



HAL
open science

Understanding metric-related pitfalls in image analysis validation

Annika Reinke, Minu D. Tizabi, Michael Baumgartner, Matthias Eisenmann, Doreen Heckmann-Nötzel, A. Emre Kavur, Tim Rädsch, Carole H. Sudre, Laura Acion, Michela Antonelli, et al.

► **To cite this version:**

Annika Reinke, Minu D. Tizabi, Michael Baumgartner, Matthias Eisenmann, Doreen Heckmann-Nötzel, et al.. Understanding metric-related pitfalls in image analysis validation. 2023. hal-04345927

HAL Id: hal-04345927

<https://hal.science/hal-04345927>

Preprint submitted on 14 Dec 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Understanding metric-related pitfalls in image analysis validation

ANNIKA REINKE^{*†}, German Cancer Research Center (DKFZ) Heidelberg, Division of Intelligent Medical Systems and HI Helmholtz Imaging, Germany and Faculty of Mathematics and Computer Science, Heidelberg University, Heidelberg, Germany

MINU D. TIZABI[†], German Cancer Research Center (DKFZ) Heidelberg, Division of Intelligent Medical Systems, Germany and National Center for Tumor Diseases (NCT), NCT Heidelberg, a partnership between DKFZ and University Medical Center Heidelberg, Germany

MICHAEL BAUMGARTNER, German Cancer Research Center (DKFZ) Heidelberg, Division of Medical Image Computing, Germany

MATTHIAS EISENMANN, German Cancer Research Center (DKFZ) Heidelberg, Division of Intelligent Medical Systems, Germany

DOREEN HECKMANN-NÖTZEL, German Cancer Research Center (DKFZ) Heidelberg, Division of Intelligent Medical Systems, Germany and National Center for Tumor Diseases (NCT), NCT Heidelberg, a partnership between DKFZ and University Medical Center Heidelberg, Germany

A. EMRE KAVUR, HI Applied Computer Vision Lab, Division of Medical Image Computing; German Cancer Research Center (DKFZ) Heidelberg, Division of Intelligent Medical Systems, Germany

TIM RÄDSCH, German Cancer Research Center (DKFZ) Heidelberg, Division of Intelligent Medical Systems and HI Helmholtz Imaging, Germany

CAROLE H. SUDRE, MRC Unit for Lifelong Health and Ageing at UCL and Centre for Medical Image Computing, Department of Computer Science, University College London, London, UK and School of Biomedical Engineering and Imaging Science, King's College London, London, UK

LAURA ACION, Instituto de Cálculo, CONICET – Universidad de Buenos Aires, Buenos Aires, Argentina

MICHELA ANTONELLI, School of Biomedical Engineering and Imaging Science, King's College London, London, UK and Centre for Medical Image Computing, University College London, London, UK

TAL ARBEL, Centre for Intelligent Machines and MILA (Quebec Artificial Intelligence Institute), McGill University, Montreal, Canada

SPYRIDON BAKAS, Division of Computational Pathology, Dept of Pathology & Laboratory Medicine, Indiana University School of Medicine, IU Health Information and Translational Sciences Building, Indianapolis, USA and Center for Biomedical Image Computing and Analytics (CBICA), University of Pennsylvania, Richards Medical Research Laboratories FL7, Philadelphia, PA, USA

ARRIEL BENIS, Department of Digital Medical Technologies, Holon Institute of Technology, Holon, Israel and European Federation for Medical Informatics, Le Mont-sur-Lausanne, Switzerland

MATTHEW B. BLASCHKO, Center for Processing Speech and Images, Department of Electrical Engineering, KU Leuven, Leuven, Belgium

FLORIAN BUETTNER, German Cancer Consortium (DKTK), partner site Frankfurt/Mainz, a partnership between DKFZ and UCT Frankfurt-Marburg, Germany, German Cancer Research Center (DKFZ) Heidelberg, Germany, Goethe University Frankfurt, Department of Medicine, Germany, Goethe University Frankfurt, Department of Informatics, Germany, and Frankfurt Cancer Institute, Germany

M. JORGE CARDOSO, School of Biomedical Engineering and Imaging Science, King's College London, London, UK

VERONIKA CHEPLYGINA, Department of Computer Science, IT University of Copenhagen, Copenhagen, Denmark

JIANXU CHEN, Leibniz-Institut für Analytische Wissenschaften – ISAS – e.V., Dortmund, Germany

EVANGELIA CHRISTODOULOU, German Cancer Research Center (DKFZ) Heidelberg, Division of Intelligent Medical Systems, Germany

BETH A. CIMINI, Imaging Platform, Broad Institute of MIT and Harvard, Cambridge, Massachusetts, USA

GARY S. COLLINS, Centre for Statistics in Medicine, University of Oxford, Oxford, UK

KEYVAN FARAHANI, Center for Biomedical Informatics and Information Technology, National Cancer Institute, Bethesda, MD, USA

LUCIANA FERRER, Instituto de Investigación en Ciencias de la Computación (ICC), CONICET-UBA, Ciudad Universitaria, Ciudad Autónoma de Buenos Aires, Argentina

ADRIAN GALDRAN, Universitat Pompeu Fabra, Barcelona, Spain and University of Adelaide, Adelaide, Australia

BRAM VAN GINNEKEN, Fraunhofer MEVIS, Bremen, Germany and Radboud Institute for Health Sciences, Radboud University Medical Center, Nijmegen, The Netherlands

BEN GLOCKER, Department of Computing, Imperial College London, London, UK

PATRICK GODAU, German Cancer Research Center (DKFZ) Heidelberg, Division of Intelligent Medical Systems, Germany, Faculty of Mathematics and Computer Science, Heidelberg University, Heidelberg, Germany, and National Center for Tumor Diseases (NCT), NCT Heidelberg, a partnership between DKFZ and University Medical Center Heidelberg, Germany

ROBERT HAASE, Now with: Center for Scalable Data Analytics and Artificial Intelligence (ScaDS.AI), Leipzig University, Leipzig, Germany, DFG Cluster of Excellence "Physics of Life", Technische Universität (TU) Dresden, Dresden, Germany, and Center for Systems Biology, Dresden, Germany

DANIEL A. HASHIMOTO, Department of Surgery, Perelman School of Medicine, Philadelphia, PA, USA and General Robotics Automation Sensing and Perception Laboratory, School of Engineering and Applied Science, University of Pennsylvania, Philadelphia, PA, USA

MICHAEL M. HOFFMAN, Princess Margaret Cancer Centre, University Health Network, Toronto, Canada, Department of Medical Biophysics, University of Toronto, Toronto, Canada, Department of Computer Science, University of Toronto, Toronto, Canada, and Vector Institute for Artificial Intelligence, Toronto, Canada

MEREL HUISMAN, Department of Radiology and Nuclear Medicine, Radboud University Medical Center, Nijmegen, The Netherlands

FABIAN ISENSEE, German Cancer Research Center (DKFZ) Heidelberg, Division of Medical Image Computing and HI Applied Computer Vision Lab, Germany

PIERRE JANNIN, Laboratoire Traitement du Signal et de l'Image – UMR_S 1099, Université de Rennes 1, Rennes, France and INSERM, Paris Cedex, France

CHARLES E. KAHN, Department of Radiology and Institute for Biomedical Informatics, University of Pennsylvania, Philadelphia, PA, USA

DAGMAR KAINMUELLER, Max-Delbrück Center for Molecular Medicine in the Helmholtz Association (MDC), Biomedical Image Analysis and HI Helmholtz Imaging, Berlin, Germany and University of Potsdam, Digital Engineering Faculty, Potsdam, Germany

BERNHARD KAINZ, Department of Computing, Faculty of Engineering, Imperial College London, London, UK and Department AIBE, Friedrich-Alexander-Universität (FAU), Erlangen-Nürnberg, Germany

ALEXANDROS KARARGYRIS, IHU Strasbourg, Strasbourg, France

ALAN KARTHIKESALINGAM, Google Health DeepMind, London, UK

HANNES KENNGOTT, Department of General, Visceral and Transplantation Surgery, Heidelberg University Hospital, Heidelberg, Germany

JENS KLEESIEK, Translational Image-guided Oncology (TIO), Institute for AI in Medicine (IKIM), University Medicine Essen, Essen, Germany

FLORIAN KOFLER, Helmholtz AI, München, Germany

THIJS KOOI, Lunit, Seoul, South Korea

ANNETTE KOPP-SCHNEIDER, German Cancer Research Center (DKFZ) Heidelberg, Division of Biostatistics, Germany

MICHAL KOZUBEK, Centre for Biomedical Image Analysis and Faculty of Informatics, Masaryk University, Brno, Czech Republic

ANNA KRESHUK, Cell Biology and Biophysics Unit, European Molecular Biology Laboratory (EMBL), Heidelberg, Germany

TAHSIN KURC, Department of Biomedical Informatics, Stony Brook University, Stony Brook, NY, USA

BENNETT A. LANDMAN, Electrical Engineering, Vanderbilt University, Nashville, TN, USA

GEERT LITJENS, Department of Pathology, Radboud University Medical Center, Nijmegen, The Netherlands

AMIN MADANI, Department of Surgery, University Health Network, Philadelphia, PA, Canada

KLAUS MAIER-HEIN, German Cancer Research Center (DKFZ) Heidelberg, Division of Medical Image Computing and HI Helmholtz Imaging, Germany and Pattern Analysis and Learning Group, Department of Radiation Oncology, Heidelberg University Hospital, Heidelberg, Germany

ANNE L. MARTEL, Physical Sciences, Sunnybrook Research Institute, Toronto, Canada and Department of Medical Biophysics, University of Toronto, Toronto, ON, Canada
PETER MATTSON, Google, Mountain View, USA
ERIK MEIJERING, School of Computer Science and Engineering, University of New South Wales, Sydney, Australia
BJOERN MENZE, Department of Quantitative Biomedicine, University of Zurich, Zurich, Switzerland
KAREL G.M. MOONS, Julius Center for Health Sciences and Primary Care, UMC Utrecht, Utrecht University, Utrecht, The Netherlands
HENNING MÜLLER, Information Systems Institute, University of Applied Sciences Western Switzerland (HES-SO), Sierre, Switzerland and Medical Faculty, University of Geneva, Geneva, Switzerland
BRENNAN NICHYPORUK, MILA (Quebec Artificial Intelligence Institute), Montréal, Canada
FELIX NICKEL, Department of General, Visceral and Thoracic Surgery, University Medical Center Hamburg-Eppendorf, Hamburg, Germany
JENS PETERSEN, German Cancer Research Center (DKFZ) Heidelberg, Division of Medical Image Computing, Germany
SUSANNE M. RAFELSKI, Allen Institute for Cell Science, Seattle, WA, USA
NASIR RAJPOOT, Tissue Image Analytics Laboratory, Department of Computer Science, University of Warwick, Coventry, UK
MAURICIO REYES, ARTORG Center for Biomedical Engineering Research, University of Bern, Bern, Switzerland and Department of Radiation Oncology, University Hospital Bern, University of Bern, Bern, Switzerland
MICHAEL A. RIEGLER, Simula Metropolitan Center for Digital Engineering, Oslo, Norway and UiT The Arctic University of Norway, Tromsø, Norway
NICOLA RIEKE, NVIDIA GmbH, München, Germany
JULIO SAEZ-RODRIGUEZ, Institute for Computational Biomedicine, Heidelberg University, Heidelberg, Germany and Faculty of Medicine, Heidelberg University Hospital, Heidelberg, Germany
CLARA I. SÁNCHEZ, Informatics Institute, Faculty of Science, University of Amsterdam, Amsterdam, The Netherlands
SHRAVYA SHETTY, Google Health, Google, CA, USA
RONALD M. SUMMERS, National Institutes of Health Clinical Center, Bethesda, MD, USA
ABDEL A. TAHA, Institute of Information Systems Engineering, TU Wien, Vienna, Austria
ALEKSEI TIULPIN, Research Unit of Health Sciences and Technology, Faculty of Medicine, University of Oulu, Oulu, Finland and Neurocenter Oulu, Oulu University Hospital, Oulu, Finland
SOTIRIOS A. TSAFTARIS, School of Engineering, The University of Edinburgh, Edinburgh, Scotland
BEN VAN CALSTER, Department of Development and Regeneration and EPI-centre, KU Leuven, Leuven, Belgium and Department of Biomedical Data Sciences, Leiden University Medical Center, Leiden, The Netherlands
GAËL VAROQUAUX, Parietal project team, INRIA Saclay-Île de France, Palaiseau, France
ZIV R. YANIV, National Institute of Allergy and Infectious Diseases, National Institutes of Health, Bethesda, MD, USA
PAUL F. JÄGER[‡], German Cancer Research Center (DKFZ) Heidelberg, Interactive Machine Learning Group and HI Helmholtz Imaging, Germany
LENA MAIER-HEIN[‡], German Cancer Research Center (DKFZ) Heidelberg, Division of Intelligent Medical Systems and HI Helmholtz Imaging, Germany, Faculty of Mathematics and Computer Science and Medical Faculty, Heidelberg University, Heidelberg, Germany, and National Center for Tumor Diseases (NCT), NCT Heidelberg, a partnership between DKFZ and University Medical Center Heidelberg, Germany

***Corresponding authors:** Annika Reinke: a.reinke@dkfz-heidelberg.de; Minu D. Tizabi: Lena Maier-Hein, l.maier-hein@dkfz-heidelberg.de; Paul F. Jäger: p.jaeger@dkfz-heidelberg.de; Lena Maier-Hein, l.maier-hein@dkfz-heidelberg.de.

[†]Shared first authors: Annika Reinke and Minu D. Tizabi

[‡]Shared last authors: Paul F. Jäger and Lena Maier-Hein

Abstract: Validation metrics are key for the reliable tracking of scientific progress and for bridging the current chasm between artificial intelligence (AI) research and its translation into practice. However, increasing evidence shows that particularly in image analysis, metrics are often chosen inadequately in relation to the underlying research problem. This could be attributed to a lack of accessibility of metric-related knowledge: While taking into account the individual strengths, weaknesses, and limitations of validation metrics is a critical prerequisite to making educated choices, the relevant knowledge is currently scattered and poorly accessible to individual researchers. Based on a multi-stage Delphi process conducted by a multidisciplinary expert consortium as well as extensive community feedback, the present work provides the first reliable and comprehensive common point of access to information on pitfalls related to validation metrics in image analysis. Focusing on biomedical image analysis but with the potential of transfer to other fields, the addressed pitfalls generalize across application domains and are categorized according to a newly created, domain-agnostic taxonomy. To facilitate comprehension, illustrations and specific examples accompany each pitfall. As a structured body of information accessible to researchers of all levels of expertise, this work enhances global comprehension of a key topic in image analysis validation.

Keywords: Validation, Evaluation, Pitfalls, Metrics, Good Scientific Practice, Biomedical Image Processing, Challenges, Computer Vision, Classification, Segmentation, Instance Segmentation, Semantic Segmentation, Detection, Localization, Medical Imaging, Biological Imaging

MAIN

Measuring performance and progress in any given field critically depends on the availability of meaningful outcome metrics. In a field such as athletics, this process is straightforward because the performance measurements (e.g., the time it takes an athlete to run a given distance) exactly reflect the underlying interest (e.g., which athlete runs a given distance the fastest?). In image analysis, the situation is much more complex as, depending on the underlying research question, vastly different aspects of an algorithm's performance might be of interest (Fig. 1) and meaningful in determining its future practical, for example clinical, applicability. If the performance of an image analysis algorithm is not measured according to relevant validation metrics, no reliable statement can be made about the suitability of this algorithm in solving the proposed task, and the algorithm is unlikely to ever reach the stage of real-life application. Moreover, unsuitable algorithms could be wrongly regarded as the best-performing ones, sparking entirely futile resource investment and follow-up research while obscuring true scientific advancements. In determining new state-of-the-art methods and informing future directions, the use of validation metrics actively shapes the evolution of research. In summary, *validation metrics are the key for both measuring and informing scientific progress, as well as bridging the current chasm between image analysis research and its translation into practice.*

In image analysis, while for some applications it might, for instance, be sufficient to draw a box around the structure of interest (e.g., a polyp in colonoscopic polyp detection), other applications (e.g., tumor volume delineation for radiotherapy planning) could require determining the exact structure boundaries. The suitability of any individual validation metric thus depends crucially on the properties of the driving image analysis problem. As a result, numerous metrics have so far been proposed in the field of image processing. In our previous work, we analyzed all biomedical image analysis competitions conducted within a period of about 15 years [57]. We found a total of 97 different metrics reported in the field of biomedicine alone, each with its own individual strengths, weaknesses, and limitations, and hence varying degrees of suitability for meaningfully measuring algorithm performance on any given research problem. Such a vast lake of options makes tracking all related information impossible for any individual researcher and consequently renders the process of metric selection error-prone. Thus, the frequent reliance on flawed, historically grown validation practices in current literature comes as no surprise. To make matters worse, there is currently no comprehensive resource that can provide an overview of the relevant definitions, (mathematical) properties, limitations, and pitfalls pertaining to a metric of interest. *While taking into account the individual properties and limitations of metrics is imperative for choosing adequate validation metrics, the required knowledge is thus largely inaccessible.*

As a result, numerous flaws and pitfalls are prevalent in image analysis validation, with researchers often being unaware of them due to a lack of knowledge of intricate metric properties and limitations. Accordingly, increasing evidence shows that metrics are often selected inadequately in image analysis (e.g., [34, 48, 83]). In the absence of a central information resource, it is common for researchers to resort to popular validation metrics, which, however, can be entirely unsuitable, for instance due to a mismatch of the metric's inherent mathematical properties with the underlying research question and specifications of the data set at hand (see Fig. 1).

The present work addresses this important roadblock in image analysis research with a crowd-sourcing-based approach that involved both a Delphi process undergone by a multidisciplinary expert consortium as well as a social media campaign. It represents the *first comprehensive collection, visualization, and detailed discussion of pitfalls, drawbacks, and limitations regarding validation metrics commonly used in image analysis.* Our work provides researchers with a *reliable, single point*

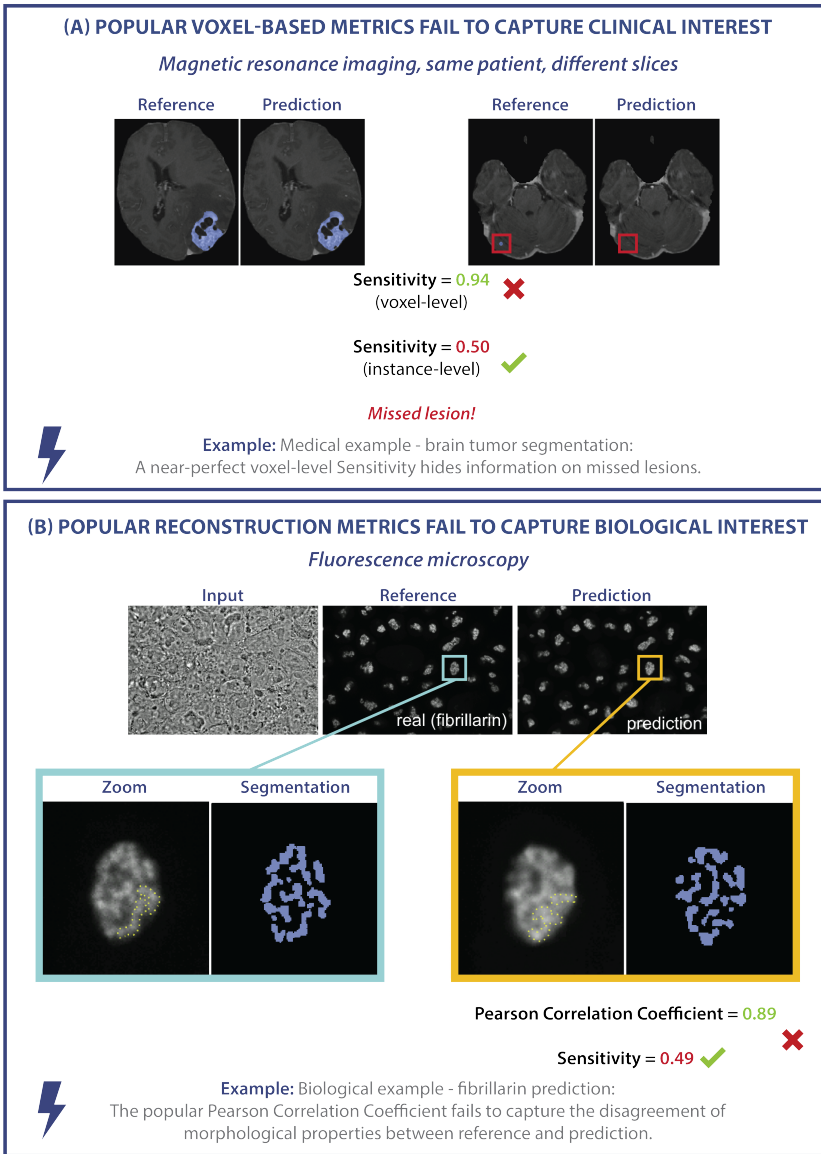


Fig. 1. Examples of metric-related pitfalls in image analysis validation. (A) Medical image analysis example: Voxel-based metrics are not appropriate for detection problems. Measuring the voxel-level performance of a prediction yields a near-perfect Sensitivity. However, the Sensitivity at the instance level reveals that lesions are actually missed by the algorithm. (B) Biological image analysis example: The task of predicting fibrillar in the dense fibrillary component of the nucleolus should be phrased as a segmentation task, for which segmentation metrics reveal the low quality of the prediction. Phrasing the task as image reconstruction instead and validating it using metrics such as the Person Correlation Coefficient yields misleadingly high metric scores [12, 67, 73, 87, 87].

of access to this critical and yet, until now, poorly retrievable or outright unavailable information. Owing to the enormous complexity of the matter, the metric properties and pitfalls are discussed in the specific context of classification problems, i.e., image analysis problems that can be considered classification tasks at either the image, object, or pixel level. Specifically, these encompass the four problem categories of image-level classification, semantic segmentation, object detection, and instance segmentation. Our contribution includes a dedicated profile for each metric (Suppl. Note 3) as well as the creation of a new common taxonomy that categorizes pitfalls in a domain-agnostic manner (Fig. 2). Depicted for individual metrics in tables provided in this paper (see Extended Data Tabs. 1-5), the taxonomy enables researchers to quickly grasp whether using a certain metric comes with pitfalls in a given use case. While our work grew out of image analysis research and practice in the field of biomedicine, a field of high complexity and particularly high stakes due to its direct impact on human health, we believe the identified pitfalls to be transferable to other application areas of imaging research. It should be noted that this work focuses on identifying, categorizing, and illustrating metric pitfalls, while the sister publication of this work gives specific recommendations on which metrics to apply under which circumstances [58].

RESULTS

Information on metric pitfalls is largely inaccessible

Researchers and algorithm developers seeking to validate image analysis algorithms frequently face the problem of choosing adequate validation metrics while at the same time navigating a range of potential pitfalls. Following common practice is often not the best option, as evidenced by a number of recent publications [34, 48, 57, 83]. Making an educated choice from a vast array of possibilities requires a researcher to be aware of not only the definitions and mathematical properties of different metrics but also their strengths and weaknesses, as well as limitations related to their use under certain conditions. The endeavor is notably complicated by the absence of any comprehensive databases or reviews covering the topic and thus the lack of a central resource for reliable information on validation metrics.

This lack of accessibility is considered by experts to be a major bottleneck in image analysis validation [57]. To illustrate this point, we searched the literature for available information on commonly used validation metrics. The search was conducted on the platform Google Scholar using search strings that combined different notations of the metric name, including synonyms and acronyms, with search terms indicating problems, such as “pitfall” or “limitation”. The mean and median number of hits for the metrics addressed in the present work were 159,329 and 22,100, respectively, and ranged between 49 for centerline Dice Similarity Coefficient (cDice) and 962,000 for Sensitivity. Moreover, despite valuable literature on individual relevant aspects (e.g., [14, 15, 36, 48, 79, 80, 83]), we did not find a common point of entry to metric-related pitfalls in image analysis in the form of a review paper or other credible source. It is thus unfeasible for any individual researcher to, within reasonable time and effort, retrieve comprehensive information on properties and pitfalls pertaining to one or multiple metrics of interest from the current body of research literature. We conclude that *the key knowledge required for making educated decisions and avoiding pitfalls related to the use of validation metrics is highly scattered and not accessible by individuals.*

Historically grown practices are not always justified

To obtain an initial insight into current common practice regarding validation metrics, we prospectively captured the designs of challenges organized by the IEEE Society of the International Symposium of Biomedical Imaging (ISBI), the Medical Image Computing and Computer Assisted Interventions (MICCAI) Society and the Medical Imaging with Deep Learning (MIDL) foundation. The organizers of the respective competitions were asked to provide a rationale for the choice of metrics in their competition. An analysis of a total of 138 competitions conducted between 2018 and 2022 revealed that metrics are frequently (in 24% of the competitions) based on common practice in the community. We found, however, that common practices are often not well-justified, and poor practices may even be propagated from one generation to the next.

One remarkable example for this issue is the widespread adoption of an incorrect naming and inconsistent mathematical formulation of a metric proposed for cell instance segmentation. The term "mean Average Precision (mAP)" usually refers to one of the most common metrics in object detection (object-level classification) [56, 72]. Here, Precision denotes the Positive Predictive Value (PPV), which is "averaged" over varying thresholds on the predicted class scores of an object detection algorithm. The "mean" Average Precision (AP) is then obtained by taking the mean over classes [29, 72]. Despite the popularity of mAP, a widely known challenge on cell instance segmentation¹ introduced a new "Mean Average Precision" in 2018. Although the task matches the task of the original "mean" AP, object detection, all terms in the newly proposed metric (mean, average, and precision) refer to entirely different concepts. For instance, the common definition of Precision from literature $TP/(TP + FP)$ was altered to $TP/(TP + FP + FN)$, where TP, FP, and FN refer to the cardinalities of the confusion matrix (i.e., the true/false positives/negatives). The latter formula actually defines the Intersection over Union (IoU) metric. Despite these problems, the terminology was adopted by subsequent influential works [46, 76, 78], indicating widespread propagation and usage within the community.

A multidisciplinary Delphi process reveals numerous pitfalls in biomedical image analysis validation

With the aim of creating a comprehensive, reliable collection and future point of access to biomedical image analysis metric definitions and limitations, we formed an international multidisciplinary consortium of 62 experts from various biomedical image analysis-related fields that engaged in a multi-stage Delphi process [9] for consensus building. Further pitfalls were crowdsourced through the publication of a dynamic preprint of this work [72] as well as a social media campaign, both of which asked the scientific community for contributions. This approach allowed us to integrate distributed, cross-domain knowledge on metric-related pitfalls within a single resource. In total, the process revealed 37 distinct sources of pitfalls (see Fig. 2). Notably, these pitfall sources (e.g., class imbalances, uncertainties in the reference, or poor image resolution) can occur irrespective of a specific imaging modality or application. As a result, many pitfalls generalize across different problem categories in image processing (image-level classification, semantic segmentation, object detection, and instance segmentation), as well as imaging modalities and domains. A detailed discussion of all pitfalls can be found in Suppl. Note 2.

¹<https://www.kaggle.com/competitions/data-science-bowl-2018/overview/evaluation>

A common taxonomy enables domain-agnostic categorization of pitfalls

One of our key objectives was to facilitate information retrieval and provide structure within this vast topic. Specifically, we wanted to enable researchers to identify at a glance which metrics are affected by which types of pitfalls. To this end, we created a comprehensive taxonomy that categorizes the different pitfalls in a semantic fashion. The taxonomy was created in a domain-agnostic manner to reflect the generalization of pitfalls across different imaging domains and modalities. An overview of the taxonomy is presented in Fig. 2, and the relations between the pitfall categories and individual metrics can be found in Extended Data Tabs. 1-5. We distinguish the following three main categories:

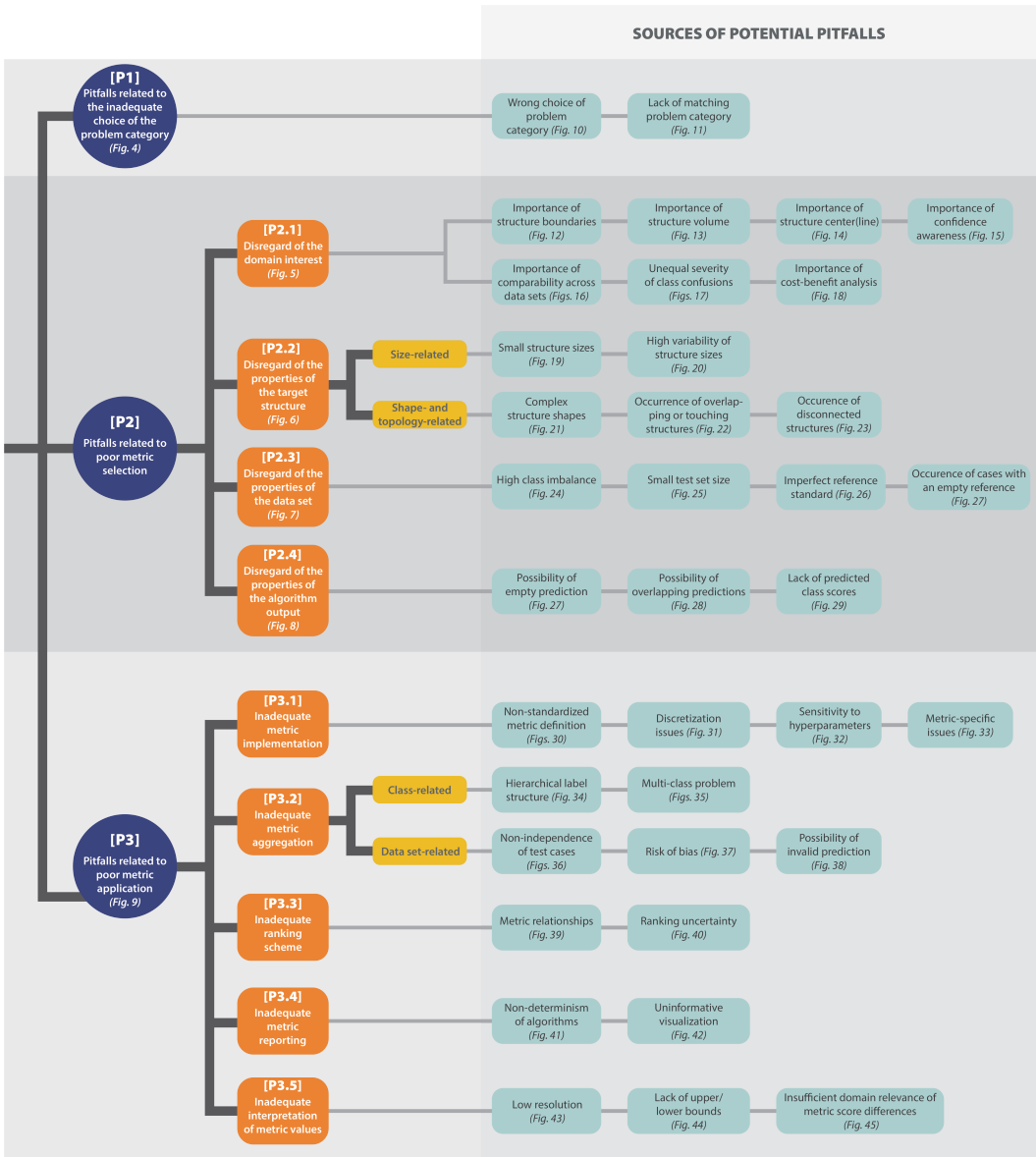


Fig. 2. Overview of the taxonomy for metric-related pitfalls. Pitfalls can be grouped into three main categories: [P1] Pitfalls related to the inadequate choice of the problem category, [P2] pitfalls related to poor metric selection, and [P3] pitfalls related to poor metric application. [P2] and [P3] are further split into subcategories. For all categories, pitfall sources are presented (green), with references to corresponding illustrations of representative examples. Note that the order in which the pitfall sources are presented does not correlate with importance.

[P1] Pitfalls related to the inadequate choice of the problem category. A common pitfall lies in the use of metrics for a problem category they are not suited for because they fail to fulfill crucial requirements of that problem category, and hence do not reflect the domain interest (Fig. 1). For instance, popular voxel-based metrics, such as the Dice Similarity Coefficient (DSC) or Sensitivity, are widely used in image analysis problems, although they do not fulfill the critical requirement of detecting all objects in a data set. In a cancer monitoring application they fail to measure instance progress, i.e., the potential increase in number of lesions (Fig. 1), which can have serious consequences for the patient. For some problems, there may even be a lack of matching problem category (Fig. SN 2.4), rendering common metrics inadequate. We present further examples of pitfalls in this category in Suppl. Note 2.1.

[P2] Pitfalls related to poor metric selection. Pitfalls of this category occur when a validation metric is selected while disregarding specific properties of the given research problem or method used that make this metric unsuitable in the particular context. [P2] can be further divided into the following four subcategories:

[P2.1] Disregard of the domain interest. Commonly, several requirements arise from the domain interest of the underlying research problem that may clash with particular metric limitations. For example, if there is particular interest in the structure boundaries, it is important to know that overlap-based metrics such as the DSC do not take the correctness of an object's boundaries into account, as shown in Fig. 4(a). Similar issues may arise if the structure volume (Fig. SN 2.6) or center(line) (Fig. SN 2.7) are of particular interest. Other domain interest-related properties may include an unequal severity of class confusions. This may be important in an ordinal grading use case, in which the severity of a disease is categorized by different scores. Predicting a low severity for a patient that actually suffers from a severe disease should be substantially penalized. Common classification metrics do not fulfill this requirement. An example is provided in Fig. 4(b). On pixel level, this property relates to an unequal severity of over- vs. undersegmentation. In applications such as radiotherapy, it may be highly relevant whether an algorithm tends to over- or undersegment the target structure. Common overlap-based metrics, however, do not represent over- and undersegmentation equally [95]. Further pitfalls may occur if confidence awareness (Fig. SN 2.8), comparability across data sets (Fig. SN 2.9), or a cost-benefit analysis (Fig. SN 2.11) are of particular importance, as illustrated in Suppl. Note 2.2.1.

[P2.2] Disregard of the properties of the target structures. For problems that require capturing local properties (object detection, semantic or instance segmentation), the properties of the target structures to be localized and/or segmented may have important implications for the choice of metrics. Here, we distinguish between *size-related* and *shape- and topology-related* pitfalls. Common metrics, for example, are sensitive to structure sizes, such that single-pixel differences may hugely impact the metric scores, as shown in Extended Data Fig. 1(a). Shape- and topology-related pitfalls may relate to the fact that common metrics disregard complex shapes (Extended Data Fig. 1(b)) or that bounding boxes do not capture the disconnectedness of structures (Fig. SN 2.16). A high variability of structure sizes (Fig. SN 2.13) and overlapping or touching structures (Fig. SN 2.15) may also influence metric values. We present further examples of [P2.2] pitfalls in Suppl. Note 2.2.2.

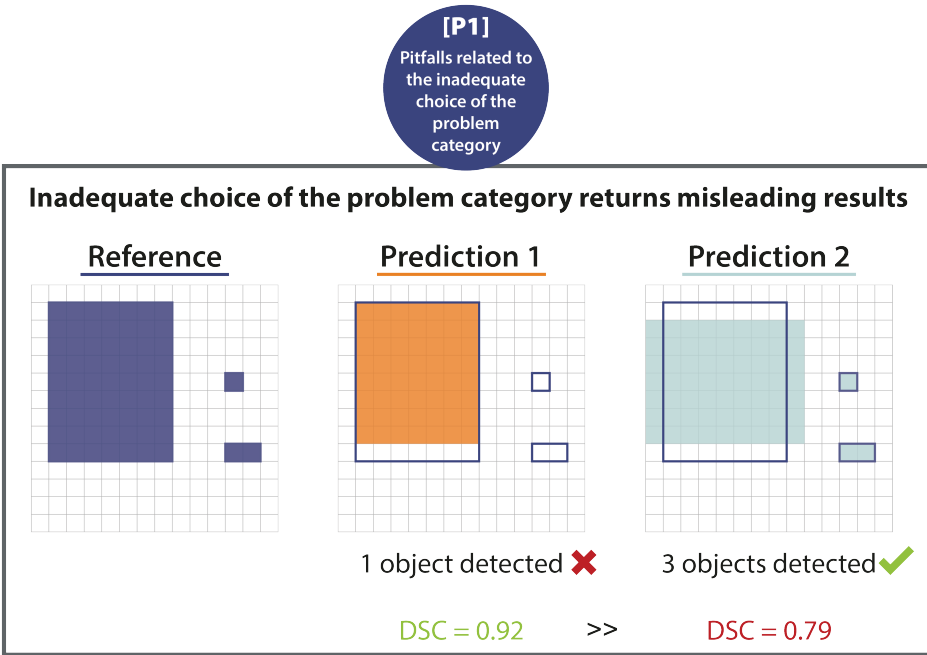


Fig. 3. [P1] Pitfalls related to the inadequate choice of the problem category. **Wrong choice of problem category.** Effect of using segmentation metrics for object detection problems. The pixel-level Dice Similarity Coefficient (DSC) of a prediction recognizing every structure (*Prediction 2*) is lower than that of a prediction that only recognizes one of the three structures (*Prediction 1*).

[P2.3] *Disregard of the properties of the data set.* Various properties of the data set such as class imbalances (Fig. 5(a)), small sample sizes (Fig. 5(b)), or the quality of the reference annotations, may directly affect metric values. Common metrics such as the Balanced Accuracy (BA), for instance, may yield a very high score for a model that predicts many False Positive (FP) samples in an imbalanced setting (see Fig. 5(a)). When only small test data sets are used, common calibration metrics (which are typically biased estimators) either underestimate or overestimate the true calibration error of a model (Fig. 5(b)) [37]. On the other hand, metric values may be impacted by reference annotations (Fig. SN 2.19). Spatial outliers in the reference may have a huge impact on distance-based metrics such as the Hausdorff Distance (HD) (Fig. 5(c)). Additional pitfalls may arise from the occurrence of cases with an empty reference (Extended Data Fig. 2(b)), causing division by zero errors. We present further examples of [P2.3] pitfalls in Suppl. Note 2.2.3.

[P2.4] *Disregard of the properties of the algorithm output.* Reference-based metrics compare the algorithm output to a reference annotation to compute a metric score. Thus, the content and format of the prediction are of high importance when considering metric choice. Overlapping predictions in segmentation problems, for instance, may return misleading results. In Extended Data Fig. 2(a), the predictions only overlap to a certain extent, not representing that the reference instances actually overlap substantially. This is not detected by common metrics. Another example are empty predictions that may cause division by zero errors in metric calculations, as illustrated in Extended Data Fig. 2(b), or the lack of predicted class scores (Fig. SN 2.22). We present further examples of [P2.4] pitfalls in Suppl. Note 2.2.3.

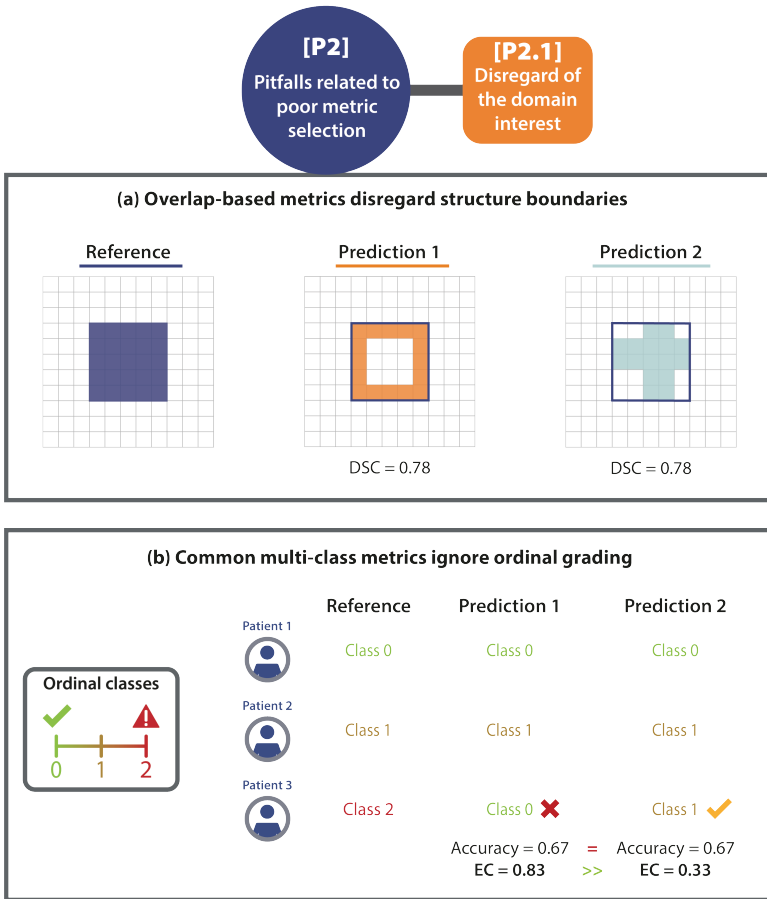


Fig. 4. [P2.1] **Disregard of the domain interest.** (a) **Importance of structure boundaries.** The predictions of two algorithms (*Prediction 1/2*) capture the boundary of the given structure substantially differently, but lead to the exact same Dice Similarity Coefficient (DSC), due to its boundary unawareness. This pitfall is also relevant for other overlap-based metrics such as centerline Dice Similarity Coefficient (cIDice), pixel-level F_β Score, and Intersection over Union (IoU), as well as localization criteria such as Box/Approx/Mask IoU, Center Distance, Mask IoU > 0, Point inside Mask/Box/Approx, and Intersection over Reference (IoR). (b) **Unequal severity of class confusions.** When predicting the severity of a disease for three patients in an ordinal classification problem, *Prediction 1* assumes a much lower severity for *Patient 3* than actually observed. This critical issue is overlooked by common metrics (here: Accuracy), which measure no difference to *Prediction 2*, which assesses the severity much better. Metrics with pre-defined weights (here: Expected Cost (EC)) correctly penalize *Prediction 1* much more than *Prediction 2*. This pitfall is also relevant for other counting metrics, such as Balanced Accuracy (BA), F_β Score, Positive Likelihood Ratio (LR+), Matthews Correlation Coefficient (MCC), Net Benefit (NB), Negative Predictive Value (NPV), Positive Predictive Value (PPV), Sensitivity, and Specificity.

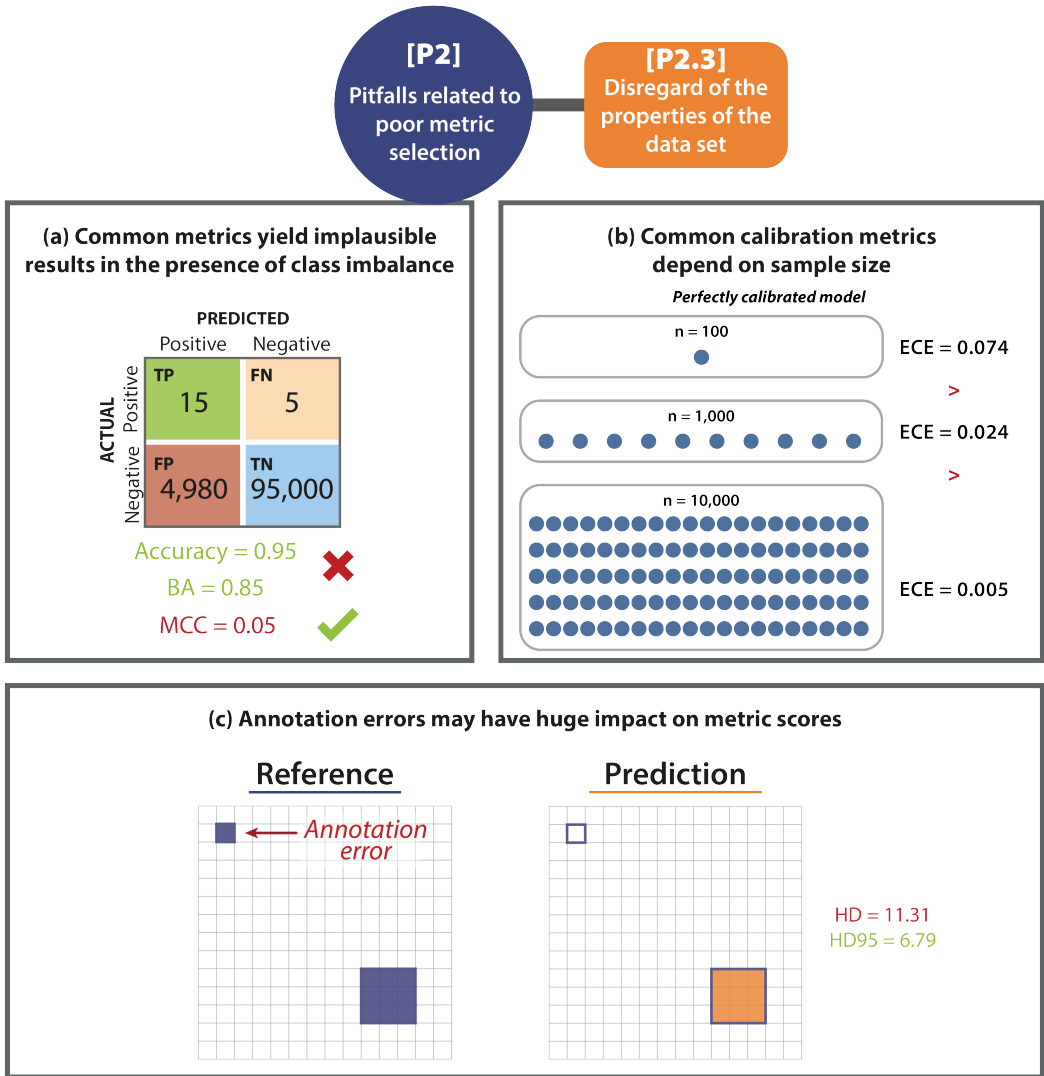


Fig. 5. [P2.3] **Disregard of the properties of the data set.** (a) **High class imbalance.** In the case of underrepresented classes, common metrics may yield misleading values. In the given example, Accuracy and Balanced Accuracy (BA) have a high score despite the high amount of False Positive (FP) samples. The class imbalance is only uncovered by metrics considering predictive values (here: Matthews Correlation Coefficient (MCC)). This pitfall is also relevant for other counting and multi-threshold metrics such as Area under the Receiver Operating Characteristic Curve (AUROC), Expected Cost (EC) (depending on the chosen costs), Positive Likelihood Ratio (LR+), Net Benefit (NB), Sensitivity, Specificity, and Weighted Cohen’s Kappa (WCK). (b) **Small test set size.** The values of the Expected Calibration Error (ECE) depend on the sample size. Even for a simulated perfectly calibrated model, the ECE will be substantially greater than zero for small sample sizes [37]. (c) **Imperfect reference standard.** A single erroneously annotated pixel may lead to a large decrease in performance, especially in the case of the Hausdorff Distance (HD) when applied to small structures. The Hausdorff Distance 95th Percentile (HD95), on the other hand, was designed to deal with spatial outliers. This pitfall is also relevant for localization criteria such as Box/Approx Intersection over Union (IoU) and Point inside Box/Approx. Further abbreviations: True Positive (TP), False Negative (FN), True Negative (TN).

[P3] Pitfalls related to poor metric application. Once selected, the metrics need to be applied to an image or an entire data set. This step is not straightforward and comes with several pitfalls. For instance, when aggregating metric values over multiple images or patients, a common mistake is to ignore the hierarchical data structure, such as data from several hospitals or a varied number of images per patient. We present three examples of [P3] pitfalls in Fig. 6; for more pitfalls in this category, please refer to Suppl. Note 2.3. [P3] can further be divided into five subcategories that are presented in the following paragraphs.

[P3.1] Inadequate metric implementation. Metric implementation is, unfortunately, not standardized. As shown by [35], different researchers typically employ various different implementations for the same metric, which may yield a substantial variation in the metric scores. While some metrics are straightforward to implement, others require more advanced techniques and offer different possibilities. In the following, we provide some examples for inadequate metric implementation:

- The method of how identical confidence scores are handled in the computation of the AP metric may lead to substantial differences in the metric scores. Microsoft Common Objects in Context (COCO) [56], for instance, processes each prediction individually, while CityScapes [18] processes all predictions with the same score in one joint step. Fig. 6(a) provides an example with two predictions having the same confidence score, in which the final metric scores differ depending on the chosen handling strategy for identical confidence scores. Similar issues may arise with other curve-based metrics, such as AUROC, AP, or Free-Response Receiver Operating Characteristic (FROC) scores (see e.g., [62]).
- Metric implementation may be subject to discretization issues such as the chosen discretization of continuous variables, which may cause differences in the metric scores, as exemplary illustrated in Fig. SN 2.24.
- For metrics assessing structure boundaries, such as the Average Symmetric Surface Distance (ASSD), the exact boundary extraction method is not standardized. Thus, for example, the boundary extraction method implemented by the Liver Tumor Segmentation (LiTS) challenge [7] and that implemented by Google DeepMind² may produce different metric scores for the ASSD. This is especially critical for metrics that are sensitive to small contour changes, such as the HD.
- Suboptimal choices of hyperparameters may also lead to metric scores that do not reflect the domain interest. For example, the choice of a threshold on a localization criterion (see Fig. SN 2.25) or the chosen hyperparameter for the F_β Score will heavily influence the subsequent metric scores [82].

More [P3.1] pitfalls can be found in Suppl. Note 2.3.1.

