

Overview of the CLEF–2021 CheckThat! Lab on Detecting Check-Worthy Claims, Previously Fact-Checked Claims, and Fake News

Preslav Nakov¹, Giovanni Da San Martino², Tamer Elsayed³,
Alberto Barrón-Cedeño⁴, Rubén Míguez⁵, Shaden Shaar¹, Firoj Alam¹,
Fatima Haouari³, Maram Hasanain³, Watheq Mansour³, Bayan Hamdan¹¹,
Zien Sheikh Ali³, Nikolay Babulkov⁶, Alex Nikolov⁶, Gautam Kishore Shahi⁷,
Julia Maria Struß⁸, Thomas Mandl⁹, Mucahid Kutlu¹⁰, Yavuz Selim Kartal¹⁰

¹ Qatar Computing Research Institute, HBKU, Qatar

² University of Padova, Italy

³ Qatar University, Qatar

⁴ DIT, Università di Bologna, Italy

⁵ Newtral Media Audiovisual, Spain

⁶ Sofia University, Bulgaria

⁷ University of Duisburg-Essen, Germany

⁸ University of Applied Sciences Potsdam, Germany

⁹ University of Hildesheim, Germany

¹⁰ TOBB University of Economics and Technology, Turkey

¹¹ Independent Researcher

{pnakov, sshaar, fialam}@hbku.edu.qa,
{telsayed, 200159617, maram.hasanain, wm1900793, zs1407404}@qu.edu.qa,
dasan@math.unipd.it, a.barron@unibo.it, ruben.miguez@newtral.es,
{nbabulkov, alexnickolow}@gmail.com, gautam.shahi@uni-due.de,
struss@fh-potsdam.de, mandl@uni-hildesheim.de,
{m.kutlu, ykartal}@etu.edu.tr, bayan.hamdan995@gmail.com

Abstract. We describe the fourth edition of the CheckThat! Lab, part of the 2021 Conference and Labs of the Evaluation Forum (CLEF). The lab evaluates technology supporting tasks related to factuality, and covers Arabic, Bulgarian, English, Spanish, and Turkish. Task 1 asks to predict which posts in a Twitter stream are worth fact-checking, focusing on COVID-19 and politics (in all five languages). Task 2 asks to determine whether a claim in a tweet can be verified using a set of previously fact-checked claims (in Arabic and English). Task 3 asks to predict the veracity of a news article and its topical domain (in English). The evaluation is based on mean average precision or precision at rank k for the ranking tasks, and macro- F_1 for the classification tasks. This was the most popular CLEF-2021 lab in terms of team registrations: 132 teams. Nearly one-third of them participated: 15, 5, and 25 teams submitted official runs for tasks 1, 2, and 3, respectively.

Keywords: Fact-Checking, Disinformation, Misinformation, Check-Worthiness Estimation, Verified Claim Retrieval, Fake News Detection, COVID-19.

1 Introduction

The mission of the **CheckThat!** lab is to foster the development of technology to enable the (semi-)automatic verification of claims. Systems for claim identification and verification can be very useful as supportive technology for investigative journalism, as they could provide help and guidance, thus saving time [34,45,47,97,54]. A system could automatically identify check-worthy claims, make sure they have not been fact-checked already by a reputable fact-checking organization, and then present them to a journalist for further analysis in a ranked list. Additionally, the system could identify documents that are potentially *useful* for humans to perform manual fact-checking of a claim, and it could also estimate a *veracity score* supported by evidence to increase the journalist’s understanding and trust in the system’s decision.

CheckThat! at CLEF 2021 is the fourth edition of the lab. The 2018 edition [65] focused on the identification and verification of claims in political debates. The 2019 edition [31,32] featured political debates and isolated claims, in conjunction with a closed set of Web documents to retrieve evidence from.

In 2020 [15], the focus was on social media—in particular on *Twitter*—as information posted on this platform is not checked by an authoritative entity before posting and such posts tend to disseminate very quickly. Moreover, social media posts lack context due to their short length and conversational nature; thus, identifying a claim’s context is sometimes key for effective fact-checking [23].

In the 2021 edition of the **CheckThat!** lab, we feature three tasks: 1. check-worthiness estimation, 2. detecting previously fact-checked claims, and 3. predicting the veracity of news articles and their domain. In these tasks, we focus on (i) *tweets*, (ii) *political debates and speeches*, and (iii) *news articles*. Moreover, besides Arabic and English, we extend our language coverage to Bulgarian, Spanish, and Turkish. We further add a new task (task 3) on multi-class fake news detection for news articles and topical domain identification, which can help direct the article to the right fact-checking expert[68].

2 Previously on CheckThat!

Three editions of the **CheckThat!** lab have been held so far, and some of the tasks in the 2021 edition are reformulated from previous editions. Below, we discuss some relevant tasks from previous years.

2.1 CheckThat! 2020

Task 1₂₀₂₀. Given a topic and a stream of potentially related tweets, rank the tweets by check-worthiness for the topic [43,82]. The most successful runs adopted state-of-the-art transformer models. The top-ranked teams for the English version of this task used BERT [24] and RoBERTa [70,98]. For the Arabic version, the top systems used AraBERT [52,98] and the multilingual BERT [42].

Task 2₂₀₂₀. Given a check-worthy claim and a dataset of verified claims, rank the verified claims, so that those that verify the input claim (or a sub-claim in it) are ranked on top of the list [82]. The most effective approaches fine-tuned large-scale pre-trained transformers such as BERT and RoBERTa. In particular, the top-ranked run fine-tuned RoBERTa [18].

Task 4₂₀₂₀. Given a check-worthy claim on a specific topic and a set of potentially-relevant Web pages, predict the veracity of the claim [43]. Two runs were submitted for the task [94], using a scoring function that computes the degree of concordance and negation between a claim and all input text snippets for that claim.

Task 5₂₀₂₀. Given a political debate or a speech, segmented into sentences, together with information about who the speaker of each sentence is, prioritize the sentences for fact-checking [82]. For this task, only one out of eight runs outperformed a strong bi-LSTM baseline [59].

2.2 CheckThat! 2019

Task 1₂₀₁₉. Given a political debate, an interview, or a speech, segmented into sentences, rank the sentences by the priority with which they should be fact-checked [10]. The most successful approaches used neural networks for the classification of the individual instances. For example, Hansen et al. [40] learned domain-specific word embeddings and syntactic dependencies and used an LSTM with a classification layer on top of it.

Task 2₂₀₁₉. Given a claim and a set of potentially relevant Web pages, identify which of the pages (and passages thereof) are useful for assisting a human to fact-check that claim. There was also a second subtask, asking to determine the factuality of the claim [44]. The most effective approach for this task used textual entailment and external data [35].

2.3 CheckThat! 2018

Task 1₂₀₁₈ [9] was identical to Task 1₂₀₁₉. The best approaches used *pseudo-speeches* as a concatenation of all interventions by a debater [104], and represented the entries with embeddings, part-of-speech tags, and syntactic dependencies [39].

Task 2₂₀₁₈. Given a check-worthy claim in the form of a (transcribed) sentence, determine whether the claim is likely to be true, half-true, or false [17]. The best approach retrieved relevant information from the Web, and fed the claim with the most similar Web-retrieved text to a convolutional neural network [39].

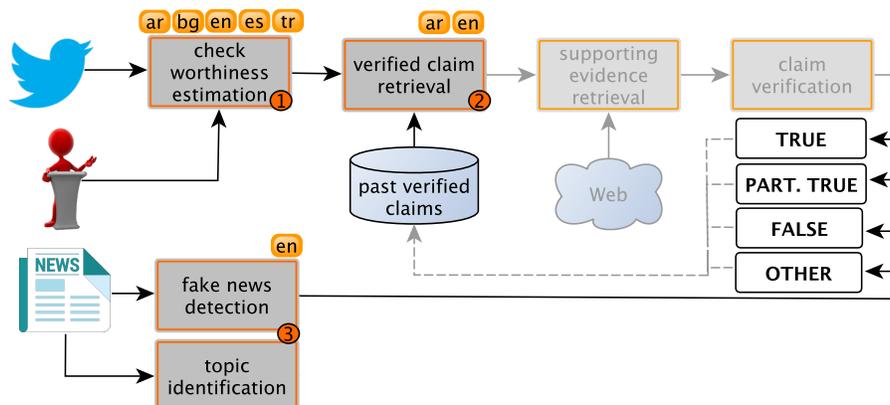


Fig. 1. The full verification pipeline. The 2021 lab covers three tasks from that pipeline: (i) check-worthiness estimation, (ii) verified claim retrieval, and (iii) fake news detection. The gray tasks were addressed in previous editions of the lab [16,32].

3 Description of the Tasks

The lab is organized around three tasks, each of which in turn has several sub-tasks. Figure 1 shows the full **CheckThat!** verification pipeline, and the three tasks we target this year are highlighted.

3.1 Task 1: Check-Worthiness Estimation

The aim of Task 1 is to determine whether a piece of text is worth fact-checking. In order to do that, we either resort to the judgments of professional fact-checkers or we ask human annotators to answer several auxiliary questions [3,4], such as “does it contain a verifiable factual claim?”, “is it harmful?” and “is it of general interest?”, before deciding on the final check-worthiness label.

Subtask 1A: Check-worthiness of tweets. Given a tweet, produce a ranked list of tweets, ordered by their check-worthiness. This is a ranking task, focusing either on COVID-19 or politics. It was offered in Arabic, Bulgarian, English, Spanish, and Turkish. The participants were free to work on any language(s) of their choice, and they could also use multilingual approaches that make use of all datasets for training.

Subtask 1B: Check-worthiness of debates or speeches. Given a political debate/speech, return a ranked list of its sentences, ordered by their check-worthiness. This is a ranking task, and it was offered in English.

3.2 Task 2: Detecting Previously Fact-Checked Claims

Given a check-worthy claim in the form of a tweet, and a set of previously fact-checked claims, rank these previously fact-checked claims in order of their usefulness to fact-check that new claim.

