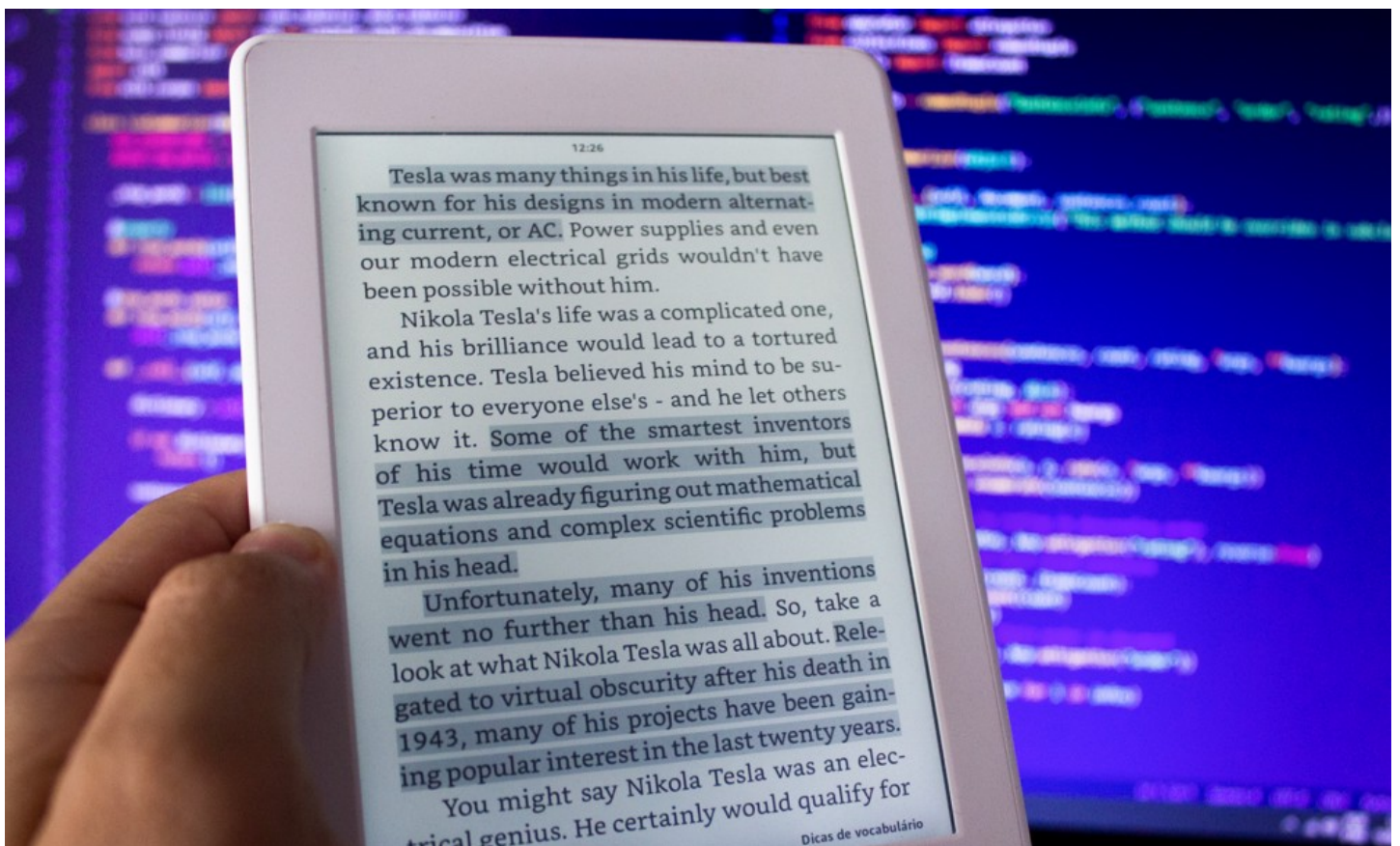


Automatic Text Summarization with Machine Learning — An overview



Luís Gonçalves

Apr 11 · 11 min read



Summarization is the task of condensing a piece of text to a shorter version, reducing the size of the initial text while at the same time preserving key informational elements and the meaning of content. Since manual text summarization is a time expensive and generally laborious task, the automatization of the task is gaining increasing popularity and therefore constitutes a strong motivation for academic research.