[P3.2] Inadequate metric aggregation. A common pitfall with respect to metric application is to simply aggregate metric values over the entire data set and/or all classes. As detailed in Fig. 6(b) and Suppl. Note 2.3.2, important information may get lost in this process, and metric results can be misleading. For example, the popular TorchMetrics framework calculates the DSC metric by default as a global average over all pixels in the data set without considering their image or class of origin³. Such a calculation eliminates the possibility of interpreting the final metric score with respect to individual images and classes. For example, errors in small structures may be suppressed by correctly segmented larger structures in other images (see e.g., Fig. SN 2.28). An adequate

²<https://github.com/deepmind/surface-distance>

³<https://torchmetrics.readthedocs.io/en/stable/classification/dice.html?highlight=dice>

aggregation scheme is also crucial for handling hierarchical class structure (Fig. SN 2.29), missing values (Fig. SN 2.31), and potential biases (Fig. SN 2.30) of the algorithm. Further [P3.2] pitfalls are shown in Suppl. Note 2.3.2.

[P3.3] Inadequate ranking scheme. Rankings are often created to compare algorithm performances. In this context, several pitfalls pertain to either metric relationships or ranking uncertainty. For example, to assess different properties of an algorithm, it is advisable to select multiple metrics and determine their values. However, the chosen metrics should assess complementary properties and should not be mathematically related. For example, the DSC and IoU are closely related, so using both in combination would not provide any additional information over using either of them individually (Fig. SN 2.32). Note in this context that unawareness of metric synonyms can equally mislead. Metrics can be known under different names; for instance, Sensitivity and Recall refer to the same mathematical formula. Despite this fact potentially appearing trivial, an analysis of 138 biomedical image analysis challenges [58] found three challenges that unknowingly used two versions of the same metric to calculate their rankings. Moreover, rankings themselves may be unstable (Fig. SN 2.33). [57] and [93] demonstrated that rankings are highly sensitive to altering the metric aggregation operators, the underlying data set, or the general ranking method. Thus, if the robustness of rankings is disregarded, the winning algorithm may be identified by chance rather than true superiority.

[P3.4] Inadequate metric reporting. A thorough reporting of metric values and aggregates is important both in terms of transparency and interpretability. However, several pitfalls are to be avoided in this regard. Notably, different types of visualization may vary substantially in terms of interpretability, as shown in Figs 6(c). For example, while a box plot provides basic information, it does not depict the distribution of metric values. This may conceal important information, such as specific images on which an algorithm performed poorly. Other pitfalls in this category relate to the non-determinism of algorithms, which introduces a natural variability to the results of a neural network, even with fixed seeds (Fig. SN 2.34). This issue is aggravated by inadequate reporting, for instance, reporting solely the results from the best run instead of proper cross-validation and reporting of the variability across different runs. Generally, shortcomings in reporting, such as providing no standard deviation or confidence intervals in the presented results, are common. Concrete examples of [P3.4] pitfalls can be found in Suppl. Note 2.3.4.

[P3.5] Inadequate interpretation of metric values. Interpreting metric scores and aggregates is an important step for the analysis of algorithm performances. However, several pitfalls can arise from the interpretation. In rankings, for example, minor differences in metric scores may not be relevant from an application perspective but may still yield better ranks (Fig. SN 2.38). Furthermore, some metrics do not have upper or lower bounds, or the theoretical bounds may not be achievable in practice, rendering interpretation difficult (Fig. SN 2.37). More information on interpretation-based pitfalls can be found in Suppl. Note 2.3.5.

[P3]
Pitfalls related to poor metric application

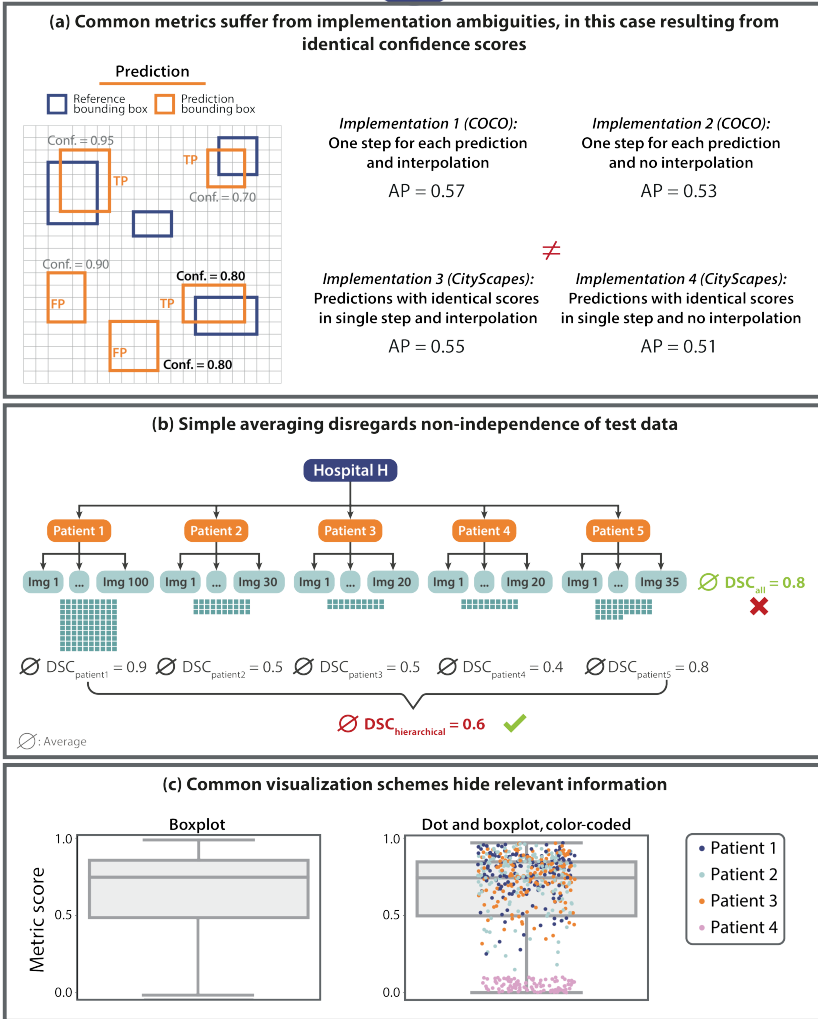


Fig. 6. [P3] **Pitfalls related to poor metric application.** (a) **Non-standardized metric implementation.** In the case of the Average Precision (AP) metric and the construction of the Precision-Recall (PR)-curve, the strategy of how identical scores (here: confidence score of 0.80 is present twice) are treated has a substantial impact on the metric scores. Microsoft Common Objects in Context (COCO) [56] and CityScapes [18] are used as examples. (b) **Non-independence of test cases.** The number of images taken from *Patient 1* is much higher compared to that acquired from *Patients 2-5*. Averaging over all Dice Similarity Coefficient (DSC) values, denoted by \emptyset , results in a high averaged score. Aggregating metric values per patient reveals much higher scores for *Patient 1* compared to the others, which would have been hidden by simple aggregation. (c) **Uninformative visualization.** A single box plot (left) does not give sufficient information about the raw metric value distribution. Adding the raw metric values as jittered dots on top (right) adds important information (here: on clusters). In the case of non-independent validation data, color/shape-coding helps reveal data clusters.

The first illustrated common access point to metric definitions and pitfalls

To underline the importance of a common access point to metric pitfalls, we conducted a search for individual metric-related pitfalls on the platforms Google Scholar and Google, with the purpose of determining how many of the pitfalls identified through our work could be located in existing resources. We were only able to locate a portion of the pitfalls identified by our approach in existing research literature (68%) or online resources such as blog posts (11%; 8% were found in both). Only 27% of the located pitfalls were presented visually.

Our work now provides this key resource in a highly structured and easily understandable form. Suppl. Note 2, contains a dedicated illustration for each of the pitfalls discussed, thus facilitating reader comprehension and making the information accessible to everyone regardless of their level of expertise. A further core contribution of our work are the metric profiles presented in Suppl. Note 2, which, for each metric, summarize the most important information deemed of particular relevance by the *Metrics Reloaded* consortium of the sister work to this publication [58]. The profiles provide the reader with a compact, at-a-glance overview of each metric and an enumeration of the limitations and pitfalls identified in the Delphi process conducted for this work.

DISCUSSION

Flaws in the validation of biomedical image analysis algorithms significantly impede the translation of methods into (clinical) practice and undermine the assessment of scientific progress in the field [55]. They are frequently caused by poor choices due to disregarding the specific properties and limitations of individual validation metrics. The present work represents the first comprehensive collection of pitfalls and limitations to be taken into account when using validation metrics in image-level classification, semantic segmentation, instance segmentation, and object detection tasks. Our work enables researchers to gain a deep understanding of and familiarity with both the overall topic and individual metrics by providing a common access point to previously largely scattered and inaccessible information — key knowledge they can resort to when conducting validation of image analysis algorithms. This way, our work aims to disrupt the current common practice of choosing metrics based on their popularity rather than their suitability to the underlying research problem. This practice, which, for instance, often manifests itself in the unreflected and inadequate use of the DSC, is concerningly prevalent even among prestigious, high-quality biomedical image analysis competitions (challenges) [19, 34, 43, 48, 49, 57, 59, 83]. The educational aspect of our work is complemented by dedicated 'metric profiles' which detail the definitions and properties of all metrics discussed. Notably, our work pioneers the examination of artificial intelligence (AI) validation pitfalls in the biomedical domain, a domain in which they are arguably more critical than in many others as flaws in biomedical algorithm validation can directly affect patient wellbeing and safety.

We posited that shortcomings in current common practice are marked by the low accessibility of information on the pitfalls and limitations of commonly used validation metrics. A literature search conducted from the point of view of a researcher seeking information on individual metrics confirmed that the number of search results far exceeds any amount that could be overseen within reasonable time and effort, as well as the lack of a common point of entry to reliable metric information. Even when knowing the specific pitfalls and related keywords uncovered by our consortium, only a fraction of those pitfalls could be found in existing literature, indicating the novelty and added value of our work.

For transparency, several constraints regarding our literature search must be noted. First, it must be acknowledged that the remarkably high search result numbers inevitably include duplicates of papers (e.g., the same work in a conference paper and on arXiv) as well as results that are out of scope (e.g., [11], [26]), in the cited examples for instance due to a metric acronym (AUC) simultaneously being an acronym for another entity (a trinucleotide) in a different domain, or the word "sensitivity" being used in its common, non-metric meaning. Moreover, common words used to describe pitfalls such as "problem" or "issue" are by nature present in many publications discussing any kind of research, rendering them unusable for a dedicated search, which could, in turn, account for missing publications that do discuss pitfalls in these terms. Similarly, when searching for specific pitfalls, many of the returned results containing the appropriate keywords did not actually refer to metrics or algorithm validation but to other parts of a model or biomedical problem (e.g., the need for stratification is commonly discussed with regard to the design of clinical studies but not with regard to their validation). Character limits in the Google Scholar search bar further complicate or prevent the use of comprehensive search strings. Finally, it is both possible and probable that our literature search did not retrieve all publications or non-peer-reviewed online resources that mention a particular pitfall, since even extensive search strings might not cover the particular words used for a pitfall description.

None of these observations, however, detracts from our hypothesis. In fact, all of the above observations reinforce our finding that, for any individual researcher, retrieving information on metrics of interest is difficult to impossible. In many cases, finding information on pitfalls only appears feasible if the specific pitfall and its related keywords are exactly known, which, of course, is not the situation most researchers realistically find themselves in. Overall accessibility of such vital information, therefore, currently leaves much to be desired.

Compiling this information through a multi-stage Delphi process allowed us to leverage distributed knowledge from experts across different biomedical imaging domains and thus ensure that the resulting illustrated collection of metric pitfalls and limitations turned out to be both comprehensive and of maximum practical relevance. Continued proximity of our work to issues occurring in practical application was achieved through sharing the first results of this process as a dynamic preprint [71] with dedicated calls for feedback, as well as crowdsourcing further suggestions on social media.

Although their severity and practical consequences might differ between applications, we found that the pitfalls generalize across different imaging modalities and application domains. By categorizing them solely according to their underlying sources, we were able to create an overarching taxonomy that goes beyond domain-specific concerns and thus enjoys broad applicability. Given the large number of identified pitfalls, our taxonomy crucially establishes structure in the topic. Moreover, by relating types of pitfalls to the respective metrics they apply to and illustrating them, it enables researchers to gain a deeper, systemic understanding of the causes of metric failure.

Our complementary *Metrics Reloaded* recommendation framework, which guides researchers towards the selection of appropriate validation metrics for their specific tasks and is introduced in a sister publication to this work [58], shares the same principle of domain independence. Its recommendations are based on the creation of a 'problem fingerprint' that abstracts from specific domain knowledge and, informed by the pitfalls discussed here, captures all properties relevant to metric selection for a specific biomedical problem. In this sister publication, we present recommendations to avoid the pitfalls presented in this work. Importantly, the finding that pitfalls generalize and can be categorized in a domain-independent manner opens up avenues for future expansion of

our work to other fields of ML-based imaging, such as general computer vision (see below), thus freeing it from its major constraint of exclusively focusing on biomedical problems.

It is worth mentioning that we only examined pitfalls related to the tasks of image-level classification, semantic segmentation, instance segmentation, and object detection, as these can all be considered classification tasks at different levels (image/object/pixel) and hence share similarities in their validation. While including a wider range of biomedical problems not considered classification tasks, such as regression or registration, would have gone beyond the scope of the present work, we envision this expansion in future work. Moreover, our work focused on pitfalls related to reference-based metrics. Including pitfalls pertaining to non-reference-based metrics, such as metrics that assess speed, memory consumption, or carbon footprint, could be a future direction to take. Finally, while we aspired to be as comprehensive as possible in our compilation, we cannot exclude that there are further pitfalls to be taken into account that the consortium and the participating community have so far failed to recognize. Should this be the case, our dynamic *Metrics Reloaded* online platform, which is currently under development and will continuously be updated after release, will allow us to easily and transparently append missed pitfalls. This way, our work can remain a reliable point of access, reflecting the state of the art at any given moment in the future. In this context, we note that we explicitly welcome feedback and further suggestions from the readership of *Nature Methods*.

The expert consortium was primarily compiled in a way to cover the required expertise from various fields but also consisted of researchers of different countries, (academic) ages, roles, and backgrounds (details can be found in the Methods). It mainly focused on biomedical applications. The pitfalls presented here are therefore of the highest relevance for biological and clinical use cases. Their clear generalization across different biomedical imaging domains, however, indicates broader generalizability to fields such as general computer vision. Future work could thus see a major expansion of our scope to AI validation well beyond biomedical research. Regardless of this possibility, we strongly believe that by raising awareness of metric-related pitfalls, our work will kick off a necessary scientific debate. Specifically, we see its potential in inducing the scientific communities in other areas of AI research to follow suit and investigate pitfalls and common practices impairing progress in their specific domains.

In conclusion, our work presents the first comprehensive and illustrated access point to information on validation metric properties and their pitfalls. We envision it to not only impact the quality of algorithm validation in biomedical imaging and ultimately catalyze faster translation into practice, but to raise awareness on common issues and call into question flawed AI validation practice far beyond the boundaries of the field.

METHODS

Literature search

The literature search of metric pitfalls and limitations was conducted on the platform Google Scholar. The checkbox "include patents" was activated and the checkbox "include citations" was deactivated; other default settings were left unchanged. For each metric, a specific search string using the Boolean operators OR and AND was generated as follows:

- (Different notations of the metric name, including synonyms and acronyms, enclosed in quotation marks, respectively, and combined with OR)
- AND "metric"

- AND (different expressions pertaining to the concept of pitfalls, limitations and flaws, enclosed in quotation marks, respectively, and combined with OR)

For example, the following search string was used for the literature search of DSC pitfalls: ("DSC" OR "Dice Similarity Coefficient" OR "Sørensen-Dice coefficient" OR "F1 score" OR "DCE") AND "metric" AND ("pitfall" OR "limitation" OR "caveat" OR "drawback" OR "shortcoming" OR "weakness" OR "flaw" OR "disadvantage" OR "suffer").

A second literature search dedicated to the pitfalls collected during the Delphi process was conducted on the platforms Google Scholar and Google. This search served the purpose of determining how many of the proposed pitfalls could be found in either existing research literature or online resources such as blogs, assuming that the issue is already roughly known to the person conducting the search. We further determined whether or not a found pitfall was presented in a visual manner. We analyzed the first three results pages (corresponding to thirty results) from each search platform and excluded our own previous work on metric pitfalls from the analysis.

Delphi process

The collection of pitfalls was achieved via a multi-stage Delphi process conducted among an international expert consortium comprised of more than 60 biomedical image analysis experts, as well as community feedback. A Delphi process is a structured group communication process that serves to pool opinions from an expert panel via a series of individual interrogations, usually in the form of questionnaires, interspersed with feedback from the respondents [9]. The technique is widely used for building consensus among experts in medicine, particularly in the development of best practices in areas where evidence may be limited, conflicting, or absent [65]. Expert selection was initially based on membership in major relevant societies such as the Biomedical Image Analysis ChallengeS (BIAS) initiative, the Medical Open Network for Artificial Intelligence (MONAI) Working Group for Evaluation, Reproducibility and Benchmarks, and the MICCAI Special Interest Group for Challenges (previously MICCAI board working group), as well as a track record of expertise in the areas of metrics, challenges and/or best practices. To reflect as broad a range of application areas and metric pitfalls as possible, the number of consortium members was increased throughout the process to a final number of 62 members. The Delphi process comprised four surveys. Each survey was developed by the coordinating team of the process and sent out to the remaining members of the consortium. Upon completion, the coordinating team then analyzed the results and iteratively refined the list of pitfalls. The main stages of the compilation and consensus building process are detailed in the following:

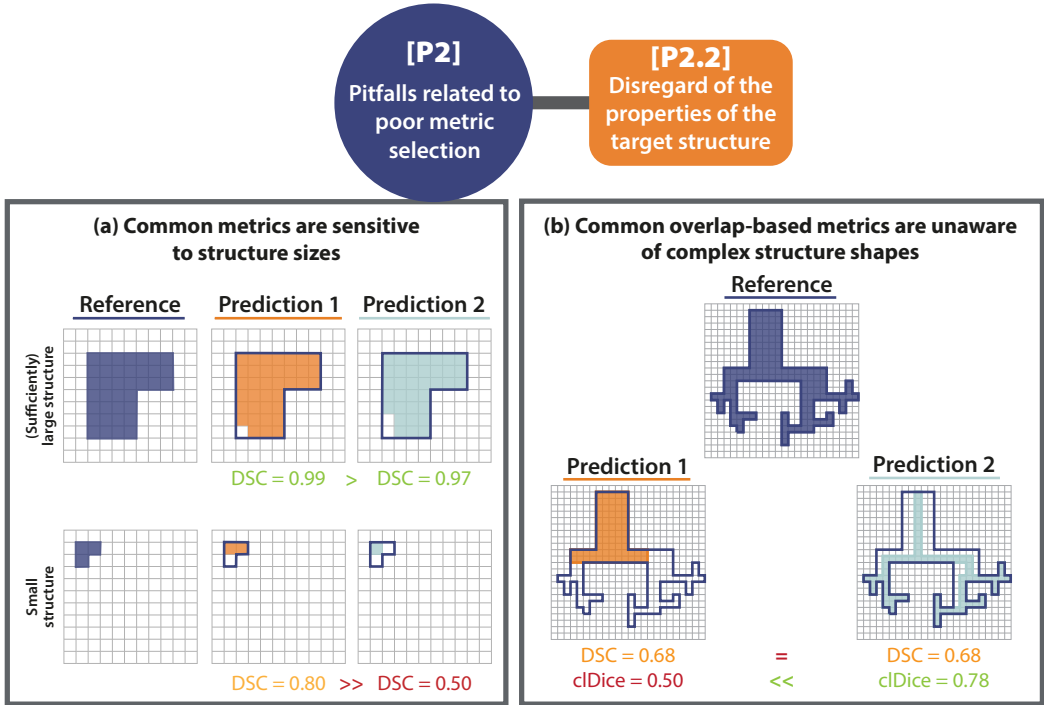
- (1) *Compilation of pitfall sources*: The primary purpose of the first survey was obtaining agreement on sources of pitfalls.
- (2) *Collection of pitfalls*: The following survey specifically asked for concrete pitfalls in the presence of those problem characteristics.
- (3) *Community feedback*: The proposed list of pitfalls was further complemented by social media-based feedback from the general scientific community.
- (4) *Final agreement on pitfalls*: The subsequent survey served to obtain consensus agreement on which pitfalls to include. For each pitfall, it asked whether the pitfall should be included. In addition, the experts were given the opportunity to provide feedback on each pitfall and to suggest further pitfalls. The final collection of pitfalls was illustrated and all metric values were verified by two independent observers.

- (5) *Creation of taxonomy*: The collected pitfalls were analyzed and a taxonomy was created. In the final survey, approval of the consortium for the structure and phrasing of the taxonomy and the assignment of specific pitfalls to the taxonomy was obtained.

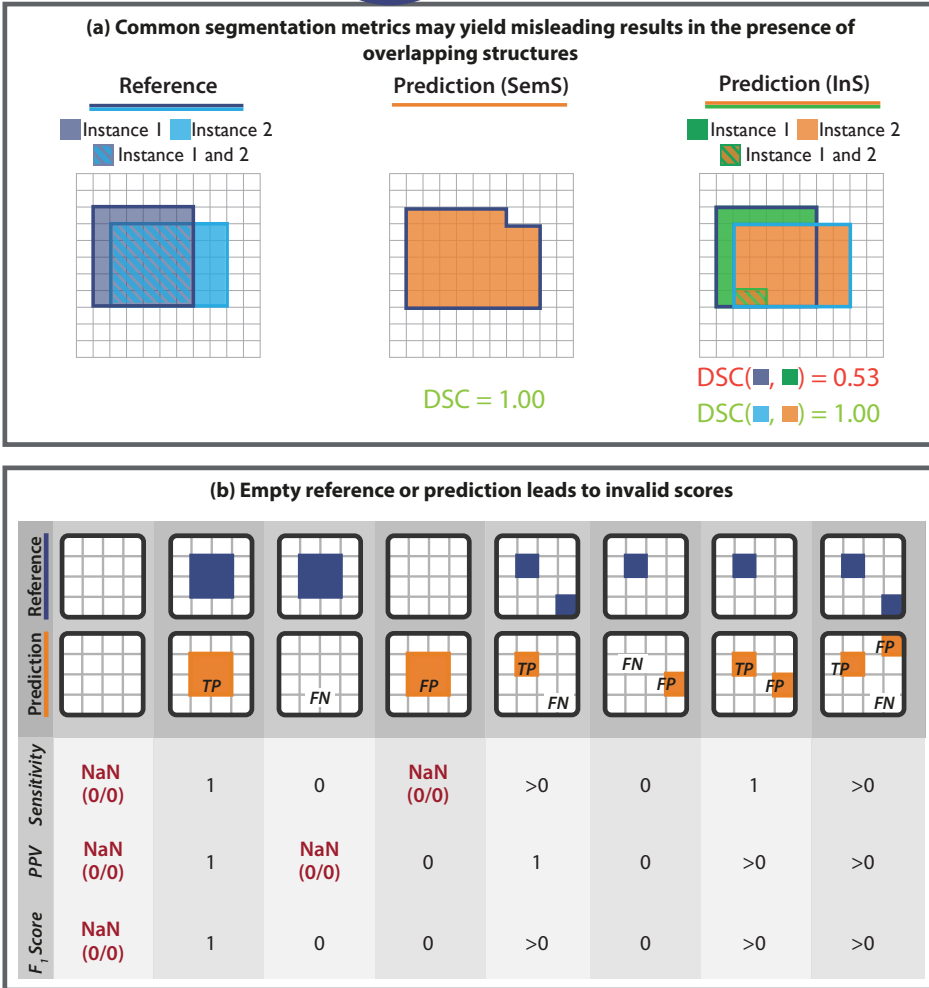
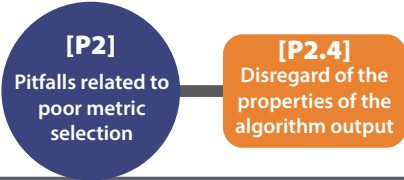
Expert consortium

The expert consortium consisted of a total of 70 researchers (70% male, 30% female) from a total of 65 institutions. The majority of experts (50%) were professors, followed by postdoctoral researchers (39%). The median h-index of the consortium was 31.5 (mean: 36; minimum: 6; maximum: 113) and the median academic age was 18 years (mean: 19; minimum: 3; max: 42). Experts were from 19 countries and 5 continents. 60% of experts had a technical, 6% a clinical, 3% a biological, and 23% a mixed background. Of the 65 institutions, we could identify the number of employees for 89%. Of those, the majority of institutions had a size between 1,000 and 10,000 employees (57%), followed by even larger institutions between 10,000 and 100,000 employees (22%), and smaller institutions below 1,000 employees (20%). Only a small portion of institutions were above 100,000 employees (2%).

EXTENDED DATA



Extended Data Fig. 1. [P2.2] **Disregard of the properties of the target structures.** (a) **Small structure sizes.** The predictions of two algorithms (*Prediction 1/2*) differ in only a single pixel. In the case of the small structure (bottom row), this has a substantial effect on the corresponding Dice Similarity Coefficient (DSC) metric value (similar for the Intersection over Union (IoU)). This pitfall is also relevant for other overlap-based metrics such as the centerline Dice Similarity Coefficient (cDice), and localization criteria such as Box/Approx/Mask IoU and Intersection over Reference (IoR). (b) **Complex structure shapes.** Common overlap-based metrics (here: DSC) are unaware of complex structure shapes and treat *Predictions 1* and *2* equally. The cDice uncovers the fact that *Prediction 1* misses the fine-granular branches of the reference and favors *Prediction 2*, which focuses on the center line of the object. This pitfall is also relevant for other overlap-based such as metrics IoU and pixel-level F_β Score as well as localization criteria such as Box/Approx/Mask IoU, Center Distance, Mask IoU > 0, Point inside Mask/Box/Approx, and IoR.



Extended Data Fig. 2. [P2.4] **Disregard of the properties of the algorithm output.** (a) **Possibility of overlapping predictions.** If multiple structures of the same type can be seen within the same image (here: reference objects R1 and R2), it is generally advisable to phrase the problem as instance segmentation (InS; right) rather than semantic segmentation (SemS; left). This way, issues with boundary-based metrics resulting from comparing a given structure boundary to the boundary of the wrong instance in the reference can be avoided. In the provided example, the distance of the red boundary pixel to the reference, as measured by a boundary-based metric in SemS problems, would be zero, because different instances of the same structure cannot be distinguished. This problem is overcome by phrasing the problem as InS. In this case, (only) the boundary of the matched instance (here: R2) is considered for distance computation. (b) **Possibility of empty prediction or reference.** Each column represents a potential scenario for per-image validation of objects, categorized by whether True Positives (TPs), False Negatives (FNs), and False Positives (FPs) are present ($n > 0$) or not ($n = 0$) after matching/assignment. The sketches on the top showcase each scenario when setting " $n > 0$ " to " $n = 1$ ". For each scenario, Sensitivity, Positive Predictive Value (PPV), and the F₁ Score are calculated. Some scenarios yield undefined values (Not a Number (NaN)).

Extended Data Tab. 1. **Overview of pitfall sources for *image-level classification metrics*** ((a): counting metrics, (b): multi-threshold metrics) related to poor metric selection [P2]. A warning sign indicates a potential pitfall for the metric in the corresponding column, in case the property represented by the respective row holds true. Comprehensive illustrations of pitfalls are available in Suppl. Note 2. A comprehensive list of pitfalls is provided separately for each metrics in the metrics cheat sheets (Suppl. Note 3). Note that we only list sources of pitfalls relevant to the considered metrics. Other sources of pitfalls are neglected for this table.

(a) **Counting metrics.** Considered metrics: Accuracy (Fig. SN 3.40), Balanced Accuracy (BA) (Fig. SN 3.41), Expected Cost (EC) (Fig. SN 3.44), F_β Score (Fig. SN 3.45), Matthews Correlation Coefficient (MCC) (Fig. SN 3.48), Net Benefit (NB) (Fig. SN 3.49), Negative Predictive Value (NPV) (Fig. SN 3.50), Positive Likelihood Ratio (LR+) (Fig. SN 3.52), Positive Predictive Value (PPV) (Fig. SN 3.53), Sensitivity (Sens) (Fig. SN 3.54), Specificity (Spec) (Fig. SN 3.55), Weighted Cohen’s Kappa (WCK) (Fig. SN 3.56).

Source of potential pitfall	Accuracy	BA	EC	F_β Score	LR+	MCC	NB	PPV/ NPV	Sens/ Spec	WCK
Importance of confidence awareness	⚠*	⚠*	⚠*	⚠*	⚠*	⚠*	⚠*	⚠*	⚠*	⚠*
Importance of comparability across data sets	⚠ (Fig. SN 2.9)		⚠** (Fig. SN 2.9)	⚠ (Fig. SN 2.9)		⚠ (Fig. SN 2.9)	⚠ (Fig. SN 2.9)	⚠ (Fig. SN 2.9)		⚠ (Fig. SN 2.9)
Unequal severity of class confusions	⚠ (Fig. 4b)	⚠ (Fig. 4b)		⚠*** (Fig. 4b)	⚠ (Fig. 4b)	⚠ (Fig. 4b)		⚠ (Fig. 4b)	⚠ (Fig. 4b)	
Importance of cost-benefit analysis	⚠ (Fig. SN 2.11)	⚠ (Fig. SN 2.11)		⚠*** (Fig. SN 2.11)	⚠ (Fig. SN 2.11)	⚠ (Fig. SN 2.11)		⚠ (Fig. SN 2.11)	⚠ (Fig. SN 2.11)	
High class imbalance	⚠ (Figs. 5a, SN 2.17)	⚠ (Fig. 5a)	⚠** (Fig. 5a)		⚠ (Fig. 5a)		(Figs. 5a, SN 2.17)	NPV: ⚠ (Figs. 5a, SN 2.17)	⚠ (Sens: Fig. 5a, Spec: Figs. 5a, SN 2.17)	⚠ (Figs. 5a, SN 2.17)
Small test set size	⚠ (Fig. SN 2.18)	⚠ (Fig. SN 2.18)	⚠ (Fig. SN 2.18)	⚠ (Fig. SN 2.18)	⚠ (Fig. SN 2.18)	⚠ (Fig. SN 2.18)	⚠ (Fig. SN 2.18)	⚠ (Fig. SN 2.18)	⚠ (Fig. SN 2.18)	⚠ (Fig. SN 2.18)

* Discrimination metrics do not assess whether the predicted class scores reflect the confidence of the classifier. This is typically achieved with additional calibration metrics, which come with their own pitfalls (see Figs. SN 2.8 and SN 2.24, Extended Data Fig. 1b and the metric profiles in Suppl. Note 3.2).

** The weights in EC can be adjusted to avoid this pitfall.

*** The hyperparameter β can be used as a penalty for class confusions in the binary case. This property is not applicable to multi-class problems.

(b) **Multi-threshold metrics.** Considered metrics: Area under the Receiver Operating Characteristic Curve (AUROC) (Fig. SN 3.57) and Average Precision (AP) (Fig. SN 3.58).

Source of potential pitfall	AP	AUROC
Importance of confidence awareness	⚠*	⚠*
Importance of comparability across data sets	⚠ (Fig. SN 2.9)	
High class imbalance		⚠ (Fig. 5a)
Small test set size	⚠ (Fig. SN 2.18)	⚠ (Fig. SN 2.18)
Lack of predicted class scores	⚠ (Fig. SN 2.22)	⚠ (Fig. SN 2.22)

* Discrimination metrics do not assess whether the predicted class scores reflect the confidence of the classifier. This is typically achieved with additional calibration metrics, which come with their own pitfalls (see Figs. SN 2.8 and SN 2.24, Extended Data Fig. 1b and the metric profiles in Suppl. Note 3.2).

Extended Data Tab. 2. **Overview of pitfall sources for *semantic segmentation metrics*** ((a): overlap-based metrics, (b): boundary-based metrics) related to poor metric selection [P2]. A warning sign indicates a potential pitfall for the metric in the corresponding column, in case the property represented by the respective row holds true. Comprehensive illustrations of pitfalls are available in Suppl. Note 2. A comprehensive list of pitfalls is provided separately for each metrics in the metrics cheat sheets (Suppl. Note 3). Note that we only list sources of pitfalls relevant to the considered metrics. Other sources of pitfalls are neglected for this table.

(a) **Overlap-based metrics.** Considered metrics: Considered metrics: centerline Dice Similarity Coefficient (cDice) (Fig. SN 3.42), Dice Similarity Coefficient (DSC) (Fig. SN 3.43), F_β Score (Fig. SN 3.45), Intersection over Union (IoU) (Fig. SN 3.47).

Source of potential pitfall	cDice	DSC/IoU	F_β Score
Importance of structure boundaries	⚠ (Fig. 4a)	⚠ (Fig. 4a)	⚠ (Fig. 4a)
Importance of structure center(line)		⚠ (Fig. SN 2.7, Extended Data Fig. 1b)	⚠ (Fig. SN 2.7, Extended Data Fig. 1b)
Unequal severity of class confusions	⚠ (Fig. SN 2.10)	⚠ (Fig. SN 2.10)	⚠ (Fig. SN 2.12, Extended Data Fig. 1a)
Small structure sizes	⚠ (Fig. SN 2.12, Extended Data Fig. 1a)	⚠ (Fig. SN 2.12, Extended Data Fig. 1a)	⚠ (Fig. SN 2.12, Extended Data Fig. 1a)
High variability of structure sizes	⚠ (Fig. SN 2.13)	⚠ (Fig. SN 2.13)	⚠ (Fig. SN 2.13)
Complex structure shapes		⚠ (Fig. SN 2.14)	⚠ (Fig. SN 2.14)
Occurrence of overlapping or touching structures	⚠ (Fig. SN 2.15)	⚠ (Fig. SN 2.15)	⚠ (Fig. SN 2.15)
Imperfect reference standard		⚠ (Fig. SN 2.19)	⚠ (Fig. SN 2.19)
Occurrence of cases with an empty reference	⚠ (Fig. SN 2.20)	⚠ (Fig. SN 2.20)	⚠ (Fig. SN 2.20)
Possibility of empty prediction	⚠ (Fig. SN 2.20)	⚠ (Fig. SN 2.20)	⚠ (Fig. SN 2.20)
Possibility of overlapping predictions	⚠ (Fig. SN 2.21, Extended Data Fig. 2a)	⚠ (Fig. SN 2.21, Extended Data Fig. 2a)	⚠ (Fig. SN 2.21, Extended Data Fig. 2a)

(b) **Boundary-based metrics.** Considered metrics: Average Symmetric Surface Distance (ASSD) (Fig. SN 3.60), Boundary Intersection over Union (Boundary IoU) (Fig. SN 3.61), Hausdorff Distance (HD) (Fig. SN 3.62), Hausdorff Distance 95th Percentile (HD95) (Fig. SN 3.65), Mean Average Surface Distance (MASD) (Fig. SN 3.63), Normalized Surface Distance (NSD) (Fig. SN 3.64).

Source of potential pitfall	ASSD	Boundary IoU	HD	HD95	MASD	NSD
Importance of structure volume	⚠ (Fig. SN 2.6)	⚠ (Fig. SN 2.6)	⚠ (Fig. SN 2.6)	⚠ (Fig. SN 2.6)	⚠ (Fig. SN 2.6)	⚠ (Fig. SN 2.6)
Importance of structure center(line)	⚠ (Fig. SN 2.7, Extended Data Fig. 1b)	⚠ (Fig. SN 2.7, Extended Data Fig. 1b)	⚠ (Fig. SN 2.7, Extended Data Fig. 1b)	⚠ (Fig. SN 2.7, Extended Data Fig. 1b)	⚠ (Fig. SN 2.7, Extended Data Fig. 1b)	⚠ (Fig. SN 2.7, Extended Data Fig. 1b)
Occurrence of overlapping or touching structures	⚠ (Fig. SN 2.15)	⚠ (Fig. SN 2.15)	⚠ (Fig. SN 2.15)	⚠ (Fig. SN 2.15)	⚠ (Fig. SN 2.15)	⚠ (Fig. SN 2.15)
Imperfect reference standard	⚠ (Figs. 5c, SN 2.17)	⚠ (Figs. 5c, SN 2.17)	⚠ (Figs. 5c, SN 2.17)	⚠ (Figs. 5c*, SN 2.17)	⚠ (Figs. 5c, SN 2.17)	⚠ (Figs. 5c, SN 2.17)
Occurrence of cases with an empty reference	⚠ (Fig. SN 2.20)	⚠ (Fig. SN 2.20)	⚠ (Fig. SN 2.20)	⚠ (Fig. SN 2.20)	⚠ (Fig. SN 2.20)	⚠ (Fig. SN 2.20)
Possibility of empty prediction	⚠ (Fig. SN 2.20)	⚠ (Fig. SN 2.20)	⚠ (Fig. SN 2.20)	⚠ (Fig. SN 2.20)	⚠ (Fig. SN 2.20)	⚠ (Fig. SN 2.20)
Possibility of overlapping predictions	⚠ (Fig. SN 2.21, Extended Data Fig. 2a)	⚠ (Fig. SN 2.21, Extended Data Fig. 2a)	⚠ (Fig. SN 2.21, Extended Data Fig. 2a)	⚠ (Fig. SN 2.21, Extended Data Fig. 2a)	⚠ (Fig. SN 2.21, Extended Data Fig. 2a)	⚠ (Fig. SN 2.21, Extended Data Fig. 2a)

* Can be mitigated by the choice of the percentile.

Extended Data Tab. 3. **Overview of sources of pitfalls for *object detection metrics*** ((a): detection metrics, (b): localization criteria) related to poor metric selection [P2]. A warning sign indicates a potential pitfall for the metric in the corresponding column, in case the property represented by the respective row holds true. Comprehensive illustrations of pitfalls are available in Suppl. Note 2. A comprehensive list of pitfalls is provided separately for each metrics in the metrics cheat sheets (Suppl. Note 3). Note that we only list sources of pitfalls relevant to the considered metrics. Other sources of pitfalls are neglected for this table.

(a) **Detection metrics.** Considered counting metrics: F_β Score (Fig. SN 3.45), Positive Predictive Value (PPV) (Fig. SN 3.53), Sensitivity (Sens) (Fig. SN 3.54). Considered multi-threshold metrics: Average Precision (AP) (Fig. SN 3.58) and Free-Response Receiver Operating Characteristic (FROC) (Fig. SN 3.59).

Source of potential pitfall	F_β Score	PPV	Sens	AP	FROC Score
Unequal severity of class confusions	▲* (Fig. 4b)	▲ (Fig. 4b)	▲ (Fig. 4b)	▲ (Fig. 4b)	▲ (Fig. 4b)
High class imbalance			▲ (Fig. 5a)		
Small test set size	▲ (Fig. SN 2.18)	▲ (Fig. SN 2.18)	▲ (Fig. SN 2.18)	▲ (Fig. SN 2.18)	▲ (Fig. SN 2.18)
Occurrence of cases with an empty reference	▲ (Fig. SN 2.20, Extended Data Fig. 2b)	▲ (Fig. SN 2.20, Extended Data Fig. 2b)	▲ (Fig. SN 2.20, Extended Data Fig. 2b)	▲ (Fig. SN 2.20, Extended Data Fig. 2b)	▲ (Fig. SN 2.20, Extended Data Fig. 2b)
Possibility of empty prediction	▲ (Fig. SN 2.20, Extended Data Fig. 2b)	▲ (Fig. SN 2.20, Extended Data Fig. 2b)	▲ (Fig. SN 2.20, Extended Data Fig. 2b)	▲ (Fig. SN 2.20, Extended Data Fig. 2b)	▲ (Fig. SN 2.20, Extended Data Fig. 2b)
Lack of predicted class scores				▲ (Fig. SN 2.22)	▲ (Fig. SN 2.22)

* The hyperparameter β can be used as a penalty for class confusions in the binary case. This property is not applicable to multi-class problems.

(b) **Localization criteria.** Considered localization criteria: Box/Approx Intersection over Union (IoU) (Fig. SN 3.76), Center Distance (Fig. SN 3.74), Mask IoU > 0 (Fig. SN 3.77), and Point inside Mask/ Box/ Approx (Fig. SN 3.78).

Source of potential pitfall	Box/ Approx IoU	Center Distance	Mask IoU > 0	Point inside Mask/ Box/ Approx
Importance of structure boundaries	▲ (Fig. 5a)	▲ (Fig. 5a)	▲ (Fig. 5a)	▲ (Fig. 5a)
Importance of structure volume		▲ (Fig. SN 2.6)	▲ (Fig. SN 2.6)	▲ (Fig. SN 2.6)
Importance of structure center(line)	▲ (Fig. SN 2.7, Extended Data Fig. 1b)		▲ (Fig. SN 2.7, Extended Data Fig. 1b)	▲ (Fig. SN 2.7, Extended Data Fig. 1b)
Unequal severity of class confusions	▲ (Fig. SN 2.10)	▲ (Fig. SN 2.10)*	▲ (Fig. SN 2.10)	▲ (Fig. SN 2.10)*
Small structure sizes	▲ (Fig. SN 2.12, Extended Data Fig. 1a)			
Complex structure shapes	▲ (Figs. SN 2.13, SN 2.16)	▲ (Fig. SN 2.13)	▲ (Fig. SN 2.13)	▲ (Fig. SN 2.13)
Occurrence of disconnected structures	▲ (Fig. SN 2.16)			Point inside Box: ▲ (Fig. SN 2.16)
Imperfect reference standard	▲ (Fig. 5c)			

* Criterion implies point prediction, thus overlap assessment is not applicable.

Extended Data Tab. 4. **Overview of sources of pitfalls for *instance segmentation metrics (Part 1)*** ((a): detection metrics, (b): localization criteria) related to poor metric selection [P2]. A warning sign indicates a potential pitfall for the metric in the corresponding column, in case the property represented by the respective row holds true. Comprehensive illustrations of pitfalls are available in Suppl. Note 2. A comprehensive list of pitfalls is provided separately for each metrics in the metrics cheat sheets (Suppl. Note 3). Note that we only list sources of pitfalls relevant to the considered metrics. Other sources of pitfalls are neglected for this table.

(a) **Detection metrics.** Considered counting metrics: F_β Score (Fig. SN 3.45), Positive Predictive Value (PPV) (Fig. SN 3.53), Panoptic Quality (PQ) (Fig. SN 3.51), Sensitivity (Sens) (Fig. SN 3.54). Considered multi-threshold metrics: Average Precision (AP) (Fig. SN 3.58) and Free-Response Receiver Operating Characteristic (FROC) (Fig. SN 3.59).

Source of potential pitfall	F_β Score	PPV	PQ	Sens	AP	FROC Score
Unequal severity of class confusions	⚠* (Fig. 4b)	⚠ (Fig. 4b)	⚠ (Fig. 4b)	⚠ (Fig. 4b)		
High class imbalance				⚠ (Fig. 5a)		
Small test set size	⚠ (Fig. SN 2.18)	⚠ (Fig. SN 2.18)	⚠ (Fig. SN 2.18)	⚠ (Fig. SN 2.18)	⚠ (Fig. SN 2.18)	⚠ (Fig. SN 2.18)
Lack of predicted class scores					⚠ (Fig. SN 2.22)	⚠ (Fig. SN 2.22)

* The hyperparameter β can be used as a penalty for class confusions in the binary case. This property is not applicable to multi-class problems.

(b) **Localization criteria.** Considered localization criteria: Boundary Intersection over Union (Boundary IoU) (Fig. SN 3.73), Intersection over Reference (IoR) (Fig. SN 3.75), Mask IoU (Fig. SN 3.47).

Source of potential pitfall	Boundary IoU	IoR	Mask IoU
Importance of structure boundaries		⚠ (Fig. 4a)	⚠ (Fig. 4a)
Importance of structure volume	⚠ (Fig. SN 2.4)		
Importance of structure center(line)	⚠ (Fig. SN 2.7, Extended Data Fig. 1b)	⚠ (Fig. SN 2.7, Extended Data Fig. 1b)	⚠ (Fig. SN 2.7, Extended Data Fig. 1b)
Unequal severity of class confusions	⚠ (Fig. SN 2.10)	⚠ (Fig. SN 2.10)	⚠ (Fig. SN 2.10)
Small structure sizes		⚠ (Fig. SN 2.12, Extended Data Fig. 1a)	⚠ (Fig. SN 2.12, Extended Data Fig. 1a)
Complex structure shapes		⚠ (Fig. SN 2.14)	⚠ (Fig. SN 2.12)
Imperfect reference standard	⚠ (Fig. SN 2.19)	⚠ (Fig. SN 2.19)	⚠ (Fig. SN 2.19)

Extended Data Tab. 5. **Overview of sources of pitfalls for *instance segmentation metrics (Part 2)*** ((a) per instance segmentation overlap-based metrics, (b) per instance segmentation boundary-based metrics) related to poor metric selection [P2]. A warning sign indicates a potential pitfall for the metric in the corresponding column, in case the property represented by the respective row holds true. Comprehensive illustrations of pitfalls are available in Suppl. Note 2. Note that we only list sources of pitfalls relevant to the considered metrics. Other sources of pitfalls are neglected for this table.

(a) **Per instance segmentation overlap-based metrics.** Considered metrics: centerline Dice Similarity Coefficient (cDice) (Fig. SN 3.42), Dice Similarity Coefficient (DSC) (Fig. SN 3.43), F_β Score (Fig. SN 3.45), Intersection over Union (IoU) (Fig. SN 3.47).

Source of potential pitfall	cDice	DSC/IoU	F_β Score
Importance of structure boundaries	⚠ (Fig. 4a)	⚠ (Fig. 4a)	⚠ (Fig. 4a)
Importance of structure center(line)		⚠ (Fig. SN 2.7, Extended Data Fig. 1b)	⚠ (Fig. SN 2.7, Extended Data Fig. 1b)
Unequal severity of class confusions	⚠ (Fig. SN 2.10)	⚠ (Fig. SN 2.10)	
Small structure sizes	⚠ (Fig. SN 2.12, Extended Data Fig. 1a)	⚠ (Fig. SN 2.12, Extended Data Fig. 1a)	⚠ (Fig. SN 2.12, Extended Data Fig. 1a)
Complex structure shapes		⚠ (Fig. SN 2.14)	⚠ (Fig. SN 2.14)
Imperfect reference standard		⚠ (Fig. SN 2.19)	⚠ (Fig. SN 2.19)

(b) **Per instance segmentation boundary-based metrics.** Considered metrics: Average Symmetric Surface Distance (ASSD) (Fig. SN 3.60), Boundary Intersection over Union (Boundary IoU) (Fig. SN 3.61), Hausdorff Distance (HD) (Fig. SN 3.62), Hausdorff Distance 95th Percentile (HD95) (Fig. SN 3.65), Mean Average Surface Distance (MASD) (Fig. SN 3.63), Normalized Surface Distance (NSD) (Fig. SN 3.64).

Source of potential pitfall	ASSD	Boundary IoU	HD	HD95	MASD	NSD
Importance of structure volume	⚠ (Fig. SN 2.6)	⚠ (Fig. SN 2.6)	⚠ (Fig. SN 2.6)	⚠ (Fig. SN 2.6)	⚠ (Fig. SN 2.6)	⚠ (Fig. SN 2.6)
Importance of structure center(line)	⚠ (Fig. SN 2.7, Extended Data Fig. 1b)	⚠ (Fig. SN 2.7, Extended Data Fig. 1b)	⚠ (Fig. SN 2.7, Extended Data Fig. 1b)	⚠ (Fig. SN 2.7, Extended Data Fig. 1b)	⚠ (Fig. SN 2.7, Extended Data Fig. 1b)	⚠ (Fig. SN 2.7, Extended Data Fig. 1b)
Imperfect reference standard	⚠ (Figs. 5c, SN 2.19)	⚠ (Figs. 5c, SN 2.19)	⚠ (Figs. 5c, SN 2.19)	⚠ (Figs. 5c*, SN 2.19)	⚠ (Figs. 5c, SN 2.19)	

* Can be mitigated by the choice of the percentile.

CODE AVAILABILITY STATEMENT

We provide reference implementations for all *Metrics Reloaded* metrics within the MONAI open-source framework. They are accessible at <https://github.com/Project-MONAI/MetricsReloaded>.

ACKNOWLEDGEMENTS

This work was initiated by the Helmholtz Association of German Research Centers in the scope of the Helmholtz Imaging Incubator (HI), the MICCAI Special Interest Group for biomedical image analysis challenges, and the benchmarking working group of the MONAI initiative. It has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement No. [101002198], NEURAL SPICING) and the Surgical Oncology Program of the National Center for Tumor Diseases (NCT) Heidelberg. It was further supported in part by the Intramural Research Program of the National Institutes of Health Clinical Center as well as by the National Cancer Institute (NCI) and the National Institute of Neurological Disorders and Stroke (NINDS) of the National Institutes of Health (NIH), under award numbers NCI:U01CA242871 and NINDS:R01NS042645. The content of this publication is solely the responsibility of the authors and does not represent the official

views of the NIH. T.A. acknowledges the Canada Institute for Advanced Research (CIFAR) AI Chairs program, the Natural Sciences and Engineering Research Council of Canada. F.B. was co-funded by the European Union (ERC, TAPO, 101088594). Views and opinions expressed are however those of the authors only and do not necessarily reflect those of the European Union or the European Research Council. Neither the European Union nor the granting authority can be held responsible for them. M.J.C. acknowledges funding from Wellcome/EPSRC Centre for Medical Engineering (WT203148/Z/16/Z), the Wellcome Trust (WT213038/Z/18/Z), and the InnovateUK funded London AI Centre for Value-Based Healthcare. J.C. is supported by the Federal Ministry of Education and Research (BMBF) under the funding reference 161L0272. V.C. acknowledges funding from NovoNordisk Foundation (NNF21OC0068816) and Independent Research Council Denmark (1134-00017B). B.A.C. was supported by NIH grant P41 GM135019 and grant 2020-225720 from the Chan Zuckerberg Initiative DAF, an advised fund of the Silicon Valley Community Foundation. G.S.C. was supported by Cancer Research UK (programme grant: C49297/A27294). M.M.H. is supported by the Natural Sciences and Engineering Research Council of Canada (RGPIN-2022-05134). A.Ka. is supported by French State Funds managed by the “Agence Nationale de la Recherche (ANR)” - “Investissements d’Avenir” (Investments for the Future), Grant ANR-10-IAHU-02 (IHU Strasbourg). M.K. was funded by the Ministry of Education, Youth and Sports of the Czech Republic (Project LM2018129). Ta.K. was supported in part by 4UH3-CA225021-03, 1U24CA180924-01A1, 3U24CA215109-02, and 1UG3-CA225-021-01 grants from the National Institutes of Health. G.L. receives research funding from the Dutch Research Council, the Dutch Cancer Association, HealthHolland, the European Research Council, the European Union, and the Innovative Medicine Initiative. S.M.R. wishes to acknowledge the Allen Institute for Cell Science founder Paul G. Allen for his vision, encouragement and support. M.R. is supported by Innosuisse grant number 31274.1 and Swiss National Science Foundation Grant Number 205320_212939. C.H.S. is supported by an Alzheimer’s Society Junior Fellowship (AS-JF-17-011). R.M.S. is supported by the Intramural Research Program of the NIH Clinical Center. A.T. acknowledges support from Academy of Finland (Profi6 336449 funding program), University of Oulu strategic funding, Finnish Foundation for Cardiovascular Research, Wellbeing Services County of North Ostrobothnia (VTR project K62716), and Terttu foundation. S.A.T. acknowledges the support of Canon Medical and the Royal Academy of Engineering and the Research Chairs and Senior Research Fellowships scheme (grant RCSRF1819\8\25). B.V.C. was supported by Research Foundation Flanders (FWO grant G097322N) and Internal Funds KU Leuven (grant C24M/20/064).

