

Multilingual Previously Fact-Checked Claim Retrieval

Anonymous ACL submission

Abstract

Fact-checkers are often hampered by the sheer amount of online content that needs to be fact-checked. NLP can help them by retrieving already existing fact-checks relevant to the content being investigated. This paper introduces a new multilingual dataset – *MultiClaim* – for previously fact-checked claim retrieval. We collected 28k posts in 27 languages from social media, 206k fact-checks in 39 languages written by professional fact-checkers, as well as 31k connections between these two groups. This is the most extensive and the most linguistically diverse dataset of this kind to date. We evaluated how different unsupervised methods fare on this dataset and its various dimensions. We show that evaluating such a diverse dataset has its complexities and proper care needs to be taken before interpreting the results. We also evaluated a supervised fine-tuning approach, improving upon the unsupervised method significantly.

1 Introduction

Fact-checking organizations have made progress in recent years in manually and professionally fact-checking viral content (Micallef et al., 2022; Full Fact, 2020). To reduce some of the fact-checkers’ manual efforts and make their work more effective, several studies have recently examined their needs and identified tasks that could be automated (Nakov et al., 2021; Full Fact, 2020; Micallef et al., 2022; Dierickx et al., 2022; Hrkova et al., 2022). These tasks include searching for the source of evidence for verification, searching for other versions of misinformation, and searching within existing fact-checks. These tasks were identified as particularly challenging for fact-checkers working in low-resource languages (Hrkova et al., 2022).

In this work, we focus on *previously fact-checked claim retrieval* (PFCR) (Shaar et al., 2020). Given a text making an *input claim* (e.g., a social media post) and a set of *fact-checked claims*, our task is to

rank the *fact-checked claims* so that those that are the most relevant w.r.t. the *input claim* (and thus the most useful from the fact-checker’s perspective) are ranked as high as possible.

Previously, this task was mostly done in English. Other languages that have been considered include Arabic (Nakov et al., 2022), Bengali, Hindi, Malayalam, and Tamil (Kazemi et al., 2021). However, many other languages or even entire major language families have not been considered at all. Additionally, so far only *monolingual PFCR* has been tackled, when the input claim and the fact-checked claims are in the same language. To address these shortcomings, we introduce in this paper a new extensive multilingual dataset. Our two main contributions are:

1. *MultiClaim* – Multilingual dataset for PFCR.

We collected and made available¹ a novel multilingual dataset for PFCR. The dataset consists of 205,751 fact-checks in 39 languages and 28,092 social media posts (from now on just *posts*) in 27 languages. For most of these languages, this is the first time this task has been considered at all. This is also the biggest dataset of fact-checks released to date.

All the posts were previously reviewed by professional fact-checkers who also assigned appropriate fact-checks to them. We collected these assignments and gathered 31,305 pairs consisting of a post and a fact-check reviewing the claim made in the post. 4,212 of these pairs are crosslingual (i.e., the language of the fact-check and the language of the post are different). This dataset introduces *crosslingual PFCR* as a new task that has not been tackled before. This is the biggest collection of such pairs that were confirmed by professional fact-checkers. The dataset also includes OCR transcripts of the images attached to the posts

¹The dataset and code are available at Zenodo. Data are available upon request *for research purposes only: anonymized*.

and machine translation of all the data into English.

2. In-depth multilingual evaluation. We evaluated the performance of various text embedding models and BM25 for both the original multilingual data and their English translations. We describe several pitfalls related to the complexity of evaluating such a linguistically diverse dataset. We also explore the performance across several other data dimensions, such as post length or publication date. Finally, we show that we can improve text embedding methods further by using supervised training with our data.

2 Related Work

Other names are used for PFCR or similar tasks for various reasons, e.g., fact-checking URL recommendation (Vo and Lee, 2018), fact-checked claims detection (Shaar et al., 2020), verified claim retrieval (Barrón-Cedeño et al., 2020), searching for fact-checked information (Vo and Lee, 2020), or claim matching (Kazemi et al., 2021).

Datasets. *CheckThat!* datasets (Barrón-Cedeño et al., 2020; Shaar et al., 2021) have the most similar collection approach to ours. They collect English and Arabic tweets mentioned in fact-checks to create preliminary pairs and then manually filter them. Compared to this work, we broaden the scope of data collection and omit the manual cleaning in favor of using fact-checkers’ reports. Shaar et al. (2020) collected data from fact-checking of English political debates done by fact-checkers. The *CrowdChecked* dataset (Hardalov et al., 2022) was created by searching for fact-check URLs on Twitter and collecting English tweets from retrieved threads. The process is inherently noisy and, the authors propose different noise filtering techniques.

Kazemi et al. (2021) collected several million chat messages from public chat groups and tiplines in English, Bengali, Hindi, Malayalam, and Tamil and 150k fact-checks. Then they sampled roughly 2,300 pairs based on their embedding similarity and manually annotated them. In the end, they obtained only roughly 250 positive pairs. Jiang et al. (2021) matched COVID-19 tweets and 90 COVID-19 claims in a similar manner. Their data could be used for PFCR, but the authors worked on classification instead.

PFCR datasets are summarized in Table 1. Our dataset has the highest number of fact-checked

| | Input claims | FC claims | Pairs | Languages |
|-------------------------|--------------|-----------|---------|-----------|
| Kazemi et al., 2021 | NA | 150,000 | 258 | 5 |
| Jiang et al., 2021 | NA | 90 | 1,573 | 1 |
| Shaar et al., 2020 | NA | 27,032 | 1,768 | 1 |
| Shaar et al., 2021 | 2,259 | 44,164 | 2,440 | 2 |
| Hardalov et al., 2022 | 316,564 | 10,340 | 332,660 | 1 |
| <i>MultiClaim</i> (our) | 28,092 | 205,751 | 31,305 | 27/39 |

Table 1: PFCR datasets. FC claims are *fact-checked*. NA means that we were not able to identify the correct number of input claims. The number should be similar to the number of pairs in most cases.

claims. It also has the second-highest number of input claims and pairs after *CrowdChecked*, but that dataset is significantly noisier. Finally, our dataset has by far the most languages, while the second biggest dataset in this regard has 5 language with only 50 samples per language.

Methods. Methods used for PFCR are usually either BM25 (and other similar information retrieval algorithms) or various text embedding-based approaches (Vo and Lee, 2018; Shaar et al., 2022a,b, i.a.). Reranking is often used to combine several methods to side-step compute requirements or as a sort of ensembling (Shaar et al., 2020, i.a.). PFCR task is also a target of the *CLEF’s CheckThat!* challenge, with many teams contributing with their solutions (Nakov et al., 2022). Other methods use visual information from images (Mansour et al., 2022; Vo and Lee, 2020), abstractive summarization (Bhatnagar et al., 2022), or key sentence identification (Sheng et al., 2021) to improve the results.

3 Our Dataset

Our dataset *MultiClaim* consists of fact-checks, social media posts and pairings between them.

Fact-checks. We have collected the majority of fact-checks listed in the Google Fact Check Explorer, as well as fact-checks from additional manually identified major sources (e.g., Snopes) that were missing. Overall, we have collected 205,751 fact-checks from 142 fact-checking organizations covering 39 languages. We publish the *claim*, *title*, *publication date*, and *URL* of each fact-check. We do not publish the full body of the articles. The claim is usually (in 88.2% of the cases) a one sentence long summarization of the information being fact-checked.

Social media posts. We used two ways to find relevant social media posts from Facebook, Instagram and Twitter. In both cases, it was professional

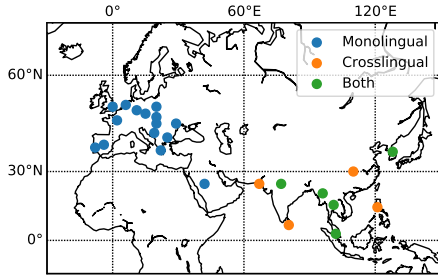


Figure 1: Major languages from our dataset. Crosslingual languages all have English fact-checks.

fact-checkers that assigned the fact-checks to the posts. (1) Some fact-checks use the *ClaimReview* schema², which has a field for reviewed items. All the links to the three social media platforms from this field are used to collect the posts and form the *pairs*. (2) We searched for URLs to Facebook and Instagram in the main body of the fact-checks. This is our pool of potentially relevant posts. Then we use the fact-checking warnings these two platforms provide. These warnings contain links to relevant fact-checking articles. We use these links to establish additional *pairs*.

In total, we collected 28,092 posts from 27 languages. There are 31,305 *fact-check-to-post pairs*, each post in our dataset is paired with at least one fact-check. 26,774 of these pairs are monolingual and 4,212 are crosslingual (as predicted by the language identification, see below). Figure 1 shows the major (more than 100 samples) languages. All the crosslingual cases have the visualized language for posts and English for fact-checks. We can see that there is a clear distinction between these two groups, probably caused by different fact-checking cultures in different regions.

We publish the *text*, *OCR of the attached images (if any)*, *publication date*, *social media platform*, and *fact-checker’s rating* of each post. The *rating* is the reason why the post was flagged (see Section 4.2 for more details). We do not publish URLs in an effort to protect the users and their privacy as much as possible. For detailed information about the implementation of this dataset collection pipeline, see Appendix B. For a more detailed breakdown of dataset statistics (by languages and sources), see Appendix C. Examples from our dataset can be seen in Appendix G.

²<https://schema.org/ClaimReview>

Dataset versions. We machine-translated all the published texts into English, resulting in two parallel versions of our dataset: the *original version* and the *English version*. We also used automatic language identification on all the texts. Both translations and language identifications are published as well.

Noise ratio. We manually checked 100 randomly selected pairs from our dataset and evaluated their validity. Three authors rated these pairs and assessed whether the claim from the fact-check was made in the post. In case of disagreement, they discussed the annotation until an agreement was reached. Based on our assessment, 87 out of 100 pairs were correct. The remaining 13 pairs were not errors made by social media platforms or fact-checkers, but rather posts that required visual information (either from video or image) to fully match the assigned fact-check. The 95% Agresti-Coull confidence interval (Agresti and Coull, 1998) for correct samples in our dataset is 79-92%.

4 Unsupervised Evaluation

We formulate the task we are solving with our dataset as a ranking task, i.e., for each post, the methods rank all the fact-checks. Then, we evaluate the performance based on the rank of the desired fact-checks by using success-at-K (S@K) as the main evaluation metric. We define as the percentage of pairs when the desired fact-check ends up in the top K. Throughout the paper, we report this metric with the 95% Agresti-Coull confidence interval.

For unsupervised evaluation, we evaluated text embedding models and the BM25 algorithm to understand how they are able to handle pairs in different languages or even crosslingual pairs. Fact-checks are represented with their claims only. Posts are represented with their main texts concatenated with the OCR transcripts. We use either the original texts or their English translations, depending on the version of the dataset that is reported.

Text embedding models (TEMs). We use various neural TEMs (Reimers and Gurevych, 2019) that encode texts into a vector space. These are usually based on pre-trained transformer language models fine-tuned as Siamese networks to generate well-performing text embeddings. We use these models to embed both social media posts and fact-checked claims into a common vector space. The

retrieval is then reduced to calculating and sorting distances between vectors.

BM25. With BM25 (Robertson and Zaragoza, 2009), we use the posts as queries and fact-checked claims as documents. The score is then calculated based on the lexical overlap between the query and all the documents.

4.1 Main Results

We compare the performance of 15 English TEMs, 5 multilingual TEMs, and BM25. The English TEMs were only evaluated with the *English* version. The multilingual TEMs and BM25 were evaluated with both the *original* and the *English* versions. BM25 with different versions will be denoted as BM25-Original and BM25-English, respectively.

In this section, we use different strategies to evaluate monolingual and crosslingual pairs. For monolingual pairs, we only search within the pool of fact-checks written in the same language as the post (e.g., for a French post we only rank the French fact-checks). For crosslingual pairs, we search in all the fact-checks³. In both cases, we report the average performance for individual languages. We only report for languages with more than 100 pairs. For crosslingual pairs, we also consider a separate *Other* category for all the leftover pairs.

We present the main results in Table 2 and we visualize them in Figure 2. We conclude that: **(1)** English TEMs are the best performing option for both monolingual and crosslingual claim retrieval. **(2)** Machine translation significantly improved the performance of both BM25 and TEMs. The difference between the best performing *English* version method and the best performing *original* version method is 35% for crosslingual and 14% for monolingual S@10. Currently, machine translation systems also have better language coverage than multilingual TEMs. **(3)** TEMs have a strong correlation between monolingual and crosslingual performance (Pearson’s $\rho = 0.98$, $P = 4e-10$ for English TEMs). These two capabilities do not conflict. **(4)** There is almost no correlation (Pearson’s $\rho = 0.03$, $P = 0.89$ for English TEMs) between model size and performance. The training procedure is much more important. GTR is an exceptionally well-performing family, with all three models being Pareto optimal w.r.t. model size and performance. Another notable model is MiniLM – a surprisingly powerful model for its size (33M).

³The index created for BM25 is multilingual as well.

| Method | Size [M] | Ver. | Mono | Cross | SLB |
|-----------------------------|----------|------|-------------|-------------|------|
| BM25 | | | | | |
| | | En | 0.78 | 0.39 | 0.18 |
| | | Og | 0.62 | 0.06 | 0.69 |
| English TEMs | | | | | |
| DistilRoBERTa | 82 | En | 0.76 | 0.43 | 0.18 |
| GTR-T5-Base | 110 | En | 0.81 | 0.51 | 0.19 |
| GTR-T5-Large | 336 | En | 0.83 | 0.56 | 0.20 |
| GTR-T5-XL | 1242 | En | 0.83 | 0.56 | 0.20 |
| MPNet-Base | 109 | En | 0.78 | 0.47 | 0.18 |
| MSMARCO-BERT-Base | 109 | En | 0.78 | 0.46 | 0.18 |
| MiniLM-L12 | 33 | En | 0.80 | 0.48 | 0.18 |
| MultiQA-MPNet-Base | 109 | En | 0.80 | 0.50 | 0.18 |
| SGPT-125M | 125 | En | 0.63 | 0.25 | 0.14 |
| SGPT-2.7B | 2700 | En | 0.77 | 0.50 | 0.19 |
| Sentence-T5-Base | 110 | En | 0.73 | 0.37 | 0.14 |
| Sentence-T5-Large | 336 | En | 0.75 | 0.41 | 0.15 |
| Sentence-T5-XL | 1242 | En | 0.78 | 0.47 | 0.16 |
| Multilingual TEMs | | | | | |
| DistilUSE-Base-Multilingual | 135 | En | 0.74 | 0.40 | 0.15 |
| | | Og | 0.66 | 0.20 | 0.16 |
| LaBSE | 472 | En | 0.63 | 0.22 | 0.13 |
| | | Og | 0.69 | 0.22 | 0.17 |
| MPNet-Base-Multilingual | 278 | En | 0.75 | 0.41 | 0.16 |
| | | Og | 0.70 | 0.21 | 0.17 |
| MiniLM-L2-Multilingual | 118 | En | 0.74 | 0.38 | 0.16 |
| | | Og | 0.63 | 0.15 | 0.17 |
| XLM-R | 278 | En | 0.72 | 0.33 | 0.15 |
| | | Og | 0.66 | 0.15 | 0.16 |

Table 2: Results for methods showing both *monolingual* and *crosslingual* S@10. Ver. denotes either the *original* (Og) or the *English* (En) version of our dataset. The best results for these two versions are bolded. SLB denotes *same language bias*.

Languages. Performance for individual languages is shown in Figure 3. We show the results for the best performing TEMs for both versions (**GTR-T5-Large** for the *English* and **MPNet-Base-Multilingual** for the *original*, which are denoted as GTR-T5 and MPNet from now on) and both BM25s. We cannot directly compare the performance numbers across different monolingual languages, since they use different pools of fact-checks with different sizes. This is also why smaller languages seem to have better scores.

BM25-Original, despite its seemingly weak overall performance, is actually competitive in some languages, e.g., Spanish, Portuguese, or Malay. It is better than multilingual TEMs for 7 out of 20 monolingual cases. Its overall monolingual performance is significantly decreased by Thai and Myanmar, due to their use of *scriptio continua*. On the other hand, unlike multilingual TEMs, BM25-Original is by design not capable of any crosslingual retrieval and the results are shown only for completeness.