There are important applications for text summarization in various NLP related tasks such as text classification, question answering, legal texts summarization, news summarization, and headline generation. Moreover, the generation of summaries can

be integrated into these systems as an intermediate stage which helps to reduce the length of the document.

In the big data era, there has been an explosion in the amount of text data from a variety of sources. This volume of text is an inestimable source of information and knowledge which needs to be effectively summarized to be useful. This increasing availability of documents has demanded exhaustive research in the NLP area for automatic text summarization. Automatic text summarization is the task of producing a concise and fluent summary without any human help while preserving the meaning of the original text document.

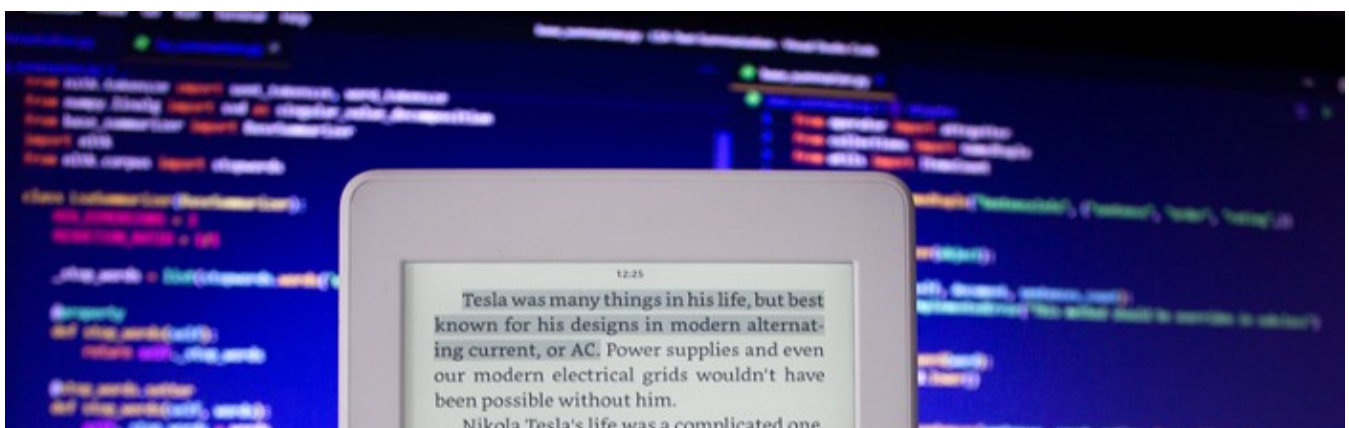
It is very challenging, because when we as humans summarize a piece of text, we usually read it entirely to develop our understanding, and then write a summary highlighting its main points. Since computers lack human knowledge and language capability, it makes automatic text summarization a very difficult and non-trivial task.

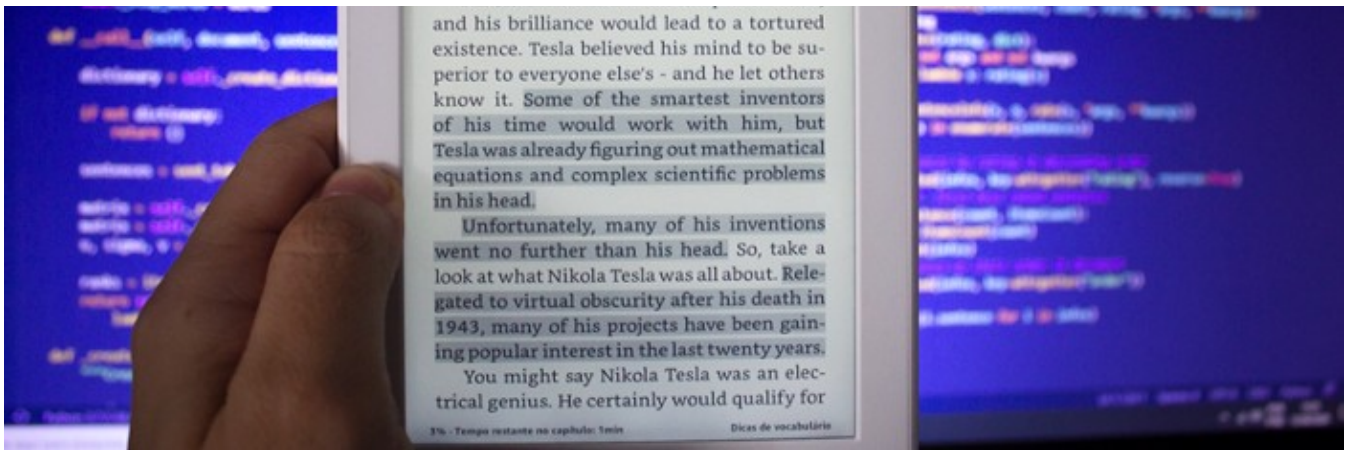
Various models based on machine learning have been proposed for this task. Most of these approaches model this problem as a classification problem which outputs whether to include a sentence in the summary or not. Other approaches have used topic information, Latent Semantic Analysis (LSA), Sequence to Sequence models, Reinforcement Learning and Adversarial processes.