Subtask 2A: Detect previously fact-checked claims from tweets. Given a tweet, detect whether the claim it makes was previously fact-checked with respect to a collection of fact-checked claims. This is a ranking task, offered in Arabic and English, where the systems need to return a list of top- n candidates.

Subtask 2B: Detect previously fact-checked claims in political debates or speeches. Given a claim in a political debate or a speech, detect whether the claim has been previously fact-checked with respect to a collection of previously fact-checked claims. This is a ranking task, and it was offered in English.

3.3 Task 3: Fake News Detection

Task 3 was offered for the first time, as a pilot task. It includes two subtasks.

Subtask 3A: Multi-class fake news detection of news articles. Given the text of a news article, determine whether the claims made in the article are *true*, *partially true*, *false*, or *other*. This is a classification task, offered in English.

Subtask 3B: Given the text of a news article, determine the topical domain of the article. This is a classification task to determine the topical domain of a news article [86]. It involves six categories (health, crime, climate, election, and education), and was offered in English.

4 Datasets

Here, we briefly describe the datasets for each of the three tasks. For more details, refer to the task description paper for each individual task [80,81,88].

4.1 Task 1: Check-Worthiness Estimation

Subtask 1A: Check-worthiness for tweets. We produced datasets in five languages with tweets covering COVID-19, politics, and other topics. We refer to these datasets as the CT-CWT-21 corpus, which stands for **C**heck**T**hat! check-worthiness for tweets 2021. Table 1 shows statistics about the corpus.

For **Arabic**, the training set is sampled from the corpus used in the 2020 edition of the **C**heck**T**hat! lab [43]; we only kept tweets with full agreement between the annotators. The tweets mainly cover politics and COVID-19. The newly collected testing set covers two political events: Gulf reconciliation and US Capitol riots. They were labelled by two expert annotators, and the disagreements were resolved by discussion between the annotators.

Table 1. Task 1A (Check-worthiness in tweets): Statistics about the CT-CWT-21 corpus for all five languages. The bottom part of the table shows the main topics.

Partition	Arabic	Bulgarian	English	Spanish	Turkish	Total
Training	3,444	3,000	822	2,495	1,899	11,660
Development	661	350	140	1,247	388	2,786
Testing	600	357	350	1,248	1,013	3,568
Total	4,705	3,707	1,312	4,990	3,300	18,014

Main topics						
COVID-19	■	■	■		■	
Politics	■			■	■	

For **Bulgarian**, we created a new dataset focusing on COVID-19. The tweets were annotated by three annotators, and disagreements were resolved by majority voting, and then by a consolidator.

For **English**, the dataset also focused on COVID-19. For training, we released the data used in the **CheckThat!** lab of 2020 [82]. For testing, we annotated new instances, where we had three annotators per example, and we resolved the disagreements by majority voting, and then by a consolidator.

For **Spanish**, we had a new dataset. The tweets were manually annotated by journalists from **Newtral**—a Spanish fact-checking organization—and came from the Twitter accounts of 300 Spanish politicians.

For **Turkish**, the training set came from the **TrClaim-19** dataset [53], whereas the testing set was labelled for this task by three annotators. We applied majority voting for aggregation. The training set covers important events in Turkey in 2019 (e.g., the earthquake in Istanbul, and the military operation in Syria), whereas the test set focuses on COVID-19.

The datasets for Arabic, Bulgarian, and English have annotations for some auxiliary questions. For example, annotators were asked question such as “Is the claim of interest to the public?” and “Would the claim cause harm?”

Subtask 1B: Check-worthiness for debates/speeches. For training, we collected 57 debates/speeches from 2012–2018, and we selected sentences from the transcript that were checked by human fact-checkers. After a political debate/speech, **PolitiFact** journalists publish an article fact-checking some of the claims made in it. We collected all such sentences and considered them check-worthy, and the rest non check-worthy. However, as **PolitiFact** journalists only fact-check a few claims made in the claims, there is an abundance of false negative examples in the dataset. To address this issue at test time, we manually looked over the debates from the test set and we attempted to check whether each sentence contains a verified claim using **BM25** suggestions. Table 2 shows some statistics about the data. Note the higher proportion of positive examples in the test set compared to the training and the development sets.

Further details about the CT-CWT-21 corpus for Task 1 can be found in [81].

Table 2. Task 1B (Check-worthiness in Debates/Speeches): Statistics about the CT-CWT-21 corpus for subtask 1B.

Dataset	# of debates	# of sentences	
		Check-worthy	Non-check-worthy
Training	40	429	41,604
Development	9	69	3,517
Test	8	298	5,002
Total	57	796	50,123

4.2 Task 2: Detecting Previously Fact-Checked Claims

Subtask 2A: Detecting previously fact-checked claims from tweets. For **English**, we have 1,401 annotated tweets, each matching a single claim in a set of 13,835 verified claims from Snopes.

For **Arabic**, we have 858 tweets, matching 1,039 verified claims (some tweets match more than one verified claim) in a collection of 30,329 previously fact-checked claims. The latter include 5,921 Arabic claims from AraFacts [5] and 24,408 English claims from ClaimsKG [92], translated to Arabic using the Google translate API (<http://cloud.google.com/translate>).

Subtask 2B: Detecting previously fact-checked claims in political debates/speeches. We have 669 claims from political debates [79], matched against 804 verified claims (some input claims match more than one verified claim) in a collection of 19,250 verified claims in PolitiFact.

Table 3 shows statistics about the CT-VCR-21 corpus for Task 2, including both subtasks and languages. CT-VCR-21 stands for **C**heck**T**hat! verified claim retrieval 2021. *Input-VerClaim* pairs represent input claims with their corresponding verified claims by a fact-checking source. The input for subtask 2A (2B) is a tweet (sentence from a political debate or a speech). More details about the corpus construction can be found in [80].

4.3 Task 3: Fake News Detection

The process of corpus creation for Task 3 extends the AMUSED framework [83]. Starting with articles written by fact-checking organizations, we scraped the links to the original articles they verified, together with the factuality judgments. This process was done in two steps. First, in an automatic filtering step, all links with posts from social media channels or to multimedia documents were filtered out. In a second step, the remaining links were subjected to a manual checking process. During this step, we additionally made sure that the scraped link actually pointed to the checked document and that the document still existed (thus, eliminating error pages, articles with other content, etc.). After successful verification for each article, we scraped its title and full text.

Table 3. Task 2: Statistics about the CT-VCR-21 corpus, including the number of *Input-VerClaim* pairs and the number of *VerClaim* claims to match the input claim against.

	2A-Arabic	2A-English	2B-English
Input claims	858	1,401	669
Training	512	999	472
Development	85	200	119
Test	261	202	78
Input-<i>VerClaim</i> pairs	1,039	1,401	804
Training	602	999	562
Development	102	200	139
Test	335	202	103
Verified claims (to match against)	30,329	13,835	19,250

Subtask 3A: Multi-class fake news categorization of news articles.

This subtask was offered in English only. We collected a total of 900 news articles for training and 354 news articles for testing from 11 fact-checking websites such as PolitiFact. The label for the original fact-checking site was given as a rating. However, due to the heterogeneous labeling schemes of different fact-checking organizations (e.g., *false*: incorrect, inaccurate, misinformation; *partially false*: mostly false, half false), we merged labels with shared meaning according to [84], resulting in the following four classes: *false*, *partially false*, *true* and *other*. We provided an ID, the title of the article, the text of the article, and our rating as data to the participants. No further metadata about the article was made available in the dataset. The ID is a unique identifier created for the dataset, the title is the title given in the target article, the text is the full-text content of the article, and our rating is the normalized rating provided in one of the above four label categories.

Subtask 3B: Topical domain identification of news articles. This subtask is also offered in English only. We annotated a subset of the articles from subtask 3A with their topic: 318 articles for training, and 137 articles for testing in six different classes as shown in Table 4 based on [85]. We refer to the corpus as CT-FAN-21, which stands for **C**heck**T**hat! 2021 Fake News. We provided the ID, the title, the text, and our rating as the metadata for the dataset. Here, ID is the unique ID, title is the title of the fake news article, the text is the full-text content of the article, and domain is the domain, expressed in terms of one of the above six categories.

The datasets for subtasks 3A and 3B are available in Zenodo [87]. We did not provide any other information (e.g., a link to the article, a publication date, eventual tags, authors, location of publication, etc.).

Table 4. Task 3: Statistics about the number of documents and class distribution for the CT-FAN-21 corpus for fake news detection (left) and for topic identification (right).

Class	Training	Test	Topic	Training	Test
False	465	111	Health	127	54
True	142	65	Climate	49	21
Partially false	217	138	Economy	43	19
Other	76	40	Crime	39	17
Total	900	354	Elections	32	14
			Education	28	12
			Total	318	137

5 Evaluation

For the ranking tasks, as in the two previous editions of the **CheckThat!** lab, we used *Mean Average Precision* (MAP) as the official evaluation measure. We further calculated and reported reciprocal rank, and $P@k$ for $k \in \{1, 3, 5, 10, 20, 30\}$, as unofficial measures. For the classification tasks, we used accuracy and macro- F_1 score.

6 Results for Task 1: Check-Worthiness Estimation

Below, we report the evaluation results for task 1 and its two subtasks for all five languages.

6.1 Task 1A. Check-Worthiness of Tweets

Fifteen teams took part in this task, with English and Arabic being the most popular languages. Four out of the fifteen teams submitted runs for all five languages —most of them having trained independent models for each language (yet, team UPV trained a single multilingual model). For all five languages, we had a monolingual baseline based on n -gram representations. Table 5 shows the performance of the official submissions on the test set, in addition to the n -gram baseline. The official run was the last valid blind submission by each team. The table shows the runs ranked on the basis of the official MAP measure and includes all five languages.

Arabic Eight teams participated for Arabic, submitting a total of 17 runs (yet, recall that only the last submission counts). All participating teams fine-tuned existing pre-trained models, such as AraBERT, and multilingual BERT models. We can see that the top two systems additionally worked on improved training datasets. Team **Accenture** used a label augmentation approach to increase the number of positive examples, while team **bigIR** augmented the training set with the Turkish training set (which they automatically translated to Arabic).

