

Grand Challenge On Detecting Cheapfakes

Duc-Tien Dang-Nguyen^{1,3}, Sohail Ahmed Khan¹, Cise Midoglu²,
Michael Riegler², Pål Halvorsen², Minh-Son Dao⁴

¹University of Bergen (UiB), ²Simula Metropolitan Center for Digital Engineering (SimulaMet),

³Kristiania University College, ⁴National Institute of Information and Communications Technology (NICT)
{ductien.dangnguyen, sohail.khan}@uib.no, {cise, michael, paalh}@simula.no, dao@nict.go.jp

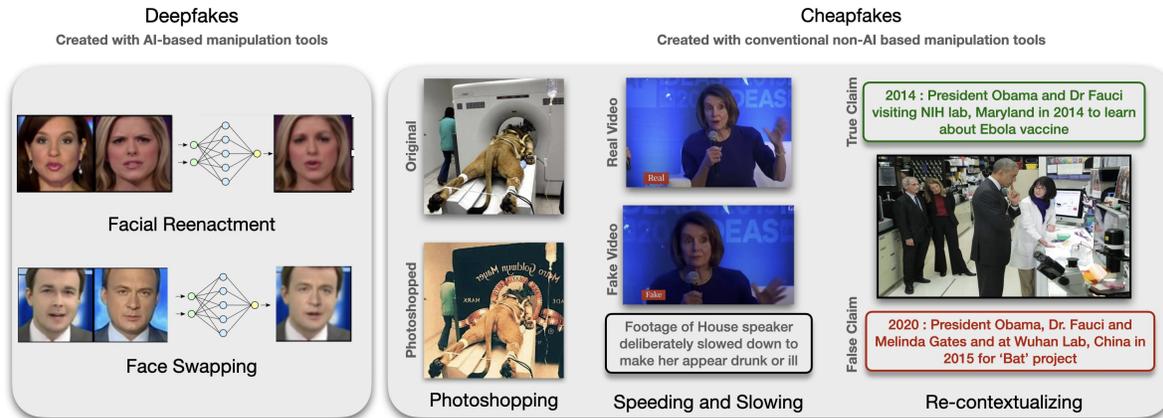


Figure 1: Deepfakes (left) are defined as falsified media created using sophisticated AI-based media manipulation tools and techniques. Cheapfakes (right) include falsified media created with/without contemporary non-AI based editing tools which are easily accessible. Photoshopping tools can be used to tamper with images. Videos can be sped up or slowed down to change the intent or misrepresent the person in the video. Re-contextualizing includes associating falsified or unrelated claims with a genuine image to misrepresent events or persons. This challenge is focused on detecting re-contextualized cheapfakes. Image sources: [1, 2, 3, 4, 5]

1 Motivation and Background

Cheapfake is a recently coined general term that encompasses non-AI (“cheap”) manipulations of multimedia content, created without using deep learning methods. Although a lot of attention has been paid to the creation, detection, and misuse of deepfakes in the last years, cheapfakes are actually known to be more prevalent than deepfakes [6, 7].

Cheapfakes can be created by using contemporary editing tools, which are non-AI based and are easily accessible, such as Adobe Photoshop or PremierePro, or even without using any software. Image manipulations, speeding/slowing of videos, deliberate alteration of the context of the multimedia asset in, e.g., news captions, by sharing the media alongside misleading claims, are some of the methods that are currently in use (see Figure 1). The latter is referred to as context alteration or out-of-context (OOC) misuse of media. OOC media are much harder to detect than fake media, since the images and/or videos are not tampered. An overview of different types of cheapfakes surfacing the Internet are reported by [8].

Depending on the type of cheapfakes, different detection tools can be used. Methods to detect image manipulations such as photoshopping and image splicing have been investigated [9, 10, 11, 12]. Re-contextualization or OOC misuse, which include associating falsified or unrelated claims with a genuine image in order to misrepresent events or persons is, however, relatively niche and unexplored. Aneja et al. [13] have recently introduced this task, provided a dataset of real-world news posts called COSMOS, and proposed a method for detecting cheapfakes, which was benchmarked using the COSMOS dataset.