In general, there are two different approaches for automatic summarization: **extraction** and **abstraction**.

The extractive approach

Extractive summarization picks up sentences directly from the document based on a scoring function to form a coherent summary. This method work by identifying important sections of the text cropping out and stitch together portions of the content to produce a condensed version.





Extractive summarization work by identifying important sections of the text cropping out and stitch together portions of the content to produce a condensed version. Thus, they depend only on the extraction of sentences from the original text.

Thus, they depend only on the extraction of sentences from the original text. Most of the summarization research today has focused on extractive summarization, once it is easier and yields naturally grammatical summaries requiring relatively little linguistic analysis. Moreover, extractive summaries contain the most important sentences of the input, which can be a single document or multiple documents.

A typical flow of extractive summarization systems consists of:

1. Constructs an intermediate representation of the input text intending to find salient content. Typically, it works by computing TF metrics for each sentence in the given matrix.
2. Scores the sentences based on the representation, assigning a value to each sentence denoting the probability with which it will get picked up in the summary.
3. Produces a summary based on the top k most important sentences. Some studies have used Latent semantic analysis (LSA) to identify semantically important sentences.





For a good starting point to the LSA models in summarization, [check this paper](#) and [this one](#). An implementation of LSA for extractive text summarization in Python is available in this [github repo](#). For example, I used this code to make the following summary:

Original text:

De acordo com o especialista da Certsys (empresa que tem trabalhado na implementação e alteração de fluxos desses robôs), Diego Howës, as empresas têm buscado incrementar os bots de atendimento ao público interno com essas novas demandas de prevenção, para que os colaboradores possam ter à mão informações sobre a doença, tipos de cuidado, boas práticas de higiene e orientações gerais sobre a otimização do home office. Já os negócios que buscam se comunicar com o público externo enxergam outras necessidades. “Temos clientes de varejo que pediram para que fossem criados novos fluxos abordando o tema, e informando aos consumidores que as entregas dos produtos adquiridos online podem sofrer algum atraso”, comenta Howës, da Certsys, que tem buscado ampliar o escopo desses canais para se adequar ao momento de atenção. Ainda segundo o especialista, em todo o mercado é possível observar uma tendência de automatização do atendimento à população, em busca de chatbots que trabalhem em canais de alto acesso, como o WhatsApp, no caso de órgãos públicos. Na área de saúde, a disseminação de informação sobre a pandemia do vírus tem sido um esforço realizado.

Summarized text:

De acordo com o especialista da Certsys (empresa que tem trabalhado na implementação e alteração de fluxos desses robôs), Diego Howës, as empresas têm buscado incrementar os bots de atendimento ao público interno com essas novas demandas de prevenção, para que os colaboradores possam ter à mão informações sobre a doença, tipos de cuidado, boas práticas de higiene e orientações gerais sobre a otimização do home office. Já os negócios que buscam se comunicar com o público externo enxergam outras necessidades. Na área de saúde, a disseminação de informação sobre a pandemia do vírus tem sido um esforço realizado.

Recent studies have applied deep learning in extractive summarization as well. For instance, Sukriti proposes an extractive text summarization approach for factual reports using a deep learning model, exploring various features to improve the set of sentences selected for the summary.

Yong Zhang proposed a document summarization framework based on convolutional neural networks to learn sentence features and perform sentence ranking jointly using a CNN model for sentence ranking. The authors adapt the original classification model of Y. Kim to address a regression process for sentence ranking. The neural architecture used in that paper is compound by one single convolution layer that is built on top of the pre-trained word vectors followed by a max-pooling layer. The author carried experiments on both single and multi-document summarization tasks to evaluate the proposed model. Results have shown the method achieved competitive or even better performance compared with baselines. The source code used in experiments can be found here.

