

# Streamlined Summarization: Pointer-Generator Networks in Focus for Efficient Content Condensation

*Nirmal Koshy*

## Table

1) Introduction.....	3
2) Abstract.....	4
3) Model Overview.....	5
3.1) Fundamental Sequence-to-Sequence Model.....	5
3.2) Pointer Generator Network.....	7
3.3) Coverage Mechanism.....	8
4) Dataset Used.....	11
5) Experiment.....	11
6) Result.....	12
7) Conclusion.....	13
8) References.....	14

## 1. Introduction

Abstractive text summarization represents a paradigm shift in natural language processing, aiming to distil the core essence of a document into a concise summary through interpretation and paraphrasing rather than mere extraction and reordering of sentences. Among the cutting-edge methodologies within this domain, pointer-generator networks have emerged as a pivotal innovation, bridging the gap between the advantages of extractive and abstractive techniques.

At its essence, abstractive text summarization using pointer-generator networks involves the creation of summaries that go beyond the mere replication or rearrangement of sentences, instead, they generate novel sentences that capture the essential information while maintaining coherence and readability. Unlike traditional extractive methods, which may struggle with understanding context and generating fluent summaries, abstractive approaches have the potential to produce more human-like summaries by synthesizing information from the source text in a more flexible and creative manner.

Pointer-generator networks leverage a sophisticated architecture that combines the strengths of both abstractive and extractive approaches. These networks typically consist of an encoder-decoder framework augmented with attention mechanisms, allowing the model to effectively capture and comprehend the salient features of the source text. Importantly, pointer-generator networks introduce mechanisms that dynamically determine whether to generate new words or directly copy words from the source text, thus providing the model with the flexibility to balance between fidelity to the original content and the creation of novel, informative summaries.

The integration of pointer mechanisms into the abstractive summarization process offers several key advantages. Firstly, it enables the model to accurately reproduce factual details from the source text, addressing one of the primary challenges faced by purely abstractive models. Additionally, pointer-generator networks mitigate the risk of generating repetitive or irrelevant content by allowing the model to selectively copy information from the source text as needed, thereby enhancing the coherence and relevance of the generated summaries.

## 2. Abstract

In the landscape of abstractive text summarization, traditional approaches have often struggled with accurately reproducing factual details from the source text and avoiding repetitive content in generated summaries. This paper presents a pioneering architecture that seeks to address these inherent challenges within neural sequence-to-sequence models. The proposed model introduces two key advancements to enhance the summarization process. Firstly, a hybrid pointer-generator network is incorporated, allowing the model to directly copy words from the source text using pointing mechanisms. This innovative approach aims to improve the accuracy of information reproduction while still preserving the model's ability to generate novel words through the generator component. By leveraging this hybrid approach, the model can better capture essential details from the source material, thereby enhancing the overall quality of the generated summaries. Secondly, coverage mechanisms are integrated into the model to monitor the summarization progress and mitigate the issue of repetitive content commonly observed in existing approaches. These mechanisms enable the model to keep track of which parts of the source text have been summarized, thereby reducing redundancy and ensuring a more coherent and concise summary output. Empirical evaluation on the challenging CNN/Daily Mail summarization task demonstrates the effectiveness of the proposed architecture, showcasing a notable improvement over the current state-of-the-art methods by a substantial margin of at least 2 ROUGE points. This research not only pushes the boundaries of abstractive summarization techniques but also equips them with practical solutions for addressing issues of factual accuracy and content repetition. Overall, this work represents a significant advancement in the field, making text summarization smarter, more refined, and more adept at capturing the essence of the source material.

### 3. Model Overview

This segment provides an overview of our **fundamental sequence-to-sequence model**, our **pointer-generator model**, and the **coverage mechanism** that can be implemented into either of the preceding models. Our abstractive summarization is achieved using these three models one by one where each model is designed to address specific challenges or enhance certain aspects of performance in the overall task. The combination of these models may result in a system that is more robust, context-aware, and capable of handling various nuances in the input data.