The aim of this challenge is to develop models that can be used to detect OOC images, and more specifically the misuse of real photographs with conflicting image captions in news items, based on a version of the COSMOS dataset. We have organized similar challenges on the detection of cheapfakes at the ACM Multimedia Systems Conference 2021 (MMSys’21) and the ACM Multimedia Conference 2022 (ACMMM 2022) [14, 15].

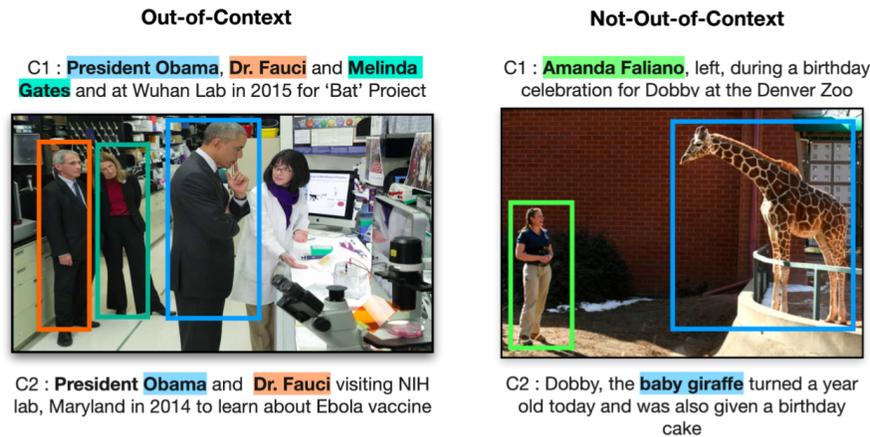


Figure 2: Each image in the COSMOS dataset is accompanied by two captions that the image was circulated together with on the Internet. Left: one of the two captions is misleading due to the alteration of context, indicating out-of-context (OOC) misuse. Right: none of the captions are misleading, hence not-out-of-context (NOOC). Image source: [13]

2 Host Organizations and Coordinator Contacts

This challenge is a collaboration between the University of Bergen (UiB) in Norway, Simula Metropolitan Center for Digital Engineering (SimulaMet) in Norway, and the National Institute of Information and Communications Technology (NICT) in Japan. Organization committee:

- Duc Tien Dang Nguyen, UiB, ductien.dangnguyen@uib.no
- Sohail Ahmed Khan, SFI-MediaFutures - UiB, sohail.khan@uib.no
- Cise Midoglu, SimulaMet, cise@simula.no
- Michael Riegler, SimulaMet, michael@simula.no
- Pål Halvorsen, SimulaMet, paalh@simula.no
- Minh-Son Dao, NICT, dao@nict.go.jp

3 Challenge Description

An image serves as evidence of an event described by a news caption. Presenting an image as evidence of untrue and/or unrelated events is defined as out-of-context (OOC) use of the image (see Figure 2). The aim of this challenge is to develop and benchmark models that can be used to detect OOC misuse of images in news items.

Task 1: Detection of Conflicting Image-Caption Triplets

If two captions associated with an image are valid, then they should describe the same event. If they refer to same object(s) in the image, but are semantically different, i.e., associate the same subject to different events, this indicates OOC use of the image. However, if the captions correspond to the same event, irrespective of the object(s) they describe, this is defined as not-out-of-context (NOOC) use of the image.

In this task, participants are asked to come up with methods to detect conflicting image-caption triplets, which indicate miscontextualization. More specifically, given $\langle \text{Image}, \text{Caption1}, \text{Caption2} \rangle$ triplets as input, their proposed model should predict corresponding class labels 1 (OOC) or 0 (NOOC). The goal is not to identify individual captions as true/false, but rather to detect the existence of miscontextualization. Such methods are considered particularly useful for assisting fact checkers, as highlighting conflicting image-caption triplets allows them to narrow down their search space.

Task 2: Detection of Fake Captions

In a practical scenario, multiple captions might not be available for a given image, and the challenge boils down to figuring out whether an individual caption associated with an image is genuine or not.