Table 5. Task 1A: results for the official submissions in all five languages.

Team	MAP	MRR	RP	P@1	P@3	P@5	P@10	P@20	P@30
Arabic									
1 Accenture [99]	0.658	1.000	0.599	1.000	1.000	1.000	1.000	0.950	0.840
2 bigIR	0.615	0.500	0.579	0.000	0.667	0.600	0.600	0.800	0.740
3 SCUoL [6]	0.612	1.000	0.599	1.000	1.000	1.000	1.000	0.950	0.780
4 iCompass	0.597	0.333	0.624	0.000	0.333	0.400	0.400	0.500	0.640
4 QMUL-SDS [1]	0.597	0.500	0.603	0.000	0.667	0.600	0.700	0.650	0.720
6 TOBB ETU [100]	0.575	0.333	0.574	0.000	0.333	0.400	0.400	0.500	0.680
7 DamascusTeam	0.571	0.500	0.558	0.000	0.667	0.600	0.800	0.700	0.640
8 UPV [14]	0.548	1.000	0.550	1.000	0.667	0.600	0.500	0.400	0.580
9 ngram-baseline	0.428	0.500	0.409	0.000	0.667	0.600	0.500	0.450	0.440
Bulgarian									
1 bigIR	0.737	1.000	0.632	1.000	1.000	1.000	1.000	1.000	0.800
2 UPV [14]	0.673	1.000	0.605	1.000	1.000	1.000	1.000	0.800	0.700
3 ngram-baseline	0.588	1.000	0.474	1.000	1.000	1.000	0.900	0.750	0.640
4 Accenture [99]	0.497	1.000	0.474	1.000	1.000	0.800	0.700	0.600	0.440
5 TOBB ETU [100]	0.149	0.143	0.039	0.000	0.000	0.000	0.200	0.100	0.060
English									
1 NLP&IR@UNED [49]	0.224	1.000	0.211	1.000	0.667	0.400	0.300	0.200	0.160
2 Fight for 4230 [102]	0.195	0.333	0.263	0.000	0.333	0.400	0.400	0.250	0.160
3 UPV [14]	0.149	1.000	0.105	1.000	0.333	0.200	0.200	0.100	0.120
4 bigIR	0.136	0.500	0.105	0.000	0.333	0.200	0.100	0.100	0.120
5 GPLSI [77]	0.132	0.167	0.158	0.000	0.000	0.000	0.200	0.150	0.140
6 csum112	0.126	0.250	0.158	0.000	0.000	0.200	0.200	0.150	0.160
7 abaruah	0.121	0.200	0.158	0.000	0.000	0.200	0.200	0.200	0.140
8 NLytics [75]	0.111	0.071	0.053	0.000	0.000	0.000	0.000	0.050	0.120
9 Accenture [99]	0.101	0.143	0.158	0.000	0.000	0.000	0.200	0.200	0.100
10 TOBB ETU [100]	0.081	0.077	0.053	0.000	0.000	0.000	0.000	0.050	0.080
11 ngram-baseline	0.052	0.020	0.000	0.000	0.000	0.000	0.000	0.000	0.020
Spanish									
1 TOBB ETU [100]	0.537	1.000	0.525	1.000	1.000	0.800	0.900	0.700	0.680
2 GPLSI [77]	0.529	0.500	0.533	0.000	0.667	0.600	0.800	0.750	0.620
3 bigIR	0.496	1.000	0.483	1.000	1.000	0.800	0.800	0.600	0.620
4 NLP&IR@UNED [49]	0.492	1.000	0.475	1.000	1.000	1.000	0.800	0.800	0.620
5 Accenture [99]	0.491	1.000	0.508	1.000	0.667	0.800	0.900	0.700	0.620
6 ngram-baseline	0.450	1.000	0.450	1.000	0.667	0.800	0.700	0.700	0.660
7 UPV	0.446	0.333	0.475	0.000	0.333	0.600	0.800	0.650	0.580
Turkish									
1 TOBB ETU [100]	0.581	1.000	0.585	1.000	1.000	0.800	0.700	0.750	0.660
2 SU-NLP [22]	0.574	1.000	0.585	1.000	1.000	1.000	0.800	0.650	0.680
3 bigIR	0.525	1.000	0.503	1.000	1.000	1.000	0.800	0.700	0.720
4 UPV [14]	0.517	1.000	0.508	1.000	1.000	1.000	1.000	0.850	0.700
5 Accenture [99]	0.402	0.250	0.415	0.000	0.000	0.400	0.400	0.650	0.660
6 ngram-baseline	0.354	1.000	0.311	1.000	0.667	0.600	0.700	0.600	0.460

Bulgarian Four teams took part for Bulgarian, submitting a total of 11 runs. The top-ranked team was **bigIR**. They did not submit a task description paper, and thus we cannot give much detail about their system. Team **UPV** is the second best system, and they used multilingual sentence transformer representation (SBERT) with knowledge distillation. They also introduced an auxiliary language identification task, aside from the downstream check-worthiness task.

English Ten teams took part in task 1A for English, with a total of 21 runs. The top-ranked team was **NLP&IR@UNED**, and they fine-tuned several pre-trained transformers models. They reported BERTweet was best on the development set. The model was trained using RoBERTa on 850 million English tweets and 23 million COVID-19 related English tweets. The second best system (Team **Fight for 4230**) also used BERTweet with a dropout layer. It also included pre-processing and data augmentation.

Spanish Six teams took part for Spanish, with a total of 13 runs. The top team **TOBB ETU** explored different data augmentation strategies, including machine translation and weak supervision. However, they submitted a fine-tuned BETO model without any data augmentation. The first runner up **GPLSI** opted for using the BETO Spanish transformer together with a number of hand-crafted features, such as the presence of numbers or words in the LIWC lexicon.

Turkish Five teams participated for Turkish, submitting a total of 9 runs. All participants used BERT-based models. The top ranked team **TOBB ETU** fine-tuned BERTurk after removing user mentions and URLs. The runner up team **SU-NLP** applied a pre-processing step that includes removing hashtags, emojis, and replacing URLs and mentions with special tokens. Subsequently, they used an ensemble of BERTurk models fine-tuned with different seed values. The third-ranked team **bigIR** machine-translated the Turkish text to Arabic and then fine-tuned AraBERT on the translated text.

All languages. Table 6 summarizes the MAP performance of all the teams that submitted predictions for all languages in Task 1A. We can see that team **BigIR** performed best overall.

Table 6. MAP performance for the official submissions to **Task 1A** in all five languages. μ shows a standard mean of the five MAP scores; μ_w shows a weighed mean, where each MAP is multiplied by the size of the testing set.

Team	ar	bg	en	es	tr	μ	μ_w
1 bigIR	0.615	0.737	0.136	0.496	0.525	0.502	0.513
2 UPV [14]	0.548	0.673	0.149	0.446	0.517	0.467	0.477
3 TOBB ETU [100]	0.575	0.149	0.081	0.537	0.581	0.385	0.472
4 Accenture [99]	0.658	0.497	0.101	0.491	0.402	0.430	0.456
5 ngram-baseline	0.428	0.588	0.052	0.450	0.354	0.374	0.394

Table 7. Task 1B (English): Official evaluation results, in terms of MAP, MRR, R-Precision, and Precision@ k . The teams are ranked by the official evaluation measure: MAP.

Rank	Team	MAP	MRR	RP	P@1	P@3	P@5	P@10	P@20	P@30
1	Fight for 4230 [102]	0.402	0.917	0.403	0.875	0.833	0.750	0.600	0.475	0.350
2	ngram-baseline	0.235	0.792	0.263	0.625	0.583	0.500	0.400	0.331	0.217
3	NLytics [75]	0.135	0.345	0.130	0.250	0.125	0.100	0.137	0.156	0.135

6.2 Task 1B. Check-Worthiness of Debates/Speeches

Two teams took part in this subtask, submitting a total of 3 runs. Table 7 shows the performance of the official submissions on the test set, in addition to the ngram baseline. Similarly to Task 1A, the official run was the last valid blind submission by each team. The table shows the runs ranked on the basis of the official MAP measure.

The top-ranked team, **Fight for 4230**, fine-tuned BERTweet after normalizing the claims, augmenting the data using WordNet-based substitutions and removal of punctuation. They were able to beat the ngram baseline by 18 MAP points absolute.

7 Results for Task 2: Verified Claim Retrieval

7.1 Subtask 2A: Detecting Previously Fact-Checked Claims in Tweets

Table 8 shows the official results for Task 2A in both Arabic and English. A total of four teams participated in this task, and they all managed to improve over the Elastic Search (ES) baseline.

Arabic One team, bigIR, submitted a run for this subtask. They used AraBERT to rerank a list of candidates retrieved by a BM25 model. Their approach consists of three main steps. First, constructing a balanced training dataset, where the positive examples correspond to the query relevances (qrels) provided by the organizers, while the negative examples were selected from the top retrieved candidates by BM25 such that they were not already labeled as positive. Second, they fine-tuned AraBERT to predict the relevance score for a given tweet-VerClaim pair. They added two neural network layers on top of AraBERT to perform the classification task. Finally, at inference time, they first used BM25 to retrieve the top-20 candidate verified claims. Then, they fed each tweet-VerClaim pair to the fine-tuned model to get a relevance score and to rerank the candidate claims accordingly. As Table 8 shows, team **bigIR** outperformed the Elastic Search baseline by a good margin achieving a MAP@5 of 0.908 versus 0.794 for the baseline.

Table 8. Task 2A: Official evaluation results, in terms of MRR, MAP@ k , and Precision@ k . The teams are ranked by the official evaluation measure: MAP@5. Here, *ES-baseline* refers to the Elastic Search baseline.

