# Science Checker: Extractive-Boolean Question Answering For Scientific Fact Checking

Loïc Rakotoson
loic.rakotoson@opscidia.com
Opscidia
Paris, France

Charles Letaillieur
charles.letaillieur@opscidia.com
Opscidia
Paris, France

Sylvain Massip
sylvain.massip@opscidia.com
Opscidia
Paris, France

Fréjus A. A. Laleye
frejus.laleye@opscidia.com
Opscidia
Paris, France

## ABSTRACT

With the explosive growth of scientific publications, making the synthesis of scientific knowledge and fact checking becomes an increasingly complex task.

In this paper, we propose a multi-task approach for verifying the scientific questions based on a joint reasoning from facts and evidence in research articles. We propose an intelligent combination of (1) an automatic information summarization and (2) a Boolean Question Answering which allows to generate an answer to a scientific question from only extracts obtained after summarization.

Thus on a given topic, our proposed approach conducts structured content modeling based on paper abstracts to answer a scientific question while highlighting texts from paper that discuss the topic. We based our final system on an end-to-end Extractive Question Answering (EQA) combined with a three outputs classification model to perform in-depth semantic understanding of a question to illustrate the aggregation of multiple responses. With our light and fast proposed architecture, we achieved an average error rate of 4% and a F1-score of 95.6%. Our results are supported via experiments with two QA models (BERT, RoBERTa) over 3 Million Open Access (OA) articles in the medical and health domains on Europe PMC.

## CCS CONCEPTS

• **Information systems** → **Question answering**; **Information extraction**; *Summarization*; *Language models*; Trust.

## KEYWORDS

Scientific Fact checking, Question Answering, Medical information extraction, Query Analysis

## 1 INTRODUCTION

For many years, public trust in science has been one of the concerns, and sometimes fears, of the scientific community [1, 7, 23]. From simple misconception to distrust, from discredit to the creation of an alternative science, opposition to science is increasingly reinforced nowadays.

Hence, the development of Open Access should be a solution to improve the communication of scientific result to the general public. Yet, even with completely open access, two important difficulties remain: domain knowledge is needed to understand the literature, and the volume of research article is such that it is difficult to read and analyse every article on a topic.

Modern communication means and the recent evolution in science communication gave the possibility to experts to popularize knowledge, notably in public conferences and social networks [2, 11, 29]. However, this has also given a voice to non-experts and pseudo-science. Science being a slow process, this gives time for obscurantism to take hold. This results in mistrust which sometimes adds to ignorance and contributes to the scientific fake news expansion.

Therefore, automatic means to source and verify the information is more and more needed, and can be a nice complement to the work of human expertise. This work is part of a project called Science Checker, which aims to help non experts to navigate in the scientific literature and improve its access to the best scientific information. The general motivation and approach of the project is described elsewhere [18, 28]. In this article, we focus on the description and analysis of the classification pipeline which aims to classify articles as supporting or contradicting a scientific affirmation.

Overall, our science checker system operates on two linked tasks: (1) an extractive question answering component that, by neural semantic matching between the representations of a query and all available abstracts of scientific articles, gives an optimal representation at the granularity of text chunks for fact retrieval; (2) a multi

Question: Does Hydroxychloroquine cure COVID-19 ?

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

Abstract 1: Background Some disease-modifying agents [...] No significant difference was found in terms of rates of usage of hydroxychloroquine or colchicine between those who were found positive for SARS-CoV-2 and those who were found negative (0.23% versus 0.25% for hydroxychloroquine, and 0.53% versus 0.48% for colchicine, respectively) Conclusion These findings raise doubts regarding the protective role of these medications in the battle against SARS-CoV-2 infection.
Answer 1: No

Abstract 2: [...] Conclusions Hydroxychloroquine has received worldwide attention as a potential treatment for covid-19 because of positive results from small studies. However, the results of this study do not support its use in patients admitted to hospital with covid-19 who require oxygen.
Answer 2: No

**Figure 1: Example of task. Given the question, the system is required to find the candidates pieces of texts and to give an answer. (Abstract 1 [6], Abstract 2 [17])**

label classification model that, from the facts found, gives the probabilities that theses are **False**, **True** or **Neutral** with a weighted average of the results as output. Figure 1 perfectly illustrates the task tackled in this work behind which one can imagine the difficulties for the accurate selection of potential texts containing scientific facts from a huge knowledge base.

To contribute to Open Science, our results, data and code will be made publicly available under a Free/Libre Open Source license.

## 2 RELATED WORKS

*Fact checking.* The SemEval-2019 fact-checking competition [19] in Task 8B highlighted approaches by classifying terms and authors in a forum to obtain the veracity of information in a community. The platform built by Miranda et al. [20] allows the fact-checking of a claim by selecting the closest sentences to it within a threshold, from about ten thousand of newspaper articles, and then classifying them if they refute or support the claim. The classifier was trained on the Wikipedia-based fact-checking dataset FEVER (Thorne et al. [31]). Similar work in [21] proposed an hybrid CNN-RNN model to detect the sentence that may be the fake news. Karadzhov et al. [10]'s work starts with a sentence which is preprocessed to create a query on Google and Bing search engines. The snippets in the results are then compared to the source sentence to determine whether it is factual or not. The model which compute the comparison is a combination of SVM and a recurrent neural network, and achieves excellent performance based on a dataset built from Snopes. The approach proposed by Yang et al. [34] aims at extracting information by combining text and image which is processed by a convolutional neural network, based on their explicit and latent variables and which performs an F1 score of 92%.