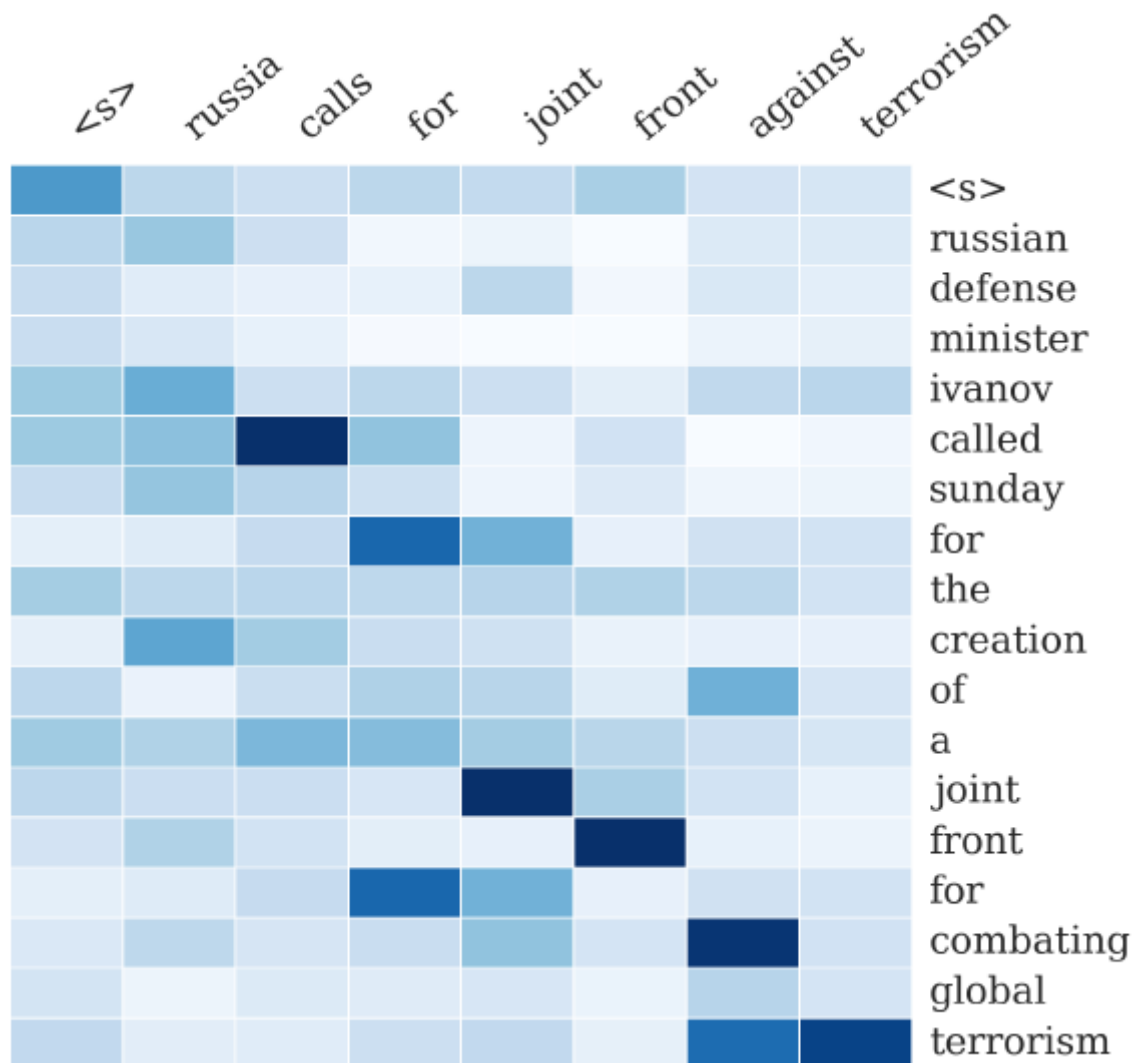
Abstractive summarization

Abstractive summarization methods aim at producing summary by interpreting the text using advanced natural language techniques in order to generate a new shorter text — parts of which may not appear as part of the original document, that conveys the most critical information from the original text, requiring rephrasing sentences and incorporating information from full text to generate summaries such as a human-written abstract usually does. Thus, they are not restricted to simply selecting and rearranging passages from the original text.