We would like to thank Peter Bankhead, Gary S. Collins, Robert Haase, Fred Hamprecht, Alan Karthikesalingam, Hannes Kenngott, Peter Mattson, David Moher, Bram Stieltjes, and Manuel Wiesenfarth for the fruitful discussions on this work.

We would like to thank Sandy Engelhardt, Sven Koehler, M. Alican Noyan, Gorkem Polat, Hassan Rivaz, Julian Schroeter, Anindo Saha, Lalith Sharan, Peter Hirsch, and Matheus Viana for suggesting additional illustrations that can be found in [71].

COMPETING INTERESTS

The authors declare the following competing interests: F.B. is an employee of Siemens AG (Munich, Germany). B.v.G. is a shareholder of Thirona (Nijmegen, NL). B.G. is an employee of HeartFlow Inc (California, USA) and Kheiron Medical Technologies Ltd (London, UK). M.M.H. received an Nvidia GPU Grant. Th. K. is an employee of Lunit (Seoul, South Korea). G.L. is on the advisory board of Canon Healthcare IT (Minnetonka, USA) and is a shareholder of Aiosyn BV (Nijmegen, NL). Na.R. is the founder and CSO of Histofy (New York, USA). Ni.R. is an employee of Nvidia GmbH (Munich,

Germany). J.S.-R. reports funding from GSK (Heidelberg, Germany), Pfizer (New York, USA) and Sanofi (Paris, France) and fees from Traveo Therapeutics (California, USA), Stadapharm (Bad Vilbel, Germany), Astex Therapeutics (Cambridge, UK), Pfizer (New York, USA), and Grunenthal (Aachen, Germany). R.M.S. receives patent royalties from iCAD (New Hampshire, USA), ScanMed (Nebraska, USA), Philips (Amsterdam, NL), Translation Holdings (Alabama, USA) and PingAn (Shenzhen, China); his lab received research support from PingAn through a Cooperative Research and Development Agreement. S.A.T. receives financial support from Canon Medical Research Europe (Edinburgh, Scotland).

REFERENCES

- [1] Saeid Asgari Taghanaki, Kumar Abhishek, Joseph Paul Cohen, Julien Cohen-Adad, and Ghassan Hamarneh. Deep semantic segmentation of natural and medical images: a review. *Artificial Intelligence Review*, 54(1):137–178, 2021.
- [2] John Attia. Moving beyond sensitivity and specificity: using likelihood ratios to help interpret diagnostic tests. *Australian prescriber*, 26(5):111–113, 2003.
- [3] Min Bai and Raquel Urtasun. Deep watershed transform for instance segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5221–5229, 2017.
- [4] D Bamira and MH Picard. Imaging: Echocardiology—assessment of cardiac structure and function. *Elsevier*, 2018.
- [5] Andriy I Bandos, Howard E Rockette, Tao Song, and David Gur. Area under the free-response roc curve (froc) and a related summary index. *Biometrics*, 65(1):247–256, 2009.
- [6] Miroslav Beneš and Barbara Zitová. Performance evaluation of image segmentation algorithms on microscopic image data. *Journal of microscopy*, 257(1):65–85, 2015.
- [7] Patrick Bilic, Patrick Christ, Hongwei Bran Li, Eugene Vorontsov, Avi Ben-Cohen, Georgios Kaissis, Adi Szeskin, Colin Jacobs, Gabriel Efrain Humpire Mamani, Gabriel Chartrand, et al. The liver tumor segmentation benchmark (lits). *Medical Image Analysis*, 84:102680, 2023.
- [8] Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*, volume 4. Springer, 2006.
- [9] Bernice B Brown. Delphi process: a methodology used for the elicitation of opinions of experts. Technical report, Rand Corp Santa Monica CA, 1968.
- [10] Chang Cao, Davide Chicco, and Michael M Hoffman. The mcc-f1 curve: a performance evaluation technique for binary classification. *arXiv preprint arXiv:2006.11278*, 2020.
- [11] Alberto Carbonell, Marcos De la Pena, Ricardo Flores, and Selma Gago. Effects of the trinucleotide preceding the self-cleavage site on eggplant latent viroid hammerheads: differences in co- and post-transcriptional self-cleavage may explain the lack of trinucleotide auc in most natural hammerheads. *Nucleic acids research*, 34(19):5613–5622, 2006.
- [12] Jianxu Chen, Liya Ding, Matheus P Viana, HyeonWoo Lee, M Filip Sluezwski, Benjamin Morris, Melissa C Hendershott, Ruian Yang, Irina A Mueller, and Susanne M Rafelski. The allen cell and structure segmenter: a new open source toolkit for segmenting 3d intracellular structures in fluorescence microscopy images. *BioRxiv*, page 491035, 2020.
- [13] Bowen Cheng, Ross Girshick, Piotr Dollár, Alexander C Berg, and Alexander Kirillov. Boundary iou: Improving object-centric image segmentation evaluation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15334–15342, 2021.
- [14] Davide Chicco and Giuseppe Jurman. The advantages of the matthews correlation coefficient (mcc) over f1 score and accuracy in binary classification evaluation. *BMC genomics*, 21(1):1–13, 2020.
- [15] Davide Chicco, Niklas Tötsch, and Giuseppe Jurman. The matthews correlation coefficient (mcc) is more reliable than balanced accuracy, bookmaker informedness, and markedness in two-class confusion matrix evaluation. *BioData mining*, 14(1):1–22, 2021. The manuscript addresses the challenge of evaluating binary classifications. It compares MCC to other metrics, explaining their mathematical relationships and providing use cases where MCC offers more informative results.
- [16] Nancy Chinchor. Muc-4 evaluation metrics. In *Proceedings of the 4th Conference on Message Understanding, MUC4 '92*, page 22–29, USA, 1992. Association for Computational Linguistics. ISBN 1558602739. doi: 10.3115/1072064.1072067. URL <https://doi.org/10.3115/1072064.1072067>.
- [17] Jacob Cohen. A coefficient of agreement for nominal scales. *Educational and psychological measurement*, 20(1):37–46, 1960.
- [18] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Scharwächter, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset. In *CVPR Workshop on The Future of Datasets in Vision*, 2015.

- [19] Paulo Correia and Fernando Pereira. Video object relevance metrics for overall segmentation quality evaluation. *EURASIP Journal on Advances in Signal Processing*, 2006:1–11, 2006.
- [20] George Cybenko, Dianne P O’Leary, and Jorma Rissanen. *The Mathematics of Information Coding, Extraction and Distribution*, volume 107. Springer Science & Business Media, 1998.
- [21] Jesse Davis and Mark Goadrich. The relationship between precision-recall and roc curves. In *Proceedings of the 23rd international conference on Machine learning*, pages 233–240, 2006.
- [22] Bert De Brabandere, Davy Neven, and Luc Van Gool. Semantic instance segmentation with a discriminative loss function. *arXiv preprint arXiv:1708.02551*, 2017.
- [23] Jeffrey De Fauw, Joseph R LedSAM, Bernardino Romera-Paredes, Stanislav Nikolov, Nenad Tomasev, Sam Blackwell, Harry Askham, Xavier Glorot, Brendan O’Donoghue, Daniel Visentin, et al. Clinically applicable deep learning for diagnosis and referral in retinal disease. *Nature medicine*, 24(9):1342–1350, 2018.
- [24] Rosario Delgado and Xavier-Andoni Tibau. Why cohen’s kappa should be avoided as performance measure in classification. *PLoS one*, 14(9):e0222916, 2019.
- [25] Elizabeth R DeLong, David M DeLong, and Daniel L Clarke-Pearson. Comparing the areas under two or more correlated receiver operating characteristic curves: a nonparametric approach. *Biometrics*, pages 837–845, 1988.
- [26] Antonio Di Sabatino and Gino Roberto Corazza. Nonceliac gluten sensitivity: sense or sensibility?, 2012.
- [27] Lee R Dice. Measures of the amount of ecologic association between species. *Ecology*, 26(3):297–302, 1945.
- [28] Mark Everingham, Andrew Zisserman, Christopher KI Williams, Luc Van Gool, Moray Allan, Christopher M Bishop, Olivier Chapelle, Navneet Dalal, Thomas Deselaers, Gyuri Dorkó, et al. The 2005 pascal visual object classes challenge. In *Machine Learning Challenges. Evaluating Predictive Uncertainty, Visual Object Classification, and Recognising Tectual Entailment: First PASCAL Machine Learning Challenges Workshop, MLCW 2005, Southampton, UK, April 11-13, 2005, Revised Selected Papers*, pages 117–176. Springer, 2006.
- [29] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2):303–338, 2010.
- [30] Mark Everingham, SM Ali Eslami, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes challenge: A retrospective. *International journal of computer vision*, 111(1):98–136, 2015.
- [31] Luciana Ferrer. Analysis and comparison of classification metrics. *arXiv preprint arXiv:2209.05355*, 2022.
- [32] Naiyu Gao, Yanhu Shan, Yupei Wang, Xin Zhao, Yinan Yu, Ming Yang, and Kaiqi Huang. Ssap: Single-shot instance segmentation with affinity pyramid. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 642–651, 2019.
- [33] Tilmann Gneiting and Adrian E Raftery. Strictly proper scoring rules, prediction, and estimation. *Journal of the American statistical Association*, 102(477):359–378, 2007.
- [34] Mark J Gooding, Annamaria J Smith, Maira Tariq, Paul Aljabar, Devis Peressutti, Judith van der Stoep, Bart Reymen, Daisy Emans, Djoya Hattu, Judith van Loon, et al. Comparative evaluation of autocontouring in clinical practice: a practical method using the turing test. *Medical physics*, 45(11):5105–5115, 2018.
- [35] Mark J Gooding, Djamel Boukerroui, Eliana Vasquez Osorio, René Monshouwer, and Ellen Brunenberg. Multicenter comparison of measures for quantitative evaluation of contouring in radiotherapy. *Physics and Imaging in Radiation Oncology*, 24:152–158, 2022.
- [36] Margherita Grandini, Enrico Bagli, and Giorgio Visani. Metrics for multi-class classification: an overview. *arXiv preprint arXiv:2008.05756*, 2020.
- [37] Sebastian Gruber and Florian Buettner. Trustworthy deep learning via proper calibration errors: A unifying approach for quantifying the reliability of predictive uncertainty. *arXiv preprint arXiv:2203.07835*, 2022.
- [38] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On Calibration of Modern Neural Networks. *ICML*, page 10, 2017.
- [39] Metin N Gurcan, Anant Madabhushi, and Nasir Rajpoot. Pattern recognition in histopathological images: An icpr 2010 contest. In *International Conference on Pattern Recognition*, pages 226–234. Springer, 2010.
- [40] James A Hanley and Barbara J McNeil. The meaning and use of the area under a receiver operating characteristic (roc) curve. *Radiology*, 143(1):29–36, 1982.
- [41] Trevor Hastie, Robert Tibshirani, Jerome H Friedman, and Jerome H Friedman. *The elements of statistical learning: data mining, inference, and prediction*, volume 2. Springer, 2009.
- [42] Peter Hirsch, Lisa Mais, and Dagmar Kainmueller. Patchperpix for instance segmentation. *arXiv preprint arXiv:2001.07626*, 2020.
- [43] Katrin Honauer, Lena Maier-Hein, and Daniel Kondermann. The hci stereo metrics: Geometry-aware performance analysis of stereo algorithms. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2120–2128, 2015.
- [44] Daniel P Huttenlocher, Gregory A Klanderma, and William J Rucklidge. Comparing images using the hausdorff distance. *IEEE Transactions on pattern analysis and machine intelligence*, 15(9):850–863, 1993.

- [45] Paul Jaccard. The distribution of the flora in the alpine zone. 1. *New phytologist*, 11(2):37–50, 1912.
- [46] Kaggle. Satorius Cell Instance Segmentation 2021. <https://www.kaggle.com/c/sartorius-cell-instance-segmentation>, 2021. [Online; accessed 25-April-2022].
- [47] Alexander Kirillov, Kaiming He, Ross Girshick, Carsten Rother, and Piotr Dollár. Panoptic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9404–9413, 2019.
- [48] Florian Kofler, Ivan Ezhov, Fabian Isensee, Christoph Berger, Maximilian Korner, Johannes Paetzold, Hongwei Li, Suprosanna Shit, Richard McKinley, Spyridon Bakas, et al. Are we using appropriate segmentation metrics? Identifying correlates of human expert perception for CNN training beyond rolling the DICE coefficient. *arXiv preprint arXiv:2103.06205v1*, 2021.
- [49] Ender Konukoglu, Ben Glocker, Dong Hye Ye, Antonio Criminisi, and Kilian M Pohl. Discriminative segmentation-based evaluation through shape dissimilarity. *IEEE transactions on medical imaging*, 31(12):2278–2289, 2012.
- [50] Jonathan Krause, Varun Gulshan, Ehsan Rahimy, Peter Karth, Kasumi Widner, Greg S Corrado, Lily Peng, and Dale R Webster. Grader variability and the importance of reference standards for evaluating machine learning models for diabetic retinopathy. *Ophthalmology*, 125(8):1264–1272, 2018.
- [51] Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955.
- [52] Victor Kulikov and Victor Lempitsky. Instance segmentation of biological images using harmonic embeddings. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3843–3851, 2020.
- [53] Meelis Kull, Miquel Perello Nieto, Markus Kängsepp, Telmo Silva Filho, Hao Song, and Peter Flach. Beyond temperature scaling: Obtaining well-calibrated multi-class probabilities with dirichlet calibration. *Advances in neural information processing systems*, 32, 2019.
- [54] Ananya Kumar, Percy S Liang, and Tengyu Ma. Verified uncertainty calibration. *Advances in Neural Information Processing Systems*, 32, 2019.
- [55] Jochen K Lennerz, Ursula Green, Drew FK Williamson, and Faisal Mahmood. A unifying force for the realization of medical ai. *npj Digital Medicine*, 5(1):1–3, 2022.
- [56] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.
- [57] Lena Maier-Hein, Matthias Eisenmann, Annika Reinke, Sinan Onogur, Marko Stankovic, Patrick Scholz, Tal Arbel, Hrvoje Bogunovic, Andrew P Bradley, Aaron Carass, et al. Why rankings of biomedical image analysis competitions should be interpreted with care. *Nature communications*, 9(1):1–13, 2018. With this comprehensive analysis of biomedical image analysis competitions (challenges), the authors initiated a shift in how such challenges are designed, performed, and reported in the biomedical domain. Its concepts and guidelines have been adopted by reputed organizations such as MICCAI.
- [58] Lena Maier-Hein, Annika Reinke, Evangelia Christodoulou, Ben Glocker, Patrick Godau, Fabian Isensee, Jens Kleesiek, Michal Kozubek, Mauricio Reyes, Michael A Riegler, et al. Metrics reloaded: Pitfalls and recommendations for image analysis validation. *arXiv preprint arXiv:2206.01653*, 2022.
- [59] Ran Margolin, Lihi Zelnik-Manor, and Ayellet Tal. How to evaluate foreground maps? In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 248–255, 2014.
- [60] Martin Maška, Vladimír Ulman, David Svoboda, Pavel Matula, Petr Matula, Cristina Ederra, Ainhoa Urbiola, Tomás España, Subramanian Venkatesan, Deepak MW Balak, et al. A benchmark for comparison of cell tracking algorithms. *Bioinformatics*, 30(11):1609–1617, 2014.
- [61] Brian W Matthews. Comparison of the predicted and observed secondary structure of t4 phage lysozyme. *Biochimica et Biophysica Acta (BBA)-Protein Structure*, 405(2):442–451, 1975.
- [62] John Muschelli. Roc and auc with a binary predictor: a potentially misleading metric. *Journal of classification*, 37(3): 696–708, 2020.
- [63] Mahdi Pakdaman Naeini, Gregory Cooper, and Milos Hauskrecht. Obtaining well calibrated probabilities using bayesian binning. In *Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015.
- [64] Ying-Hwey Nai, Bernice W Teo, Nadya L Tan, Sophie O’Doherty, Mary C Stephenson, Yee Liang Thian, Edmund Chiong, and Anthonin Reilhac. Comparison of metrics for the evaluation of medical segmentations using prostate mri dataset. *Computers in Biology and Medicine*, 134:104497, 2021.
- [65] Prashant Nasa, Ravi Jain, and Deven Juneja. Delphi methodology in healthcare research: how to decide its appropriateness. *World Journal of Methodology*, 11(4):116, 2021.
- [66] Stanislav Nikolov, Sam Blackwell, Alexei Zverovitch, Ruheena Mendes, Michelle Livne, Jeffrey De Fauw, Yojan Patel, Clemens Meyer, Harry Askham, Bernadino Romera-Paredes, et al. Clinically applicable segmentation of head and neck anatomy for radiotherapy: deep learning algorithm development and validation study. *Journal of Medical Internet Research*, 23(7):e26151, 2021.

- [67] Chawin Ounkomol, Sharmishta Seshamani, Mary M Maleckar, Forrest Collman, and Gregory R Johnson. Label-free prediction of three-dimensional fluorescence images from transmitted-light microscopy. *Nature methods*, 15(11):917–920, 2018.
- [68] Hung Viet Pham, Shangshu Qian, Jiannan Wang, Thibaud Lutellier, Jonathan Rosenthal, Lin Tan, Yaoliang Yu, and Nachiappan Nagappan. Problems and opportunities in training deep learning software systems: An analysis of variance. In *Proceedings of the 35th IEEE/ACM international conference on automated software engineering*, pages 771–783, 2020.
- [69] Teodora Popordanoska, Raphael Sayer, and Matthew B Blaschko. A consistent and differentiable lp canonical calibration error estimator. In *Advances in Neural Information Processing Systems*, 2022.
- [70] David Martin Ward Powers. The problem with kappa. In *Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics*, pages 345–355, 2012.
- [71] Annika Reinke, Matthias Eisenmann, Minu D Tizabi, Carole H Sudre, Tim Rädtsch, Michela Antonelli, Tal Arbel, Spyridon Bakas, M Jorge Cardoso, Veronika Cheplygina, Keyvan Farahani, Ben Glocker, Doreen Heckmann-Nötzel, Fabian Isensee, Pierre Jannin, Charles Kahn, Jens Kleesiek, Tahsin Kurc, Michal Kozubek, Bennett A Landman, Geert Litjens, Klaus Maier-Hein, Anne L Martel, Henning Müller, Jens Petersen, Mauricio Reyes, Nicola Rieke, Bram Stieltjes, Ronald M Summers, Sotirios A Tsafaris, Bram van Ginneken, Annette Kopp-Schneider, Paul Jäger, and Lena Maier-Hein. Common limitations of image processing metrics: A picture story. *arXiv preprint arXiv:2104.05642*, 2021.
- [72] Annika Reinke, Matthias Eisenmann, Minu D Tizabi, Carole H Sudre, Tim Rädtsch, Michela Antonelli, Tal Arbel, Spyridon Bakas, M Jorge Cardoso, Veronika Cheplygina, et al. Common limitations of image processing metrics: A picture story. *arXiv preprint arXiv:2104.05642*, 2021.
- [73] Brock Roberts, Amanda Haupt, Andrew Tucker, Tanya Grancharova, Joy Arakaki, Margaret A Fuqua, Angeliq Nelson, Caroline Hookway, Susan A Ludmann, Irina A Mueller, et al. Systematic gene tagging using crispr/cas9 in human stem cells to illuminate cell organization. *Molecular biology of the cell*, 28(21):2854–2874, 2017.
- [74] Azriel Rosenfeld and John L Pfaltz. Sequential operations in digital picture processing. *Journal of the ACM (JACM)*, 13(4):471–494, 1966.
- [75] Anindo Saha, Joeran Bosma, Jasper Linmans, Matin Hosseinzadeh, and Henkjan Huisman. Anatomical and diagnostic bayesian segmentation in prostate mri – should different clinical objectives mandate different loss functions? *arXiv preprint arXiv:2110.12889*, 2021.
- [76] Uwe Schmidt, Martin Weigert, Coleman Broaddus, and Gene Myers. Cell detection with star-convex polygons. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 265–273. Springer, 2018.
- [77] Suprosanna Shit, Johannes C Paetzold, Anjany Sekuboyina, Ivan Ezhov, Alexander Unger, Andrey Zhylyka, Josien PW Pluim, Ulrich Bauer, and Bjoern H Menze. cldice-a novel topology-preserving loss function for tubular structure segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16560–16569, 2021.
- [78] Carsen Stringer, Tim Wang, Michalis Michaelos, and Marius Pachitariu. Cellpose: a generalist algorithm for cellular segmentation. *Nature methods*, 18(1):100–106, 2021.
- [79] Abdel Aziz Taha and Allan Hanbury. Metrics for evaluating 3d medical image segmentation: analysis, selection, and tool. *BMC medical imaging*, 15(1):1–28, 2015. The paper discusses the importance of effective metrics for evaluating the accuracy of 3D medical image segmentation algorithms. The authors analyze existing metrics, propose a selection methodology, and develop a tool to aid researchers in choosing appropriate evaluation metrics based on the specific characteristics of the segmentation task.
- [80] Abdel Aziz Taha, Allan Hanbury, and Oscar A Jimenez del Toro. A formal method for selecting evaluation metrics for image segmentation. In *2014 IEEE international conference on image processing (ICIP)*, pages 932–936. IEEE, 2014.
- [81] Alaa Tharwat. Classification assessment methods. *Applied Computing and Informatics*, 2020.
- [82] Thuy Nuong Tran, Tim Adler, Amine Yamlihi, Evangelia Christodoulou, Patrick Godau, Annika Reinke, Minu Dietlinde Tizabi, Peter Sauer, Tillmann Persicke, Jörg Gerhard Albert, et al. Sources of performance variability in deep learning-based polyp detection. *arXiv preprint arXiv:2211.09708*, 2022.
- [83] Femke Vaassen, Colien Hazelaar, Ana Vaniqui, Mark Gooding, Brent van der Heyden, Richard Canters, and Wouter van Elmpt. Evaluation of measures for assessing time-saving of automatic organ-at-risk segmentation in radiotherapy. *Physics and Imaging in Radiation Oncology*, 13:1–6, 2020.
- [84] Juozas Vaicenavicius, David Widmann, Carl Andersson, Fredrik Lindsten, Jacob Roll, and Thomas Schön. Evaluating model calibration in classification. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 3459–3467. PMLR, 2019.
- [85] Bram Van Ginneken, Samuel G Armato III, Bartjan de Hoop, Saskia van Amelsvoort-van de Vorst, Thomas Duindam, Meindert Niemeijer, Keelin Murphy, Arnold Schilham, Alessandra Retico, Maria Evelina Fantacci, et al. Comparing and combining algorithms for computer-aided detection of pulmonary nodules in computed tomography scans: the anode09 study. *Medical image analysis*, 14(6):707–722, 2010.

- [86] C Van Rijsbergen. Information retrieval: theory and practice. In *Proceedings of the Joint IBM/University of Newcastle upon Tyne Seminar on Data Base Systems*, volume 79, 1979.
- [87] Matheus P Viana, Jianxu Chen, Theo A Knijnenburg, Ritvik Vasan, Calysta Yan, Joy E Arakaki, Matte Bailey, Ben Berry, Antoine Borensztein, Eva M Brown, et al. Integrated intracellular organization and its variations in human ipscs. *Nature*, pages 1–10, 2023.
- [88] Andrew J Vickers and Elena B Elkin. Decision curve analysis: a novel method for evaluating prediction models. *Medical Decision Making*, 26(6):565–574, 2006.
- [89] Andrew J Vickers, Ben Van Calster, and Ewout W Steyerberg. Net benefit approaches to the evaluation of prediction models, molecular markers, and diagnostic tests. *bmj*, 352, 2016.
- [90] David S Wack, Michael G Dwyer, Niels Bergsland, Carol Di Perri, Laura Ranza, Sara Hussein, Deepa Ramasamy, Guy Poloni, and Robert Zivadinov. Improved assessment of multiple sclerosis lesion segmentation agreement via detection and outline error estimates. *BMC medical imaging*, 12(1):1–10, 2012.
- [91] Matthijs J Warrens. Some paradoxical results for the quadratically weighted kappa. *Psychometrika*, 77(2):315–323, 2012.
- [92] David Widmann, Fredrik Lindsten, and Dave Zachariah. Calibration tests in multi-class classification: A unifying framework. *Advances in Neural Information Processing Systems*, 32, 2019.
- [93] Manuel Wiesenfarth, Annika Reinke, Bennett A Landman, Matthias Eisenmann, Laura Aguilera Saiz, M Jorge Cardoso, Lena Maier-Hein, and Annette Kopp-Schneider. Methods and open-source toolkit for analyzing and visualizing challenge results. *Scientific Reports*, 11(1):1–15, 2021.
- [94] Varduhi Yeghiazaryan and Irina Voiculescu. An overview of current evaluation methods used in medical image segmentation. *Department of Computer Science, University of Oxford*, 2015.
- [95] Varduhi Yeghiazaryan and Irina D Voiculescu. Family of boundary overlap metrics for the evaluation of medical image segmentation. *Journal of Medical Imaging*, 5(1):015006, 2018.
- [96] Qiuming Zhu. On the performance of matthews correlation coefficient (mcc) for imbalanced dataset. *Pattern Recognition Letters*, 136:71–80, 2020.

SUPPLEMENTARY NOTES

1 METRIC FUNDAMENTALS

The present work focuses on biomedical image analysis problems that can be interpreted as classification tasks at the image, object, or pixel level. The vast majority of metrics for these problem categories are directly or indirectly based on epidemiological principles of True Positive (TP), False Negative (FN), False Positive (FP), True Negative (TN), i.e., the *cardinalities* of the so-called confusion matrix. The TP/FN/FP/TN are henceforth referred to as cardinalities. In the case of more than two classes C , we also refer to the entries of the $C \times C$ confusion matrix as cardinalities. For simplicity and clarity in notation, we restrict ourselves to the binary case in most examples. Cardinalities can be computed at the image (segment), object, or pixel level. They are typically computed by comparing the prediction of the algorithm to a reference annotation. Modern neural network-based approaches commonly require a threshold to be set in order to convert the algorithm output comprising predicted class scores (also referred to as continuous class scores) to a confusion matrix. For the purpose of metric recommendation, the available metrics can be broadly classified as follows (see also [10]):

- **Counting metrics** operate directly on the confusion matrix and express the metric value as a function of the cardinalities. In the context of segmentation, they are typically referred to as **overlap-based** metrics [79]. We distinguish **multi-class counting metrics**, which are defined for an arbitrary number of classes and invariant under class order, from **per-class counting metrics**, which are computed by treating one class as foreground/positive class and all other classes as background. Popular examples for the former include MCC or Accuracy, while examples for the latter are Sensitivity, Specificity and DSC.
- **Multi-threshold metrics** operate on a dynamic confusion matrix, reflecting the conflicting properties of interest, such as high Sensitivity and high Specificity. Popular examples include the AUROC and AP.
- **Distance-based metrics** have been designed for semantic and instance segmentation tasks. They operate exclusively on the TPs and rely on the explicit definition of object boundaries. Popular examples are the HD and the NSD.

Depending on the context (e.g., image-level classification vs. semantic segmentation task) and the community (e.g., medical imaging community vs. computer vision community), identical metrics are referred to with different terminology. For example, Sensitivity, True Positive Rate (TPR) and Recall refer to the same concept. The same holds true for the DSC and the F_1 Score. The most relevant metrics for the problem categories in the scope of this paper are introduced in the following.

Most metrics are recommended to be applied per class (except for the multi-class counting metrics), meaning that a potential multi-class problem is converted to multiple binary classification problems, such that each relevant class serves as the positive class once. This results in different confusion matrices depending on which class is used as the positive class.

1.1 Image-level Classification

Image-level classification refers to the process of assigning one or multiple labels, or *classes*, to an image. Modern algorithms usually output **predicted class scores** (or continuous class scores) between 0 and 1 for every image and class, indicating the probability of the image belonging to a specific class. By introducing a threshold (e.g., 0.5), predictions are considered as positive (e.g., cancer = true) if they are above the threshold, or negative if they are below the threshold.

Subsequently, predictions are assigned to the cardinalities (e.g., a cancer patient with prediction cancer = true is considered as TP) [21]. The most popular classification metrics are counting metrics, operating on a confusion matrix with fixed threshold on the class probabilities, and multi-threshold metrics, as detailed in the following.

Counting metrics. As stated previously, counting metrics rely on the confusion matrix. We distinguish between per-class and multi-class counting metrics. Popular multi-class counting metrics include:

Accuracy [81]: Fig. SN 3.40

Balanced Accuracy (BA) [81]: Fig. SN 3.41

Expected Cost (EC) (also referred to as Expected Prediction Error or Expected Loss) [8, 31, 41]: Fig. SN 3.44

Matthews Correlation Coefficient (MCC) (also referred to as Phi Coefficient) [61]: Fig. SN 3.48

Weighted Cohen's Kappa (WCK) (also referred to as Weighted Cohen's Kappa Coefficient, Weighted Kappa Statistic or Weighted Kappa Score) [17]: Fig. SN 3.56

Popular per-class counting metrics include:

F_β Score [16, 86]: Fig. SN 3.45

Net Benefit (NB) [88]: Fig. SN 3.49

Negative Predictive Value (NPV) [81]: Fig. SN 3.50

Positive Predictive Value (PPV) (also referred to as Precision) [81]: Fig. SN 3.53

Sensitivity (also referred to as Recall, TPR or Hit Rate) [81]: Fig. SN 3.54

Specificity (also referred to as Selectivity or True Negative Rate (TNR)) [81]: Fig. SN 3.55

Multi-threshold metrics. The classical counting metrics presented above rely on fixed thresholds to be set on the predicted class probabilities (if available), resulting in them being based on the cardinalities of the confusion matrix. **Multi-threshold metrics** overcome this limitation by calculating metric scores based on multiple thresholds. Popular examples are:

Area under the Receiver Operating Characteristic Curve (AUROC) (also referred to as Area under the Curve (AUC), AUC - ROC (Area under the Curve - Receiver Operating Characteristics), C-Index, C-Statistics) [40]: Fig. SN 3.57

Average Precision (AP) [30, 56]: Fig. SN 3.58

Calibration metrics. While most research in biomedical image analysis focuses on the discrimination capabilities of classifiers, a complementary property of relevance is the *calibration* of predicted class scores (also known as *confidence scores*). Intuitively speaking, a system is well-calibrated if the predicted class scores (i.e., the output of the model) reflect the true probabilities of the outcome. In practice, this means that calibrated scores match the empirical success rate of associated predictions. For a binary classification task, calibration implies that of all the data samples assigned a predicted score of 0.8 for the positive class, empirically, 80% belong to this class. Popular examples are:

Brier Score (BS) [33]: Fig. SN 3.66

Class-Wise Calibration Error (CWCE) [53, 54]: Fig. SN 3.67

Expected Calibration Error (ECE) [63]: Fig. SN 3.68

Expected Calibration Error Kernel Density Estimate (ECE^{KDE}) [69]: Fig. SN 3.69

Kernel Calibration Error (KCE) [37, 92]: Fig. SN 3.70

Negative Log Likelihood (NLL) [20]: Fig. SN 3.71

Root Brier Score (RBS) [37]: Fig. SN 3.72

1.2 Semantic Segmentation

Semantic segmentation is commonly defined as the process of partitioning an image into multiple segments/regions. To this end, one or multiple labels are assigned to every pixel such that pixels with the same label share certain characteristics. Semantic segmentation can therefore also be regarded as pixel-level classification. As in image-classification problems, predicted class probabilities are typically calculated for each pixel, deciding on the class affiliation based on a threshold over the class scores [1]. In semantic segmentation problems, the pixel-level classification is typically followed by a post-processing step, in which connected components are defined as objects, and object boundaries are created accordingly. Semantic segmentation metrics can roughly be classified into: (1) counting metrics or overlap-based metrics, for measuring the overlap between the reference annotation and the prediction of the algorithm, (2) distance-based or boundary-based metrics, for measuring the distance between object boundaries, and (3) problem-specific metrics, measuring, for example, object volumes.

Counting metrics. The most frequently used segmentation metrics are **counting metrics**. In the context of segmentation they are also referred to as **overlap-based metrics**, as they essentially measure the overlap between a reference mask and the algorithm prediction. Popular examples of overlap-based metrics include:

Dice Similarity Coefficient (DSC) (also referred to as Sørensen–Dice Coefficient, F_1 Score, Balanced F Score) [27]: Fig. SN 3.43

Intersection over Union (IoU) (also referred to as Jaccard Index, Tanimoto Coefficient) [45]: Fig. SN 3.47

centerline Dice Similarity Coefficient (cIDice) [77]: Fig. SN 3.42

Distance-based metrics. Overlap-based metrics are often complemented by **distance-based metrics** that operate exclusively on the TPs and compute one or several distances between the reference and the prediction. Besides few exceptions, distance-based metrics are often **boundary-based metrics** which focus on assessing the accuracy of object boundaries. Popular examples include:

Average Symmetric Surface Distance (ASSD) (also referred to as Weighted Bilateral Mean Contour Distance) [94]: Fig. SN 3.60

Boundary Intersection over Union (Boundary IoU) [13]: Fig. SN 3.61

Hausdorff Distance (HD) (also referred to as Maximum Symmetric Surface Distance, Hausdorff Metric, Pompeiu–Hausdorff Distance) [44]: Fig. SN 3.62

Hausdorff Distance 95th Percentile (HD95) [44]: Fig. SN 3.65

Mean Average Surface Distance (MASD) (also referred to as Mean Surface Distance) [6]: Fig. SN 3.63

Normalized Surface Distance (NSD) (also referred to as Normalized Surface Dice, Surface Distance, Surface Dice, Surface DSC) [66]: Fig. SN 3.64

Problem-specific segmentation metrics. While overlap- and distance-based metrics are the standard metrics used by the general computer vision community, biomedical applications often have special domain-specific requirements. In medical imaging, for example, the actual volume of

an object (e.g., a tumor) may be of particular interest. In this case, **volume metrics** such as the *Absolute* or *Relative Volume Error* and the *Symmetric Relative Volume Difference* can be computed [64].

1.3 Object Detection

Object detection refers to the detection of one or multiple objects (or: instances) of a particular class (e.g., lesion) in an image [56]. The following description assumes single-class problems, but translation to multi-class problems is straightforward, as validation for multiple classes on object level is performed individually per class. Notably, as multiple predictions and reference instances may be present in one image, the predictions need to include localization information, such that reference and predicted objects can be matched. Important design choices with respect to the validation of object detection methods include:

- (1) *How to represent an object?* Representation is typically composed of location information and a class affiliation. The former may for example take the form of a bounding box (i.e., a list of coordinates), a pixel mask, or the object's center point. Additionally, modern algorithms typically assign a confidence value to each object, representing the probability of a prediction corresponding to an actual object of the respective class. Note that a confusion matrix is later computed for a fixed threshold on the predicted class probabilities.⁴
- (2) *How to decide whether a reference instance was correctly detected?* This step is achieved by applying the *localization criterion*. A localization criterion may, for example, be based on comparing the object centers of the reference and prediction or computing their overlap.
- (3) *How to resolve assignment ambiguities?* The above step might lead to ambiguous matchings, such as two predictions being assigned to the same reference object. Several strategies exist for resolving such cases.

The following sections provide details on (1) applying the localization criterion, (2) applying the assignment strategy, and (3) computing the actual performance metrics.

Localization criterion. As one image may contain multiple objects or no object at all, the **localization criterion** or **hit criterion** measures the (spatial) similarity between a prediction (represented by a bounding box, pixel mask, center point or similar) and a reference object. It defines whether the prediction *hit/detected* (TP) or *missed* (FP) the reference. Any reference object not detected by the algorithm is defined as FN. Please note that TNs are not defined for object detection tasks. Popular localization criteria include:

Box/Approx Intersection over Union (IoU) [45]: Fig. SN 3.76

Mask IoU > 0 [45, 90]: Fig. SN 3.77

Center Distance [39]: Fig. SN 3.74

Point inside Mask/ Box/ Approx⁵: Fig. SN 3.78

Assignment strategy. The localization criterion alone is not sufficient to extract the final confusion matrix based on a fixed threshold for the predicted class probabilities (confidence scores), as ambiguities can occur. For example, two predictions may have been assigned to the same reference object in the localization step, or vice versa. These ambiguities need to be resolved in a further

⁴Please note that we will use the term confidence scores analogously to predicted class probabilities in the context of object detection and instance segmentation.

⁵<https://cada.grand-challenge.org/Assessment/>

assignment step. This assignment and thus the resolving of potential assignment ambiguities can be done via different strategies:

- Greedy (by Score) Matching** [30]: Fig. SN 3.79
- Optimal (Hungarian) Matching** [51]: Fig. SN 3.81
- Matching via Overlap > 0.5** [28]: Fig. SN 3.82
- Greedy (by Localization Criterion) Matching** [58]: Fig. SN 3.80

Metric computation. Similar to image-level classification and semantic segmentation algorithms, object detection algorithms are commonly assessed with counting metrics, assuming a fixed confusion matrix. Popular examples include:

- F_β Score** [16, 86]: Fig. SN 3.45
- False Positives per Image (FPPI)** [5, 85]: Fig. SN 3.46
- Positive Predictive Value (PPV)** (also referred to as Precision) [81]: Fig. SN 3.53
- Sensitivity** (also referred to as Recall, TPR or Hit Rate) [81]: Fig. SN 3.54

Similarly, multi-threshold metrics rely on a range of thresholds. Popular examples are:

- Average Precision (AP)** [30, 56]: Fig. SN 3.58
- Free-Response Receiver Operating Characteristic (FROC) Score** [85]: Fig. SN 3.59

1.4 Instance Segmentation

In contrast to semantic segmentation, **instance segmentation** problems distinguish different instances of the same class (e.g., different lesions). Similarly to object detection problems, the task is to detect individual instances of the same class, but detection performance is measured by pixel-level correspondences (as in semantic segmentation problems). Optionally, instances can be applied to one of multiple classes. Validation metrics in instance segmentation problems often combine common detection metrics with segmentation metrics applied per instance. For instance, segmentation problems, we consider different localization criteria, namely:

Localization criteria:

- Boundary Intersection over Union (Boundary IoU)** [13]: Fig. SN 3.73
- Mask IoU** [45]: Fig. SN 3.76
- Intersection over Reference (IoR)** [60]: Fig. SN 3.75

Additional counting metric: If detection and segmentation performance should be assessed simultaneously in a single score, the **PQ** metric can be utilized [47]: Fig. SN 3.51.

It should be noted that instance segmentation problems are often phrased as semantic segmentation problems with an additional post-processing step, such as connected component analysis [74].

2 METRIC PITFALLS

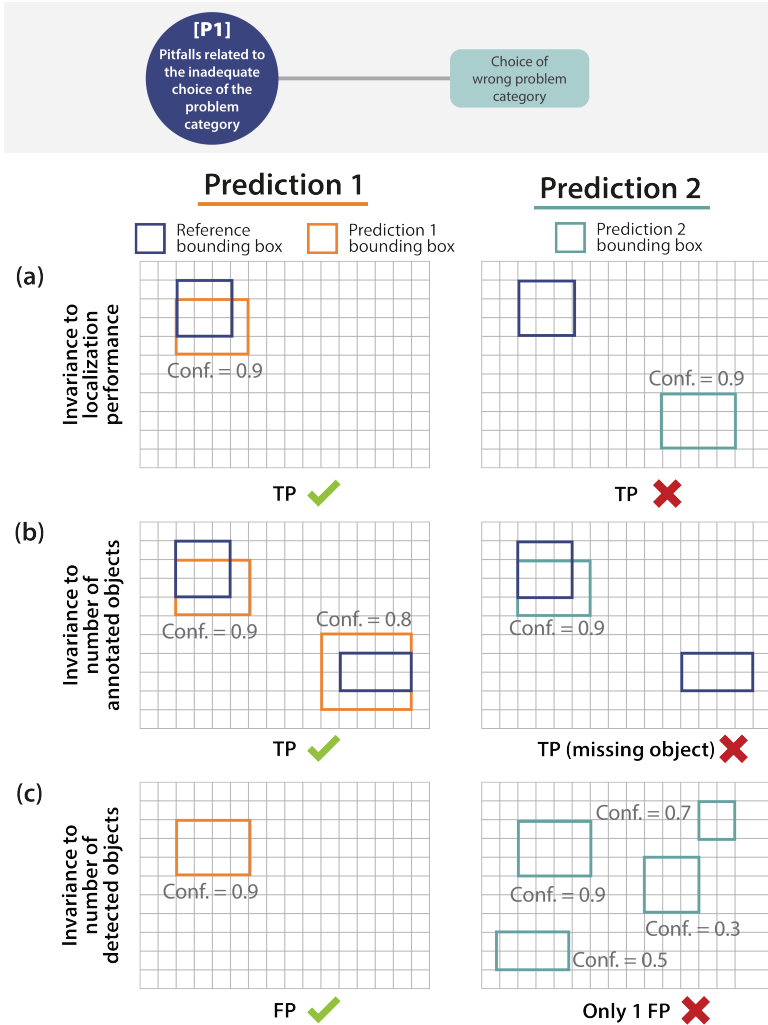
This section presents common limitations of image processing metrics related to [P1] an inadequate choice of problem category (Suppl. Note 2.1), [P2] poor metric selection (Suppl. Note 2.2) and [P3] poor metric application (Suppl. Note 2.3) in an illustrated manner.

To preserve visual clarity, the most important of the presented metric values may be highlighted with color. Green metric values correspond to a "good" metric value (e.g. a high *Sensitivity* score), whereas red values correspond to a "bad" value (e.g. a low *Sensitivity*). Green check marks indicate desirable behavior of metrics, red crosses indicate undesirable behavior. Please note that a low metric value is not automatically a "bad" score. A metric value should always be put into perspective and compared to inter-rater variability. For simplicity, we still use the terms "good" and "bad/poor" throughout the section. Finally, our illustrations do not provide the concrete class probabilities of the presented classifiers.

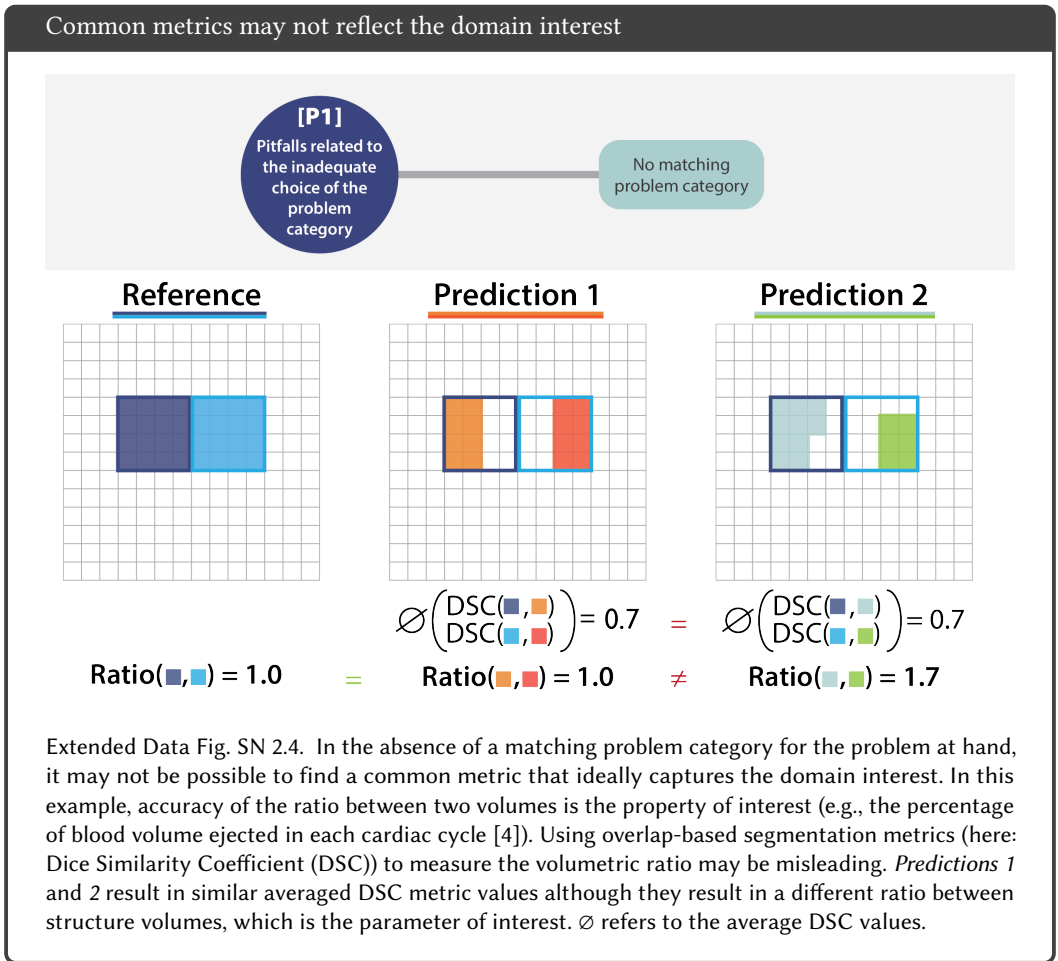
2.1 Pitfalls related to an inadequate choice of the problem category

Performance metrics are typically expected to reflect a domain-specific (e.g., clinical) validation goal. Previous research, however, suggests that this is often not the case [75]. Before choosing validation metrics, the correct problem category needs to be defined. In the following, we present pitfalls related to metrics not being applied to the appropriate problem category. These can either be associated with a wrong choice of the problem category (here: Figs. 3 and SN 2.3; more examples are provided in [71]) or the lack of a matching problem category (Fig. SN 2.4).

Assessing object detection performance at image level yields misleading results



Extended Data Fig. SN 2.3. Image-level classification metrics such as the Area under the Receiver Operating Characteristic Curve (AUROC) curve can be used to validate object detection models by first aggregating predictions to one image-level score (per class). This validation scheme discards the information on the object matching (localization, number of objects etc.). This leads to several problems: **(a)** The image-level Receiver Operating Characteristic (ROC) curve does not measure the localization performance. Both *Prediction 1* and *2* are considered as True Positive (TP) due to their score being very high, although *Prediction 2* does not hit the annotated object. **(b)** The image-level ROC is invariant to the number of annotated objects in an image. The curve does not discriminate between a model detecting all positives (*Prediction 1*) and a model detecting only one of the positives (*Prediction 2*), as long as the maximum score is the same. **(c)** The image-level ROC is invariant to the number of detections in an image. The curve does not discriminate between a model with many False Positives (FP) (*Prediction 2*), and a model with just one FP (*Prediction 1*), as long as the maximum score is the same. The class probabilities are represented by confidence scores (Conf.).