#### 3.1) Fundamental Sequence-To-Sequence Model

In a sequence-to-sequence (seq2seq) model, repetition of words can indeed occur, particularly if the model isn't explicitly guided to avoid it. This repetition might stem from various factors such as the absence of constraints during training, characteristics of the training data itself, or the inherent architecture of the model, especially if it employs autoregressive decoding mechanisms. Our foundational model closely follows the structure introduced by Nallapati et al. (2016) which is clearly depicted in **Figure 1**. The individual tokens of the article are sequentially inputted into the encoder, implemented as a single-layer bidirectional LSTM. This process generates a sequence of encoder hidden states. Encoder hidden states are the intermediate representations of the source sentence generated by the encoder component of the model. Each element in the sequence is processed, and the encoder hidden states capture the information and context associated with that specific position in the sequence. In the context of recurrent neural networks (RNNs) or LSTMs, each hidden state is computed based on the input at the corresponding position and the previous hidden state. At each time step, the decoder, a single-layer unidirectional LSTM, takes as input the word embedding of the preceding word. During training, this is the prior word from the reference summary and during testing, it corresponds to the previous word generated by the decoder. The decoder maintains a state. The computation of the attention distribution at a given time step follows the methodology outlined by Bahdanau et al. (2015). Formula is expressed in the form below:

$$a_i^t = \frac{e^{t_i}}{\sum_{j=1}^n e^{t_j}}$$

Where:

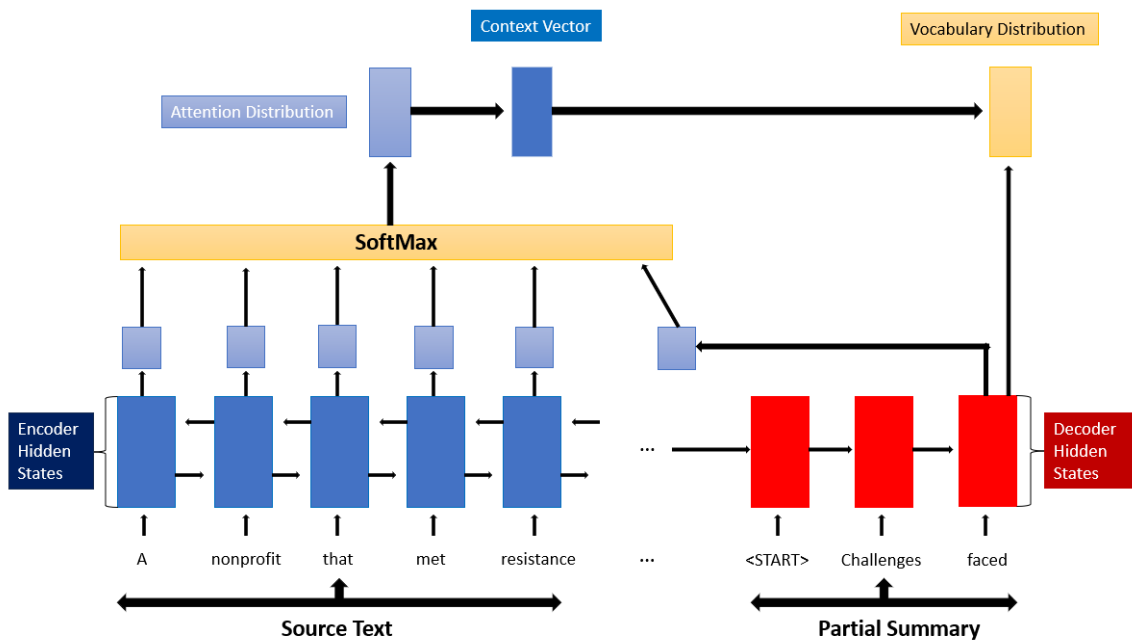
- $a_i^t$  represents the  $i$ -th element of the attention distribution at time step  $t$ .
- $e^{t_i}$  is the  $i$ -th element of the vector  $e^t$ , which is calculated as  $v^t * \tanh(x * \omega h h_i * \omega_{ss_t} * b_{attn})$ .
- The denominator  $\sum_{j=1}^n e^{t_j}$  ensures that the attention distribution sums to 1, as it is computed using the SoftMax function.

