



Published in final edited form as:

Ophthalmology. 2019 November ; 126(11): 1475–1479. doi:10.1016/j.ophtha.2019.09.014.

An Ophthalmologist's Guide to Deciphering Studies in Artificial Intelligence

Daniel SW TING, MD PhD¹, Aaron Y Lee, MD MSCI², Tien Y WONG, MD PhD¹

¹Singapore Eye Research Institute, Singapore National Eye Center, Duke-NUS Medical School, National University of Singapore

² Department of Ophthalmology, University of Washington, School of Medicine, Seattle WA

Over the past few years, there has been an influx of artificial intelligence (AI) articles in medicine¹ and ophthalmology.^{2–4} Deep learning, a recently described AI machine learning technique, when applied to image analysis allows the algorithm to analyze data using multiple processing layers to extract different image features,¹ with the lower processing layers recognizing basic features (e.g. number and arrangement of edges of an image) and higher layers identifying items more meaningful to human observers (e.g. nose, faces, disease lesions). In ophthalmology, many groups have reported exceptional diagnostic performance using deep learning algorithms to detect various ocular conditions based on anterior segment topography (e.g. keratoconus),⁵ surgical videos (e.g. identification of phases in cataract surgeries),⁶ fundus photographs (e.g. diabetic retinopathy,^{7–11} glaucoma,¹² age-related macular degeneration^{13–16} and retinopathy of prematurity^{17,18}), and anterior and posterior segment optical coherence tomography (OCT) (e.g. glaucoma¹⁹ and multiple retinal diseases¹⁶).

A common concern, however, is that these articles have varying standards, and often lack an agreed, standardized format with respect to presenting methods, statistics, reporting metrics, and clinical translational value. Furthermore, many readers may not understand the nuances and technical details of this relatively new area of research in ophthalmology. How should readers read the AI research papers more effectively and apply them in the context of prior and current work? As editors of different journals, we have now been able to review a number of high-quality papers. In this article, we share some pointers and insights to help readers better understand, critically appraise and evaluate the data presented in AI papers in ophthalmology.

Corresponding author: Daniel SW TING, Assistant Professor, Duke-NUS Medical School, Consultant, Vitreo-retinal Department, Singapore National Eye Center, Adjunct Professor, State Key Laboratory of Ophthalmology, Zhongshan Ophthalmic Center, Sun Yat-sen University, Guangzhou, J. William Fulbright Scholar (2017-2018), Johns Hopkins University, daniel.ting.s.w@singhealth.com.sg, Address: 11 Third Hospital Avenue, Singapore 168751.

Conflict of interest: Drs Ting and Wong are the co-inventors and patent holders of a deep learning system for retinal diseases. Dr Ting is an editor of *Ophthalmology* and section editor for *British Journal of Ophthalmology*, Dr Aaron Lee is an editor of *American Journal of Ophthalmology and Translational Vision Science & Technology*, and Dr Wong is an editor of *JAMA Ophthalmology* and on International Advisory Board of *Lancet Digital Health*.

Is this study answering a question that matters?

Firstly, in the Introduction, high quality studies usually state clearly the unmet clinical or public health need, and the primary research question that needs to be tackled by AI technology. A clear example is the role of AI to screen for diabetic retinopathy (DR). At present, DR screening is commonly performed by the humans (e.g. ophthalmologists, optometrists, orthoptists or non-medical trained personnel), and often in countries with well-established and existing DR screening programs such as United States, United Kingdom and Singapore. With a projected rise in diabetes prevalence,²⁰ many countries do not have established DR screening programs (e.g. China, Indonesia, many countries in South America, Africa and Asia). To help address the limited and difficult to sustain resource in screening capacity, manpower costs and/or DR grading expertise, AI may be a helpful solution for countries, even in those with existing screening programs. Given the different regulations for health information technology in each country, the variety of deployment models, and varying thresholds for referable diabetic retinopathy, careful consideration is required before adoption.

As part of the introduction, a comprehensive literature search is often available on the information of similar technologies relevant to the specific diseases and the 'value-add' of the proposed AI system. Based on these information, the readers could better understand the clinical gap that the proposed AI algorithm may bridge, for example its application on specific target populations and clinical settings, using different classification outputs on fundus photographs (e.g. different ocular diseases - DR, glaucoma or AMD,⁸) different severity levels of a specific disease (e.g. early vs late glaucoma²¹) or same imaging device with different diagnostic outcomes (e.g. detection of specific retinal lesions^{22,23} vs triaging referral decisions using OCT¹⁶).

What are the core components in an AI system?