Abstractive methods take advantage of recent developments in deep learning. Since it can be regarded as a sequence mapping task where the source text should be mapped to the target summary, abstractive methods take advantage of the recent success of the sequence to sequence models. These models consist of an encoder and a decoder, where a neural network reads the text, encodes it and then generates target text.

In general, building abstract summaries is a challenging task, which is relatively harder than data-driven approaches such as sentence extraction and involves complex language modeling. Thus, they are still far away from reaching human-level quality in summary generation, despite recent progress using neural networks inspired by the progress of neural machine translation and sequence to sequence models.

An example is the work of [Alexander et al](#), which proposed a neural attention model for abstractive sentence summarization (**NAMAS**) by exploring a fully data-driven approach for generating abstractive summaries using an attention-based encoder-decoder method. [Attention mechanism](#) has been broadly used in sequence to sequence models where the decoder extracts information from the encoder based on the attention scores on the source-side information. The code to reproduce the experiments from the **NAMAS** paper [can be found here](#).



Example output of the attention-based summarization of [Alexander et al](#). The heatmap represents a soft alignment between the input (right) and the generated summary (top). The columns represent the distribution over the input after generating each word.

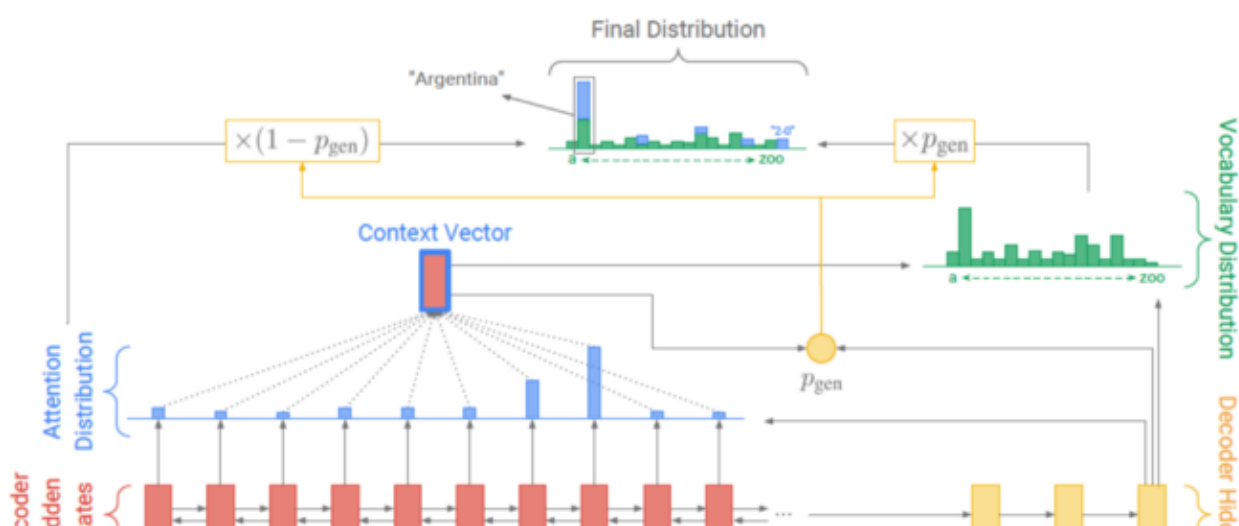
Recent studies have argued attention-based sequence to sequence models for abstractive summarization can suffer from repetition and semantic irrelevance, causing grammatical errors and insufficient reflection of the main idea of the source text. [Junyang Lin et al](#) propose to implement a gated unit on top of the encoder outputs at each time step, which is a CNN that convolves all the encoder outputs, in order to tackle this problem.

Based on the convolution and self-attention of [Vaswani et al.](#), a convolutional gated unit sets a gate to filter the source annotations from the RNN encoder, in order to select information relevant to the global semantic meaning. In other words, it refines the representation of the source context with a CNN to improve the connection of the word representation with the global context. Their model is capable of reducing repetition compared with the sequence to sequence model outperforming the state-of-the-art methods. The source code of paper [can be found here](#).

Other methods for abstractive summarization have borrowed the concepts from the pointer network of [Vinyals et al](#) to addresses the undesirable behavior of sequence to sequence models. Pointer Network is a neural attention-based sequence-to-sequence architecture that learns the conditional probability of an output sequence with elements that are discrete tokens corresponding to positions in an input sequence.