The attention distribution serves as a probability distribution across the source words, guiding the decoder on where to focus in generating the subsequent word. Subsequently, this attention distribution is utilized to generate a weighted sum of the encoder hidden states. The weighted sum, often denoted as the context vector  $h_t^*$ , is calculated by taking the sum of the encoder hidden states, each multiplied by its corresponding attention weight. Mathematically, it can be expressed as:

$$h_t^* = a_1^t h_1 + a_2^t h_2 + \dots + a_n^t h_n$$

Where:

- $a_i^t$  is the attention weight associated with the  $i$ -th encoder hidden state at time step  $t$ .
- $h_n$  is the number of hidden states of the encoder.



**Figure 1: Fundamental Sequence-To-Sequence Model with Attention Mechanism**

The context vector, acting as a compact representation of the information extracted from the source at the current step, is combined with the decoder layer state. This combined representation undergoes processing through two linear layers, resulting in the generation of the vocabulary distribution denoted as  $P_{Vocab}$ .

$$P_{Vocab} = \text{SoftMax} (v_0 * (\text{concat} (S_t, h_t^*)) + b_0)$$

Where:

- $v_0$  is a weight matrix for the linear layer.
- $S_t$  is the decoder state.
- $\text{concat} (S_t, h_t^*)$  represents the concatenation of the decoder state and the context vector.

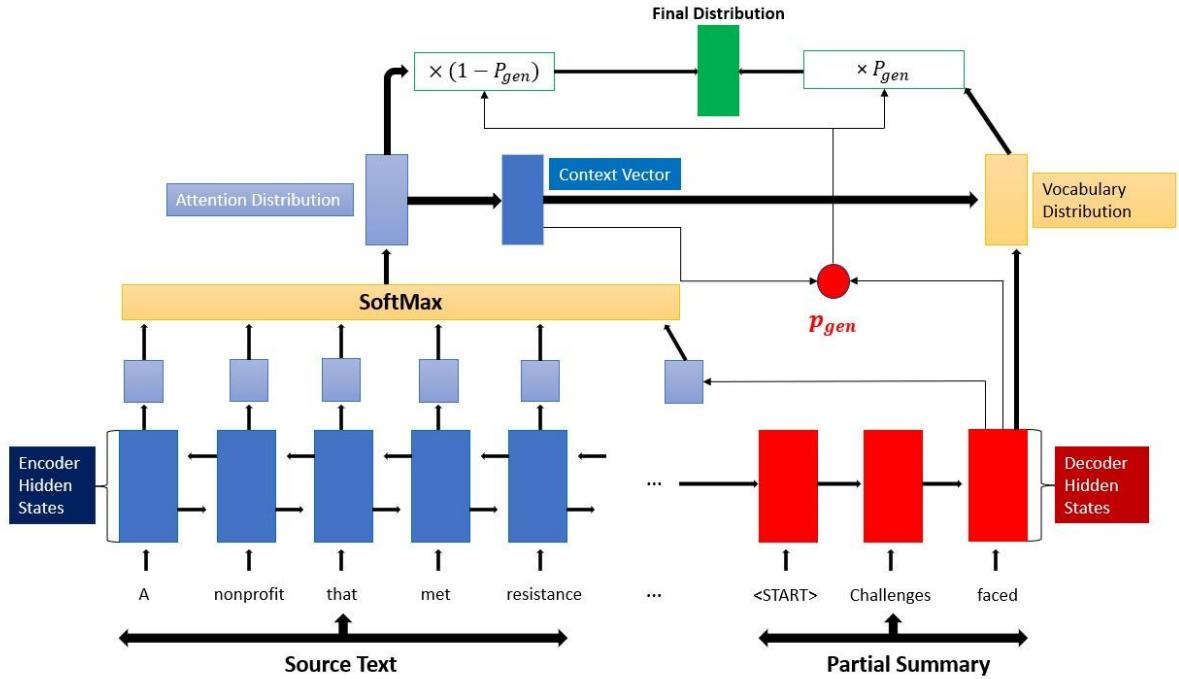
- $b_0$  is the bias vector. The bias vector allows the model to learn an offset or a shift in the output space.