Many of the ophthalmology and clinician readers, without prior computer science background, may feel overwhelmed with the Methods section of AI studies with the usage of sophisticated technical terminologies. While it may not be necessary for readers to fully understand the entire technical architecture or mathematical formulas, it is important to appreciate the key components of the AI systems that are essential for clinical translation.

Broadly speaking, an AI system consists of two phases: 1) training/validation and; 2) testing. For training/validation, the AI system requires two main components: 1) training dataset consisting of clinical data/images and; 2) selection of a technical network (also known as convolutional neural network – CNN). More commonly, the majority of the dataset will be used for training and validation, followed by testing (e.g. 60%/20%/20%; 70%/20%/10% or 80%/10%/10% respectively), ideally at patients' level with no overlaps of same data/image in any phases to avoid the images from the same patient are used in the training and testing phase. Irrespective of the split in the training/validation datasets, it is more important to have sufficiently powered, preferably several independent testing datasets (which will be discussed in the later subsection – “How to clinically assess the diagnostic performance?”). A similar analogy would be that a basketball player could have different ways/styles of

training, but he/she would make a great basketball player if consistently performs well against players he/she has never seen before.

What makes a good training/testing datasets?

For the training/testing, it is important to have a sufficiently large and robust clinical dataset. Many readers may ask what would be the sufficient number or sample size required to train an AI algorithm, but these numbers may vary depending on how obvious or unique a condition is. In principle, conditions that have obvious, unique and easily distinguishable characteristics (e.g. neovascular AMD and proliferative DR on fundus photographs) may require less training examples as compared to those with abstract or subtle changes (early AMD with drusens, mild non-proliferative DR). Similarly, smaller OCT datasets are needed for training large hemorrhagic pigment epithelial detachment or full thickness macula hole compared to a subtle silver of subretinal fluid in AMD or mild diabetic macular edema.

In contrast, conditions that may have similar features would require larger number for training (e.g. differentiation between hard exudates versus drusen, microaneurysms vs small dot hemorrhages or pigments, severe non-proliferative DR vs central retinal vein occlusions on fundus photographs; cystoid macular edema secondary to different causes related to DR, AMD, retinal vein occlusions on OCTs). Thus, there is no simple answer to the question of “minimum sample size” as apart from the above-mentioned factors, it also depends on the specific research questions, technical methods and desired clinical output. Based on our observations, high quality AI papers published in the major journals use large data and image bank for AI algorithms development,^{7,8,16} although some technical approaches (e.g. transfer learning) have been shown to reduce the size of training datasets (this will be discussed further in the later subsection).²⁴

Apart from the actual number of data points or images, it is also important to understand exactly the number of eyes compared to patients involved in the training. This is because some AI algorithms may be trained on thousands of data points or images, but only from several hundreds of patients as these images or data were collected repetitively from the same patients. Such “large” datasets may lack diversity and potential generalizability to other populations. As in other studies, the study design and methods of recruitment (retrospective or prospective, inclusion and exclusion criteria, randomized trial versus cohort study) are also essential components for understanding the robustness of the training datasets.

What makes a good reference standard?

In particular, the reference standard (or also known as “ground truth” in technical terminology) acts as the “brain” for the AI algorithms. Thus, it is important to find out how the reference standard was derived. Many reference standards are based on human diagnosis and assessment. So for image-based research, the professional background and training, the experience and number of the human assessors should be specified. High-quality reference standard could consist of board-certified ophthalmologists and fellowship trained subspecialists, or certified non-medical professional graders or optometrists in reading

centers who have undertaken intensive training and accreditation with reproducible and consistent outcomes. This essential step avoids the “garbage-in, garbage-out” circumstance. There is some debate as to whether ophthalmologists or sub-specialists should be considered as the reference standard as increasingly reading centers have utilized professional graders as the gold standard for many clinical trials.^{25–27}

What are the appropriate machine learning or deep learning techniques?