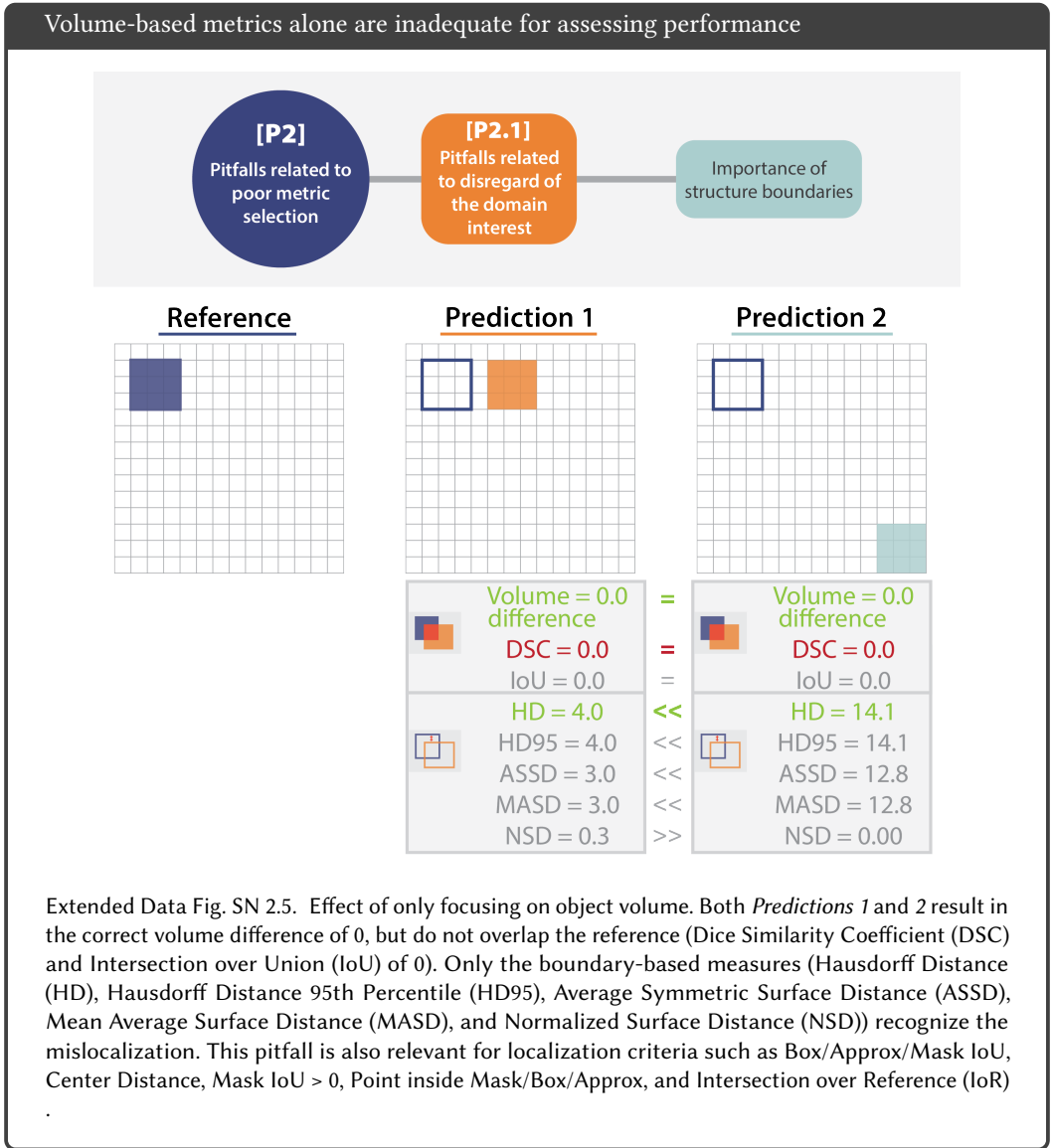


2.2 Pitfalls related to poor metric selection

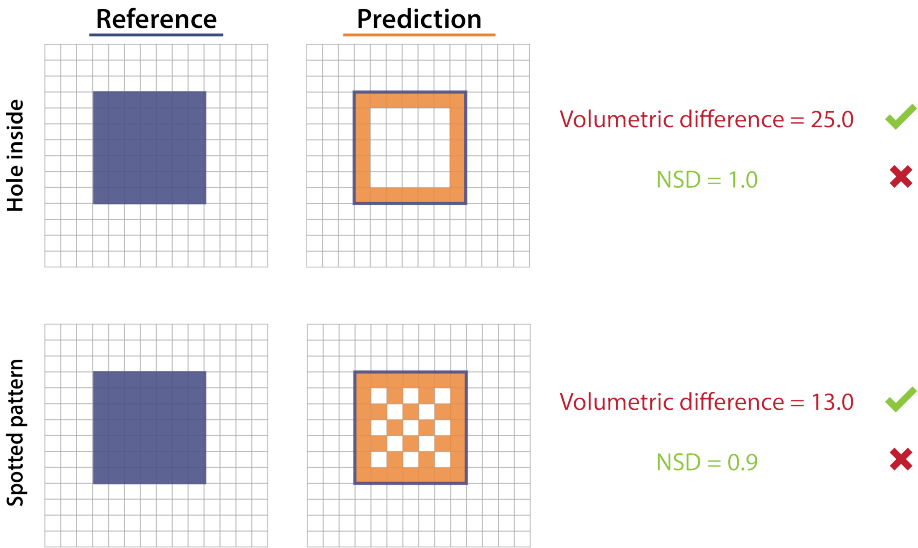
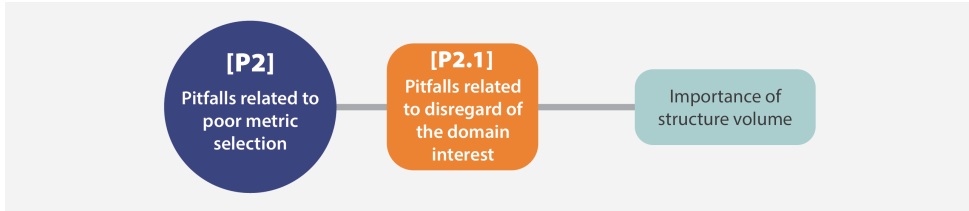
Validation metrics typically assess a specific property of interest. Thus, a metric designed for a particular purpose often cannot be used to appropriately validate another property. This is due to both the limitations as well as the mathematical properties of individual metrics, both of which are often neglected. In this section, we present pitfalls related to poor metric selection.

2.2.1 Pitfalls related to disregard of the domain interest. Several requirements for metric selection arise from the domain interest, which may clash with particular metric limitations. In the following, we present pitfalls related to disregard of the domain interest, stemming from the following sources:

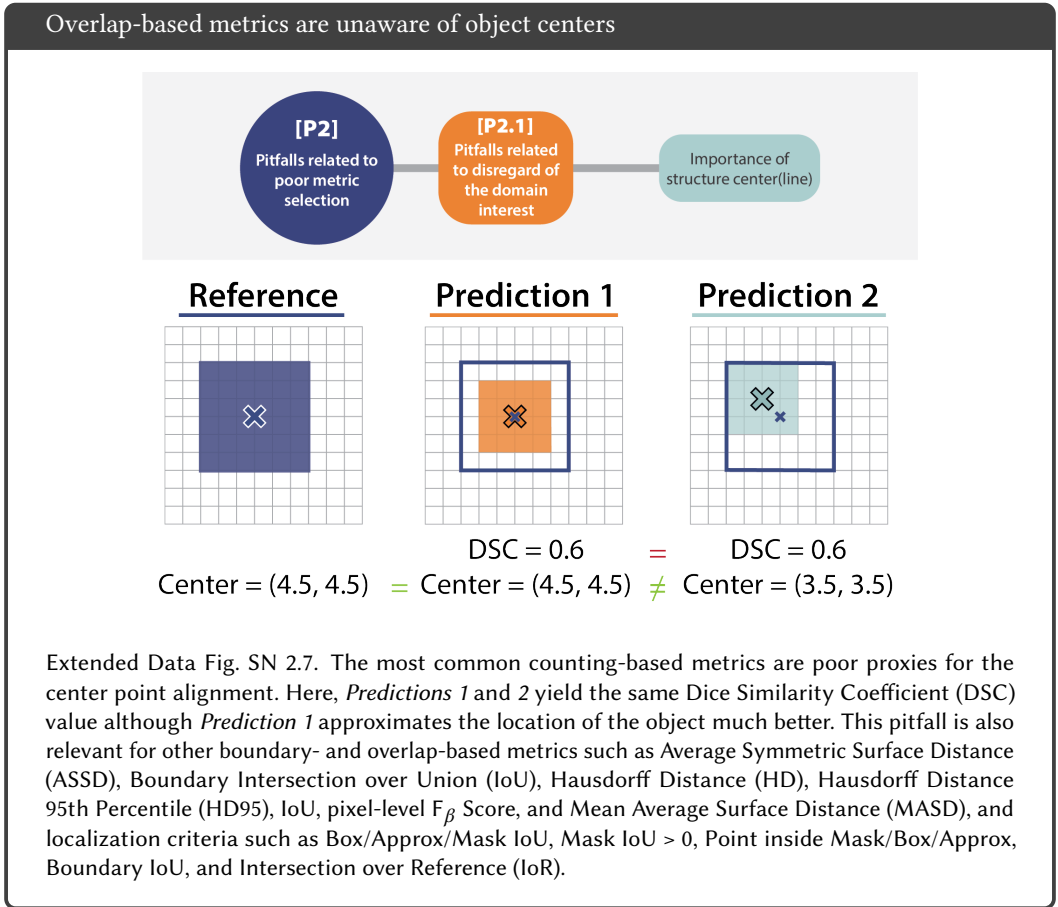
- Importance of structure boundaries (Figs. 4a and SN 2.5)
- Importance of structure volume (Fig. SN 2.6)
- Importance of structure center(line) (Fig. SN 2.7)
- Importance of confidence awareness (Fig. SN 2.8)
- Importance of comparability across data sets (Figs. SN 2.9)
- Unequal severity of class confusions (Figs. 4b and SN 2.10)
- Importance of cost-benefit analysis (Fig. SN 2.11)



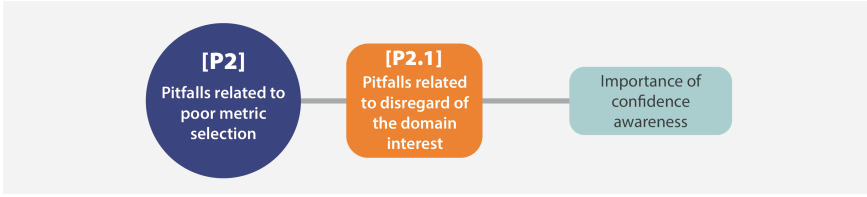
Boundary-based metrics disregard holes in the segmentation



Extended Data Fig. SN 2.6. Boundary-based metrics commonly ignore the overlap between structures and are thus insensitive to holes in structures. In the examples, the Prediction respectively features a hole or spotted pattern within the object. Boundary-based metrics (here: Normalized Surface Distance (NSD)) do not recognize this problem, yielding (near) perfect metric scores of 1.0 and 0.9, whereas the volumetric difference reflects the fact that the inner area is inadequately predicted. NSD was calculated for $\tau = 2$. This pitfall is also relevant for other boundary-based metrics such as Average Symmetric Surface Distance (ASSD), Boundary Intersection over Union (Boundary IoU), Hausdorff Distance (HD), Hausdorff Distance 95th Percentile (HD95), and Mean Average Surface Distance (MASD), as well as localization criteria such as Center Distance, Mask IoU > 0, Point inside Mask/Box/Approx, Boundary IoU, Intersection over Reference (IoR), and Mask IoU.



Common calibration metrics falsely imply perfect calibration



Top-label calibration

$g(X)$	$P[Y \in \cdot \mid g(X)]$
(0.1, 0.3, 0.6)	(0.2, 0.2, 0.6)
(0.1, 0.6 , 0.3)	(0.0, 0.7 , 0.3)
(0.3, 0.1, 0.6)	(0.2, 0.2, 0.6)
(0.3, 0.6 , 0.1)	(0.4, 0.5 , 0.1)
(0.6 , 0.1, 0.3)	(0.7 , 0.0, 0.3)
(0.6 , 0.3, 0.1)	(0.5 , 0.4, 0.1)

top-label ECE = 0

Class-wise calibration

$g(X)$	$P[Y \in \cdot \mid g(X)]$
(0.1 , 0.3, 0.6)	\emptyset [(0.2, 0.2, 0.6)]
(0.1 , 0.6, 0.3)	\emptyset [(0.0, 0.7, 0.3)]
(0.3 , 0.1, 0.6)	\emptyset [(0.2, 0.2, 0.6)]
(0.3 , 0.6, 0.1)	\emptyset [(0.4, 0.5, 0.1)]
(0.6 , 0.1, 0.3)	\emptyset [(0.7, 0.0, 0.3)]
(0.6 , 0.3, 0.1)	\emptyset [(0.5, 0.4, 0.1)]

\emptyset
over all
classes
(example
shown for
class 1)

class-wise ECE = 0

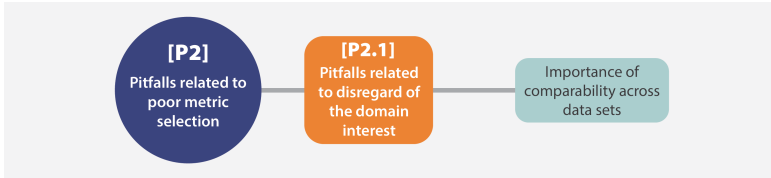
Canonical calibration

$g(X)$	$P[Y \in \cdot \mid g(X)]$
(0.1 , 0.3 , 0.6)	(0.2 , 0.2 , 0.6)
(0.1 , 0.6 , 0.3)	(0.0 , 0.7 , 0.3)
(0.3 , 0.1 , 0.6)	(0.2 , 0.2 , 0.6)
(0.3 , 0.6 , 0.1)	(0.4 , 0.5 , 0.1)
(0.6 , 0.1 , 0.3)	(0.7 , 0.0 , 0.3)
(0.6 , 0.3 , 0.1)	(0.5 , 0.4 , 0.1)

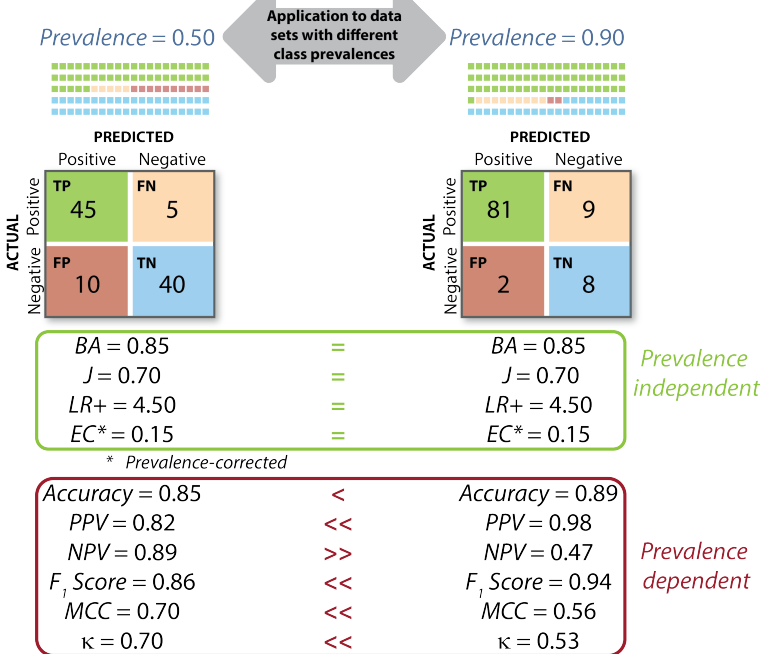
canonical ECE > 0

Extended Data Fig. SN 2.8. Effect of different definitions of calibration on the Expected Calibration Error (ECE) when focusing on confidence or predicted class scores (confidence awareness). For top-label calibration, only the maximum values of the predicted class scores $g(X)$ are considered, while all other values are neglected, resulting in a perfect calibration for this example. Similarly, for class-wise calibration, the predicted class scores are compared class-wise per value, also yielding a perfect score. Only canonical calibration considers all components of the predicted class score vectors, showing that the model is not perfectly calibrated [37, 84]. A more detailed insight in different definitions of calibration is given in [58]. It should be noted that discrimination metrics generally do not assess calibration performance, i.e., perfect discrimination does not imply good calibration performance.

Comparison of metric scores across data sets may be misleading

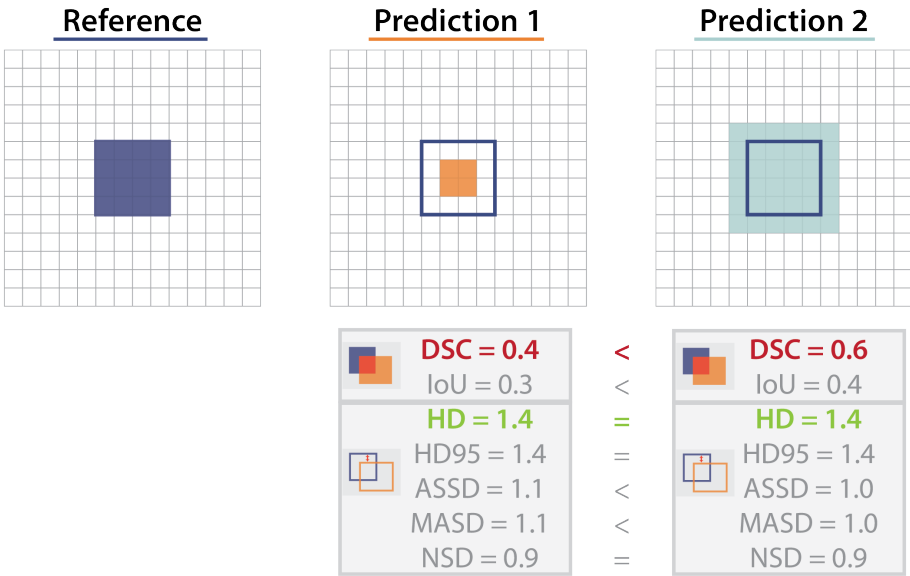
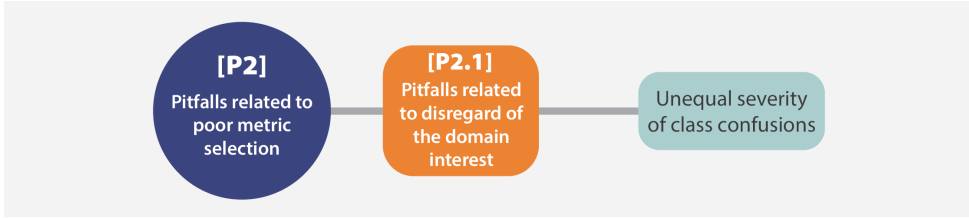


Inherent properties of a method: *Sensitivity = 0.90, Specificity = 0.80*



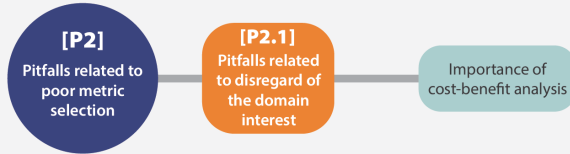
Extended Data Fig. SN 2.9. Effect of prevalence dependency. An algorithm with specific inherent properties (here: Sensitivity of 0.9 and Specificity of 0.8) may perform completely differently on different data sets if the prevalences differ (here: 50% (left) and 90% (right)) and prevalence-dependent metrics are used for validation (here: Accuracy, Positive Predictive Value (PPV), Negative Predictive Value (NPV), F_1 Score, Matthews Correlation Coefficient (MCC), Cohen's Kappa κ). In contrast, prevalence-independent metrics (here: Balanced Accuracy (BA), Youden's Index J, Positive Likelihood Ratio (LR+), and Expected Cost (EC)) can be used to compare validation results across different data sets. Used abbreviations: True Positive (TP), False Negative (FN), False Positive (FP) and True Negative (TN). This pitfall is also relevant for other counting metrics such as Net Benefit (NB).

Overlap-based metrics prefer oversegmentation over undersegmentation



Extended Data Fig. SN 2.10. Effect of undersegmentation vs. oversegmentation. The outlines of the predictions of two algorithms (*Prediction 1/2*) differ in only a single layer of pixels (*Prediction 1*: undersegmentation – smaller structure compared to reference, *Prediction 2*: oversegmentation – larger structure compared to reference). This has no (or only a minor) effect on the Hausdorff Distance (HD)/(95%), the Normalized Surface Distance (NSD), MASD, and the Average Symmetric Surface Distance (ASSD), but yields a substantially different Dice Similarity Coefficient (DSC) or Intersection over Union (IoU) score [79, 95]. If penalizing of either over- or undersegmentation is desired (unequal severity of class confusions), other metrics such as the F_{β} Score provide specific penalties for either depending on the chosen hyperparameter β . This pitfall is also relevant for other overlap-based metrics such as centerline Dice Similarity Coefficient (cIDice) and localization criteria such as Box/Approx/Mask IoU, Boundary IoU, and Intersection over Reference (IoR).

Common metrics disregard cost-benefit analysis



Cost-benefit analysis:
 ~9 unnecessary biopsies for one detected lesion are acceptable.

1) BIOPSY IN ALL PATIENTS:

		PREDICTED	
		Positive	Negative
ACTUAL	Positive	TP 30	FN 0
	Negative	FP 75	TN 0

Accuracy = 0.29

2) MARKER-BASED DECISION ON BIOPSY:

		PREDICTED	
		Positive	Negative
ACTUAL	Positive	TP 20	FN 15
	Negative	FP 60	TN 10

Accuracy = 0.29

75 unnecessary biopsies (FP)



60 unnecessary biopsies (FP)

NB metrics relate the benefit of TPs with the cost of FPs
 (here: 1/9 based on benefit-cost analysis)

NB = 0.21

>

NB = 0.13

Same Accuracy, but better clinical utility

Same Accuracy, but poorer clinical utility

Extended Data Fig. SN 2.11. Effect of neglecting a cost-benefit analysis. In a cost-benefit analysis, clinicians are able to define a risk-specific exchange rate that is used in the computation of the Net Benefit (NB) metric. Common metrics such as Accuracy do not consider this analysis and would favor the marker-based decision on biopsy, while NB indicates that biopsies of all patients actually yield a better clinical outcome [89]. This pitfall is also relevant for other counting metrics such as Balanced Accuracy (BA), Positive Likelihood Ratio (LR+), Matthews Correlation Coefficient (MCC), Negative Predictive Value (NPV), Positive Predictive Value (PPV), Sensitivity, and Specificity. For binary problems, the hyperparameter β of the F_β Score can be used as a dynamic penalty for class confusions.

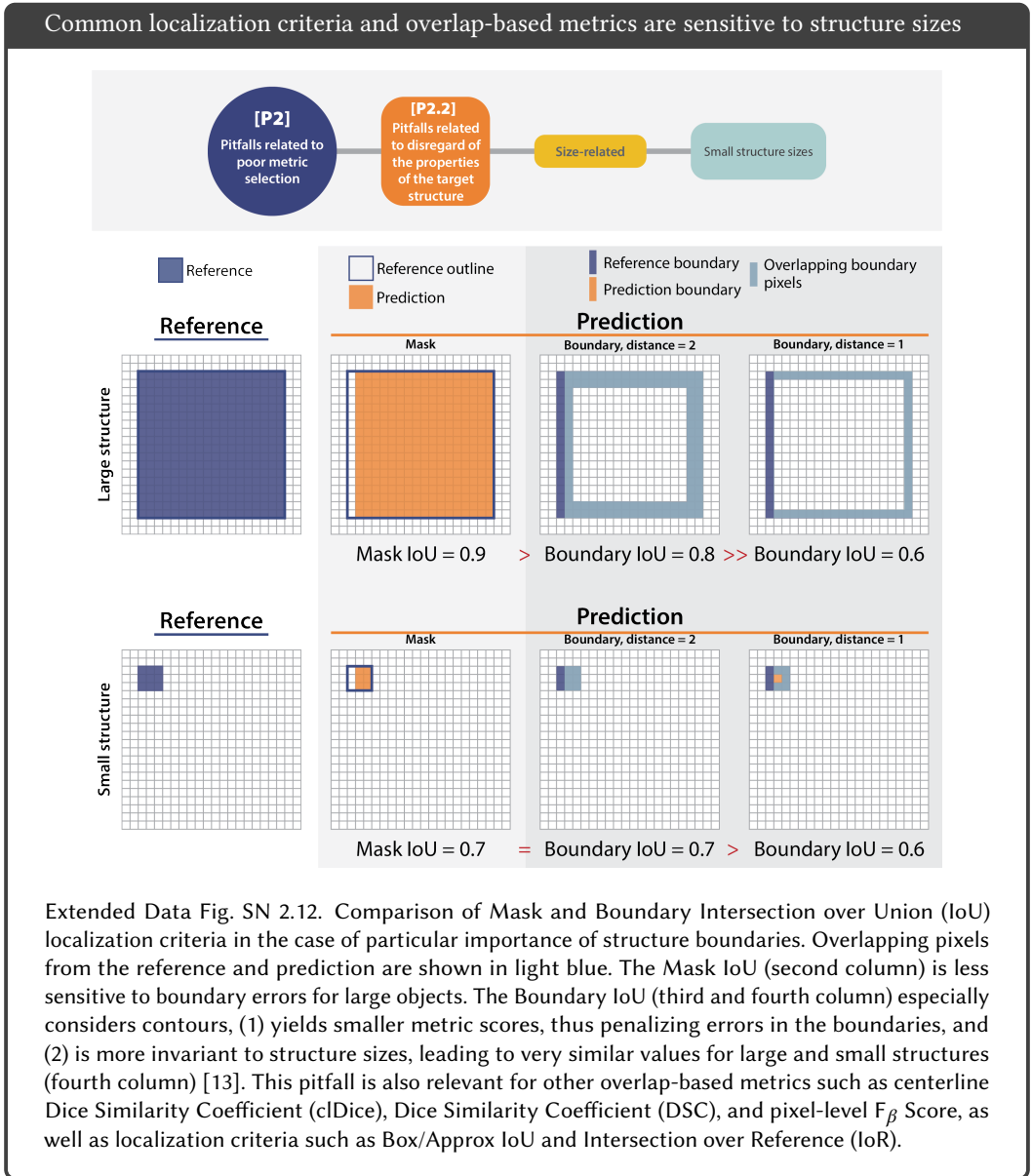
2.2.2 Pitfalls related to disregard of the properties of the target structure. For problems that require capturing local properties (object detection, semantic or instance segmentation), the properties of the target structures to be localized and/or segmented may have severe implications for metric choice. Pitfalls can be further subdivided into *size-related* and *shape- and topology-related* pitfalls. In the following, we present pitfalls stemming from the following sources:

Size-related pitfalls:

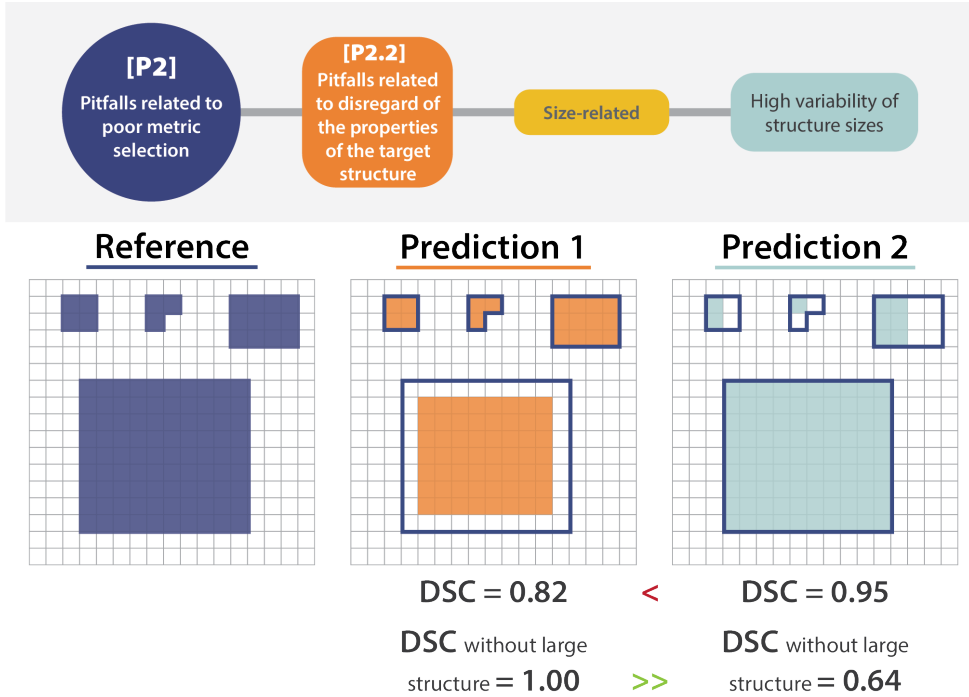
- Small structure sizes (Extended Data Fig. 1a and Fig. SN 2.12)
- High variability of structure sizes (Fig. SN 2.13)

Shape- and topology-related pitfalls

- Complex structure shapes (Extended Data Fig. 1b and Fig. SN 2.14)
- Occurrence of overlapping or touching structures (Fig. SN 2.15)
- Occurrence of disconnected structures (Fig. SN 2.16)

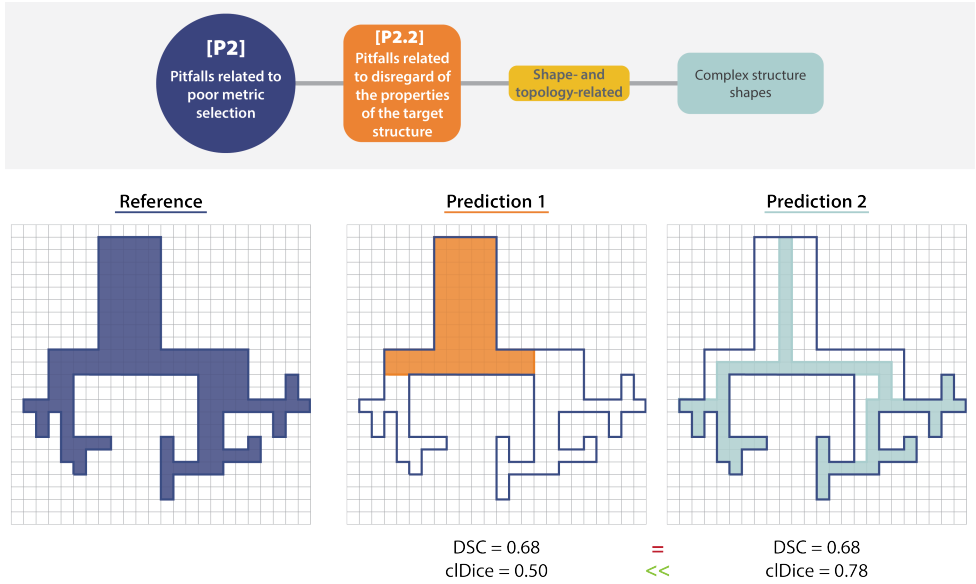


Effect of high variability of structure sizes



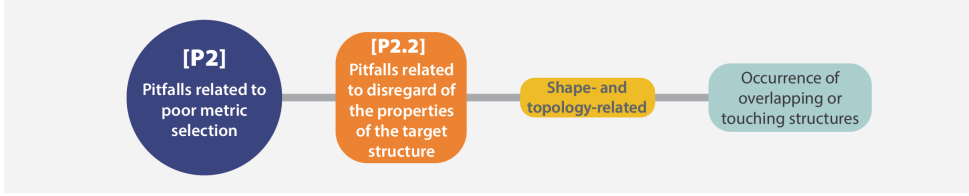
Extended Data Fig. SN 2.13. Large structures completely dominate overlap-based metrics in semantic segmentation problems. While *Prediction 1* perfectly segments all three small structures, the metric score (here: Dice Similarity Coefficient (DSC)) is much worse compared to the score of *Prediction 2*, with only one perfect prediction for the large structure. This is highlighted by only computing the metric without the large structure. This pitfall is also relevant for other overlap-based metrics such as centerline Dice Similarity Coefficient (cIDice), Dice Similarity Coefficient (DSC), and pixel-level F_β Score, as well as localization criteria such as Mask/Box/Approx Intersection over Union (IoU) and Intersection over Reference (IoR).

Common metrics are unaware of object shapes

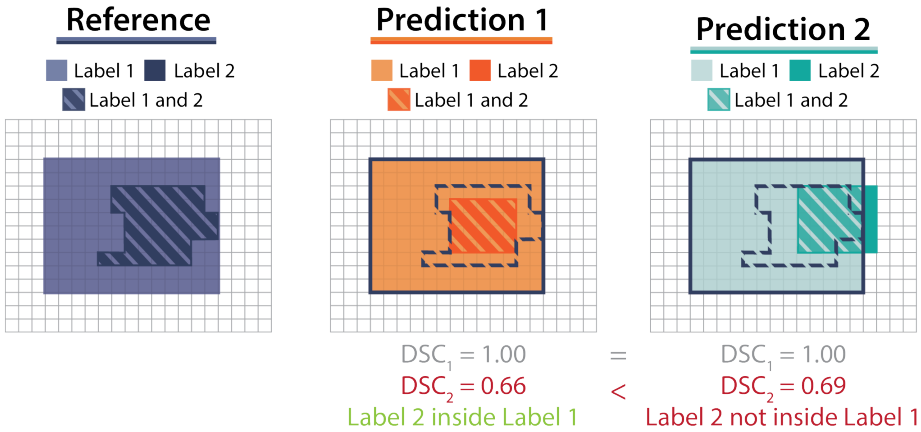


Extended Data Fig. SN 2.14. Effect of complex shapes. Common overlap-based metrics such as the Dice Similarity Coefficient (DSC) are unaware of complex structure shapes and treat *Predictions 1* and *2* equally. The centerline Dice Similarity Coefficient (cDice) uncovers that *Prediction 1* misses the fine-granular branches of the reference and favors *Prediction 2*, which focuses on the object's center line and better captures its fine branches. This pitfall is also relevant for other overlap-based metrics such as Intersection over Union (IoU) and pixel-level F_β Score, and localization criteria such as Box/Approx/Mask IoU, Center Distance, Mask IoU > 0, Point inside Mask/Box/Approx, and Intersection over Reference (IoR).

Common metrics do not account for hierarchical label structure

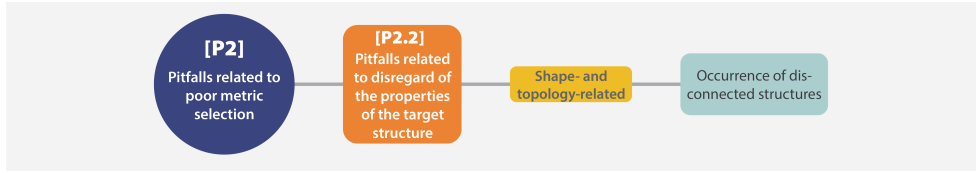


Required: Label 2 is inside Label 1

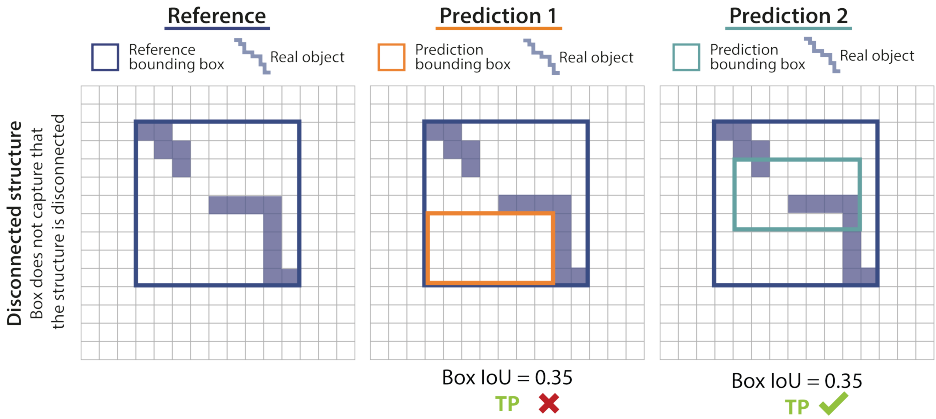


Extended Data Fig. SN 2.15. Effect of nested multi-label structures. The requirement of *Label 2* being inside *Label 1* is violated by *Prediction 2*. Nevertheless, *Prediction 2* has a higher Dice Similarity Coefficient (DSC) score compared to *Prediction 1*, which adheres to the requirement. This pitfall is also relevant for other boundary- and overlap-based metrics such as Average Symmetric Surface Distance (ASSD), Boundary Intersection over Union (IoU), centerline Dice Similarity Coefficient (clDice), Hausdorff Distance (HD), Hausdorff Distance 95th Percentile (HD95), IoU, pixel-level F_β Score, Mean Average Surface Distance (MASD), and Normalized Surface Distance (NSD).

Bounding boxes are inadequate for representing complex shapes and disconnected structures



Box IoU > 0.3: True positive (TP)
 Box IoU ≤ 0.3: False positive (FP)



Extended Data Fig. SN 2.16. Bounding boxes are not well-suited for representing disconnected shapes, in particular multi-component structures. *Predictions 1* and *2* both yield a True Positive (TP) detection, as the Box Intersection over Union (IoU) is larger than the threshold 0.3. However, *Prediction 1* does not hit the real object at all.

2.2.3 *Pitfalls related to disregard of the properties of the data set and algorithm output.* Properties of the data set such as class imbalances or high inter-rater variability may directly affect metric values. Pitfalls can be further subdivided into *class-related* and *reference-related* pitfalls. For reference-based metrics, the algorithm output will be compared against the reference annotation to compute a metric score. Thus, the content and format of the prediction is of high relevance for metric choice. In the following, we present pitfalls stemming from the following sources:

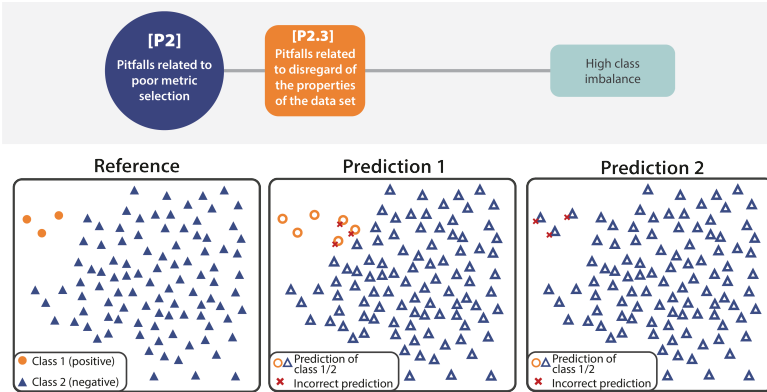
[P2.3] Disregard of the properties of the data set

- High class imbalance (Figs. 5a and SN 2.17)
- Small test set size (Figs. 5b and SN 2.18)
- Imperfect reference standard (Figs. 5c and SN 2.19)

[P2.4] Disregard of the properties of the algorithm output

- Possibility of empty prediction (Extended Data Fig. 2b and Fig. SN 2.20)
- Possibility of overlapping predictions (Extended Data Fig. 2a and Fig. SN 2.21)
- Lack of predicted class scores (Fig. SN 2.22)

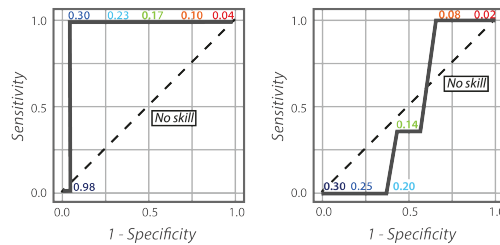
Common metrics yield implausible results in the presence of class imbalance



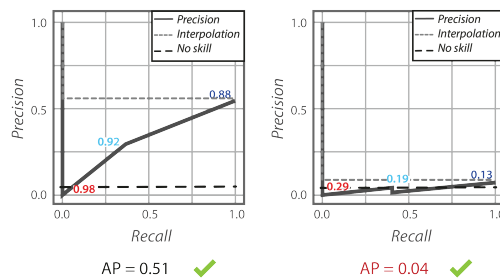
(A) COUNTING METRICS

Accuracy = 0.97	=	Accuracy = 0.97
BA = 0.98	>>	BA = 0.50
Sensitivity = 1.00	>>	Sensitivity = 0.00
PPV = 0.50	<	PPV = NaN
Specificity = 0.97	>>	Specificity = 1.00
F ₁ Score = 0.67	>	F ₁ Score = 0.00
NPV = 1.00	>>	NPV = 0.97
MCC = 0.70	>>	MCC = 0.00
WCK = 0.65	>>	WCK = 0.00

(B) ROC-CURVE

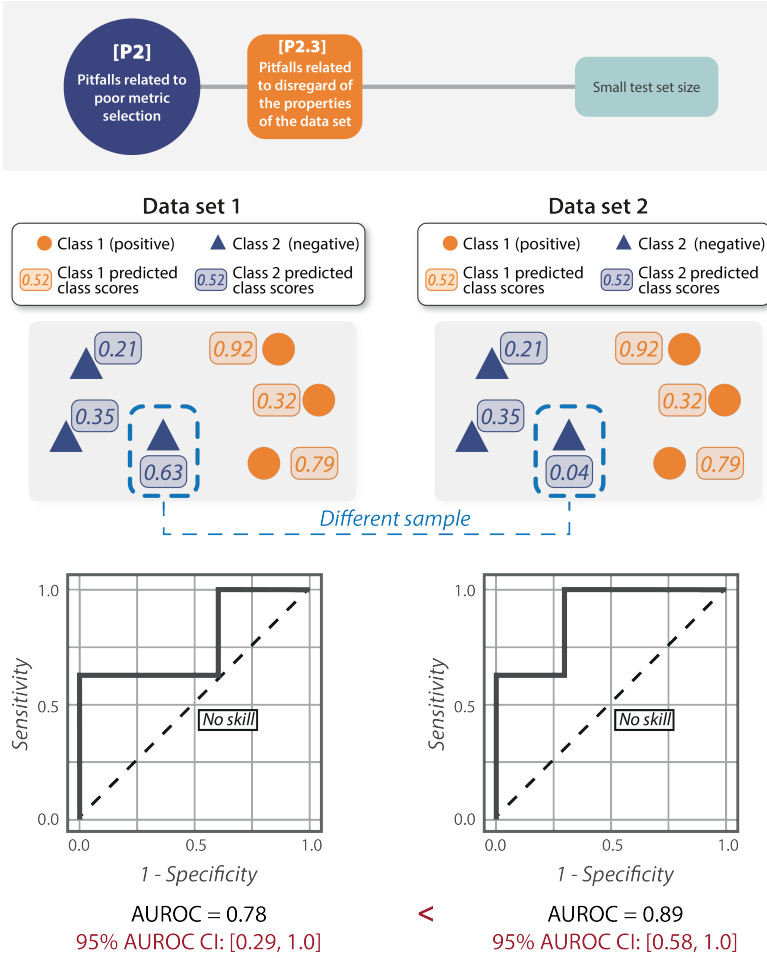


(C) PR-CURVE

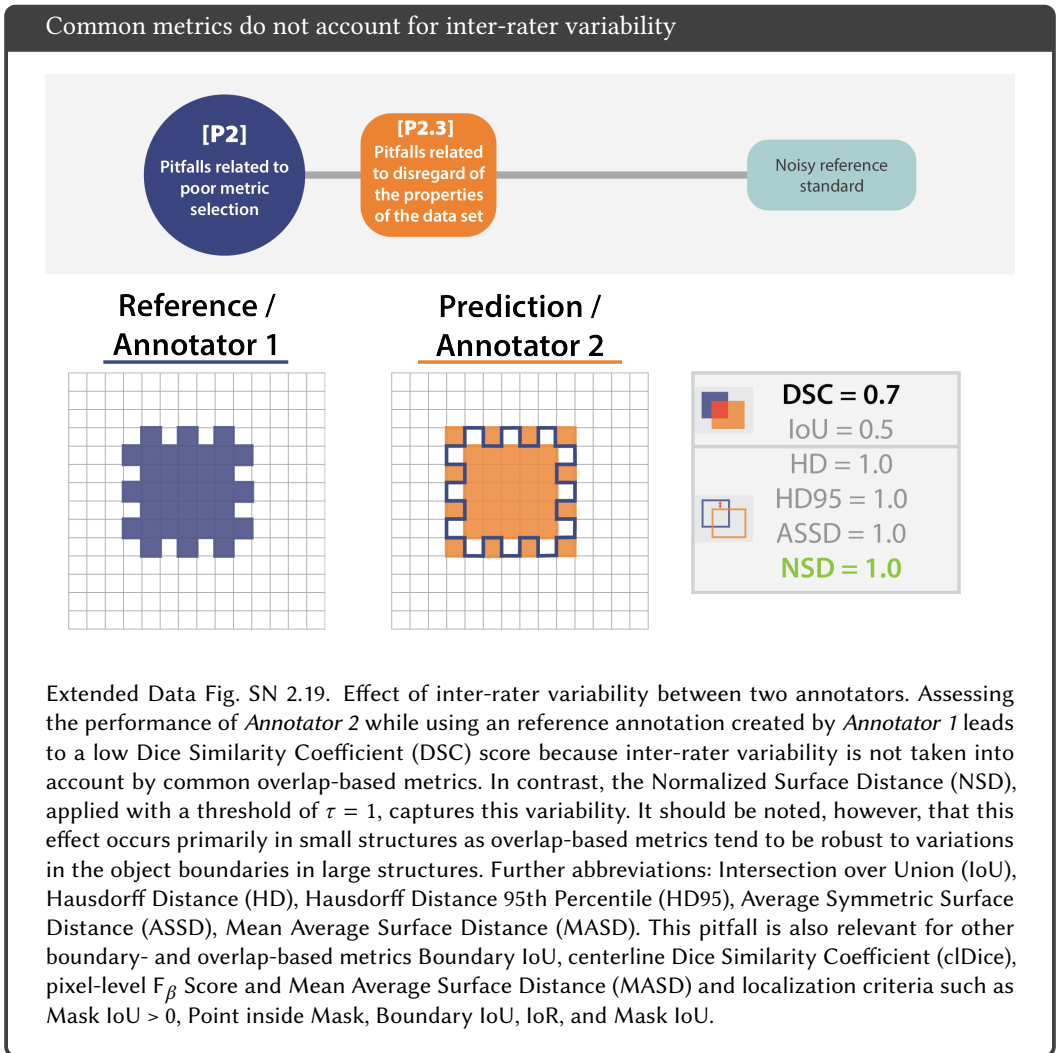


Extended Data Fig. SN 2.17. Effect of class imbalance. Not every metric is designed to reflect class imbalance [14]. In the case of underrepresented classes, an unsuitable metric, such as Accuracy, yields a high value even if the classifier performs very poorly for one of the classes (here: *Prediction 2*). Multi-threshold metrics, such as the Area under the Receiver Operating Characteristic Curve (AUROC) and the Average Precision (AP), reveal the weakness, indicating that *Prediction 2* is not better than random guessing. For comparison, a no-skill classifier (random guessing) is shown as a black dashed line. For the Precision-Recall (PR) curves, the interpolation applied to compute the AP metric is shown as a dashed grey line. Thresholds used for curve generation are provided as small numbers above the curve. Further abbreviations: Positive Predictive Value (PPV), Negative Predictive Value (NPV), Matthews Correlation Coefficient (MCC), Weighted Cohen's Kappa (WCK). This pitfall is also relevant for other counting metrics such as Net Benefit (NB).

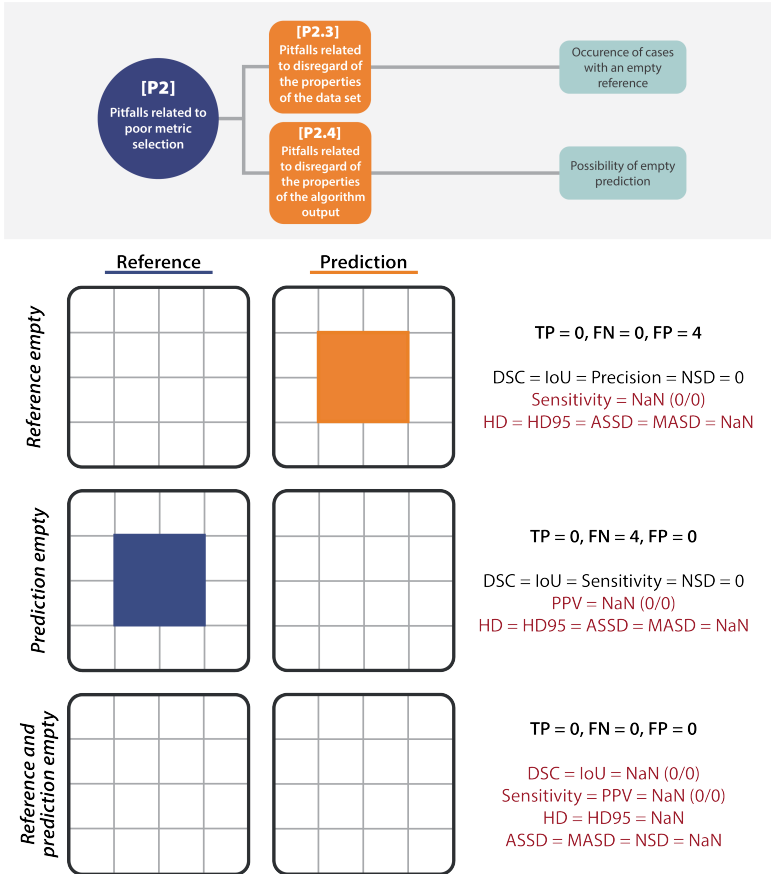
Common multi-threshold metrics are not well-suited for small sample sizes



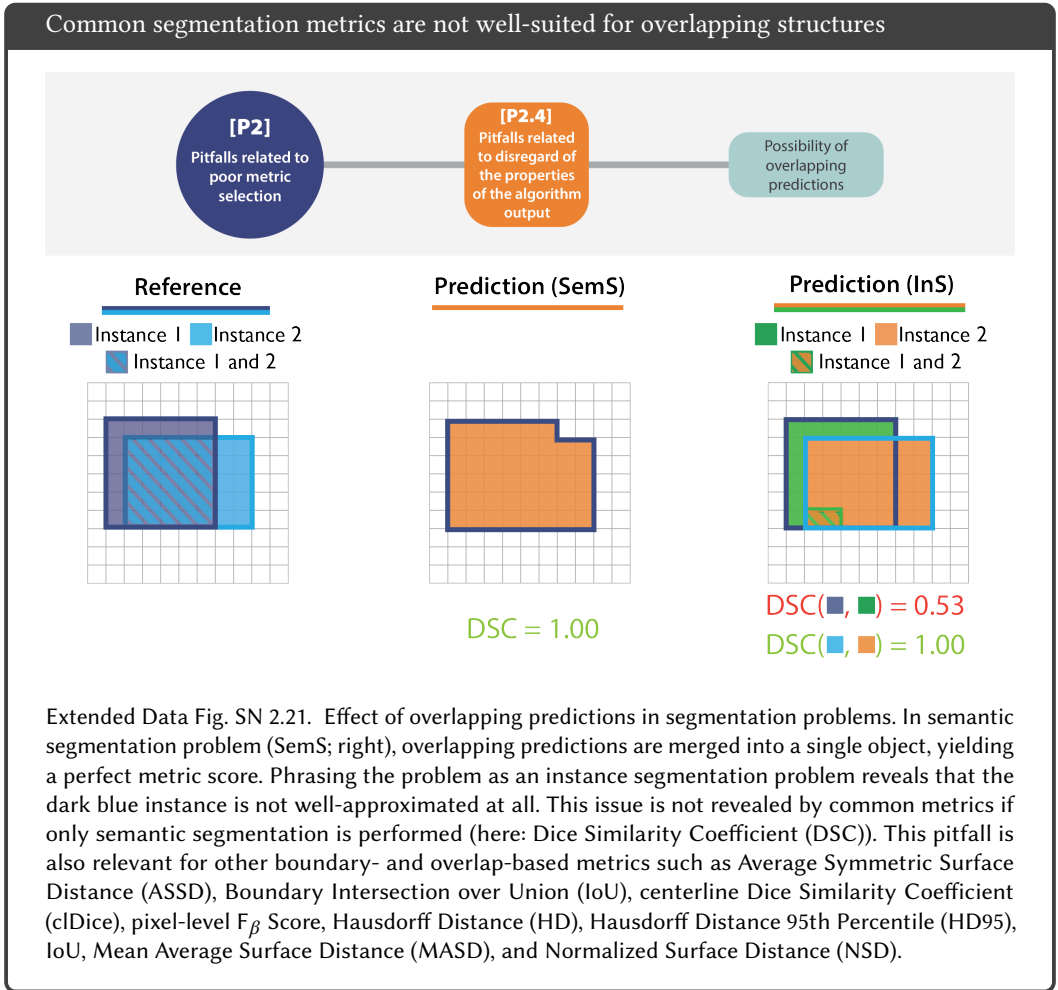
Extended Data Fig. SN 2.18. Effect of calculating the Area under the Receiver Operating Characteristic Curve (AUROC) for very small sample sizes. The AUROC is very unstable for small sample sizes. *Data sets* 1 and 2 only contain six samples each, for which only one predicted score differs between sets. Drawing the Receiver Operating Characteristic (ROC) curve and calculating the AUROC leads to a large difference in scores between both data sets. The 95% Confidence Interval (CI) reveals that there is a large range of possible AUROC values. CIs were calculated based on [25]. This pitfall is also relevant for other counting metrics such as Accuracy, Average Precision (AP), Balanced Accuracy (BA), Expected Cost (EC), F_β Score, Free-Response Receiver Operating Characteristic (FROC) Score, Positive Likelihood Ratio (LR+), Matthews Correlation Coefficient (MCC), Net Benefit (NB), Negative Predictive Value (NPV), Positive Predictive Value (PPV), Sensitivity, Specificity, and Weighted Cohen’s Kappa (WCK).