False positive rate. We noticed that BM25-Original seems to perform better for languages with larger fact-check pools. We conducted an experiment to measure how pool size affects the results. We randomly selected 100 pairs for 7 of our languages with the largest fact-check pools. We then measured the performance for these 100 pairs

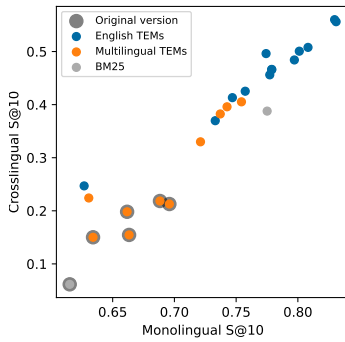


Figure 2: Comparison of different method families. Unless stated otherwise, the methods use *English* version of our dataset.

while increasing the pool size from 100 to 2,100 by gradually adding random fact-checks.

We found that our initial observation was correct and that BM25-Original performs better than the MPNet model as the pool size increases (especially for Spanish, Portuguese, and French). The relative comparison between BM25 methods and TEMs is shown in Figure 4. This suggests, that MPNet has a higher *false positive rate*, i.e., it is more likely to assign high scores to irrelevant fact-checks. As the number of fact-checks grows, the risk of selecting irrelevant fact-checks also grows. **Different methods may be appropriate for different languages based on the number of fact-checks available.** We did not find the same pattern when comparing the methods using the *English* version.

Same language bias. The fact that we reduce the fact-checks pool to one language in monolingual evaluation is motivated by what we call *same language bias* (SLB) – a tendency of methods to retrieve fact-checks that have the same language as the post. We approximate SLB by calculating the percentage of top 10 fact-checks that have the same language as the input post when we use the full pool. This number is reported in Table 2.

BM25-Original has the highest SLB score of all the methods, as it has an implicit language filter that effectively removes fact-checks from other languages from the pool. This reduction makes the task easier, but it violates our requirement that the method should take fact-checks in all the languages into consideration. We used language-filtered fact-checks in monolingual evaluation to reduce the effect the SLB has on the results. Without this filtering, BM25-Original would clearly outperform MPNet (S@10 51.9 vs 38.5), even though our re-

sults in Figure 3 show that for many languages, its language understanding capabilities are actually worse.

However, it is not necessarily true that a higher SLB leads to worse crosslingual performance. As shown in Figure 5, TEMs with the highest SLB actually have the best performance for crosslingual evaluation. Even more strikingly, the relative crosslingual performance compared to monolingual performance increases with SLB as well. We theorize that a certain amount of SLB is healthy, as long as the methods focus on meaningful similarities in texts written in the same language, such as local topics, named entities, and events, rather than on superfluous lexical overlaps. SLB can also be useful to localize claims that are not specific enough. For example, it is impossible to identify the country of origin for the following claim translated to English: *Educational institutions are re-opening from January 18*. However, as soon as we know that the original language was Bengali, we can guess that it is about Bangladeshi institutions.

4.2 Other Dimensions

In this section, we report results for various data splits. Since we often work with small splits, we are not able to report the results as an average per language as in the previous section. Instead, we report the average score across the samples. This will give more weight to the more common languages, penalizing the methods with high *false positive rate* (e.g., multilingual TEMs).

Time. We grouped the posts for which we were able to obtain the publication date (N = 26,337) into 20-quantiles and measured the performance of individual methods. The results are shown in Figure 6. There is a visible drop-off for all the methods at the start of 2020, largely caused by the COVID-19 pandemic. We confirmed this by measuring how well the methods worked on posts with the substrings *corona*, *covid* or *korona*.⁴ The results are shown in Table 3 (top panel).

The relative differences between individual methods seem stable. We hypothesized that TEMs might have problems with aging, since many of the foundation language models were originally trained before 2020. We correlated the average post time for each quantile with the difference be-

⁴We chose this as a very simple, high-precision filtering technique. Many other COVID-19-related posts were not retrieved.

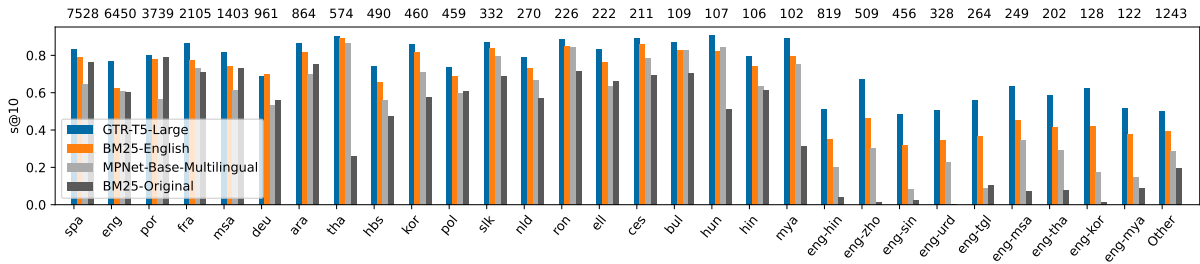


Figure 3: Performance of selected methods for individual languages. For crosslingual pairs (e.g., *eng-hin*), the first language is for the fact-checks and the second is for the posts. The number of pairs is shown at the top.

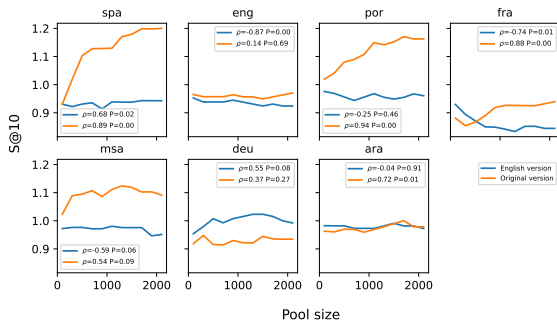


Figure 4: Relative performance ($S@10$) between BM25 methods and TEMs for different fact-check pool sizes. For both versions we compare the best performing TEMs (GTR-T5 and MPNet) with BM25. Positive ρ means that BM25 gets better with the growing pool size.

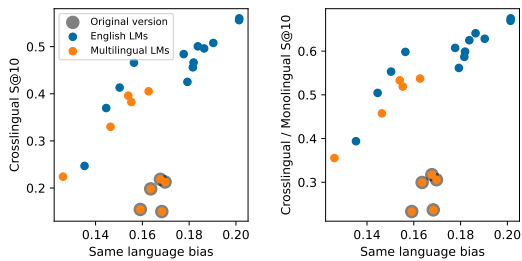


Figure 5: Relation between *same language bias* and performance for TEMs.

tween GTR-T5 and BM25-English performance and found a negative, but statistically insignificant correlation (Pearson $\rho = -0.33$, $P = 0.17$ for monolingual $S@10$). Similar results were measured for crosslingual performance. In both cases, the direction signals that the GTR model is indeed getting worse over time. We found no such signal comparing methods using the *original* version.

There is a risk that the fact-check was written based on the very post we are using, and an information leak might have happened (e.g., the fact-checker might have used parts of the post verbatim). To test this, we compared pairs where the post is

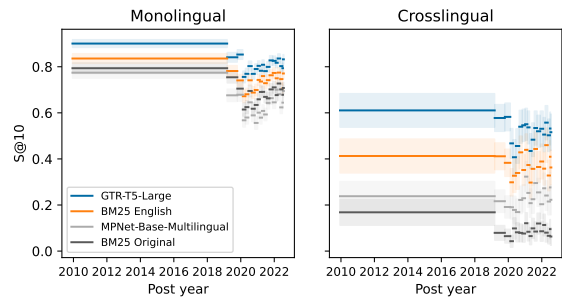


Figure 6: Performance of selected methods for posts from different time intervals. Shaded areas are confidence intervals.

newer with the pairs where the post is older. We found that the two groups have virtually the same performance for all the methods (e.g., 80.02 vs 80.04 monolingual $S@10$ for GTR-T5). If there is an information leak happening, we were not able to measure it.

Post rating. In the case of Facebook and Instagram posts, fact-checkers use the so-called *ratings* to describe the type of fallacy present. We show the results for the most common ratings in Table 3 (middle panel). *Missing context* has a slightly lower score than *(Partially) False information*. This might be caused by the fact that the rating is defined by what is *not* written in the post, making it harder to match with an appropriate fact-check. *Altered photo / video* rating has an even lower score. This is an expected behavior, since our purely text-based models cannot handle cases when the crux of the post is in its visual aspect.

Post length. We show how the length of the posts influence the results in Figure 8. In general, the performance peaks at around 500 characters. Posts that are too short are too difficult to match (and extremely short posts may even indicate noise in the data). On the other hand, for posts longer than 500 characters, the methods gradually lose their ef-

| | Monolingual | | | | | Crosslingual | | | | |
|--------------------------|-------------|-------------|-------------|-------------|-------------|--------------|-------------|-------------|-------------|-------------|
| | N | GTR-T5 | BM25-En | MPNet | BM25-Og | N | GTR-T5 | BM25-En | MPNet | BM25-Og |
| COVID-related | 4159 | 0.72 ± 0.01 | 0.68 ± 0.01 | 0.50 ± 0.02 | 0.60 ± 0.01 | 514 | 0.40 ± 0.04 | 0.29 ± 0.04 | 0.17 ± 0.03 | 0.06 ± 0.02 |
| Otherwise | 22615 | 0.83 ± 0.00 | 0.75 ± 0.01 | 0.66 ± 0.01 | 0.70 ± 0.01 | 3698 | 0.55 ± 0.02 | 0.39 ± 0.02 | 0.23 ± 0.01 | 0.08 ± 0.01 |
| False information | 14812 | 0.82 ± 0.01 | 0.75 ± 0.01 | 0.65 ± 0.01 | 0.69 ± 0.01 | 2155 | 0.52 ± 0.02 | 0.37 ± 0.02 | 0.22 ± 0.02 | 0.09 ± 0.01 |
| Partly false information | 4498 | 0.82 ± 0.01 | 0.75 ± 0.01 | 0.63 ± 0.01 | 0.70 ± 0.01 | 669 | 0.53 ± 0.04 | 0.39 ± 0.04 | 0.21 ± 0.03 | 0.08 ± 0.02 |
| Missing context | 1993 | 0.77 ± 0.02 | 0.70 ± 0.02 | 0.61 ± 0.02 | 0.63 ± 0.02 | 268 | 0.53 ± 0.06 | 0.35 ± 0.06 | 0.19 ± 0.05 | 0.05 ± 0.03 |
| Altered photo/video | 753 | 0.73 ± 0.03 | 0.66 ± 0.03 | 0.52 ± 0.04 | 0.64 ± 0.03 | 142 | 0.47 ± 0.08 | 0.34 ± 0.08 | 0.17 ± 0.06 | 0.12 ± 0.05 |
| Facebook | 24668 | 0.81 ± 0.00 | 0.74 ± 0.01 | 0.64 ± 0.01 | 0.68 ± 0.01 | 3927 | 0.52 ± 0.02 | 0.37 ± 0.02 | 0.22 ± 0.01 | 0.08 ± 0.01 |
| Instagram | 1473 | 0.78 ± 0.02 | 0.74 ± 0.02 | 0.56 ± 0.03 | 0.75 ± 0.02 | 44 | 0.56 ± 0.14 | 0.37 ± 0.14 | 0.19 ± 0.11 | 0.19 ± 0.11 |
| Twitter | 682 | 0.84 ± 0.03 | 0.74 ± 0.03 | 0.69 ± 0.03 | 0.70 ± 0.03 | 244 | 0.64 ± 0.06 | 0.49 ± 0.06 | 0.38 ± 0.06 | 0.06 ± 0.03 |
| Total | 26774 | 0.81 ± 0.00 | 0.74 ± 0.01 | 0.64 ± 0.01 | 0.68 ± 0.01 | 4212 | 0.53 ± 0.02 | 0.38 ± 0.01 | 0.23 ± 0.01 | 0.08 ± 0.01 |

Table 3: Performance (S@10) with confidence intervals for various splits and methods.

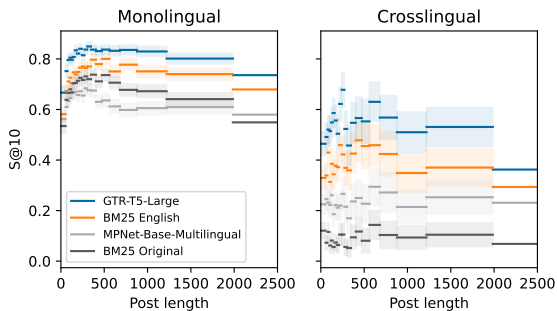


Figure 7: Performance of selected methods for posts with different lengths. Shaded areas are confidence intervals.

fectiveness. The relative performance of methods seems to be relatively stable.

Social media platforms. The results for social media platforms are in Table 3 (bottom pannel). We can see that Twitter has the best performance overall. We believe that this is, to a large extent, caused by the limited length of the Twitter posts.

5 Supervised Training

To validate that our dataset can be used as a training set, we fine-tuned TEMs and evaluated their performance. We split the posts randomly into 80:10:10% train, development, and test sets. We used *cosine* or *contrastive* training losses to fine-tune the models. In both cases, both positive and negative pairs are required for training. We used our data as positive samples and random pairs as negative samples. We performed a hyperparameter search with GTR-T5 and MPNet TEMs (see §D). Here, we report the best performing fine-tuned model we were able to achieve for both TEMs.

The overall results for the test set are reported in Table 4. We can see that GTR-T5 achieved only modest improvements⁵. On the other hand,

⁵This is largely caused by the fact that we had to use

MPNet improved significantly in both monolingual and crosslingual performance, even surpassing the performance of BM25-English. We observed that the improvements were global across all languages.

We also observed that the TEMs were able to saturate the training set quite quickly, achieving 99.5%+ average precision after only a few epochs. This shows that our naive random selection of negative samples was too easy. The model can learn only a limited amount of information from such samples, and we would need a more elaborate scheme for generating more challenging negative samples. This could lead to further performance improvements.

6 Post-Hoc Results Analysis

The pairs, we obtained from the fact-checks, are only a subset of all the potentially valid pairs. This incompleteness limits our understanding of the dataset and also our evaluation. We decided to manually annotate a subset of the results generated by the methods to better understand what is missing from our data. We generated the top 10 fact-checks for the 87 test set posts that we knew had valid fact-checks (see §3). We used the 4 unsupervised and 2 supervised methods from Section 5.

These methods generated 3,390 unique pair predictions for these 87 posts. Three authors went through each prediction and marked, whether they agreed with it, i.e., whether they found the fact-check to be valid and useful for the post. The agreement rates between the annotators were sufficiently high: 82.2%, 85.5% and 92.9%. We consider pairs where at least two annotators agreed to be *correct*. In total, the methods were able to find 719 correct pairs. 96 of these were present

smaller batch size due to (1) this model being larger, and (2) this model not supporting *average mixed precision*. We believe that with larger batch size the performance could be even better.

| Model | Section 5 (S@10) | | Section 6 (S@10) | | Section 6 (R@10) | |
|-------------------------|------------------|--------------|------------------|-------------|------------------|-------------|
| | Monolingual | Crosslingual | Our dataset | Annotated | Our dataset | Annotated |
| Unsupervised | | | | | | |
| GTR-T5-Large | 0.82 ± 0.01 | 0.55 ± 0.05 | 0.70 ± 0.09 | 0.93 ± 0.05 | 0.69 ± 0.09 | 0.59 ± 0.04 |
| BM25-English | 0.74 ± 0.02 | 0.40 ± 0.05 | 0.67 ± 0.10 | 0.85 ± 0.07 | 0.67 ± 0.09 | 0.48 ± 0.04 |
| MPNet-Base-Multilingual | 0.63 ± 0.02 | 0.23 ± 0.04 | 0.51 ± 0.10 | 0.70 ± 0.09 | 0.47 ± 0.09 | 0.32 ± 0.03 |
| BM25 Original | 0.68 ± 0.02 | 0.09 ± 0.03 | 0.60 ± 0.10 | 0.71 ± 0.09 | 0.58 ± 0.09 | 0.26 ± 0.03 |
| Supervised | | | | | | |
| GTR-T5-Large | 0.84 ± 0.01 | 0.59 ± 0.05 | 0.71 ± 0.09 | 0.92 ± 0.05 | 0.70 ± 0.09 | 0.65 ± 0.03 |
| MPNet-Base-Multilingual | 0.76 ± 0.02 | 0.42 ± 0.05 | 0.62 ± 0.10 | 0.85 ± 0.07 | 0.60 ± 0.09 | 0.45 ± 0.04 |

Table 4: Test set performance (§5) and annotated results performance (§6) of unsupervised and supervised methods.

in our original dataset. This suggests that there is roughly $7\times$ more pairs in our dataset than we had previously identified. No method was able to find 9 fact-checks out of 105 that were already in our dataset. Of the 719 correct pairs, only 247 were monolingual, 136 were crosslingual with an English fact-check, and 336 were crosslingual with a non-English fact-check. The last category in particular is almost completely missing from our dataset.