In AI studies, there are many potential technical methods available to train an algorithm, and this largely depends again on the research question and the type of training dataset. Generally, the training dataset can be based on clinical data only, image only, or multi-modal models (clinical data + image + genetic and serum biomarkers). To analyse clinical data (e.g. electronic health records, population-based studies), many research groups have utilized conventional machine learning methods (e.g. random forest, support vector machines) and statistical methods (e.g. multi-variable logistic regression, generalized linear mixed models) to train the model, as these conventional methods appears sufficient to generate robust predictive algorithms. It is of course possible to adopt the newer deep learning methods for longitudinal and predictive clinical problems that have temporal sequence.^{28,29}

In contrast, for image-based training data (fundus photographs or OCT), deep learning is the most popular technique thus far because of increased diagnostic performance. Readers may come across a term called convolutional neural network (CNN), which are deep learning algorithms consisting of “neurons-like” computational layers that can have varying numbers of layers. The popular CNNs consists of AlexNet, VGGNet, Inception V4, ResNet and DenseNet,³ and these CNNs are commonly found off-the-shelf and can be downloaded from public domains. To further enhance the diagnostic performance, a technique known as “transfer learning” has been used. In this method, CNNs are usually pre-trained with the ImageNet database (consisting of millions of images such as cars, animals and etc), before being applied to the specific dataset in question. The “value-add” of this approach is reduced if there is sufficient number in the training datasets. It is usually not necessary for readers to fully understand how these CNNs are constructed mathematically, but to understand the next important step in the Methodology: the overall operational flow of the CNN within an AI system.

What makes an excellent operational flow of an AI system?

An easy to understand and user-friendly AI system would ideally describe the operational workflow to the readers. This is similar to describing “a patient journey” from registration, followed by clinical consultation and diagnosis. First, in terms of AI registration, the readers can first find out the ability of an AI system in detecting the types of ocular images (anterior segment, fundus and OCT photographs), field of view (macula-centered versus optic disc centered) and gradeability (gradable vs non-gradable). Second, once the AI input has been registered, it would be seen by the ‘doctors’ (in this case, the AI algorithm) for diagnosis and consultation. During this step, the image can be preprocessed (cropping, contrast enhancement and etc), followed by the analysis of the AI system (CNNs) and generation of diagnosis (output classifications, sometimes binary but sometimes multi-class outputs).

Lastly, some groups have also reported using visualization maps to highlight the abnormal areas detected by the AI system,³⁰ serving as good clinical decision support tools for clinical implementation. However, similar steps may not be applicable to other non-image modalities (e.g. visual field, clinical data).

How to clinically assess the diagnostic performance?

Earlier this year, the US Food and Drug Administration (FDA) published a framework on Artificial Intelligence/Machine Learning - Based Software as a Medical Device (SaMD), using the definition by the International Medical Device Regulators Forum (IMDRF) to consider AI-based software as a medical device (SaMD).³¹ Apart from advocating good quality systems and machine learning practices, it also stated the importance of evaluating the AI system for validity in having correct AI input data to generate accurate output data and between the AI output and target condition/intended purposes in clinical care.

For clinical adoption, it is useful to look at the disease prevalence on the described testing datasets, and to better understand whether it is conducted in a population-based or clinic-based settings. The sample size power calculation for testing datasets may include the disease prevalence, type 1 and 2 errors with the 95% confidence intervals, similar to other clinical trials or diagnostic tests. Showing the reproducibility of an AI system in an out-of-sample, geographically distinct population would also increase the quality of the presented work. Apart from that, it is also important for the readers to pay attention to the disease classification systems adopted in the AI studies (e.g. functional versus structural diagnosis for glaucoma,³² the UK National Health Service DR classifications³³ vs the International Classification DR severity scales).³⁴

Many AI studies report the area under the receivers' operating curve (AUC), a popular statistical method that most clinicians are familiar with, although a precision-recall curve (a method that is commonly used in the technical or machine learning world) may be more appropriate for imbalanced datasets.³⁵ From what we have observed, many high-quality AI papers describe the rationale on how an operating threshold is determined in the training dataset, and demonstrated the sensitivity and specificity on the independent datasets performed on same operating threshold. We also observed that the accuracy, positive predictive, negative predictive values, Cohen's kappa may be reported, although this may not be as popular as the traditional metrics of AUC, sensitivity and specificity. The recommended reporting metrics were sometimes included as tables and figures format, especially if there are multiple independent testing datasets.⁸ Given the variability of the reporting metrics, a consensus on the best metrics for ophthalmology related AI research may be warranted.

What to watch out for in an AI system for clinical adoption?

To help the readers' understanding on how the AI system can be integrated and adopted within clinical practice, high-quality AI papers generally will include the limitations of the specific AI systems (e.g. the AI system can only detect high quality images, or lack ability to diagnose other conditions). Many AI research groups conduct due diligence about their

system, especially when they were only tested in the research settings. In fact, we encourage transparency in reporting to expedite clinical translation of AI technology.

How can we better guide AI research in the field?

