



OPEN ACCESS

EDITED BY

Alessandro Bria,
University of Cassino, Italy

REVIEWED BY

Edgardo Manuel Felipe Riverón,
National Polytechnic Institute (IPN), Mexico
Gábor Szűcs,
Budapest University of Technology and
Economics, Hungary

*CORRESPONDENCE

Róbert Lakatos
✉ lakatos.robert@inf.unideb.hu

RECEIVED 25 October 2023

ACCEPTED 06 February 2024

PUBLISHED 28 February 2024

CITATION

Lakatos R, Pollner P, Hajdu A and Joó T (2024)
A multimodal deep learning architecture for
smoking detection with a small data
approach. *Front. Artif. Intell.* 7:1326050.
doi: 10.3389/frai.2024.1326050

COPYRIGHT

© 2024 Lakatos, Pollner, Hajdu and Joó. This
is an open-access article distributed under the
terms of the [Creative Commons Attribution
License \(CC BY\)](#). The use, distribution or
reproduction in other forums is permitted,
provided the original author(s) and the
copyright owner(s) are credited and that the
original publication in this journal is cited, in
accordance with accepted academic practice.
No use, distribution or reproduction is
permitted which does not comply with these
terms.

A multimodal deep learning architecture for smoking detection with a small data approach

Róbert Lakatos^{1,2,3*}, Péter Pollner⁴, András Hajdu² and Tamás Joó^{3,4}

¹Doctoral School of Informatics, University of Debrecen, Debrecen, Hungary, ²Department of Data Science and Visualization, Faculty of Informatics, University of Debrecen, Debrecen, Hungary, ³Neumann Technology Platform, Neumann Nonprofit Ltd., Budapest, Hungary, ⁴Data-Driven Health Division of National Laboratory for Health Security, Health Services Management Training Centre, Semmelweis University, Budapest, Hungary

Covert tobacco advertisements often raise regulatory measures. This paper presents that artificial intelligence, particularly deep learning, has great potential for detecting hidden advertising and allows unbiased, reproducible, and fair quantification of tobacco-related media content. We propose an integrated text and image processing model based on deep learning, generative methods, and human reinforcement, which can detect smoking cases in both textual and visual formats, even with little available training data. Our model can achieve 74% accuracy for images and 98% for text. Furthermore, our system integrates the possibility of expert intervention in the form of human reinforcement. Using the pre-trained multimodal, image, and text processing models available through deep learning makes it possible to detect smoking in different media even with few training data.

KEYWORDS

AI supported preventive healthcare, pre-training with generative AI, multimodal deep learning, automated assessment of covert advertisement, few-shot learning, smoking detections

1 Introduction

The WHO currently estimates that smoking causes around 8 million deaths a day. It is the leading cause of death from a wide range of diseases, for example, heart attacks, obstructive pulmonary disease, respiratory diseases, and cancers. 15% of people aged 15 years and over smoke in the OECD countries and 17% in the European Union ([Economic Co-operation and Development, 2023](#)). Moreover, of the 8 million daily deaths, 15% result from passive smoking ([World Health Organization, 2022](#)). The studies ([Pechmann and Shih, 1996](#); [Chapman and Davis, 1997](#)) below highlight the influence of smoking portrayal in movies and the effectiveness of health communication models. However, quantifying media influence is complex. For internet media like social sites, precise ad statistics are unavailable. Furthermore, calculating incited and unmarked ads poses a significant difficulty as well. Therefore, accurate knowledge of the smoking-related content appearing in individual services can be an effective tool in reducing the popularity of smoking. Methods for identifying content include continuous monitoring of advertising

intensity (Kong et al., 2022), structured data generated by questionnaires (Fielding et al., 2004), and AI-based solutions that can effectively support these goals. The authors of the article “Machine learning applications in tobacco research” (Fu et al., 2023) point out in their review that artificial intelligence is a powerful tool that can advance tobacco control research and policy-making. Therefore, researchers are encouraged to explore further possibilities.

Nonetheless, these methods are highly data-intensive. In the case of image processing, an excellent example of this is the popular ResNet (He et al., 2016) image processing network, which was trained on the ImageNet dataset (Deng et al., 2009) containing 14,197,122 images. Regarding text processing, we can mention the popular and pioneering BERT network (Devlin et al., 2018b) trained by the Toronto BookCorpus (Zhu et al., 2015) was trained by the 4.5 GB of Toronto BookCorpus. Generative text processing models such as GPT (Radford et al., 2018) are even larger and were trained with significantly more data than BERT. For instance, the training set of GPT 3.0 was the Common Crawl dataset (<https://commoncrawl.org/>), which has a size of 570 GB.

