

# Are Explainability Tools Gender Biased? A Case Study on Face Presentation Attack Detection

Marco Huber<sup>1,2</sup>, Meiling Fang<sup>1,2</sup>, Fadi Boutros<sup>1</sup>, Naser Damer<sup>1</sup>

<sup>1</sup>Fraunhofer Institute for Computer Graphics Research IGD, Darmstadt, Germany

<sup>2</sup>Department of Computer Science, TU Darmstadt, Darmstadt, Germany

**Abstract**—Face recognition (FR) systems continue to spread in our daily lives with an increasing demand for higher explainability and interpretability of FR systems that are mainly based on deep learning. While bias across demographic groups in FR systems has already been studied, the bias of explainability tools has not yet been investigated. As such tools aim at steering further development and enabling a better understanding of computer vision problems, the possible existence of bias in their outcome can lead to a chain of biased decisions. In this paper, we explore the existence of bias in the outcome of explainability tools by investigating the use case of face presentation attack detection. By utilizing two different explainability tools on models with different levels of bias, we investigate the bias in the outcome of such tools. Our study shows that these tools show clear signs of gender bias in the quality of their explanations.

**Index Terms**—Face Recognition, Bias, Explainability, Face PAD

## I. INTRODUCTION

Face recognition (FR) is increasingly present in our everyday lives, whether it is crossing borders or unlocking our smartphones. Current FR systems [1], [2] achieve outstanding performances that can even exceed those of humans [3], but are difficult for humans to understand and analyse due to the opacity of the deep learning methods used [4]. To increase the understanding of the deep learning models' performance and their behavior in computer vision tasks, several explainability methods have been proposed, such as GradCAM [5] or GradCAM++ [6], to highlight important areas for a given task on an image. These methods are gaining increasing attention in the field of biometrics [4], [7], [8]. Explainability tools aim at enhancing trust in biometrics technologies and can also lead to new solutions for challenges facing biometric systems, such as differential performance and bias. Bias refers to relative performance differences towards certain demographic or non-demographic subgroups [9] that might enable unfair behavior of the system or systematic discrimination.

While recent works investigated the demographic bias and the fairness of FR [9], [10] and its related tasks such as face image quality [11], [12] and face presentation attack detection [13], [14], the bias that might be present in the explanations provided by the emerging explainability tools

This research work has been funded by the German Federal Ministry of Education and Research and the Hessian Ministry of Higher Education, Research, Science and the Arts within their joint support of the National Research Center for Applied Cybersecurity ATHENE. This work has been partially funded by the German Federal Ministry of Education and Research through the Software Campus Project.

used to increase the interpretability of the used models has not been investigated so far. Different explainability tools have already been used in biometric systems [8], [15], [16], mainly focusing on visual explanations. However, none of these works investigated the bias that these explanations may contain and present to the user without further notice or discussion of the possible present bias. As the explainability outcome is used to direct the design choices of algorithm developers as well as the decisions of system operators, the existence of bias in these explanations might lead to a chain of biased decisions.

In this work, we elaborate on the question "Are explainability tools gender-biased?" by performing a case study on face presentation attack detection (face PAD). Face PAD refers to an attack, in which an attacker attempts to impersonate another identity while using an FR system by, for example, using a video, a print, or a 3D mask. This task was chosen as it represents a simple binary classification task, which eliminates the influence of higher complexities. We investigate the bias in explainability outcomes of face PAD in terms of gender bias as it is one of the most well-known and discussed biases. For the case study, we utilize an arbitrary face presentation attack detector [13] and two different, widely used explainability tools, GradCAM [5] and GradCAM++ [6] and evaluate their explanations based on the presented gender using a deletion-and-insertion evaluation scheme [17]. This is therefore the first work to investigate demographic bias in the outcomes of explainability tools.