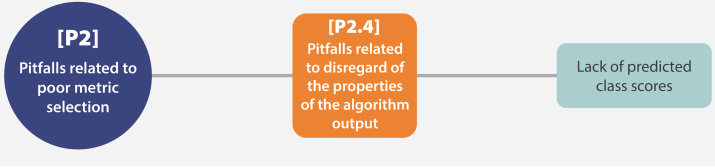
Empty reference or prediction leads to invalid scores



Extended Data Fig. SN 2.20. Effect of empty references or predictions when applying common metrics per image (here for semantic segmentation). Empty images lead to division by zero for many common metrics as the numbers of the TPs, FPs, FNPs turn zero. Used abbreviations: Average Symmetric Surface Distance (ASSD), Dice Similarity Coefficient (DSC), Hausdorff Distance (HD), Hausdorff Distance 95th Percentile (HD95), Intersection over Union (IoU), Mean Average Surface Distance (MASD), Not a Number (NaN), Normalized Surface Distance (NSD). This pitfall is also relevant for other boundary-based, overlap-based and counting metrics such as Boundary IoU, centerline Dice Similarity Coefficient (cIDice), F_{β} Score, Negative Predictive Value (NPV), Positive Predictive Value (PPV), Sensitivity, and Specificity.

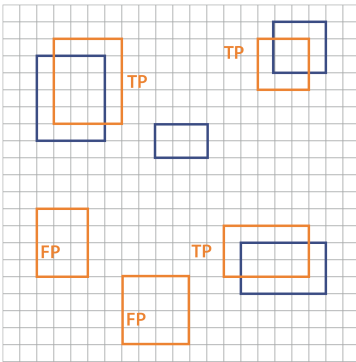


Selection of multi-threshold metrics in the absence of predicted class scores



Prediction

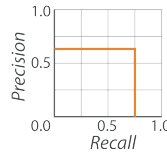
□ Reference bounding box □ Prediction bounding box



Absence of predicted class scores/confidence scores:
Multi-threshold metrics cannot be computed without hacks that influence the metric scores

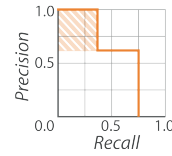
F_1 Score = 0.67 ✓

AP Implementation 1:
 Start curve at (0, Precision)



AP = 0.45

AP Implementation 2:
 Add extra point at (0, 1)



AP = 0.60

<

✗

Extended Data Fig. SN 2.22. Multi-threshold metrics should only be computed if predicted class scores are available, although an increasing body of work computes multi-threshold metrics such as AP in the absence of class scores (e.g., [3, 22, 32, 42, 52]). Otherwise, the strategy chosen for compensating the lack of class scores (here reflected by *Implementations 1* and *2*) leads to metric scores that are less well interpretable than those of established counting metrics working on a fixed confusion matrix (here: F_1 Score). This pitfall is also relevant for other multi-threshold metrics such as Area under the Receiver Operating Characteristic Curve (AUROC) and Free-Response Receiver Operating Characteristic (FROC) Score.

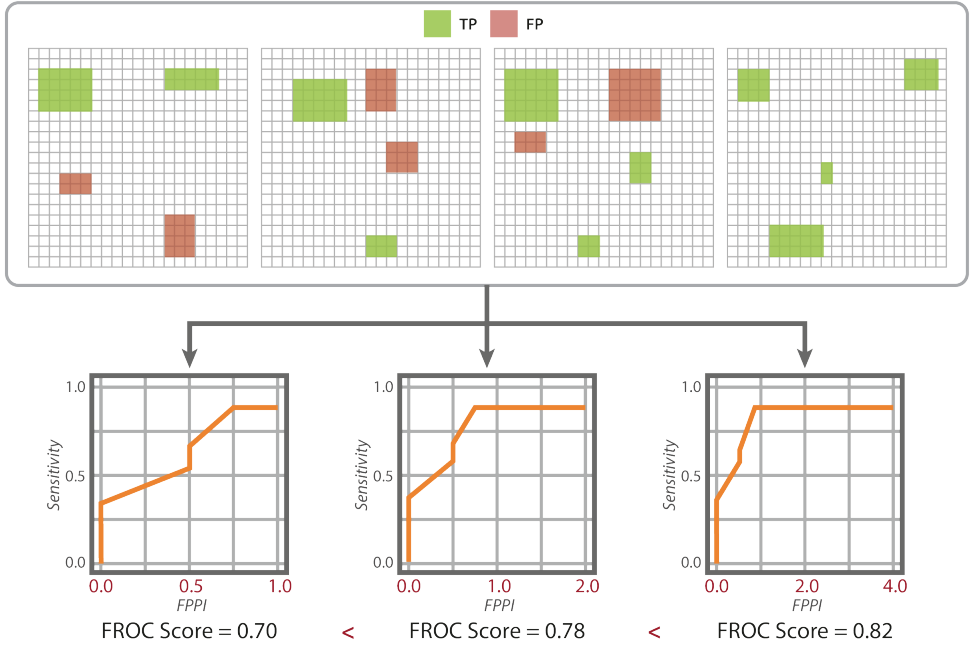
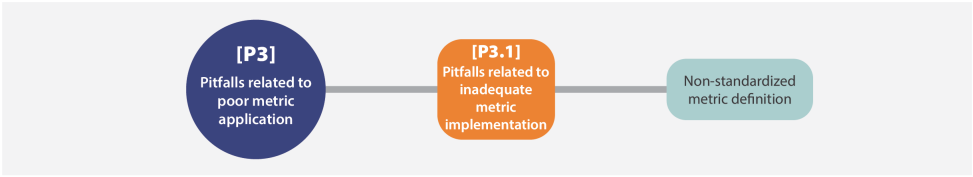
2.3 Pitfalls related to poor metric application

A data set typically contains several hundreds or thousands of images. When analyzing, aggregating and combining metric values, a number of factors need to be taken into account.

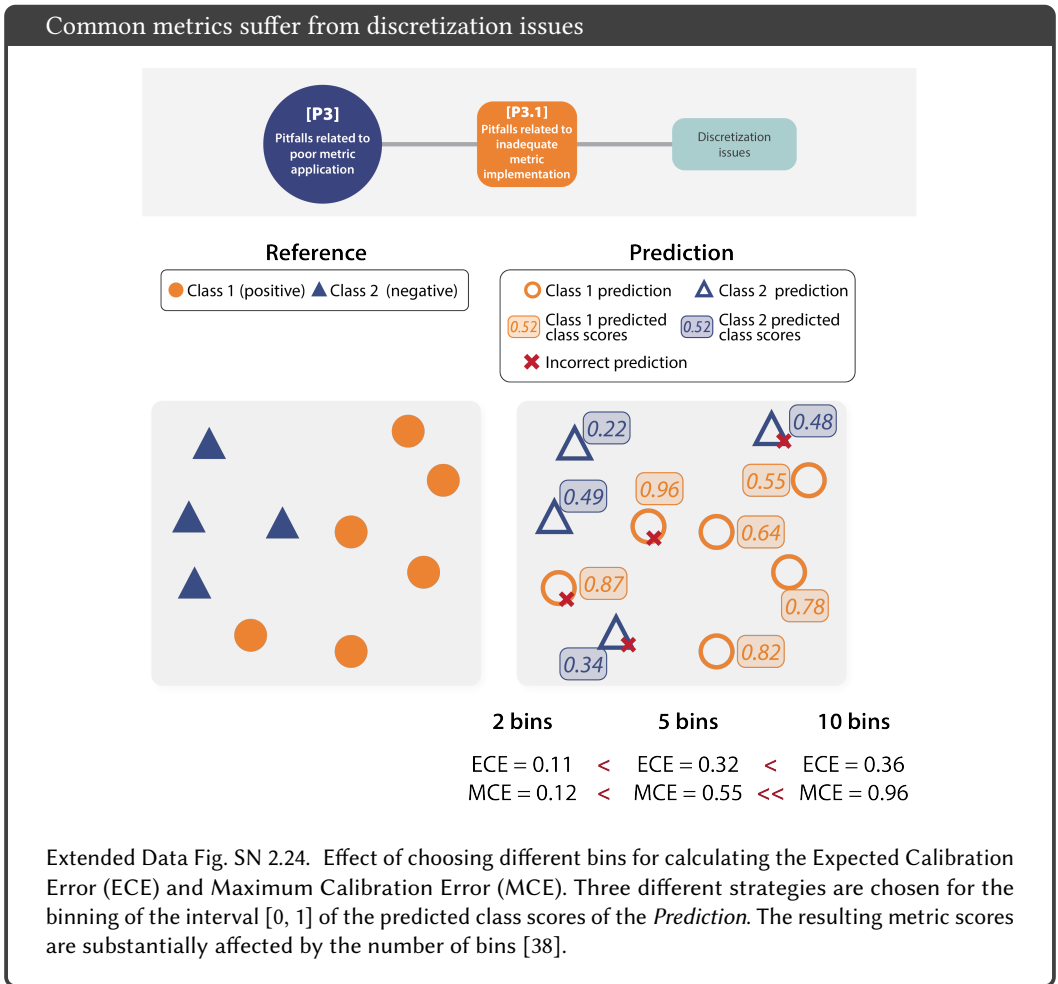
2.3.1 Pitfalls related to inadequate metric implementation. The implementation of metrics is, unfortunately, not standardized. While some metrics are straightforward to implement, others require more advanced techniques and offer a variety of implementation possibilities. Sources of metric implementation pitfalls include:

- Non-standardized definitions (Figs. 6a and SN 2.23)
- Discretization issues (Fig. SN 2.24)
- Sensitivity to hyperparameters (Fig. SN 2.25)
- Metric-specific issues (Fig. SN 2.26)

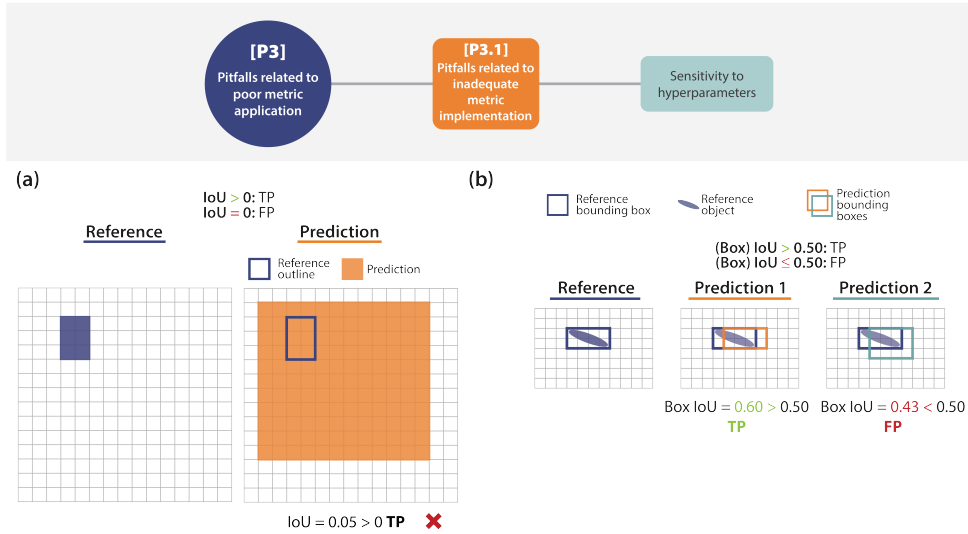
Lack of standardization leads to variation in metric scores



Extended Data Fig. SN 2.23. Effect of defining different ranges for the False Positives per Image (FPPI) (which are unbounded to the top) used to draw the Free-Response Receiver Operating Characteristic (FROC) curve for the same prediction (top). The resulting FROC Scores differ for different boundaries of the x-axis used for the FPPI ([0, 1], [0, 2] and [0, 4]). Publications make use of different ranges for the x-axis, complicating comparison between works.

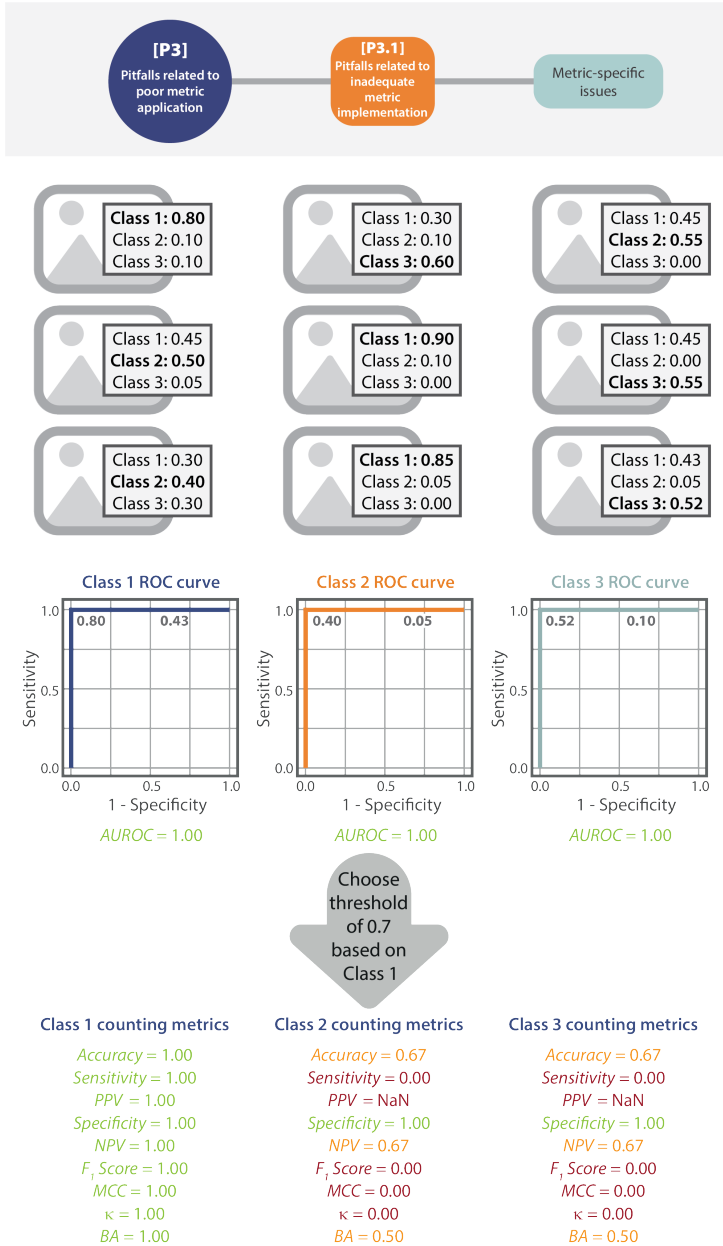


Choice of hyperparameters may have largely impact metric scores



Extended Data Fig. SN 2.25. Effect of the Intersection over Union (IoU) threshold on the localization (here Box IoU). **(a)** When defining a True Positive (TP) by a very loose $IoU > 0$, the resulting localizations may be deceived by very large predictions. **(b)** On the other hand, a strict IoU criterion may be problematic when the bounding box does not approximate the target structure shapes well. Although *Predictions 1* and *2* are very similar (differing in one pixel in one dimension), only *Prediction 1* is a TP because the number of bounding box pixels increases quadratically with the size of diagonal narrow structures. Further abbreviation: False Positive (FP).

Per-class tuning of the decision threshold yields misleading results



Extended Data Fig. SN 2.26. Effect of the determination of a global threshold for all classes based on a single class. In a data set of three classes and nine images, the Area under the Receiver Operating Characteristic Curve (AUROC) score is 1.0 for every class. In practice, however, a global decision threshold needs to be set in multi-class problems, which typically renders substantially worse results. Here, the optimal threshold for *Class 1* yields poor results for *Classes 2* and *3* (see e.g., [23, 50]). Used abbreviations: Positive Predictive Value (PPV), Negative Predictive Value (NPV), Matthews Correlation Coefficient (MCC), Cohen’s Kappa κ , and Balanced Accuracy (BA).

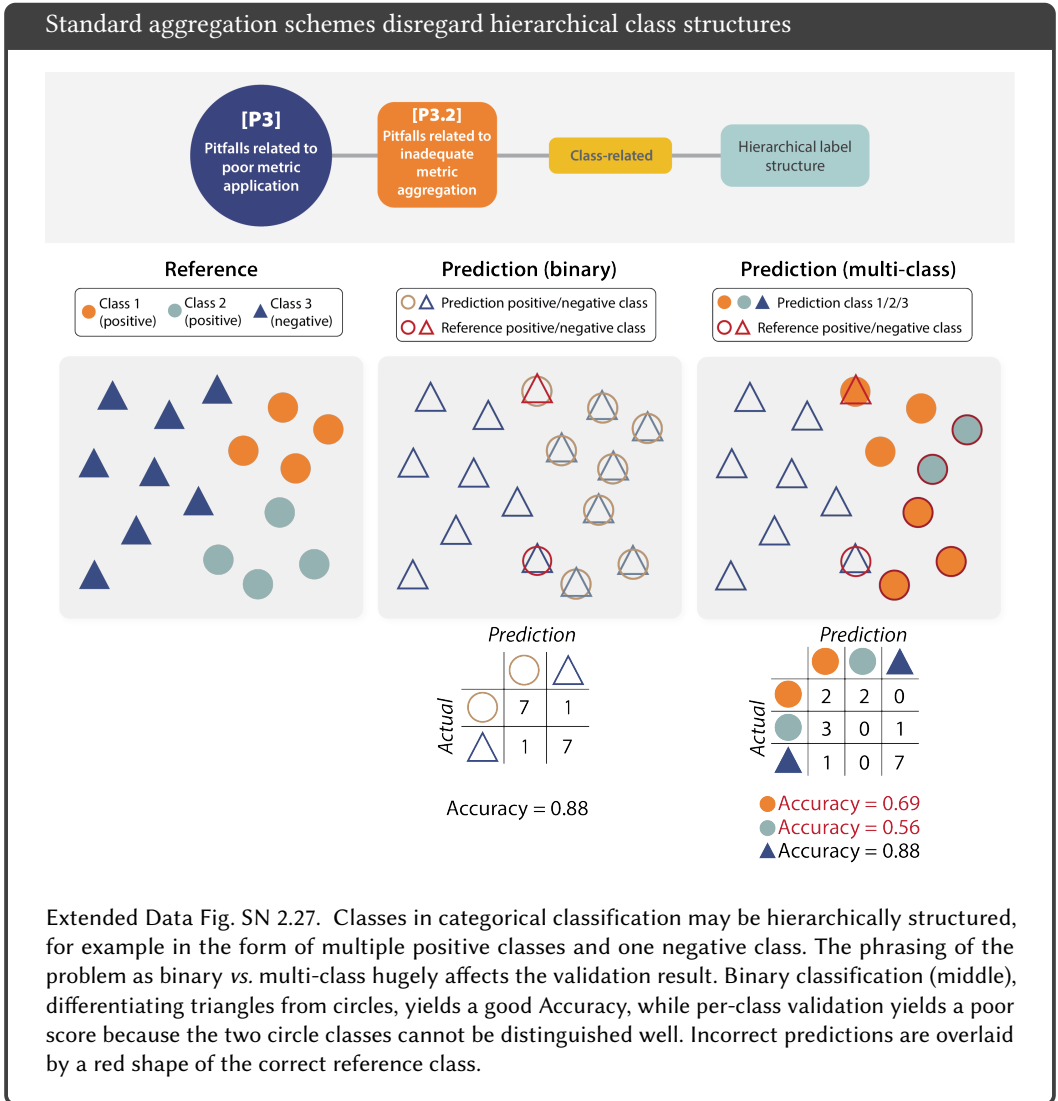
2.3.2 Pitfalls related to inadequate metric aggregation. When aggregating metric values over multiple cases (data points), the method of metric aggregation should be clearly defined and reported including details for example on the aggregation operator (e.g., mean or median) and missing value handling. In addition, special care should be taken when aggregating across classes or different hierarchy levels. Pitfalls can be further subdivided into *class-related* and *data set-related* pitfalls. In the following, we present pitfalls stemming from the following sources:

Class-related pitfalls

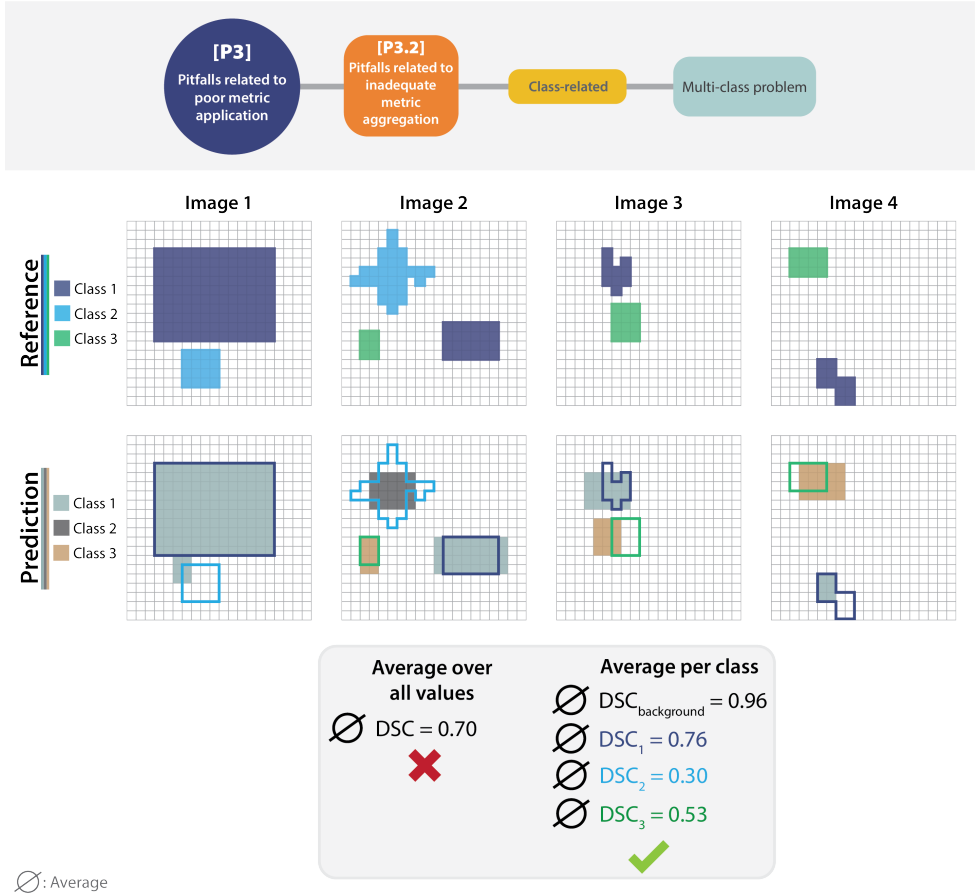
- Hierarchical label structure (Fig. SN 2.27)
- Multi-class problem (Fig. SN 2.28)

Data set-related pitfalls

- Non-independence of test cases (Figs. 6b and SN 2.29)
- Risk of bias (Fig. SN 2.30)
- Possibility of invalid prediction (Fig. SN 2.31)

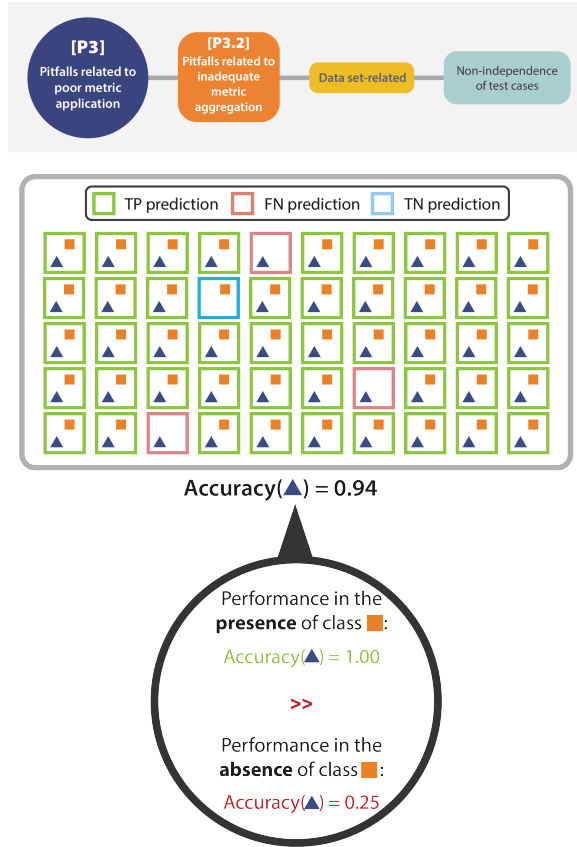


Lack of per-class validation conceals important information



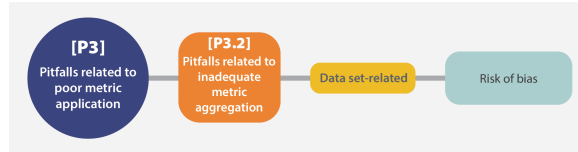
Extended Data Fig. SN 2.28. Effect of ignoring the presence of multiple classes when aggregating metric values (here: using the mean). The overall average of all Dice Similarity Coefficient (DSC) scores for the four images is 0.7. Averaging per class reveals a very low performance for *Classes 2 and 3*. \emptyset refers to the average DSC values.

Inter-class dependencies are concealed in standard aggregation schemes

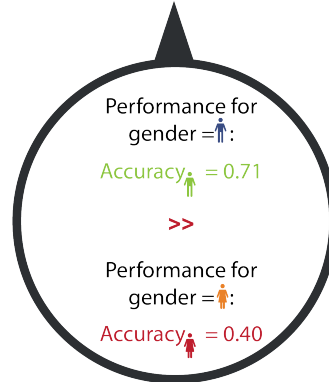


Extended Data Fig. SN 2.29. Effect of interdependencies between classes. A prediction may show a near-perfect Accuracy score of 0.94 for the dark blue triangle as it frequently appears in conjunction with the orange square. By calculating the Accuracy in the *presence* and *absence* of the orange square class, it can be seen that the algorithm only works well in the presence of the orange square class.

Lack of stratification conceal biases



Accuracy = 0.58



Extended Data Fig. SN 2.30. Effect of disregarding relevant meta-information (here: gender). When ignoring the available meta-information of the patient's gender per image, any metric (here: *Accuracy*) fails to reveal that the algorithm performs much better for men compared to women. In this example, correct predictions are marked by a green check mark, incorrect predictions by a red cross.

Lack of missing data handling strategy yields misleading results

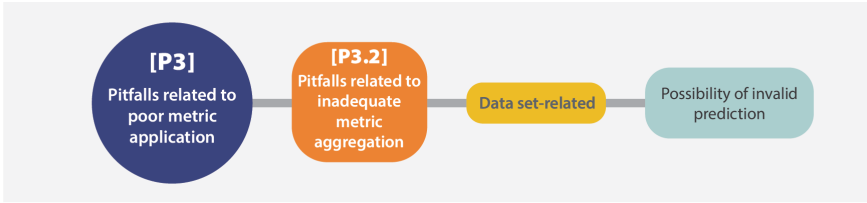


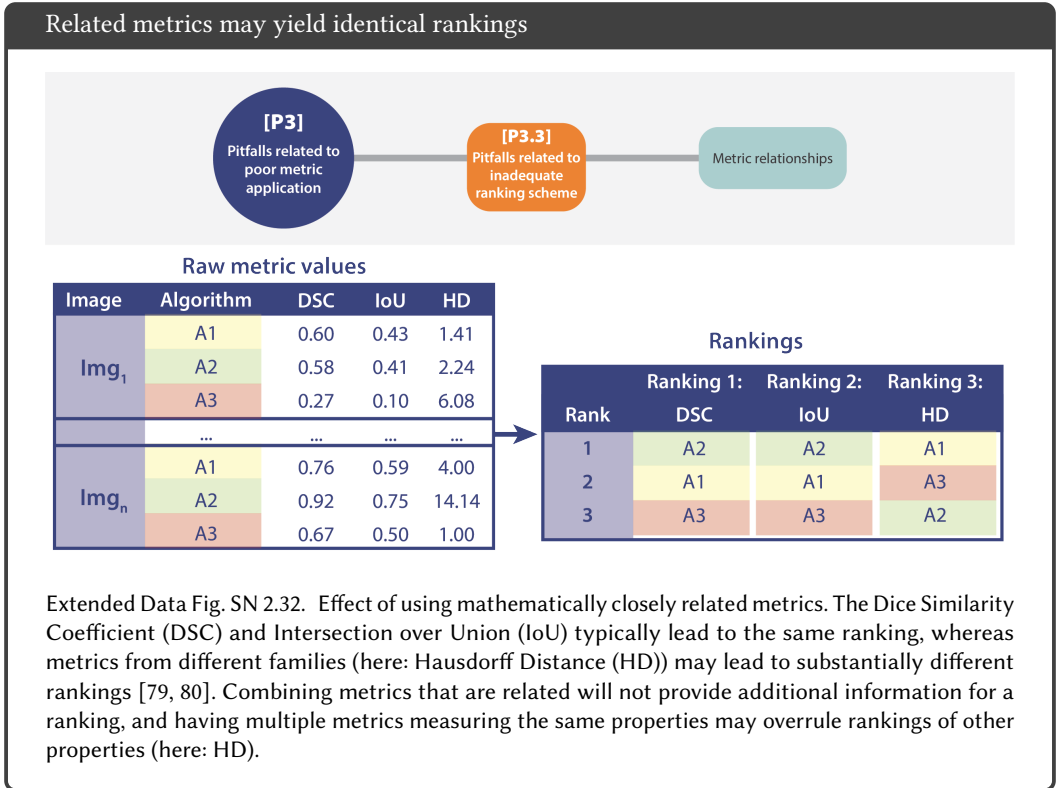
Image	I ₁	I ₂	I ₃	I ₄	I ₅	I ₆
DSC	0.94	NaN	0.87	0.90	NaN	0.89



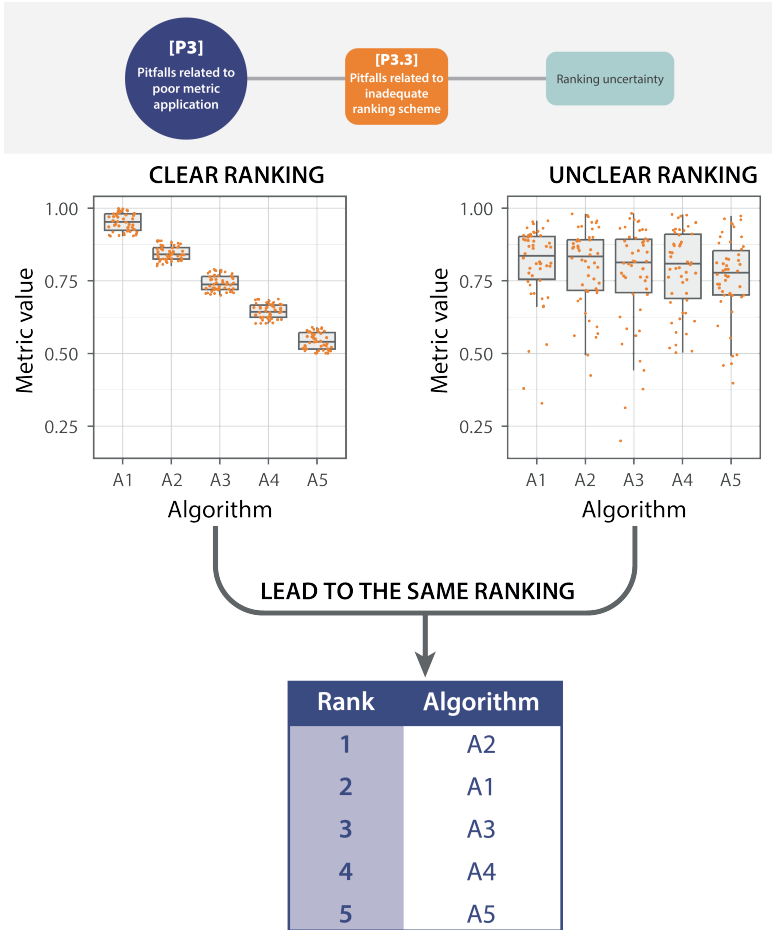
Extended Data Fig. SN 2.31. Effect of invalid predictions (missing values) when aggregating metric values. In this example, ignoring missing values leads to a substantially higher Dice Similarity Coefficient (DSC) compared to setting missing values to the worst possible value (here: 0).

2.3.3 Pitfalls related to inadequate ranking scheme. Rankings are often created to compare algorithm performances. In this context, we present pitfalls stemming from the following sources:

- Metric relationships (Fig. SN 2.32)
- Ranking uncertainty (Fig. SN 2.33)



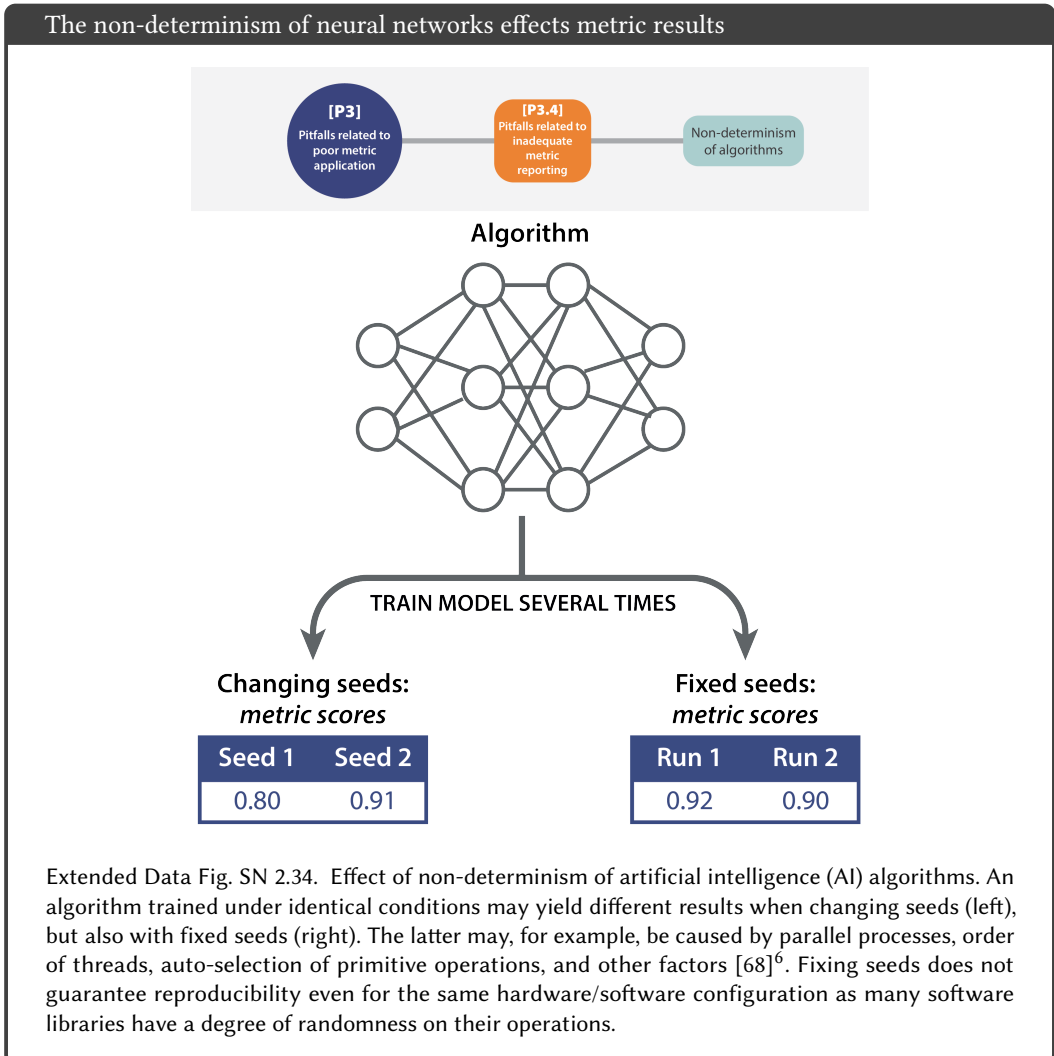
Ranking tables do not reflect ranking uncertainty



Extended Data Fig. SN 2.33. Effect of ranking uncertainty. The results of two benchmarking experiments with five algorithms A1-A5 differ substantially, as shown by the boxplots of the metric values for every algorithm. While the left situation introduces a clear ranking visible from the boxplots, the right use case is not clear as performance is very similar across algorithms. However, both situations lead to the same ranking [57, 93]. Thus, solely providing ranking tables conceals information on ranking uncertainty.

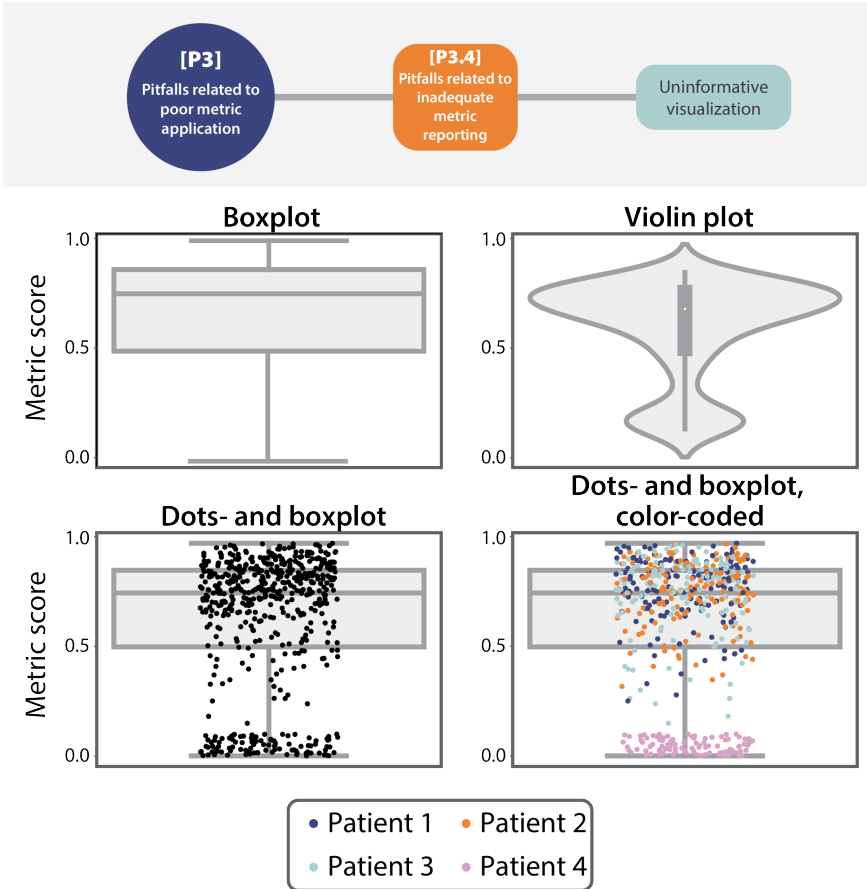
2.3.4 *Pitfalls related to inadequate metric reporting.* A thorough reporting of metric values and aggregates is important both in terms of transparency and interpretability. However, several pitfalls are to be avoided in this regard. Sources of metric reporting pitfalls include:

- Non-determinism of algorithms (Fig. SN 2.34)
- Uninformative visualization (Figs. 6c and SN 2.35)



⁶See for example: <https://pytorch.org/docs/stable/notes/randomness.html>

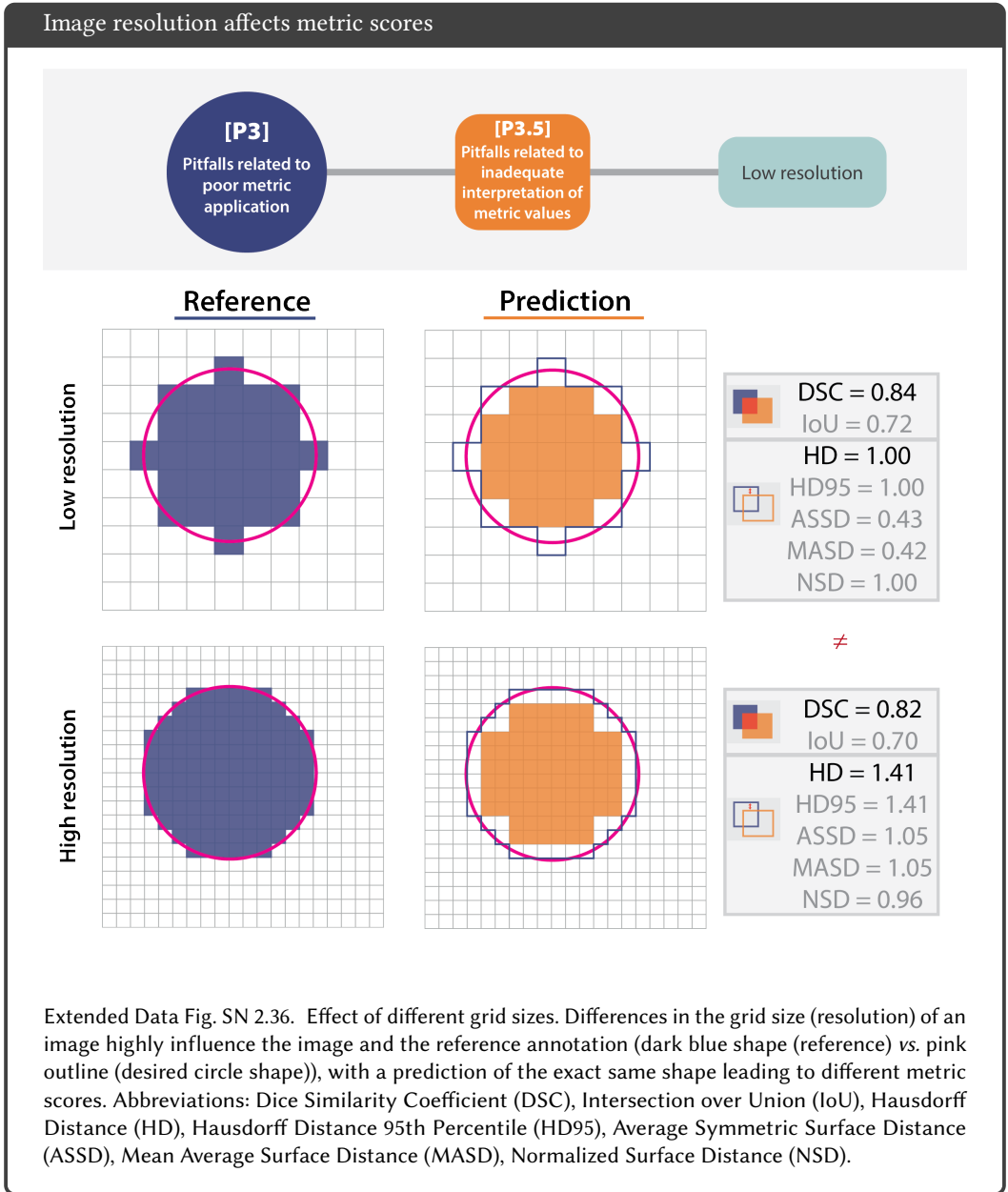
Common visualization schemes conceal relevant information



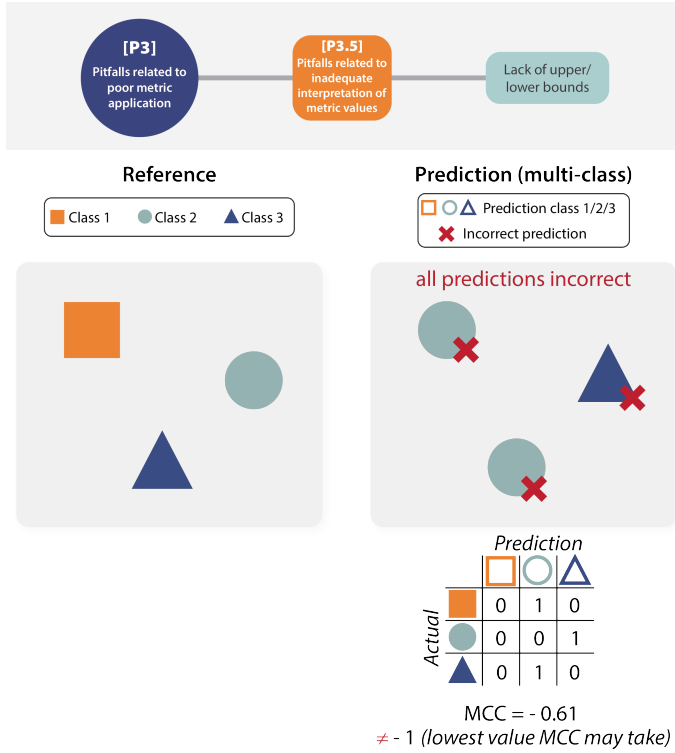
Extended Data Fig. SN 2.35. Effect of different visualization types. A single boxplot (top left) does not provide sufficient information about the raw metric value distribution (here: Dice Similarity Coefficient (DSC)). Using a violin plot (top right) or adding the raw metric values as jittered dots on top (bottom left) adds important information. In the case of non-independent validation data, color/shape-coding helps reveal data clusters (bottom right).

2.3.5 *Pitfalls related to inadequate interpretation of metric values.* Interpreting metric scores and aggregates is an important step in algorithm performance analysis. However, several pitfalls can arise from interpretation. In the following, we present pitfalls related to:

- Low resolution (Fig. SN 2.36)
- Lack of upper/lower bounds (Fig. SN 2.37)
- Insufficient domain relevance of metric score differences (Fig. SN 2.38)

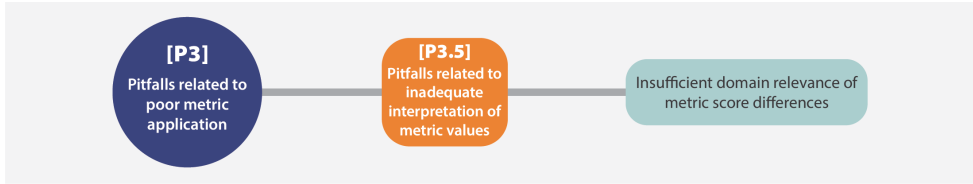


Lower bounds of metrics may not be achievable in practice



Extended Data Fig. SN 2.37. Effect of theoretical bounds that may not be achievable in practice. In this multi-class example, all samples were predicted incorrectly. However, the theoretical lowest value for the Matthews Correlation Coefficient (MCC) metric (-1) cannot be achieved in this situation, rendering interpretation difficult.