Throughout the training process, the loss at a given time step  $t$  is determined by the negative logarithm of the probability assigned to the correct target word for that particular time step. The cumulative or overall loss for the entire sequence is then computed by summing the individual losses across all time steps which is given below:

$$\text{loss} = \frac{1}{T} \sum_{t=0}^T \text{loss}_t$$

### 3.2) Pointer Generator Network

The pointer-generator network we've developed is an amalgamation of our baseline model and a pointer network. This unique architecture enables our model to seamlessly switch between copying words directly from the input sequence through pointing and generating words from a predefined vocabulary.



**Figure 2: Pointer Generator Network**

In our pointer-generator model, illustrated in **Figure 2**, we calculate the attention distribution and the context vector following the procedures outlined in Section 3.1. Beyond that, we introduce a generation probability for each timestep. This probability is computed based on the context vector, the decoder state, and the decoder input. It is expressed as:

$$P_{gen} = \sigma(w_T h_t^* + \omega_T s_t + \omega_T x_t + b_{ptr})$$

Where:

- $w_T h_t^*$  is dot product of the weight vector associated with the context vector  $h_t^*$ . The subscript  $T$  denotes the transpose of the weight vector, and the dot product involves multiplying corresponding elements of  $w$  and  $h_t^*$  and summing up the results.
- $w_T s_t$  involves the dot product of the weight vector  $w$  associated with the decoder state  $s_t$ .
- $w_T x_t$  is similar to the previous terms, this involves the dot product of the weight vector  $w$  associated with the decoder input  $x_t$ .
- $b_{ptr}$  is the bias term associated with the pointer mechanism. It's a constant value that is added to the weighted sum of the other terms.

For each document, let the extended vocabulary be the union of the vocabulary and all words appearing in the source document. We define the probability distribution over this extended vocabulary as:

$$P(w) = p_{gen} * P_{Vocab}(w) + (1 - p_{gen}) \sum_{i:w_i^a = t_i} a_{t_i}$$

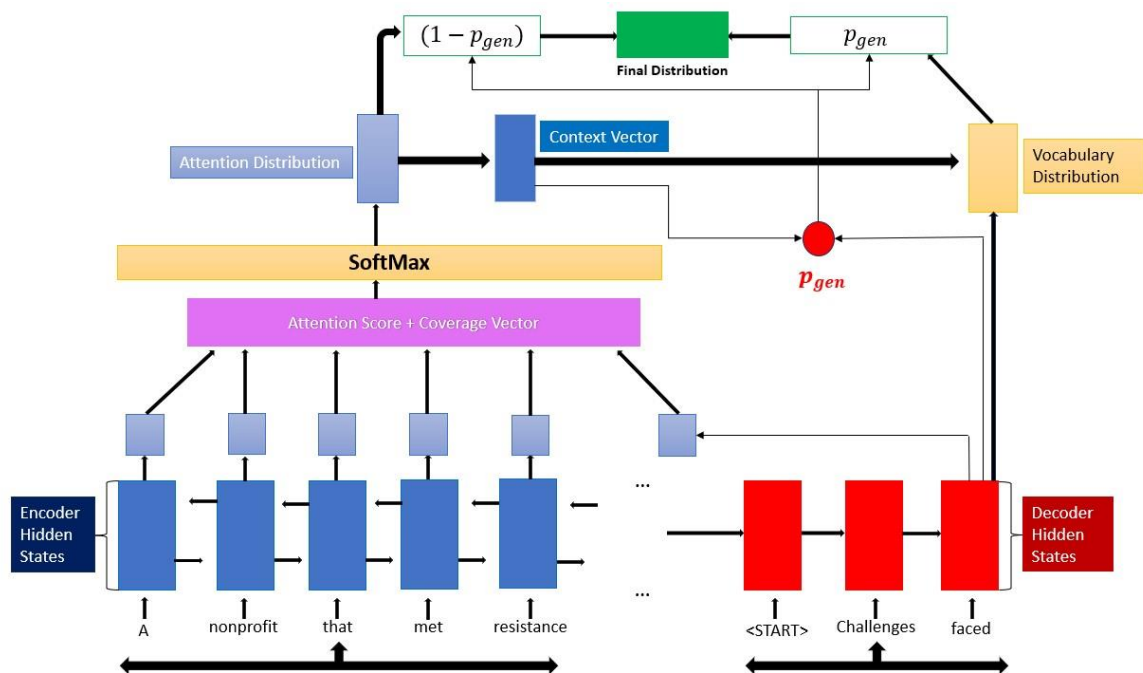
In this equation,  $P_{Vocab}(w)$  represents the probability of encountering word  $w$  in the fixed vocabulary, and  $\sum_{i:w_i^a = t_i} a_{t_i}$  denotes the sum of attention weights assigned to positions where word  $w$  appears in the source document. The described probability distribution and its associated model play a crucial role in addressing challenges related to vocabulary limitations. The probability distribution  $P(w)$  is tailored to handle out-of-vocabulary (OOV) words, a common issue in traditional models with fixed vocabularies. The pointer-generator model elegantly resolves this challenge by incorporating an attention mechanism, enabling the generation of OOV words. The model showcases adaptability by dynamically adjusting its probability distribution based on the vocabulary found in the source document. This feature proves invaluable when dealing with diverse content, allowing the model to effectively navigate variations in word usage across different documents. Utilizing the parameter  $p_{gen}$ , the model gains flexibility in word generation. It can intelligently decide whether to generate words from a fixed vocabulary or to point to words in the source document. This adaptability enhances the model's ability to produce coherent and contextually relevant summaries. In contrast to traditional models constrained by fixed vocabularies, the pointer-generator model triumphs over such limitations. By dynamically incorporating words from the source document, it ensures a more accurate representation of diverse content, showcasing the model's versatility in capturing the intricacies of different textual contexts.

### 3.3) Coverage Mechanism

The coverage mechanism is an important component in pointer-generator networks, designed to address the issue of repetitiveness in abstractive text summarization. In abstractive summarization, a model generates a concise summary that may not be a verbatim extraction from the source text. Repetitiveness is a common challenge in this task, where the model tends to repeatedly select and use the same words or phrases in the



generated summary. The coverage mechanism helps mitigate this problem by keeping track of the attention weights assigned to different positions in the source document over time. The coverage mechanism interacts with the attention mechanism by augmenting the computation of attention scores at each decoding step. Initially, the attention mechanism calculates attention scores based on the decoder's current state and the encoder's hidden states, indicating the relevance of each encoder state to the decoding process. The coverage mechanism introduces a coverage vector that tracks the attention history, representing which parts of the input sequence have been attended to. At each decoding step, this coverage vector is incorporated into the computation of attention scores, typically by adding it to the attention scores obtained from the attention mechanism. The resulting combined scores reflect both the relevance of the encoder states and the coverage of the input sequence. Subsequently, these combined scores are passed through a SoftMax function to normalize them into a probability distribution, known as the attention distribution as shown in **Figure 3**. This distribution represents the attention weights assigned to each encoder state, taking into account both the current relevance and the coverage of the input sequence.



**Figure 3: Pointer Generator Network with Coverage Mechanism**

The idea is to encourage the model to attend to parts of the source document that have not been attended to before. This is done by introducing a coverage vector, which is updated at each decoding step. Plus, remember the coverage vector  $c^0$  is initialized as a zero vector. This is because, at the onset of the decoding process during the first timestep, none of the positions in the source document have been covered. Consequently,  $c^0$  serves as the starting point, representing the absence of cumulative attention on any position in the input sequence. As the decoding progresses, the coverage vector  $c^t$  accumulates the attention weights assigned to different positions, capturing the evolving coverage status and influencing the model's attention pattern during subsequent timesteps. In our coverage

model, the coverage vector is calculated as the cumulative sum of attention distributions from all preceding decoder timesteps:

$$c^t = \sum_{t'=0}^{t-1} a^{t'}$$

Where:

- $c^t$  is the coverage vector at timestep  $t$ . This is a vector of values that keeps track of the cumulative attention assigned to each position in the input sequence up to the current decoding step.
- $a^{t'}$  signifies the attention weight assigned to the position in the input sequence at timestep  $t'$ . In the context of a sequence-to-sequence model, attention weights are computed for each position in the input sequence at each decoding step.