## II. OUR CASE STUDY

In this case study, we investigate the question if explainability tools provide explanations that are gender biased based on face PAD systems. These face PAD solutions are commonly trained to solve the PAD problem using binary classification between bona fide and presentations attack images. Recent works in the literature have demonstrated that state-of-the-art PADs are gender biased [13]. In our experiments, we utilize three different models, one trained on a balanced gender dataset, and two models solely trained on male or female data. We then apply two explainability tools to generate visual explanations of the models' decisions. These explanations are then evaluated statistically based on insertion and deletion evaluation curves [17]. By inserting and deleting the important pixels as identified by the explainability tool, we quantify the explainability performance - and performance differences

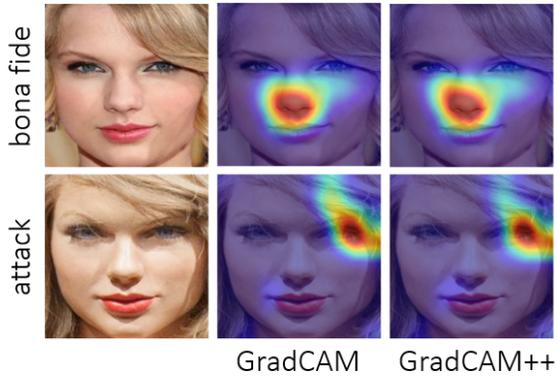


Fig. 1. Example of explanation maps generated using GradCAM and GradCAM++ on a bona fide and an attack image. Both methods highlight a similar area in the image.

based on gender, which would indicate gender bias presented in the explanations.

#### A. Experimental Setup

For the experiments we utilize three different models  $PAD_B$ ,  $PAD_F$ , and  $PAD_M$ . All models share the same ResNet-50 [18] architecture. We chose this architecture as it serves as a backbone in many state-of-the-art PAD methods [19], [20] and achieved good PAD performance [13]. The  $PAD_B$  model is trained on the training set of the CAAD-PAD dataset [13] on images of both, males and females. We test on the testing split of the CAAD-PAD [13] dataset, which provides a testing dataset consisting only of female or male images. The testing data consists of 53.827 male and 19.042 female images. The  $PAD_M$  and the  $PAD_F$  models are trained only on the male and female training sets of the CAAD-PAD dataset, respectively. We follow the implementation details provided in [13]. Using models trained on different gender-based subsets allows us to investigate explanations obtained from models with different levels of bias. The used  $PAD_B$  model achieved an Equal Error Rate (EER) of 2.54% on the male test set and 3.00% on the female test set. The  $PAD_F$  and the  $PAD_M$  achieved an EER of 2.96% and 13.13% on the male set and an EER of 5.90% and 10.62% on the female test set, respectively, and thus clearly show bias as discussed in [13].

As explainability tools, we use GradCAM [5] and GradCAM++ [6] in our experiments. Both approaches produce saliency maps that highlight the important regions in an image for the predicted value. Examples of explanation maps produced by GradCAM and GradCAM++ are provided in Figure 1.

#### B. Evaluation Metrics

To measure the bias of the explanations produced by the explainability tools, we utilize an insertion and deletion curve evaluation, following the trend in evaluating the performance

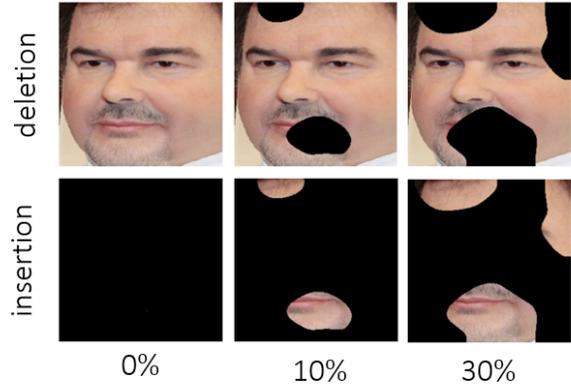


Fig. 2. Visualisation of the insertion and deletion. Based on the calculated explanation map a fraction of the most important pixels are either removed or inserted.

of explainability outcomes [17], [21]. In the insertion evaluation, we iteratively insert pixels from the input image into a black canvas. The pixels, in this case, are selected based on the importance scores produced by the explainability tool. In the deletion evaluation, we iteratively delete pixels from the input image based on the calculated explanation map by setting their values to zero. After inserting or removing a certain amount of pixels based on their assigned importance to the decision, we evaluate the performance of the models on these newly generated images with the identified important pixels inserted or removed. In our experiments, we limited the amount of removed or added pixels, starting from 5% up to 30% with steps of 5%, as the explanation maps often only indicate a small area as important to the decision. A visualization of the insertion and deletion evaluation procedure is shown in Figure 2.