In Table 4, we show the results for individual methods. We compare S@10 (now defined as how many posts have at least one correct fact-check produced in the top 10) as approximated with our dataset and the true S@10 obtained by the annotation. We can see that the score for our dataset is significantly lower and the true performance of our methods is better than what was measured previously. We also compare recall-at-10 (R@10), defined as the percentage of expected pairs a method was able to produce in the top 10. In this case, both our dataset and manual annotation are only estimates, since they do not contain *all* the valid pairs, they both contain only a subset obtained by different methods. Here we can see that our dataset actually provides higher estimates. We assume that our annotation is more precise, so we conclude that the recall calculated from our dataset is overinflated (possibly due to selection bias). **It also seems that our dataset has a bias in favor of BM25**, compared to the results obtained from annotated data.

7 Discussion

Complexity of crosslingual evaluation. Phenomena such as *same language bias* or *false positive rate* make the evaluation of multilingual and crosslingual datasets inherently complex. If we were to abstract the whole evaluation into a single number, as is often done in practice, we would have completely missed these pitfalls. Without an in-depth evaluation, we might have been misled while applying our methods in practice, e.g., while developing helpful tools for fact-checkers. Our

evaluation procedures were previously impossible to develop in the absence of linguistically diverse PFCR datasets.

Machine translation beats multilingual TEMs.

These two technologies represent the two main multilingual and crosslingual learning paradigms – label transfer and parameter transfer (Pikuliak et al., 2021). Machine translation is a clear winner in our case. English TEMs significantly outperform multilingual approaches for both monolingual and crosslingual retrieval.

COVID-19. As shown in Table 3, it seems that the performance for COVID-19 is significantly worse than for the rest of the dataset. However, this might not necessarily mean that the methods are having issues with the domain shift. The sheer amount of fact-checks written about COVID-19 makes it hard for the methods to pick the desired fact-check in the presence of thousands of other very similar ones. This is evident considering that BM25 also has worse results, even though it should be less prone to domain shift based on its design.

8 Conclusions

In this paper, we introduced a new multilingual *previously fact-checked claim retrieval* dataset. Our collection process yielded a unique and diverse dataset with a relatively small amount of noise in it. We believe that the evaluation of various methods is also insightful and can lead to the development of better fact-checking tools in the future.

We believe that our dataset opens up many interesting research directions. We have, for example, barely scraped the surface of crosslingual learning in this work. Applying various transfer learning methods (especially for low-resource languages) is an important future direction.

9 Limitations

9.1 Dataset

Noise. Based on our annotation (see §3), we expect that around 13% of posts in our dataset do not contain the claim in the textual form. These are the cases when the claim being made on the social media is based on visual information. Note, that the methods might still be able to retrieve correct fact-checks for some of these posts, based on spurious correlations, e.g., overlaps in named entities.

AI APIs. We use out-of-the-box AI services to perform optical character recognition, machine translation to English and language detection. All of these have limited precision and might inject noise into our data.

- *OCR* was too sensitive and was often reading imaginary character, watermarks, etc. We had to address this by a more aggressive text cleaning.
- *Machine translation to English* is not perfect and the quality of translations depends on source language, particular topics or even the writing style.
- *Language detection* is an important component in our pipeline as we use it to group samples by language and then reason about these languages. Noise in language detection might have influenced our results and insights.

Selection bias. There is a possibility that selection bias influences our results. First, sometimes fact-checkers writing the fact-checks base their writing on a particular post and the fact-check might contain parts of it verbatim. We tried to measure the size of this effect by comparing cases when the fact-checks are newer and older than posts (see §4.2), but we did not find a signal that this is the case. However, we know that there are at least few samples with this problem.

Second, there might be a bias towards social media posts that the social media platform or fact-checkers are already able to detect. Other, more difficult cases might still elude us.

Linguistic bias. Although our dataset is quite diverse, compared to most published datasets, there is still a bias towards major languages and Indo-European language family in particular. Crosslingual pairs consist mostly of East or South Asian

posts with non-Latin script mapped to English fact-checks. It is hard to estimate how our results would generalize to other language pairs. We visualize the languages in Figure 1. The annotation efforts in Section 6 shows that there are many crosslingual pairs that our data collection methodology was not able to collect.

9.2 Methods

Language support. The methods we use have different degrees of support for different languages. BM25 requires a proper tokenization to work. We have languages that use *scriptio continua* – Thai and Myanmar – where this is a problem. BM25-Original performance for these two is subpar, but could be improved by implementing custom tokenization models.

Multilingual TEMs we use do not support Sinhala and Tagalog languages, i.e., they were not trained with their data. The performance for these two languages is again subpar. Additionally, all methods depending on machine translation are naturally only able to handle languages that have a machine translation system available, although we believe that this was not a significant problem in our dataset.

Hidden positive pairs. The results we report might be deflated from the practical point of view because of unmarked correct pairs that are in the dataset. We have information only about a small subset of all the pairs. Our attempt to approximate true performance is provided in Section 6.

Supervised learning overfitting. It is possible that our supervised training yielded model that is overfitted on the particular languages and time frame that are represented in our dataset. The increase in performance might not transfer to out-of-domain pairs.

10 Ethical Considerations

We analyzed the likelihood and impact of ethical and societal risks for the most affected stakeholders, such as social media users and profile owners, fact-checkers, researchers, or social media platforms. For the most severe risks, we proposed respective countermeasures, following the guidelines and arguments in (Franzke et al., 2020; Townsend and Wallace, 2016; Mancosu and Vegetti, 2020).

Data collection process. While Twitter posts were collected using a publicly available API, the

Terms of Service (ToS) of Facebook and Instagram do not currently allow for the accessing or collecting of data using automated means. To minimize the harm to these social media platforms and their users, we made sure to only collect publicly available posts that are accessible even without logging in. This complies with the ToS.

Even if we admit the risk that such research activities could potentially violate the ToS, we argue that ignoring posts from Facebook and Instagram would prohibit research that seeks to address key current issues such as disinformation on these platforms (Bruns, 2019). These are some of the main platforms for disinformation dissemination in many countries. We consider the collection of such public data and its usage for research purposes to be in the public interest, especially considering the status of disinformation as a hybrid security threat (ENISA, 2022), which could justify minor harms to social media platforms.

Other considerable risks include the risk of accessibility privacy intrusion (Tavani, 2016) of social media users by observing them in an environment where they do not want to be observed. We did not obtain explicit consent from social media users to collect their posts. However, the criteria for considering social media data private or public depend on the assumption of whether social media users can reasonably expect to be observed by strangers (Townsend and Wallace, 2016). Twitter is considered an open platform. The collected posts on Facebook or Instagram are not only public, but the users can also expect that their posts will be widely shared, commented or reacted to and they can end up being fact-checked if it is the case.

Data publication. To minimize the risk of third-party misuse, the dataset is available only to researchers for research purposes. The full texts of the fact-checks are not published to avoid possible copyright violations.

Automatic translation has the risk of unintentional harm from misinterpretation of the original claims. To counter this risk, we always provide the original text as well.

We assessed the risk of re-identification, as well as the risk of revealing incorrect, highly sensitive or offensive content regarding social media users. At the same time, we had to take into account the fact that social media platforms remove some posts after they have been flagged as disinformation. Therefore, we decided to include the original texts of

the posts in the dataset to prevent it from decaying. Otherwise, it would become progressively less usable and research based on it less reproducible. This also allows us avoid publishing the URLs of posts, which would directly reveal the identities of the users. It is not possible to guarantee complete anonymity, since the posts are still linked in the fact-checks. The posts could also theoretically be found by full-text search.

On the other hand, all the posts released in our dataset are already mentioned in a publicly available space in the context of fact-checking efforts. Our publication of these posts does not significantly increase their already existing public exposure, especially considering the limited access options of our dataset.

To support users' rights to rectification and erasure in case of the publication of incorrect or sensitive information, we provide a procedure for them to request the removal of their posts from the dataset. However, we assess that the risk of wrongfully assigned fact-checks has a low probability (see §3).

As the dataset can also be used for supervised training (see §5), there is a risk of propagating biases present in the data (see §9). We recommend performing a proper linguistic analysis of any supervised model w.r.t. all the languages for which the model is intended. The results shown in this paper may not reflect the performance of the methods on other languages. We are also aware of the risk of propagating the biases of the fact-checkers, as it is they who decide what to fact-check. Although they should generally follow principles of fact-checking ethics (see, e.g., the IFCN's Code of Principles), there may still be present some human or systemic biases (Schwartz et al., 2022) that could affect the results when using the dataset for other purposes.

References

- Alan Agresti and Brent A Coull. 1998. Approximate is better than "exact" for interval estimation of binomial proportions. *The American Statistician*, 52(2):119–126.
- Alberto Barrón-Cedeño, Tamer Elsayed, Preslav Nakov, Giovanni Da San Martino, Maram Hasanain, Reem Suwaileh, Fatima Haouari, Nikolay Babulkov, Bayan Hamdan, Alex Nikolov, Shaden Shaar, and Zien Sheikh Ali. 2020. [Overview of CheckThat! 2020: Automatic Identification and Verification of Claims in Social Media](#). In *Experimental IR Meets Multilinguality, Multimodality, and Interaction*, vol-

| | | | |
|-----|--|---|-----|
| 783 | ume 12260 of <i>Lecture Notes in Computer Science</i> , | 59th Annual Meeting of the Association for Compu- | 837 |
| 784 | pages 215–236, Cham. Springer International Pub- | tational Linguistics and the 11th International Joint | 838 |
| 785 | lishing. | Conference on Natural Language Processing (Vol- | 839 |
| 786 | Varad Bhatnagar, Diptesh Kanojia, and Kameswari Che- | ume 1: Long Papers), pages 4504–4517, Online. As- | 840 |
| 787 | brolu. 2022. Harnessing abstractive summarization | sociation for Computational Linguistics. | 841 |
| 788 | for fact-checked claim detection . In <i>Proceedings of</i> | Moreno Mancosu and Federico Vegetti. 2020. What | 842 |
| 789 | <i>the 29th International Conference on Computational</i> | You Can Scrape and What Is Right to Scrape: A | 843 |
| 790 | <i>Linguistics</i> , pages 2934–2945, Gyeongju, Republic | Proposal for a Tool to Collect Public Facebook Data . | 844 |
| 791 | of Korea. International Committee on Computational | <i>Social Media + Society</i> , 6(3). SAGE Publications | 845 |
| 792 | Linguistics. | Ltd. | 846 |
| 793 | David M Blei, Andrew Y Ng, and Michael I Jordan. | Watheq Mansour, Tamer Elsayed, and Abdulaziz Al-Ali. | 847 |
| 794 | 2003. Latent Dirichlet allocation. <i>Journal of Ma-</i> | 2022. Did I See It Before? Detecting Previously- | 848 |
| 795 | <i>chine Learning Research</i> , 3(Jan):993–1022. | Checked Claims over Twitter . In <i>Advances in In-</i> | 849 |
| 796 | Axel Bruns. 2019. After the ‘APIcalypse’: social media | <i>formation Retrieval</i> , Lecture Notes in Computer Sci- | 850 |
| 797 | platforms and their fight against critical scholarly | ence, pages 367–381, Cham. Springer International | 851 |
| 798 | research . <i>Information, Communication & Society</i> , | Publishing. | 852 |
| 799 | 22(11):1544–1566. | Nicholas Micallef, Vivienne Armacost, Nasir Memon, | 853 |
| 800 | Laurence Dierickx, Ghazaal Sheikhi, Duc Tien | and Sameer Patil. 2022. True or false: Studying the | 854 |
| 801 | Dang Nguyen, and Carl-Gustav Lindén. 2022. Re- | work practices of professional fact-checkers . <i>Proc.</i> | 855 |
| 802 | port on the user needs of fact-checkers . Technical | <i>ACM Hum.-Comput. Interact.</i> , 6(CSCW1). | 856 |
| 803 | report, NORDIS – NORdic observatory for digital | Preslav Nakov, Alberto Barrón-Cedeño, Giovanni | 857 |
| 804 | media and information DISorders. | da San Martino, Firoj Alam, Julia Maria Struß, | 858 |
| 805 | ENISA. 2022. ENISA Threat Landscape 2022 . | Thomas Mandl, Rubén Míguez, Tommaso Caselli, | 859 |
| 806 | Aline Shakti Franzke, Anja Bechmann, Michael Zim- | Mucahid Kutlu, Wajdi Zaghoulani, Chengkai Li, | 860 |
| 807 | mer, Charles Ess, and the Association of Internet Re- | Shaden Shaar, Gautam Kishore Shahi, Hamdy | 861 |
| 808 | searchers. 2020. Internet research: Ethical guidelines | Mubarak, Alex Nikolov, Nikolay Babulkov, | 862 |
| 809 | 3.0. | Yavuz Selim Kartal, Michael Wiegand, Melanie | 863 |
| 810 | Full Fact. 2020. The challenges of online fact checking . | Siegel, and Juliane Köhler. 2022. Overview | 864 |
| 811 | Maarten Grootendorst. 2022. BERTopic: Neural topic | of the CLEF–2022 CheckThat! Lab on Fighting | 865 |
| 812 | modeling with a class-based TF-IDF procedure . | the COVID-19 infodemic and fake news detection . | 866 |
| 813 | arXiv:2203.05794 [cs.CL]. | In <i>Experimental IR Meets Multilinguality, Multi-</i> | 867 |
| 814 | Momchil Hardalov, Anton Chernyavskiy, Ivan Koychev, | <i>modality, and Interaction</i> , pages 495–520, Cham. | 868 |
| 815 | Dmitry Ilvovsky, and Preslav Nakov. 2022. Crowd- | Springer International Publishing. | 869 |
| 816 | Checked: Detecting previously fact-checked claims | Preslav Nakov, David Corney, Maram Hasanain, Firoj | 870 |
| 817 | in social media . In <i>Proceedings of the 2nd Confer-</i> | Alam, Tamer Elsayed, Alberto Barrón-Cedeño, Paolo | 871 |
| 818 | <i>ence of the Asia-Pacific Chapter of the Association</i> | Papotti, Shaden Shaar, and Giovanni Da San Martino. | 872 |
| 819 | <i>for Computational Linguistics and the 12th Interna-</i> | 2021. Automated Fact-Checking for Assisting Hu- | 873 |
| 820 | <i>tional Joint Conference on Natural Language Pro-</i> | man Fact-Checkers . In <i>Proceedings of the Thirtieth</i> | 874 |
| 821 | <i>cessing (Volume 1: Long Papers)</i> , pages 266–285, | <i>International Joint Conference on Artificial Intelli-</i> | 875 |
| 822 | Online only. Association for Computational Linguis- | <i>gence (IJCAI-21)</i> , pages 4551–4558. International | 876 |
| 823 | tics. | Joint Conferences on Artificial Intelligence Organi- | 877 |
| 824 | Andrea Hrckova, Robert Moro, Ivan Srba, Jakub | zation. | 878 |
| 825 | Simko, and Maria Bielikova. 2022. Automated, | Matúš Pikuliak, Marián Šimko, and Mária Bieliková. | 879 |
| 826 | not automatic: Needs and practices in European | 2021. Cross-lingual learning for text processing: A | 880 |
| 827 | fact-checking organizations as a basis for design- | survey . <i>Expert Systems with Applications</i> , 165. | 881 |
| 828 | ing human-centered AI systems . arXiv:2211.12143 | Nils Reimers and Iryna Gurevych. 2019. Sentence- | 882 |
| 829 | [cs.CY]. | BERT: Sentence embeddings using Siamese BERT- | 883 |
| 830 | Ye Jiang, Xingyi Song, Carolina Scarton, Ahmet Aker, | networks . In <i>Proceedings of the 2019 Conference on</i> | 884 |
| 831 | and Kalina Bontcheva. 2021. Categorising Fine-to- | <i>Empirical Methods in Natural Language Processing</i> | 885 |
| 832 | Coarse Grained Misinformation: An Empirical Study | <i>and the 9th International Joint Conference on Natu-</i> | 886 |
| 833 | of COVID-19 Infodemic . arXiv:2106.11702 [cs]. | <i>ral Language Processing (EMNLP-IJCNLP)</i> , pages | 887 |
| 834 | Ashkan Kazemi, Kiran Garimella, Devin Gaffney, and | 3982–3992, Hong Kong, China. Association for Com- | 888 |
| 835 | Scott Hale. 2021. Claim matching beyond English | putational Linguistics. | 889 |
| 836 | to scale global fact-checking . In <i>Proceedings of the</i> | Stephen Robertson and Hugo Zaragoza. 2009. The | 890 |
| | | probabilistic relevance framework: Bm25 and be- | 891 |
| | | yond . <i>Found. Trends Inf. Retr.</i> , 3(4):333–389. | 892 |