In this work and unlike the works [19–21], we focus on verifying the facts by collecting information from the scientific literature and rendering it with an intelligible answer. In contrast to newspaper information and more oriented document, the answer to a query is very rarely explicit in scientific article. Closed questions do not get a definite binary solution in the scientific literature. So, unlike the work in [33] which similarly worked on the same task, our approach does not aim only to report the information but also to identify the interesting elements that allow for a clear answer.

In [36], the authors proposed a graph-based information aggregation approaches from a claim to select the relative information from Wikipedia by representing it on a connected evidence graph (GEAR) for verifying a query.

On the same basis, Liu et al. [16]'s proposed an extension of the graph-based approach (GEAR) for a finer-grained highlighting of claim responses. On FEVER dataset, they outperform GEAR approach of 3.3%.

We deviate, in this work, from these approaches in the sense that the interest of our task is to bring out the information by scientific document so as to allow an easy tracing of the information sources.

*Question-answering.* In this paragraph, we cite some question answering works that have adopted an extrative approach.

Lewis et al. [12] and Puri et al. [25] developed an approach that consists in creating an extractive question answering system without data annotations. This consists in a first step of question generation by masking or in zero-shot [24], followed by a second step of question answering on the previously generated data. This method avoids the need to search the dataset for a specific subtask but requires the existence of the data in which the models will evolve. The downside is the quality dependency between the models of each phase. Yoon et al. [35] obtained the best scores during the BioASQ competition [22] on question-answering task with their approach based on modifying BioBERT followed by a task specific and a post-processing layer.

Many state-of-the-art approaches have been developed to address the question answering task since the long form question answering dataset [4] has been proposed.

These approaches and methods, unlike our work, focus exclusively on extracting information without answering the question in a straightforward and concise way for the non-specialist. SciFact [33] partially addresses this problem by collecting sentences that explicitly respond to a claim with an additional field that states whether or not that group of sentences supports it. However, our work aims to extract less obvious information from larger sources that are intended to be unbiased. The necessary information is found on fragmented parts all over the document and the purpose is to gather it. Finally, in addition to Yes or No responses, we want to capture the neutral conclusions frequently encountered in the scientific literature.

## 3 TASK FORMALIZATION

As shown in the example in figure 1, for a query in a specific domain (medical and health in this work), we aim first to collect the open access scientific articles related to it.

This large corpus of documents will serve as a basis for representing the scientific consensus on this question. This representation must be detailed, motivated and easily intelligible for non-experts in the field.

Let $\mathbb{D}$ be the set of available scientific articles, $D_i$ be an individual abstract and $S_i^j$ be the $j$-th text chunk in the $i$-th. For the given fact $c_i$, and the set of abstracts $\mathbb{D} = \{D_0, D_1, ..., D_n\}$ where $D_i = \{s_i^0, s_i^1, ..., s_i^n\}$, the task is to provide the predicted response $(\hat{R}_i, \hat{y}_i)$ where $\hat{R}_i$ represents the set of sentences of an individual abstract and $\hat{y}_i \in \{T, F, N\}$ gives the probability that the given fact is True ($T$), False ($F$) and Neutral ($N$).

## 4 OUR APPROACH

### 4.1 Method Background

The purpose would be to summarize multiple articles and come up with an answer (as described in the algorithm 1). However, we want to exploit each document to be able to detail their contribution, the aggregation of which will constitute the final answer. We then propose to process each article individually, whose content can respond to the query, to extract the information needed for answering.

With the mass of open access articles on health and the thousands of possible questions, in this work we focused on closed questions worded as «**Does agent X *prevent / cure / cause / increase* disease Y?**».

To retrieve the set of articles considered in our work as candidate abstracts, we perform a semantic search in the mass of open access articles using the expression for agent X and disease Y.

---

**Algorithm 1:** Short summary of our proposed approach

**Input:** $X$: Does agent X verb disease Y?
**Output:** Predicted response: $(\hat{R}_i, \hat{y}_i)$

Find the set of candidate abstracts $\mathbb{D}$.
**for** $D_i$ **in** $\mathbb{D}$ **do**
  Create $n + 1$ windows of $t$ sentences and stride
  $0 \le p \le t - 1$ where:
  $d = t - p$
  $w_0 = \{s_1, ..., s_t\}$
  $w_n = \{s_{dn+1}, ...s_{dn+t}\}$
  where $w_j$ is truncated or padded to have 350 tokens.
  **for** $j \leftarrow 0$ **to** $n$ **do**
    EQA$_j(C, w_j) = \{s_b*, ..., s_e*\}$
    where $s*$ are the answers, $s_b$ the beginning of
      sentence highlight, $s_e$ the highlight end
  **end**
  $\hat{R}_i = \{\text{EQA}_1, \text{EQA}_j, ..., \text{EQA}_n\}$
  $\hat{y}_i = \text{BQA}(C, \hat{R}_i)$
**end**

---

The task is divided into two steps: an information gathering phase followed by a question answering phase. We used in each step Transformers-based structural textual representations [32]. The first step is to summarize an abstract based on the extraction of pieces of texts and the second step is purely a multi-outputs Boolean question-answering.