As the field evolves, we must determine and reach a consensus in deciding how best to evaluate and critically appraise AI research in ophthalmology to ensure robustness and reproducibility of algorithms and generalizability of study findings to real-world clinical settings. For diagnostic tests, the Standards for Reporting of Diagnostic Accuracy (STARD) steering committee have proposed a list of 25 items to standardize the reporting standard in diagnostic tests studies.³⁶ AI algorithms could also be considered as novel diagnostic tests. Hence, a good-quality study could include most of the STARD criteria, with additional information on the training methodologies. Furthermore, a number of generic AI guidelines (not specific to ophthalmology) are currently under development, such as the AI extension of the TRIPOD guidelines (TRIPOD-AI).³⁷

Furthermore, it is worth noting that in the traditional biomarker research, the term “validation” is generally used as the way to confirm the diagnostic value of a biomarker, which is dissimilar to what it is used in the AI settings (cross-validation is used as part of the training phase as mentioned earlier). Nevertheless, many AI groups have used different terminologies as the “testing” phase and this may potentially create some confusions to the readers when they browse through the articles.

Our editorial summarizes some observations which may help readers navigate, review and appreciate AI papers in ophthalmology (Table 1). The rapid emergence of this area of research and the lack of standardized reporting format and standards suggests the need for a formation of a working taskforce to develop consensus and guidelines for AI research in ophthalmology, possibly done by relevant experts using a Dephi approach. This framework may help guide researchers, clinicians and policy makers better evaluate, apply and translate AI technologies in ophthalmology, and harness the power of the big data revolution that will transform healthcare.

References

1. LeCun Y, Bengio Y, Hinton G. Deep learning. *Nature*. 2015;521(7553):436–444. [PubMed: 26017442]
2. Schmidt-Erfurth U, Sadeghipour A, Gerendas BS, Waldstein SM, Bogunovic H. Artificial intelligence in retina. *Prog Retin Eye Res*. 2018;67:1–29. [PubMed: 30076935]
3. Ting DSW, Peng L, Varadarajan AV, et al. Deep learning in ophthalmology: The technical and clinical considerations. *Prog Retin Eye Res*. 2019.
4. Lee A, Taylor P, Kalpathy-Cramer J, Tufail A. Machine Learning Has Arrived! *Ophthalmology*. 2017;124(12):1726–1728. [PubMed: 29157423]
5. Hwang ES, Perez-Straziota CE, Kim SW, Santhiago MR, Randleman JB. Distinguishing Highly Asymmetric Keratoconus Eyes Using Combined Scheimpflug and Spectral-Domain OCT Analysis. *Ophthalmology*. 2018;125(12):1862–1871. [PubMed: 30055838]
6. Yu F, Croso GS, Kim TSea. Assessment of Automated Identification of Phases in Videos of Cataract Surgery Using Machine Learning and Deep Learning Techniques. *JAMA Netw Open*. 2019;2(e191860).