The effective tools for identifying the content of natural language texts are topic modeling (Blei et al., 2003) and the embedding of words (Mikolov et al., 2013; Pennington et al., 2014; Bojanowski et al., 2017), tokens, sentences (Reimers and Gurevych, 2019), or characters (Clark et al., 2022) clustering (Arthur and Vassilvitskii, 2007). For a more precise identification of the content elements of the texts, we can use the named-entity recognition (Ali et al., 2022) techniques. In image processing, we can highlight classification and object detection to detect smoking. The most popular image processing models are VGG (Simonyan and Zisserman, 2014), ResNet (He et al., 2016), Xception (Chollet, 2017), EfficientNet (Tan and Le, 2019), Inception (Szegedy et al., 2016), and YOLO (Redmon et al., 2016). There also are architectures such as CAMFFNet (Lin et al., 2022) that are proposed for detecting smoking-related diseases. The development of multimodal models also is gaining increasing focus (Liu et al., 2019, 2022), which can use texts and images to solve the tasks at the same time. For movies, scene recognition is particularly challenging compared to images (Rao et al., 2020). Scene recognition is also linked to sensitive events such as fire, smoke, or other disaster detection systems (Gagliardi et al., 2021), but there are attempts to investigate point-of-sale and tobacco marketing practices (Bianco et al., 2021) as well.

We concluded that there is currently no publicly available specific smoking-related dataset that would be sufficient to train a complex model from scratch. Hence, we propose a multimodal architecture that uses pre-trained image and language models to detect smoking-related content in text and images. By combining image processing networks with multimodal architectures and language models, we leverage textual and image data simultaneously. This offers a data-efficient and robust solution that can be further improved with expert input. This paper demonstrates the remarkable potential of artificial intelligence, especially deep learning, for the detection of covert advertising, alongside its capacity to provide unbiased, replicable, and equitable quantification of tobacco-related media content.

2 Methods

2.1 Model architecture

As illustrated in Figure 1 by a schematic flow diagram, our solution relies on pre-trained language and image processing models and can handle both textual and image data.

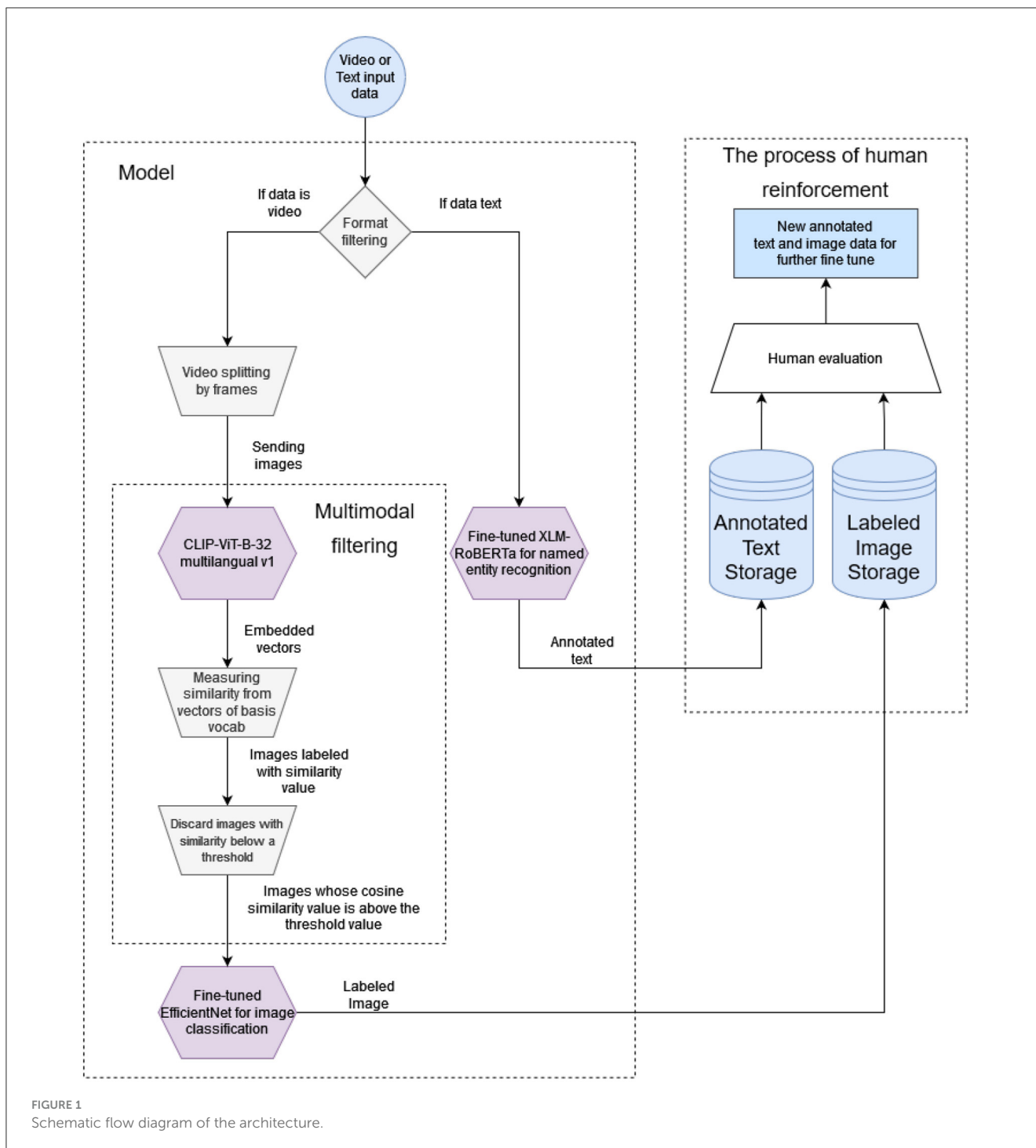
The first step of our pipeline is to define the incoming data format because need to direct the data to the appropriate model for its format. The video recordings are analyzed with multimodal and image processing models, while the texts are analyzed with a large language model. In the case of video recordings, we applied the CLIP-ViT-B-32 multilingual (Reimers and Gurevych, 2020; Radford et al., 2021) model. The model has been developed for over 50 languages with a special training technique (Reimers and Gurevych, 2020). The model supports Hungarian, which was our target language. We use the CLIP-ViT-B-32 model as a filter. After filtering, to achieve more accurate results, we recommend using the pre-trained EfficientNet B5 model, which we fine-tuned with smoking images for the classification task. To process texts, we use name entity recognition to identify smoking-related terms. For this purpose, we have integrated into our architecture an XLM-RoBERTa model (Conneau et al., 2019) that is pre-trained, multilingual, and also supports the Hungarian language, which is important to us.