Metric score differences leading to different rankings may be irrelevant



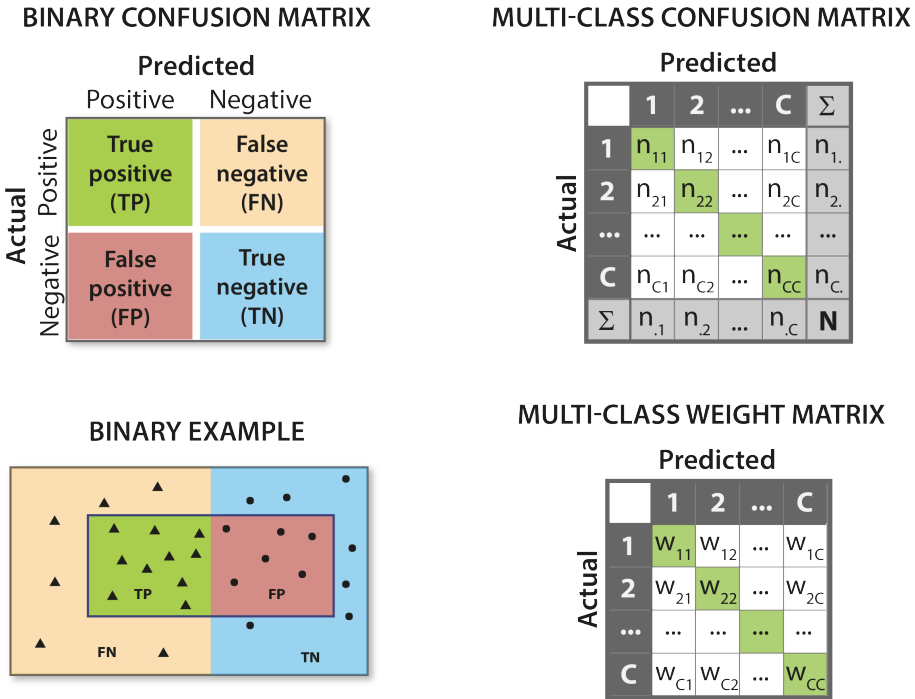
Rank	Algorithm	Aggregated metric score
1	A2	0.9543
2	A1	0.9542
3	A3	0.8703

} Difference: 0.0001 ❌
} Difference: 0.0839 ✅

Extended Data Fig. SN 2.38. Effect of irrelevant metric score differences in rankings. The difference of the metric score aggregates of algorithms A1 and A2 is extremely low and not of biomedical relevance. However, the numerical difference would assign them different ranks.

3 METRIC PROFILES

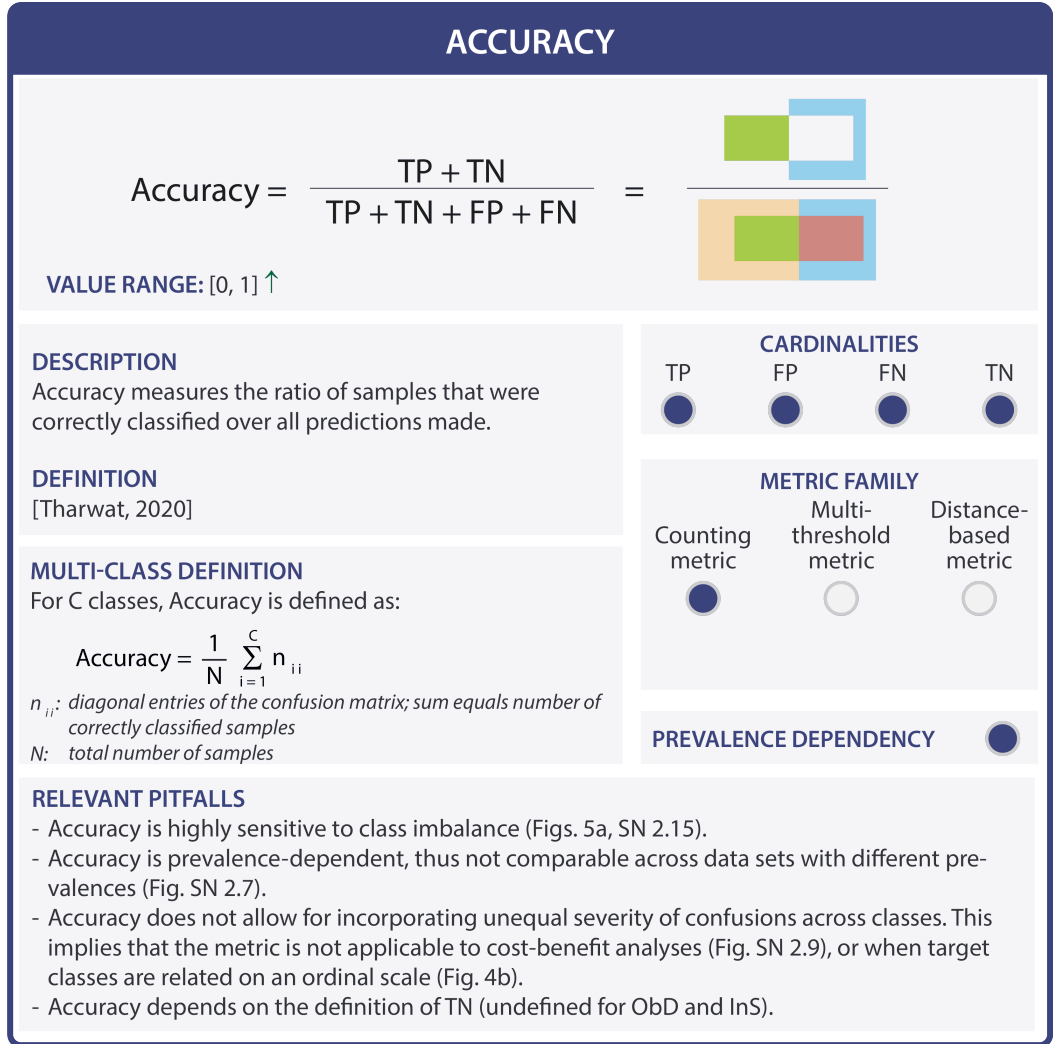
This section presents profiles for the metrics deemed particularly relevant by the *Metrics Reloaded* consortium [58]. For each metric, the respective description, formula, and value range (upward arrow: higher values better than lower values; downward arrow: lower values are better than higher values) are provided, along with further important characteristics, such as the used cardinalities of a confusion matrix, or potential prevalence dependency. Finally, relevant pitfalls are highlighted. Many of the presented metrics rely on the confusion matrix, which is illustrated in Fig. SN 3.39.



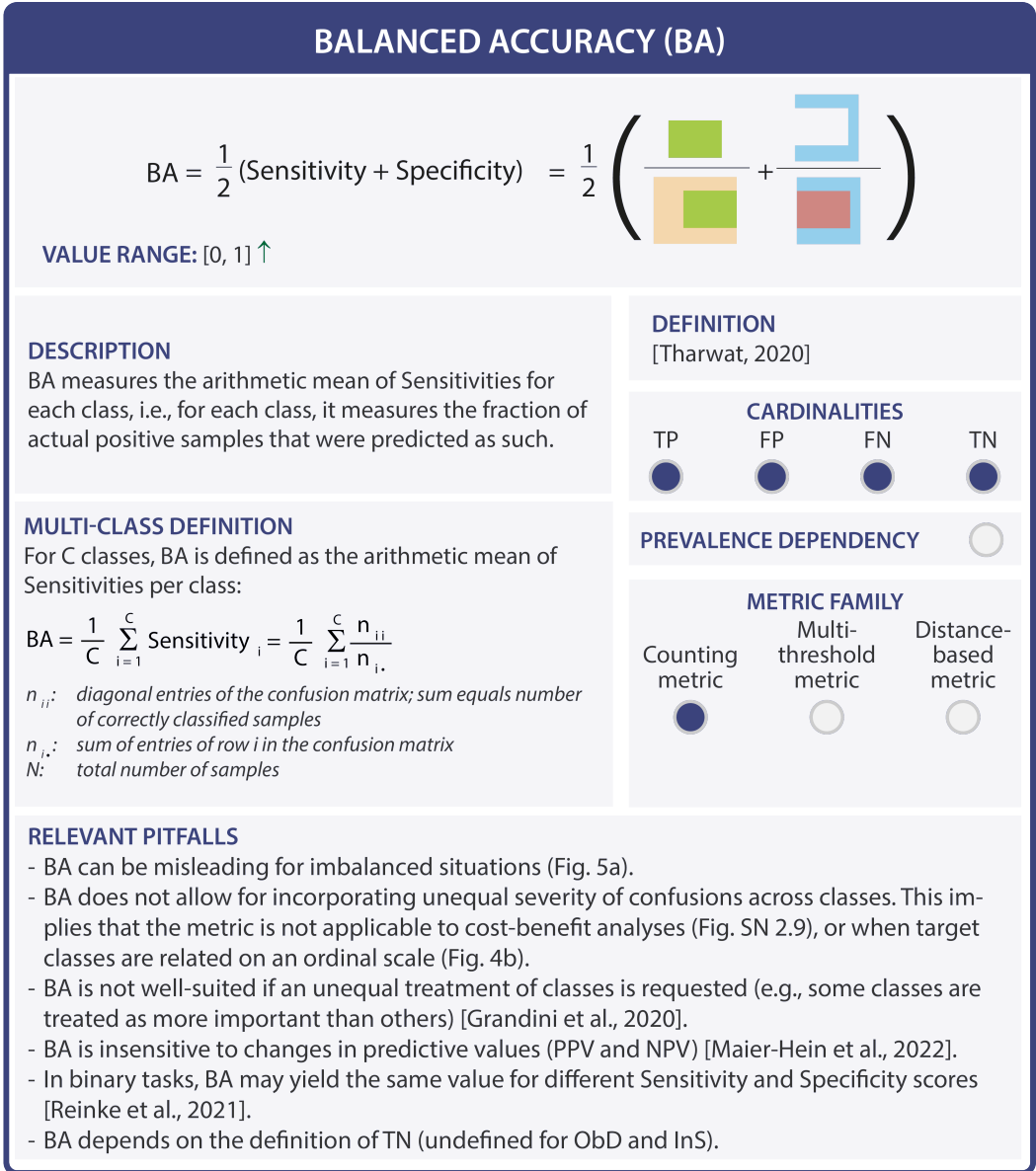
Extended Data Fig. SN 3.39. Schematic example of the confusion matrix for two and for C classes. For the latter case, we also present a weight or cost matrix with weights $w_{ij} > 0$ without loss of generality. For the binary confusion matrix, we show an example illustrating the cardinalities for a prediction of triangles and circles.

3.1 Discrimination metrics

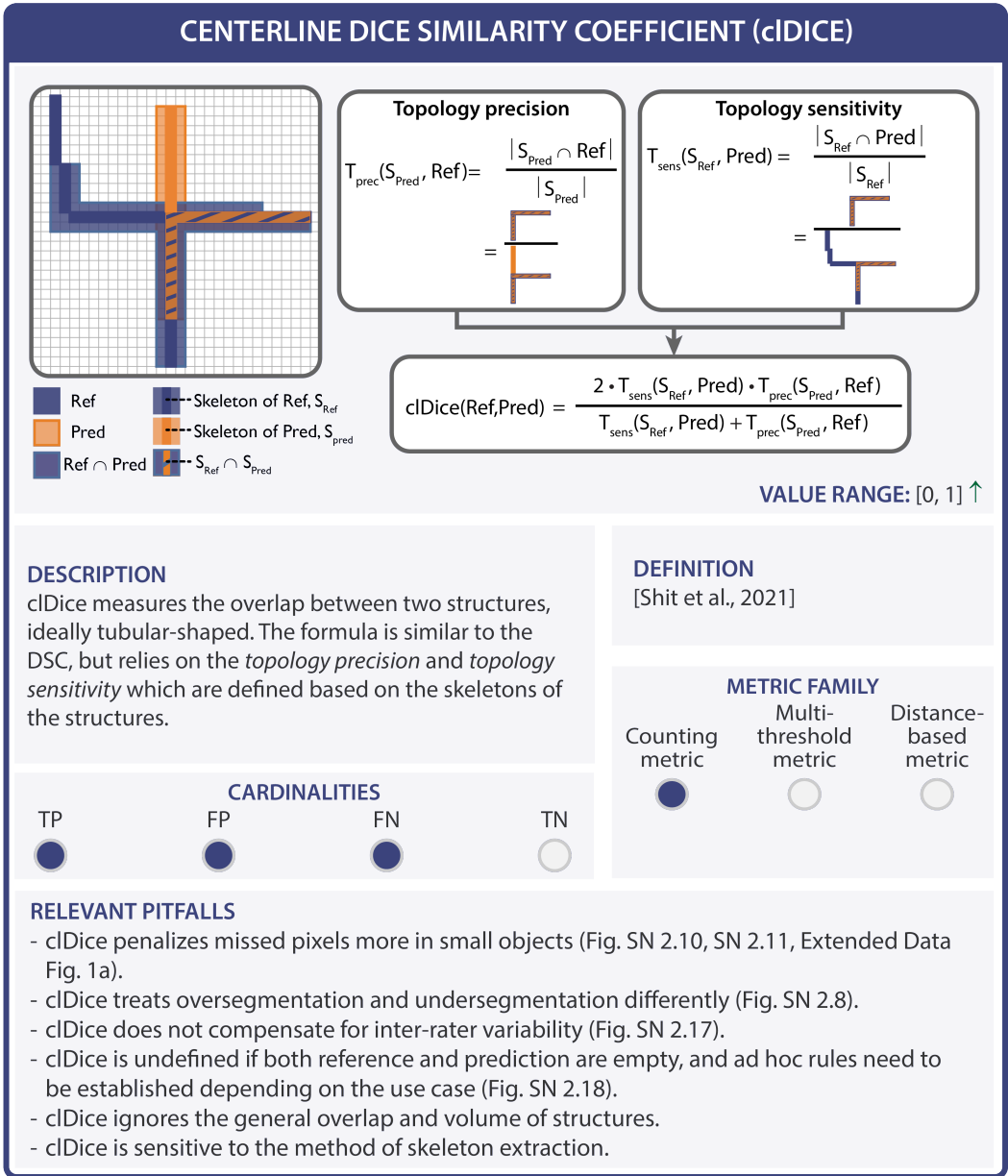
3.1.1 Counting metrics.



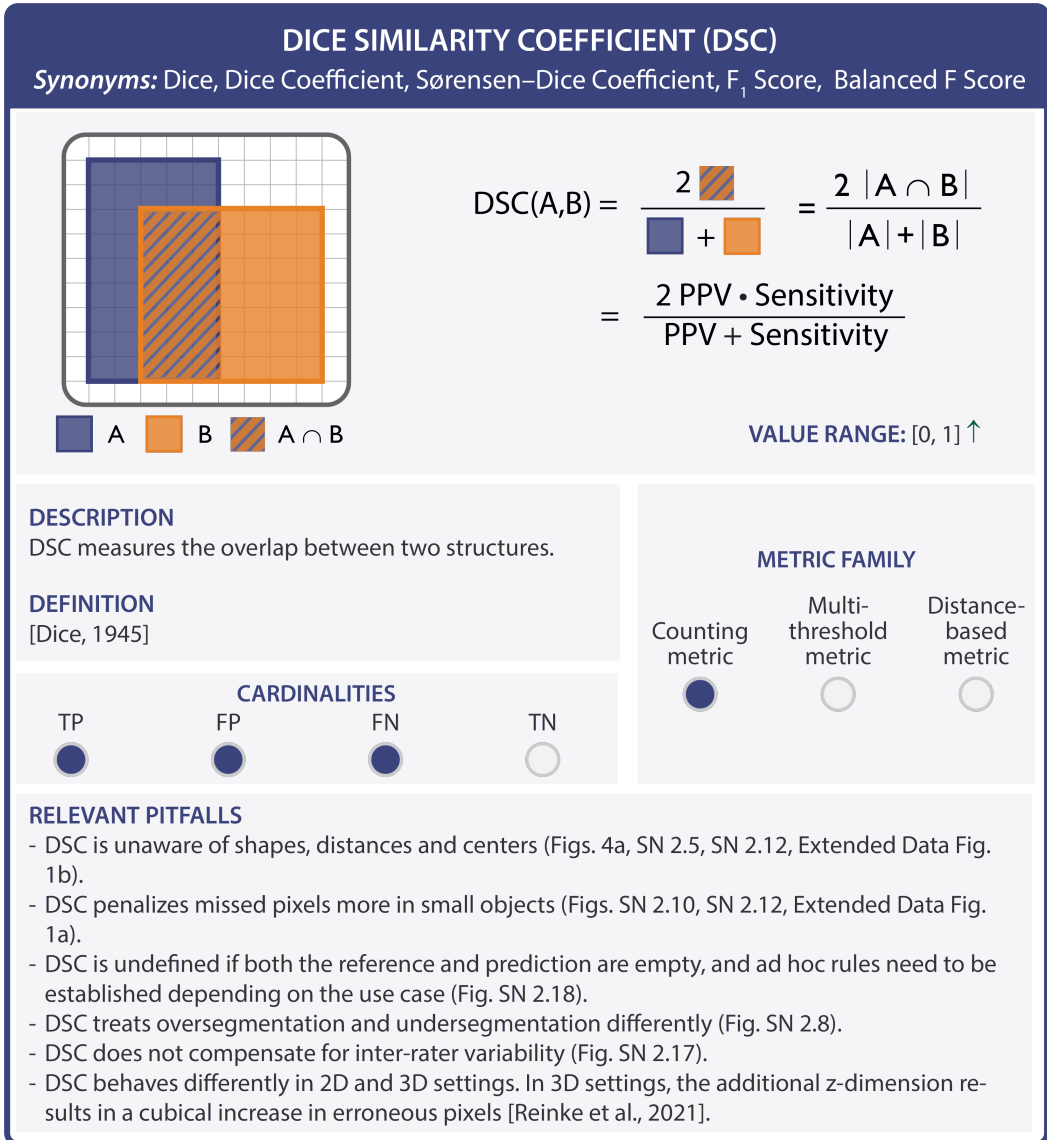
Extended Data Fig. SN 3.40. Metric profile of Accuracy. The upward arrow in the value range indicates that higher values are better than lower values. Abbreviations: False Negative (FN), False Positive (FP), Instance Segmentation (InS), Object Detection (ObD), True Negative (TN), True Positive (TP). Reference: Tharwat, 2020: [81]. Mentioned figures: Figs. 4b, 5a, SN 2.9, SN 2.11, SN 2.17.



Extended Data Fig. SN 3.41. Metric profile of Balanced Accuracy (BA). The upward arrow in the value range indicates that higher values are better than lower values. Abbreviations: False Negative (FN), False Positive (FP), Instance Segmentation (InS), Object Detection (ObD), Positive Predictive Value (PPV), True Negative (TN), True Positive (TP). References: Grandini et al., 2020: [36], Maier-Hein et al., 2022: [58], Reinke et al., 2021: [71], Tharwat, 2020: [81]. Mentioned figures: Figs. 4b, 5a, SN 2.11.



Extended Data Fig. SN 3.42. Metric profile of centerline Dice Similarity Coefficient (cIDice). The upward arrow in the value range indicates that higher values are better than lower values. Abbreviations: False Negative (FN), False Positive (FP), True Negative (TN), True Positive (TP). Reference: Shit et al., 2021: [77]. Mentioned figures: Extended Data Fig. 1a, Figs. SN 2.10, SN 2.12, SN 2.13, SN 2.19, SN 2.20.



Extended Data Fig. SN 3.43. Metric profile of Dice Similarity Coefficient (DSC). The upward arrow in the value range indicates that higher values are better than lower values. Abbreviations: False Negative (FN), False Positive (FP), Positive Predictive Value (PPV), True Negative (TN), True Positive (TP). References: Dice, 1945: [27], Reinke et al., 2021: [71]. Mentioned figures: Figs. 4a, SN 2.7, SN 2.10, SN 2.12, SN 2.13, SN 2.14, SN 2.19, SN 2.20, Extended Data Fig. 1a-b.

EXPECTED COST (EC)/NORMALIZED EC (ECN)

Synonyms: Expected prediction error, Expected loss

$$EC = w_{miss} \cdot \frac{FN}{TP + FN} \cdot \frac{TP + FN}{TP + TN + FP + FN} + w_{FA} \cdot \frac{FP}{TN + FP} \cdot 1 - \frac{TP + FN}{TP + TN + FP + FN}$$

$$= w_{miss} \cdot \frac{\text{[Diagram: Orange U-shape]} \cdot \text{[Diagram: Orange and Green Box]}}{\text{[Diagram: Orange and Green Box]}} + w_{FA} \cdot \frac{\text{[Diagram: Red Box]} \cdot \text{[Diagram: Red and Blue Box]}}{\text{[Diagram: Red and Blue Box]}} \cdot 1 - \frac{\text{[Diagram: Orange and Green Box]}}{\text{[Diagram: Orange, Green, and Red Box]}}$$

P_{miss} : FN (miss) rate, P_{FA} : FP (false alarm) rate
 P_{tar} : prior probability (prevalence)
 w_{miss}/w_{FA} : (estimation of) costs of the respective errors; can be adjusted as a weighting of them.

VALUE RANGE: $[0, \infty)$ ↓
EC can be assumed to be positive if costs are non-negative, which can be done without loss of generality.

DESCRIPTION

EC is a generalization of the probability of error (which is, in turn, 1 - Accuracy) for cases in which errors cannot all be considered to have equally severe consequences. It is defined as the expectation of the cost, where the cost incurred on a certain sample depends on the sample's class and the decision made for that sample. In practice, the expectation can be estimated as a simple average of the costs over the evaluation samples. EC describes the weighted sum of error rates. It can be used to measure discrimination and calibration in one score.

VARIANT

Normalized EC (ECN): normalizes EC by the EC of a naive system.

DEFINITION
[Bishop and Nasrabadi, 2006; Hastie et al., 2009; Ferrer, 2022]

PREVALENCE DEPENDENCY?

Both options are possible depending on how the priors are set in the definition of the metric.

MULTI-CLASS DEFINITION
For C classes, EC is defined as:

$$EC = \sum_{i=1}^C \sum_{j=1}^C P_i \cdot w_{ij} \cdot \frac{n_{ij}}{n_i}$$

n_{ij} : entry of the confusion matrix for row i and column j, i.e., samples of actual class i that have been predicted as class j
 n_i : sum of entries of row i of the confusion matrix
 w_{ij} : costs for the entry of the confusion matrix for row i and column j, i.e., the cost for predicting a sample of actual class i that was predicted as class j
 P_i : prevalence of class i; usually (n_i / N) , but in some cases one might want to plug in P_i directly from a target application

CARDINALITIES

TP FP FN TN

METRIC FAMILY

Counting metric Multi-threshold metric Distance-based metric

RELEVANT PITFALLS

- EC is rather uncommon and can therefore not be used for comparison with other publications.
- EC can be close to optimal even for poor predictive values (PPV and NPV) [Maier-Hein et al., 2022]
- ECN cannot be configured to ensure equal class contribution without losing its ability to ensure high predictive values [Maier-Hein et al., 2022].

Extended Data Fig. SN 3.44. Metric profile of Expected Cost (EC). The downward arrow in the value range indicates that lower values are better than higher values. Abbreviations: False Negative (FN), False Positive (FP), True Negative (TN), True Positive (TP). References: Bishop and Nasrabadi, 2006: [8], Ferrer 2022: [31], Hastie et al., 2009: [41], Maier-Hein et al., 2022: [58].

F_β SCORE

$$F_{\beta} \text{ Score} = (1+\beta^2) \frac{\text{Precision} \cdot \text{Sensitivity}}{\beta^2 \cdot \text{Precision} + \text{Sensitivity}}$$

$$= \frac{(1+\beta^2) \cdot \text{TP}}{(1+\beta^2) \cdot \text{TP} + \beta^2 \cdot \text{FN} + \text{FP}} = \frac{(1+\beta^2) \cdot \text{TP}}{(1+\beta^2) \cdot \text{TP} + \beta^2 \cdot \text{FN} + \text{FP}}$$

VALUE RANGE: [0, 1] ↑

DESCRIPTION

The F_β Score weights PPV (FP) and Sensitivity (FN) with the parameter β.

The special case of β = 1 is the harmonic mean of PPV and Sensitivity and is a common metric in segmentation problems (here usually referred to as DSC). In segmentation problems, F_β Score weights the penalization of oversegmentation (FP) and undersegmentation (FN) with the parameter β.

DEFINITION
[Van Rijsbergen, 1979; Chinchor 1992]

CARDINALITIES

TP	FP	FN	TN
<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>

PREVALENCE DEPENDENCY

METRIC FAMILY

Counting metric	Multi-threshold metric	Distance-based metric
<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>

RELEVANT PITFALLS

F_β Score for classification/detection assessment:

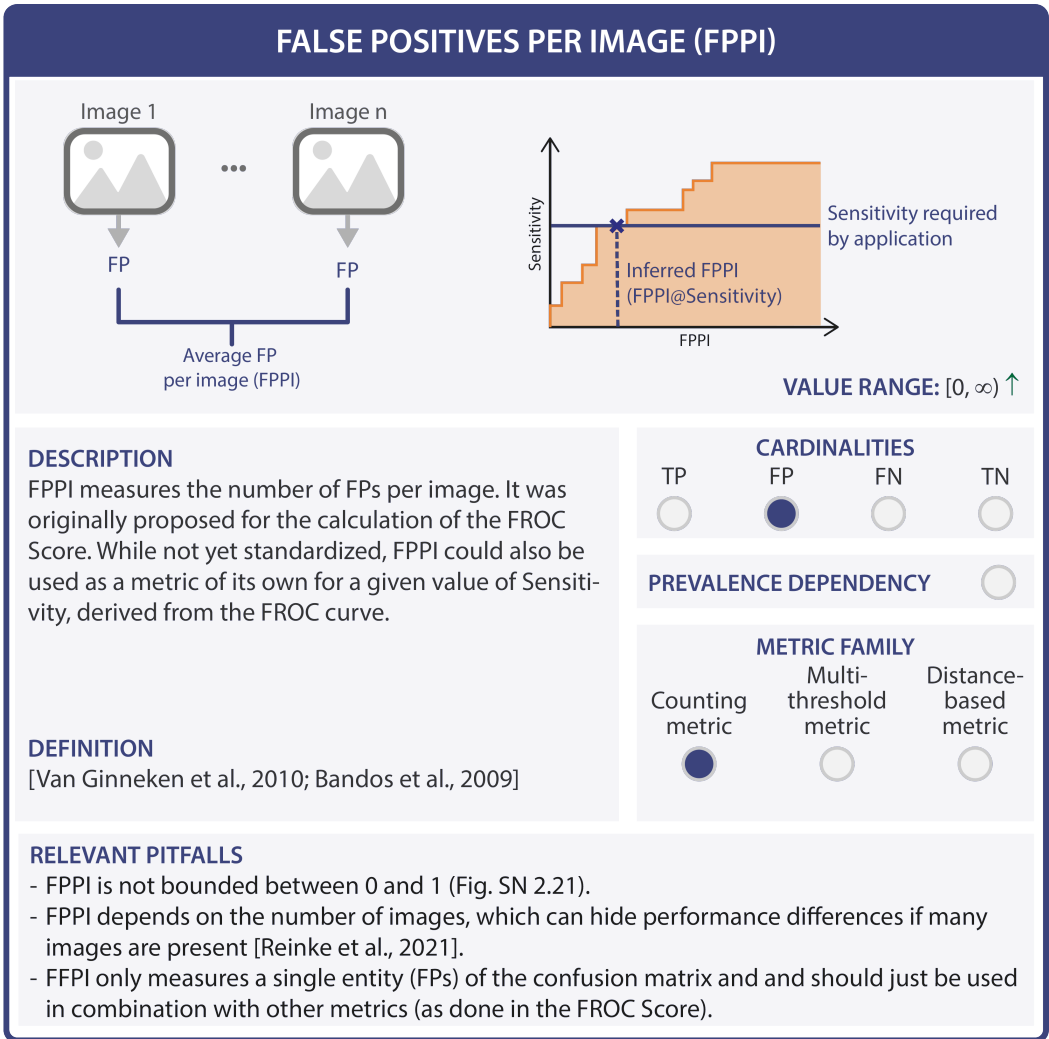
- F_β Score is prevalence-dependent, thus not comparable across data sets with different prevalences (Figs. SN 2.7, SN 2.15).
- Compared to other per-class counting metrics (e.g., LR+) it lacks the interpretability with respect to a naive classifier
- F_β Score depends on the definition of the positive class [Reinke et al., 2021].

F_β Score for segmentation assessment:

- F_β Score is unaware of the structure shape and center (Figs. 4a, SN 2.5, SN 2.12, Extended Data Fig. 1b).
- F_β Score penalizes missed pixels more in small objects (Fig. SN 2.10, Extended Data Fig. 1a).
- F_β Score does not compensate for inter-rater variability (Fig. 2.17).
- F_β Score behaves differently in 2D and 3D settings. In 3D settings, the additional z-dimension results in a cubical increase in erroneous pixels [Reinke et al., 2021].

F_β Score is undefined if both reference and prediction are empty, and ad hoc rules need to be established depending on the use case (Fig. SN 2.18, Extended Data Fig. 2b).

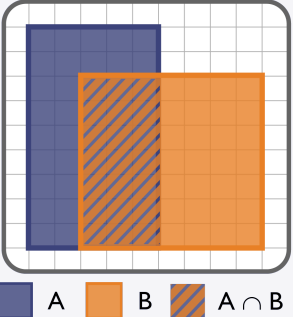
Extended Data Fig. SN 3.45. Metric profile of F_β Score.[16, 86]. The upward arrow in the value range indicates that higher values are better than lower values. Abbreviations: Dice Similarity Coefficient (DSC), False Negative (FN), False Positive (FP), Positive Predictive Value (PPV), True Negative (TN), True Positive (TP).References: Chinchor 1992: [16], Reinke et al., 2021: [71], Van Rijsbergen, 1979: [86]. Mentioned figures: Figs. 4a, SN 2.7, SN 2.9, SN 2.12, SN 2.14, SN 2.17, SN 2.19, SN 2.20, Extended Data Figs. 1a-b and 2b.



Extended Data Fig. SN 3.46. Metric profile of False Positives per Image (FPPI). The upward arrow in the value range indicates that higher values are better than lower values. Abbreviations: False Negative (FN), False Positive (FP), Free-Response Receiver Operating Characteristic (FROC), True Negative (TN), True Positive (TP). References: Bandos et al., 2009: [5], Reinke et al., 2021: [71], Van Ginneken et al., 2010: [85]. Mentioned figure: Fig. SN 2.23.

INTERSECTION OVER UNION (IoU)

Synonyms: Jaccard Index, Tanimoto Coefficient



$$IoU(A,B) = \frac{\text{Intersection}}{\text{Union}}$$

$$= \frac{|A \cap B|}{|A| + |B| - |A \cap B|} = \frac{|A \cap B|}{|A \cup B|}$$

$$= \frac{PPV \cdot Sensitivity}{PPV + Sensitivity - PPV \cdot Sensitivity}$$

VALUE RANGE: [0, 1] ↑

DESCRIPTION

IoU measures the overlap between two structures. It is often referred to as **Box IoU** when comparing bounding boxes, **Mask IoU** when comparing segmentation masks, or **Approx IoU** when comparing approximations of objects beyond bounding boxes.

DEFINITION

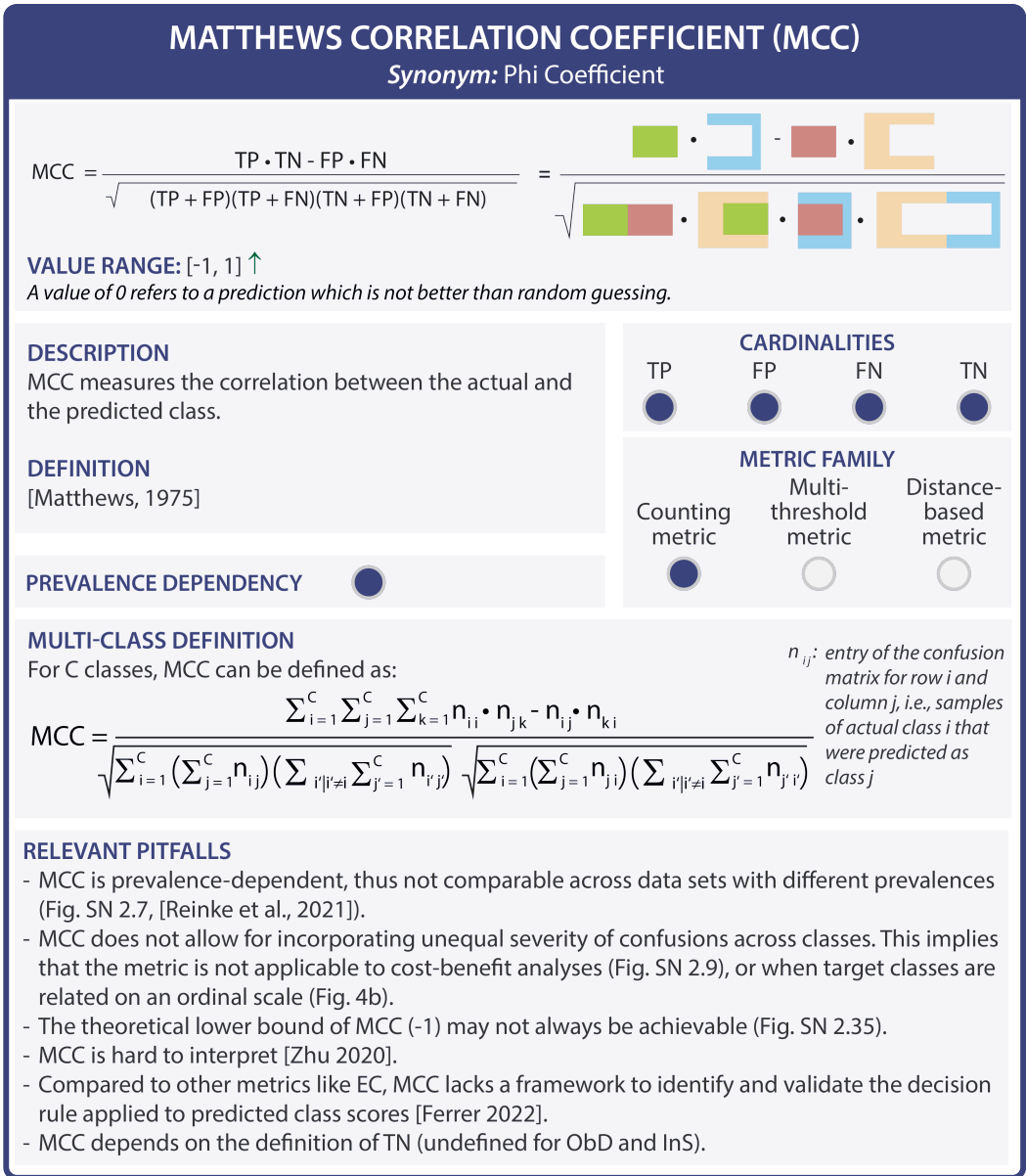
[Jaccard, 1912]

RELEVANT PITFALLS

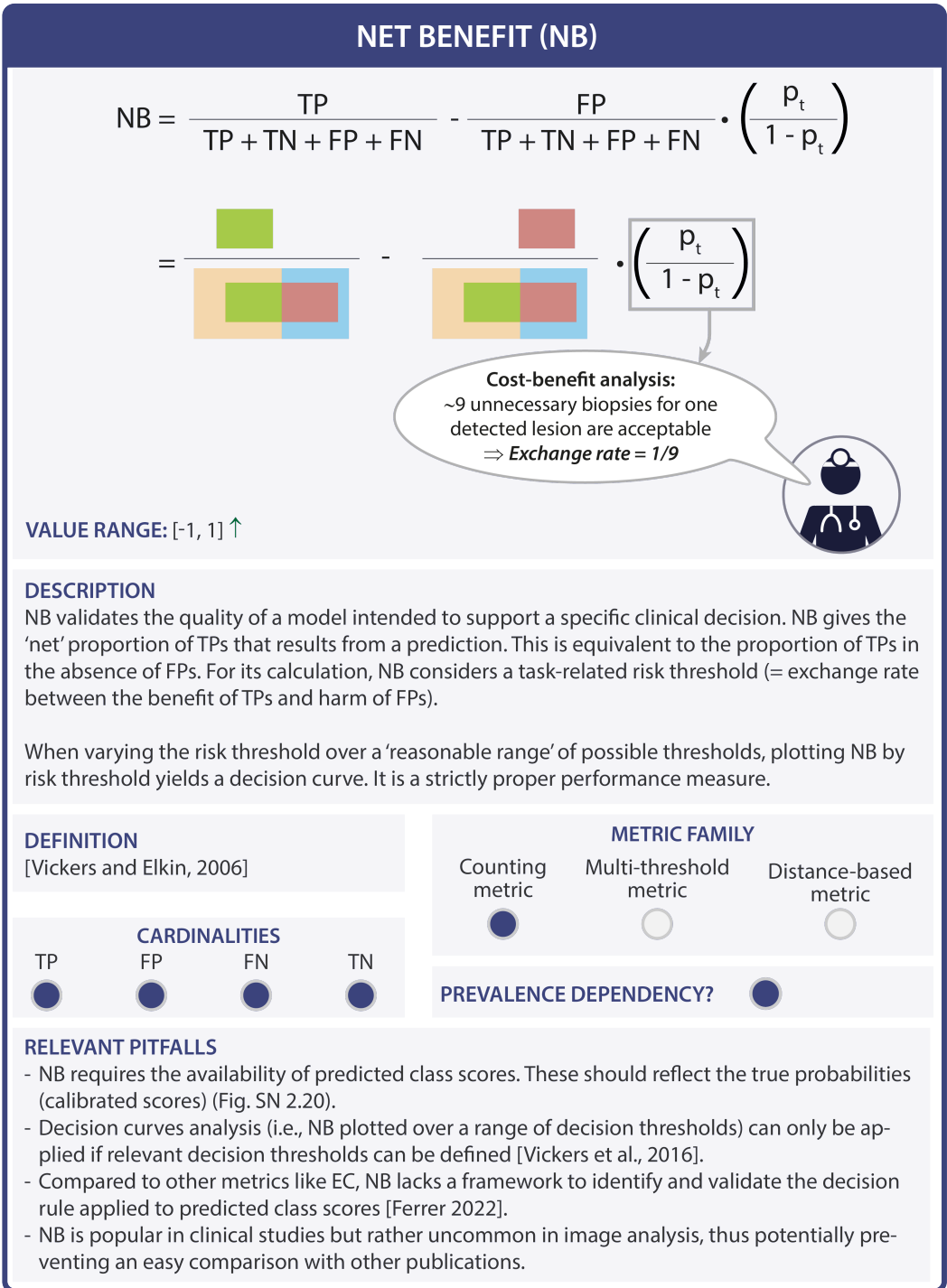
- IoU is unaware of shapes, boundaries, distances and centers (Figs. 4a, SN 2.5, SN 2.12, Extended Data Fig. 1b).
- IoU penalizes missed pixels more in small objects (Figs. SN 2.10, SN 2.11, Extended Data Fig. 1a).
- IoU is undefined if both the reference and prediction are empty, and ad hoc rules need to be established depending on the use case (Fig. SN 2.18).
- IoU treats oversegmentation and undersegmentation differently (Fig. SN 2.8).
- IoU does not compensate for inter-rater variability (Fig. SN 2.17).
- IoU behaves differently in 2D and 3D settings. In 3D settings, the additional z-dimension results in a cubical increase in erroneous pixels [Reinke et al., 2021].

CARDINALITIES			
TP	FP	FN	TN
<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>
METRIC FAMILY			
Counting metric	Multi-threshold metric	Distance-based metric	
<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	

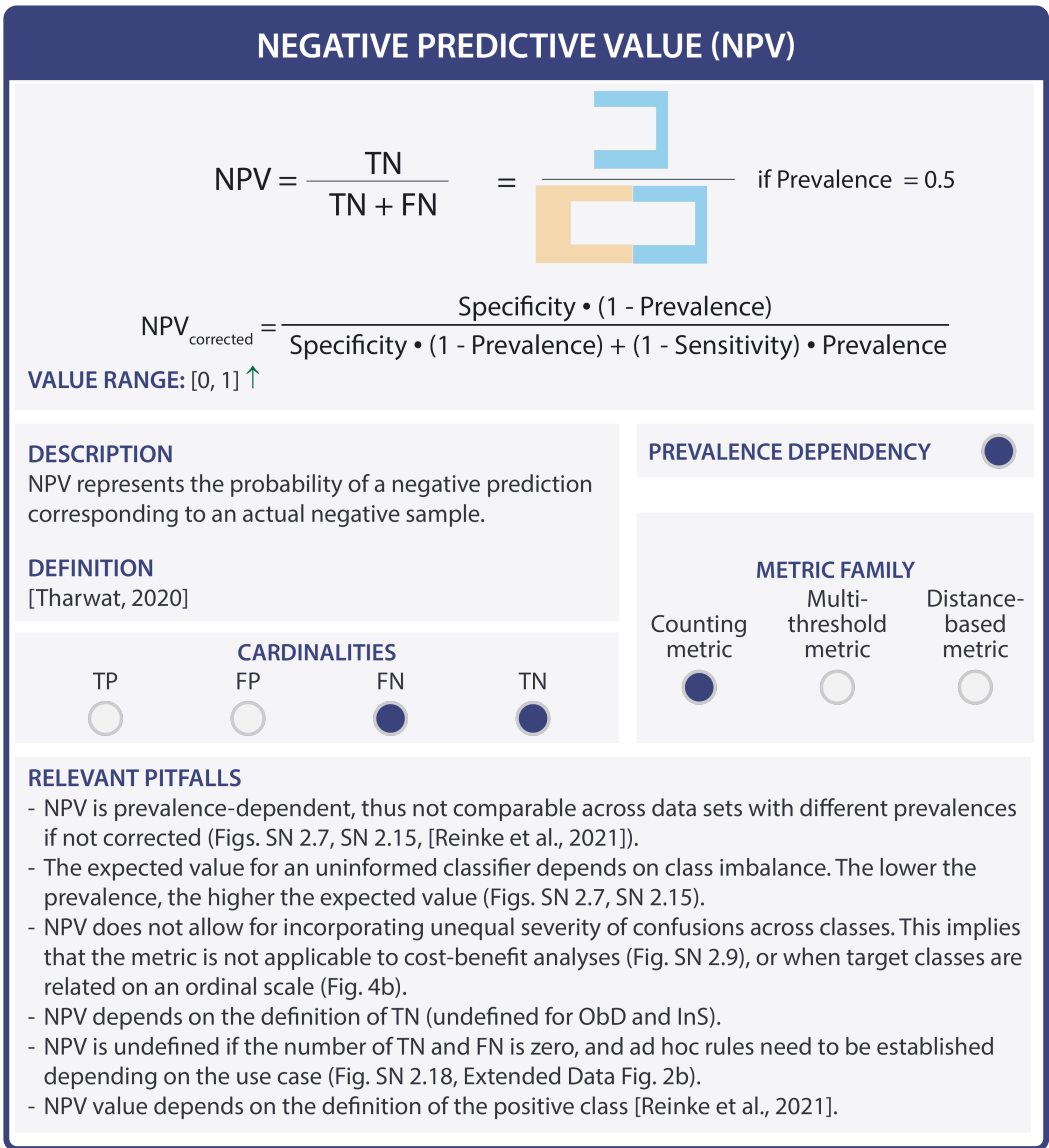
Extended Data Fig. SN 3.47. Metric profile of Intersection over Union (IoU). The upward arrow in the value range indicates that higher values are better than lower values. Abbreviations: False Negative (FN), False Positive (FP), Positive Predictive Value (PPV), True Negative (TN), True Positive (TP). References: Jaccard, 1912: [45], Reinke et al., 2021: [71]. Mentioned figures: Figs. 4a, SN 2.7, SN 2.10, SN 2.12, SN 2.13, SN 2.14, SN 2.19, SN 2.20, Extended Data Fig. 1a-b.



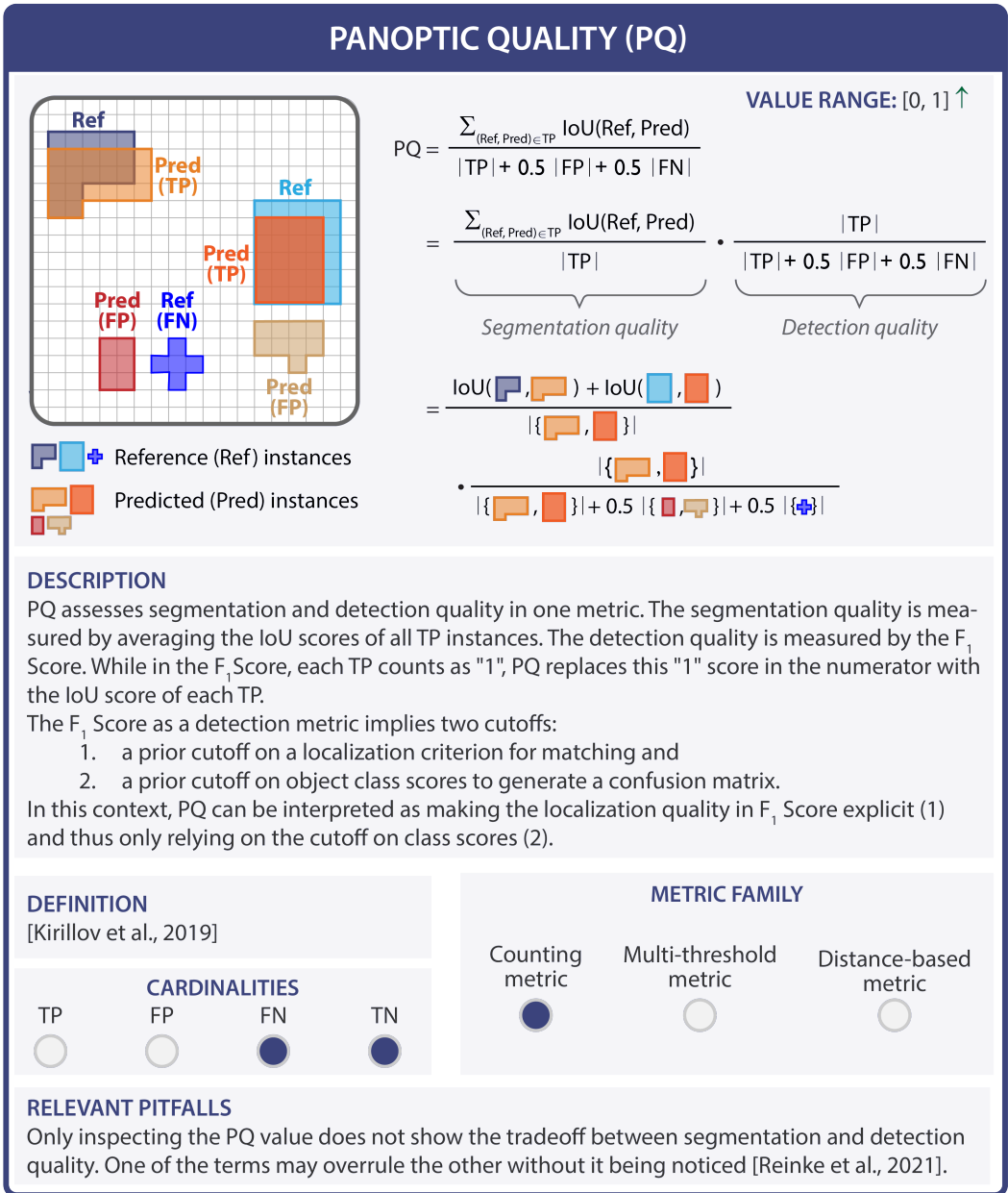
Extended Data Fig. SN 3.48. Metric profile of Matthews Correlation Coefficient (MCC). The upward arrow in the value range indicates that higher values are better than lower values. Abbreviations: Expected Cost (EC), False Negative (FN), False Positive (FP), Instance Segmentation (InS), Object Detection (ObD), True Negative (TN), True Positive (TP). References: Ferrer, 2022: [31], Matthews, 1975: [61], Reinke et al., 2021: [71], Zhu, 2020: [96]. Mentioned figures: Figs. 4b, SN 2.9, SN 2.11, SN 2.37.



Extended Data Fig. SN 3.49. Metric profile of Net Benefit (NB). The upward arrow in the value range indicates that higher values are better than lower values. Abbreviations: Expected Cost (EC), False Negative (FN), False Positive (FP), True Negative (TN), True Positive (TP). References: Ferrer, 2022: [31], Vickers and Elkin, 2006: [88], Vickers et al., 2016: [89]. Mentioned figure: Fig. SN 2.22.



Extended Data Fig. SN 3.50. Metric profile of Negative Predictive Value (NPV). The upward arrow in the value range indicates that higher values are better than lower values. Abbreviations: False Negative (FN), False Positive (FP), Instance Segmentation (InS), Object Detection (ObD), True Negative (TN), True Positive (TP). References: Reinke et al., 2021: [71], Tharwat, 2020: [81]. Mentioned figures: Figs. 4b, SN 2.9, SN 2.11, SN 2.17, SN 2.20, Extended Data Fig. 2b.



Extended Data Fig. SN 3.51. Metric profile of Panoptic Quality (PQ). The upward arrow in the value range indicates that higher values are better than lower values. Abbreviations: Average Precision (AP), False Negative (FN), False Positive (FP), Free-Response Receiver Operating Characteristic (FROC), Intersection over Union (IoU), True Negative (TN), True Positive (TP). References: Kirillov et al., 2019: [47], Reinke et al., 2021: [71].

POSITIVE LIKELIHOOD RATIO (LR+)

Synonyms: Likelihood ratio positive, Likelihood ratio for positive results

$$LR+ = \frac{\text{Sensitivity}}{1 - \text{Specificity}} = \frac{\text{Green Box}}{\text{Orange Box} - \text{Green Box}} \Bigg/ \left(1 - \frac{\text{Blue Box} - \text{Red Box}}{\text{Blue Box}} \right)$$

VALUE RANGE: $[0, \infty)$ ↑

DESCRIPTION

LR+ indicates the factor by which a positive prediction occurs more frequently among actual positive samples than among actual negative samples. In a clinical example where the quality of a diagnostic test is to be assessed, this could be interpreted as how much more likely a positive test result is for a diseased person compared to a healthy person (the higher the better).