| | | | |
|-----|---|--|-----|
| 893 | Reva Schwartz, Apostol Vassilev, Kristen Greene, Lori | Nguyen Vo and Kyumin Lee. 2020. Where are the | 949 |
| 894 | Perine, Andrew Burt, and Patrick Hall. 2022. Towards a standard for identifying and managing bias | facts? Searching for fact-checked information to alle- | 950 |
| 895 | in artificial intelligence . NIST Special Publication | viate the spread of fake news . In <i>Proceedings of the</i> | 951 |
| 896 | 1270, National Institute of Standards and Technology. | <i>2020 Conference on Empirical Methods in Natural</i> | 952 |
| 897 | | <i>Language Processing (EMNLP)</i> , pages 7717–7731, | 953 |
| | | Online. Association for Computational Linguistics. | 954 |
| 898 | Shaden Shaar, Firoj Alam, Giovanni Da San Martino, | | |
| 899 | and Preslav Nakov. 2022a. The role of context in | A Computational Resources | 955 |
| 900 | detecting previously fact-checked claims . In <i>Find-</i> | We calculated all the results on an AWS-based vir- | 956 |
| 901 | <i>ings of the Association for Computational Linguis-</i> | tual machine located in the Ohio AWS data center. | 957 |
| 902 | <i>tics: NAACL 2022</i> , pages 1619–1631, Seattle, United | The machine has one NVIDIA Tesla T4 GPU in- | 958 |
| 903 | States. Association for Computational Linguistics. | stalled. The unsupervised experiments would take | 959 |
| 904 | Shaden Shaar, Nikolay Babulkov, Giovanni Da San Mar- | approximately 2 GPU days to replicate. The su- | 960 |
| 905 | tino, and Preslav Nakov. 2020. That is a known lie: | pervised experiments would take approximately 3 | 961 |
| 906 | Detecting previously fact-checked claims . In <i>Pro-</i> | GPU days to replicate. Additional roughly 4 GPU | 962 |
| 907 | <i>ceedings of the 58th Annual Meeting of the Asso-</i> | days were spent on other experiments that were | 963 |
| 908 | <i>ciation for Computational Linguistics</i> , pages 3607– | discarded or are not reported in this paper. | 964 |
| 909 | 3618, Online. Association for Computational Lin- | | |
| 910 | guistics. | | |
| 911 | Shaden Shaar, Nikola Georgiev, Firoj Alam, Giovanni | B Dataset Pipeline Details | 965 |
| 912 | Da San Martino, Aisha Mohamed, and Preslav Nakov. | | |
| 913 | 2022b. Assisting the human fact-checkers: Detect- | B.1 Dataset Collection | 966 |
| 914 | ing all previously fact-checked claims in a document . | | |
| 915 | In <i>Findings of the Association for Computational</i> | Crawling. We use a <i>Selenium</i> -based web crawler | 967 |
| 916 | <i>Linguistics: EMNLP 2022</i> , pages 2069–2080, Abu | that visits the links, extracts the HTML content and | 968 |
| 917 | Dhabi, United Arab Emirates. Association for Com- | parses it with the <i>Beautiful Soup</i> ⁶ library. | 969 |
| 918 | putational Linguistics. | | |
| 919 | Shaden Shaar, Fatima Haouari, Watheq Mansour, | Source of fact-checks. We only processed fact- | 970 |
| 920 | Maram Hasanain, Nikolay Babulkov, Firoj Alam, | checks written by the AFP news agency. We chose | 971 |
| 921 | Giovanni Da San Martino, Tamer Elsayed, and | them because they are an established fact-checking | 972 |
| 922 | Preslav Nakov. 2021. Overview of the CLEF-2021 | organization with high editing standards and are | 973 |
| 923 | CheckThat! Lab Task 2 on Detecting Previously Fact- | also a part of Meta’s <i>Third-Party Fact-Checking</i> | 974 |
| 924 | Checked Claims in Tweets and Political Debates . In | <i>Program</i> . Pairs with fact-checks from other or- | 975 |
| 925 | <i>Proceedings of the Working Notes of CLEF 2021 -</i> | ganizations might have been established from the | 976 |
| 926 | <i>Conference and Labs of the Evaluation Forum</i> , vol- | warnings. | 977 |
| 927 | ume 2936. CEUR-WS. | | |
| 928 | Qiang Sheng, Juan Cao, Xueyao Zhang, Xirong Li, | Archiving services. Since the content from so- | 978 |
| 929 | and Lei Zhong. 2021. Article reranking by memory- | cial media networks may disappear in time, fact- | 979 |
| 930 | enhanced key sentence matching for detecting | checkers tend to use various content archiving ser- | 980 |
| 931 | previously fact-checked claims . In <i>Proceedings of the</i> | vices (e.g., <code>perma.cc</code>). We extract the content | 981 |
| 932 | <i>59th Annual Meeting of the Association for Compu-</i> | from these services as well. | 982 |
| 933 | <i>tational Linguistics and the 11th International Joint</i> | | |
| 934 | <i>Conference on Natural Language Processing (Vol-</i> | AI APIs. We use following services to process | 983 |
| 935 | <i>ume 1: Long Papers)</i> . Association for Computational | our samples: | 984 |
| 936 | Linguistics. | | |
| 937 | H.T. Tavani. 2016. <i>Ethics and Technology: Controvers-</i> | • <i>Google Vision API</i> . We use Google Vision | 985 |
| 938 | <i>ies, Questions, and Strategies for Ethical Comput-</i> | API to extract text from images attached to the | 986 |
| 939 | <i>ing</i> , 5th edition. Wiley. | post. The API also returns a list of languages | 987 |
| 940 | Leanne Townsend and Claire Wallace. 2016. Social | found in each image with their percentage. | 988 |
| 941 | media research: A guide to ethics . | | |
| 942 | Nguyen Vo and Kyumin Lee. 2018. The Rise of | • <i>Google Translate API</i> . We use Google Trans- | 989 |
| 943 | Guardians: Fact-checking URL Recommendation | late API to translate all the texts into English. | 990 |
| 944 | to Combat Fake News . In <i>The 41st International</i> | The API also returns a most probable lan- | 991 |
| 945 | <i>ACM SIGIR Conference on Research & Development</i> | guage. | 992 |
| 946 | <i>in Information Retrieval, SIGIR ’18</i> , pages 275–284, | | |
| 947 | New York, NY, USA. Association for Computing | | |
| 948 | Machinery. | | |

⁶<https://www.crummy.com/software/BeautifulSoup/>

B.2 Dataset Pre-Processing

We performed several cleaning and pre-processing steps with our dataset. All the pre-processing is available in the released code repository.

Removing noisy claims. We removed fact-checks that had no claim or where the claim was shorter than 10 characters.

Fact-check deduplication. We unified fact-checks with identical claims.

Noise in social media posts. We removed texts or OCR transcripts that we deemed noisy (shorter than 25 characters or more than 50% non-alphabetical characters). We then only kept posts where at least one text was considered not noisy. We also removed noisy lines from OCR transcripts (Lines shorter than 5 characters or with more than 50% non-alphabetical characters). We also removed URLs.

Post deduplication. We unified posts that ended up with identical text contents after the cleaning process.

Machine translation. We translated all the texts into English. The only exceptions were fact-check claims coming from English-language providers (e.g., Snopes) that we considered English by default, and fact-check claims where CLD3⁷ identified English language. We confirmed experimentally that CLD3 has a high precision on English texts.

Language identification normalization. We observed that there are some systematic errors in the language identification models we used. We found out that the model often selected less common languages based on spurious patterns, e.g., mentions of Filipino politicians sometimes led to Ilocano language prediction. Based on data analysis, we changed some predictions automatically, e.g., all Ilocano predictions were changed into English. Sometimes we only did it when the script did not match the language, e.g., for posts with Latin script identified as Oromo. We do not recommend using this process automatically on any data. In other contexts, the generated predictions might be less noisy. Even in our case, we have different rules for posts and fact-checks based on the characteristics of these two domains. If the predictions proved to be too noisy, we unified several

⁷<https://github.com/google/cld3>

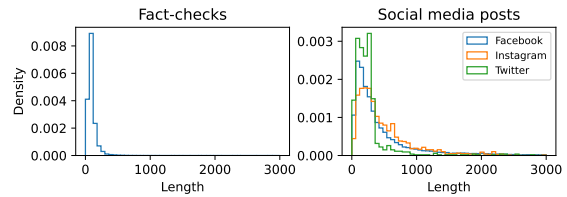


Figure 8: Density plots for the character lengths of the fact-checked claims and the social media posts in our dataset.

languages or language varieties into one. This is the case of Croatian, Bosnian and Serbian, as well as Indonesian and Malay.

C Dataset Statistics

We show the number of fact-checks and posts per language in Table 5. For fact-checks, we only take into consideration the language of claim, since we mostly only work with claims in this work.

Posts can have more than one language detected based on its overall compositions. We calculated percentage for each language based on the language prediction methods. We consider all languages with at least 20% to be relevant. 25,482 posts have only one language detected, while 2,549 has two, 59 has three, 1 has four and 1 has zero.

Table 6 shows the sources of our fact-checks. Here we only show the statistics for the fact-checks we actually used in our experiments. There are additional 6k fact-checks that we have not used because they we were not able to fill their *claim* field.

Table 7 show the number of fact-check-to-post pairs for different language combinations.

Figure 8 show the density of lengths for both the fact-checked claims and the SMPs. Both have long tail distributions, but the claims are in general much shorter. 99% of claims are shorter than 379 characters. For social media posts, it is 4129 characters.

D Hyperparameters

D.1 BM25

We use the default PyTerrier values for BM25 algorithm: $k_1 = 1.2$, $b = 0.75$. Our preliminary results show that the results are not very sensitive towards these two hyperparameters, probably because of the relatively short length of the documents that we retrieve. Most claims in our dataset have only one sentence.

| Code | Language | # fact-checks | # posts |
|------|----------------|---------------|---------|
| ara | Arabic | 14201 | 931 |
| asm | Assamese | 60 | 5 |
| aze | Azerbaijani | 178 | 2 |
| bul | Bulgarian | 162 | 114 |
| ben | Bengali | 4143 | 113 |
| cat | Catalan | 574 | 100 |
| ces | Czech | 254 | 265 |
| dan | Danish | 648 | 6 |
| deu | German | 4996 | 932 |
| ell | Greek | 1821 | 175 |
| eng | English | 85814 | 7307 |
| spa | Spanish | 14082 | 7319 |
| fas | Farsi | 418 | 17 |
| fin | Finnish | 109 | 103 |
| tgl | Tagalog | 462 | 439 |
| fra | French | 4355 | 2146 |
| hbs | Serbo-Croatian | 2451 | 481 |
| hin | Hindi | 7149 | 833 |
| hun | Hungarian | 139 | 113 |
| ita | Italian | 3047 | 65 |
| heb | Hebrew | 202 | 2 |
| jpn | Japanese | 62 | 7 |
| khm | Khmer | 144 | 6 |
| kor | Korean | 510 | 474 |
| mkd | Macedonian | 1125 | 1 |
| mal | Malayalam | 1206 | 4 |
| msa | Malay | 8424 | 1389 |
| mya | Myanmar | 92 | 172 |
| nld | Dutch | 1232 | 257 |
| nor | Norwegian | 440 | 5 |
| pol | Polish | 6912 | 453 |
| por | Portuguese | 21569 | 3366 |
| ron | Romanian | 204 | 238 |
| rus | Russian | 2715 | 28 |
| sin | Sinhala | 825 | 534 |
| slk | Slovak | 260 | 363 |
| sqi | Albanian | 726 | 1 |
| tam | Tamil | 1612 | 29 |
| tel | Telugu | 2450 | 11 |
| tha | Thai | 382 | 626 |
| tur | Turkish | 6676 | 7 |
| ukr | Ukrainian | 68 | 6 |
| urd | Urdu | 0 | 378 |
| zho | Chinese | 2586 | 595 |
| | Others | 266 | 343 |

Table 5: List of languages with at least 50 fact-checks or 50 posts.

D.2 Supervised Training

Table 8 show the range of hyperparameters in our hyperparameter search done for supervised training in Section 5, as well as the best performing hyperparameters.

E Additional Results

E.1 Additional Metrics

Table 9 show additional IR metrics calculated for the experiments done in Section 4.1. There is a strong correlation between all these metrics, as shown in Table 10. This is caused by the fact that most of the SMPs have only one fact-check assigned and the calculations for such cases are very similar for different metrics. We ultimately decided to use S@10 as our main evaluation metric in this work as we find it to be the most interpretable measure (*for how many pairs the expected fact-checked claim ended up in the top 10*).

E.2 Detailed Per-language Results

Table 11 shows additional language-specific results for all the methods from Section 4 including confidence intervals.

F Other Ideas

Here we discuss some additional ideas that were tried and that we decided not to include in the main text for various reasons.

Sliding window embedding. Figure 8 shows that the performance for methods decreases for posts with certain length. The decrease is generally starting at around 500 characters. We experimented with using sliding windows with various sizes (both based on the number of characters and the number of sentences) and strides. TEMs then encode only this sliding window and the final vector similarity is calculated as the maximum similarity of any of the windows. We found out that this technique can slightly (+0.01 – 0.02 S@10) improve the results for TEMs.

Using fact-check titles alongside claims. We represent fact-checks with the *claim* field obtained from the data in our main text experiments. We also experimented with the *title* field that we were able to obtain for the majority of the fact-checks. We found out that representing the fact-check as a concatenation of a claim and a title improves the results slightly (+0.00 – 0.01 S@10) for BM25 methods.