By combining all the Boolean outputs, we can therefore provide a final answer using a simple majority vote. The final output is as follows:

- Affirmative: the *yes* dominates.
- Negative: the *no* prevails.
- Balanced: the scientific opinions are divided on the fact.
- Neutral: the selected OA articles do not allow to answer the question.

This combination is made to allow traceability from the response generated so that it can easily be traced back to the source articles. Figure 2 presents an overview of our approach and shows how the two parts are linked from input representation to produce the responses.

### 4.2 Input representation

Let $D_i$ and $N$ be respectively an abstract in the set of candidate abstracts $\mathbb{D}$ and a fixed number of word windows. We build the inputs by splitting $D_i$ into sliding windows of 350 tokens, thus forming for each abstract $D_i$, so that the $jth$ window $w_{j*}$ denotes an input feeding the models separately. $D_i$ is obtained as follows.

$$D_i = \sum_{j=0}^{N-1} w_{j*} \tag{1}$$

Each window is fed separately to the EQA block which processes them all in parallel with as outputs: 0 when no window allows to respond and 1 otherwise with the spans of the resulting windows. The outputs are then concatenated (Figure 3) to feed the BQA block. It should be noted that the clear answer to a query can be located in various places of the abstract (neighboring sentences or not) and that for an abstract, the answers can vary according to the input question. Our approach therefore makes it possible to address this problem by covering the information throughout the entire abstract on the one hand and by aggregating it on the other hand for a precise answer.

### 4.3 Model

In this section, we describe the core of our approach, ie the extractive and boolean models and how they have been combined. Our information summarization model is based on an extractive approach which consists in quoting a part of the text that carries relevant information unlike the abstractive approach which consists in generating a shorter and more structured alternative text. This choice is motivated by the fact that we need the raw pieces of the text to feed the boolean model without bringing it external knowledge which could add biases to the main information and impact the decision making by the boolean model. We nevertheless experimented with the two approaches for the purposes of performance comparisons only on the information extraction.

*4.3.1 Extractive highlights.* The Extractive model consists in selecting only the relevant parts of a long text to generate a new content to be used for answering the question. Let $(C, D)$ be the set of question and text that is divided into sentences $s^*$ which are divided into tokens $\omega_N^* = \omega_1, \omega_2, ..., \omega_N$. We then assign to each token the probabilities ($\mathbb{P}$, a matrix $\mathcal{M}$ of $N$ tokens and 2 indexes)
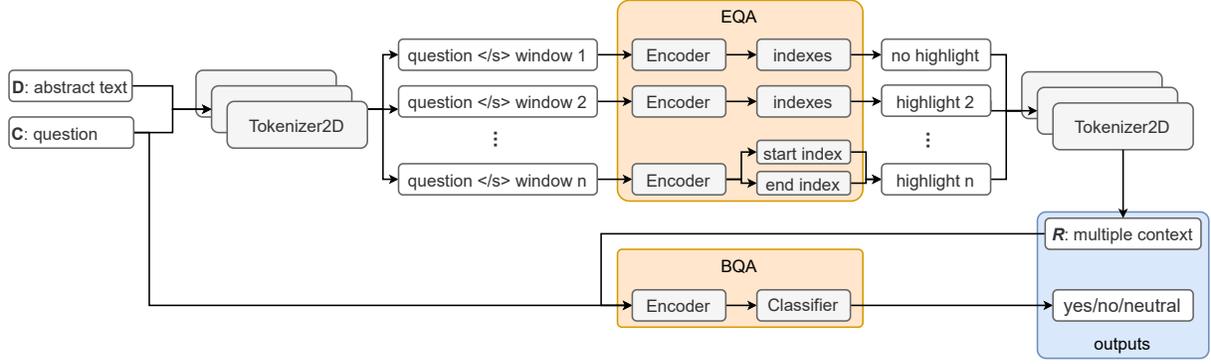
**Figure 2: Combined system for science checker.** It consist of two two parts: abstract summarization based on Extractive Question Answering (EQA) that exploits information contained in sliding pieces of an abstract (Eq. 1) and a multi-outputs Boolean Question-Answering (BQA) that combines all extracted contents to generate multiple output.

that they are the beginning and the end of an important part. For one window, $R$ is computed such that:

$$(C, D) = \omega_N^*, \mathbb{P}(C, D) \in \mathcal{M}_{(N, 2)}$$
$$(m, M) = \text{argmax}\left[\mathbb{P}(C, D)\right] \in \mathbb{R}^2 \qquad (2)$$
$$R = h = \omega_{m:M}$$

For each window $w_j$ of the Tokenizer2D, the couple $(C, D_i)$ gives a highlight $h_j$. Their concatenation without repetition gives the context $R$ (Figure 3). To achieve such solution, we built a model which, starting from a text representation, outputs the start and end positions of an important part of it. We used BERT and RoBERTa [15], each with their base and large versions, to calculate the text embeddings. The positions are given by two dense layers of the same number of units as the input layer dimension (Figure 4).