Team	MRR	MAP					Precision				
		@1	@3	@5	@10	@20	@1	@3	@5	@10	@20
Arabic											
1 bigIR	0.924	0.787	0.905	0.908	0.910	0.912	0.908	0.391	0.237	0.120	0.061
2 ES-baseline	0.835	0.682	0.782	0.794	0.799	0.802	0.793	0.344	0.217	0.113	0.058
English											
1 Aschern [25]	0.884	0.861	0.880	0.883	0.884	0.884	0.861	0.300	0.182	0.092	0.046
2 NLytics [75]	0.807	0.738	0.792	0.799	0.804	0.806	0.738	0.289	0.179	0.093	0.048
3 DIPS [60]	0.795	0.728	0.778	0.787	0.791	0.794	0.728	0.282	0.177	0.092	0.048
4 ES-baseline	0.761	0.703	0.741	0.749	0.757	0.759	0.703	0.262	0.164	0.088	0.046

English Three teams participated for English, submitting a total of ten runs. All of them managed to improve over the Elastic Search (ES) baseline by a large margin. Team **Aschern** had the top-ranked system, which used TF.IDF, fine-tuned pre-trained sentence-BERT, and the reranking LambdaMART model. The system is 13.4 (MAP@5) points absolute above the baseline. The second best system is the **NLytics**, which used RoBERTa to train their model and this system was 5 (MAP@5) point above the baseline.

7.2 Subtask 2B: Detecting Previously Fact-Checked Claims in Political Debates and Speeches

Table 9 shows the official results for Task 2B, which was offered in English only. We can see that only three teams participated in this subtask, submitting a total of five runs, and no team managed to beat the Elastic Search (ES) baseline, which was based on BM25.

Among the three participating teams, Team **DIPS** was the top-ranked one. They used sentence BERT (S-BERT) embeddings for all claims, and computed the cosine similarity for each pair of an input claim and a verified claim from the dataset of previously fact-checked claims. Their prediction was made by passing a sorted list of cosine similarities to a neural network. Team **BeaSkU** was the second-best team, which used a triplet loss training method to perform fine-tuning of the S-BERT model. Then, they used the scores predicted by the fine-tuned model along with BM25 scores as features to train a reranker based on rankSVM. In addition, they discussed the impact of applying online mining of triplets. They also performed some experiments aiming at augmenting the training dataset with additional examples.

Table 9. Task 2B (English): Official evaluation results, in terms of MAP, MAP@ k , and Precision@ k . The teams are ranked by the official evaluation measure: MAP@5.

Team	MRR	MAP					Precision				
		@1	@3	@5	@10	@20	@1	@3	@5	@10	@20
1 ES-baseline	0.350	0.304	0.339	0.346	0.351	0.353	0.304	0.143	0.091	0.052	0.027
2 DIPS [60]	0.336	0.278	0.313	0.328	0.338	0.342	0.266	0.143	0.099	0.059	0.032
3 Beasku [90]	0.320	0.266	0.308	0.327	0.332	0.332	0.253	0.139	0.101	0.056	0.028
4 NLytics [75]	0.216	0.171	0.210	0.215	0.219	0.222	0.165	0.101	0.068	0.038	0.022

8 Overview of Task 3: Fake News Detection

In this section, we present an overview of all task submissions for tasks 3A and 3B. Overall, there were 88 submissions by 27 teams for Task 3A and 49 submissions by 20 teams for task 3B. For task 3, unlike the other tasks, each participant could submit up to 5 runs. After evaluation, we found that two teams from task 3A and seven teams from task 3B submitted the wrong files, and thus we have not considered them for evaluation; we report the ranking for 25 teams for task 3A and 13 teams for task 3B. In Tables 10 and 11, we report the best submission of each team for task 3A and 3B, respectively. In the following sections, we report the results for each of the subtasks.

8.1 Task 3A. Multi-Class Fake News Detection of News Articles

Most teams used deep learning models and in particular the transformer architecture for this pilot task. There have been no attempts to model knowledge with semantic technology, e.g., argument processing [30].

The best submission (team **NoFake**) was ahead of the rest by a rather large margin and achieved a macro-F1 score of 0.838. They applied BERT and made extensive use of external resources and in particular downloaded collections of misinformation datasets from fact-checking sites. The second best submission (team **Saud**) achieved a macro-F1 score of 0.503 and used lexical features, traditional weighting methods as features, and standard machine learning algorithms. This shows, that traditional approaches can still outperform deep learning models for this task. Many teams used BERT and its newer variants. Such systems are ranked after the second position. The most popular model was RoBERTa, which was used by seven teams. Team **MUCIC** used a majority voting ensemble with three BERT variants [12]. The participating teams that used BERT had to find solutions for handling the length of the input: BERT and its variants have limitations on the length of their input, but the length of texts in the CT-FAN-21 dataset, which consists of newspaper articles, is much longer. In most cases, heuristics were used for the selection of part of the text. Overall, most submissions achieved a macro-F1 score below 0.5.

Table 10. Task 3A: Performance of the best run per team based on F_1 score for individual classes, and accuracy and macro- F_1 for the overall measure.

Team	True	False	Partially False	Other	Accuracy	Macro-F1
1 NoFake*[56]	0.824	0.862	0.879	0.785	0.853	0.838
2 Saud*	0.321	0.615	0.502	0.618	0.537	0.514
3 DLRG* [50]	0.250	0.588	0.519	0.656	0.528	0.503
4 NLP&IR@UNED [49]	0.247	0.629	0.536	0.459	0.528	0.468
5 NITK_NLP [57]	0.196	0.617	0.523	0.459	0.517	0.449
6 UAICS [26]	0.442	0.470	0.482	0.391	0.458	0.446
7 CIVIC-UPM [48]	0.268	0.577	0.472	0.340	0.463	0.414
8 Uni. Regensburg [41]	0.231	0.489	0.497	0.400	0.438	0.404
9 Pathfinder* [95]	0.277	0.517	0.451	0.360	0.452	0.401
10 CIC* [8]	0.205	0.542	0.490	0.319	0.410	0.389
11 Black Ops [91]	0.231	0.518	0.327	0.453	0.427	0.382
12 NLytics*	0.130	0.575	0.522	0.318	0.475	0.386
13 Nkovachevich [55]	0.237	0.643	0.552	0.000	0.489	0.358
14 talhaanwar*	0.283	0.407	0.435	0.301	0.367	0.357
15 abaruah	0.165	0.531	0.552	0.125	0.455	0.343
16 Team GPLSI[77]	0.293	0.602	0.226	0.092	0.356	0.303
17 Sigmoid [76]	0.222	0.345	0.323	0.154	0.291	0.261
18 architap	0.154	0.291	0.394	0.187	0.294	0.257
19 MUCIC [12]	0.143	0.446	0.275	0.070	0.331	0.233
20 Probity	0.163	0.401	0.335	0.033	0.302	0.233
21 M82B [7]	0.130	0.425	0.241	0.094	0.305	0.223
22 Spider	0.046	0.482	0.145	0.069	0.316	0.186
23 Qword [96]	0.108	0.458	0.000	0.033	0.277	0.150
24 ep*	0.060	0.479	0.000	0.000	0.319	0.135
25 azaharudue*	0.060	0.479	0.000	0.000	0.319	0.135
Majority class baseline	0.000	0.477	0.000	0.000	0.314	0.119

* Runs submitted after the deadline, but before the release of the results.

The second most popular neural network model was the recurrent neural network, which was used by six teams. Many participants experimented also with traditional text processing methods as they were commonly used for knowledge representation in information retrieval. For example, team **Kovachevich** used a Naïve Bayes classifier with TF.IDF features for the 500 most frequent stems in the dataset [55]. Some lower-ranked teams used additional techniques and resources. These include LIWC [49], data augmentation by inserting artificially created similar documents [8], semantic analysis with the Stanford Empath Tool [26], and the reputation of the sites of a search engine result after searching with the title of the article [49].

Table 11. Task 3B: Performance of the best run per team based on F1-measure for individual classes, and accuracy and macro-F₁ for overall measure.

Team	Climate	Crime	Economy	Education	Elections	Health	Acc	Macro F1
1 NITK_NLP [57]	0.950	0.872	0.824	0.800	0.897	0.946	0.905	0.881
2 NoFake*	0.800	0.875	0.900	0.957	0.692	0.907	0.869	0.855
3 Nkovachevich [55]	0.927	0.872	0.743	0.737	0.857	0.911	0.869	0.841
4 DLRG	0.952	0.743	0.688	0.800	0.828	0.897	0.847	0.818
5 CIC* [8]	0.952	0.750	0.688	0.588	0.889	0.871	0.832	0.790
6 architap	0.900	0.711	0.774	0.609	0.815	0.907	0.825	0.786
7 NLytics	0.826	0.714	0.710	0.500	0.769	0.867	0.788	0.731
8 CIVIC-UPM* [48]	0.864	0.700	0.645	0.421	0.609	0.821	0.745	0.677
9 ep*	0.727	0.476	0.222	0.343	0.545	0.561	0.511	0.479
10 Pathfinder* [95]	0.900	0.348	0.250	0.000	0.526	0.667	0.599	0.448
11 M82B [7]	0.294	0.000	0.000	0.000	0.000	0.576	0.409	0.145
12 MUCIC [12]	0.294	0.000	0.000	0.000	0.000	0.576	0.409	0.145
13 azaharudue*	0.129	0.000	0.000	0.125	0.000	0.516	0.321	0.128
<i>Majority class baseline</i>	0.000	0.000	0.000	0.000	0.000	0.565	0.394	0.094

* Runs submitted after the deadline, but before the release of the results.

8.2 Task 3B. Topical Domain Identification of News Articles

The performance of the systems for task 3B was overall higher than for task 3A. The first three submissions were close together and all used transformer-based architectures. The best submission, by team **NITK_NLP**, used an ensemble of three transformers [57]. The second best submission (by team **NoFake**) and the third best submission (by team **Nkovachevich**) used BERT.

9 Related Work

There has been work on checking the factuality/credibility of a claim, of a news article, or of an information source [11,13,51,58,64,69,73,103]. Claims can come from different sources, but special attention has been paid to those from social media [37,62,66,78,79,89,101]. Check-worthiness estimation is still a fairly-new problem especially in the context of social media [34,45,46,47]. A lot of research was performed on fake news detection for news articles, which is mostly approached as a binary classification problem [71].