| Name | Lang. | N | Name | Lang. | N | Name | Lang. | N |
|----------------------|-------|-------|---------------------------|-------|------|----------------------------|-------|-----|
| snopes.com | eng | 18376 | washingtonpost.com | eng | 1413 | agi.it | ita | 246 |
| politifact.com | eng | 9029 | dogrulukpayi.com | tur | 1360 | verify-sy.com | ara | 242 |
| misbar.com | ara | 9027 | stopfake.org | rus | 1307 | cbsnews.com | eng | 242 |
| boomlive.in | eng | 7949 | colombiacheck.com | spa | 1271 | factchecknederland.afp.com | nld | 234 |
| factcheck.afp.com | eng | 6853 | tempo.co | id | 1143 | butac.it | ita | 220 |
| cekfakta.com | ind | 6523 | vistinomer.mk | mkd | 1141 | efe.com | spa | 219 |
| altnews.in | eng | 6199 | faktograf.hr | hbs | 1094 | br.de | deu | 214 |
| factly.in | eng | 5818 | dubawa.org | eng | 1066 | annielab.org | eng | 204 |
| leadstories.com | eng | 5319 | factcheck.kz | rus | 1044 | globes.co.il | heb | 202 |
| sapo.pt | por | 5200 | istinomer.rs | hbs | 958 | factcheckhub.com | eng | 200 |
| demagog.org.pl | pol | 4292 | boombd.com | ben | 937 | ghanafact.com | eng | 199 |
| fullfact.org | eng | 4260 | bufale.net | ita | 928 | telemundo.com | spa | 195 |
| factual.afp.com | spa | 4051 | apublica.org | por | 915 | apa.at | deu | 185 |
| uol.com.br | por | 3908 | rappler.com | tgl | 874 | verificat.afp.com | ron | 177 |
| checkyourfact.com | eng | 3620 | verificat.cat | spa | 821 | efetococuyo.com | spa | 170 |
| teyit.org | tur | 3289 | kallxo.com | sqi | 728 | factcheckni.org | eng | 157 |
| newsmobile.in | eng | 3265 | aap.com.au | eng | 687 | proveri.afp.com | bul | 152 |
| newtral.es | spa | 3256 | projetocomprova.com.br | por | 686 | icirnigeria.org | eng | 142 |
| dpa-factchecking.com | nld | 2839 | tjekdet.dk | dan | 648 | tenykerdes.afp.com | hun | 138 |
| indiatoday.in | eng | 2799 | dogrula.org | tur | 634 | liberation.fr | fra | 134 |
| factcheck.org | eng | 2716 | faktencheck.afp.com | deu | 629 | factcheckgreek.afp.com | ell | 129 |
| aosfatos.org | por | 2596 | thip.media | ben | 598 | radio-canada.ca | fra | 123 |
| boatos.org | por | 2553 | dailyo.in | ben | 591 | maharat-news.com | ara | 121 |
| aahtak.in | hin | 2493 | univision.com | spa | 582 | factcheckmyanmar.afp.com | mya | 119 |
| dabegad.com | ara | 2342 | periksafakta.afp.com | ind | 563 | jachai.org | ben | 113 |
| factcheck.afp.com/ar | ara | 2292 | lemonde.fr | fra | 558 | nieuwscheckers.nl | nld | 111 |
| estadao.com.br | por | 2197 | check4spam.com | eng | 524 | europapress.es | spa | 108 |
| factuel.afp.com | fra | 2178 | healthfeedback.org | eng | 499 | faktantarkistus.afp.com | fin | 107 |
| thequint.com | eng | 2058 | mygopen.com | zho | 494 | tagesschau.de | deu | 103 |
| tfc-taiwan.org.tw | zho | 1960 | sprawdzam.afp.com | pol | 458 | scroll.in | eng | 100 |
| observador.pt | por | 1930 | faktisk.no | nor | 444 | thelallantop.com | hin | 99 |
| usatoday.com | eng | 1901 | presseportal.de | deu | 439 | theferret.scot | eng | 96 |
| oko.press | pol | 1872 | 20minutes.fr | fra | 419 | france24.com | fra | 92 |
| fatabyyano.net | ara | 1844 | cinjenice.afp.com | hbs | 387 | voachinese.com | zho | 92 |
| factrescendo.com | ben | 1808 | factcheckthailand.afp.com | tha | 382 | comprovem.afp.com | cat | 90 |
| correctiv.org | deu | 1783 | factcheckkorea.afp.com | kor | 382 | factandfurious.com | fra | 82 |
| maldita.es | spa | 1748 | asianetnews.com | mal | 365 | factchecker.in | eng | 74 |
| ellinikahoaxes.gr | ell | 1688 | newsweek.com | eng | 364 | telugupost.com | tel | 73 |
| checamos.afp.com | por | 1672 | factnameh.com | fas | 356 | zimfact.org | eng | 72 |
| facta.news | ita | 1652 | fakenews.pl | pol | 320 | factcheckbangla.afp.com | ben | 62 |
| youturn.in | tam | 1609 | fastcheck.cl | spa | 313 | buzzfeed.com | eng | 56 |
| malumatfurus.org | tur | 1572 | newsmeter.in | eng | 290 | verificado.com.mx | spa | 55 |
| polygraph.info | eng | 1527 | factrakers.org | eng | 276 | ripplesnigeria.com | eng | 52 |
| metafact.io | eng | 1526 | semakanfakta.afp.com | msa | 267 | poynter.org | eng | 52 |
| africacheck.org | eng | 1468 | fakty.afp.com | slk | 260 | globo.com | por | 52 |
| animalpolitico.com | spa | 1468 | napravoumiru.afp.com | ces | 255 | radiofarda.com | fas | 51 |
| verafiles.org | tgl | 1414 | factograph.info | rus | 253 | stern.de | deu | 50 |

Table 6: Fact-checking sources with at least 50 fact-checks in our dataset.

| | Fact-check language | | | | | | | | | | | | | | | | | | | | | | | | | |
|-------|---------------------|-----|-----|-----|-----|-----|-----|------|------|-----|------|-----|-----|-----|-----|------|-----|-----|-----|------|-----|-----|-----|-----|-------|---|
| | ara | ben | bul | cat | ces | deu | ell | eng | fin | fra | hbs | hin | hun | kor | msa | mya | nld | pol | por | ron | sin | slk | spa | tha | Other | |
| ara | 864 | 0 | 1 | 0 | 0 | 0 | 0 | 28 | 1 | 0 | 23 | 0 | 0 | 0 | 0 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| ben | 0 | 109 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | |
| bul | 0 | 0 | 56 | 0 | 0 | 0 | 0 | 75 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| cat | 0 | 0 | 0 | 79 | 0 | 0 | 0 | 0 | 10 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| ces | 0 | 0 | 0 | 0 | 211 | 2 | 1 | 5 | 1 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 28 | 0 | 0 | |
| deu | 1 | 0 | 0 | 0 | 1 | 961 | 1 | 5 | 1 | 0 | 0 | 4 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 1 | 0 | 0 | 3 | 1 | 1 | |
| ell | 0 | 0 | 0 | 0 | 0 | 0 | 222 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| eng | 2 | 1 | 2 | 4 | 12 | 37 | 4 | 6450 | 77 | 5 | 50 | 8 | 16 | 3 | 25 | 67 | 0 | 12 | 20 | 34 | 3 | 5 | 14 | 12 | 28 | |
| fin | 3 | 0 | 0 | 43 | 0 | 1 | 2 | 39 | 7528 | 0 | 5 | 0 | 0 | 0 | 0 | 4 | 0 | 0 | 6 | 28 | 1 | 0 | 1 | 0 | 13 | |
| fra | 0 | 0 | 0 | 0 | 0 | 3 | 0 | 0 | 0 | 95 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | |
| hbs | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 264 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| hin | 1 | 0 | 3 | 0 | 0 | 1 | 2 | 20 | 3 | 0 | 2105 | 0 | 0 | 0 | 0 | 2 | 0 | 0 | 1 | 4 | 2 | 1 | 1 | 0 | 7 | |
| hun | 0 | 2 | 0 | 0 | 0 | 2 | 0 | 4 | 0 | 0 | 0 | 490 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 6 | |
| kor | 0 | 0 | 3 | 0 | 0 | 0 | 0 | 819 | 0 | 0 | 0 | 0 | 106 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 2 | |
| msa | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 107 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 0 | 0 | |
| mya | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 128 | 0 | 0 | 0 | 0 | 0 | 0 | 460 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| nld | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 249 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1403 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 1 | |
| pol | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 122 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 102 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | |
| por | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 270 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| ron | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 7 | 0 | 0 | 1 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 459 | 0 | 0 | 0 | 1 | 0 | 1 | |
| sin | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 16 | 32 | 0 | 3 | 1 | 2 | 0 | 0 | 2 | 0 | 0 | 1 | 3739 | 0 | 0 | 1 | 0 | 4 | |
| slk | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 226 | 0 | 0 | 0 | 0 | 1 | |
| spa | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 456 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 57 | 0 | 0 | 3 | |
| tgl | 0 | 0 | 0 | 0 | 3 | 0 | 0 | 3 | 2 | 0 | 1 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 332 | 0 | 0 |
| tha | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 202 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 574 | 0 | |
| urd | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 328 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| zho | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 509 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 17 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 15 | |
| Other | 5 | 1 | 0 | 0 | 2 | 5 | 3 | 216 | 20 | 1 | 8 | 3 | 1 | 0 | 1 | 7 | 0 | 3 | 0 | 5 | 1 | 0 | 2 | 2 | 17 | |

Table 7: Number of fact-check-to-post pairs for different language combinations. Note that one SMP can have more than one language assigned.

| Hyperparameter | Range | GTR-T5-Large | MPNet-Base-Multilingual |
|---------------------------------------|--|--------------------|-------------------------|
| Loss | contrastive cosine online-contrastive | online-contrastive | online-contrastive |
| Learning rate | $[1e-3, 1e-7]$ | $1e-5$ | $5e-6$ |
| Learning rate schedule | cosine linear | cosine | cosine |
| Warmup steps | $[100, 3200]$ | 800 | 1600 |
| Weight decay rate | $[1e-7, 1e-4]$ | $1e-5$ | $8e-5$ |
| Ratio of positive to negative samples | $[10, 50\%]$ | 10% | 30% |
| Margin | $[0.1, 0.5]$ | 0.5 | 0.4 |
| Batch size | Maximum possible | 2 | 8 |

Table 8: Range of hyperparameters used in our supervised hyperparameter search and the hyperparameters of our most successful models. The ranges adjusted during the experimentation according to the preliminary results.

| Method | Ver. | S@10 | | MRR | | MAP | | NDCG | | MAP@10 | |
|-----------------------------|------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | | Mono | Cross | Mono | Cross | Mono | Cross | Mono | Cross | Mono | Cross |
| BM25 | | | | | | | | | | | |
| BM25 | En | 0.61 | 0.22 | 0.78 | 0.39 | 0.61 | 0.22 | 0.68 | 0.33 | 0.82 | 0.40 |
| BM25 | Og | 0.48 | 0.04 | 0.62 | 0.06 | 0.47 | 0.04 | 0.56 | 0.11 | 0.65 | 0.07 |
| English TEMs | | | | | | | | | | | |
| DistilRoBERTa | En | 0.59 | 0.24 | 0.76 | 0.43 | 0.59 | 0.24 | 0.67 | 0.36 | 0.80 | 0.44 |
| GTR-T5-Base | En | 0.65 | 0.30 | 0.81 | 0.51 | 0.64 | 0.30 | 0.71 | 0.41 | 0.85 | 0.53 |
| GTR-T5-Large | En | 0.67 | 0.33 | 0.83 | 0.56 | 0.67 | 0.33 | 0.73 | 0.45 | 0.88 | 0.58 |
| GTR-T5-XL | En | 0.67 | 0.34 | 0.83 | 0.56 | 0.67 | 0.33 | 0.74 | 0.45 | 0.87 | 0.58 |
| MPNet-Base | En | 0.61 | 0.27 | 0.78 | 0.47 | 0.60 | 0.27 | 0.68 | 0.38 | 0.82 | 0.48 |
| MSMARCO-BERT-Base | En | 0.62 | 0.26 | 0.78 | 0.46 | 0.61 | 0.26 | 0.69 | 0.37 | 0.82 | 0.48 |
| MiniLM-L12 | En | 0.64 | 0.29 | 0.80 | 0.48 | 0.63 | 0.29 | 0.71 | 0.40 | 0.84 | 0.50 |
| MultiQA-MPNet-Base | En | 0.64 | 0.29 | 0.80 | 0.50 | 0.63 | 0.29 | 0.71 | 0.41 | 0.84 | 0.52 |
| SGPT-125M | En | 0.47 | 0.14 | 0.63 | 0.25 | 0.47 | 0.14 | 0.56 | 0.25 | 0.66 | 0.26 |
| SGPT-2.7B | En | 0.60 | 0.29 | 0.77 | 0.50 | 0.60 | 0.29 | 0.68 | 0.40 | 0.81 | 0.52 |
| Sentence-T5-Base | En | 0.57 | 0.21 | 0.73 | 0.37 | 0.56 | 0.21 | 0.65 | 0.33 | 0.77 | 0.38 |
| Sentence-T5-Large | En | 0.58 | 0.24 | 0.75 | 0.41 | 0.58 | 0.23 | 0.66 | 0.35 | 0.78 | 0.43 |
| Sentence-T5-XL | En | 0.61 | 0.27 | 0.78 | 0.47 | 0.60 | 0.26 | 0.68 | 0.38 | 0.82 | 0.48 |
| Multilingual TEMs | | | | | | | | | | | |
| DistilUSE-Base-Multilingual | En | 0.58 | 0.22 | 0.74 | 0.40 | 0.58 | 0.22 | 0.66 | 0.34 | 0.78 | 0.41 |
| | Og | 0.50 | 0.10 | 0.66 | 0.20 | 0.49 | 0.10 | 0.59 | 0.21 | 0.69 | 0.21 |
| LaBSE | En | 0.47 | 0.13 | 0.63 | 0.22 | 0.46 | 0.13 | 0.56 | 0.23 | 0.66 | 0.23 |
| | Og | 0.53 | 0.12 | 0.69 | 0.22 | 0.53 | 0.12 | 0.61 | 0.22 | 0.72 | 0.23 |
| MPNet-Base-Multilingual | En | 0.60 | 0.24 | 0.75 | 0.41 | 0.59 | 0.23 | 0.67 | 0.35 | 0.79 | 0.42 |
| | Og | 0.53 | 0.11 | 0.70 | 0.21 | 0.53 | 0.11 | 0.61 | 0.22 | 0.73 | 0.22 |
| MiniLM-L2-Multilingual | En | 0.58 | 0.22 | 0.74 | 0.38 | 0.57 | 0.22 | 0.66 | 0.34 | 0.77 | 0.40 |
| | Og | 0.48 | 0.08 | 0.63 | 0.15 | 0.47 | 0.08 | 0.57 | 0.18 | 0.67 | 0.16 |
| XLM-R | En | 0.55 | 0.19 | 0.72 | 0.33 | 0.55 | 0.19 | 0.63 | 0.30 | 0.76 | 0.34 |
| | Og | 0.50 | 0.08 | 0.66 | 0.15 | 0.50 | 0.08 | 0.59 | 0.18 | 0.70 | 0.16 |

Table 9: The results for different ranking methods. This table shows the same experiment as Table 2, but also calculates additional information retrieval metrics: MRR, MAP, NDCG, MAP@10.

| | S@10 | MRR | MAP | NDCG | MAP@10 |
|--------|-------|-------|-------|-------|--------|
| S@10 | 1.000 | 0.986 | 0.986 | 0.993 | 1.000 |
| MRR | 0.986 | 1.000 | 1.000 | 0.998 | 0.988 |
| MAP | 0.986 | 1.000 | 1.000 | 0.998 | 0.988 |
| NDCG | 0.993 | 0.998 | 0.998 | 1.000 | 0.995 |
| MAP@10 | 1.000 | 0.988 | 0.988 | 0.995 | 1.000 |

Table 10: Pearson correlation coefficient between different metrics as calculated in Table 9. Both monolingual and crosslingual scores are taken into consideration.

Topic detection. We attempted to run a topic detection over our posts to better understand how different methods handle different topics and themes in our data. We experimented with both *original* and *English* versions, with both multilingual and monolingual topic detection models, such as LDA (Blei et al., 2003) or BERTopic (Grootendorst, 2022). Ultimately we were not content with the quality of topic detection, as the models failed to reliably identify even the most frequent topics in our data, such as the COVID-19 pandemic or Russo-Ukrainian war. We believe that this is caused by the short length of the majority of the posts, as well as their relatively noisy nature.

Mixing original and English versions. We experimented with representing both fact-checks and posts as a concatenation of both the original language texts and the English translations, so that the multilingual methods can use both sources of information. However, this increased the *same language bias* significantly while the performance decreased significantly across the board.

G Examples

This Appendix contains 5 randomly selected fact-check-post pairs from our dataset. We show here all the information present in our dataset for these samples.