For example, [Abigail See et al.](#) presented an architecture called Pointer-Generator, which allows copying words from the input sequence via pointing of specific positions, whereas a generator allows generating words from a fixed vocabulary of 50k words. The architecture can be viewed as a balance between extractive and abstractive approaches.

In order to overcome the repetition problems, the paper adapts the coverage model of [Tu et al.](#), which was proposed to overcome the lacking coverage of source words in neural machine translation models. Specifically, Abigail See et al. defined a flexible coverage loss to penalize repeatedly attending to the same locations, only penalizing the overlap between each attention distribution and the coverage up to the current time step helping to prevents repeated attention. The source code for the model [can be found here](#).

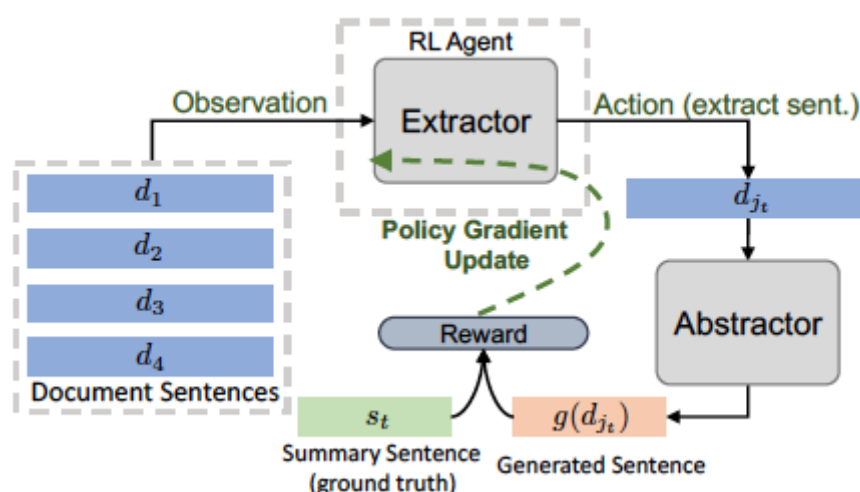




The Pointer-generator model. For each timestep in the decoder, the probability of generating words from the fixed vocabulary, versus copying words from source using a pointer is weighted by a generation probability p_{gen} . The vocabulary distribution and attention distribution are weighted and summed to obtain the final distribution. The attention distribution can be viewed as a probability distribution over the source words, that tells the decoder where to look to generate the next word. It is used to produce a weighted sum of the encoder hidden states, known as the context vector.

Other studies have in abstractive summarization have borrowed the concepts from the reinforcement learning (RL) field to improve model accuracy. For example, [Chen et al.](#) proposed a hybrid extractive-abstractive architecture using two neural networks in a hierarchical way, that selects salient sentences using an RL guided extractor from the source and then rewrites them abtractively to generate a summary.

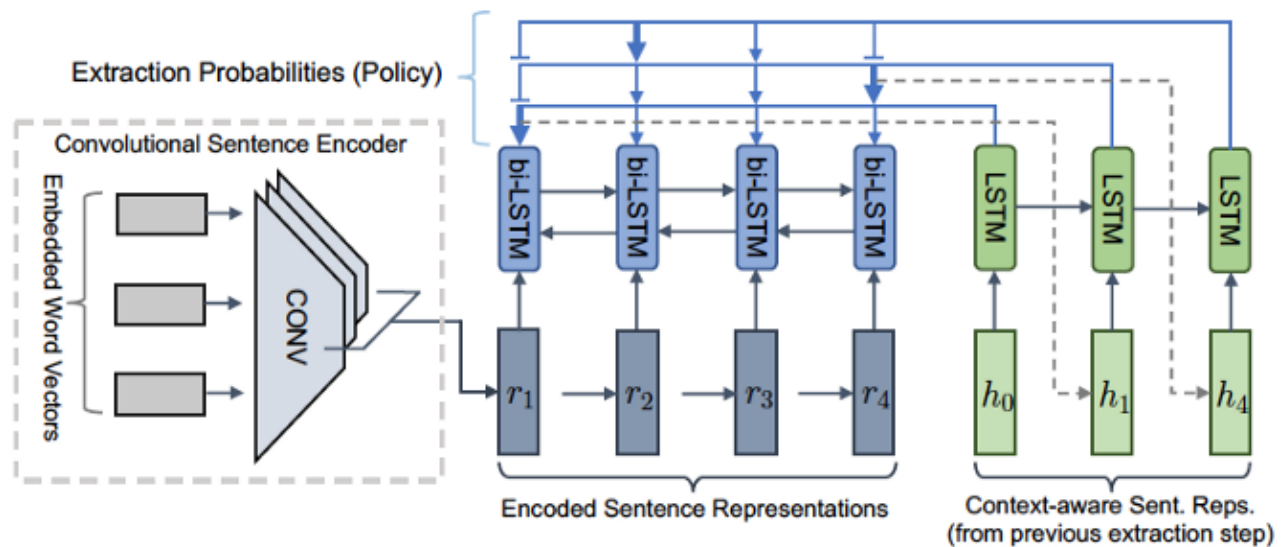
In other words, the model simulates how humans summarize long documents first using an extractor agent to select salient sentences or highlights, and then employs an abstractor — an encoder-aligner-decoder model — network to rewrite each of these extracted sentences. To train the extractor on available document-summary pairs, the model uses a policy-based reinforcement learning (RL) with sentence-level metric rewards to connect both extractor and abstractor networks and to learn sentence saliency.