If the explainability methods do not produce gender bias, a similar performance should be observable, i.e. the accuracy in selecting the most important parts of the image to make the decision is similar for both male and female samples. If this is not the case, the explainability tools provide explanations that are gender biased.

For the evaluation of the PAD system, we report the Half Total Error Rate (HTER) at the fixed threshold of the EER on the unaltered images. The HTER, which is widely used to report PAD performance [22]–[24], is the half of the sum of the Attack Presentation Classification Error Rate (APCER) and the Bona Fide Presentation Classification Error Rate (BPCER) and allows us to report APCER and BPCER in a single curve. We also keep the threshold fixed for the evaluation of the different degrees of insertion or deletion as the threshold is also fixed in practice.

As the different models do not perform similarly on genders, we normalize the insertion and deletion curves with respect to the performance without manipulated images to provide comparable evaluation results. This allows a better visual interpretation and also allows us to calculate and report the

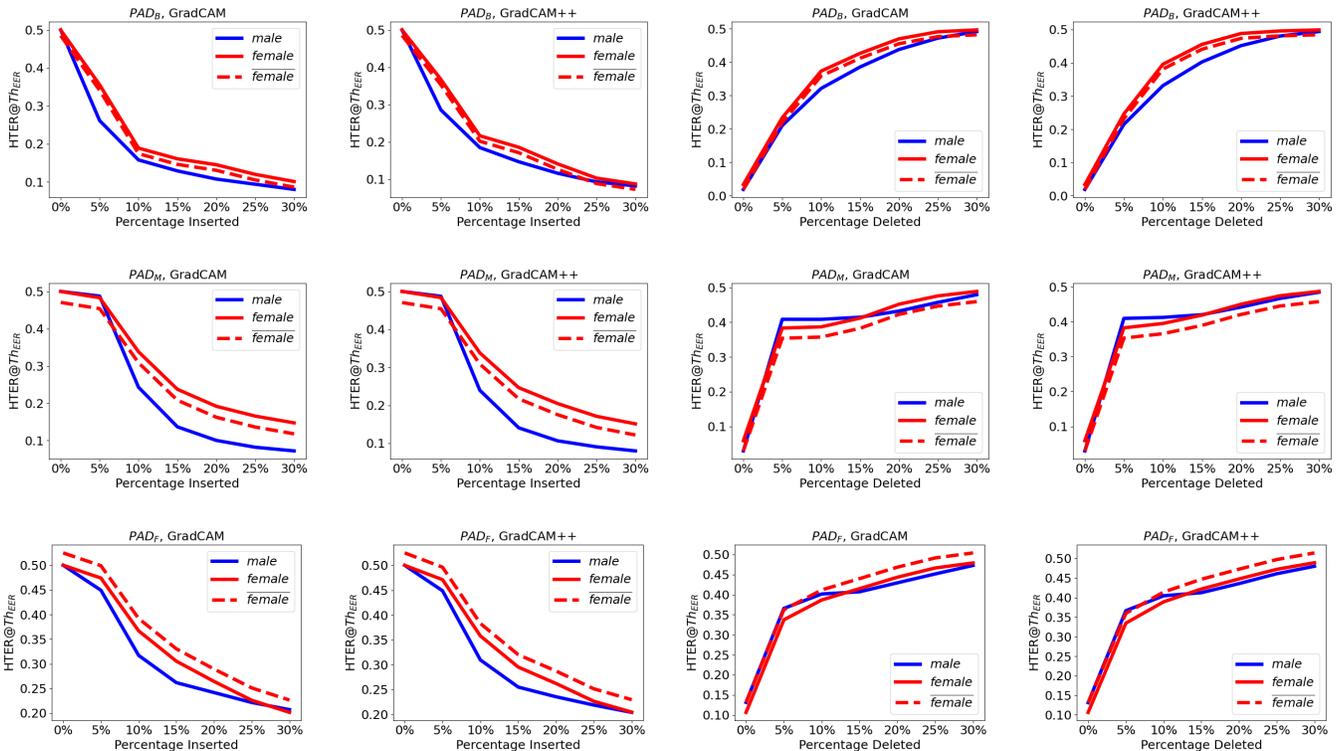


Fig. 3. Insertion and Deletion Curves for the three models,  $PAD_B$  (balanced),  $PAD_M$  (male),  $PAD_F$  (female) using GradCAM [5] and GradCAM++ [6]. The red-dotted line (*female*) shows the normalized female performance with the same starting point as the male performance. Especially in the deletion evaluation curve, the gender bias of the model trained only on one gender ( $PAD_M$  and  $PAD_F$ ) can be observed as there is a performance difference between the error of evaluation of the male (blue) and female (red) dataset.