The coverage vector is combined with attention mechanism, changing the equation of sequence-to-sequence model to the one displayed below:

$$e_i^t = v^T \tanh(\omega_h h_i + w_s s_t + \omega_c c_i^t + b_{attn})$$

Where:

- $e_i^t$  represents the attention score for the  $i$ -th element at time step  $t$ . It is computed as a weighted sum of the input elements, with weights determined by the attention mechanism.
- $h_i$  represents the hidden state at position  $i$ . In the context of an LSTM (Long Short-Term Memory) or a similar recurrent neural network (RNN),  $h_i$  is the output or hidden state of the network at position  $i$ .
- $s_t$  represents the cell state at time step  $t$ . In an LSTM, the cell state is updated based on the previous cell state, the input, and the forget and input gates.
- $c_i^t$  represents the context or memory at position  $i$  at time step  $t$ . In the context of an LSTM, it is the cell state  $c_i$  at position  $i$  at time  $t$ .
- All the  $w$ 's are weight matrices associated with the hidden state, context, and input, respectively. They are learnable parameters that the model optimizes during training.
- $b_{attn}$  is a bias term associated with the attention mechanism.
- $\tanh$  is the hyperbolic tangent activation function. It squashes the input values between -1 and 1, introducing non-linearity to the model.

We find it imperative to introduce an additional definition, namely a coverage loss. This is essential to impose a penalty for recurrently attending to identical locations. This strategic addition becomes crucial to address a recurring issue where the model tends to concentrate on identical positions within the input sequence. The coverage loss serves as a regularization term integrated into the training process, aiming to penalize the model for repeatedly attending to the same locations. By doing so, we fix the risk of the model overlooking crucial information and generating repetitive outputs. This regularization mechanism encourages a more balanced distribution of attention across the entire input sequence, promoting a

comprehensive understanding and preventing the model from sticking onto specific positions throughout the generation process. It is expressed as:

$$covloss_t = \sum_i \min(a_i^t, c_i^t)$$

Our loss function exhibits greater flexibility, particularly tailored for summarization tasks where uniform coverage is not a strict requirement. Instead of penalizing uniformity, we focus on mitigating overlap between each attention distribution and the coverage accumulated thus far, effectively preventing the model from repeatedly attending to the same locations. The coverage loss, adjusted by a hyperparameter  $\lambda$ , is then integrated into the primary loss function. This results in the creation of a novel composite loss function, providing a balanced approach that caters specifically to the nuances of summarization requirements:

$$loss_t = -\log(Pw_t^*) + \lambda \sum_i \min(a_i^t, c_i^t)$$

#### 4. Dataset Used

The dataset utilized for this research is the CNN/Daily Mail dataset, as introduced by Hermann et al. in 2015 and further refined by Nallapati et al. in 2016. This dataset comprises a large collection of news articles paired with corresponding human-generated summaries, sourced from the websites of CNN and the Daily Mail. With its diverse range of topics and high-quality summaries, the CNN/Daily Mail dataset serves as a standard benchmark for evaluating abstractive text summarization models. By leveraging this widely used dataset, our research ensures consistency and comparability with existing literature, facilitating robust evaluation and benchmarking of our proposed approach against state-of-the-art methods.

#### 5. Experiment

In our experiments, we maintain a consistent model configuration with hidden states of 150 dimensions and word embeddings of 100 dimensions across all models. When employing pointer-generator models, we set 64000 rows of trainset and 3200 for validation set. It's worth noting that the pointer network's capability to address out-of-vocabulary (OOV) words allows us to utilize a smaller vocabulary size compared to prior works, such as Nallapati et al. (2016), which used vocabularies of 150,000 for the source and 60,000 for the target. Additionally, for the seq2seq model, we experiment with trainset containing 50k articles. Unlike Nallapati et al.'s approach (2016), we eschew pre-training the word embeddings—opting instead to learn them from scratch during training. Our optimization method of choice is RMSProp (Duchi et al., 2011), employing a learning rate of 0.15 and initializing the accumulator value at 0.1. Apart from that we have used Adam optimizer.