G.1 Example #1

G.1.1 Fact-check

ID: 104315

Published at: 2021-09-27 factcheck.afp.com

Claim

Original text: *Photo shows meeting of five international intelligence agencies in Delhi*

Translated text: *Photo shows meeting of five international intelligence agencies in Delhi*

Detected languages: eng: 100.0%

Title

Original text: *This photo shows a delegation-level meeting between Indian and Russian national security advisors*

Translated text: *This photo shows a delegation-level meeting between Indian and Russian national security advisors*

Detected languages: eng: 100.0%

G.1.2 Social Media Post

ID: 16806

Published at: Facebook 2021-09-16

Verdicts: Partly false information

Main text

Original text: *Post Giri IyerNow in Delhi !!India RAWIsrael MOSSADAmerica CIARussia KGBEngland MI6First time ever that the top five intelligence agency of the world are sitting together for a high level meeting in Delhi. This is the power of new India*

Translated text: *Post Giri IyerNow in Delhi !!India RAWIsrael MOSSADAmerica CIARussia KGBEngland MI6First time ever that the top five intelligence agency of the world are sitting together for a high level meeting in Delhi. This is the power of new India*

Detected languages: eng: 100.0%

G.2 Example #2

G.2.1 Fact-check

ID: 34296

Published at: 2019-10-29 factuel.afp.com

Claim

Original text: *COMMENT TRAITER L'HÉPATITE B PAR LES PLANTES*

Translated text: *HOW TO TREAT HEPATITIS B WITH HERBS*

Detected languages: fra: 100.0%

Title

Original text: *Non, cette boisson à base de papaye, de citron, de racines de cocotier et de moringa bouillis, ne guérit pas l'hépatite B*

Translated text: *No, this drink made from boiled papaya, lemon, coconut palm roots and moringa does not cure hepatitis B*

Detected languages: fra: 100.0%

| | Ver. | spa | eng | por | fra | msa | deu | ara | tha | hbs | kor |
|-----------------------------|------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| BM25 | En | 0.79 ± 0.01 | 0.62 ± 0.01 | 0.78 ± 0.01 | 0.77 ± 0.02 | 0.74 ± 0.02 | 0.70 ± 0.03 | 0.82 ± 0.03 | 0.89 ± 0.03 | 0.66 ± 0.04 | 0.81 ± 0.04 |
| BM25 | Og | 0.76 ± 0.01 | 0.60 ± 0.01 | 0.79 ± 0.01 | 0.71 ± 0.02 | 0.73 ± 0.02 | 0.56 ± 0.03 | 0.75 ± 0.03 | 0.26 ± 0.04 | 0.47 ± 0.04 | 0.58 ± 0.04 |
| DistilRoBERTa | En | 0.72 ± 0.01 | 0.64 ± 0.01 | 0.64 ± 0.02 | 0.79 ± 0.02 | 0.75 ± 0.02 | 0.58 ± 0.03 | 0.79 ± 0.03 | 0.89 ± 0.03 | 0.65 ± 0.04 | 0.82 ± 0.04 |
| GTR-T5-Base | En | 0.80 ± 0.01 | 0.72 ± 0.01 | 0.75 ± 0.01 | 0.84 ± 0.02 | 0.80 ± 0.02 | 0.66 ± 0.03 | 0.85 ± 0.02 | 0.90 ± 0.02 | 0.71 ± 0.04 | 0.85 ± 0.03 |
| GTR-T5-Large | En | 0.84 ± 0.01 | 0.77 ± 0.01 | 0.80 ± 0.01 | 0.86 ± 0.01 | 0.82 ± 0.02 | 0.69 ± 0.03 | 0.86 ± 0.02 | 0.90 ± 0.02 | 0.74 ± 0.04 | 0.86 ± 0.03 |
| GTR-T5-XL | En | 0.83 ± 0.01 | 0.77 ± 0.01 | 0.80 ± 0.01 | 0.86 ± 0.01 | 0.83 ± 0.02 | 0.70 ± 0.03 | 0.87 ± 0.02 | 0.91 ± 0.02 | 0.73 ± 0.04 | 0.86 ± 0.03 |
| MPNet-Base | En | 0.75 ± 0.01 | 0.68 ± 0.01 | 0.67 ± 0.02 | 0.81 ± 0.02 | 0.79 ± 0.02 | 0.57 ± 0.03 | 0.80 ± 0.03 | 0.88 ± 0.03 | 0.70 ± 0.04 | 0.80 ± 0.04 |
| MSMARCO-BERT-Base | En | 0.78 ± 0.01 | 0.63 ± 0.01 | 0.75 ± 0.01 | 0.80 ± 0.02 | 0.79 ± 0.02 | 0.59 ± 0.03 | 0.82 ± 0.03 | 0.87 ± 0.03 | 0.66 ± 0.04 | 0.84 ± 0.03 |
| MiniLM-L12 | En | 0.78 ± 0.01 | 0.70 ± 0.01 | 0.71 ± 0.01 | 0.82 ± 0.02 | 0.78 ± 0.02 | 0.64 ± 0.03 | 0.84 ± 0.02 | 0.89 ± 0.03 | 0.72 ± 0.04 | 0.83 ± 0.03 |
| MultiQA-MPNet-Base | En | 0.79 ± 0.01 | 0.70 ± 0.01 | 0.71 ± 0.01 | 0.84 ± 0.02 | 0.81 ± 0.02 | 0.64 ± 0.03 | 0.82 ± 0.03 | 0.90 ± 0.02 | 0.71 ± 0.04 | 0.84 ± 0.03 |
| SGPT-125M | En | 0.54 ± 0.01 | 0.40 ± 0.01 | 0.52 ± 0.02 | 0.60 ± 0.02 | 0.55 ± 0.03 | 0.45 ± 0.03 | 0.70 ± 0.03 | 0.82 ± 0.03 | 0.52 ± 0.04 | 0.71 ± 0.04 |
| SGPT-2.7B | En | 0.77 ± 0.01 | 0.65 ± 0.01 | 0.72 ± 0.01 | 0.79 ± 0.02 | 0.79 ± 0.02 | 0.60 ± 0.03 | 0.82 ± 0.03 | 0.91 ± 0.02 | 0.70 ± 0.04 | 0.83 ± 0.03 |
| Sentence-T5-Base | En | 0.71 ± 0.01 | 0.60 ± 0.01 | 0.63 ± 0.02 | 0.73 ± 0.02 | 0.77 ± 0.02 | 0.43 ± 0.03 | 0.76 ± 0.03 | 0.89 ± 0.03 | 0.59 ± 0.04 | 0.83 ± 0.03 |
| Sentence-T5-Large | En | 0.74 ± 0.01 | 0.64 ± 0.01 | 0.65 ± 0.02 | 0.76 ± 0.02 | 0.76 ± 0.02 | 0.47 ± 0.03 | 0.80 ± 0.03 | 0.89 ± 0.02 | 0.59 ± 0.04 | 0.81 ± 0.04 |
| Sentence-T5-XL | En | 0.77 ± 0.01 | 0.68 ± 0.01 | 0.70 ± 0.01 | 0.80 ± 0.02 | 0.79 ± 0.02 | 0.53 ± 0.03 | 0.82 ± 0.03 | 0.91 ± 0.02 | 0.68 ± 0.04 | 0.81 ± 0.04 |
| DistilUSE-Base-Multilingual | En | 0.69 ± 0.01 | 0.57 ± 0.01 | 0.63 ± 0.02 | 0.75 ± 0.02 | 0.74 ± 0.02 | 0.53 ± 0.03 | 0.80 ± 0.03 | 0.88 ± 0.03 | 0.65 ± 0.04 | 0.79 ± 0.04 |
| DistilUSE-Base-Multilingual | Og | 0.64 ± 0.01 | 0.56 ± 0.01 | 0.58 ± 0.02 | 0.69 ± 0.02 | 0.60 ± 0.03 | 0.50 ± 0.03 | 0.74 ± 0.03 | 0.72 ± 0.04 | 0.57 ± 0.04 | 0.74 ± 0.04 |
| LaBSE | En | 0.56 ± 0.01 | 0.44 ± 0.01 | 0.53 ± 0.02 | 0.60 ± 0.02 | 0.59 ± 0.03 | 0.34 ± 0.03 | 0.68 ± 0.03 | 0.82 ± 0.03 | 0.44 ± 0.04 | 0.72 ± 0.04 |
| LaBSE | Og | 0.64 ± 0.01 | 0.44 ± 0.01 | 0.66 ± 0.02 | 0.72 ± 0.02 | 0.67 ± 0.02 | 0.48 ± 0.03 | 0.77 ± 0.03 | 0.79 ± 0.03 | 0.57 ± 0.04 | 0.77 ± 0.04 |
| MPNet-Base-Multilingual | En | 0.72 ± 0.01 | 0.62 ± 0.01 | 0.64 ± 0.02 | 0.79 ± 0.02 | 0.73 ± 0.02 | 0.58 ± 0.03 | 0.83 ± 0.03 | 0.89 ± 0.03 | 0.64 ± 0.04 | 0.80 ± 0.04 |
| MPNet-Base-Multilingual | Og | 0.64 ± 0.01 | 0.61 ± 0.01 | 0.57 ± 0.02 | 0.73 ± 0.02 | 0.61 ± 0.03 | 0.53 ± 0.03 | 0.70 ± 0.03 | 0.86 ± 0.03 | 0.56 ± 0.04 | 0.71 ± 0.04 |
| MiniLM-L2-Multilingual | En | 0.69 ± 0.01 | 0.59 ± 0.01 | 0.60 ± 0.02 | 0.75 ± 0.02 | 0.71 ± 0.02 | 0.54 ± 0.03 | 0.82 ± 0.03 | 0.86 ± 0.03 | 0.65 ± 0.04 | 0.79 ± 0.04 |
| MiniLM-L2-Multilingual | Og | 0.57 ± 0.01 | 0.57 ± 0.01 | 0.51 ± 0.02 | 0.66 ± 0.02 | 0.54 ± 0.03 | 0.49 ± 0.03 | 0.49 ± 0.03 | 0.82 ± 0.03 | 0.55 ± 0.04 | 0.61 ± 0.04 |
| XML-R | En | 0.64 ± 0.01 | 0.53 ± 0.01 | 0.60 ± 0.02 | 0.72 ± 0.02 | 0.69 ± 0.02 | 0.56 ± 0.03 | 0.81 ± 0.03 | 0.86 ± 0.03 | 0.58 ± 0.04 | 0.79 ± 0.04 |
| XML-R | Og | 0.58 ± 0.01 | 0.52 ± 0.01 | 0.55 ± 0.02 | 0.70 ± 0.02 | 0.57 ± 0.03 | 0.54 ± 0.03 | 0.51 ± 0.03 | 0.82 ± 0.03 | 0.54 ± 0.04 | 0.63 ± 0.04 |

| | Ver. | pol | slk | nld | ron | ell | ces | bul | hun | hin | mya |
|-----------------------------|------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| BM25 | En | 0.69 ± 0.04 | 0.83 ± 0.04 | 0.73 ± 0.05 | 0.84 ± 0.05 | 0.76 ± 0.06 | 0.85 ± 0.05 | 0.81 ± 0.07 | 0.81 ± 0.07 | 0.73 ± 0.08 | 0.78 ± 0.08 |
| BM25 | Og | 0.60 ± 0.04 | 0.69 ± 0.05 | 0.57 ± 0.06 | 0.71 ± 0.06 | 0.66 ± 0.06 | 0.69 ± 0.06 | 0.70 ± 0.08 | 0.51 ± 0.09 | 0.61 ± 0.09 | 0.32 ± 0.09 |
| DistilRoBERTa | En | 0.65 ± 0.04 | 0.85 ± 0.04 | 0.72 ± 0.05 | 0.86 ± 0.05 | 0.73 ± 0.06 | 0.88 ± 0.04 | 0.86 ± 0.06 | 0.82 ± 0.07 | 0.65 ± 0.09 | 0.77 ± 0.08 |
| GTR-T5-Base | En | 0.68 ± 0.04 | 0.87 ± 0.04 | 0.74 ± 0.05 | 0.88 ± 0.04 | 0.78 ± 0.05 | 0.90 ± 0.04 | 0.86 ± 0.06 | 0.87 ± 0.06 | 0.76 ± 0.08 | 0.85 ± 0.07 |
| GTR-T5-Large | En | 0.74 ± 0.04 | 0.87 ± 0.04 | 0.78 ± 0.05 | 0.88 ± 0.04 | 0.83 ± 0.05 | 0.88 ± 0.04 | 0.86 ± 0.06 | 0.89 ± 0.06 | 0.79 ± 0.08 | 0.88 ± 0.06 |
| GTR-T5-XL | En | 0.70 ± 0.04 | 0.88 ± 0.03 | 0.80 ± 0.05 | 0.87 ± 0.04 | 0.85 ± 0.05 | 0.88 ± 0.04 | 0.87 ± 0.06 | 0.86 ± 0.07 | 0.80 ± 0.08 | 0.84 ± 0.07 |
| MPNet-Base | En | 0.67 ± 0.04 | 0.84 ± 0.04 | 0.74 ± 0.05 | 0.86 ± 0.04 | 0.77 ± 0.05 | 0.89 ± 0.04 | 0.87 ± 0.06 | 0.87 ± 0.06 | 0.72 ± 0.08 | 0.80 ± 0.08 |
| MSMARCO-BERT-Base | En | 0.64 ± 0.04 | 0.85 ± 0.04 | 0.66 ± 0.06 | 0.83 ± 0.05 | 0.77 ± 0.06 | 0.86 ± 0.05 | 0.87 ± 0.06 | 0.85 ± 0.07 | 0.80 ± 0.08 | 0.82 ± 0.07 |
| MiniLM-L12 | En | 0.72 ± 0.04 | 0.86 ± 0.04 | 0.73 ± 0.05 | 0.86 ± 0.04 | 0.80 ± 0.05 | 0.86 ± 0.05 | 0.90 ± 0.05 | 0.86 ± 0.07 | 0.77 ± 0.08 | 0.80 ± 0.08 |
| MultiQA-MPNet-Base | En | 0.70 ± 0.04 | 0.86 ± 0.04 | 0.75 ± 0.05 | 0.87 ± 0.04 | 0.79 ± 0.05 | 0.87 ± 0.04 | 0.89 ± 0.06 | 0.87 ± 0.06 | 0.73 ± 0.08 | 0.84 ± 0.07 |
| SGPT-125M | En | 0.54 ± 0.05 | 0.71 ± 0.05 | 0.60 ± 0.06 | 0.75 ± 0.06 | 0.54 ± 0.06 | 0.73 ± 0.06 | 0.76 ± 0.08 | 0.78 ± 0.08 | 0.50 ± 0.09 | 0.75 ± 0.08 |
| SGPT-2.7B | En | 0.65 ± 0.04 | 0.83 ± 0.04 | 0.70 ± 0.05 | 0.83 ± 0.05 | 0.76 ± 0.06 | 0.85 ± 0.05 | 0.87 ± 0.06 | 0.83 ± 0.07 | 0.69 ± 0.09 | 0.82 ± 0.07 |
| Sentence-T5-Base | En | 0.57 ± 0.05 | 0.84 ± 0.04 | 0.62 ± 0.06 | 0.87 ± 0.04 | 0.76 ± 0.06 | 0.82 ± 0.05 | 0.87 ± 0.06 | 0.84 ± 0.07 | 0.70 ± 0.09 | 0.76 ± 0.08 |
| Sentence-T5-Large | En | 0.60 ± 0.04 | 0.86 ± 0.04 | 0.62 ± 0.06 | 0.86 ± 0.04 | 0.77 ± 0.05 | 0.86 ± 0.05 | 0.85 ± 0.07 | 0.82 ± 0.07 | 0.68 ± 0.09 | 0.80 ± 0.08 |
| Sentence-T5-XL | En | 0.66 ± 0.04 | 0.87 ± 0.04 | 0.70 ± 0.05 | 0.87 ± 0.04 | 0.81 ± 0.05 | 0.86 ± 0.05 | 0.87 ± 0.06 | 0.83 ± 0.07 | 0.73 ± 0.08 | 0.80 ± 0.08 |
| DistilUSE-Base-Multilingual | En | 0.68 ± 0.04 | 0.82 ± 0.04 | 0.70 ± 0.05 | 0.84 ± 0.05 | 0.76 ± 0.06 | 0.78 ± 0.06 | 0.89 ± 0.06 | 0.83 ± 0.07 | 0.63 ± 0.09 | 0.80 ± 0.08 |
| DistilUSE-Base-Multilingual | Og | 0.60 ± 0.04 | 0.76 ± 0.05 | 0.61 ± 0.06 | 0.80 ± 0.05 | 0.60 ± 0.06 | 0.71 ± 0.06 | 0.81 ± 0.07 | 0.78 ± 0.08 | 0.53 ± 0.09 | 0.62 ± 0.09 |
| LaBSE | En | 0.45 ± 0.05 | 0.73 ± 0.05 | 0.51 ± 0.06 | 0.77 ± 0.05 | 0.67 ± 0.06 | 0.80 ± 0.05 | 0.78 ± 0.08 | 0.83 ± 0.07 | 0.55 ± 0.09 | 0.76 ± 0.08 |
| LaBSE | Og | 0.57 ± 0.05 | 0.74 ± 0.05 | 0.61 ± 0.06 | 0.78 ± 0.05 | 0.70 ± 0.06 | 0.81 ± 0.05 | 0.84 ± 0.07 | 0.82 ± 0.07 | 0.56 ± 0.09 | 0.77 ± 0.08 |
| MPNet-Base-Multilingual | En | 0.64 ± 0.04 | 0.84 ± 0.04 | 0.71 ± 0.05 | 0.86 ± 0.05 | 0.75 ± 0.06 | 0.82 ± 0.05 | 0.85 ± 0.07 | 0.82 ± 0.07 | 0.70 ± 0.09 | 0.77 ± 0.08 |
| MPNet-Base-Multilingual | Og | 0.60 ± 0.04 | 0.79 ± 0.04 | 0.66 ± 0.06 | 0.84 ± 0.05 | 0.63 ± 0.06 | 0.78 ± 0.06 | 0.81 ± 0.07 | 0.83 ± 0.07 | 0.63 ± 0.09 | 0.75 ± 0.08 |
| MiniLM-L2-Multilingual | En | 0.66 ± 0.04 | 0.83 ± 0.04 | 0.74 ± 0.05 | 0.81 ± 0.05 | 0.75 ± 0.06 | 0.79 ± 0.05 | 0.87 ± 0.06 | 0.81 ± 0.07 | 0.62 ± 0.09 | 0.79 ± 0.08 |
| MiniLM-L2-Multilingual | Og | 0.61 ± 0.04 | 0.77 ± 0.04 | 0.64 ± 0.06 | 0.79 ± 0.05 | 0.58 ± 0.06 | 0.75 ± 0.06 | 0.83 ± 0.07 | 0.79 ± 0.08 | 0.49 ± 0.09 | 0.58 ± 0.09 |
| XML-R | En | 0.59 ± 0.04 | 0.82 ± 0.04 | 0.64 ± 0.06 | 0.84 ± 0.05 | 0.71 ± 0.06 | 0.77 ± 0.06 | 0.89 ± 0.06 | 0.78 ± 0.08 | 0.71 ± 0.09 | 0.81 ± 0.07 |
| XML-R | Og | 0.59 ± 0.04 | 0.80 ± 0.04 | 0.65 ± 0.06 | 0.79 ± 0.05 | 0.65 ± 0.06 | 0.78 ± 0.06 | 0.82 ± 0.07 | 0.80 ± 0.07 | 0.62 ± 0.09 | 0.75 ± 0.08 |

