready to select the questions in the blueprint. We first filter out questions that cannot be answered using an answerability classifier (see Appendix B and AF column in the table). We then select, for each sentence in the summary, the question from with the highest lexical overlap (SS column in the table). Finally, we collated these questions together with diacritic – to create the blueprint.

D Grounding and Controllability

We further examine whether the *planning* aspect of the blueprint models remains intact when they are also expected to generate summaries with references to the input. We measure the extent to which the output is grounded to the blueprint plan using the answerability classifier. Specifically, we check if the questions in the blueprint are answerable by the summary as a whole and by individual sentences (and report the average). Table 8 shows our results for blueprint models with and without attribution. As can be seen, extractive blueprint models are less grounded compared to their abstractive counterparts, even though the former are better at attribution (see Table 2).

We further illustrate how to improve the attribution quality of abstractive models, thereby showcasing the controllability aspect of blueprint models.

We simply remove (post-hoc) generated blueprint questions that are unanswerable based on the input passages, and then generate the summary using the filtered blueprint. This post-processing substantially improves attribution by 7.09 AutoAIS points (compared to Table 2) while retaining or marginally improving output quality on other dimensions (63.12 ROUGE-L, 89.45 ANLI, 93.08, Answerability, and 81.25 AutoAIS).

E Experimental Instructions

Figure 4 presents the experimental instructions used in our human elicitation study. To ensure high quality annotations, we created a screener test to determine the raters’ suitability for the task. We also limited the maximum tasks per annotator per study to six to avoid annotation fatigue. A total of 108 annotators participated in the elicitation study for AQuAMuSe, and a total of 88 and 249 annotators participated in the studies for ASQA and ELI5, respectively. The annotator agreement scores are as follows: 69.87% on AQuAMuSe, 93.00% on ASQA, and 89.04% on ELI5. Our annotators were paid adequately by our suppliers adhering to the supplier code of conduct.

F Human Evaluation Results on ASQA and ELI5

Table 10 reports our human evaluation results for ASQA and ELI5, respectively. Results mostly align with the combined results in Table 4, showing that our model significantly produces more accurate citations.

G The ALCE Benchmark

ALCE (Gao et al., 2023) is a collection of datasets aimed at evaluating LLM citation capabilities. They contain factual questions which require long-form answers aggregating information over multiple sources (i.e., 100-word passages). We summarize various statistics on these datasets in Table 9 and describe them below.

ASQA (Stelmakh et al., 2022) focuses on ambiguous questions which have multiple interpretations (see Table 11 for an example). The questions were taken from AmbigQA (Min et al., 2020), while long-form answers were crowdsourced by synthesizing information from multiple documents. For ALCE (Gao et al., 2023), ASQA questions were additionally paired with Wikipedia passages

Instructions

In this task, you will be provided with a long Machine-Generated Response to a user's question(s) along with several Target Sentences and their cited Evidence. The Target Sentences should be based on information found within the cited Evidence, however it will often either misrepresent the information or will provide additional information not found in the Evidence.

The purpose of this task is to detect those misrepresentations and/or the additional information as you compare the Target Sentence to the Evidence. For task purposes, consider the Evidence reliable. Your job is to read through the Response and then rate the task sentence-by-sentence:

1. **To what extent can the Target Sentence be verified by its Evidence?**
 - Only use the cited Evidence.
 - Look for information in the highlighted sentence that is "**Unsupported**" (not contained in the evidence) or "**Contradictory**" (misrepresents the evidence).
 - The Target Sentence can be "**Fully**" supported by its Evidence even if
 - the sentence rephrases and/or combines information from several pieces of Evidence;
 - the sentence leaves out some information from its Evidence.
2. **Mark whether there is irrelevant Evidence cited by the Target Sentence.**
 - An irrelevant piece of Evidence does not support or contradict the Target Sentence in any way.

Note that:

- The two questions above are purely about a single, highlighted Target Sentence and the Evidence it cites each time. The other sentences in the Response should be used as context to understand the Target Sentence.
- The Evidence may be in a language different from that of the Target Sentence. Please evaluate normally in this case if possible, otherwise select "Unclear".
- The Evidence may contain symbolic math or chemistry that is similar, but not identical to that of the Target Sentence. Please navigate to the blue link associated with the evidence to evaluate for correctness and equivalence.