Afterwards, the best of the 4 architectures with a richer representation, and by nature more efficient than its basic version, is distilled to 2 simpler versions, MobileBERT proposed by Sun et al. [30] and TinyBERT proposed by Jiao et al. [8]. We experimented with our version of TinyBERT by keeping the same number of Encoder layers as $\text{BERT}_{\text{BASE}}$ but reducing the attention heads to one pair for each block. A quantization is finally applied to the optimal model.

*4.3.2 Abstractive model.* As indicated above, we also experimented with Abstractive model for summarizing information. The idea behind the Abstractive model is to take a pair of question and long document $(C, D)$, and then produce a shorter summary $R$, a sequence of tokens $\omega$ of length $T$. We then perform a directed conditional text generation such that:

$$R = \omega_1, \omega_2, ..., \omega_T$$
$$R^t = \omega_{1:t} \text{ and } R^0 = \varnothing \qquad (3)$$
$$\mathbb{P}(R \mid C, D) = \prod_{t=1}^{T} \mathbb{P}\left(\omega_t \mid R^{t-1}, C, D\right)$$

$(C, D_j)$ are contained in each window $w_j$ of the first block of the Tokenizer2D, and all $R_j$ are concatenated in the second block of the Tokenizer2D. The found optimal value of $T$ is 80 which does not

prevent our early stopping strategy to stop the generation before $t = 80$ when the token EOS is part of the most likely branch.

For this type of sequence-to-sequence model, we used the auto-regressive $T5$ configuration proposed in [26] by focusing on the text summarization task. We have fine-tuned a base version with 220 million parameters and a lighter version with about 60 million parameters.

*4.3.3 Boolean answer.* The Boolean model is the final block which receives the extractive model outputs $R_i$, and gives a Boolean answer for each abstract $D_i$. We fed the classifier with the text representation using the CLS token; and for the latter we evaluated the performance on a direct learning with the BERT architecture in its TinyBERT version and compare it with $\text{RoBERTa}_{\text{BASE}}$. The model has three outputs for *affirmative*, *negative* and *neutral* answers only if the context does not allow to answer the question.

## 5 EXPERIMENTS

First, we present the experiments with the extractive and abstractive summarization approaches before analyzing the performance of the extractive approach that best meets our task. The goal is to reduce an abstract by keeping only the relevant information.
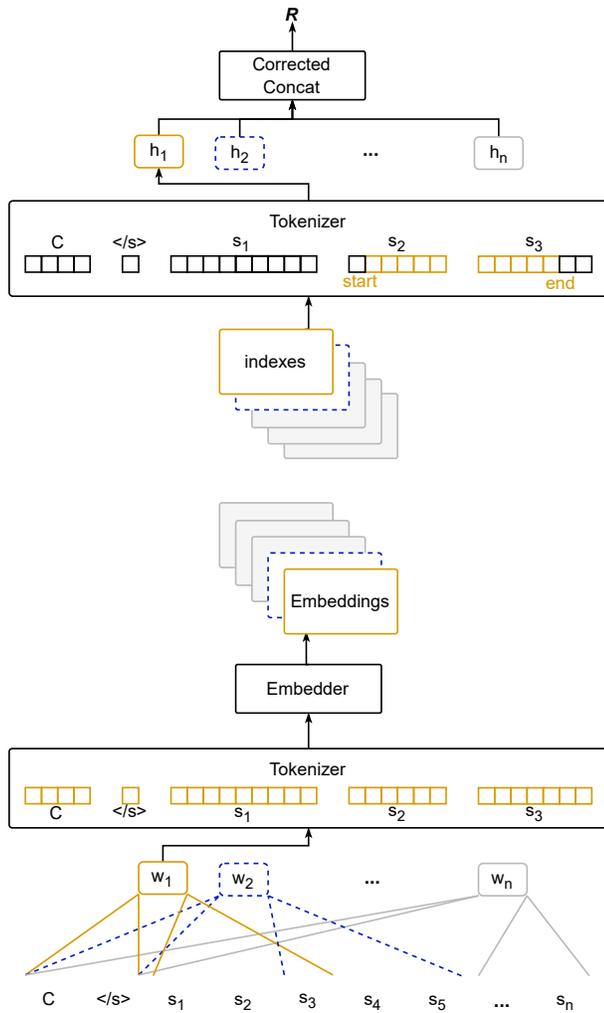
### 5.1 Experimental settings

*5.1.1 Datasets.* To train our extractive models, we formatted the dataset to have a final context which is the concatenation of the question, the separation token and the initial context. The outputs are replaced by the answer start and end tokens positions after the final context tokenization.

Two datasets were constituted and formatted. The first one is SQuAD v2 [27] on which the initial training was performed, the second one is the merge of BioASQ (Factoid and Lists) [22] and PubMedQA [9] in which the responses are fully included in the context.

Our Abstractive training dataset is a fusion of PubMedQA and BioASQ for the *Factoid* task, where the question and context block are concatenated with the separation token. We used Tokenizer2D with a size of $t = 7$ sentences and stride of $p = 0$ to capture the whole text, i.e., 2 windows per article abstract at most. During the

**C**: Does hydroxychloroquine cure sars-cov-2 ?

**R**: Some disease-modifying agents [...] such as hydroxychloroquine [...]. However, the role of such agents as prophylactic tools is still not clear. [...] An overall sample of 14,520 subjects were screened for SARS-CoV-2 [...] No significant difference was found in terms of rates of usage of hydroxychloroquine [...] These findings raise doubts regarding the protective role of these medications in the battle against SARS-CoV-2 infection.