| | Ver. | eng-hin | eng-zho | eng-sin | eng-urd | eng-tgl | eng-msa | eng-tha | eng-kor | eng-mya | Other |
|-----------------------------|------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| BM25 | En | 0.35 ± 0.03 | 0.46 ± 0.04 | 0.32 ± 0.04 | 0.35 ± 0.05 | 0.37 ± 0.06 | 0.45 ± 0.06 | 0.42 ± 0.07 | 0.42 ± 0.08 | 0.38 ± 0.08 | 0.36 ± 0.03 |
| BM25 | Og | 0.04 ± 0.01 | 0.02 ± 0.01 | 0.03 ± 0.01 | 0.01 ± 0.01 | 0.11 ± 0.04 | 0.08 ± 0.03 | 0.09 ± 0.04 | 0.03 ± 0.03 | 0.10 ± 0.05 | 0.18 ± 0.02 |
| DistilRoBERTa | En | 0.33 ± 0.03 | 0.50 ± 0.04 | 0.34 ± 0.04 | 0.38 ± 0.05 | 0.37 ± 0.06 | 0.53 ± 0.06 | 0.56 ± 0.07 | 0.48 ± 0.09 | 0.43 ± 0.09 | 0.34 ± 0.03 |
| GTR-T5-Base | En | 0.44 ± 0.03 | 0.61 ± 0.04 | 0.41 ± 0.05 | 0.47 ± 0.05 | 0.46 ± 0.06 | 0.62 ± 0.06 | 0.50 ± 0.07 | 0.61 ± 0.08 | 0.52 ± 0.09 | 0.43 ± 0.03 |
| GTR-T5-Large | En | 0.51 ± 0.03 | 0.67 ± 0.04 | 0.48 ± 0.05 | 0.51 ± 0.05 | 0.56 ± 0.06 | 0.63 ± 0.06 | 0.59 ± 0.07 | 0.62 ± 0.08 | 0.52 ± 0.09 | 0.46 ± 0.03 |
| GTR-T5-XL | En | 0.51 ± 0.03 | 0.66 ± 0.04 | 0.45 ± 0.05 | 0.56 ± 0.05 | 0.53 ± 0.06 | 0.66 ± 0.06 | 0.60 ± 0.07 | 0.60 ± 0.08 | 0.53 ± 0.09 | 0.47 ± 0.03 |
| MPNet-Base | En | 0.41 ± 0.03 | 0.52 ± 0.04 | 0.36 ± 0.04 | 0.42 ± 0.05 | 0.44 ± 0.06 | 0.60 ± 0.06 | 0.55 ± 0.07 | 0.54 ± 0.09 | 0.42 ± 0.09 | 0.40 ± 0.03 |
| MSMARCO-BERT-Base | En | 0.42 ± 0.03 | 0.52 ± 0.04 | 0.38 ± 0.04 | 0.38 ± 0.05 | 0.41 ± 0.06 | 0.55 ± 0.06 | 0.50 ± 0.07 | 0.51 ± 0.09 | 0.50 ± 0.09 | 0.40 ± 0.03 |
| MiniLM-L12 | En | 0.45 ± 0.03 | 0.54 ± 0.04 | 0.40 ± 0.04 | 0.46 ± 0.05 | 0.45 ± 0.06 | 0.60 ± 0.06 | 0.55 ± 0.07 | 0.54 ± 0.09 | 0.44 ± 0.09 | 0.41 ± 0.03 |
| MultiQA-MPNet-Base | En | 0.45 ± 0.03 | 0.57 ± 0.04 | 0.41 ± 0.04 | 0.52 ± 0.05 | 0.46 ± 0.06 | 0.61 ± 0.06 | 0.56 ± 0.07 | 0.55 ± 0.08 | 0.46 ± 0.09 | 0.42 ± 0.03 |
| SGPT-125M | En | 0.24 ± 0.03 | 0.30 ± 0.04 | 0.24 ± 0.04 | 0.23 ± 0.05 | 0.19 ± 0.05 | 0.29 ± 0.06 | 0.30 ± 0.06 | 0.26 ± 0.08 | 0.24 ± 0.07 | 0.21 ± 0.02 |
| SGPT-2.7B | En | 0.42 ± 0.03 | 0.58 ± 0.04 | 0.40 ± 0.04 | 0.52 ± 0.05 | 0.45 ± 0.06 | 0.59 ± 0.06 | 0.55 ± 0.07 | 0.53 ± 0.09 | 0.52 ± 0.09 | 0.40 ± 0.03 |
| Sentence-T5-Base | En | 0.35 ± 0.03 | 0.43 ± 0.04 | 0.29 ± 0.04 | 0.35 ± 0.05 | 0.36 ± 0.06 | 0.52 ± 0.06 | 0.40 ± 0.07 | 0.37 ± 0.08 | 0.32 ± 0.08 | 0.32 ± 0.03 |
| Sentence-T5-Large | En | 0.38 ± 0.03 | 0.48 ± 0.04 | 0.34 ± 0.04 | 0.38 ± 0.05 | 0.44 ± 0.06 | 0.59 ± 0.06 | 0.44 ± 0.07 | 0.39 ± 0.08 | 0.37 ± 0.08 | 0.35 ± 0.03 |
| Sentence-T5-XL | En | 0.43 ± 0.03 | 0.52 ± 0.04 | 0.38 ± 0.04 | 0.45 ± 0.05 | 0.49 ± 0.06 | 0.63 ± 0.06 | 0.53 ± 0.07 | 0.46 ± 0.09 | 0.38 ± 0.08 | 0.38 ± 0.03 |
| DistilUSE-Base-Multilingual | En | 0.38 ± 0.03 | 0.49 ± 0.04 | 0.33 ± 0.04 | 0.38 ± 0.05 | 0.28 ± 0.05 | 0.55 ± 0.06 | 0.44 ± 0.07 | 0.46 ± 0.09 | 0.35 ± 0.08 | 0.32 ± 0.03 |
| DistilUSE-Base-Multilingual | Og | 0.23 ± 0.03 | 0.35 ± 0.04 | 0.02 ± 0.01 | 0.25 ± 0.05 | 0.11 ± 0.04 | 0.38 ± 0.06 | 0.16 ± 0.05 | 0.20 ± 0.07 | 0.10 ± 0.05 | 0.23 ± 0.02 |
| LaBSE | En | 0.33 ± 0.03 | 0.23 ± 0.04 | 0.18 ± 0.04 | 0.30 ± 0.05 | 0.13 ± 0.04 | 0.28 ± 0.06 | 0.17 ± 0.05 | 0.19 ± 0.07 | 0.24 ± 0.07 | 0.23 ± 0.02 |
| LaBSE | Og | 0.24 ± 0.03 | 0.30 ± 0.04 | 0.20 ± 0.04 | 0.22 ± 0.04 | 0.15 ± 0.04 | 0.28 ± 0.06 | 0.13 ± 0.05 | 0.23 ± 0.07 | 0.22 ± 0.07 | 0.26 ± 0.02 |
| MPNet-Base-Multilingual | En | 0.34 ± 0.03 | 0.49 ± 0.04 | 0.32 ± 0.04 | 0.39 ± 0.05 | 0.34 ± 0.06 | 0.57 ± 0.06 | 0.50 ± 0.07 | 0.42 ± 0.08 | 0.33 ± 0.08 | 0.37 ± 0.03 |
| | | | | | | | | | | | |

1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264

G.2.2 Social Media Post

ID: 11569

Published at: Facebook 2019-07-06

Verdicts: False information

Main text

Original text: *HÉPATITE B ET Cet très rassurant faite cette expérience et rêvené témoigné !!! L hépatite n'es qu'un vieux souvenirs après !!REMÈDES POUR TRAITER ET ÉRADIQUER L'HÉPATITE B DU CORPSL'hépatite B est une infection virale qui s'attaque au foie.Le virus se transmet par le sang ou lors des rapports sexuels. En effet, les seules sécrétions ou liquides corporels qui permettent de transmettre le virus sont le sang, le sperme,les sécrétions vaginales, la salive et les liquides issus d'une plaieIngrédients Une Papaye non mur Les RACINES de Papayer Femelle Les feuilles fraîches de Papayer femelle racines de Moringa feuilles fraîche Moringa 4 citrons à couper en deux. racines de cocotierPréparation-Mettez les tous dans la marmite, les feuilles en dernier position. Ajoutez de l'eau et faites bouillir le mélange.Mode d'emploiBoire 2 à 3 verres par jour.Ajoutez de l'eau à chaque fois et faites bouillir une fois par jour . Suivez le traitement pendant un mois.Faites vous examiné par un médecin et revenez témoigner .Bonne guérison...Aimes ton prochain par le partage de ce messageLa Boutique du Naturopathe Vous soigne de toutes vos maladies à l'aide des plantes naturelles moins chère et plus sure sans effets secondaire.*

Translated text: *HEPATITIS B AND Cand very reassuring made this experience and dreamed witnessed !!!Hepatitis is just an old memory afterwards!!REMEDIES TO TREAT AND ERADICATE HEPATITIS B FROM THE BODYHepatitis B is a viral infection that attacks the liver. The virus is transmitted through blood or during sexual intercourse. Indeed, the only secretions or bodily fluids that can transmit the virus are blood, semen, vaginal secretions, saliva and fluids from a wound.Ingredients An unripe Papaya The ROOTS of Female Papaya Fresh female papaya leaves Moringa roots fresh Moringa leaves 4 lemons to be cut in half. coconut rootsPreparationPut them all in the pot, the leaves last. Add water and boil the mixture.ManualDrink 2-3 glasses a day. Add water each time and boil once a day. Follow the treatment for a month. Get examined by a doctor and come back to testify.Good recovery...Love*

your neighbor by sharing this messageLa Boutique du Naturopathe Treats you to all your illnesses using cheaper and safer natural plants without side effects.

Detected languages: fra: 100.0%

OCR transcripts

Original text: *PoymeraseVirus del'hépatite BParticule filamenteuseADNAntigèneHBSParticule sphérique*

Translated text: *Polymerasevirushepatitis BFilamentous particleDNAAntigenHBSspherical particle*

Detected languages: fra: 72.4%, lb: 9.2%

G.3 Example #3

G.3.1 Fact-check

ID: 93800

Published at: 2021-12-07 factual.afp.com

Claim

Original text: *Nicolás Maduro se fotografió con una camiseta del candidato chileno Gabriel Boric*

Translated text: *Nicolás Maduro was photographed with a shirt of the Chilean candidate Gabriel Boric*

Detected languages: spa: 100.0%

Title

Original text: *El tuit de Maduro con una camiseta del candidato chileno Gabriel Boric es un doble montaje*

Translated text: *Maduro's tweet with a t-shirt of the Chilean candidate Gabriel Boric is a double montage*

Detected languages: spa: 100.0%

G.3.2 Social Media Post

ID: 20617

Published at: Facebook 2021-12-03

Verdicts: Altered photo

Main text

Original text: *Vamos con esos apoyos Gabrielito*

Translated text: *Let's go with those support Gabrielito*

Detected languages: spa: 100.0%

OCR transcripts

Original text: *Nicolás Maduro@Nicolas Maduro-*

1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314

1315 *Por la patria grande, nuestro totalapoyo desde*
1316 *Venezuela al compañeroGabriel Boric.Tik Tok*
1317 *#horicchanta #horiccorrupto*
1318 **Translated text:** *Nicholas*
1319 *Maduro@NicolasMaduroFor the great coun-*
1320 *try, our totalsupport from Venezuela to the*
1321 *comradeGabriel Boric.tik tok #horicchanta*
1322 *#horiccorrupt*
1323 **Detected languages:** spa: 50.4%, eng: 31.6%, qu:
1324 10.0%

1326 G.4 Example #4

1327 G.4.1 Fact-check

1328 **ID:** 26926

1329 **Published at:** 2022-03-16 checamos.afp.com

1330 Claim

1332 **Original text:** *As fronteiras da Ucrânia não*
1333 *foram registradas na ONU e não são reconhecidas*
1334 *internacionalmente*

1335 **Translated text:** *Ukraine's borders have not been*
1336 *registered with the UN and are not internationally*
1337 *recognized*

1338 **Detected languages:** por: 100.0%

1340 Title

1341 **Original text:** *As fronteiras da Ucrânia são*
1342 *reconhecidas e não é preciso que sejam registradas*
1343 *na ONU*

1344 **Translated text:** *Ukraine's borders are recognized*
1345 *and do not need to be registered with the UN*

1346 **Detected languages:** por: 100.0%

1348 G.4.2 Social Media Post

1349 **ID:** 8853

1350 **Published at:** Facebook 2022-02-25

1351 **Verdicts:** False information.