In this task, participants are asked to come up with methods to determine whether a given image-caption pair is genuine (real) or falsely generated (fake). More specifically, given an $\langle \text{Image}, \text{Caption} \rangle$ pair as input, their proposed model should predict corresponding class labels 0 (real) or 1 (fake).

We acknowledge that this is a challenging task without prior knowledge of the image origin, even for human moderators. In fact, Luo et al [16] have verified this challenge with a study on human evaluators, who were instructed not to use search engines, where the average human accuracy was only around 65%.

4 Dataset

For this challenge, an augmented version of the COSMOS dataset [13] will be used. A part of this dataset is sampled and assigned as the public dataset. The public dataset, consisting of the training, validation and public test splits, is provided openly to participants for training and testing their algorithms. The remaining part of the COSMOS dataset is augmented with new samples and modified to create the hidden test split, similar to [14, 15]. The hidden test split is not made publicly available, and will be used by the challenge organizers to evaluate the submissions. Details of the challenge dataset are summarized in Table 1.

Table 1: Challenge dataset statistics.

Dataset Split	Number of Images	Number of Captions	Context Annotation
Training	161752	360749	No
Validation	41006	90036	No
Public Test	1000	2000	Yes
Hidden Test	1000	2000	Yes

5 Evaluation Criteria

Participant models will be evaluated and ranked according to two aggregate scores, composed of 5 and 3 metrics respectively.

- *Effectiveness*: accuracy, precision, recall, F1-score, and Matthews correlation coefficient (MCC). Participants are asked to calculate these 5 metrics for their model and include the values in the results section of their submission.
- *Efficiency*: latency, number of parameters, and model size. Participants are asked to calculate these 3 metrics for their model and include the values in the results section of their submission.

After the participants evaluate their own models on the public test split, they are asked to provide code and trained model weights to the organization committee, in order for their models to be evaluated on the hidden test split. Participants are allowed to submit their solutions in three alternative ways as described in Section 7, provided that they abide by the deadlines listed in Section 6.

6 Important Dates

- Dataset release (public training): Monday, 16 January 2023
- Dataset release (public test): Monday, 20 February 2023
- Model submission deadline: Wednesday, 29 March 2023
- Paper submission deadline: Monday, 17 March 2023
- Model evaluation results announcement: Monday, 03 April 2023
- Paper acceptance announcement: Monday, 24 April 2023

7 Submission Guidelines

Docker Container

We recommend challenge participants to submit their solutions as a Docker container, since it will make sure that we don't get any errors resulting from software incompatibility issues or any other similar reason. In this case, we recommend them to follow the instructions given under <https://github.com/detecting-cheapfakes/detecting-cheapfakes-code>.

Standard Python Executable

If the participants face any difficulties in submitting their solutions as a Docker container, or if they feel more comfortable submitting their solution as a standard Python project, they can do so by following the instructions below. It would also

be helpful for us if the participants use PyTorch as the main library if they would like to submit their Python projects, however, this is not compulsory.

- We expect that the submitted code will be executable by a single command, for example:

```
python solution.py <path to folder containing the hidden test split file
private_test.json>
```

- Participants should expect the same format for both the `private_test.json` file and the `public_test.json` file.
- To cope with any software incompatibility issues, we request the participants to provide a `requirements.txt` file along with their solutions, containing the names and specific version numbers of software packages used. This is fairly easy to do with both **Conda** (`conda list`), and **Pip** (`pip list`). We recommend that the participants use **Conda** and create a fresh environment before starting to write the code for their challenge solution.

Jupyter Notebook

Participants can also submit their solutions using Jupyter notebooks, by following the instructions below.

- Participants should structure their code so that it allows us to change the input path by updating a single line in the main file, i.e., the submitted notebook should be executable after changing the `INPUT_FOLDER` parameter from the path for the public test split file, to the path for the private test split file, as shown below:

```
INPUT_FOLDER = <path to folder containing the hidden test split file
private_test.json>
```

- Participants should expect the same format for both the `private_test.json` file and the `public_test.json` file.