2.2 Format check

The first step in processing is deciding whether the model has to process video recordings or text data. Since there are many formats for videos and texts, we chose the simple solution of only supporting mp4 and txt file formats. The mp4 is a popular video format, and practically all other video recording formats can be converted to mp4. We consider txt files utf8-encoded raw text files that are ideally free of various metadata. It is important to emphasize that here we ignore the text cleaning processes required to prepare raw text files. The reason is that we did not deal with faulty or txt files requiring further cleaning during the trial.

2.3 Processing of videos and images

The next step in the processing of video footage is to break it down into frames by sampling every second. The ViT image encoder of the CLIP-ViT-B-32 model was trained by its creators for various image sizes. For this, they used the ImageNet (Deng et al., 2009) dataset in which the images have an average size of 469×387 pixels. The developers of CLIP-ViT-B-32 do not recommend an exact size for the image encoder. The model specification only specifies a minimum size of 224×224 . In the case of EfficientNetB5, the developers have optimized an image size of 224×224 . For these reasons, we have taken this image size as a reference and transformed the images sampled from the video recordings to this image size.



2.4 Multimodal filtering

The images sampled from the video recordings were filtered using the CLIP-ViT-B-32 multilingual v1 model. The pre-trained CLIP-ViT-B-32 multilingual v1 model consists of two main components from a ViT (Dosovitskiy et al., 2020) image processing model and a DistilBERT-based (Sanh et al., 2019) multilingual language model. We convert into a 512-long embedded vector (Mikolov et al., 2013) the images and texts with CLIP-ViT-B-32. The embedded vectors for texts and images can be compared

based on their content meaning if we measure cosine similarities between the vectors. The cosine similarity is a value falling in the interval $[-1,1]$, and the similarity of two vectors will be larger the closer their cosine similarity is to 1. Since we aimed to find smoking-related images, we defined a smoking-related term. We converted it to a vector and measured it against the embedded vectors generated from the video images. The term we chose was the word “smoking.” We can use more complex expressions, which could complicate the measurement results interpretation.