DEFINITION
[Attia, 2003]

PREVALENCE DEPENDENCY

METRIC FAMILY

Counting metric	Multi-threshold metric	Distance-based metric
<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>

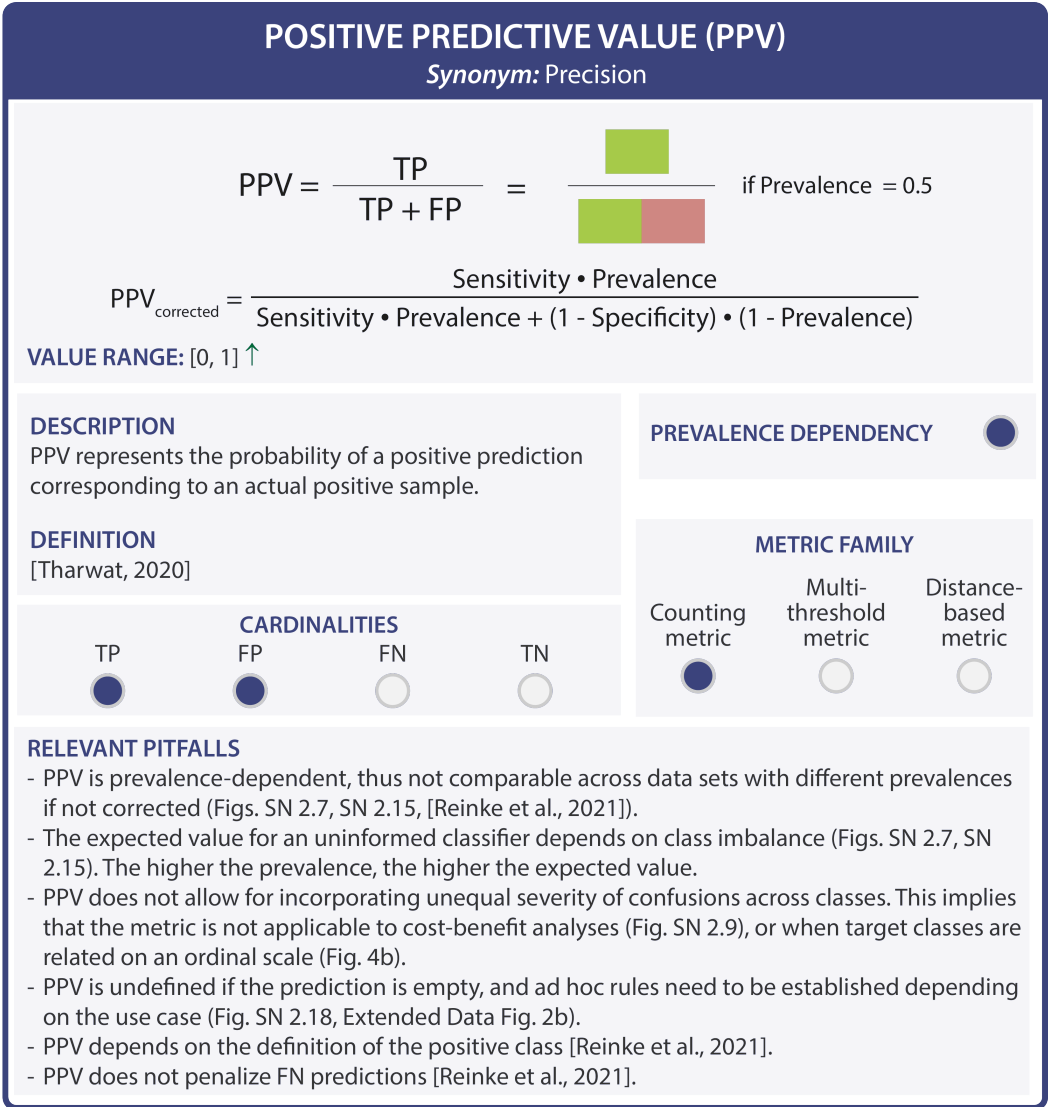
CARDINALITIES

TP	FP	FN	TN
<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>

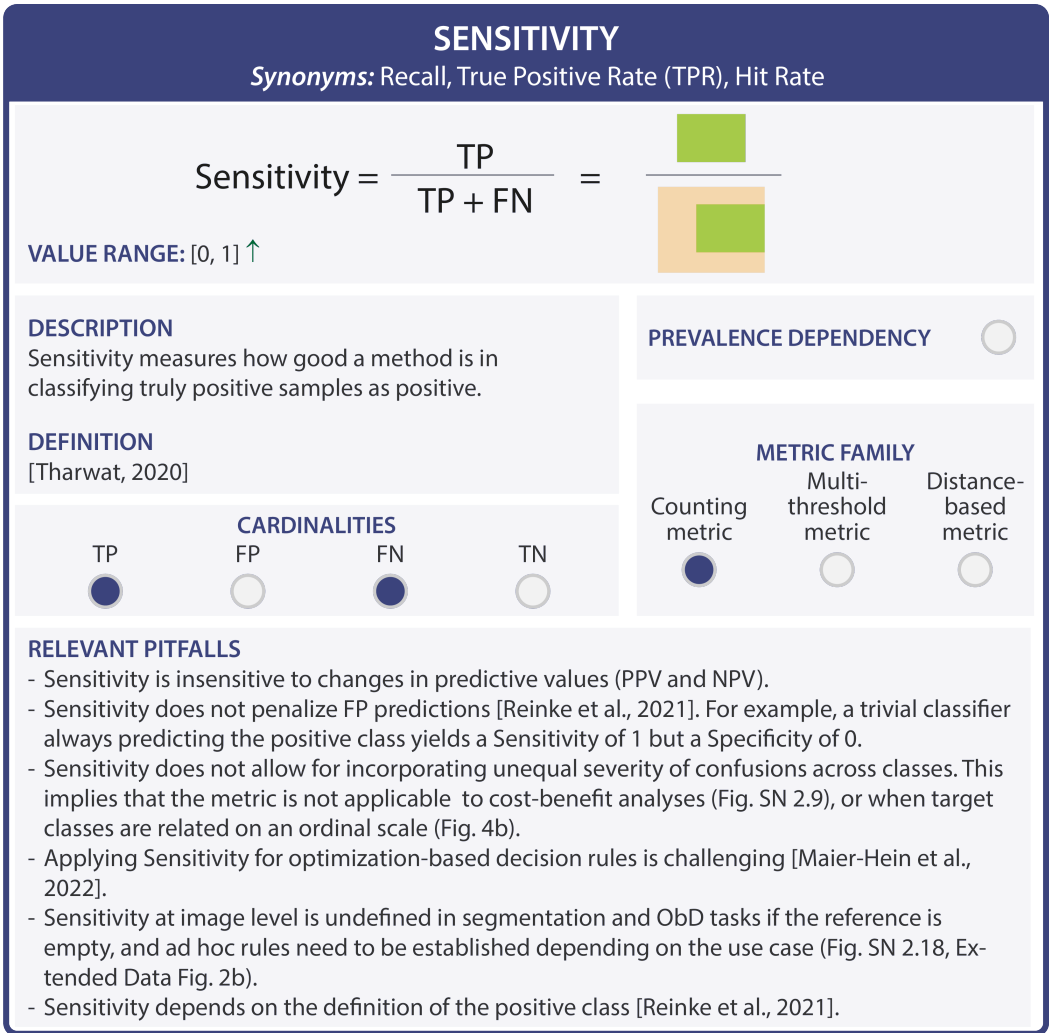
RELEVANT PITFALLS

- LR+ is insensitive to changes in predictive values (PPV and NPV; Fig. 5a).
- LR+ is undefined if the Specificity is 1, and ad hoc rules need to be established depending on the use case (Fig. SN 2.18, Extended Data Fig. 2b).
- LR+ does not allow for incorporating unequal severity of confusions across classes. This implies that the metric is not applicable to cost-benefit analyses (Fig. SN 2.9), or when target classes are related on an ordinal scale (Fig. 4b).
- LR+ depends on the definition of the positive class [Reinke et al., 2021].
- LR+ may yield the same value for different Sensitivity and Specificity scores [Reinke et al., 2021].
- LR+ depends on the definition of TN (undefined for ObD and InS).

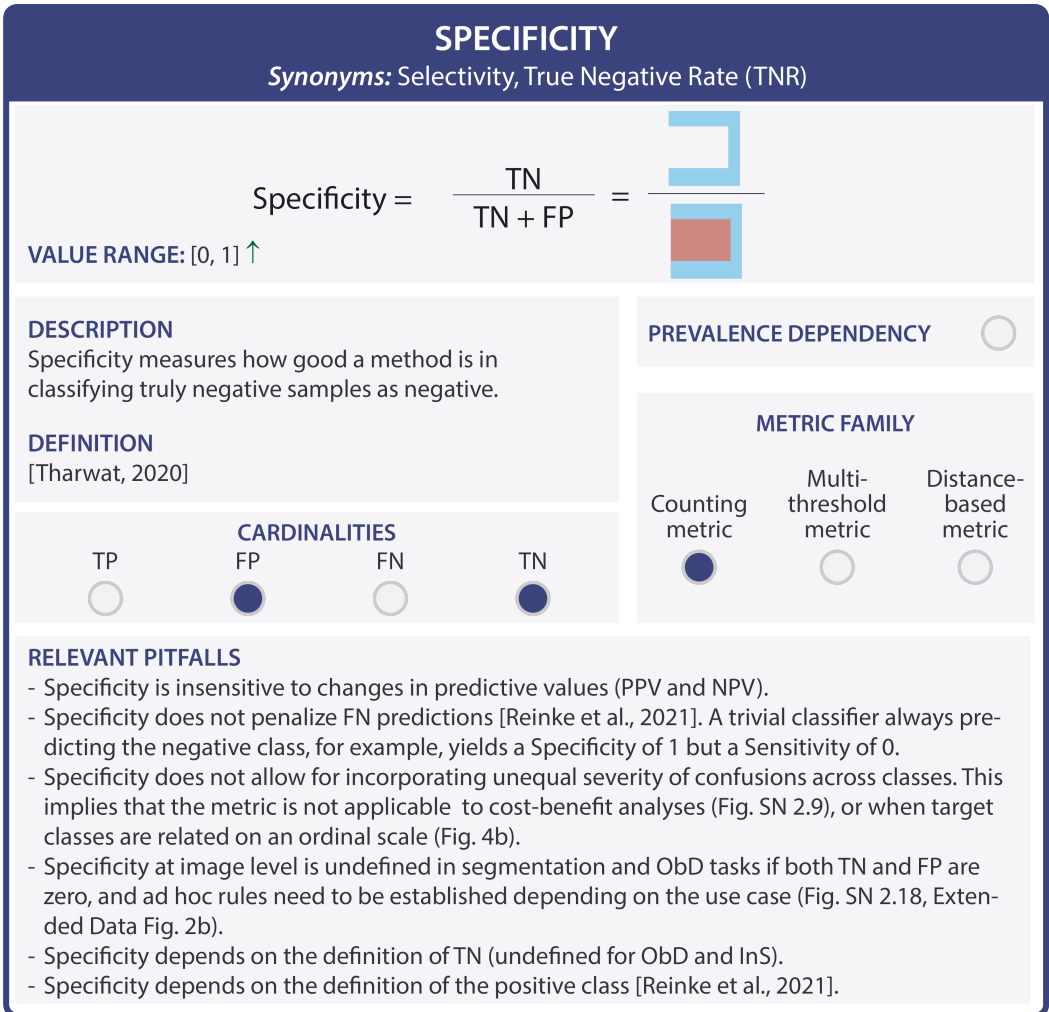
Extended Data Fig. SN 3.52. Metric profile of Positive Likelihood Ratio (LR+). The upward arrow in the value range indicates that higher values are better than lower values. Abbreviations: False Negative (FN), False Positive (FP), Instance Segmentation (InS), Object Detection (ObD), Positive Predictive Value (PPV), True Negative (TN), True Positive (TP). References: Attia, 2003: [2], Reinke et al., 2021: [71]. Mentioned figures: Figs. 4b, 5a, SN 2.11, SN 2.20, Extended Data Fig. 2b.



Extended Data Fig. SN 3.53. Metric profile of the Positive Predictive Value (PPV). The upward arrow in the value range indicates that higher values are better than lower values. Abbreviations used in the figure: False Negative (FN), False Positive (FP), Instance Segmentation (InS), Object Detection (ObD), True Negative (TN), True Positive (TP). References used in the figure: Reinke et al., 2021: [71], Tharwat, 2020: [81]. Mentioned figures: Figs. 4b, SN 2.9, SN 2.11, SN 2.17, SN 2.20, Extended Data Fig. 2b.



Extended Data Fig. SN 3.54. Metric profile of Sensitivity. The upward arrow in the value range indicates that higher values are better than lower values. Abbreviations: False Negative (FN), False Positive (FP), Object Detection (ObD), Positive Predictive Value (PPV), True Negative (TN), True Positive (TP). References: Maier-Hein et al., 2022: [58], Reinke et al., 2021: [71], Tharwat, 2020: [81]. Mentioned figures: Figs. 4b, SN 2.11, SN 2.20, Extended Data Fig. 2b.



Extended Data Fig. SN 3.55. Metric profile of Specificity. The upward arrow in the value range indicates that higher values are better than lower values. Abbreviations: False Negative (FN), False Positive (FP), True Negative (TN), True Positive (TP). References: Reinke et al., 2021: [71], Tharwat, 2020: [81]. Mentioned figures: Figs. 4b SN 2.11, SN 2.20, Extended Data Fig. 2b.

WEIGHTED COHEN'S KAPPA (WCK)

Synonyms: Weighted Cohen's Kappa Coefficient, Weighted Kappa Statistic, Weighted Kappa Score

$$WCK = \frac{p_0^w - p_e^w}{1 - p_e^w}, p_0^w = \frac{w_{TP}TP + w_{TN}TN + w_{FP}FP + w_{FN}FN}{TP + TN + FP + FN} = \frac{w_{TP} \cdot \text{TP} + w_{TN} \cdot \text{TN} + w_{FP} \cdot \text{FP} + w_{FN} \cdot \text{FN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

$$p_e^w = w_{TP} \frac{(TP + FP)(TP + FN)}{TP + TN + FP + FN} + w_{TN} \frac{(TN + FP)(TN + FN)}{TP + TN + FP + FN} + w_{FN} \frac{(FN + FP)(FN + TN)}{TP + TN + FP + FN} + w_{FP} \frac{(FP + TP)(FP + TN)}{TP + TN + FP + FN}$$

$$= w_{TP} \frac{\text{TP} \cdot \text{TP} + \text{FP} \cdot \text{TP} + \text{FN} \cdot \text{TP}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} + w_{TN} \frac{\text{TN} \cdot \text{TN} + \text{FP} \cdot \text{TN} + \text{FN} \cdot \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} + w_{FN} \frac{\text{FN} \cdot \text{FN} + \text{FP} \cdot \text{FN} + \text{TP} \cdot \text{FN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} + w_{FP} \frac{\text{FP} \cdot \text{FP} + \text{FP} \cdot \text{TP} + \text{FP} \cdot \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

VALUE RANGE: [-1, 1] ↑
 A value of 0 refers to a prediction which is not better than random guessing.
 $w_{TP}/w_{TN}/w_{FP}/w_{FN}$: (estimation of) costs of the respective cardinalities; can be adjusted as a weighting of them.

DESCRIPTION WCK calculates the degree of agreement between the reference and prediction while incorporating the agreement resulting from chance. WCK is a generalization of CK with 0-1 weights.	CARDINALITIES			
	TP	FP	FN	TN
	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>
DEFINITION [Cohen, 1960]	PREVALENCE DEPENDENCY			
	<input checked="" type="radio"/>			
MULTI-CLASS DEFINITION For C classes, WCK can be defined as: $WCK = 1 - \left(\frac{\sum_{i=1}^C \sum_{j=1}^C w_{ij} \cdot n_{ij}}{\sum_{i=1}^C \sum_{j=1}^C w_{ij} \cdot \frac{n_{i \cdot} \cdot n_{\cdot j}}{N^2}} \right)$	METRIC FAMILY			
	Counting metric	Multi-threshold metric	Distance-based metric	
	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	

MULTI-CLASS DEFINITION
 For C classes, WCK can be defined as: $WCK = 1 - \left(\frac{\sum_{i=1}^C \sum_{j=1}^C w_{ij} \cdot n_{ij}}{\sum_{i=1}^C \sum_{j=1}^C w_{ij} \cdot \frac{n_{i \cdot} \cdot n_{\cdot j}}{N^2}} \right)$

n_{ij} : entry of the confusion matrix for row i and column j, i.e. samples of actual class i that were predicted as class j
 N : total number of samples

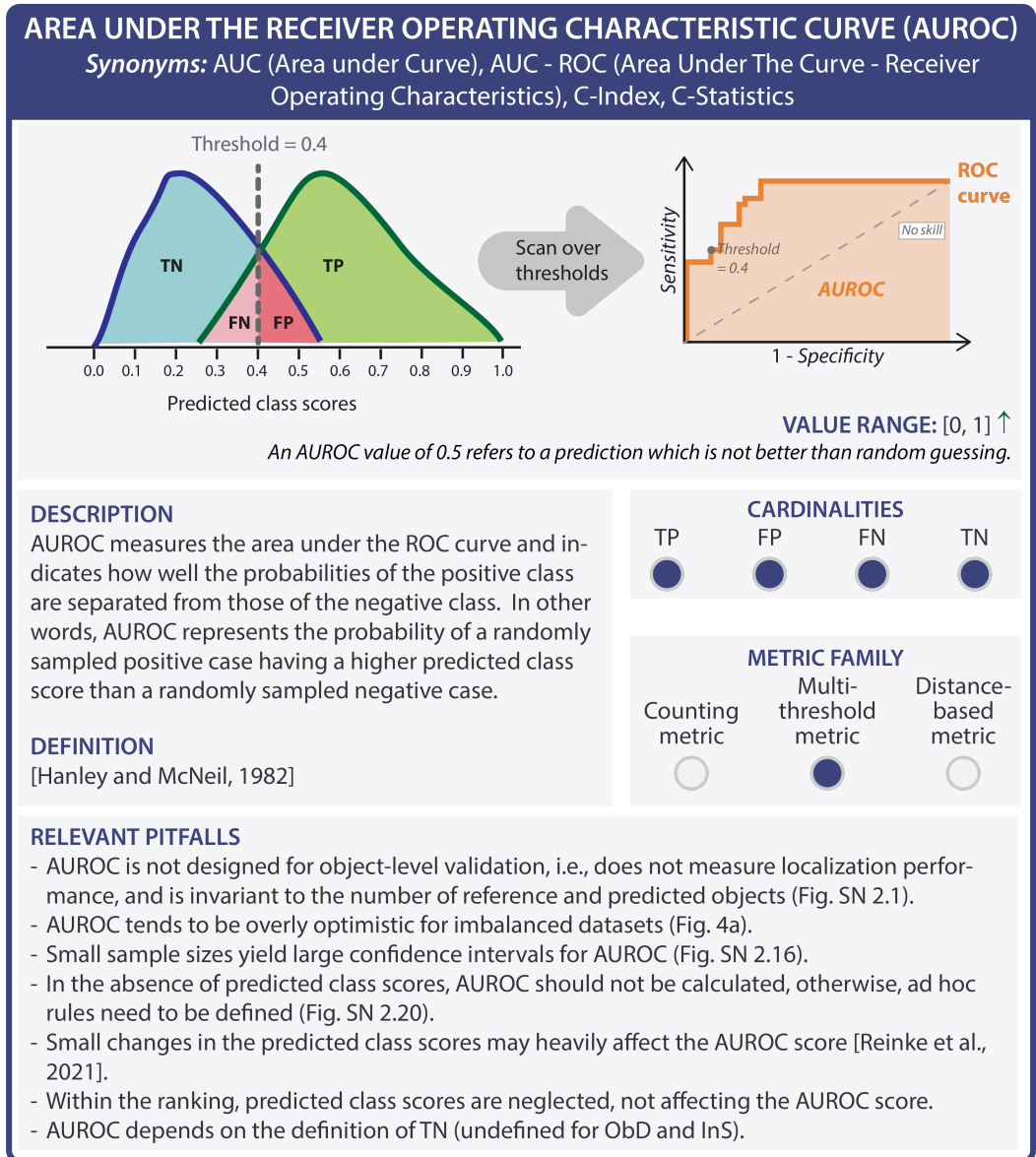
$n_{i \cdot}$: sum of entries of row i of the confusion matrix
 $n_{\cdot j}$: sum of entries of column j of the confusion matrix
 w_{ij} : costs for the entry of the confusion matrix for row i and column j, i.e., the cost for samples of actual class i that were predicted as class j

RELEVANT PITFALLS

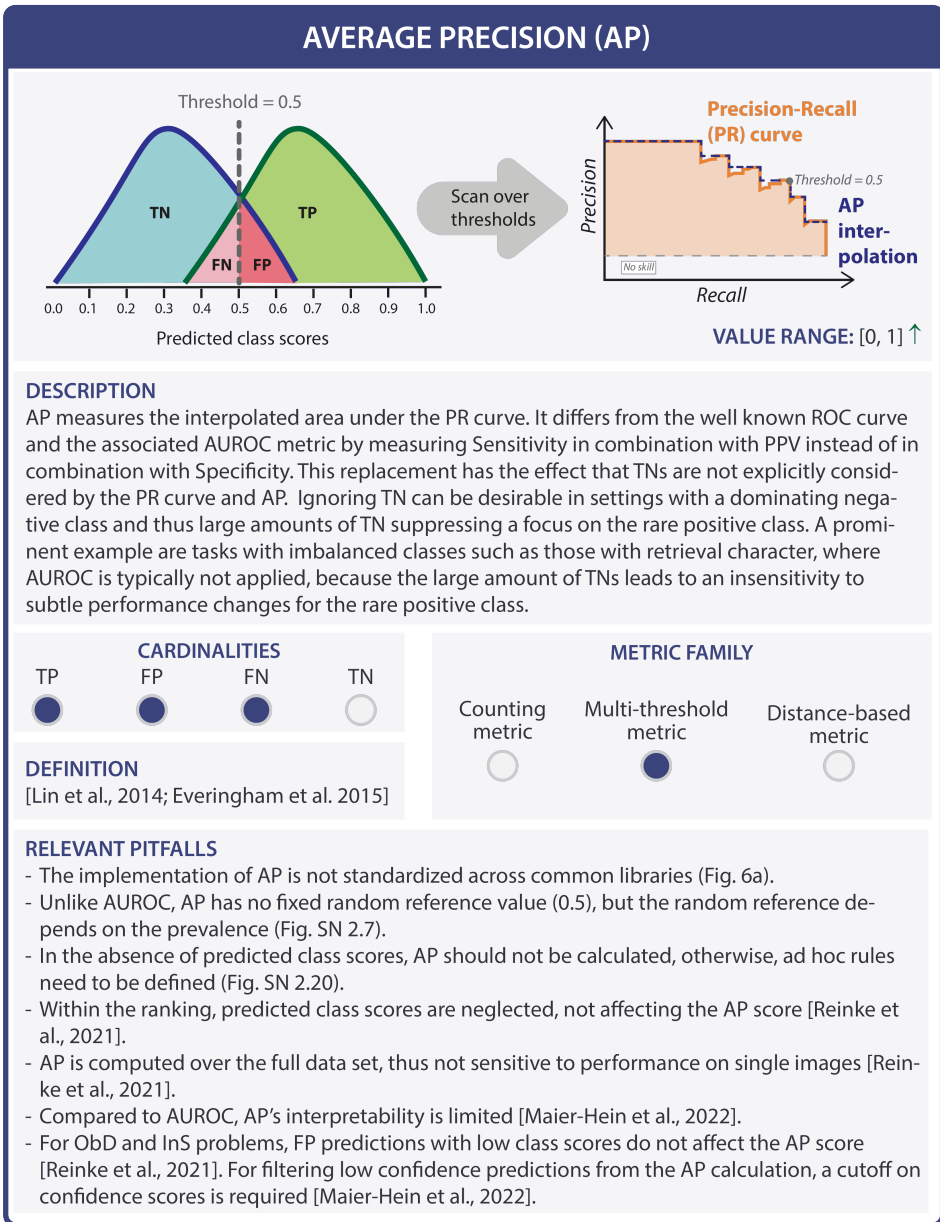
- WCK is prevalence-dependent, thus not comparable across data sets and may yield different rankings than the BA (Figs. SN 2.7, SN 2.15, [Reinke et al., 2021]).
- The theoretical lower bound of WCK (-1) may not always be achievable (Fig. SN 2.35).
- In settings where target classes lie on an ordinal scale, WCK may harshly penalize label shifts [Reinke et al., 2021].
- WCK is hard to interpret [Delgado and Tibau, 2019].
- WCK was designed for symmetric situations (guesses of two raters) [Powers, 2012].
- Compared to other multi-class metrics like EC, WCK lacks a framework to identify and validate the decision rule applied to predicted class scores [Ferrer, 2022].
- WCK with quadratic weights may yield "paradoxical results" [Warrens, 2012].
- The selection of weights to different types of mistakes is problem-dependent and requires domain knowledge.
- WCK depends on the definition of TN (undefined for ObD and InS).

Extended Data Fig. SN 3.56. Metric profile of Weighted Cohen's Kappa (WCK). The upward arrow in the value range indicates that higher values are better than lower values. Abbreviations: Balanced Accuracy (BA), Cohen's Kappa (CK), Expected Cost (EC), False Negative (FN), False Positive (FP), Instance Segmentation (InS), Object Detection (ObD), True Negative (TN), True Positive (TP). References: Cohen, 1960: [17], Delgado and Tibau, 2019: [24], Ferrer, 2022: [31], Powers, 2012: [70], Reinke et al., 2021: [71], Warrens, 2012: [91]. Mentioned figures: Figs. SN 2.9, SN 2.17, SN 2.37.

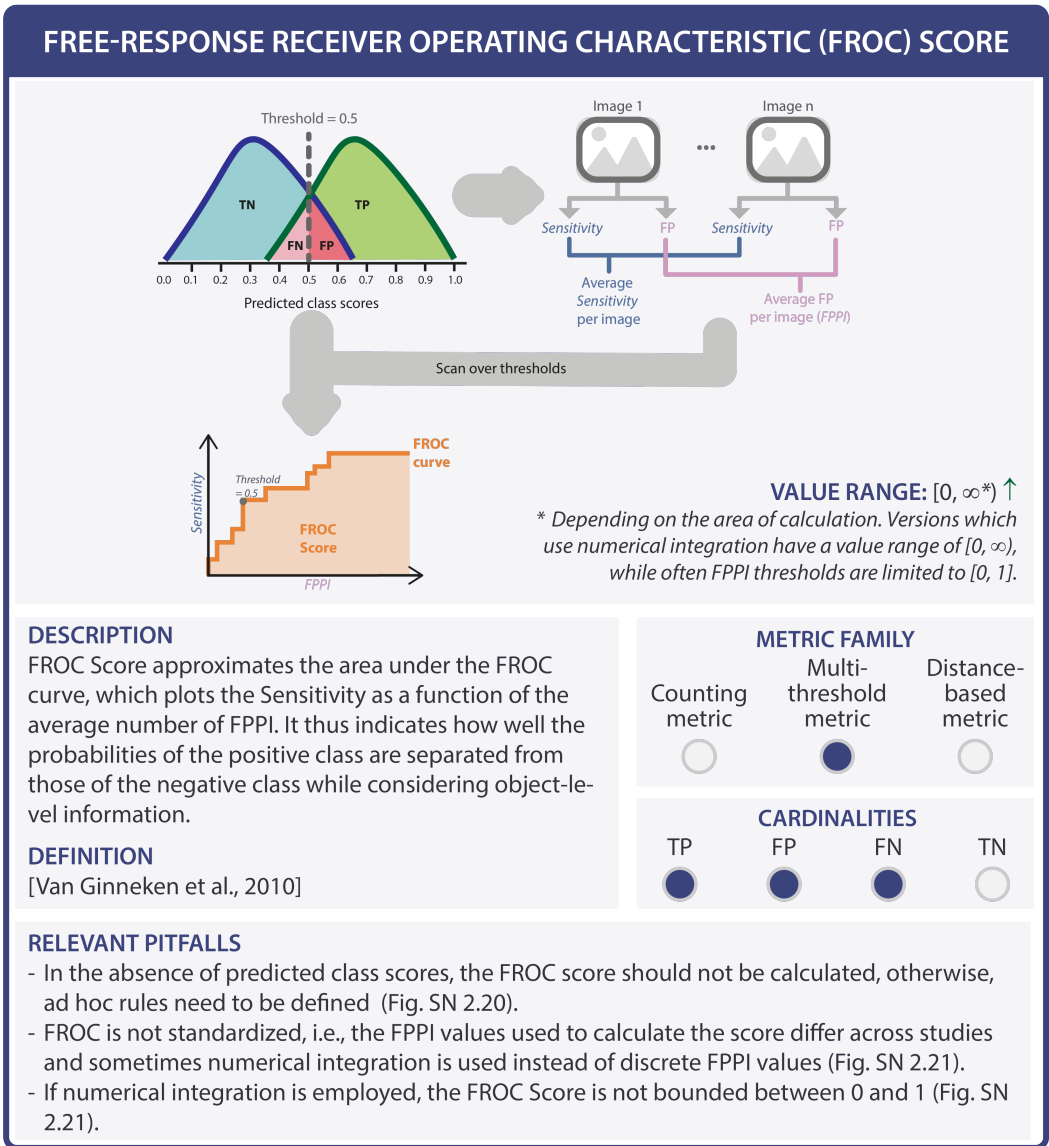
3.1.2 Multi-threshold metrics.



Extended Data Fig. SN 3.57. Metric profile of Area under the Receiver Operating Characteristic Curve (AUROC). The upward arrow in the value range indicates that higher values are better than lower values. Abbreviations: False Negative (FN), False Positive (FP), Instance Segmentation (InS), Object Detection (ObD), Receiver Operating Characteristic (ROC), True Negative (TN), True Positive (TP). References: Hanley and McNeil, 1982: [40], Reinke et al., 2021: [71]. Mentioned figures: Figs. 5a, SN 2.3, SN 2.18, SN 2.22.

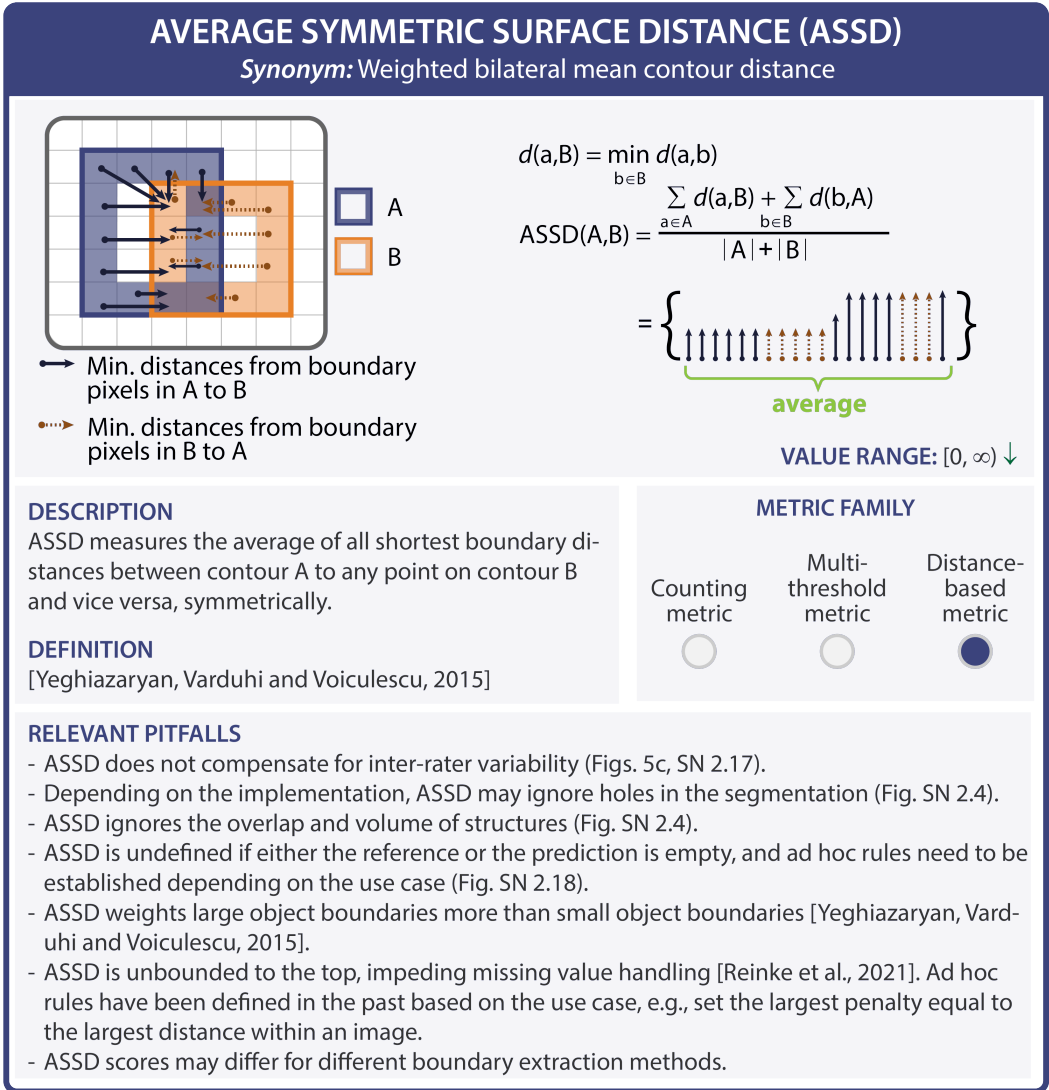


Extended Data Fig. SN 3.58. Metric profile of Average Precision (AP). The upward arrow in the value range indicates that higher values are better than lower values. Abbreviations: Area under the Receiver Operating Characteristic Curve (AUROC), False Negative (FN), False Positive (FP), Instance Segmentation (InS), Object Detection (ObD), Precision-Recall (PR), True Negative (TN), True Positive (TP). References: Everingham et al., 2015: [29], Lin et al., 2014: [56], Maier-Hein et al., 2022: [58], Reinke et al., 2021: [71]. Mentioned figures: Figs. 6a, SN 2.9, SN 2.22.

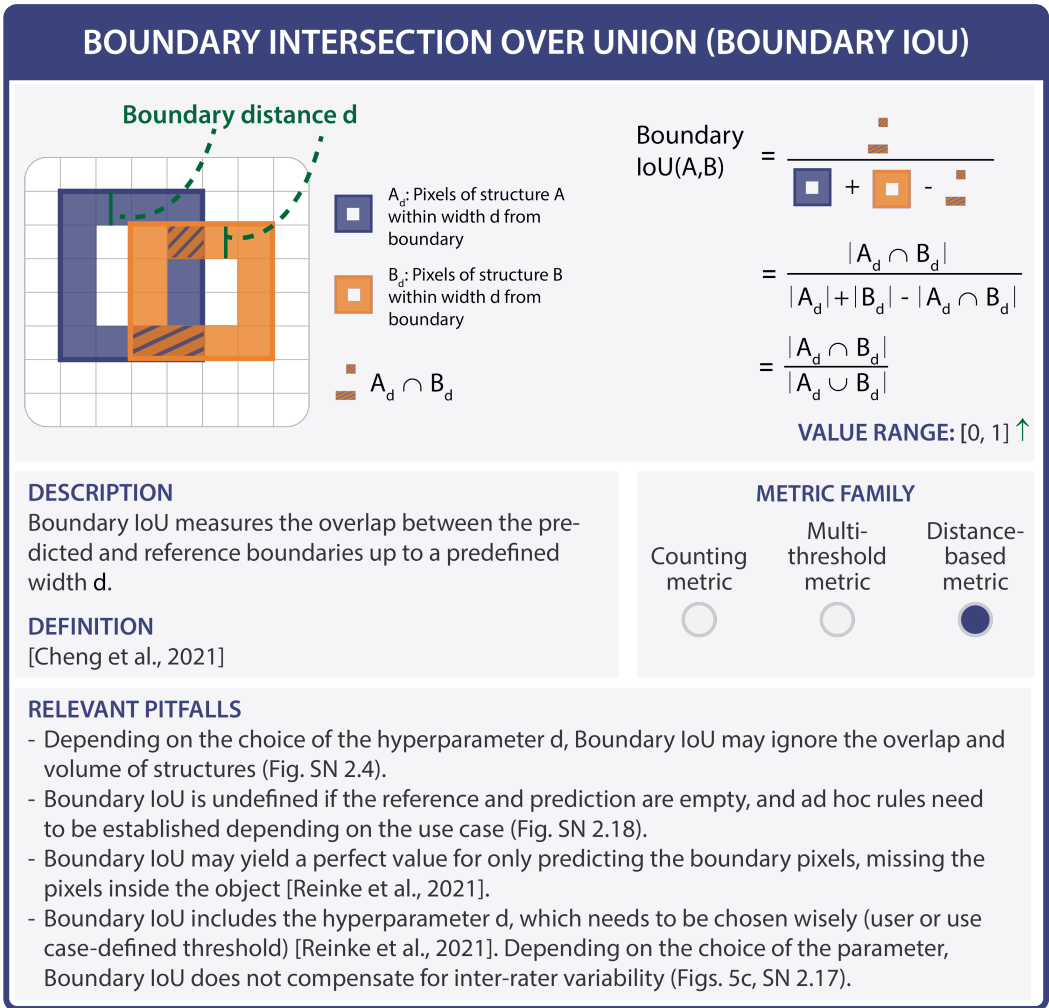


Extended Data Fig. SN 3.59. Metric profile of Free-Response Receiver Operating Characteristic (FROC). The upward arrow in the value range indicates that higher values are better than lower values. Abbreviations: False Negative (FN), False Positive (FP), False Positives per Image (FPPI), True Negative (TN), True Positive (TP). References: Van Ginneken et al., 2010: [85]. Mentioned figures: Figs. SN 2.22, SN 2.23.

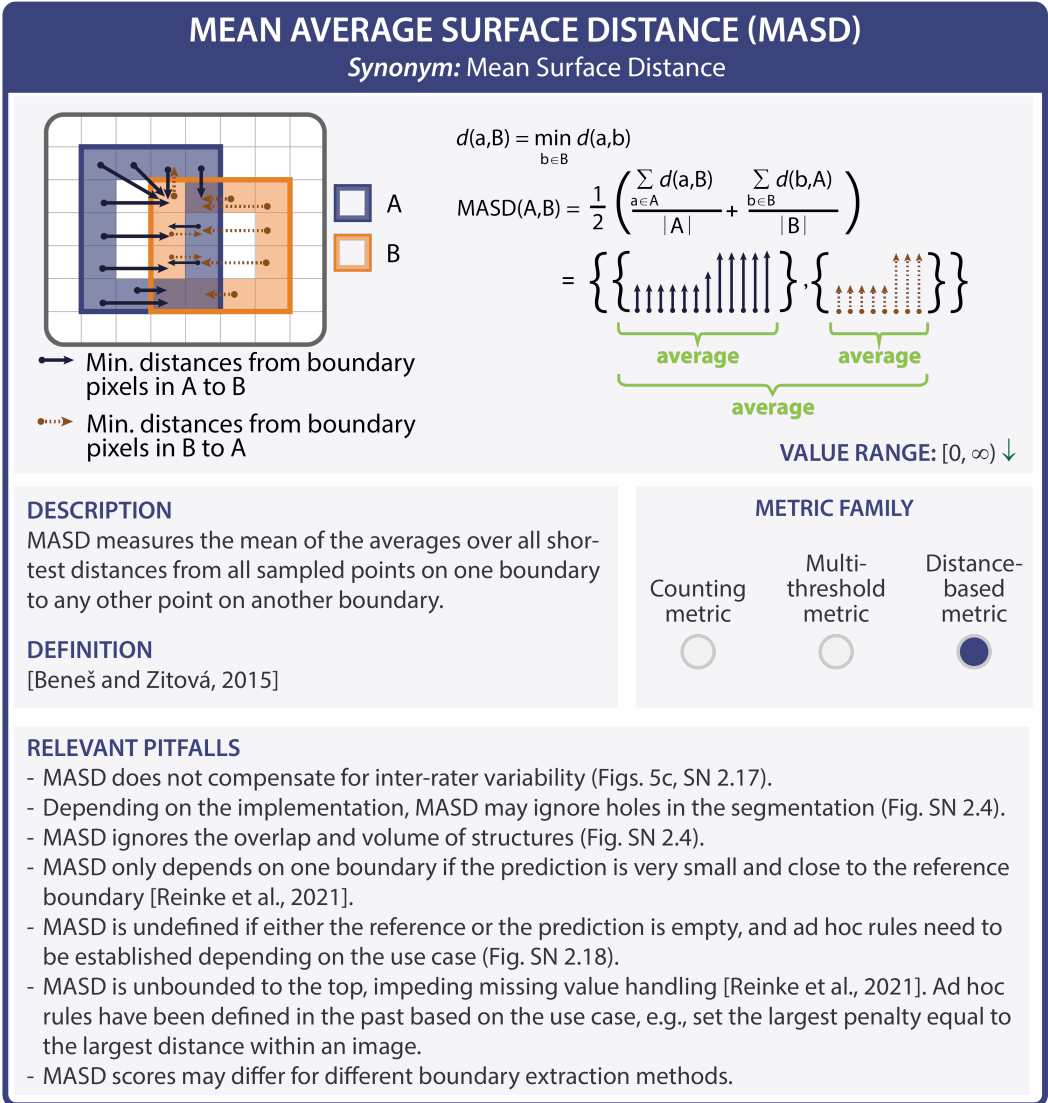
3.1.3 Distance-based metrics.



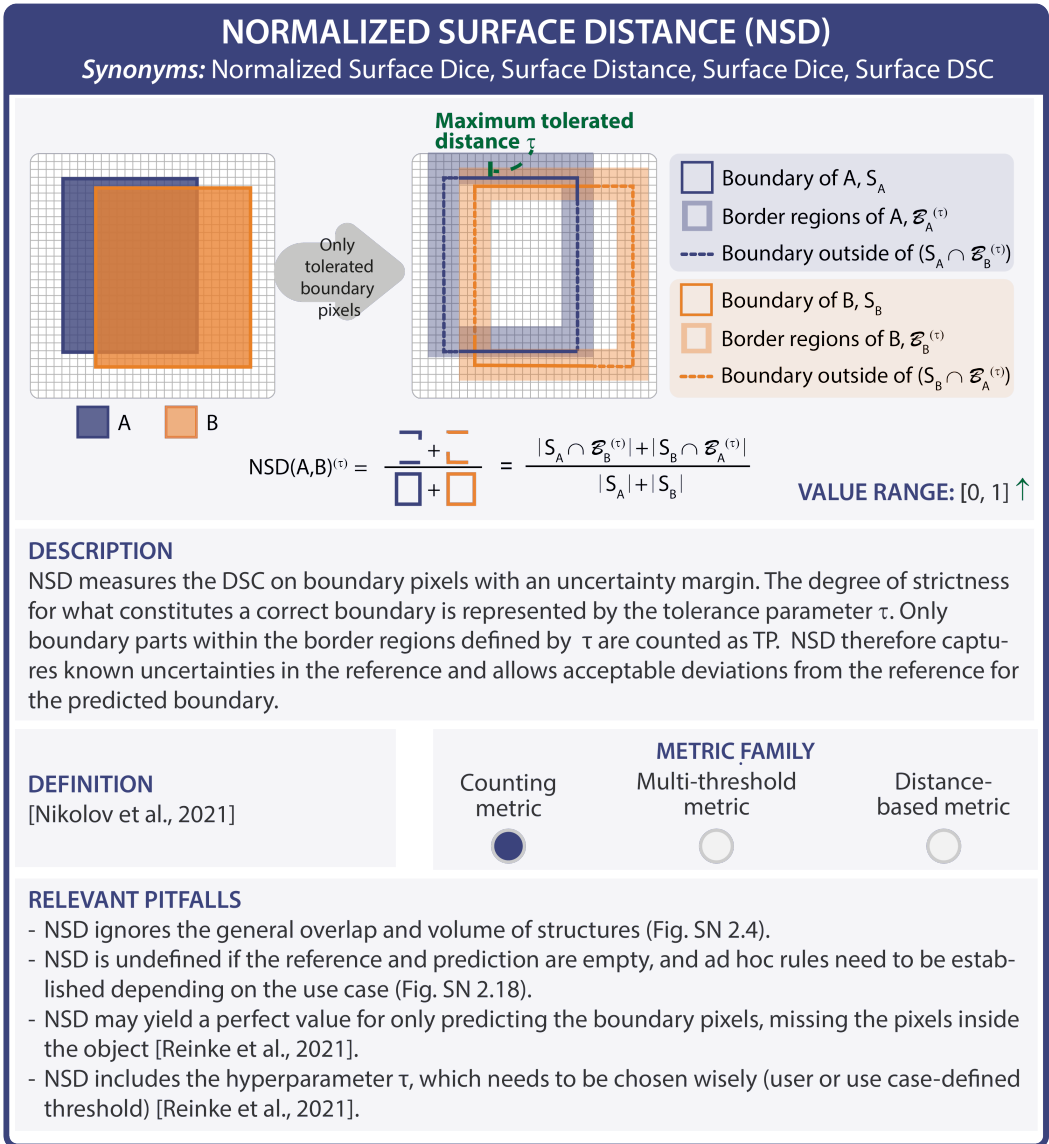
Extended Data Fig. SN 3.60. Metric profile of Average Symmetric Surface Distance (ASSD). The downward arrow in the value range indicates that lower values are better than higher values. Abbreviation: Semantic Segmentation (SemS). References: Reinke et al., 2021: [71], Yeghiazaryan, Varduhi and Voiculescu, 2015: [94]. Mentioned figures: Figs. 5c, SN 2.6, SN 2.19, SN 2.20.



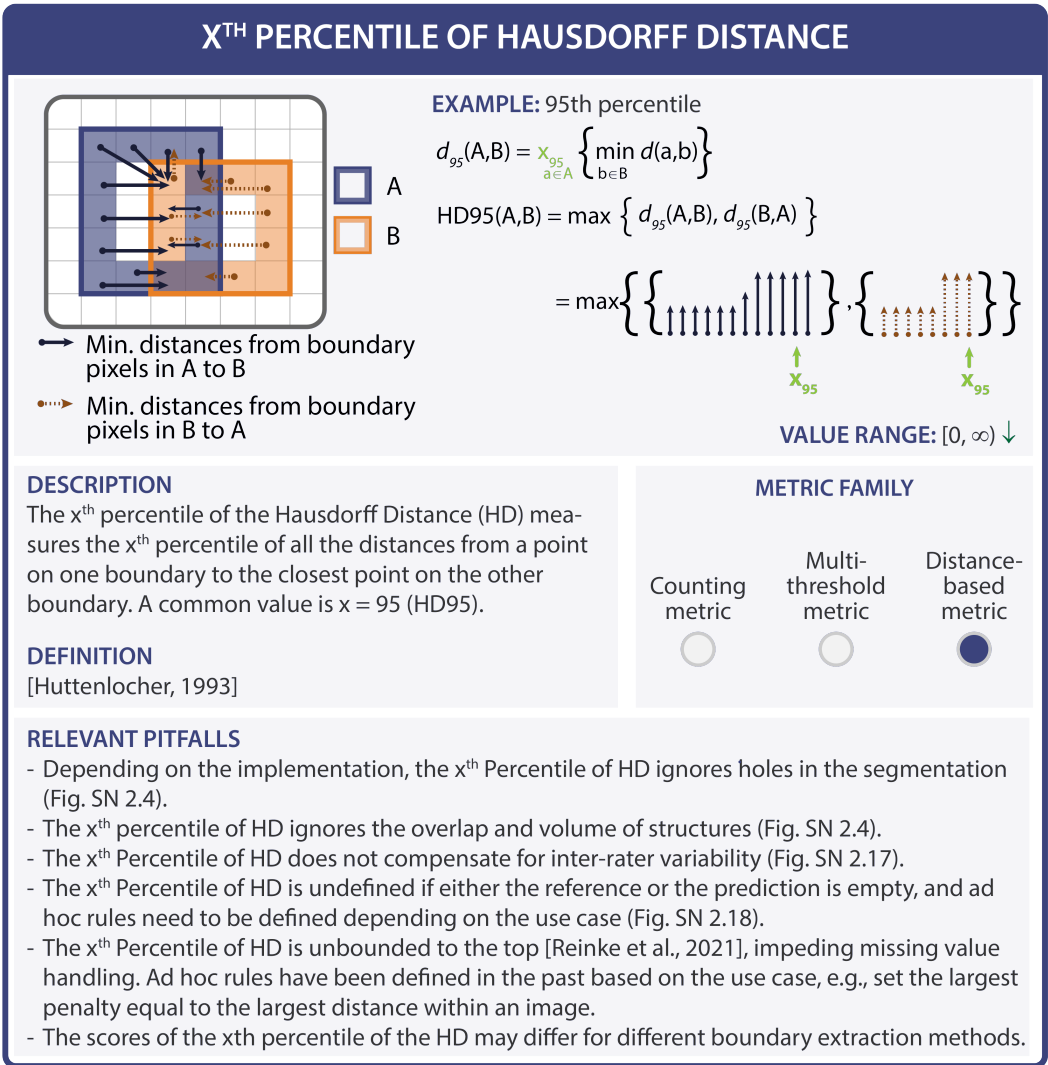
Extended Data Fig. SN 3.61. Metric profile of the Boundary Intersection over Union (IoU). The upward arrow in the value range indicates that higher values are better than lower values. References: Cheng et al., 2021: [13], Reinke et al., 2021: [71]. Mentioned figures: Figs. 5c, SN 2.6, SN 2.19, SN 2.20.



Extended Data Fig. SN 3.63. Metric profile of Mean Average Surface Distance (MASD). The downward arrow in the value range indicates that lower values are better than higher values. Abbreviation: Semantic Segmentation (SemS). References: Beneš and Zitová, 2015: [6], Reinke et al., 2021: [71]. Mentioned figures: Figs. 5c, SN 2.6, SN 2.19, SN 2.20.



Extended Data Fig. SN 3.64. Metric profile of Normalized Surface Distance (NSD). The upward arrow in the value range indicates that higher values are better than lower values. Abbreviation: Dice Similarity Coefficient (DSC). References: Nikolov et al., 2021: [66], Reinke et al., 2021: [71]. Mentioned figures: Figs. SN 2.6, SN 2.20.