During the training phase and during testing, we employ truncation strategies where we limit the article's length to 800 tokens and constrain the summary's length to 150 tokens

during training and testing. This practice aims to expedite both training and testing procedures. Remarkably, our observations indicate that truncating the article can actually enhance the model's performance. This outcome is attributed to several factors, including the facilitation of quicker computation, the reduction of memory usage, and the encouragement of the model to prioritize pertinent information, thereby fostering more concise and informative summaries.

We conducted training sessions for the seq2seq model up to 10 epochs, consuming a few days like 2 days. Notably, the pointer-generator model exhibited faster training dynamics, completing 10 epochs consuming a few hours merely a day. Noteworthy is the pointer-generator model's rapid advancement during the initial training phases. To refine our final coverage model, we introduced the coverage mechanism with coverage loss weighted at  $\lambda = 1$ , as specified in coverage mechanism section (last equation). During this period, the coverage loss converged to approximately 0.2, from its initial value of about 0.5. Experimentation with a more aggressive  $\lambda$  value of 2 decreased coverage loss but simultaneously escalated the primary loss function, rendering it unsuitable for use.

We conducted trials by omitting the dedicated loss function in training the coverage model, hypothesizing that the attention mechanism could organically adapt to minimize repetitive attention to the same locations. Nonetheless, this strategy yielded no discernible reduction in repetition, proving its ineffectiveness. Furthermore, we investigated incorporating coverage right from the onset rather than as a separate training stage. However, we noted that during the initial training phases, integrating coverage objectives prematurely interfered with the primary training objective, leading to an overall deterioration in performance.

## 6. Result

Model Used	ROUGE			
	1	2	L	
Sequence-To-Sequence Model	3.75	0.24	3.58	
Pointer Generator + Coverage Mechanism	1	2	L	
	36.56	12.92	33.10	

**Table 1: F1 ROUGE scores of both models employed**

We present our findings in Table 1. Our assessments utilize the widely-used ROUGE metric (Lin, 2004b), which measures the degree of overlap in words (ROUGE-1), bigrams (ROUGE-2), and the longest common sequence (ROUGE-L) between the generated summary and the reference summary. We calculate the F1 scores for ROUGE-1, ROUGE-2, and ROUGE-L. For computing ROUGE scores, we utilize the pyrouge package. As we generate summaries in plain text while Nallapati et al. (2016; 2017) produce anonymized summaries (refer to Section 4), direct comparisons of our ROUGE scores are not entirely equitable. It appears that datasets containing original text may yield higher ROUGE scores than anonymized datasets.

Based on the F1 ROUGE scores presented in Table 1, it is evident that the Pointer Generator with Coverage Mechanism model significantly outperforms the Sequence-To-Sequence Model. Specifically, the Pointer Generator with Coverage Mechanism achieved substantially higher scores across all ROUGE metrics: ROUGE-1, ROUGE-2, and ROUGE-L.

The Sequence-To-Sequence Model attained F1 scores of 3.75, 0.24, and 3.58 for ROUGE-1, ROUGE-2, and ROUGE-L respectively, indicating comparatively lower performance. In contrast, the Pointer Generator + Coverage Mechanism model demonstrated remarkable proficiency with F1 scores of 36.56, 12.92, and 33.10 for ROUGE-1, ROUGE-2, and ROUGE-L, showcasing its superior capability in generating more accurate and coherent summaries.

These results highlight the effectiveness of incorporating a Pointer Generator and Coverage Mechanism within the summarization architecture, enabling better content selection and synthesis, leading to enhanced summary quality.