Reinforced training of the extractor (for one extraction step) and its interaction with the abstractor.

The abstractor network is an attention-based encoder-decoder which compresses and paraphrases an extracted document sentence to a concise summary sentence.

Moreover, the abstractor has a useful mechanism to help directly copy some out-of-vocabulary (OOV) words.



The convolutional extractor agent

The extractor agent is a convolutional sentence encoder that computes representations for each sentence based on input embedded word vectors. Further, an RNN encoder computes context-aware representation and then an RNN decoder selects sentence at time step t . Once the sentence is selected, the context-aware representation will be fed into the decoder at time $t + 1$.

Thus, the method incorporates the abstractive approach advantages of concisely rewriting sentences and generating novel words from the full vocabulary, whereas adopts intermediate extractive behavior to improve the overall model's quality, speed, and stability. The author argued model training is 4x faster than the previous state-of-the-art. Both [source code](#) and best pre-trained models were released to promote future research.

Other recent studies have proposed using a combination of the adversarial processes and reinforcement learning to abstractive summarization. An example is [Liu et al. \(2017\)](#), whose work proposes an adversarial framework to jointly train a generative model and a discriminative model similar to [Goodfellow et al. \(2014\)](#). In that framework, a generative model takes the original text as input and generates the summary using reinforcement learning to optimize the generator for a highly rewarded summary. Further, a discriminator model tries to distinguish the ground truth summaries from the generated summaries by the generator.

The discriminator is implemented as a text classifier that learns to classify the generated summaries as machine or human-generated, while the training procedure of generator is to maximize the probability of discriminator making a mistake. The idea is this adversarial process can eventually let the generator to generate plausible and high-quality abstractive summaries. The author provided supplementary material [here](#). The source code is available in [this github repo](#).

In short

Automatic text summarization is an exciting research area with several applications on the industry. By condensing large quantities of information into short, summarization can aid many downstream applications such as creating news digests, report generation, news summarization, and headline generation. There are two prominent types of summarization algorithms.

First, extractive summarization systems form summaries by copying and rearranging passages from the original text. Second, abstractive summarization systems generate new phrases, rephrasing or using words that were not in the original text. Due to the difficulty of abstractive summarization, the great majority of past work has been extractive.

The extractive approach is easier because copying large chunks of text from the source document ensures good levels of grammaticality and accuracy. On the other hand, sophisticated abilities that are crucial to high-quality summarization, such as paraphrasing, generalization, or the incorporation of real-world knowledge, are possible only in an abstractive framework. Even though abstractive summarization is a more challenging task, there has been a number of advances so far, thanks to recent developments in the deep learning area.

Further reading

1. Extractive Text Summarization using Neural Networks — [Sinha et al.\(2018\)](#).
2. Extractive document summarization based on convolutional neural networks — [Y. Zhang et al. \(2016\)](#).
3. A Neural Attention Model for Abstractive Sentence Summarization — [Rush et al. \(2015\)](#).
4. Global Encoding for Abstractive Summarization — [Lin et al.\(2018\)](#).
5. Summarization with Pointer-Generator Networks — [See et al.\(2017\)](#).
6. Fast Abstractive Summarization with Reinforce-Selected Sentence Rewriting —