**C**: Does hydroxychloroquine cure sars-cov-2 ?

**s$_1$**: Some disease-modifying agents [...] diseases/autoimmune disorders.

**s$_2$**: However, the role of such agents as prophylactic tools is still not clear.

**s$_3$**: This is a retrospective study based on a large healthcare [...] 2020.

...

**s$_n$**: These findings raise doubts regarding the protective role [...] infection.

**Figure 3: Tokenizer2D.** The primary block (in the bottom) takes a pair of question $C$ and document $D$ decomposed into sentences $s_{1:n}$ and applies the sliding window strategy to return the embeddings. Here with a size $t = 3$ and a stride $p = 1$. The secondary block (on the top) returns information to initial input $D$ from the model outputs. Here it is an extractive model indexes.

training, the input data has at most 350 tokens, which corresponds to an average window size.
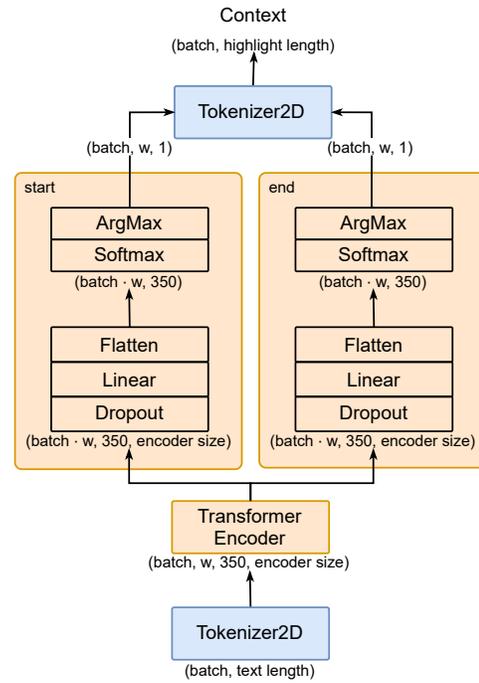
**Figure 4: Extractive Model.** Placed between the two blocks of Tokenizer2D, the model is composed of an Encoder followed by two outputs that assign weights to each input token.

All data was transformed to BoolQ dataset format [3] for Boolean model. We used PubMedQA dataset with the labels *yes/no/maybe* and BioASQ with *Yes/No* task. The training dataset was augmented with Beyond Back Translation [13] to reduce the strong imbalance between each class, then neutral contexts were additionally generated from text confusion to add material to those already existing, such that at most 30% of the data is synthetic for each class. The remaining sets did not have any augmentation.

*5.1.2 Metrics.* In text generation, semantic content can be generated with different sentence structures; thus, to evaluate the results, we compared them with the paragraph written by experts using the ROUGE metric [14] and report the $F$1-score for ROUGE-1, ROUGE-2 and ROUGE-L.

We based our extractive model evaluation on the average EM (Exact Match) score and the macro-F1 often used for this task. We also noticed that the predictions capture the true answers well but tend to add a larger or smaller margin around them; the Recall indicates the position that the relevant content takes in the predicted part.

Finally, we report Accuracy and Macro-F1 for Boolean models.

## 5.2 Results

*5.2.1 Extractive highlights.* Table 1 presents the results of the highlights $h_i$ extraction from each $(C, D_i)$ pair. Once $R$ is constructed as an output of the Tokenizer2D from these, the extractive approach builds a mapping over the whole document $D$ by highlighting all the important $h_i$ information scattered all over the large text. The

strength of the approach lies in capturing the indirect and implicit relationships between important terms. As a result, the subsequent Boolean model will no longer need to deal with irrelevant material to answer the question but will start directly from condensed information.

**Table 1: Results of extractive models.**

| Model | Scores | | | Statistics | |
|---|---|---|---|---|---|
| | EM | F1 | Recall | P | V (w/s) |
| $BERT_B$ | 40.00 | 58.35 | 72.34 | 110 | 1.49 |
| $BERT_L$ | 41.69 | 63.32 | 74.55 | 334 | 0.58 |
| $RoBERTa_B$ | 39.66 | 62.98 | 71.81 | 124 | 1.42 |
| $RoBERTa_L$ | 39.39 | 62.82 | 67.01 | 355 | 0.39 |
| MobileBERT | 41.46 | 60.14 | 71.42 | 25 | 2.03 **x4**[*] |
| TinyBERT | 37.40 | 58.35 | 71.02 | 6 | 6.50 **x11**[*] |

Parameters (P) in Millions. Inference speed (V) in windows per second (w/s), computed on CPU: Intel Xeon @ 2.20GHz
[*]Speed-Up compared to $BERT_L$

Among the architectures trained directly on the corpus, $BERT_{LARGE}$ performs best and infers faster than RoBERTa for the same level of complexity. At this stage, we only consider the encoders' performance. A direct benefit would be to take the most efficient one and lighten it to a less complex and faster version while keeping its knowledge. In favor of the inference speed, it perform at least as well as the base even if its distillations lower its scores.