Area-Under-the-Curve (AUC) as a quantified metric. For the normalization, we calculated the error of the models on the male and female data separately and then subtracted this initial performance difference from all following female error rates to get the normalized female performance. If the explanations of the models are not biased, the difference in the AUC should be zero, as the explanations would indicate with the same performance the most important pixels, independently of the gender of the presented sample.

### III. RESULTS

This section presents the results of our case study on face PAD using the considered explainability tools to investigate if the explanations provided are gender biased. First, we will look at the insertion and deletion curves of the different considered PAD models. Then we will quantify the gender bias in the explanations by comparing the AUC.

The results of the insertion and deletion evaluation are presented in Figure 3. It shows the insertion and deletion curves for all three models ( $PAD_B$ ,  $PAD_M$ , and  $PAD_F$ ) for both explainability methods. The red-dotted line (*female*) in Figure 3 indicates the normalized female performance to compensate different starting performance, as explained in Section II-B. This is needed to provide a fairer explainability comparison, as for example, the performance of the  $PAD_B$  model on the female testing data is worse than on the male

testing data, which should not influence the explanation performance comparison.

The curves in Figure 3 show the HTER at the threshold of the EER (y-axis) over insertion or deletion proportion (x-axis). In the insertion curves, the error is decreased over higher insertion rate and vice versa for the deletion curves, i.e., the error is increased over an increased deletion rate. In the insertion curve, a fast decrease indicates a better performance of the explainability tool, and a low AUC, therefore, indicates superior explainability performance. For the deletion curve, it is the other way around where a fast increase and a higher AUC indicate better performance of the explainability tool. The bias, therefore, is present if a performance gap is observable between explaining samples of different demographic groups. Therefore, the degree of bias can be indicated by the explainability performance difference when processing different groups.

The performance gaps (in term of HTER) over insertion and deletion rates (the outcomes of the explanation tools) between the case when the model is evaluated on male and female subsets, respectively, is observable in the insertion and deletion curves (Figure 3), indicating bias in the explanation outcome.

In the insertion curve, we observe bias as the decrease of the error when evaluated on male data is faster than the decrease in the error evaluated on the female dataset. This remains true even when normalizing the curve depending on

the starting performance (dotted red line). Similar behavior can be observed on the more biased models,  $PAD_M$  and  $PAD_F$ . However, the bias in the explanations was smaller ( $PAD_B$ ) than the bias present in the explanations of the  $PAD_M$  model, at least by using GradCAM, as shown by the closer red and blue curves in Figure 3.

Gender bias in explainability tools can also be observed in the deletion curves. In the  $PAD_B$  model, the increase in error tends to be faster for females than for males. On the deletion curve for the  $PAD_M$ , we can observe a steeper increase in error as pixels are removed on the male data than in comparison to the female data, which indicates a gender bias. The opposite is true for the  $PAD_F$  model, in which the error of the normalized deletion curve for the female testing data increases faster than the error on the male testing data, also indicating bias.

In addition to the insertion and deletion curves, we provide a quantified evaluation of the bias in explainability tools by reporting the AUCs and the performance difference between the male curve and the normalized female curve of each plot in Figure 3.

From the quantified results in Table I, we made the following observations: a) The explainability tools are less biased when evaluated on the less biased  $PAD_B$  model, than when they are evaluated on more biased PAD models ( $PAD_F$  and  $PAD_M$ ), b) in the explanations of the models  $PAD_M$  and  $PAD_F$ , a higher performance difference is observable, indicating larger explainability gender bias. These observations are clearly observable for both explainability tools with slightly higher values for the GradCAM++ method and they are complementary for the ones reported early in this section based on the reported curves.