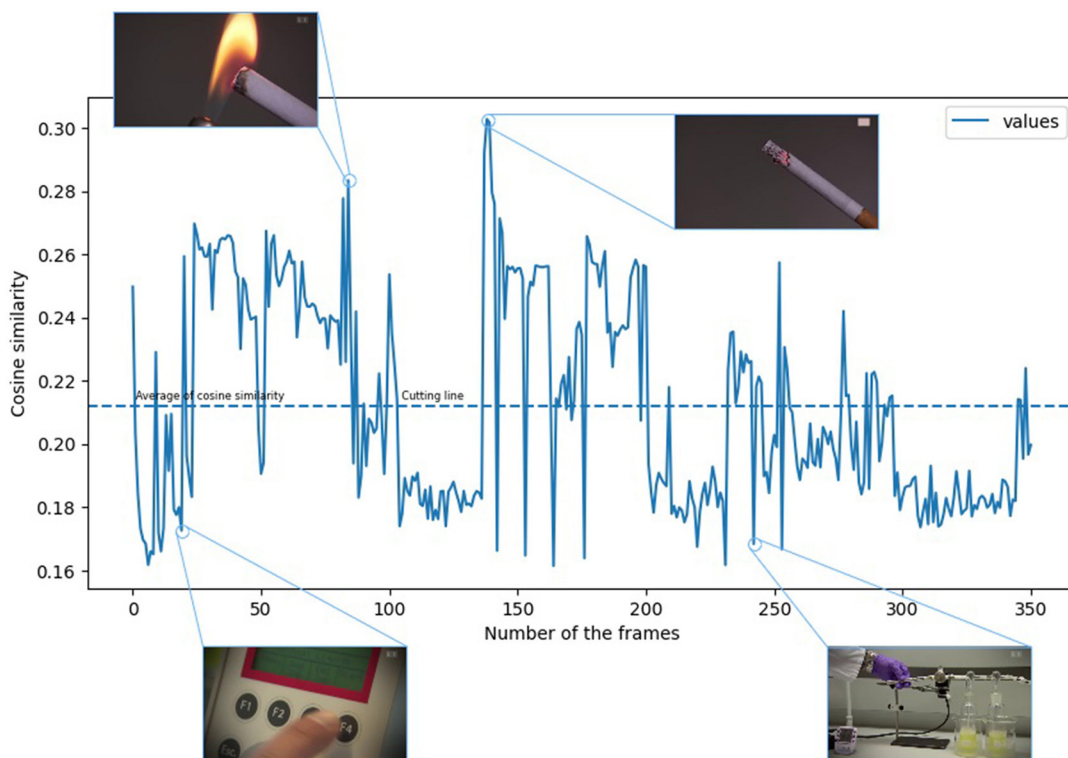


FIGURE 2
The cosine similarity of the images obtained from the video recording in chronological order. (The video footage used for the illustration was selected from the YouTube-8M [Abu-El-Hajja et al., 2016](#) dataset.)

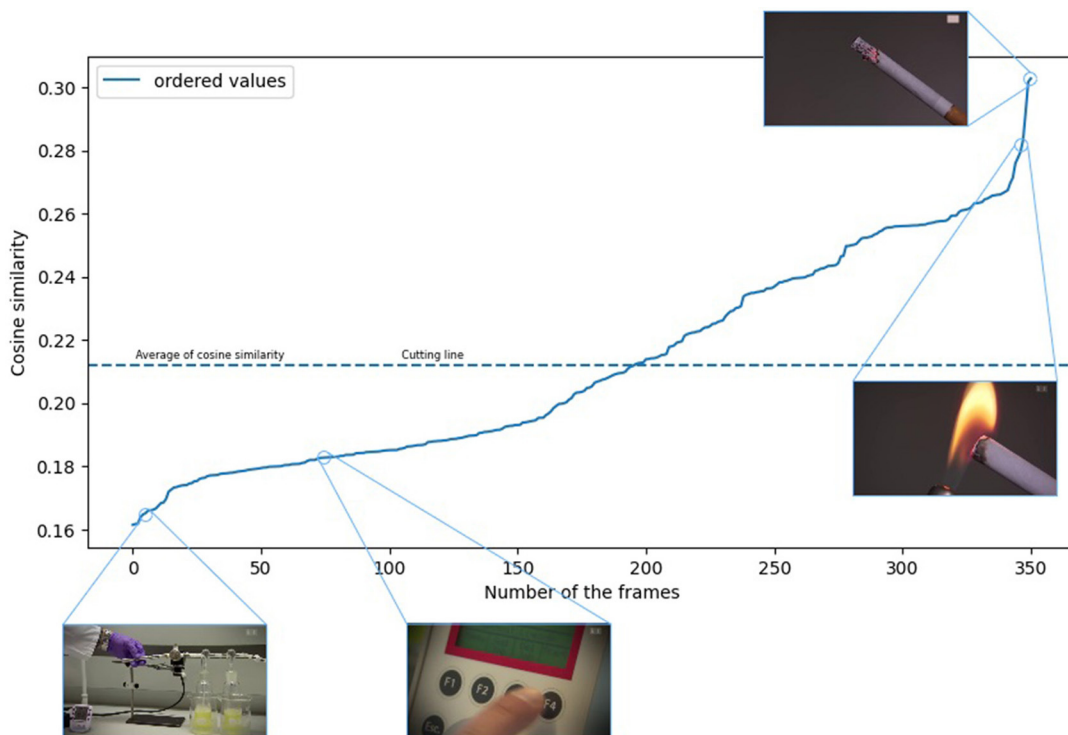


FIGURE 3
The images are in an orderly manner based on the cosine similarity values. (The video footage used for the illustration was selected from the YouTube-8M [Abu-El-Hajja et al., 2016](#) dataset.)

The cosine similarity of the vectors produced by embedding the images always results in a scalar value compared to the vector created from our expression related to “smoking.” However, the decision limit between the distances measured between the vectors produced by the CLIP-ViT-B-32 model is not always clear. Namely, even in the case of images with meanings other than “smoking,” we get a value that is not too distant. We had to understand the distribution of the smoking images to eliminate this kind of blurring of the decision boundary. To this end, we examined the characteristics of the distribution of the images. It is clear from [Figure 2](#) that because the images with a semantic meaning closer to smoking appear randomly in a video recording, it is difficult to grasp the series of images that can be useful for us. [Figure 2](#) is actually a function whose vertical axis has the cosine similarity values belonging to the individual images.

At the same time, the horizontal axis shows the position of the images in the video. To solve this problem, we introduced the following procedure. If we put the cosine similarity values in ascending order, we get a function that describes the ordered evolution of the cosine similarity values.

The ordered function generated from [Figure 2](#) can be seen in [Figure 3](#). As shown in [Figures 2, 3](#), we found that if we take the similarity value of the images sampled from the given sample to the word “smoking,” their average results in a cutting line, and we can use it as a filter.

Furthermore, considering the specifics of the video recordings, we consider that the average can be corrected with a constant value. In this mean, the constant value can thus also be defined as the hyperparameter of the model. We chose the 0 default value for the correction constant because of more apparent measurements. Because the choice of the best constant value may differ depending on the recording type and may distort the exact measurement results. We show the complete process of multimodal filtering in [Algorithm 1](#) of [Supplementary material](#).