Task

Machine-Generated Response:

Multilingual speakers have acquired and maintained at least one language during childhood, the so-called first language. [2] The first language (sometimes also referred to as the mother tongue) is acquired without formal education, by mechanisms heavily disputed. [2] Children acquiring two languages in this way are called simultaneous bilinguals. [2] Even in the case of simultaneous bilinguals, one language usually dominates the other. [1] People who know more than one language have been reported to have better skills at language learning compared to monolinguals, and have important economic advantages in their professional career over monolingual individuals. [2]

(1) Target Sentence and Evidence

Multilingual speakers have acquired and maintained at least one language during childhood, the so-called first language.

Evidence:

[2] You learned (or better: you acquired) your mother tongue or your first language as a child without formal teaching. You memorized sound and grammar so deeply, that you master the language automatically. When children acquire two languages they are "simultaneous bilinguals", generally known as bilingual kids.

Question. To what extent can the **Target Sentence** be verified by its **Evidence**?

- **Contradictory:** Some information misrepresents the cited Evidence.
- **Unsupported:** Some important information is not contained in the cited Evidence; no contradiction.
- **Partially:** All important information is supported by the cited Evidence (though some trivial/ minor information is not supported); no contradiction.
- **Fully:** Everything in the Target Sentence can be verified by the cited Evidence.
- **Unclear:** The Target Sentence or Evidence is hard to understand, or can't understand the language, or can't make a choice.

(2) Target Sentence and Evidence

...

Figure 4: Experimental instructions presented to participants during the human elicitation study. The question repeats for each sentence in the machine-generated response.

Retrieval-based Models	Summary Sentences	
+Blueprint _A − Attribution	96.67	96.49
+Blueprint _E − Attribution	91.32	80.06
+Blueprint _A + Attribution	97.46	93.64
+Blueprint _E + Attribution	91.22	79.97

Table 8: Grounding results for Blueprint Models (with and without attribution).

(2018-12-20 snapshot) which purportedly contained the answers.

QAMPARI (Rubin et al., 2022) is an open-domain QA dataset where answers are lists of entities, drawn from different passages. All questions in QAMPARI have at least 5 answers, with an average of 13 answers. Multi-answer questions were automatically generated using manually defined templates, answers were collated from Wikipedia, and examples were verified and paraphrased by crowdworkers (see Table 11 for an example). For ALCE (Gao et al., 2023), QAMPARI questions were also paired with Wikipedia passages (2018-12-20 snapshot).

ELI5 (Fan et al., 2019) mostly consists of how/why/what questions that require in-depth long answers and multiple passages as evidence. Questions (and answers) were elicited from the subreddit *Explain Like I’m Five* (ELI5) where users are encouraged to provide answers which are comprehensible by a five year old (see Table 11). For ALCE (Gao et al., 2023), ELI5 questions were paired with passages from Sphere (Piktus et al., 2021), a filtered version of Common Crawl.

ALCE (Gao et al., 2023) contains 1,000 randomly selected examples from each dataset (development set). It does not provide training data as it is aimed at assessing the citation capabilities of LLMs.

H Example Output

Tables 12, 14, and 16 present output examples for the baseline system (−Blueprint +Attribution) and our two blueprint variants (+Blueprint_{A|E} +Attribution). Blueprints created by the plan-based systems are respectively shown in Tables 13, 15, and 17. In general, we observe that the baseline system considers a smaller number of input passages when generating its response, while plan-based summaries contain more diverse passage references. As far as the blueprints are concerned, we find that the abstractive system generates (on average) more

questions compared to the extractive model and the reference. Although the questions making up the blueprints of the two models are not identical, we see that they cover similar topics. For instance, in Table 13 the blueprints emphasize topics like the mean temperatures of a tropical monsoon climate, its precipitation, dry seasons, and variance in temperature.

Dataset	Question Type	#passages	Question Length	Answer Length
ASQA	Factoid (ambiguous)	Wikipedia (21M)	12.5	65.7
QAMPARI	Factoid (list)	Wikipedia (21M)	10.5	13.0
ELI5	Why/How/What	Sphere (899M)	17.4	120.6

Table 9: Datasets in the ALCE benchmark (Gao et al., 2023). Question and Answer Length refer to the average number of words; for QAMPARI, Answer Length is the average number of answers. Examples are shown in Table 11.

ASQA	FSupp	PSupp	NSupp	Contra
Vicuna-13B	57.87*	3.40	38.30*	0.43
ChatGPT (Vanilla)	78.02*	1.61	18.77*	1.61
+ Blueprint_A + Attribution	88.58	2.75	7.40	1.27
ELI5	FSupp	PSupp	NSupp	Contra
Vicuna-13B	37.57*	1.17	61.25*	0.00*
ChatGPT (Vanilla)	75.76*	2.12	20.93*	1.19
+ Blueprint_A + Attribution	87.42	4.60	7.67	0.31

Table 10: Proportion of sentences with citations Fully Supported by the referenced passages, Partially Supported, Not Supported, and in Contradiction. Systems marked with * are significantly different from the model in **bold** ($p < 0.01$, using bootstrap resampling).