1352 Main text

1353 **Original text:** *E esta hein..??!!... haverá con-*
1354 *traditório..??.....” O secretário-geral das Nações*
1355 *Unidas afirmou que a Ucrânia não solicita registro*
1356 *de fronteira desde 1991, então o estado da Ucrânia*
1357 *não existe.... E não sabemos disso!!! 04/07/2014*
1358 *O secretário-geral da ONU, Ban Ki-moon, fez uma*
1359 *declaração impressionante, cuja distribuição na*
1360 *mídia ucraniana e na Internet está proibida. O*
1361 *conflito entre os dois países foi discutido na sessão*
1362 *do Conselho de Segurança da ONU. A partir disso,*
1363 *chegou-se à seguinte conclusão: A Ucrânia não*
1364

1365 *registra suas fronteiras desde 25/12/1991. A ONU*
1366 *não registrou as fronteiras da Ucrânia como um*
1367 *estado soberano. Portanto, pode-se supor que a*
1368 *Rússia não está cometendo nenhuma violação*
1369 *de direitos em relação à Ucrânia. De acordo*
1370 *com o Tratado da CEI, o território da Ucrânia*
1371 *é um distrito administrativo da URSS. Portanto,*
1372 *ninguém pode ser culpado pelo separatismo e pela*
1373 *mudança forçada das fronteiras da Ucrânia. Sob*
1374 *a lei internacional, o país simplesmente não tem*
1375 *fronteiras oficialmente reconhecidas. Para resolver*
1376 *esse problema, a Ucrânia precisa concluir a*
1377 *demarcação das fronteiras com os países vizinhos*
1378 *e obter o acordo dos países vizinhos, incluindo*
1379 *a Rússia, em sua fronteira comum. É necessário*
1380 *documentar tudo e assinar tratados com todos os*
1381 *estados vizinhos. A União Europeia prometeu o*
1382 *seu apoio à Ucrânia nesta importante questão e*
1383 *decidiu prestar toda a assistência técnica. Mas*
1384 *a Rússia assinará um tratado de fronteira com*
1385 *a Ucrânia? Não, claro que não Como a Rússia*
1386 *é a sucessora legal da URSS (isso é confirmado*
1387 *pelas decisões dos tribunais internacionais sobre*
1388 *disputas de propriedade entre a ex-URSS e países*
1389 *estrangeiros), as terras em que a Ucrânia, a*
1390 *Bielorrússia e a Novorossiya estão localizadas*
1391 *pertencem à Rússia, e ninguém tem o direito de*
1392 *ficar sem o consentimento da Rússia para dispor*
1393 *desta área. Basicamente, agora tudo o que a*
1394 *Rússia precisa fazer é declarar que essa área é*
1395 *rusa e que tudo o que acontece nessa área é um*
1396 *assunto interno da Rússia. Qualquer interferência*
1397 *será vista como uma medida contra a Rússia. Com*
1398 *base nisso, eles podem anular as eleições de 25*
1399 *de maio de 2014 e fazer o que o povo quiser!*
1400 *De acordo com o Memorando de Budapeste e*
1401 *outros acordos, a Ucrânia não tem fronteiras. O*
1402 *estado da Ucrânia não existe (e nunca existiu!).”*
1403 *Alexandre Panin*

1404 **Translated text:** *And this one huh..??!!...there will*
1405 *be a contradiction..??.....” The Secretary-General*
1406 *of the United Nations stated that Ukraine has not*
1407 *applied for border registration since 1991, so the*
1408 *state of Ukraine does not exists.... And we don't*
1409 *know that!!! 04/07/2014 The Secretary-General*
1410 *of the UN, Ban Ki-moon, made an impressive*
1411 *statement, whose distribution in the Ukrainian*
1412 *media and on the Internet is prohibited. The*
1413 *conflict between the two countries was discussed*
1414 *at the UN Security Council session. From this,*
1415 *the following conclusion was reached: Ukraine*

1416 *has not registered its borders since 12/25/1991.*
 1417 *The UN has not registered Ukraine's borders as*
 1418 *a sovereign state. Therefore, it can be assumed*
 1419 *that Russia is not committing any rights violations*
 1420 *in relation to Ukraine. According to the CIS*
 1421 *Treaty, the territory of Ukraine is an administrative*
 1422 *district of USSR. Therefore, no one can be blamed*
 1423 *for separatism and the forced change of Ukraine's*
 1424 *borders. Under international law, the country*
 1425 *simply has no officially recognized borders. To*
 1426 *solve this problem, Ukraine needs to complete*
 1427 *the demarcation of borders with neighboring*
 1428 *countries and get the agreement of neighboring*
 1429 *countries, including Russia, on their common*
 1430 *border. It is necessary to document everything*
 1431 *and sign treaties with all neighboring states. The*
 1432 *European Union pledged its support to Ukraine*
 1433 *on this important issue and decided to provide*
 1434 *full technical assistance. But will Russia sign a*
 1435 *border treaty with Ukraine? No of course not Since*
 1436 *Russia is the legal successor of the USSR (this is*
 1437 *confirmed by the decisions of international courts*
 1438 *on property disputes between the former USSR and*
 1439 *foreign countries), the lands on which Ukraine,*
 1440 *Belarus and Novorossiya are located belong to*
 1441 *Russia, and no one has the right to be without*
 1442 *Russia's consent to dispose of this area. Basically,*
 1443 *now all Russia has to do is declare that this area is*
 1444 *Russian and that everything that happens in this*
 1445 *area is an internal Russian affair. Any interference*
 1446 *will be seen as a measure against Russia. Based*
 1447 *on that, they can nullify the May 25, 2014 elections*
 1448 *and do whatever the people want! According to*
 1449 *the Budapest Memorandum and other agreements,*
 1450 *Ukraine has no borders. The state of Ukraine does*
 1451 *not exist (and never did!)."* Alexandre Panin
 1452 **Detected languages:** por: 100.0%

1454 G.5 Example #5

1455 G.5.1 Fact-check

1456 **ID:** 61827

1457 **Published at:** 2019-11-26 periksafakta.afp.com

1458 Claim

1460 **Original text:** *Gadis gembala di Maroko menjadi*
 1461 *menteri pendidikan Prancis setelah dewasa*

1462 **Translated text:** *Shepherd girl in Morocco*
 1463 *becomes French education minister as an adult*

1464 **Detected languages:** msa: 100.0%
 1465

Title

1466 **Original text:** *Ini adalah foto seorang anak*
 1467 *Maroko, bukan mantan menteri pendidikan Prancis*

1468 **Translated text:** *This is a photo of a Moroccan*
 1469 *child, not a former French education minister*

1470 **Detected languages:** msa: 100.0%
 1471
 1472

1473 G.5.2 Social Media Post

1474 **ID:** 10815

1475 **Published at:** Facebook

1476 **Verdicts:** None
 1477

Main text

1479 **Original text:** *Gadis yg disebelah kiri mengiring*
 1480 *domba di maroko, wanita yg disebelah kanan*
 1481 *adalah gadis yg sama 20 thn kemudian sbg mentri*
 1482 *pendidikan Prancis. Jgn pernah berhenti bermimpi*
 1483 *dan tdk pernah berhenti bekerja keras utk impian*
 1484 *anda.... *** Enerjik. Itulah gambaran sosok*
 1485 *Najat Vallaud-Belkacem. Dulunya, dia memakai*
 1486 *baju seadanya dengan rambut dikucir ekor*
 1487 *kuda, membawa tongkat, dan menggembalakan*
 1488 *domba. Sehari-hari dia adalah seorang gadis*
 1489 *gembala di sebuah desa kecil di dekat Nador,*
 1490 *Maroko. Saat itu tidak ada yang menduga bahwa*
 1491 *kehidupannya ketika dewasa akan berubah jauh*
 1492 *lebih baik. Menjadi menteri pendidikan dan*
 1493 *penelitian Prancis. Tentu saja posisi itu tidak*
 1494 *begitu saja datang dari langit. Belkacem berusaha*
 1495 *ekstrakeras untuk meraihnya. Di kamusnya, tak*
 1496 *ada yang tidak bisa diwujudkan. Dulu, ketika*
 1497 *dia ingin berkuliah di Paris Institute of Political*
 1498 *Studies, guru sekolahnya melarangnya mendaftar.*
 1499 *Alasannya, sekolah itu mahal sekaligus susah*
 1500 *untuk dimasuki. Namun, langkah anak kedua*
 1501 *di antara tujuh bersaudara tersebut tak surut.*
 1502 *Belkacem tetap mendaftar, belajar mati-matian,*
 1503 *dan akhirnya diterima. Dia juga harus bekerja*
 1504 *paro waktu di dua tempat untuk membayar biaya*
 1505 *kuliahnya. Di kampus itu pula, dia bertemu*
 1506 *dengan Boris Vallaud yang kini menjadi salah*
 1507 *seorang penasihat Presiden Prancis Francois*
 1508 *Hollande. Mereka sama-sama aktif di Partai*
 1509 *Sosialis. Keduanya menikah pada 27 Agustus*
 1510 *2005. Jauh sebelum itu, Belkacem juga sudah*
 1511 *terbiasa hidup keras. Saat berusia empat tahun,*
 1512 *ayahnya memboyong dia, ibu, dan kakak tertuanya,*
 1513 *Fatiha, ke Amiens, kawasan pinggiran Prancis.*
 1514 *"Ayah saya tak punya masalah. Tapi, kami, saya,*
 1515 *ibu, dan kakak, mati-matian beradaptasi dengan*
 1516 *kehidupan baru," katanya seperti dikutip Vogue.*

1517 Dia bahkan sempat terheran-heran saat melihat
1518 mobil. Hal langka di negara asalnya. Belum lagi
1519 diskriminasi yang datang dari lingkungan sekita-
1520 rnyanya. Bahkan saat dia sudah menjadi anggota
1521 parlemen di Rhone-Alpes. Dalam sebuah tulisan,
1522 Belkacem bercerita, waktu itu dirinya mengadakan
1523 perjamuan makan malam dan mengundang tamu
1524 yang belum terlalu mengenalnya. Ketika tamu itu
1525 datang, Belkacem menyambut dan membantunya
1526 melepaskan mantel. Tamu itu lantas bertanya
1527 di mana sang pemilik rumah. "Hingga saat ini
1528 di Prancis, kalau ada perempuan dengan kulit
1529 berwarna yang membuka pintu rumah di kawasan
1530 mewah, selalu dianggap pembantu," tulis ibu
1531 si kembar Louis-Adel Vallaud dan Nour-Chloe
1532 Vallaud tersebut. Sejak saat itu, dia semakin
1533 mantap mengabdikan hidup untuk menghilangkan
1534 diskriminasi. Sorotan terhadap karir gemilang
1535 Belkacem mulai terjadi saat Presiden Francois
1536 Hollande menunjuknya sebagai juru bicara
1537 pemerintah dan menteri hak-hak perempuan pada
1538 16 Mei 2012. Beberapa bulan setelah itu, Hollande
1539 memberinya tanggung jawab untuk memerangi ho-
1540 mofobia. Belkacem menjabat menteri pendidikan
1541 dan penelitian pada 25 Agustus 2014, dua hari
1542 sebelum ulang tahun kesembilan pernikahannya.
1543 Penunjukan itu menjadikan dia sebagai menteri
1544 pendidikan termuda yang pernah dipunyai Prancis.
1545 Terpilihnya Belkacem seakan menjadi bukti
1546 bahwa seorang imigran juga bisa menjadi aset
1547 yang berharga bagi negara. Apalagi dia adalah
1548 seorang muslim. Tentang Belkacem Saat masih
1549 kanak-kanak, momen terbaik dalam hidupnya
1550 adalah ketika bibliobus (mobil perpustakaan
1551 keliling) menyambangi kawasan tempat tinggalnya.
1552 Sebab, dia bisa membaca beragam buku. Memiliki
1553 dua kewarganegaraan. Salah satunya Maroko
1554 karena dia berasal dari sana. Selain itu, Prancis
1555 memberinya status warga negara saat masih ku-
1556 liah. Ia adalah Anak kedua dari tujuh bersaudara,
1557 Najat Belkacem lahir di negara Maroko pada
1558 1977 di Bni Chiker, sebuah desa dekat Nador di
1559 wilayah Rif. Pada 1982 ia bergabung kembali
1560 dengan ayahnya, seorang pekerja bangunan,
1561 dengan ibunya dan kakaknya Fatiha, dan tumbuh
1562 di subperkotaan Amiens.[3] Ia lulus dari Institut
1563 d'études politiques de Paris (Institut Studi-Studi
1564 Politik Paris) pada 2002. Di Institut ia bertemu
1565 Boris Vallaud, yang menikah dengannya pada 27
1566 Agustus 2005.[4] Ia masuk Partai Sosialis pada
1567 2002 dan bergabung dengan tim Gérard Collomb,

Walikota Lyon, pada 2003 untuk menjalankan
demokrasi lokal yang kuat, perlawanan melawan
diskriminasi, mempromosikan hak-hak warga
sipil, dan akses untuk pekerjaan dan perumahan.
Terpilih dalam Dewan Wilayah Rhone-Alpes pada
2004, ia mengetuai Komisi Budaya, mengundurkan
diri pada 2008. Pada 2005, ia menjadi penasihat
Partai Sosialis. Pada 2005 dan 2006 ia menjadi
kolumnis program kebudayaan C'est tout vu di
Télé Lyon Municipale bersama dengan Stéphane
Cayrol. Pada Februari 2007 ia bergabung
dalam tim kampanye Ségolène Royal sebagai
jurubicara, bersama dengan Vincent Peillon
dan Arnaud Montebourg. Pada Maret 2008 ia
terpilih menjadi conseillère générale departemen
Rhône dalam pemilihan kantonal dengan 58.52%
suara pada putaran kedua, dibawah spanduk
Partai Sosialis di kanton Lyon-XIII. Pada 16
Mei 2012, ia dilantik pada kabinet Presiden
Perancis François Hollande sebagai Menteri
Hak-Hak Wanita dan jurubicara pemerintahan.
<https://m.facebook.com/story.php...>

Translated text: The girl on the left is herding
sheep in Morocco, the woman on the right is the
same girl 20 years later as the French Minister of
Education. Never stop dreaming and never stop
working hard for your dreams.... *** Energetic.
That is the picture of Najat Vallaud-Belkacem. In
the past, he wore modest clothes with his hair
in a ponytail, carried a stick, and herded sheep.
Everyday she is a shepherd girl in a small village
near Nador, Morocco. At that time no one expected
that his life as an adult would change much for the
better. Became the French minister of education
and research. Of course that position didn't just
come from the sky. Belkacem tried extra hard to
reach it. In his dictionary, there is nothing that
cannot be realized. In the past, when he wanted
to study at the Paris Institute of Political Studies,
his school teacher forbade him to enroll. The
reason, the school is expensive and difficult to enter.
However, the step of the second child among the
seven siblings did not subside. Belkacem continued
to apply, studied hard, and was finally accepted.
He also had to work part-time at two places to
pay for his tuition. On the same campus, he met
Boris Vallaud, who is now an adviser to French
President Francois Hollande. They are both active
in the Socialist Party. The two were married on
August 27, 2005. Long before that, Belkacem was
also used to living hard. When he was four years

1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618

1619 old, his father took him, his mother and eldest
1620 sister, Fatiha, to Amiens, a suburb of France. "My
1621 father had no problems. But, we, me, mother and
1622 brother, are desperately adapting to a new life,"
1623 he was quoted as saying by Vogue. He even had
1624 time to be surprised when he saw the car. A rare
1625 thing in their home country. Not to mention the
1626 discrimination that comes from the surrounding
1627 environment. Even when he was already a member
1628 of parliament in the Rhone-Alpes. In an article,
1629 Belkacem recounted that at that time he held a
1630 dinner banquet and invited guests who did not
1631 know him well. When the guest arrived, Belkacem
1632 greeted him and helped him take off his coat. The
1633 guest then asked where the owner of the house
1634 was. "Until now in France, if a woman of color
1635 opened the door to a house in a luxury area, it was
1636 always considered a maid," wrote the mother of
1637 twins Louis-Adel Vallaud and Nour-Chloe Vallaud.
1638 Since then, he has been steadily devoting his
1639 life to eliminating discrimination. The spotlight
1640 on Belkacem's illustrious career began when
1641 President Francois Hollande appointed him as
1642 government spokesman and minister for women's
1643 rights on 16 May 2012. Months after that,
1644 Hollande gave him the responsibility to fight
1645 homophobia. Belkacem took office as minister of
1646 education and research on August 25, 2014, two
1647 days before her ninth wedding anniversary. The
1648 appointment makes him the youngest education
1649 minister France has ever had. The election of
1650 Belkacem seems to be proof that an immigrant can
1651 also be a valuable asset for the country. Moreover,
1652 he is a Muslim. About Belkacem When he was
1653 a child, the best moment in his life was when a
1654 bibliobus (mobile library car) visited the area
1655 where he lived. Because, he can read a variety
1656 of books. Have dual citizenship. One of them is
1657 Morocco because he is from there. In addition,
1658 France gave him the status of a citizen while still
1659 in college. The second of seven children, Najat
1660 Belkacem was born in Morocco in 1977 in Bni
1661 Chiker, a village near Nador in the Rif region. In
1662 1982 he rejoined his father, a construction worker,
1663 with his mother and sister Fatiha, and grew up in
1664 the suburb of Amiens.[3] He graduated from the
1665 Institut d'études politiques de Paris (Paris Institute
1666 of Political Studies) in 2002. At the Institute he
1667 met Boris Vallaud, whom he married on 27 August
1668 2005.[4] He joined the Socialist Party in 2002
1669 and joined the team of Gérard Collomb, Mayor of

Lyon, in 2003 to promote strong local democracy,
fight against discrimination, promote civil rights,
and access to jobs and housing. Elected to the
Rhone-Alpes County Council in 2004, he chaired
the Culture Commission, resigning in 2008. In
2005, he became an adviser to the Socialist Party.
In 2005 and 2006 he was columnist for the cultural
program C'est tout vu at Télé Lyon Municipale
together with Stéphane Cayrol. In February 2007
he joined the Ségolène Royal campaign team
as a spokesperson, along with Vincent Peillon
and Arnaud Montebourg. In March 2008 he
was elected conseillère générale of the Rhne
department in cantonal elections with 58.52% of
the vote in the second round, under the banner of
the Socialist Party in the canton of Lyon-XIII. On
16 May 2012, she was appointed to the cabinet of
French President François Hollande as Minister
of Women's Rights and spokesperson for the
government. <https://m.facebook.com/story.php...>

Detected languages: msa: 100.0%

1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691