CheckThat! is related to several other initiatives at SemEval on determining rumour veracity and support for rumours [28,36], on stance detection [63], on fact-checking in community question answering forums [61], on propaganda detection [27,29], and on semantic textual similarity [2,67]. It is also related to the FEVER task [93] on fact extraction and verification, as well as to the Fake News Challenge [38], and the FakeNews task at MediaEval [72].

10 Conclusion and Future Work

We have presented the 2021 edition of the **CheckThat!** Lab, which was the most popular CLEF-2021 lab in terms of team registrations (132 teams registered), and about one-third of them actually participated: 15, 5, and 25 teams submitted official runs for tasks 1, 2, and 3, respectively. The lab featured tasks that span important steps of the verification pipeline: from spotting check-worthy claims to checking whether they have been fact-checked elsewhere before. We further featured a fake news detection task, and we also checked the class and the topical domain of news articles. Together, these tasks support the technology pipeline to assist human fact-checkers. Moreover, in-line with the general mission of CLEF, we promoted multi-linguality by offering our tasks in five different languages.

In future work, we plan to extend the datasets with more examples, more information sources, and also to cover more languages.

Acknowledgments

The work of Tamer Elsayed and Maram Hasanain is made possible by NPRP grant #NPRP-11S-1204-170060 from the Qatar National Research Fund (a member of Qatar Foundation). The work of Fatima Haouari is supported by GSRA grant #GSRA6-1-0611-19074 from the Qatar National Research Fund. The statements made herein are solely the responsibility of the authors.

This research is also part of the Tanbih mega-project, developed at the Qatar Computing Research Institute, HBKU, which aims to limit the impact of “fake news”, propaganda, and media bias, thus promoting digital literacy and critical thinking.

References

1. Abumansour, A., Zubiaga, A.: QMUL-SDS at CheckThat! 2021: Enriching pre-trained language models for the estimation of check-worthiness of Arabic tweets. In: Faggioli et al. [33]
2. Agirre, E., Banea, C., Cer, D., Diab, M., Gonzalez-Agirre, A., Mihalcea, R., Rigau, G., Wiebe, J.: SemEval-2016 task 1: Semantic textual similarity, monolingual and cross-lingual evaluation. In: Proceedings of the 10th International Workshop on Semantic Evaluation. pp. 497–511. SemEval '16 (2016)
3. Alam, F., Dalvi, F., Shaar, S., Durrani, N., Mubarak, H., Nikolov, A., Da San Martino, G., Abdelali, A., Sajjad, H., Darwish, K., Nakov, P.: Fighting the COVID-19 infodemic in social media: A holistic perspective and a call to arms. In: Proceedings of the International AAAI Conference on Web and Social Media. ICWSM '21, vol. 15, pp. 913–922 (2021)
4. Alam, F., Shaar, S., Dalvi, F., Sajjad, H., Nikolov, A., Mubarak, H., Da San Martino, G., Abdelali, A., Durrani, N., Darwish, K., Nakov, P.: Fighting the COVID-19 infodemic: Modeling the perspective of journalists, fact-checkers, social media platforms, policy makers, and the society. ArXiv preprint 2005.00033 (2020)

5. Ali, Z.S., Mansour, W., Elsayed, T., Al-Ali, A.: AraFacts: The first large Arabic dataset of naturally occurring claims. In: Proceedings of the Sixth Arabic Natural Language Processing Workshop. pp. 231–236. ANLP '21 (2021)
6. Althabiti, S., Alsalka, M., Atwell, E.: An AraBERT model for check-worthiness of Arabic tweets. In: Faggioli et al. [33]
7. Ashik, S.S., Apu, A.R., Marjana, N.J., Hasan, M.A., Islam, M.S.: M82B at CheckThat! 2021: Multiclass fake news detection using BiLSTM based RNN model. In: Faggioli et al. [33]
8. Ashraf, N., Butt, S., Sidorov, G., Gelbukh, A.: Fake news detection using machine learning and data augmentation – CLEF2021. In: Faggioli et al. [33]
9. Atanasova, P., Marquez, L., Barrón-Cedeño, A., Elsayed, T., Suwaileh, R., Zaghoulani, W., Kyuchukov, S., Da San Martino, G., Nakov, P.: Overview of the CLEF-2018 CheckThat! lab on automatic identification and verification of political claims. Task 1: Check-worthiness. In: Cappellato et al. [21]
10. Atanasova, P., Nakov, P., Karadzhov, G., Mohtarami, M., Da San Martino, G.: Overview of the CLEF-2019 CheckThat! lab on automatic identification and verification of claims. Task 1: Check-worthiness. In: Cappellato et al. [20]
11. Ba, M.L., Berti-Equille, L., Shah, K., Hammady, H.M.: VERA: A platform for veracity estimation over web data. In: Proceedings of the 25th International Conference on World Wide Web. pp. 159–162. WWW '16 (2016)
12. Balouchzahi, F., Shashirekha, H., Sidorov, G.: MUCIC at CheckThat! 2021: FaDo-fake news detection and domain identification using transformers ensembling. In: Faggioli et al. [33]
13. Baly, R., Karadzhov, G., An, J., Kwak, H., Dinkov, Y., Ali, A., Glass, J., Nakov, P.: What was written vs. who read it: News media profiling using text analysis and social media context. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics. pp. 3364–3374. ACL '20 (2020)
14. Baris Schlicht, I., Magnossão de Paula, A., Rosso, P.: UPV at CheckThat! 2021: Mitigating cultural differences for identifying multilingual check-worthy claims. In: Faggioli et al. [33]
15. Barrón-Cedeño, A., Elsayed, T., Nakov, P., Da San Martino, G., Hasanain, M., Suwaileh, R., Haouari, F., Babulkov, N., Hamdan, B., Nikolov, A., Shaar, S., Ali, Z.: Overview of CheckThat! 2020 — automatic identification and verification of claims in social media. In: Proceedings of the 11th International Conference of the CLEF Association: Experimental IR Meets Multilinguality, Multimodality, and Interaction. pp. 215–236. CLEF '2020 (2020)
16. Barrón-Cedeño, A., Elsayed, T., Nakov, P., Da San Martino, G., Hasanain, M., Suwaileh, R., Haouari, F., Babulkov, N., Hamdan, B., Nikolov, A., Shaar, S., Sheikh Ali, Z.: Proceedings of the eleventh international conference of the CLEF association: Experimental IR meets multilinguality, multimodality, and interaction. In: Arampatzis, A., Kanoulas, E., Tsikrika, T., Vrochidis, S., Joho, H., Lioma, C., Eickhoff, C., Névél, A., Cappellato, L., Ferro, N. (eds.) Overview of CheckThat! 2020: Automatic Identification and Verification of Claims in Social Media. CLEF '20, Springer (2020)
17. Barrón-Cedeño, A., Elsayed, T., Suwaileh, R., Marquez, L., Atanasova, P., Zaghoulani, W., Kyuchukov, S., Da San Martino, G., Nakov, P.: Overview of the CLEF-2018 CheckThat! lab on automatic identification and verification of political claims. Task 2: Factuality. In: Cappellato et al. [21]
18. Bouziane, M., Perrin, H., Cluzeau, A., Mardas, J., Sadeq, A.: Buster.AI at CheckThat! 2020: Insights and recommendations to improve fact-checking. In: Cappellato et al. [19]

19. Cappellato, L., Eickhoff, C., Ferro, N., Névéol, A. (eds.): CLEF 2020 Working Notes. CEUR Workshop Proceedings, CEUR-WS.org (2020)
20. Cappellato, L., Ferro, N., Losada, D., Müller, H. (eds.): Working Notes of CLEF 2019 Conference and Labs of the Evaluation Forum. CEUR Workshop Proceedings, CEUR-WS.org (2019)
21. Cappellato, L., Ferro, N., Nie, J.Y., Soulier, L. (eds.): Working Notes of CLEF 2018–Conference and Labs of the Evaluation Forum. CEUR Workshop Proceedings, CEUR-WS.org (2018)
22. Carik, B., Yeniterzi, R.: SU-NLP at CheckThat! 2021: Check-worthiness of Turkish tweets. In: Faggioli et al. [33]
23. Cazalens, S., Lamarre, P., Leblay, J., Manolescu, I., Tannier, X.: A content management perspective on fact-checking. In: Proceedings of the International Conference on World Wide Web. pp. 565–574. WWW '18 (2018)
24. Cheema, G.S., Hakimov, S., Ewerth, R.: Check_square at CheckThat! 2020: Claim detection in social media via fusion of transformer and syntactic features. In: Cappellato et al. [19]
25. Chernyavskiy, A., Ilvovsky, D., Nakov, P.: Aschern at CLEF CheckThat! 2021: Lambda-calculus of fact-checked claims. In: Faggioli et al. [33]
26. Cusmuliuc, C.G., Amarandei, M.A., Pelin, I., Cociorva, V.I., Iftene, A.: UAICS at CheckThat! 2021: Fake news detection. In: Faggioli et al. [33]
27. Da San Martino, G., Barrón-Cedeno, A., Wachsmuth, H., Petrov, R., Nakov, P.: SemEval-2020 task 11: Detection of propaganda techniques in news articles. In: Proceedings of the 14th Workshop on Semantic Evaluation. pp. 1377–1414. SemEval '20 (2020)
28. Derczynski, L., Bontcheva, K., Liakata, M., Procter, R., Wong Sak Hoi, G., Zubiaga, A.: SemEval-2017 task 8: RumourEval: Determining rumour veracity and support for rumours. In: Proceedings of the 11th International Workshop on Semantic Evaluation. pp. 69–76. SemEval '17 (2017)
29. Dimitrov, D., Ali, B.B., Shaar, S., Alam, F., Silvestri, F., Firooz, H., Nakov, P., Martino, G.D.S.: SemEval-2021 task 6: Detection of persuasion techniques in texts and images. In: Proceedings of the International Workshop on Semantic Evaluation. SemEval '21 (2021)
30. Dumani, L., Neumann, P.J., Schenkel, R.: A framework for argument retrieval - ranking argument clusters by frequency and specificity. In: Advances in Information Retrieval - 42nd European Conference on IR Research (ECIR). Lecture Notes in Computer Science, vol. 12035, pp. 431–445. Springer (2020)
31. Elsayed, T., Nakov, P., Barrón-Cedeño, A., Hasanain, M., Suwaileh, R., Da San Martino, G., Atanasova, P.: CheckThat! at CLEF 2019: Automatic identification and verification of claims. In: Advances in Information Retrieval. pp. 309–315 (2019)
32. Elsayed, T., Nakov, P., Barrón-Cedeño, A., Hasanain, M., Suwaileh, R., Da San Martino, G., Atanasova, P.: Overview of the CLEF-2019 CheckThat!: Automatic identification and verification of claims. In: Experimental IR Meets Multilinguality, Multimodality, and Interaction. pp. 301–321. LNCS (2019)
33. Faggioli, G., Ferro, N., Joly, A., Maistro, M., Piroi, F. (eds.): CLEF 2021 Working Notes. Working Notes of CLEF 2021–Conference and Labs of the Evaluation Forum. CEUR-WS.org (2021)
34. Gencheva, P., Nakov, P., Màrquez, L., Barrón-Cedeño, A., Koychev, I.: A context-aware approach for detecting worth-checking claims in political debates. In: Proceedings of the International Conference Recent Advances in Natural Language Processing. pp. 267–276. RANLP'17 (2017)