2.5 Fine-tuned image classification

After filtering the image set with a multimodal model, we applied an image processing model to classify the remaining images further to improve accuracy. Among the publicly available datasets on smoking, we have used the “smoker and non-smoker” ([Khan, 2020](#)) for augmented ([Shorten and Khoshgoftaar, 2019](#)) fine-tuning. We selected the following models for the task. EfficientNet, Inception, ResNet, VGG, and Xception. The EfficientNet B5 version was the best, with an accuracy of 93.75%. [Supplementary Table 1](#) contains our detailed measurement results concerning all models.

2.6 Processing of text

In the case of detecting smoking terms in texts, we approached the problem as a NER task and focused on the Hungarian language. Since we could not find a train and validation dataset containing annotated smoking phrases available in Hungarian, therefore, to generate the annotated data, we used the generational capabilities of ChatGPT, the smoking-related words of the Hungarian synonyms

TABLE 1 Examples of Hungarian synonyms for smoking in English.

Cigarette	King deck
Smokes	Cigar
Pipe	Coffin nail
Spangles	Like the factory chimney
Tobacco	Cigarette end

TABLE 2 A three elements example prompt for ChatGPT.

Generate a short text about smoking.
The text strictly contains the following words in the different sentences: smoking, tobacco, and cigar.

and antonyms dictionary ([Viola, 2012](#)), and prompt engineering. Accordingly, we selected words related to smoking from the synonyms and antonyms dictionary and asked ChatGPT to suggest further smoking-related terms besides words from the Hungarian dictionary.

Finally, we combined the synonyms and the expressions generated by ChatGPT into a single dictionary. A simplified English sample from the dictionary can be viewed in [Table 1](#), and a complete Hungarian dictionary is contained in [Supplementary Table 5](#).

We created blocks of a maximum of five elements from the words in our dictionary. Each block contained a random combination of a maximum of five words. The blocks are disjoint, so they do not contain the same words. This mixing step was done 10 times. This means that, in one iteration, we could form eight blocks of 5-element disjunct random blocks from our 43-word dictionary. By doing all these 10 times, we produced 80 blocks. However, due to the 10 repetitions, the 80 blocks were no longer disjoint. In other words, if we string all the blocks together, we get a dictionary in which every synonym for smoking appears a maximum of 10 times.

We made a prompt template to which, by attaching each block, we instructed ChatGPT to generate texts containing the specified expressions. Since ChatGPT uses the Hungarian language well, the generated texts contained our selected words by the rules of the Hungarian language, with the correct conjugation. An example of our prompts is illustrated in [Table 2](#).

We did not specify how long texts should be generated by ChatGPT or that every word of a 5-element block should be included in the generated text. When we experimented with ChatGPT generating fixed-length texts, it failed. Therefore, we have removed the requirement for this. Using this method, we created a smoking-related corpus consisting of 80 paragraphs, 49,000 characters, and 7,160 words. An English example of a generated text is presented in [Table 3](#), and there are more Hungarian examples in [Section 2](#) of [Supplementary Tables 6–16](#).

To find the best model according to the possibilities of our computing environment and the support of the Hungarian language, we tested the following models: XLM RoBERTa base and large, DistilBERT base cased, huBERT base ([Nemeskey, 2021](#)), BERT base multilingual ([Devlin et al., 2018a](#)), and Sentence-BERT

TABLE 3 An example paragraph generated by from the prompt of Table 2.

Smoking is a widespread and addictive habit that involves inhaling and exhaling the smoke produced by burning tobacco. Whether it's a hand-rolled cigar or a manufactured cigarette, the act of smoking revolves around the consumption of tobacco. Despite the well-known health risks, many individuals continue to engage in smoking due to its addictive nature. The allure of a cigar or a cigarette can be strong, making it challenging for people to quit smoking even when they are aware of its detrimental effects. Education and support are crucial in helping individuals break free from the cycle of smoking and its associated harms.

(Reimers and Gurevych, 2019). The best model was the XLM RoBERTa large one, which achieved 98% accuracy and 96% F1-score on the validation dataset and an F1-score of 91% with an accuracy of 98% on the test dataset.

2.7 Human reinforcement

In the architecture we have outlined, the last step in dealing with the lack of data is to ensure the system's continuous development capability. For this, we have integrated human confirmation into our pipeline. The essence is that our system's hyperparameter should be adjustable and optimizable during operation and that the data generated during detection can be fed back for further fine-tuning.

The hyperparameter of our solution is the cut line used in multimodal filtering. Its value is a default value. Therefore, its value is not immutable. After the expert has reviewed the results of the data generated during the processing process, the hyperparameter can be modified. Which can optimize the performance of the model.

The tagged images and annotated texts from the processed video recordings and texts are transferred to permanent storage in the last step of the process. This dynamically growing dataset can be further validated with additional human support, and possible errors can be filtered. So, False positives and False negatives can be fed back into the training datasets.