To conclude, we observed that the explainability tools are gender biased when they are used to explain the behavior of the considered PAD solutions. The bias in the explainability tools is, to some degree, lower for the PAD model that is less biased (trained on data of both genders) than when the explainability tools are evaluated on more biased PAD models (trained on gender-biased datasets). This might be due to bias that is present in the model investigated and inherits its bias to the explainability method. However, the exact roots of the bias in the explainability performance have to be further investigated. Interestingly, we also noticed a link between the bias in the PAD models and the bias in the explanation. In the deletion curves for the  $PAD_M$  and  $PAD_F$ , the bias manifests in the same direction as the models' bias. The performance on the male samples is better than the performance on the female samples using the  $PAD_M$  model that has been trained solely on male images. The same is the case for the  $PAD_F$  with better performance for female images, while it was trained on female images.

#### IV. CONCLUSION

In this work, we investigated the research question: "Are explainability tools gender-biased?", while taking the explanation of face PAD behavior as an example. In our effort

Deletion	Evaluation Data			Insertion	Evaluation Data		
	Male	Female	$\Delta$		Male	Female	$\Delta$
GradCAM				GradCAM			
$PAD_B$	0.104	0.109	<b>0.005</b>	$PAD_B$	0.052	0.059	<b>0.007</b>
$PAD_M$	0.119	0.110	0.009	$PAD_M$	0.067	0.078	0.011
$PAD_F$	0.118	0.125	0.007	$PAD_F$	0.092	0.107	0.015
GradCAM++				GradCAM++			
$PAD_B$	0.107	0.113	<b>0.006</b>	$PAD_B$	0.056	0.061	<b>0.005</b>
$PAD_M$	0.120	0.111	0.009	$PAD_M$	0.068	0.080	0.012
$PAD_F$	0.119	0.126	0.007	$PAD_F$	0.091	0.106	0.015

TABLE I  
AUC FOR THE DIFFERENT PADS, EVALUATION DATA GENDERS, AND EXPLAINABILITY METHODS FOR BOTH CURVES: THE HIGHER DIFFERENCE BETWEEN AUC OF MALE AND FEMALE INDICATES HIGHER BIAS IN THE EXPLAINABILITY PERFORMANCE, THE LOWEST DIFFERENCE (BIAS) IS IN BOLD. AS FOR THE BIAS IN THE PAD PERFORMANCE OF THE  $PAD_B$ , THE BIAS IN ITS EXPLANATION IS LOWER THAN THE OTHER PADS.

to answer this question, we performed a case study on the problem of face PAD by using two explainability tools, GradCAM and GradCAM++, and PAD models with different levels of gender bias. Our investigation concluded that there are differences in the explainability performance when explaining male and female samples and thus there is gender bias in the explainability outcome. As the explainability outcomes are used to increase the transparency for developers and system operators, the existence of bias in these explanations is of concern and needs attention. Future research works could investigate whether other bias factors, such as ethnicity and age, or even non-demographic biases are also affecting explainability outcomes, along with investigating the bias-inducing factors and bias mitigation possibilities.

#### REFERENCES

- [1] F. Boutros, N. Damer, F. Kirchbuchner, and A. Kuijper, "Elasticface: Elastic margin loss for deep face recognition," in *CVPR Workshops*. IEEE, 2022, pp. 1577–1586.
- [2] Q. Meng, S. Zhao, Z. Huang, and F. Zhou, "Magface: A universal representation for face recognition and quality assessment," in *CVPR*. Computer Vision Foundation / IEEE, 2021, pp. 14 225–14 234.
- [3] P. J. Phillips and A. J. O'Toole, "Comparison of human and computer performance across face recognition experiments," *Image Vis. Comput.*, vol. 32, no. 1, pp. 74–85, 2014.
- [4] P. C. Neto, T. Gonçalves, J. R. Pinto, W. Silva, A. F. Sequeira, A. Ross, and J. S. Cardoso, "Explainable biometrics in the age of deep learning," *CoRR*, vol. abs/2208.09500, 2022.
- [5] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-cam: Visual explanations from deep networks via gradient-based localization," in *ICCV*. IEEE Computer Society, 2017, pp. 618–626.
- [6] A. Chattopadhyay, A. Sarkar, P. Howlader, and V. N. Balasubramanian, "Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks," in *WACV*. IEEE Computer Society, 2018, pp. 839–847.
- [7] A. F. Sequeira, T. Gonçalves, W. Silva, J. R. Pinto, and J. S. Cardoso, "An exploratory study of interpretability for face presentation attack detection," *IET Biom.*, vol. 10, no. 4, pp. 441–455, 2021.
- [8] P. C. Neto, A. F. Sequeira, and J. S. Cardoso, "Myope models - are face presentation attack detection models short-sighted?" in *WACV (Workshops)*. IEEE, 2022, pp. 390–399.
- [9] P. Terhörst, J. N. Kolf, M. Huber, F. Kirchbuchner, N. Damer, A. M. Moreno, J. Fierrez, and A. Kuijper, "A comprehensive study on face recognition biases beyond demographics," *IEEE Transactions on Technology and Society*, vol. 3, no. 1, pp. 16–30, 2022.
- [10] T. de Freitas Pereira and S. Marcel, "Fairness in biometrics: A figure of merit to assess biometric verification systems," *IEEE Trans. Biom. Behav. Identity Sci.*, vol. 4, no. 1, pp. 19–29, 2022.