35. Ghanem, B., Glavaš, G., Giachanou, A., Ponzetto, S., Rosso, P., Rangel, F.: UPV-UMA at CheckThat! lab: Verifying Arabic claims using cross lingual approach. In: Cappellato et al. [20]
36. Gorrell, G., Kochkina, E., Liakata, M., Aker, A., Zubiaga, A., Bontcheva, K., Derczynski, L.: SemEval-2019 task 7: RumourEval, determining rumour veracity and support for rumours. In: Proceedings of the 13th International Workshop on Semantic Evaluation. pp. 845–854. SemEval '19 (2019)
37. Gupta, A., Kumaraguru, P., Castillo, C., Meier, P.: TweetCred: Real-time credibility assessment of content on Twitter. In: Proceeding of the 6th International Social Informatics Conference. pp. 228–243. SocInfo '14 (2014)
38. Hanselowski, A., PVS, A., Schiller, B., Caspelherr, F., Chaudhuri, D., Meyer, C.M., Gurevych, I.: A retrospective analysis of the fake news challenge stance-detection task. In: Proceedings of the 27th International Conference on Computational Linguistics. pp. 1859–1874. COLING '18 (2018)
39. Hansen, C., Hansen, C., Simonsen, J., Lioma, C.: The Copenhagen team participation in the check-worthiness task of the competition of automatic identification and verification of claims in political debates of the CLEF-2018 fact checking lab. In: Cappellato et al. [21]
40. Hansen, C., Hansen, C., Simonsen, J., Lioma, C.: Neural weakly supervised fact check-worthiness detection with contrastive sampling-based ranking loss. In: Cappellato et al. [20]
41. Hartl, P., Kruschwitz, U.: University of Regensburg at CheckThat! 2021: Exploring text summarization for fake news detection. In: Faggioli et al. [33]
42. Hasanain, M., Elsayed, T.: bigIR at CheckThat! 2020: Multilingual BERT for ranking Arabic tweets by check-worthiness. In: Cappellato et al. [19]
43. Hasanain, M., Haouari, F., Suwaileh, R., Ali, Z., Hamdan, B., Elsayed, T., Barrón-Cedeño, A., Da San Martino, G., Nakov, P.: Overview of CheckThat! 2020 Arabic: Automatic identification and verification of claims in social media. In: Cappellato et al. [19]
44. Hasanain, M., Suwaileh, R., Elsayed, T., Barrón-Cedeño, A., Nakov, P.: Overview of the CLEF-2019 CheckThat! lab on automatic identification and verification of claims. Task 2: Evidence and factuality. In: Cappellato et al. [20]
45. Hassan, N., Li, C., Tremayne, M.: Detecting check-worthy factual claims in presidential debates. In: Proceedings of the 24th ACM International on Conference on Information and Knowledge Management. pp. 1835–1838. CIKM '15 (2015)
46. Hassan, N., Tremayne, M., Arslan, F., Li, C.: Comparing automated factual claim detection against judgments of journalism organizations. In: Computation Journalism Symposium. pp. 1–5 (2016)
47. Hassan, N., Zhang, G., Arslan, F., Caraballo, J., Jimenez, D., Gawsane, S., Hasan, S., Joseph, M., Kulkarni, A., Nayak, A.K., et al.: ClaimBuster: The first-ever end-to-end fact-checking system. Proceedings of the VLDB Endowment **10**(12), 1945–1948 (2017)
48. Álvaro Huertas-Garcia, Huertas-Tato, J., Martín, A., Camacho, D.: CIVIC-UPM at CheckThat! 2021: Integration of transformers in misinformation detection and topic classification. In: Faggioli et al. [33]
49. Juan R. Martinez-Rico, J.M.R., Araujo, L.: NLP&IR@UNED at CheckThat! 2021: Check-worthiness estimation and fake news detection using transformer models. In: Faggioli et al. [33]
50. Kannan, R., R, R.: DLRG@CLEF2021: An ensemble approach for fake detection on news articles. In: Faggioli et al. [33]

51. Karadzhov, G., Nakov, P., Màrquez, L., Barrón-Cedeño, A., Koychev, I.: Fully automated fact checking using external sources. In: Proceedings of the International Conference Recent Advances in Natural Language Processing. pp. 344–353. RANLP’ 17 (2017)
52. Kartal, Y.S., Kutlu, M.: TOBB ETU at CheckThat! 2020: Prioritizing English and Arabic claims based on check-worthiness. In: Cappellato et al. [19]
53. Kartal, Y.S., Kutlu, M.: TrClaim-19: The first collection for Turkish check-worthy claim detection with annotator rationales. In: Proceedings of the 24th Conference on Computational Natural Language Learning. pp. 386–395 (2020)
54. Kazemi, A., Garimella, K., Shahi, G.K., Gaffney, D., Hale, S.A.: Tiplines to combat misinformation on encrypted platforms: A case study of the 2019 Indian election on WhatsApp. arXiv:2106.04726 (2021)
55. Kovachevich, N.: BERT fine-tuning approach to CLEF CheckThat! fake news detection. In: Faggioli et al. [33]
56. Kumari, S.: NoFake at CheckThat! 2021: Fake news detection using BERT. In: Faggioli et al. [33]
57. L, H.R., M, A.: NITK_NLP at CLEF CheckThat! 2021: Ensemble transformer model for fake news classification. In: Faggioli et al. [33]
58. Ma, J., Gao, W., Mitra, P., Kwon, S., Jansen, B.J., Wong, K.F., Cha, M.: Detecting rumors from microblogs with recurrent neural networks. In: Proceedings of the International Joint Conference on Artificial Intelligence. pp. 3818–3824. IJCAI ’16 (2016)
59. Martinez-Rico, J., Araujo, L., Martinez-Romo, J.: NLP&IR@UNED at CheckThat! 2020: A preliminary approach for check-worthiness and claim retrieval tasks using neural networks and graphs. In: Cappellato et al. [19]
60. Mihaylova, S., Borisova, I., Chemishanov, D., Hadzhitsanev, P., Hardalov, M., Nakov, P.: DIPS at CheckThat! 2021: Verified claim retrieval. In: Faggioli et al. [33]
61. Mihaylova, T., Karadzhov, G., Atanasova, P., Baly, R., Mohtarami, M., Nakov, P.: SemEval-2019 task 8: Fact checking in community question answering forums. In: Proceedings of the 13th International Workshop on Semantic Evaluation. pp. 860–869. SemEval ’19 (2019)
62. Mitra, T., Gilbert, E.: CREDBANK: A large-scale social media corpus with associated credibility annotations. In: Proceedings of the Ninth International AAAI Conference on Web and Social Media. pp. 258–267. ICWSM ’15 (2015)
63. Mohammad, S., Kiritchenko, S., Sobhani, P., Zhu, X., Cherry, C.: SemEval-2016 task 6: Detecting stance in tweets. In: Proceedings of the 10th International Workshop on Semantic Evaluation. pp. 31–41. SemEval ’16 (2016)
64. Mukherjee, S., Weikum, G.: Leveraging joint interactions for credibility analysis in news communities. In: Proceedings of the 24th ACM International Conference on Information and Knowledge Management. pp. 353–362. CIKM’15 (2015)
65. Nakov, P., Barrón-Cedeño, A., Elsayed, T., Suwaileh, R., Màrquez, L., Zaghouani, W., Gencheva, P., Kyuchukov, S., Da San Martino, G.: Overview of the CLEF-2018 lab on automatic identification and verification of claims in political debates. In: Working Notes of CLEF 2018 – Conference and Labs of the Evaluation Forum. CLEF ’18 (2018)
66. Nakov, P., Corney, D., Hasanain, M., Alam, F., Elsayed, T., Barrón-Cedeño, A., Papotti, P., Shaar, S., Martino, G.D.S.: Automated fact-checking for assisting human fact-checkers. In: Proceedings of the 30th International Joint Conference on Artificial Intelligence. IJCAI ’21 (2021)