In our architecture, we consider both the modifiability of the hyperparameter and the collection, verifiability, and feedback of the processed data in the training process as tools that provide the possibility of human reinforcement in order to further increase performance.

3 Results

We collected video materials to test the image processing part of our architecture. The source of the video materials was the video-sharing site YouTube. Taking into account the legal rules

regarding the usability of YouTube videos, we have collected five pieces short advertising films from the Malboro and Philip Morris companies. We ensured not to download videos longer than 2 min because longer videos, such as movies, would have required a special approach and additional pre-processing. Furthermore, we downloaded the videos at 240 p resolution and divided them into frames by sampling every second. Each frame was transformed to a size of 224×224 pixels. We manually annotated all videos. The downloaded videos averaged 64 s and contained an average of 13 s of smoking.

With the multimodal filtering technique, we discarded the images that did not contain smoking. Multimodal filtering found 25 s of smoking on average in the recording. The accuracy of the identified images was 62%. The multimodal filtering could filter out more than half of the 64-s, on average, videos. We also measured the performance of the fine-tuned EfficientNet B5 model by itself. The model detected an average of 28 s of smoking with 60% accuracy. We found that the predictions of the two constructions were sufficiently diverse to connect them using the boosting ensemble (Dietterich, 2000) solution. By connecting the two models, the average duration of perceived smoking became 12 s with 4 s on average error and 74% accuracy. The ensemble solution was the best approach since the original videos contained an average of 13 s of smoking. We deleted the videos after the measurements and did not use them anywhere for any other purpose. Supplementary Table 2 contains the exact quantitative results broken down into test videos.

We created training and validation datasets from Hungarian synonyms for smoking using ChatGPT. Samples of generated data are provided in Section 2 of the Supplementary material. From the data, we created a learning and validation data set in proportions of 80 and 20%. We trained our chosen large language models until their accuracy on the validation dataset did not increase for at least 10 epochs. The XLM-RoBERTa model achieved the best performance on the validation dataset with an F1-score of 96 and 98% accuracy. For the final measurement, we created test data from an online text related to smoking by manual annotation (Health Promotion Center, 2023). The text of the entire test data is included in Section 3 in the Supplementary material. The fine-tuned XLM-RoBERTa model achieved 98% accuracy and 0.91 F1 score on the test dataset. The measurement results of the chosen models can be viewed in more detail in Supplementary Tables 3, 4, where, in addition to the accuracy and F1-score values, the recall, precision, and cross-entropy loss values can also be found.

4 Conclusions

Multimodal and image classification models are powerful for classification tasks. In return, however, they are complex and require substantial training data, which can reduce their explainability and usability. In turn, our solution showed that pre-trained multimodal and image classification models exist that allow smoking detection even with limited data and in the matter of low-resource languages if we use the potential of human reinforcement, generative, and ensemble methods. In addition, we see further development opportunities if our approach is supplemented with an object detector, which can determine the time of occurrence of objects and their position. Moreover, with

the expected optimization of the automatic generation of images in the future and the growth of the available computing power, our method used for texts can work in the case of images.

Data availability statement

The original contributions presented in the study are included in the article/[Supplementary material](#), further inquiries can be directed to the corresponding author.

Ethics statement

Written informed consent was not obtained from the individual(s) for the publication of any potentially identifiable images or data included in this article because the pictorial illustrations we use come from video YouTube recordings that are publicly available to anyone.

Author contributions

RL: Writing - original draft. PP: Supervision, Writing - review & editing. AH: Supervision, Writing - review & editing. TJ: Supervision, Writing - review & editing.

Funding

The author(s) declare financial support was received for the research, authorship, and/or publication of this article. The project no. KDP-2021 has been implemented with the support provided by the Ministry of Culture and Innovation of Hungary from the National Research, Development, and Innovation Fund, financed under the C1774095 funding scheme. Also, this work was partly funded by the project GINOP-2.3.2-15-2016-00005 supported by the European Union, co-financed by the European Social Fund, and by the project TKP2021-NKTA-34, implemented with the support provided by the National Research, Development, and Innovation Fund of Hungary under the TKP2021-NKTA funding