The TinyBERT version particularly performs scores close to the base with a Recall not very different from its MobileBERT counterpart, but has the advantage of being lighter and faster.

This Recall difference is justified by the fact that the target $h_j$ are contained in the TinyBERT $\hat{h}_j$ but irrelevant tokens are added before and/or after. This performance loss is considered negligible as long as the additions are minor.

*5.2.2 Abstractive summarization.* We generated context summarization for each version of the model by playing with the Beam search's [5] depth and the n-grams size without repetition.

The $T5_{SMALL}$ version is understandably faster and, despite its lower complexity, its performance (Table 2) is no different from the basic version especially if we observe ROUGE-L.

In the generation, we want to avoid the insertion of insights that are not originally in the document $D$ in order not to bias the interpretation of the Boolean model. In our case, the more words the summary $R$ quotes from the article, the better. Nevertheless, this better output form is equivalent to performing extractive.

Compared to an Extractive output, the Abstractive texts are structured and fully intelligible to human. However, the downside is the loss of raw information to trace back to the original text, which is still an important point in our approach.

*5.2.3 Boolean Model and Combined System.* To fully meet our goal and based on the results of the two previous models, we plugged the Boolean model to the optimal extractive model (TinyBERT).

**Table 2: Results of Abstractive model.**

| Model | B | NR | ROUGE | | | Statistics | |
|---|---|---|---|---|---|---|---|
| | | | 1 | 2 | L | P | V(w/s) |
| Seq2Seq[*] | 5 | 3 | 28.9 | 5.4 | 23.1 | - | - |
| $T5_{SMALL}$ | 3 | 2 | 31.3 | 13.8 | 25.6 | 60 | 5.67 |
| | 5 | 4 | 31.4 | 14.4 | 25.9 | | |
| | 10 | 4 | 31.5 | 14.3 | 25.8 | | |
| $T5_{BASE}$ | 3 | 2 | 30.4 | 13.3 | 24.9 | 220 | 3.60 |
| | 5 | 4 | 31.5 | 14.3 | 26.0 | | |

Beam search (B). No-Repeat ngram size (NR). Parameters (P) in Millions. Generation speed (V) in windows per second (w/s), computed with GPU: NVIDIA GeForce GTX 1050
[*]Seq2Seq Multi-Task: performance on ELI5

It takes as input the concatenation of the question and the set of extracted windows outputs.

BERT and RoBERTa base version are used as encoder and we also experimented with the previously developed TinyBERT trained this time without distillation but directly on the raw data. We additionally reported the results (Table 3) by first running the text through one of the summary models before the best Boolean model. These are $BERT_L$ for extractive and $T5_S$ for abstractive.

**Table 3: Results of Boolean model.**

| Model | Scores | | Statistics | |
|---|---|---|---|---|
| | Accuracy | F1 | P | V (a/s) |
| $BERT_B$ | 96.85 | 96.01 | 110 | 3.45 |
| $RoBERTa_B$ | 97.32 | 97.07 | 125 | 3.61 |
| TinyBERT | 96.55 | 95.60 | 6 | 6.09 |
| Ext.QA + TinyBERT | 95.07 | 86.15 | 249 | - |
| Abs.QA + TinyBERT | 98.79 | 97.11 | 345 | - |

Parameters (P) in Millions. Inference speed (V) in article per seconde (a/s), computed on CPU: Intel Xeon @ 2.20GHz

The results in the first part of Table 3 show the performance of using the multiple encoders based on Roberta ($F1 = 97.07\%$) compared to Bert. Roberta being pre-trained to better encode the distributional representation of a sentence, it therefore reinforces the performance of the boolean question answering task. The direct training with the simple text representation of TinyBERT gives quite good scores considering its gap with our best model. We chose to quantize this model directly to speed it up to 9 items per second with a minor $F1$-score degradation of a 0.27% difference.

The second part of Table 3 presents the performance of the combined system: the boolean model plugged into the model of information summarization. We observe that abstractive model gives the best $F1$ score (97.11) compared to extractive model (86.15).

This performance gap reveals a specialization effect in the abstractive model that allows the boolean model to learn from richer

summary representations by bringing it more consolidated knowledge. On the other hand, the extract modeling does not consolidate knowledge but produce raw extracts which are consumed by the boolean model. Based on the results in Table 3, the contraction of information, sometimes with noise, performed by the abstractive leads the boolean model to better answer the question, in contrast to the extractive which reduces efficiency with respect to the unsummarized inputs. However, there is a trade-off between the performance, the traceability of the gathered information and the requirements of the task formalism.

Based on all this, the extractive model remains the best approach for our task and its combination with the boolean model provides a multi-task system that efficiently produces both:

- a boolean answer to a question;
- the extracts from the texts (not like a black box) supporting an answer;

These performances were not compared to related works as our task differs slightly from the usual fact checking task as indicated above.

## 6 CONCLUSION

By combining an extraction and answering system, we are able to find the solution to a closed question while justifying it from the source excerpts. The extractive part combined with our way of consolidating information through the produced extracts with Tokenizer2D allows to translate non-explicit information into a more concrete knowledge set. It then allowed both to give an answer to a question and to identify the relevant information for the answering system part.