ASQA
<p>Q: Who played Bill Weasley in Harry Potter and the Prisoner of Azkaban?</p> <p>A: The Weasley family is a family of wizards including Molly and Arthur Weasley as well as their seven children. Richard Fish appeared as Bill briefly in the film adaptation of Harry Potter and the Prisoner of Azkaban. Domhnall Gleeson, the son of actor Brendan Gleeson (Alastor Moody in the series), played Bill Weasley in Harry Potter and the Deathly Hallows. Chris Rankin played Percy Weasley and Fred and George Weasley were played by identical twins James Phelps and Oliver Phelps. Rupert Grint played Ron Weasley.</p>
QAMPARI
<p>Q: Who won a Coke Zero Sugar 400 competition?</p> <p>A: Sam McQuagg, Bobby Allison, Fireball Roberts, Jimmie Johnson, Tony Stewart, Tim Richmond, Jeff Gordon, David Ragan, Erik Jones, Tony Stewart, Jamie McMurray, Kyle Busch, Cale Yarborough, Dale Earnhardt Jr., Brad Keselowski, David Pearson, Dale Jarrett, Jeff Gordon, Tony Stewart, Greg Biffle, A. J. Foyt, Richard Petty, Justin Haley, Tony Stewart.</p>
ELI5
<p>Q: How does so much of our trash end up in the ocean?</p> <p>A: Because water flows downhill and very often ends up in rivers which very often end up in oceans. So when it rains, trash is washed downhill and into streams and rivers and ultimately the ocean.</p>

Table 11: Examples of question-answer pairs from ALCE datasets; input passages are omitted for the sake of brevity.

describe the seasonal patterns of the tropical savanna biome	
[0]	However, unlike most tropical climates, a sizeable amount of precipitation falls during the dry season. In essence, this version of the monsoon climate generally has less pronounced dry seasons than tropical savanna climates . . .
[1]	Tropical monsoon climates have monthly mean temperatures above 18 °C in every month of the year and feature wet and dry seasons, as Tropical savanna climates do. Tropical monsoon climates however features its driest month seeing on average less than 60 mm and this latter fact is in direct contrast to a tropical savanna climate . . .
[2]	64. Tropical wet and dry or savanna climate (Aw) have a pronounced dry season, with the driest month having precipitation less than 60 mm and less than 1/25 of the total annual precipitation. Example(s) of areas with this climate is(are) . . .
[3]	Tropical monsoon climate, occasionally also known as a tropical wet climate or tropical monsoon and trade-wind littoral climate in climate classification, is a relatively rare type of climate. Tropical monsoon climates have monthly mean temperatures above 18 °C (64 °FF) . . .
[4]	The second letter indicates the precipitation pattern—w indicates dry winters (driest winter month average precipitation less than one-tenth wettest summer month average precipitation. s indicates at least three times as much rain in the wettest month of winter as in the driest month of summer . . .
–Blueprint +Attribution	
Tropical monsoon climate, occasionally also known as a tropical wet climate or tropical monsoon and trade-wind littoral climate in climate classification, is a relatively rare type of climate. [3] Tropical monsoon climates have monthly mean temperatures above 18 °C (64 °F) in every month of the year and feature wet and dry seasons, as Tropical savanna climates do. [3] Unlike tropical savanna climates however, a tropical monsoon climate’s driest month sees less than 60 mm (2.36 in) of precipitation. [3] Also a tropical monsoon climate tends to see less variance in temperatures during the course of the year than a tropical savanna climate. [3] This climate has a driest month which nearly always occurs at or soon after the winter solstice for that side of the equator. [3]	
+Blueprint _A +Attribution	
Tropical monsoon climate, occasionally also known as a tropical wet climate or tropical monsoon and trade-wind littoral climate in climate classification, is a relatively rare type of climate. [3] Tropical monsoon climates have monthly mean temperatures above 18 °C (64 °F) in every month of the year and feature wet and dry seasons, as tropical savanna climates do. [3] Unlike tropical savanna climates however, a tropical monsoon climate’s driest month sees less than 60 mm (2.36 in) of precipitation. [3] Additionally, a tropical monsoon climate tends to see less variance in temperatures during the course of the year than a tropical savanna climate. [1] This climate has a driest month which nearly always occurs at or soon after the "winter" solstice for that side of the equator. [3] There are generally two versions of a tropical monsoon climate: Less pronounced dry seasons and Extraordinarily rainy wet seasons and pronounced dry seasons. [0] Regions with this variation of the monsoon climate typically see copious amounts of rain during the wet season. [1]	
+Blueprint _E +Attribution	
Tropical monsoon climate, occasionally also known as a tropical wet climate or tropical monsoon and trade-wind littoral climate in climate classification, is a relatively rare type of climate. [3] Tropical monsoon climates have monthly mean temperatures above 18 °C (64 °F) in every month of the year and feature wet and dry seasons, as tropical savanna climates do. [1] Unlike tropical savanna climates however, a tropical monsoon climate’s driest month sees less than 60 mm (2.36 in) of precipitation. [3] Also a tropical monsoon climate tends to see less variance in temperatures during the course of the year than a tropical savanna climate. [3] This climate has a driest month which nearly always occurs at or soon after the "winter" solstice for that side of the equator. [3]	
Reference Summary	
Tropical savanna climate or tropical wet and dry climate is a type of climate that corresponds to the Köppen climate classification categories "Aw" and "As". [9] Tropical savanna climates have monthly mean temperatures above 18 °C (64 °F) in every month of the year and typically a pronounced dry season, with the driest month having precipitation less than 60 mm and also less than 100 – (total annual precipitation (mm) / 25) of precipitation. [9] This latter fact is in direct contrast to a tropical monsoon climate, whose driest month sees less than 60 mm of precipitation but has more than 100 – (total annual precipitation (mm) / 25) of precipitation. In essence, a tropical savanna climate tends to either see less rainfall than a tropical monsoon climate or have more pronounced dry seasons. [0]	