References

- Abu-El-Hajja, S., Kothari, N., Lee, J., Natsev, P., Toderici, G., Varadarajan, B., et al. (2016). Youtube-8m: a large-scale video classification benchmark. *arXiv [Preprint] arXiv:1609.08675*.
- Ali, S., Masood, K., Riaz, A., and Saud, A. (2022). "Named entity recognition using deep learning: a review," in *2022 International Conference on Business Analytics for Technology and Security (ICBATS)* (Dubai: IEEE), 1–7.
- Arthur, D., and Vassilvitskii, S. (2007). "K-Means++: the advantages of careful seeding," in *Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms (New Orleans, Louisiana) (SODA'07)* (Society for Industrial and Applied Mathematics), 1027–1035.
- Bianco, F., Moffett, C., Abunku, P., Chaturvedi, I., Chen, G., Dobler, G., et al. (2021). *Automated Detection of Street-Level Tobacco Advertising Displays*. Authorea. Available online at: <https://www.authorea.com/>
- Blei, D. M., Ng, A. Y., and Jordan, M. I. (2003). Latent dirichlet allocation. *J. Machine Learn. Res.* 3, 993–1022. Available online at: <https://www.jmlr.org/papers/volume3/blei03a/blei03a.pdf?ref=https://githubhelp.com>
- Bojanowski, P., Grave, E., Joulin, A., and Mikolov, T. (2017). Enriching word vectors with subword information. *Trans. Assoc. Comput. Linguist.* 5, 135–146. doi: 10.1162/tacl_a_00051
- Chapman, S., and Davis, R. M. (1997). Smoking in movies: is it a problem? *Tobacco Control* 6, 269–271.
- Chollet, F. (2017). "Xception: deep learning with depthwise separable convolutions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Economic Co-operation and Development*, 1251–1258. Available online at: https://openaccess.thecvf.com/content_cvpr_2017/html/Chollet_Xception_Deep_Learning_CVPR_2017_paper.html
- Clark, J. H., Garrette, D., Turc, I., and Wieting, J. (2022). Canine: pre-training an efficient tokenization-free encoder for language representation. *Trans. Assoc. Comput. Linguist.* 10, 73–91. doi: 10.1162/tacl_a_00448
- Conneau, A., Khandelwal, K., Goyal, N., Chaudhary, V., Wenzek, G., Guzmán, F., et al. (2019). Unsupervised cross-lingual representation learning at scale. *CoRR* 747:abs/1911.02116. doi: 10.18653/v1/2020.acl-main.747

scheme. In addition, the study received further funding from the National Research, Development and Innovation Office of Hungary grant (RRF-2.3.1-21-2022-00006, Data-Driven Health Division of National Laboratory for Health Security).

Acknowledgments

The authors express their gratitude for the high-level computing resources of the ELKH cloud operated by the Computer and Automation Research Institute of the Hungarian Academy of Sciences (SZTAKI), which greatly contributed to the successful implementation of the project.

Conflict of interest

RL and TJ were employed by Neumann Technology Platform, Neumann Nonprofit Ltd.

The remaining authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

Supplementary material

The Supplementary Material for this article can be found online at: <https://www.frontiersin.org/articles/10.3389/frai.2024.1326050/full#supplementary-material>

- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. (2009). "Imagenet: a large-scale hierarchical image database," in *2009 IEEE Conference on Computer Vision and Pattern Recognition* (Miami, FL: IEEE), 248–255.
- Devlin, J., Chang, M., Lee, K., and Toutanova, K. (2018a). BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR* 2018:abs/1810.04805. Available online at: <https://arxiv.org/abs/1810.04805>
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2018b). Bert: pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*. doi: 10.48550/arXiv.1810.04805
- Dietterich, T. G. (2000). "Ensemble methods in machine learning," in *Multiple Classifier Systems: First International Workshop, MCS 2000 Cagliari, Italy, June 21–23, 2000 Proceedings 1* (Berlin: Springer), 1–15.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., et al. (2020). An image is worth 16x16 words: transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*. doi: 10.48550/arXiv.2010.11929
- Economic Co-operation and Development (2023). *Daily Smokers (Indicator)*. doi: 10.1787/1ff488c2-en
- Fielding, R., Chee, Y., Choi, K., Chu, T., Kato, K., Lam, S., et al. (2004). Declines in tobacco brand recognition and ever-smoking rates among young children following restrictions on tobacco advertisements in hong kong. *J. Publ. Health* 26, 24–30. doi: 10.1093/pubmed/fdh118
- Fu, R., Kundu, A., Mitsakakis, N., Elton-Marshall, T., Wang, W., Hill, S., et al. (2023). Machine learning applications in tobacco research: a scoping review. *Tobacco Contr.* 32, 99–109. doi: 10.1136/tobaccocontrol-2020-056438
- Gagliardi, A., de Gioia, F., and Saponara, S. (2021). A real-time video smoke detection algorithm based on kalman filter and cnn. *J. Real-Time Image Process.* 18, 2085–2095. doi: 10.1007/s11554-021-01094-y
- He, K., Zhang, X., Ren, S., and Sun, J. (2016). "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 770–778. Available online at: https://openaccess.thecvf.com/content_cvpr_2016/html/He_Deep_Residual_Learning_CVPR_2016_paper.html
- Health Promotion Center (2023). *Egészség Elvitelre*. Available online at: <https://semmelweis.hu/egeszsegfejlesztes/elvitelre/dohanyzas/>
- Khan, A. (2020). *Dataset containing smoking and not-smoking images (smoker vs. non-smoker)*. *Mendeley Data 1*. Available online at: <https://data.mendeley.com/datasets/7b52hhzs3r/1>
- Kong, G., Schott, A. S., Lee, J., Dashtian, H., and Murthy, D. (2022). Understanding e-cigarette content and promotion on youtube through machine learning. *Tobacco Contr.* 2021:57243. doi: 10.1136/tobaccocontrol-2021-057243
- Lin, J., Chen, Y., Pan, R., Cao, T., Cai, J., Yu, D., et al. (2022). Camffnet: a novel convolutional neural network model for tobacco disease image recognition. *Comput. Electr. Agri.* 202:107390. doi: 10.1016/j.compag.2022.107390
- Liu, Y., Guo, Y., Liu, L., Bakker, E. M., and Lew, M. S. (2019). Cyclematch: a cycle-consistent embedding network for image-text matching. *Pat. Recogn.* 93, 365–379. doi: 10.1016/j.patcog.2019.05.008
- Liu, Z., Chen, F., Xu, J., Pei, W., and Lu, G. (2022). Image-text retrieval with cross-modal semantic importance consistency. *IEEE Trans. Circuit. Syst. Video Technol.* 2022:3220297. doi: 10.1109/TCSVT.2022.3220297
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and Dean, J. (2013). "Distributed representations of words and phrases and their compositionality," in *Advances in Neural Information Processing Systems, Vol. 26*, eds C. J. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger (Curran Associates, Inc.). Available online at: https://proceedings.neurips.cc/paper_files/paper/2013/file/9aa42b31882ec039965f3c4923ce901b-Paper.pdf
- Nemeskey, D. M. (2021). "Introducing huBERT," in *XVII. Magyar Számítógépes Nyelvészeti Konferencia (MSZNY2021)* Szeged.
- Pechmann, C., and Shih, C. (1996). *How Smoking in Movies and Anti-smoking ADS Before Movies May Affect Teenagers' Perceptions of Peers Who Smoke*. Irvine, CA: Graduate School of Management, University of California, Irvine.
- Pennington, J., Socher, R., and Manning, C. D. (2014). "Glove: global vectors for word representation," in *Empirical Methods in Natural Language Processing (EMNLP)*, 1532–1543. Available online at: <http://www.aclweb.org/anthology/D14-1162>
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., et al. (2021). "Learning transferable visual models from natural language supervision," in *International Conference on Machine Learning (PMLR)*, 8748–8763. Available online at: <http://proceedings.mlr.press/v139/radford21a>
- Radford, A., Narasimhan, K., Salimans, T., and Sutskever, I. (2018). *Improving Language Understanding by Generative Pre-training*. Available online at: <https://www.mikecaptain.com/resources/pdf/GPT-1.pdf>
- Rao, A., Xu, L., Xiong, Y., Xu, G., Huang, Q., Zhou, B., et al. (2020). "A local-to-global approach to multi-modal movie scene segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 10146–10155. Available online at: https://openaccess.thecvf.com/content_CVPR_2020/html/Rao_A_Local-to-Global_Approach_to_Multi-Modal_Movie_Scene_Segmentation_CVPR_2020_paper.html
- Redmon, J., Divvala, S., Girshick, R., and Farhadi, A. (2016). "You only look once: unified, real-time object detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 779–788. Available online at: https://www.cv-foundation.org/openaccess/content_cvpr_2016/html/Redmon_You_Only_Look_CVPR_2016_paper.html
- Reimers, N., and Gurevych, I. (2019). Sentence-BERT: sentence embeddings using siamese BERT-networks. *arXiv preprint arXiv:1908.10084*. doi: 10.48550/arXiv.1908.10084
- Reimers, N., and Gurevych, I. (2020). Making monolingual sentence embeddings multilingual using knowledge distillation. *arXiv preprint arXiv:2004.09813*. doi: 10.48550/arXiv.2004.09813
- Sanh, V., Debut, L., Chaumond, J., and Wolf, T. (2019). Distilbert, a distilled version of BERT: smaller, faster, cheaper and lighter. *ArXiv, abs/1910.01108*. doi: 10.48550/arXiv.1910.01108
- Shorten, C., and Khoshgoftaar, T. M. (2019). A survey on image data augmentation for deep learning. *J. Big Data* 6, 1–48. doi: 10.1186/s40537-019-0197-0
- Simonyan, K., and Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*. doi: 10.48550/arXiv.1409.1556
- Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., and Wojna, Z. (2016). "Rethinking the inception architecture for computer vision," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2818–2826. Available online at: https://www.cv-foundation.org/openaccess/content_cvpr_2016/html/Szegedy_Rethinking_the_Inception_CVPR_2016_paper.html
- Tan, M., and Le, Q. (2019). "Efficientnet: rethinking model scaling for convolutional neural networks," in *International Conference on Machine Learning (PMLR)*, 6105–6114. Available online at: <http://proceedings.mlr.press/v97/tan19a.html?ref=jina-ai-gmbh.ghost.io>
- Viola, T. (2012). *Ellentétes jelentésű szavak adatbázisa*. Hungary: Tinta Könyvkiadó. Available online at: https://baranyilaszozsolt.com/pciskola/TAMOP-4_2_5-09_Ellentetes_jelentesu_szavak_adatbazisa.pdf
- World Health Organization. (2022). *Tobacco*. Geneva: World Health Organization.
- Zhu, Y., Kiros, R., Zemel, R., Salakhutdinov, R., Urtasun, R., Torralba, A., et al. (2015). "Aligning books and movies: towards story-like visual explanations by watching movies and reading books," in *Proceedings of the IEEE International Conference on Computer Vision*, 19–27. Available online at: https://www.cv-foundation.org/openaccess/content_iccv_2015/html/Zhu_Aligning_Books_and_ICCV_2015_paper.html