- [11] Z. Babnik and V. Štruc, "Assessing bias in face image quality assessment," in *EUSIPCO*, 2022, pp. 1037–1041.
- [12] P. Terhörst, J. N. Kolf, N. Damer, F. Kirchbuchner, and A. Kuijper, "Face quality estimation and its correlation to demographic and non-demographic bias in face recognition," in *IJCB*. IEEE, 2020, pp. 1–11.
- [13] M. Fang, W. Yang, A. Kuijper, V. Struc, and N. Damer, "Fairness in face presentation attack detection," *CoRR*, vol. abs/2209.09035, 2022.
- [14] N. Alshareef, X. Yuan, K. Roy, and M. Atay, "A study of gender bias in face presentation attack and its mitigation," *Future Internet*, vol. 13, no. 9, p. 234, 2021.
- [15] P. Terhörst, M. Huber, N. Damer, F. Kirchbuchner, K. B. Raja, and A. Kuijper, "Pixel-level face image quality assessment for explainable face recognition," *CoRR*, vol. abs/2110.11001, 2021.
- [16] H. Mirzaalian, M. E. Hussein, L. Spinoulas, J. May, and W. Abd-Almageed, "Explaining face presentation attack detection using natural language," in *FG*. IEEE, 2021, pp. 1–8.
- [17] V. Petsiuk, A. Das, and K. Saenko, "RISE: randomized input sampling for explanation of black-box models," in *BMVC*. BMVA Press, 2018, p. 151.
- [18] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*. IEEE Computer Society, 2016, pp. 770–778.
- [19] Y. Zhang, Z. Yin, Y. Li, G. Yin, J. Yan, J. Shao, and Z. Liu, "Celebaspoo: Large-scale face anti-spoofing dataset with rich annotations," in *ECCV (12)*, ser. Lecture Notes in Computer Science, vol. 12357. Springer, 2020, pp. 70–85.
- [20] Z. Yu, X. Li, J. Shi, Z. Xia, and G. Zhao, "Revisiting pixel-wise supervision for face anti-spoofing," *IEEE Trans. Biom. Behav. Identity Sci.*, vol. 3, no. 3, pp. 285–295, 2021.
- [21] M. Huber, P. Terhörst, F. Kirchbuchner, N. Damer, and A. Kuijper, "On evaluating pixel-level face image quality assessment," in *EUSIPCO*, 2022, pp. 1052–1056.
- [22] M. Fang, N. Damer, F. Kirchbuchner, and A. Kuijper, "Learnable multi-level frequency decomposition and hierarchical attention mechanism for generalized face presentation attack detection," in *WACV*. IEEE, 2022, pp. 1131–1140.
- [23] R. Shao, X. Lan, and P. C. Yuen, "Regularized fine-grained meta face anti-spoofing," in *AAAI*. AAAI Press, 2020, pp. 11 974–11 981.
- [24] M. Fang, M. Huber, and N. Damer, "Synthaspoo: Developing face presentation attack detection based on privacy-friendly synthetic data," *CoRR*, vol. abs/2303.02660, 2023. [Online]. Available: <https://doi.org/10.48550/arXiv.2303.02660>