67. Nakov, P., Màrquez, L., Moschitti, A., Magdy, W., Mubarak, H., Freihat, A.A., Glass, J., Randeree, B.: SemEval-2016 Task 3: Community question answering. In: Proceedings of the 10th International Workshop on Semantic Evaluation. pp. 525–545. SemEval '15 (2016)
68. Nakov, P., Martino, G.D.S., Elsayed, T., Barrón-Cedeño, A., Míguez, R., Shaar, S., Alam, F., Haouari, F., Hasanain, M., Babulkov, N., Nikolov, A., Shahi, G.K., Struß, J.M., Mandl, T.: The CLEF-2021 checkthat! lab on detecting check-worthy claims, previously fact-checked claims, and fake news. In: Hiemstra, D., Moens, M., Mothe, J., Perego, R., Potthast, M., Sebastiani, F. (eds.) Advances in Information Retrieval - 43rd European Conference on IR Research, ECIR 2021, Virtual Event, March 28 - April 1, 2021, Proceedings, Part II. Lecture Notes in Computer Science, vol. 12657, pp. 639–649. Springer (2021). https://doi.org/10.1007/978-3-030-72240-1_75, https://doi.org/10.1007/978-3-030-72240-1_75
69. Nguyen, V.H., Sugiyama, K., Nakov, P., Kan, M.Y.: FANG: Leveraging social context for fake news detection using graph representation. In: Proceedings of the 29th ACM International Conference on Information & Knowledge Management. p. 1165–1174. CIKM '20 (2020)
70. Nikolov, A., Da San Martino, G., Koychev, I., Nakov, P.: Team_Alex at Check-That! 2020: Identifying check-worthy tweets with transformer models. In: Cappellato et al. [19]
71. Oshikawa, R., Qian, J., Wang, W.Y.: A survey on natural language processing for fake news detection. In: Proceedings of the 12th Language Resources and Evaluation Conference. pp. 6086–6093. LREC '20 (2020)
72. Pogorelov, K., Schroeder, D.T., Burchard, L., Moe, J., Brenner, S., Filkukova, P., Langguth, J.: FakeNews: Corona virus and 5G conspiracy task at MediaEval 2020. In: Proceedings of the MediaEval workshop. MediaEval '20 (2020)
73. Popat, K., Mukherjee, S., Strötgen, J., Weikum, G.: Credibility assessment of textual claims on the web. In: Proceedings of the 25th ACM International Conference on Information and Knowledge Management. pp. 2173–2178. CIKM '16 (2016)
74. Pritzkau, A.: NLytics at CheckThat! 2021: Check-worthiness estimation as a regression problem on transformers. In: Faggioli et al. [33]
75. Pritzkau, A.: NLytics at CheckThat! 2021: Multi-class fake news detection of news articles and domain identification with RoBERTa - a baseline model. In: Faggioli et al. [33]
76. Sardar, A.A.M., Salma, S.A., Islam, M.S., Hasan, M.A., Bhuiyan, T.: Team Sigmoid at CheckThat! 2021: Multiclass fake news detection with machine learning. In: Faggioli et al. [33]
77. Sepúlveda-Torres, R., Saquete, E.: GPLSI team at CLEF CheckThat! 2021: Fine-tuning BETO and RoBERTa. In: Faggioli et al. [33]
78. Shaar, S., Alam, F., Martino, G.D.S., Nakov, P.: The role of context in detecting previously fact-checked claims. arXiv preprint arXiv:2104.07423 (2021)
79. Shaar, S., Babulkov, N., Da San Martino, G., Nakov, P.: That is a known lie: Detecting previously fact-checked claims. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics. pp. 3607–3618. ACL '20 (2020)
80. Shaar, S., Haouari, F., Mansour, W., Hasanain, M., Babulkov, N., Alam, F., Da San Martino, G., Elsayed, T., Nakov, P.: Overview of the CLEF-2021 Check-That! lab task 2 on detect previously fact-checked claims in tweets and political debates. In: Faggioli et al. [33]

81. Shaar, S., Hasanain, M., Hamdan, B., Ali, Z.S., Haouari, F., Nikolov, A., Kutlu, M., Kartal, Y.S., Alam, F., Da San Martino, G., Barrón-Cedeño, A., Míguez, R., Beltrán, J., Elsayed, T., Nakov, P.: Overview of the CLEF-2021 CheckThat! lab task 1 on check-worthiness estimation in tweets and political debates. In: Faggioli et al. [33]
82. Shaar, S., Nikolov, A., Babulkov, N., Alam, F., Barrón-Cedeño, A., Elsayed, T., Hasanain, M., Suwaileh, R., Haouari, F., Da San Martino, G., Nakov, P.: Overview of CheckThat! 2020 English: Automatic identification and verification of claims in social media. In: Cappellato et al. [19]
83. Shahi, G.K.: AMUSED: An annotation framework of multi-modal social media data. arXiv:2010.00502 (2020)
84. Shahi, G.K., Dirkson, A., Majchrzak, T.A.: An exploratory study of COVID-19 misinformation on Twitter. *Online social networks and media* pp. 100–104 (2021)
85. Shahi, G.K., Dirkson, A., Majchrzak, T.A.: Exploring the spread of COVID-19 misinformation on Twitter. arXiv preprint (2021)
86. Shahi, G.K., Nandini, D.: FakeCovid – a multilingual cross-domain fact check news dataset for COVID-19. In: *Workshop Proceedings of the 14th International AAAI Conference on Web and Social Media* (2020)
87. Shahi, G.K., Struß, J.M., Mandl, T.: CT-FAN-21 corpus: A dataset for Fake News Detection (Apr 2021). <https://doi.org/10.5281/zenodo.4714517>, <https://doi.org/10.5281/zenodo.4714517>
88. Shahi, G.K., Struß, J.M., Mandl, T.: Overview of the CLEF-2021 CheckThat! lab: Task 3 on fake news detection. In: Faggioli et al. [33]
89. Shu, K., Sliva, A., Wang, S., Tang, J., Liu, H.: Fake news detection on social media: A data mining perspective. *SIGKDD Explor. Newsl.* **19**(1), 22–36 (2017)
90. Skuczyńska, B., Shaar, S., Spenader, J., Nakov, P.: BeaSku at CheckThat! 2021: Fine-tuning sentence BERT with triplet loss and limited data. In: Faggioli et al. [33]
91. Sohan, S., Rajon, H.S., Khusbu, A., Islam, M.S., Hasan, M.A.: Black Ops at CheckThat! 2021: User profiles analyze of intelligent detection on fake tweets notebook in shared task. In: Faggioli et al. [33]
92. Tchechmedjiev, A., Fafalios, P., Boland, K., Gasquet, M., Zloch, M., Zapilko, B., Dietze, S., Todorov, K.: ClaimsKG: A knowledge graph of fact-checked claims. In: *Proceedings of the International Semantic Web Conference*. pp. 309–324. ISWC '19, Springer (2019)
93. Thorne, J., Vlachos, A., Christodoulopoulos, C., Mittal, A.: FEVER: a large-scale dataset for fact extraction and VERification. In: *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. pp. 809–819. NAACL-HLT '18 (2018)
94. Touahri, I., Mazroui, A.: EvolutionTeam at CheckThat! 2020: Integration of linguistic and sentimental features in a fake news detection approach. In: Cappellato et al. [19]
95. Tsoplefack, W.K.: Classifier for fake news detection and topical domain of news articles. In: Faggioli et al. [33]
96. Utsha, R.S., Keya, M., Hasan, M.A., Islam, M.S.: Qword at CheckThat! 2021: An extreme gradient boosting approach for multiclass fake news detection. In: Faggioli et al. [33]
97. Vasileva, S., Atanasova, P., Márquez, L., Barrón-Cedeño, A., Nakov, P.: It takes nine to smell a rat: Neural multi-task learning for check-worthiness prediction. In: *Proceedings of the International Conference on Recent Advances in Natural Language Processing*. pp. 1229–1239. RANLP '19 (2019)

98. Williams, E., Rodrigues, P., Novak, V.: Accenture at CheckThat! 2020: If you say so: Post-hoc fact-checking of claims using transformer-based models. In: Cappellato et al. [19]
99. Williams, E., Rodrigues, P., Tran, S.: Accenture at CheckThat! 2021: Interesting claim identification and ranking with contextually sensitive lexical training data augmentation. In: Faggioli et al. [33]
100. Zengin, M.S., Kartal, Y.S., Kutlu, M.: TOBB ETU at CheckThat! 2021: Data engineering for detecting check-worthy claims. In: Faggioli et al. [33]
101. Zhao, Z., Resnick, P., Mei, Q.: Enquiring minds: Early detection of rumors in social media from enquiry posts. In: Proceedings of the 24th International Conference on World Wide Web. pp. 1395–1405. WWW'15 (2015)
102. Zhou, X., Wu, B., Fung, P.: Fight for 4230 at CLEF CheckThat! 2021: Domain-specific preprocessing and pretrained model for ranking claims by check-worthiness. In: Faggioli et al. [33]
103. Zubiaga, A., Liakata, M., Procter, R., Hoi, G.W.S., Tolmie, P.: Analysing how people orient to and spread rumours in social media by looking at conversational threads. PLoS ONE **11**(3) (2016)
104. Zuo, C., Karakas, A., Banerjee, R.: A hybrid recognition system for check-worthy claims using heuristics and supervised learning. In: Cappellato et al. [21]

Appendix

A Systems for Task 1

The positions in the task ranking appear after each team name. See Tables 5–7 for further details.

Team Accenture [99] (1A:ar:1 1A:bg:4 1A:en:9 1A:es:5 1A:tr:5) used BERT and RoBERTa with data augmentation. They further generated additional synthetic training data using lexical substitution. To find the most probable substitutions, they used BERT-based contextual embedding to create synthetic examples for the positive class. They further added a mean-pooling layer and a dropout layer on top of the model before the final classification layer.

Team Fight for 4230 [102] (1A:en:2 1B:en:1) focused its efforts mostly on two fronts: the creation of a pre-processing module able to properly normalize the tweets and the augmentation of the data by means of machine translation and WordNet-based substitutions. The pre-processing included link removal and punctuation cleaning, as well as quantities and contractions expansion. All hashtags related to COVID-19 were normalized into one and the hashtags were expanded. Their best approach was based on BERTweet with a dropout layer and the above-mentioned pre-processing.

Team GPLSI [77] (1A:en:5 1A:es:2) applied the RoBERTa and the BERT transformers together with different manually engineered features, such as the occurrence of dates and numbers or words from LIWC. A thorough exploration of parameters was made using weighting and bias techniques. They also tried to split the four-way classification into two binary classifications and one three-way classification. They further tried oversampling and undersampling.