Extended Data Fig. SN 3.65. Metric profile of X^{th} Percentile of Hausdorff Distance (HD). The downward arrow in the value range indicates that lower values are better than higher values. Abbreviations: Hausdorff Distance (HD), Semantic Segmentation (SemS). References: Huttenlocher, 1993: [44], Reinke et al., 2021: [71]. Mentioned figures: Figs. SN 2.6, SN 2.19, SN 2.20.

3.2 Calibration metrics

BRIER SCORE (BS)/BRIER SKILL SCORE (BSS)

$$BS = \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^C (p_{ik} - y_{ik})^2$$

VALUE RANGE: [0, 2] ↓

N: number of samples
C: number of classes

p_{ik}: predicted probability for sample *x_i* and class *k*
y_{ik}: outcome; *y_{ik}* = 1 if *y_i* is equal to *k* and 0 otherwise

DESCRIPTION
BS is the mean squared error of a predicted class score and the actual outcome, thus assessing discrimination and calibration in one joint score. It is a proper scoring rule.

VARIANT
Brier Skill Score (BSS): normalizes BS by the BS of a naive system.

DEFINITION
[Gneiting and Raftery, 2007]

METRIC FAMILY

Counting metric	Multi-threshold metric	Distance-based metric	Calibration metric
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>

TYPE OF CALIBRATION

Top-label	Marginal	Canonical
<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>

RELEVANT PITFALLS

- BS/BSS simultaneously assess the discrimination and calibration performance in one score and can thus only be used for relative assessment of calibration.
- BS is highly prevalence-dependent, implying that scores may drastically change when the prevalence changes (Fig. SN 2.7), i.e., predicted class scores linked to sporadic events have little effect on the score, leading to preference of naive systems in imbalanced settings.
- BS/BSS do not allow for incorporating unequal severity of confusions across classes in discrimination. This implies that these metrics are not applicable when target classes are related on an ordinal scale (Fig. 4b, [Reinke et al., 2021]).

Extended Data Fig. SN 3.66. Metric profile of Brier Score (BS). The downward arrow in the value range indicates that lower values are better than higher values. Abbreviation: Brier Skill Score (BSS). References: Gneiting and Raftery, 2007: [33], Reinke et al., 2021: [71]. Mentioned figure: Fig. SN 2.9.

CLASS-WISE CALIBRATION ERROR (CWCE)

$$CWCE = \frac{1}{C} \sum_{c=1}^C \sum_{m=1}^M \frac{|B_{c,m}|}{N} \|\text{Accuracy}_c(B_{c,m}) - \text{Confidence}_c(B_{c,m})\|_p^p$$

N: number of samples; *C*: number of classes
B_{c,m}: bin *m* for class *c*
p: determines which *L_p* calibration error is desired; typically *p* = 1

VALUE RANGE: [0, 1] ↓

DESCRIPTION

CWCE is an estimator of the marginal calibration error applying binning to estimate the observed probabilities corresponding to a confidence range. It can be reported per class or in an aggregated fashion with class-specific weights reflecting prevalence or importance of classes, for example.

DEFINITION

[Kull et al., 2019; Kumar et al., 2019]

METRIC FAMILY

Counting metric	Multi-threshold metric	Distance-based metric	Calibration metric
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>

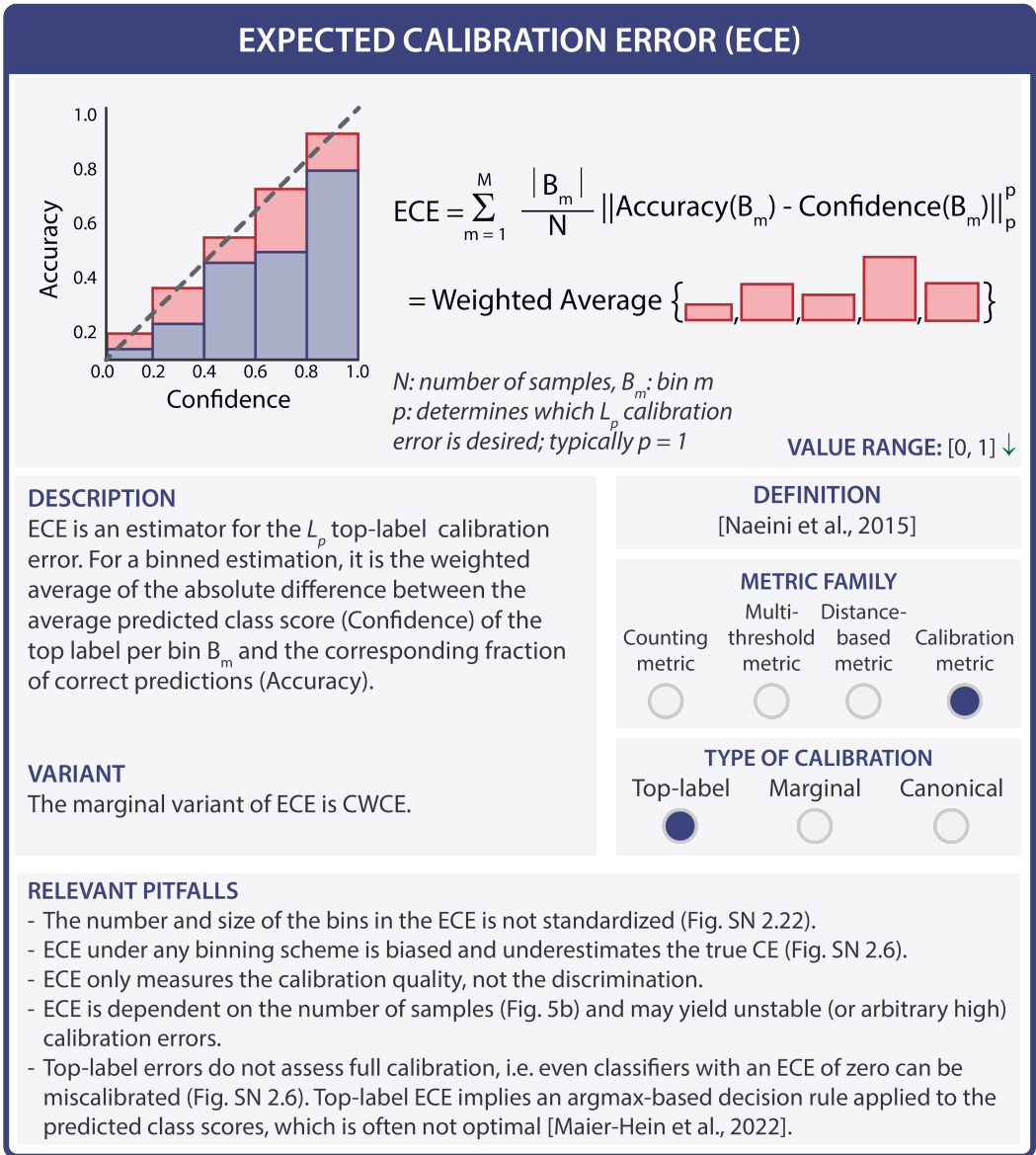
TYPE OF CALIBRATION

Top-label	Marginal	Canonical
<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>

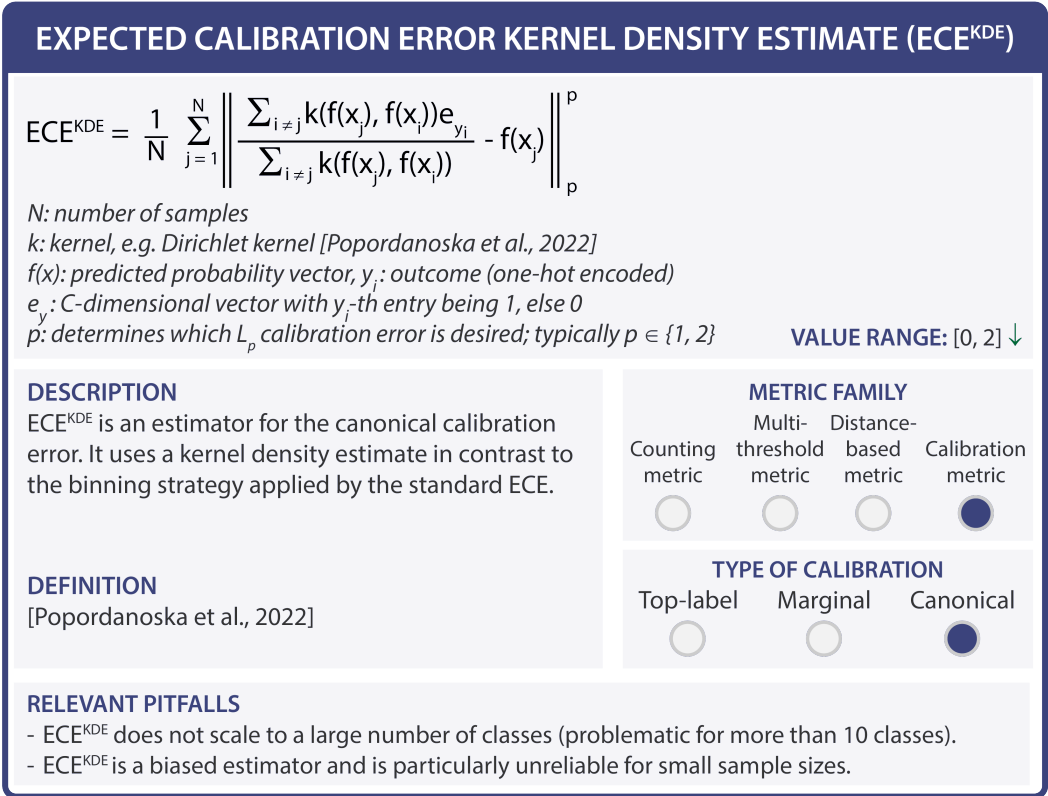
RELEVANT PITFALLS

- The number and size of the bins in the CWCE is not standardized (Fig. SN 2.22).
- CWCE under any binning scheme is biased and underestimates the true CE (Fig. SN 2.6).
- CWCE only measures the calibration quality, not the discrimination.
- Marginal errors do not assess full calibration (Fig. SN 2.6)
- CWCE is dependent on the number of samples (Fig. 5b) and may yield unstable (or arbitrary high) calibration errors. Even classifiers with an CWCE of zero can be miscalibrated (Fig. SN 2.6).
- For CWCE, there are no systematic studies on behavior for imbalanced data.

Extended Data Fig. SN 3.67. Metric profile of Class-Wise Calibration Error (CWCE). The downward arrow in the value range indicates that lower values are better than higher values. References: Kumar et al., 2019: [54], Kull et al., 2019: [53]. Mentioned figures: Figs. 5b, SN 2.8, SN 2.24.



Extended Data Fig. SN 3.68. Metric profile of Expected Calibration Error (ECE). The downward arrow in the value range indicates that lower values are better than higher values. References: Maier-Hein et al., 2022: [58], Naeini et al., 2015: [63], Reinke et al., 2021: [71]. Mentioned figures: Figs. 5b, SN 2.8, SN 2.24.



Extended Data Fig. SN 3.69. Metric profile of Expected Calibration Error Kernel Density Estimate (ECE^{KDE}). The downward arrow in the value range indicates that lower values are better than higher values. Abbreviation: Expected Calibration Error (ECE). Reference used in the figure: Popordanoska et al., 2022: [69].

KERNEL CALIBRATION ERROR (KCE)

$$KCE = \left(\mathbb{E} \left((e_y - f(x))^T k(f(x), f(x')) (e_{y'} - f(x')) \right) \right)^{1/2}$$

Example estimator: $\widehat{KCE} = \left(\binom{N}{2}^{-1} \sum_{i=1}^N \sum_{j=i+1}^N (e_{y_i} - f(x_i))^T k(f(x_i), f(x_j)) (e_{y_j} - f(x_j)) \right)^{1/2}$

N: number of samples; *k*: matrix-valued kernel; *f(x)*: predicted probability vector;
y_i: outcome; *e_{y_i}*: C-dimensional vector with *y_i*-th entry being 1, else 0

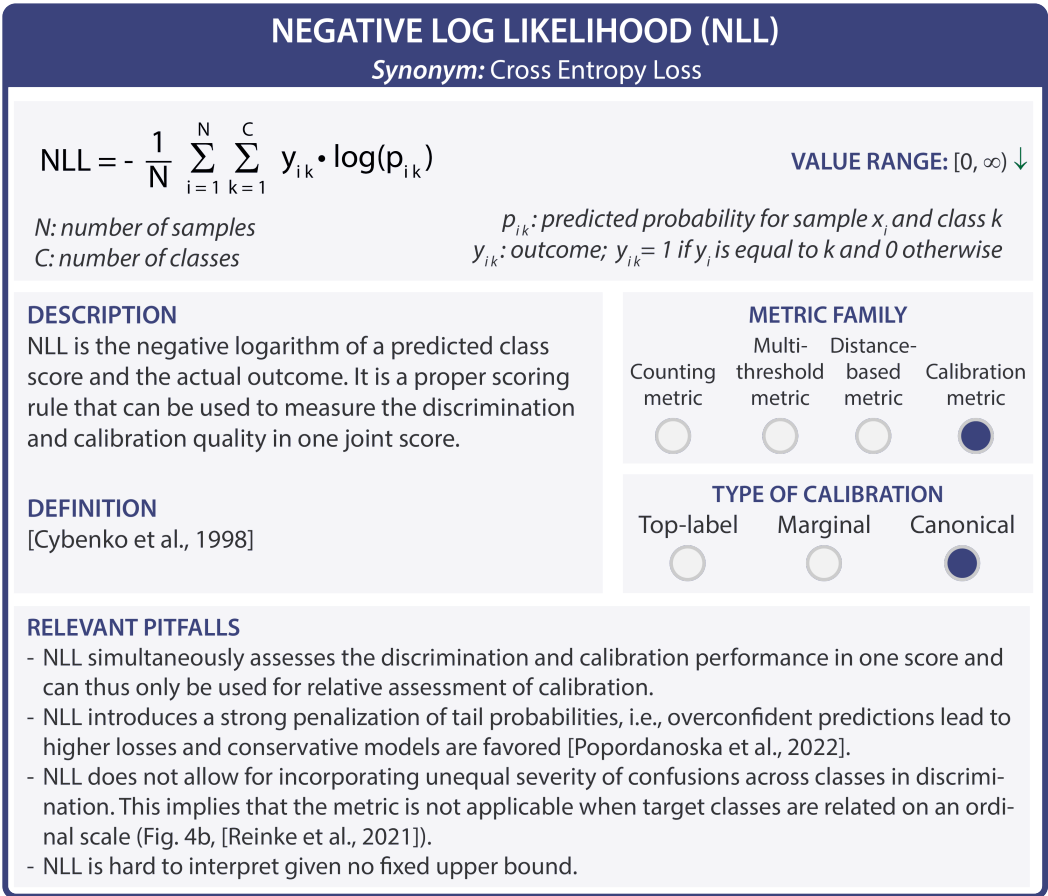
VALUE RANGE: Kernel dependent; in expectation > 0 but estimator can be arbitrarily negative

<p>DESCRIPTION</p> <p>KCE measures a canonical calibration error based on an alternative distance function, the “maximum mean discrepancy” (MMD). It is based on a matrix-valued kernel <i>k</i>.</p> <p>KCE is an unbiased estimator of the calibration error measured by MMD.</p> <p>DEFINITION</p> <p>[Widmann et al., 2019; Gruber and Buettner, 2022]</p>	<p style="text-align: center;">METRIC FAMILY</p> <table border="0" style="width: 100%; text-align: center;"> <tr> <td>Counting metric</td> <td>Multi- threshold metric</td> <td>Distance- based metric</td> <td>Calibration metric</td> </tr> <tr> <td><input type="radio"/></td> <td><input type="radio"/></td> <td><input type="radio"/></td> <td><input checked="" type="radio"/></td> </tr> </table> <p style="text-align: center;">TYPE OF CALIBRATION</p> <table border="0" style="width: 100%; text-align: center;"> <tr> <td>Top-label</td> <td>Marginal</td> <td>Canonical</td> </tr> <tr> <td><input type="radio"/></td> <td><input type="radio"/></td> <td><input checked="" type="radio"/></td> </tr> </table>	Counting metric	Multi- threshold metric	Distance- based metric	Calibration metric	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	Top-label	Marginal	Canonical	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>
Counting metric	Multi- threshold metric	Distance- based metric	Calibration metric												
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>												
Top-label	Marginal	Canonical													
<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>													

RELEVANT PITFALLS

- KCE may be hard to interpret, also due to negative output values.
- KCE cannot be used as an interpretable estimate of the calibration error and should only be used for comparative calibration assessment.
- KCE depends on nontrivial configuration choices of kernels and associated hyperparameters.
- KCE is computationally expensive.

Extended Data Fig. SN 3.70. Metric profile of Kernel Calibration Error (KCE). References: Gruber and Buettner, 2022: [37], Widmann et al., 2019: [92].



Extended Data Fig. SN 3.71. Metric profile of Negative Log Likelihood (NLL). The downward arrow in the value range indicates that lower values are better than higher values. References: Cybenko et al., 1998: [20], Popordanoska et al., 2022: [69], Reinke et al., 2021: [71]. Mentioned figure: Fig. 5b.

ROOT BRIER SCORE (RBS)

$$RBS = \sqrt{\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^C (p_{ik} - y_{ik})^2}$$

N: number of samples
C: number of classes

VALUE RANGE: $[0, \sqrt{2}]$ ↓

DESCRIPTION

RBS is the square root of the mean squared error of a predicted class score and the actual outcome.

It represents a robust upper bound of the canonical calibration error.

METRIC FAMILY

Counting metric	Multi-threshold metric	Distance-based metric	Calibration metric
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>

DEFINITION

[Gruber and Buettner, 2022]

TYPE OF CALIBRATION

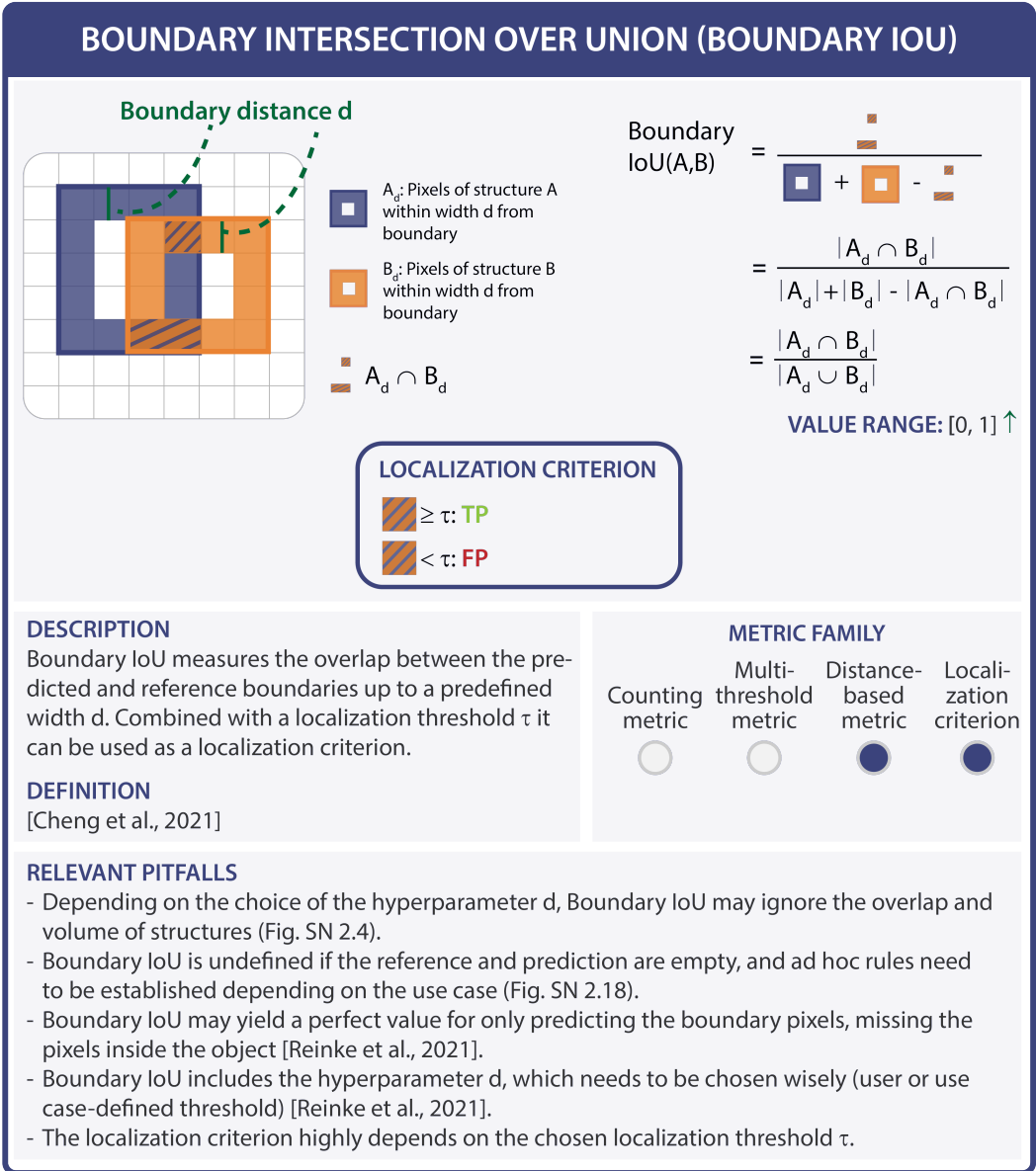
Top-label	Marginal	Canonical
<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>

RELEVANT PITFALLS

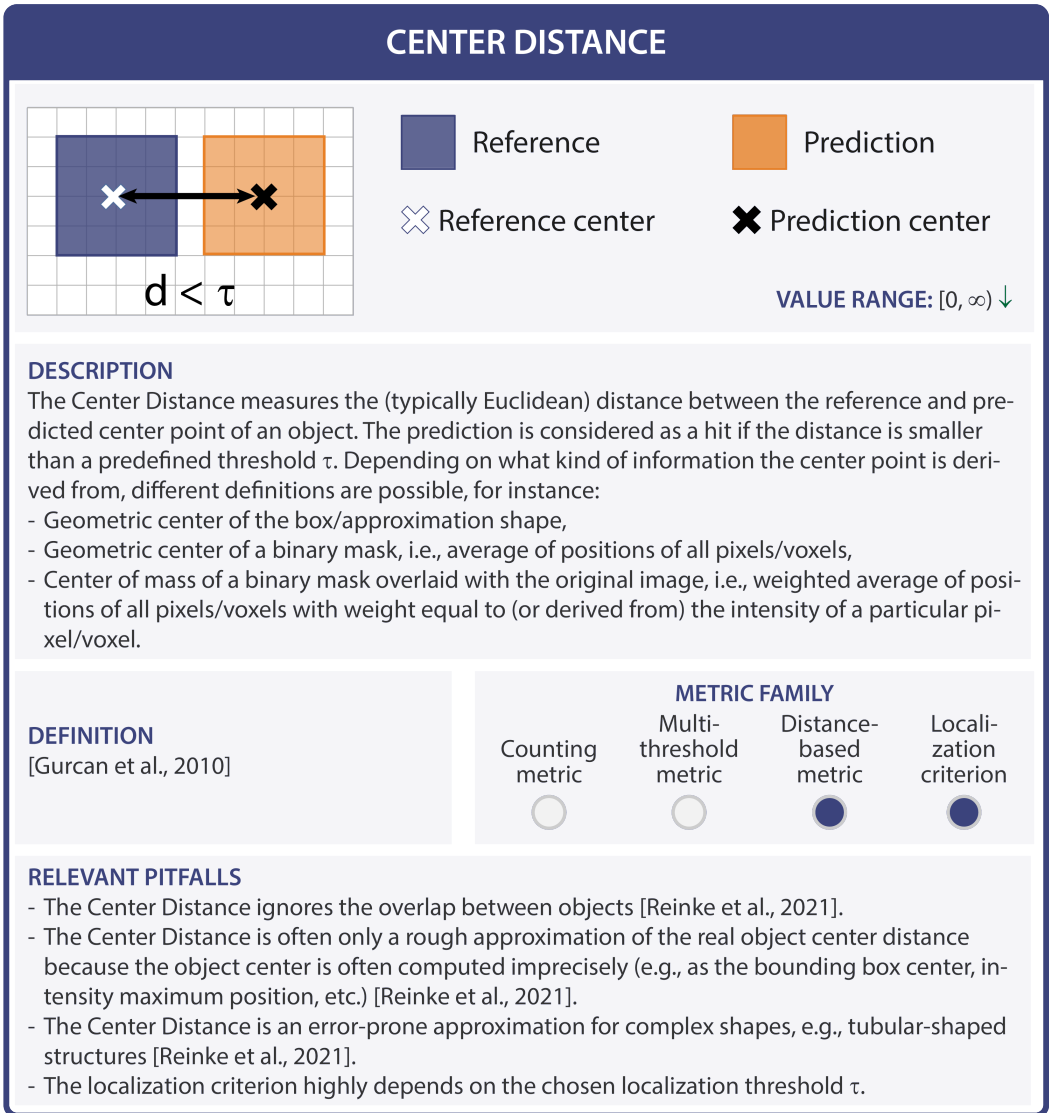
It is not clear how tight the upper bound is, especially for models with low accuracy, given that it is unclear to what extent RBS overestimates the canonical calibration error.

Extended Data Fig. SN 3.72. Metric profile of Root Brier Score (RBS). The downward arrow in the value range indicates that lower values are better than higher values. Reference: Gruber and Buettner, 2022: [37].

3.3 Localization criteria



Extended Data Fig. SN 3.73. Metric profile of the Boundary Intersection over Union (IoU) localization criterion. The upward arrow in the value range indicates that higher values of Boundary IoU are better than lower values. References: Cheng et al., 2021: [13], Reinke et al., 2021: [71]. Mentioned figures: Figs. SN 2.6, SN 2.20.



DESCRIPTION

The Center Distance measures the (typically Euclidean) distance between the reference and predicted center point of an object. The prediction is considered as a hit if the distance is smaller than a predefined threshold τ . Depending on what kind of information the center point is derived from, different definitions are possible, for instance:

- Geometric center of the box/approximation shape,
- Geometric center of a binary mask, i.e., average of positions of all pixels/voxels,
- Center of mass of a binary mask overlaid with the original image, i.e., weighted average of positions of all pixels/voxels with weight equal to (or derived from) the intensity of a particular pixel/voxel.

DEFINITION

[Gurcan et al., 2010]

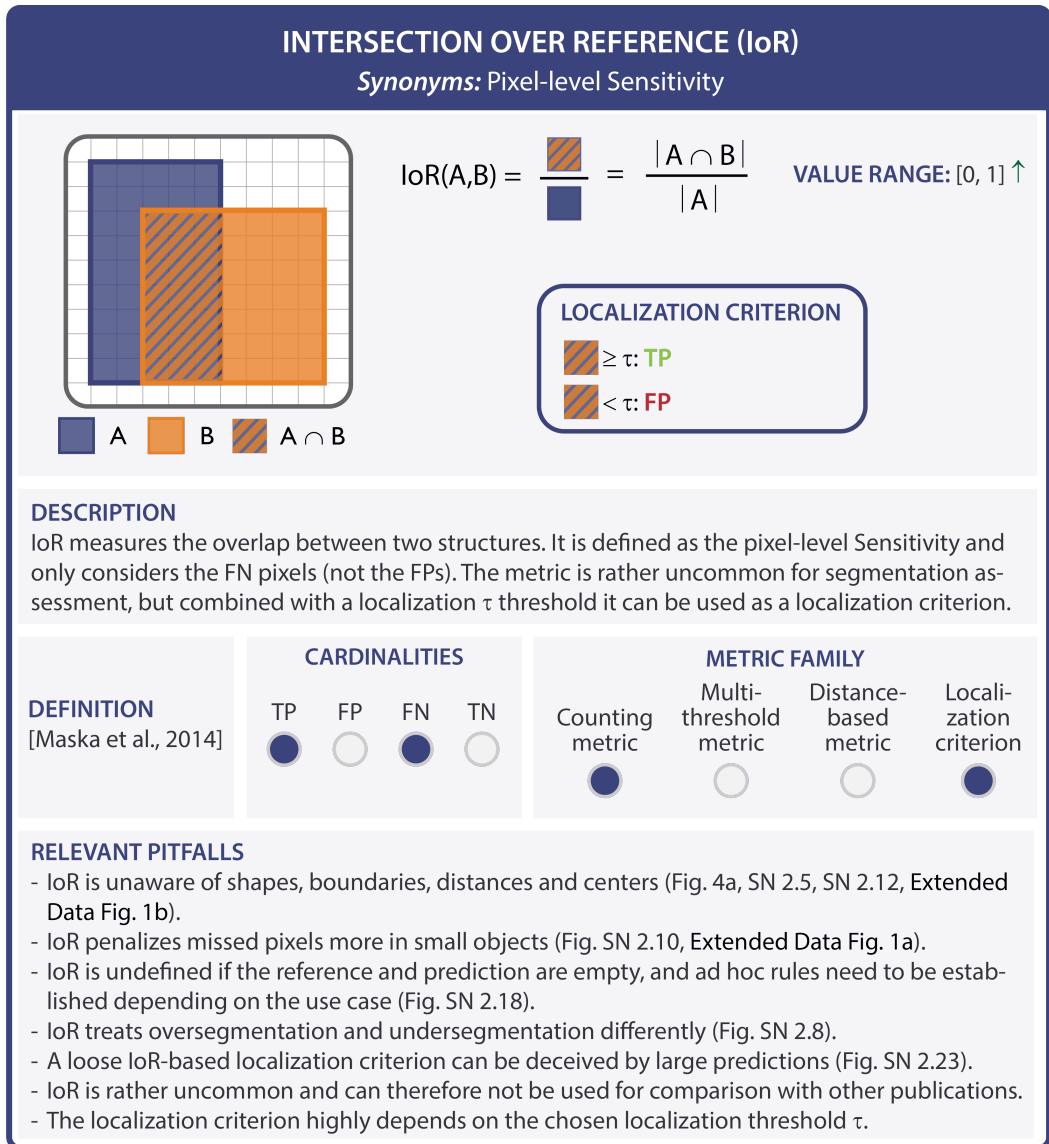
METRIC FAMILY

Counting metric	Multi-threshold metric	Distance-based metric	Localization criterion
<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>

RELEVANT PITFALLS

- The Center Distance ignores the overlap between objects [Reinke et al., 2021].
- The Center Distance is often only a rough approximation of the real object center distance because the object center is often computed imprecisely (e.g., as the bounding box center, intensity maximum position, etc.) [Reinke et al., 2021].
- The Center Distance is an error-prone approximation for complex shapes, e.g., tubular-shaped structures [Reinke et al., 2021].
- The localization criterion highly depends on the chosen localization threshold τ .

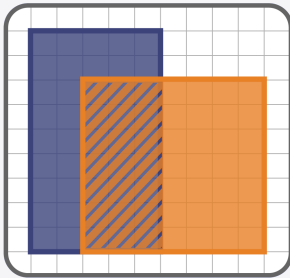
Extended Data Fig. SN 3.74. Metric profile of the Center Distance localization criterion. The downward arrow in the value range indicates that lower values of the Center Distance are better than higher values. References: Gurcan et al., 2010: [39], Reinke et al., 2021: [71].



Extended Data Fig. SN 3.75. Metric profile of the Intersection over Reference (IoR) localization criterion. The upward arrow in the value range indicates that higher values of IoR are better than lower values. Abbreviations: False Negative (FN), False Positive (FP), True Negative (TN), True Positive (TP). References: Maška et al., 2014: [60], Reinke et al., 2021: [71]. Mentioned figures: Figs. 4a, SN 2.7, SN 2.10, SN 2.12, SN 2.13, SN 2.14, SN 2.20, SN 2.25, Extended Data Fig. 1b.

MASK/BOX/APPROX INTERSECTION OVER UNION (MASK/BOX/APPROX IoU)

Synonyms: Jaccard Index, Tanimoto Coefficient



■ A ■ B ■ A ∩ B

$$\begin{aligned}
 \text{IoU}(A,B) &= \frac{\text{Area of } A \cap B}{\text{Area of } A + \text{Area of } B - \text{Area of } A \cap B} \\
 &= \frac{|A \cap B|}{|A| + |B| - |A \cap B|} = \frac{|A \cap B|}{|A \cup B|} \\
 &= \frac{\text{PPV} \cdot \text{Sensitivity}}{\text{PPV} + \text{Sensitivity} - \text{PPV} \cdot \text{Sensitivity}}
 \end{aligned}$$

LOCALIZATION CRITERION

- $\geq \tau$: TP
- $< \tau$: FP

DESCRIPTION

IoU measures the overlap between two structures (see above). Combined with a localization threshold, it is a common localization criterion. It is often referred to as **Box IoU** when comparing bounding boxes, **Mask IoU** when comparing segmentation masks, or **Approx IoU** when comparing approximations of objects beyond bounding boxes.

DEFINITION

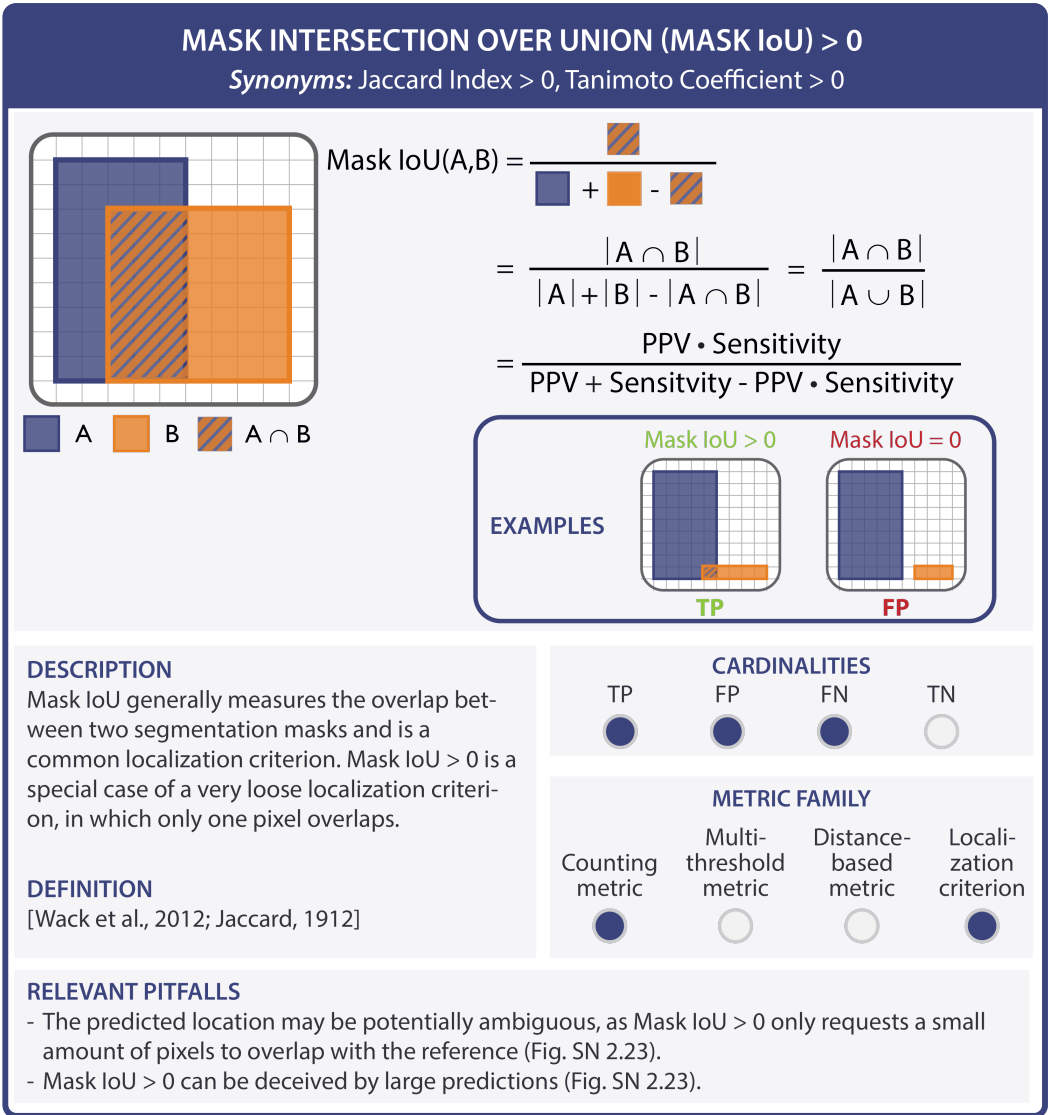
[Jaccard, 1912]

CARDINALITIES			
TP	FP	FN	TN
<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>
METRIC FAMILY			
Counting metric	Multi-threshold metric	Distance-based metric	Localization criterion
<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>

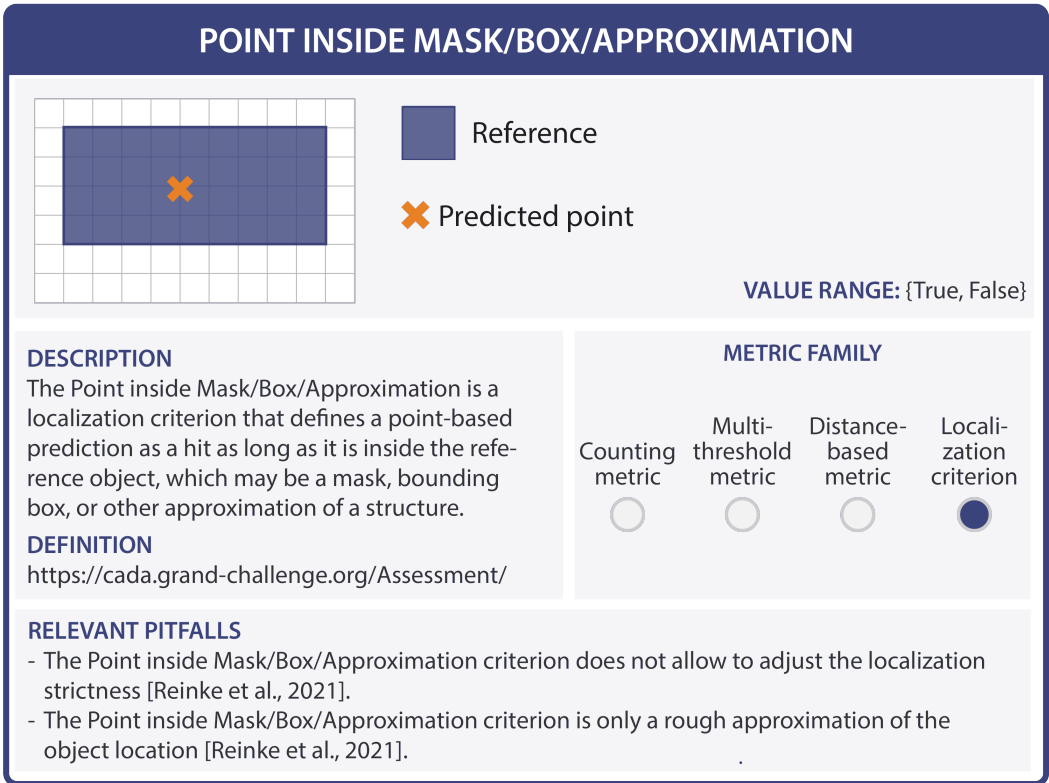
RELEVANT PITFALLS

- IoU is unaware of shapes, boundaries, distances and centers (Figs. 4a, SN 2.5, SN 2.12, Extended Data Fig. 1b).
- IoU penalizes missed pixels more in small objects (Figs. SN 2.10, SN 2.11, Extended Data Fig. 1a).
- Box IoU is not a good representation of complex or disconnected structures (Fig. SN 2.14).
- IoU is undefined if both the reference and prediction are empty, and ad hoc rules need to be established depending on the use case (Fig. SN 2.18).
- IoU treats oversegmentation and undersegmentation differently (Fig. SN 2.8).
- IoU does not compensate for inter-rater variability (Fig. SN 2.17).
- A loose IoU-based localization criterion can be deceived by large predictions (Fig. SN 2.23).
- IoU behaves differently in 2D and 3D settings. In 3D settings, the additional z-dimension results in a cubical increase in erroneous pixels [Reinke et al., 2021].
- An IoU-based localization criterion may highly penalize multiple predictions for the same reference object [Reinke et al., 2021].
- The localization criterion highly depends on the chosen localization threshold τ .

Extended Data Fig. SN 3.76. Metric profile of the Mask/Box/Approx Intersection over Union (IoU) localization criterion. Abbreviations: False Negative (FN), False Positive (FP), True Negative (TN), True Positive (TP). References: Jaccard, 1912: [45], Reinke et al., 2021: [71]. Mentioned figures: Figs. 4a, SN 2.7, SN 2.10, SN 2.12, SN 2.13, SN 2.14, SN 2.16, SN 2.19, SN 2.20, SN 2.25, Extended Data Fig. 1a-b.

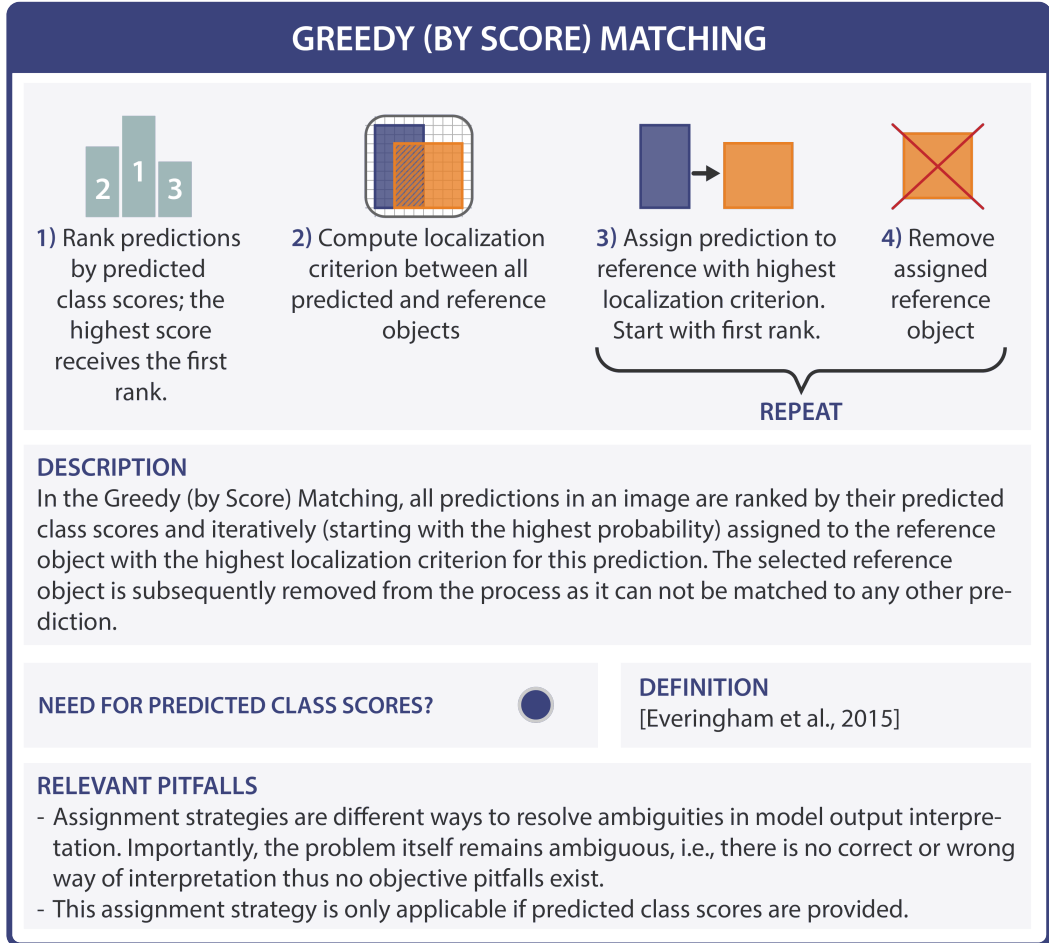


Extended Data Fig. SN 3.77. Metric profile of the Mask Intersection over Union (IoU) > 0 localization criterion. Abbreviations: False Negative (FN), False Positive (FP), True Negative (TN), True Positive (TP). References: Jaccard, 1912: [45], Wack et al., 2012: [90]. Mentioned figure: Fig. SN 2.25.




Extended Data Fig. SN 3.78. Metric profile of Point inside Mask/Box/Approximation. References: <https://cada.grand-challenge.org/Assessment/>, Reinke et al., 2021: [71].

3.4 Assignment strategies

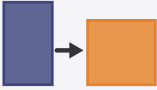


Extended Data Fig. SN 3.79. Cheat Sheet for the Greedy (by Score) Matching. Reference used in the figure: Everingham et al., 2015: [30].


GREEDY (BY LOCALIZATION CRITERION) MATCHING




1) Compute localization criterion between all predicted and reference objects



2) Assign reference to prediction with highest localization criterion.



3) Remove assigned reference object



REPEAT

DESCRIPTION
 If no predicted class scores are available, the Greedy (by Score) Matching can be replaced with the Greedy (by Localization Criterion) Matching. For this strategy, the reference with the highest localization criterion for a predicted object is matched.

NEED FOR PREDICTED CLASS SCORES?


DEFINITION
 [Maier-Hein et al., 2022]

RELEVANT PITFALLS

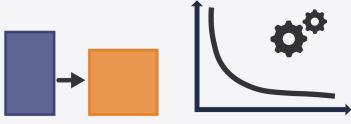
- Assignment strategies are different ways to resolve ambiguities in model output interpretation. Importantly, the problem itself remains ambiguous, i.e., there is no correct or wrong way of interpretation thus no objective pitfalls exist.
- This assignment strategy is not commonly used in the field.

Extended Data Fig. SN 3.80. Cheat Sheet for the Greedy (by Localization Criterion) Matching. Reference used in the figure: Maier-Hein et al., 2022: [58].

OPTIMAL (HUNGARIAN) MATCHING



1) Compute localization criterion between all predicted and reference objects



2) Use cost function to find the optimal assignment of predictions and references based on the localization criterion.

DESCRIPTION
 The Optimal (Hungarian) Matching is associated with a cost function, usually depending on the localization criterion, which is minimized to find the optimal assignment of predictions and reference.

NEED FOR PREDICTED CLASS SCORES?


DEFINITION
 [Kuhn, 1955]

RELEVANT PITFALLS
 The optimization may lead to overoptimistic performance results in case of ambiguous model outputs [Maier-Hein et al., 2022].

Extended Data Fig. SN 3.81. Cheat Sheet for the Optimal (Hungarian) Matching. References used in the figure: Kuhn et al., 1955: [51], Maier-Hein et al., 2022: [58].


MATCHING VIA OVERLAP > 0.5

PREREQUISITE: Overlapping predictions are not possible.




1) Compute overlap-based localization criterion between all predicted and reference objects

> 0.5:



2) If the overlap is greater than 0.5, assign prediction to the reference.



3) Remove assigned reference object

} REPEAT

DESCRIPTION

If there are no overlapping predictions, complex assignment strategies can be avoided by simply setting the localization criterion to $IoU > 0.5$. This strategy inherently avoids matching conflicts, because any secondary prediction would by definition have an overlap < 0.5 of the same reference object.

NEED FOR PREDICTED CLASS SCORES?

DEFINITION
[Everingham et al., 2006]

RELEVANT PITFALLS

- Matching via Overlap > 0.5 is unfeasible if overlapping predictions are possible [Maier-Hein et al., 2022].
- Matching via Overlap > 0.5 cannot be applied if a non-overlap based criterion is employed (e.g., Point inside Mask).

Extended Data Fig. SN 3.82. Cheat Sheet for the Matching via Overlap > 0.5 . References used in the figure: Everingham et al., 2006: [28], Maier-Hein et al., 2022: [58].

ACRONYMS

AI	artificial intelligence
AP	Average Precision
ASSD	Average Symmetric Surface Distance
AUC	Area under the Curve
AUROC	Area under the Receiver Operating Characteristic Curve
BA	Balanced Accuracy
BIAS	Biomedical Image Analysis ChallengeS
Boundary IoU	Boundary Intersection over Union
BS	Brier Score
BSS	Brier Skill Score
CI	Confidence Interval
clDice	centerline Dice Similarity Coefficient
COCO	Common Objects in Context
CK	Cohen's Kappa
CWCE	Class-Wise Calibration Error
DSC	Dice Similarity Coefficient
EC	Expected Cost
ECE	Expected Calibration Error
ECE^{KDE}	Expected Calibration Error Kernel Density Estimate
FN	False Negative
FP	False Positive
FPPI	False Positives per Image
FROC	Free-Response Receiver Operating Characteristic
HD	Hausdorff Distance
HD95	Hausdorff Distance 95th Percentile
InS	Instance Segmentation
IoU	Intersection over Union
IoR	Intersection over Reference
LR+	Positive Likelihood Ratio
KCE	Kernel Calibration Error
mAP	mean Average Precision
MASD	Mean Average Surface Distance
MCC	Matthews Correlation Coefficient
MCE	Maximum Calibration Error
MICCAI	Medical Image Computing and Computer Assisted Interventions
MONAI	Medical Open Network for Artificial Intelligence
NaN	Not a Number
NB	Net Benefit
NPV	Negative Predictive Value
NLL	Negative Log Likelihood
NSD	Normalized Surface Distance
PPV	Positive Predictive Value
ObD	Object Detection
PQ	Panoptic Quality
PR	Precision-Recall
RBS	Root Brier Score

ROC Receiver Operating Characteristic

SemS Semantic Segmentation

TN True Negative

TNR True Negative Rate

TP True Positive

TPR True Positive Rate

WCK Weighted Cohen's Kappa