Further analysis reveals the significant advantages offered by the Pointer Generator + Coverage Mechanism model over the Sequence-To-Sequence Model in the domain of text summarization. The notable performance gap between the two models underscores the importance of leveraging advanced techniques such as pointer generation and coverage mechanisms in enhancing summarization quality.

Furthermore, the Pointer Generator mechanism facilitates the incorporation of out-of-vocabulary words, thereby enhancing the model's vocabulary coverage and ensuring the fidelity of the generated summaries. Additionally, the Coverage Mechanism mitigates the issue of repetition by encouraging the model to evenly distribute its attention across the source text, thereby producing more coherent and diverse summaries.

### Output

<p><b>Article:</b> Enforcement directorate ed Friday conducted searches covering 11 states raided 35 new locations fraud accused billionaire diamond jeweller Nirav Modi according reports diamond gold worth ₹ crore seized 29 immovable properties also identified brings total value seized assets today ₹ crore.</p>
--

<p><b>Seq2Seq Model:</b> us arrested arrested found</p>
---

<p><b>Pointer Generator with Coverage Mechanism:</b> ed raids new locations fraud accused diamond jewellery...</p>
--

## 7. Conclusion

The pointer-generator network has shown significant promise in the field of natural language processing, particularly in tasks such as text summarization. Its ability to combine the strengths of both extractive and abstractive approaches, along with mechanisms like coverage, makes it a versatile and powerful tool. Our study introduces a novel hybrid pointer-generator architecture augmented with coverage mechanisms, offering a promising avenue for enhancing the quality of generated text summaries. By leveraging this architecture, we have effectively mitigated inaccuracies and repetitive patterns commonly observed in traditional abstractive summarization models. Through rigorous experimentation on a challenging long-text dataset, our proposed model has demonstrated remarkable performance gains compared to the current state-of-the-art abstractive summarization approaches. Notably, our model showcases proficient abstractive abilities, underscoring its potential for various text generation tasks. While our findings mark significant progress in the realm of abstractive summarization, there exist intriguing avenues for future exploration. One notable direction is the pursuit of higher levels of abstraction, where the model transcends mere summarization to capture deeper semantic nuances and generate more insightful and contextually rich summaries. Furthermore, the integration of multi-modal inputs, such as incorporating visual or contextual cues alongside textual information, holds promise for further enhancing the richness and relevance of generated summaries. Additionally, exploring techniques for controllable generation, allowing users to influence specific aspects of the generated summaries, could open new dimensions for personalized and tailored summarization solutions. In essence, our work underscores the potential of hybrid pointer-generator architectures in advancing the frontier of abstractive summarization and highlights exciting opportunities for future research to explore and innovate in this domain.

## 8. References

- Abigail See, Peter J. Liu, Christopher D. Manning. 2017. Get To The Point: Summarization with Pointer-Generator Networks.
- Junhao LI, Mizuho IWAHARA. 2019. Two-Encoder Pointer-Generator Network for Summarizing Segments of Long Articles.
- Ramesh Nallapati, Feifei Zhai, and Bowen Zhou. 2017. SummaRuNNer: A recurrent neural network based sequence model for extractive summarization of documents. In Association for the Advancement of Artificial Intelligence.
- Ramesh Nallapati, Bowen Zhou, Cicero dos Santos, Caglar Gulcehre, and Bing Xiang. 2016. Abstractive text summarization using sequence-to-sequence RNNs and beyond. In Computational Natural Language Learning.
- Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. In Neural Information Processing Systems.
- Vishal Gupta, Gurpreet Singh Lehal. "A Survey of Text Summarization Extractive Techniques". Journal of emerging technologies in web intelligence (Vol. 2, No. 3, August 2010).
- Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. 2015. "Pointer networks". In Neural Information Processing Systems.
- Dat P.T Nguyen, Yutaka Matsuo, and Mitsuru Ishizuka. "Exploiting Syntactic and Semantic Information for Relation Extraction from Wikipedia " Text-Mining and Link-Analysis (TextLink2007), 2007.
- Alexander M Rush, Sumit Chopra, and Jason Weston. 2015. "A neural attention model for abstractive sentence summarization". In Empirical Methods in Natural Language Processing.