Our experiments on the significant reduction of millions of model parameters made it possible to reduce the complexity of the final model, even if we can observe a slight decrease in performance. Which leads to a multi-task system that embeds the two TinyBERT models with an average error rate of 4%, which perform less than the one with a RoBERTa head at 2.68% error rate, but clearly lighter and faster.

The final system takes a query and returns two outputs for each article: the highlights and the Boolean answer.

We represent the scientific consensus on the question by using a simple majority vote of each article. Our approach is indeed to give a representation of the scientific consensus and not to give the absolute truth. However, this representation is limited by the availability of open access articles related to claims and by its quantitative aspect of the consensus with the aggregation.

We are currently work to improve the proposed approach by integrating the possibility for models to combine multiple and inevitably contradictory answers into a coherent and understandable response for humans.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Peter Achterberg, Willem de Koster, and Jeroen van der Waal. 2021. *A Science Confidence Gap: Education, Trust in Scientific Methods, and Trust in Scientific Institutions in the United States, 2014.* Springer International Publishing, Cham, 203–229.

[2] Joachim Allgaier. 2020. *Science and Medicine on YouTube.* Springer Netherlands, Dordrecht, 7–27. https://doi.org/10.1007/978-94-024-1555-1_1

[3] Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. BoolQ: Exploring the Surprising Difficulty of Natural Yes/No Questions. arXiv:1905.10044 [cs.CL]

[4] Angela Fan, Yacine Jernite, Ethan Perez, David Grangier, Jason Weston, and Michael Auli. 2019. ELI5: Long Form Question Answering. arXiv:1907.09190 [cs.CL]

[5] Markus Freitag and Yaser Al-Onaizan. 2017. Beam Search Strategies for Neural Machine Translation. *Proceedings of the First Workshop on Neural Machine Translation* (2017), 56–60. https://doi.org/10.18653/v1/w17-3207

[6] Omer Gendelman, Howard Amital, Nicola Luigi Bragazzi, Abdulla Watad, and Gabriel Chodick. 2020. Continuous hydroxychloroquine or colchicine therapy does not prevent infection with SARS-CoV-2: Insights from a large healthcare database analysis. *Autoimmunity Reviews* 19, 7 (2020), 102566. https://doi.org/10.1016/j.autrev.2020.102566 Special issue COVID19 and Autoimmunity.

[7] Benny Haerlin and Doug Parr. 1999. How to restore public trust in science. *Nature* 400, 6744 (01 Aug 1999), 499–499. https://doi.org/10.1038/22867

[8] Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. 2020. TinyBERT: Distilling BERT for Natural Language Understanding. arXiv:1909.10351 [cs.CL]

[9] Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. 2019. PubMedQA: A Dataset for Biomedical Research Question Answering. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, Vol. D19-1. 2567–2577.

[10] Georgi Karadzhov, Preslav Nakov, Lluis Marquez, Alberto Barron-Cedeno, and Ivan Koychev. 2017. Fully Automated Fact Checking Using External Sources. arXiv:1710.00341 [cs.CL]

[11] Emanuel Kulczycki. 2013. The Transformation of Science Communication in the Age of Social Media. *Teorie Vdy / Theory of Science* 35, 1 (2013), 3–28.

[12] Patrick Lewis, Ludovic Denoyer, and Sebastian Riedel. 2019. Unsupervised Question Answering by Cloze Translation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, 4896–4910. https://doi.org/10.18653/v1/P19-1484

[13] Zhenhao Li and Lucia Specia. 2019. Improving Neural Machine Translation Robustness via Data Augmentation: Beyond Back-Translation. *Proceedings of the 5th Workshop on Noisy User-generated Text (W-NUT 2019)* D19-55 (2019), 328–336. https://doi.org/10.18653/v1/d19-5543

[14] Chin-Yew Lin. 2004. ROUGE: A Package for Automatic Evaluation of Summaries. In *Text Summarization Branches Out.* Association for Computational Linguistics, Barcelona, Spain, 74–81.

[15] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, M. Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2020. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *ArXiv* abs/1907.11692 (2020).

[16] Zhenghao Liu, Chenyan Xiong, Maosong Sun, and Zhiyuan Liu. 2020. Fine-grained Fact Verification with Kernel Graph Attention Network. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, Vol. 2020.acl-main. Association for Computational Linguistics, Online, 7342–7351. https://doi.org/10.18653/v1/2020.acl-main.655

[17] Matthieu Mahévas, Viet-Thi Tran, Mathilde Roumier, Amélie Chabrol, Romain Paule, Constance Guillaud, Elena Fois, Raphael Lepeule, Tali-Anne Szwebel, François-Xavier Lescure, Frédéric Schlemmer, Marie Matignon, Mehdi Khellaf, Etienne Crickx, Benjamin Terrier, Caroline Morbieu, Paul Legendre, Julien Dang, Yoland Schoindre, Jean-Michel Pawlotsky, Marc Michel, Elodie Perrodeau, Nicolas Carlier, Nicolas Roche, Victoire de Lastours, Clément Ourghanlian, Solen Kerneis, Philippe Ménager, Luc Mouthon, Etienne Audureau, Philippe Ravaud, Bertrand Godeau, Sébastien Gallien, and Nathalie Costedoat-Chalumeau. 2020. Clinical efficacy of hydroxychloroquine in patients with covid-19 pneumonia who require oxygen: observational comparative study using routine care data. *BMJ* 369 (2020). https://doi.org/10.1136/bmj.m1844