8 Additional Information

Experience

The first two editions of this challenge have been organized within, (1) ACM Multimedia Systems Conference 2021 as the “MMSys’21 Grand Challenge on Detecting Cheapfakes” [17, 14], and (2) ACM Multimedia Conference 2022 as the “ACM Multimedia Grand Challenge on Detecting Cheapfakes” [15].

Participant Support

A challenge website has been setup under <https://detecting-cheapfakes.github.io/>, and contains information on datasets, tasks, and other resources. The GitHub organization <https://github.com/detecting-cheapfakes/> hosts all relevant repositories. A Google Group <https://groups.google.com/g/grandchallenge-cheapfakes> has also been established to support prospective participants. Interested participants can find the previously asked questions and join interactive discussions on these platforms.

Publicity

We plan to promote the challenge through the following: (i) professional networks such as MediaFutures, Norwegian Artificial Intelligence Society (NAIS), and the Norwegian AI Community (NORA), (ii) email lists, Slack workspaces, and announcement boards, including VisionList and Image World, (iii) conferences where organizers of this challenge serve as OC members, such as ACM ICMR, ACM MMSys, Global Fact, IEEE CVPR, MediaEval, MMM, NAIS, and NTCIR.

Continuation

The challenge website and other participant resources have been setup with a commitment to be maintained at least for the next 3 years. We would like to reiterate our commitment to sustain the challenge over the next years, as we believe that it is both specific, relevant and timely, as well as generic enough to evolve over time according to different needs. For instance, with the rising importance of synthetic data for research and training purposes, a possible future task could be the *generation* of fake captions.

References

- [1] A. Rössler, D. Cozzolino, L. Verdoliva, C. Riess, J. Thies, and M. Nießner. FaceForensics++: Learning to detect manipulated facial images, 2019.
- [2] Boredpanda. 30 fake viral photos people believed were real, 2019.
- [3] J. Waterson. Facebook refuses to delete fake Pelosi video spread by Trump supporters, 2019.
- [4] The White House. President Obama tours a lab at the Vaccine Research Center at the National Institutes of Health, 2014.
- [5] D. Evon. Is this Obama, Fauci, and Gates at a Wuhan Lab in 2015?, 2020.
- [6] J. S. Brennen, F. M. Simon, P. N. Howard, and R. K. Nielsen. Types, sources, and claims of COVID-19 misinformation, 2020.
- [7] N. Schick. Don't underestimate the cheapfake, 2020.
- [8] B. Paris and J. Donovan. Deepfakes and cheapfakes: The manipulation of audio and visual evidence, 2019.
- [9] C. Chen, S. McCloskey, and J. Yu. Image splicing detection via camera response function analysis. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1876–1885, 2017.
- [10] D. Cozzolino, G. Poggi, and L. Verdoliva. Splicebuster: A new blind image splicing detector. *2015 IEEE International Workshop on Information Forensics and Security (WIFS)*, pages 1–6, 2015.
- [11] M. Huh, A. Liu, A. Owens, and A. A. Efros. Fighting fake news: Image splice detection via learned self-consistency, 2018.
- [12] S. Wang, O. Wang, A. Owens, R. Zhang, and A. A. Efros. Detecting photoshopped faces by scripting photoshop. In *ICCV*, 2019.
- [13] Shivangi Aneja, Chris Bregler, and Matthias Nießner. COSMOS: Catching out-of-context misinformation with self-supervised learning, 2021.
- [14] Shivangi Aneja, Cise Midoglu, Duc-Tien Dang-Nguyen, Michael Alexander Riegler, Paal Halvorsen, Matthias Niessner, Balu Adsumilli, and Chris Bregler. MMSys'21 grand challenge on detecting cheapfakes, 2021.
- [15] Shivangi Aneja, Cise Midoglu, Duc-Tien Dang-Nguyen, Sohail Ahmed Khan, Michael Riegler, Pål Halvorsen, Chris Bregler, and Balu Adsumilli. ACM Multimedia grand challenge on detecting cheapfakes, 2022.
- [16] G. Luo, T. Darrell, and A. Rohrbach. NewsCLIPPings: automatic generation of out-of-context multimodal media, 2021.
- [17] C. Midoglu and S. Aneja. Grand challenge on detecting cheapfakes, 2021.