Table 12: Example responses to an AQuAMuSe query (top) for a sequence-to-sequence model which does not include plans, two blueprint models, (abstractive and extractive), and gold standard summary. The second block in the table shows the 5 best retrieved passages in abridged form (only the first sentence is given). Passage questions are also omitted for the sake of brevity.

Abstractive Blueprint
Q: what is a tropical monsoon climate occasionally also known as?
Q: what are the monthly mean temperatures of a tropical monsoon climate?
Q: what is the average precipitation in a tropical monsoon climate's driest month?
Q: a tropical monsoon climate tends to see less variance in temperatures during the course of the year than what other type of climate?
Q: what month nearly always occurs at or soon after the winter solstice for that side of the equator?
Q: how many versions of a tropical monsoon climate are there?
Q: what do regions with this variation of the monsoon climate typically see during the wet season?
Extractive Blueprint
Q: what type of climate has monthly mean temperatures above 18 °C (64 °F) in every month of the year?
Q: what type of climate sees less than 60 mm of precipitation in its driest month?
Q: a monsoon climate tends to see more of what than a tropical savanna climate?
Q: what type of climate has less pronounced dry seasons than tropical monsoon climates?
Q: what do regions with less pronounced dry seasons typically see during the wet season?
Reference Blueprint
Q: what is the driest month in a tropical wet and dry or savanna climate?
Q: what type of climate has monthly mean temperatures above 18 °C (64 °F) in every month of the year?
Q: a tropical monsoon climate's driest month sees less than what amount of precipitation?
Q: what do tropical wet and dry or savanna climates have a pronounced dry season?

Table 13: Abstractive and extractive blueprints for model summaries shown in Table 12. For comparison, we also include the blueprint of the reference summary.

difference between regular file and pipe file in unix
[0] Twtower, I think I know what's the reason of your problem. In old Linux (a.k.a. Old Kerbel version) named pipes weren't created by mkfifo command . . .
[1] It is not possible to change the backup destination for third-party tools. For this destination, the only objects created in the file system are named pipes. Named pipes occupy no space in the file system . . .
[2] The recovery process ends successfully, so I suppose that system correctly read all necessary recovery parameters from backup catalog. Still, it is confusing looking at the wrong log backup location in generated SQL recovery statement. . . .
[3] SAP HANA supports certified third party backup tools via the backint interface. Backups are transferred via pipe from the SAP HANA database to the 3rd party backup agent, which runs on the SAP HANA database server and then sends the backups to the 3rd party backup server. . . .
[4] This project shows you how to implement your own on-line radio receiver in Linux, using as major basis the mplayer application . . .
–Blueprint +Attribution
In Unix-like operating systems, a file is a regular file, while a pipe is a named pipe. [0]
+Blueprint _A +Attribution
Named pipes are a special type of file system object in Unix-like computer systems. [0] All Unix-like computer systems include at some point named pipes. [0]
+Blueprint _E +Attribution
In Unix-like operating systems, a file is either a regular file or a named pipe. [0] Named pipes are special files which occupy no space in the file system and can be created and destroyed like any other file. [1]
Reference Summary
In Unix-like operating systems, a file is a regular file, while a pipe is a named pipe. [0]

Table 14: Example responses to an AQUAMuSe query (top) for a sequence-to-sequence model which does not include plans, two blueprint models, (abstractive and extractive), and gold standard summary. The second block in the table shows the 5 best retrieved passages in abridged form (only the first couple of sentences are given). Passage questions are also omitted for the sake of brevity.