[18] Sylvain Massip and Charles Letaillieur. 2020. Leveraging Open Access publishing to fight fake news. Zenodo. https://doi.org/10.5281/zenodo.3776797

[19] Tsvetomila Mihaylova, Georgi Karadjov, Pepa Atanasova, Ramy Baly, Mitra Mohtarami, and Preslav Nakov. 2019. SemEval-2019 Task 8: Fact Checking in Community Question Answering Forums. *arXiv:1906.01727 [cs, stat]* (May 2019). http://arxiv.org/abs/1906.01727 arXiv: 1906.01727.

[20] Sebastião Miranda, David Nogueira, Afonso Mendes, Andreas Vlachos, Andrew Secker, Rebecca Garrett, Jeff Mitchel, and Zita Marinho. 2019. Automated Fact Checking in the News Room. arXiv:1904.02037 [cs.CL]

[21] Jamal Abdul Nasir, Osama Subhani Khan, and Iraklis Varlamis. 2021. Fake news detection: A hybrid CNN-RNN based deep learning approach. *International Journal of Information Management Data Insights* 1, 1 (2021), 100007. https://doi.org/10.1016/j.jjimei.2020.100007

[22] Anastasios Nentidis, Konstantinos Bougiatiotis, Anastasia Krithara, and Georgios Paliouras. 2020. Results of the Seventh Edition of the BioASQ Challenge.

*Communications in Computer and Information Science* 1168 (2020), 553–568. https://doi.org/10.1007/978-3-030-43887-6_51

[23] Yotam Ophir and Kathleen Hall Jamieson. 2021. The effects of media narratives about failures and discoveries in science on beliefs about and support for science. *Public Understanding of Science* 0 (2021). https://doi.org/10.1177/09636625211012630 PMID: 34000907.

[24] Raul Puri and Bryan Catanzaro. 2019. Zero-shot Text Classification With Generative Language Models. arXiv:1912.10165 [cs.CL]

[25] Raul Puri, Ryan Spring, Mostofa Patwary, Mohammad Shoeybi, and Bryan Catanzaro. 2020. Training Question Answering Models From Synthetic Data. arXiv:2002.09599 [cs.CL]

[26] 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. arXiv:1910.10683 [cs.LG]

[27] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ Questions for Machine Comprehension of Text. arXiv:1606.05250 [cs.CL]

[28] Jodi Schneider, Michele Avissar-Whiting, Caitlin Bakker, Hannah Heckner, Sylvain Massip, Randy Townsend, and Nathan D. Woods. 2021. Addressing disorder in scholarly communication: Strategies from NISO 2021. *Information Services & Use* (Sept. 2021), 1–15. https://doi.org/10.3233/ISU-210113

[29] Cassidy R. Sugimoto, Mike Thelwall, Vincent Larivière, Andrew Tsou, Philippe Mongeon, and Benoit Macaluso. 2013. Scientists Popularizing Science: Characteristics and Impact of TED Talk Presenters. *PLOS ONE* 8, 4 (04 2013), 1–8. https://doi.org/10.1371/journal.pone.0062403

[30] Zhiqing Sun, Hongkun Yu, Xiaodan Song, Renjie Liu, Yiming Yang, and Denny Zhou. 2020. MobileBERT: a Compact Task-Agnostic BERT for Resource-Limited Devices. arXiv:2004.02984 [cs.CL]

[31] James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. FEVER: a large-scale dataset for Fact Extraction and VERification. arXiv:1803.05355 [cs.CL]

[32] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention Is All You Need. arXiv:1706.03762 [cs.CL]

[33] David Wadden, Shanchuan Lin, Kyle Lo, Lucy Lu Wang, Madeleine van Zuylen, Arman Cohan, and Hannaneh Hajishirzi. 2020. Fact or Fiction: Verifying Scientific Claims. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, Online, 7534–7550. https://doi.org/10.18653/v1/2020.emnlp-main.609

[34] Yang Yang, Lei Zheng, Jiawei Zhang, Qingcai Cui, Zhoujun Li, and Philip S. Yu. 2018. TI-CNN: Convolutional Neural Networks for Fake News Detection. arXiv:1806.00749 [cs.CL]

[35] Wonjin Yoon, Jinhyuk Lee, Donghyeon Kim, Minbyul Jeong, and Jaewoo Kang. 2020. Pre-trained Language Model for Biomedical Question Answering. In *Machine Learning and Knowledge Discovery in Databases*, Peggy Cellier and Kurt Driessens (Eds.), Vol. CCIS 1168. Springer International Publishing, Cham, 727–740.

[36] Jie Zhou, Xu Han, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, and Maosong Sun. 2019. GEAR: Graph-based Evidence Aggregating and Reasoning for Fact Verification. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, Vol. 0. Association for Computational Linguistics, Florence, Italy, 892–901. https://doi.org/10.18653/v1/P19-1085