Team iCompass (ar:4) used several preprocessing steps, including (i) English word removal, (ii) removing URLs and mentions, and (iii) data normalization, removing tashkeel and the letter madda from texts, as well as duplicates, and replacing some characters to prevent mixing. They proposed a simple ensemble of two BERT-based models, which include AraBERT and Arabic-ALBERT.

Team NLP&IR@UNED [49] (1A:en:1 1A:es:4) used several transformer models, such as BERT, ALBERT, RoBERTa, DistilBERT, and Funnel-Transformer, for the experiments to compare the performance. For English, they obtained better results using BERT trained with tweets. For Spanish, they used Electra.

Team NLytics [74] (1A:en:8 1B:en:3) used RoBERTa with a regression function in the final layer, approaching the problem as a ranking task.

Team QMUL-SDS [1] (1A:ar:4) used the AraBERT preprocessing function to (i) replace URLs, email addressees, and user mentions with standard words, (ii) removed line breaks, HTML markup, repeated characters, and unwanted characters, such as emotion icons, and (iii) handled white spaces between words and digits (non-Arabic, or English), and/or a combination of both, and before and after two brackets, and also (iv) removed unnecessary punctuation. They addressed the task as a ranking problem, and fine-tuned an Arabic transformer (AraBERTv0.2-base) on a combination of the data from this year and the data from the CheckThat! lab 2020 (the CT20-AR dataset).

Team SCUoL [6] (1A:ar:3) used typical pre-processing steps, including cleaning the text, segmentation, and tokenization. Their experiments consists of fine-tuning different AraBERT models, and their final results were obtained using AraBERTv2-base.

Team SU-NLP [22] (1A:tr:2) also used several pre-processing steps, including (i) removing emojis, hashtags, and (ii) replacing all mentions with a special token (@USER), and all URLs with the respective website’s domain. If the URL is for a tweet, they replaced the URL with TWITTER and the respective user account name. They reported that this URL expansion method improved the performance. Subsequently, they used an ensemble of BERTurk models fine-tuned using different seed values.

Team TOBB ETU [100] (1A:ar:6 1A:bg:5 1A:en:10 1A:es:1 1A:tr:1) investigated different approaches to fine-tune transformer models including data augmentation using machine translation, weak supervision, and cross-lingual training. For their submission, they removed URLs and user mentions from the tweets, and fine-tuned a separate BERT-based models for each language. In particular, they fine-tuned BERTurk¹, AraBERT, BETO², and the BERT-base model for Turkish, Arabic, Spanish, and English, respectively. For Bulgarian, they fine-tune a RoBERTa model pre-trained with Bulgarian documents.³

Team UPV [14] (1A:ar:8 1A:bg:2 1A:en:3 1A:es:6 1A:tr:4) used a multilingual sentence transformer representation (S-BERT) with knowledge distillation, originally intended for question answering. They further introduced an auxiliary language identification task, aside the downstream check-worthiness task.

B Systems for Task 2

Team Aschern [25] (2A:en:1) used TF.IDF, fine-tuned pre-trained S-BERT, and the reranking LambdaMART model.

Team BeaSku [90] (2B:en:3) used triplet loss training to fine-tune S-BERT. Then, they used the scores predicted by the fine-tuned model along with BM25 scores as features to train a rankSVM re-ranker. They further discussed the impact of applying online mining of triplets. They also experimented with data augmentation.

Team DIPS [60] (2A:en:3 2B:en:2) calculated S-BERT embeddings for all claims, then computed a cosine similarity for each pair of an input claim and a verified claim. The prediction is made by passing a sorted list of cosine similarities to a neural network.

Team NLytics (2A:en:2 2B:en:4) approached the problem as a regression task, and used RoBERTa with a regression function in the final layer.

¹ <http://huggingface.co/dbmdz/bert-base-turkish-cased>

² <http://huggingface.co/dccuchile/bert-base-spanish-wwm-cased>

³ <http://huggingface.co/iarfmoose/roberta-base-bulgarian>

C Systems for Task 3

Team Black Ops [91] (3A:11) performed data pre-processing by removing stop-words and punctuation marks. Then, they experimented with decision trees, random forest, and gradient boosting classifiers for Task 3A, and found the latter to perform best.

Team CIC [8] (3A:10 3B:5) experimented with logistic regression, multi-layer perceptron, support vector machines, and random forest. Their experiments consisted of using stratified 5-fold cross-validation on the training data. Their best results were obtained using logistic regression for task 3A, and a multi-layer perceptron for task 3B.

Team CIC 3A:11 experimented with a decision tree, a random forest, and a gradient boosting algorithms. They found the latter to perform best.

Team CIVIC-UPM [48] (3A:7 3B:8) participated in the two subtasks of task 3. They performed pre-processing, using a number of tools: (i) `ftfy` to repair Unicode and emoji errors, (ii) `ekphrasis` to perform lower-casing, normalizing percentages, time, dates, emails, phones, and numbers, (iii) `contractions` for abbreviation expansion, and (iv) NLTK for word tokenization, stop-words removal, punctuation removal and word lemmatization. Then, they combined `doc2vec` with transformer representations (Electra base, T5 small and T5 base, Longformer base, RoBERTa base and DistilRoBERTa base). They further used additional data from Kaggle’s Ag News task, Kaggle’s KDD2020, and Clickbait news detection competitions. Finally, they experimented with a number of classifiers such as Naïve Bayes, Random Forest, Logistic Regression with L1 and L2 regularization, Elastic Net, and SVMs. The best system for subtask 3A used DistilRoBERTa-base on the text body with oversampling and a sliding window for dealing with long texts. Their best system for task 3B used RoBERTa-base on the title+body text with oversampling but no sliding window.

Team DLRG (3A:3 3B:4) experimented with a number of traditional approaches like Random Forest, Naïve Bayes and Logistic Regression as well as an online passive-aggressive classifier and different ensembles thereof. The best result was achieved by an ensemble of Naïve Bayes, Logistic Regression, and the Passive Aggressive classifier for task 3A. For task 3B, the Online Passive-Aggressive classifier outperformed all other approaches, including the considered ensembles.

Team GPLSI [77] (3A:16) applied the RoBERTa transformer together with different manually-engineered features, such as the occurrence of dates and numbers or words from LIWC. Both the title and the body were concatenated as a single sequence of words. Rather than going for a single multi-class setting, they used two binary models considering the most frequent classes: false vs. other, and true vs. other, followed by one three-class model.

Team MUCIC [12] (3A:19 3B:12) used a majority voting ensemble with three BERT variants. They applied BERT, Distilbert, and RoBERTa, and fine-tuned the pre-trained models.

Team NITK_NLP[57] (3A:5 3B:1) proposed an approach, that included pre-processing and tokenization of the news article, and then experimented with multiple transformer models. The final prediction was made by an ensemble.

Team NKovachevich [55] (3A:13 3B:3) created lexical features. They extracted the 500 most frequent word stems in the dataset, and calculated the TF.IDF values, which they used in a multinomial Naïve Bayes classifier. A much better performance was achieved with an LSTM model that used GloVe embeddings. A little lower F1 value was achieved using BERT. They further found RoBERTa to perform worse than BERT.

Team NLP&IR@UNED [49] (3A:4) experimented with four transformer architectures and input sizes of 150 and 200 words. In the preliminary tests, the best performance was achieved by ALBERT with 200 words. They also experimented with combining TF.IDF values from the text, all the features provided by the LIWC tool, and the TF.IDF values from the first 20 domain names returned by a query to a search engine. Unlike what was obtained in the dev dataset, in the official competition, the best results were obtained with the approach based on TF.IDF, LIWC, and domain names.

Team NLytics (3A:12 3B:7) fined-tuned RoBERTa on the dataset for each of the sub-tasks. Since the data is unbalanced, they used under-sampling. They also truncated the documents to 512 words to fit into the RoBERTa input size.

Team NoFake [56] (3A:1 3B:2) applied BERT without fine-tuning, but used an extensive amount of additional data for training, downloaded from various fact-checking websites.

Team Pathfinder [95] (3A:9 3A:10) participated in both tasks and used multinomial Naïve Bayes and random forest. The former performed better for both tasks. For task 3A, they merged the classed *false* and *partially false* into one class, which boosted the model performance by 41% (a non-official score mentioned in the paper).

Team Probity (3A:20) addressed the multiclass fake news detection subtask, they used a simple LSTM architecture where they adopted word2vec embeddings to represent the news articles.

Team Qword [96] (3A:23) applied pre-processing techniques, which included stop-word removal, punctuation removal and lemmatization using a Porter stemmer. The TF.IDF values were calculated for the words. For these features, four classification algorithms were applied. The best result was given by Extreme Gradient Boosting.

Team SAUD (3A:2) used an SVM with TF.IDF. They tried Logistic Regression, Multinomial Naïve Bayes, and Random Forest, and found SVM to work best.

Team Sigmoid [76] (3A:17) experimented with different traditional machine learning approaches, with multinomial Naïve Bayes performing best, and one deep learning approach, namely an LSTM with the Adam optimizer. The latter outperformed the more traditional approaches.

Team Spider (3A:22) applies an LSTM, after a pre-processing consisting of stop-word removal and stemming.

Team UAICS [26] (3A:6) experimented with various models including BERT, LSTM, Bi-LSTM, and feature-based models. Their submitted model is a Gradient Boosting with a weighted combination of three feature groups: bi-grams, POS tags, and lexical categories of words.

Team University of Regensburg [41] (3A:8) used different fine-tuned variants of BERT with a linear layer on top and applied different approaches to address the maximum sequence length of BERT. Besides hierarchical transformer representations, they also experimented with different summarization techniques like extractive and abstractive summarization. They performed oversampling to address the class imbalance, as well as extractive (using DistilBERT) and abstractive summarization (using distil-BART-CNN-12-6), before performing classification using fine-tuned BERT with a hierarchical transformer representation.