Abstractive Blueprint
Q: what are named pipes a special type of in unix-like computer systems?
Q: what do all unix-like computer systems include at some point?
Extractive Blueprint
Q: what type of files are created when a backup is made to a pipe?
Q: named pipes occupy how much space in the file system?
Reference Blueprint
Q: named pipes are special files that can exist where on the file system?
Q: named pipe special files are made with the command mkfifo as in what ?

Table 15: Abstractive and extractive blueprints for model summaries shown in Table 14. For comparison, we also include the blueprint of the reference summary.

briefly explain html. what are the essential features of html
[0] HTML is Hyper Text Markup Language – the standard markup language for creating web applications and pages, and a cornerstone of internet programming. HTML denotes the building blocks of a web page, such as text, paragraphs, hyperlinks, images, and other items. . .
[1] HTML: Hyper Text Markup Language is the standard markup language used to create webpages and web applications. HTML represents the structure of a web page with semantic information, images, interactive forms and other objects. . .
[2] HyperText Markup Language or simply HTML is the markup language used by web developers to create web pages. It is written in the format of HTML elements containing tags bracketed in angle brackets like <html> . . .
[3] HTML: HyperText Markup Language, commonly referred to as HTML, is the standard markup language used to create web pages. HTML is a standardized system for tagging text files to achieve font, color, graphic, and hyperlink effects on World Wide Web pages. . .
[4] HTML - HyperText Markup Language, commonly referred to as HTML, is the standard markup language used to create web pages. Web browsers can read HTML files and render them into visible or audible web pages. HTML describes the structure of a website semantically along with cues for. . .
+Blueprint _A +Attribution
HyperText Markup Language, commonly referred to as HTML, is the standard markup language used to create web pages. [4] Web browsers can read HTML files and render them into visible or audible web pages. [3] Web browsers do not display HTML tags and scripts, but use them to interpret the content of the page. [2] HTML describes the structure of a website semantically along with cues for presentation, making it a markup language, rather than a programming language. [4] HTML elements form the structure blocks of all web pages. [2] It allows images and objects to be embedded and can be used to create interactive forms. [2]
+Blueprint _E +Attribution
HTML is HyperText Markup Language – a standard markup computer language for representing the contents of a page of hypertext. [0] Web browsers receive HTML documents from a server and render them into visible or audible web pages. [4] HTML describes the structure of a web page semantically and originally included cues for the appearance of the document. [8]
Reference Summary
Hypertext Markup Language (HTML) is the standard markup language for creating web pages and web applications. [8] With Cascading Style Sheets (CSS) and JavaScript it forms a triad of cornerstone technologies for the World Wide Web. [5] Web browsers receive HTML documents from a web server or from local storage and render them into multimedia web pages. [2] HTML describes the structure of a web page semantically and originally included cues for the appearance of the document. [8]

Table 16: Example responses to an AQuAMuSe query (top) for a sequence-to-sequence model which does not include plans, two blueprint models, (abstractive and extractive), and gold standard summary. The second block in the table shows the 5 best retrieved passages in abridged form (only the first couple of sentences are given). Passage questions are also omitted for the sake of brevity.

Abstractive Blueprint
Q: what is the standard markup language used to create web pages?
Q: what can web browsers read and render into visible or audible web pages?
Q: what do web browsers use to interpret the content of a page?
Q: what describes the structure of a website semantically along with cues for presentation?
Q: what do html elements form the structure blocks of?
Q: what can be embedded in html and used to create interactive forms?
Extractive Blueprint
Q: what is the standard markup language for creating web applications and pages?
Q: what is a standardized system for tagging text files to achieve font, color, graphic, and hyperlink effects on world wide web pages?
Q: web browsers receive html documents from a server and render them into what?
Q: html describes the structure of a website semantically along with cues for what?
Reference Blueprint
Q: what is the standard markup language for creating web pages and web applications?
Q: along with html and css, what is a cornerstone technology for the world wide web?
Q: web browsers receive html documents from where?
Q: html originally included cues for the appearance of what?

Table 17: Abstractive and extractive blueprints for model summaries shown in Table 16. For comparison, we also include the blueprint of the reference summary.