7. Gulshan V, Peng L, Coram M, et al. Development and Validation of a Deep Learning Algorithm for Detection of Diabetic Retinopathy in Retinal Fundus Photographs. *JAMA*. 2016;316(22):2402–2410. [PubMed: 27898976]
8. Ting DSW, Cheung CY, Lim G, et al. Development and Validation of a Deep Learning System for Diabetic Retinopathy and Related Eye Diseases Using Retinal Images From Multiethnic Populations With Diabetes. *JAMA*. 2017;318(22):2211–2223. [PubMed: 29234807]
9. Gargeya R, Leng T. Automated Identification of Diabetic Retinopathy Using Deep Learning. *Ophthalmology*. 2017;124(7):962–969. [PubMed: 28359545]
10. Abramoff MD, Lou Y, Erginay A, et al. Improved Automated Detection of Diabetic Retinopathy on a Publicly Available Dataset Through Integration of Deep Learning. *Investigative ophthalmology & visual science*. 2016;57(13):5200–5206. [PubMed: 27701631]
11. Abramoff MD, Lavin PT, Birch M, Shah N, Folk JC. Pivotal trial of an autonomous AI-based diagnostic system for detection of diabetic retinopathy in primary care offices. *NPJ Digital Medicine*. 2018;39:1–8.
12. Li Z, He Y, Keel S, Meng W, Chang RT, He M. Efficacy of a Deep Learning System for Detecting Glaucomatous Optic Neuropathy Based on Color Fundus Photographs. *Ophthalmology*. 2018;125(8):1199–1206. [PubMed: 29506863]
13. Burlina PM, Joshi N, Pekala M, Pacheco KD, Freund DE, Bressler NM. Automated Grading of Age-Related Macular Degeneration From Color Fundus Images Using Deep Convolutional Neural Networks. *JAMA ophthalmology*. 2017;135(11):1170–1176. [PubMed: 28973096]
14. Burlina P, Joshi N, Pacheco KD, Freund DE, Kong J, Bressler NM. Utility of Deep Learning Methods for Referability Classification of Age-Related Macular Degeneration. *JAMA ophthalmology*. 2018;136(11):1305–1307. [PubMed: 30193354]
15. Peng Y, Dharssi S, Chen Q, et al. DeepSeeNet: A Deep Learning Model for Automated Classification of Patient-based Age-related Macular Degeneration Severity from Color Fundus Photographs. *Ophthalmology*. 2019;126(4):565–575. [PubMed: 30471319]
16. De Fauw J, Ledsam JR, Romera-Paredes B, et al. Clinically applicable deep learning for diagnosis and referral in retinal disease. *Nat Med*. 2018;24(9):1342–1350. [PubMed: 30104768]
17. Redd TK, Campbell JP, Brown JM, et al. Evaluation of a deep learning image assessment system for detecting severe retinopathy of prematurity. *Br J Ophthalmol*. 2018.
18. Brown JM, Campbell JP, Beers A, et al. Automated Diagnosis of Plus Disease in Retinopathy of Prematurity Using Deep Convolutional Neural Networks. *JAMA ophthalmology*. 2018;136(7):803–810. [PubMed: 29801159]
19. Fu H, Baskaran M, Xu Y, et al. A Deep Learning System for Automated Angle-Closure Detection in Anterior Segment Optical Coherence Tomography Images. *American journal of ophthalmology*. 2019;203:37–45. [PubMed: 30849350]
20. Yau JW, Rogers SL, Kawasaki R, et al. Global prevalence and major risk factors of diabetic retinopathy. *Diabetes Care*. 2012;35(3):556–564. [PubMed: 22301125]
21. Ahn JM, Kim S, Ahn KS, Cho SH, Lee KB, Kim US. A deep learning model for the detection of both advanced and early glaucoma using fundus photography. *PLoS One*. 2018;13(11):e0207982. [PubMed: 30481205]
22. Lee CS, Baughman DM, Lee AY. Deep Learning Is Effective for Classifying Normal versus Age-Related Macular Degeneration OCT Images. *Ophthalmol Retina*. 2017;1(4):322–327. [PubMed: 30693348]
23. Schmidt-Erfurth U, Bogunovic H, Sadeghipour A, et al. Machine Learning to Analyze the Prognostic Value of Current Imaging Biomarkers in Neovascular Age-Related Macular Degeneration. *Ophthalmol Retina*. 2018;2(1):24–30. [PubMed: 31047298]
24. Ting DSW, Liu Y, Burlina P, Xu X, Bressler NM, Wong TY. AI for medical imaging goes deep. *Nat Med*. 2018;24(5):539–540. [PubMed: 29736024]
25. Diabetic Retinopathy Clinical Research N, Wells JA, Glassman AR, et al. Aflibercept, bevacizumab, or ranibizumab for diabetic macular edema. *The New England journal of medicine*. 2015;372(13):1193–1203. [PubMed: 25692915]
26. Sivaprasad S, Prevost AT, Vasconcelos JC, et al. Clinical efficacy of intravitreal aflibercept versus panretinal photocoagulation for best corrected visual acuity in patients with proliferative diabetic

- retinopathy at 52 weeks (CLARITY): a multicentre, single-blinded, randomised, controlled, phase 2b, non-inferiority trial. *Lancet*. 2017;389(10085):2193–2203. [PubMed: 28494920]
27. Writing Committee for the Diabetic Retinopathy Clinical Research N, Gross JG, Glassman AR, et al. Panretinal Photocoagulation vs Intravitreal Ranibizumab for Proliferative Diabetic Retinopathy: A Randomized Clinical Trial. *JAMA*. 2015;314(20):2137–2146. [PubMed: 26565927]
28. Rajkomar A, Oren E, Chen K, et al. Scalable and accurate deep learning with electronic health records. *NPJ Digit Med*. 2018;1:18. [PubMed: 31304302]
29. Wen JC, Lee CS, Keane PA, et al. Forecasting future Humphrey Visual Fields using deep learning. *PLoS One*. 2019;14(4):e0214875. [PubMed: 30951547]
30. Sayres R, Taly A, Rahimy E, et al. Using a Deep Learning Algorithm and Integrated Gradients Explanation to Assist Grading for Diabetic Retinopathy. *Ophthalmology*. 2019;126(4):552–564. [PubMed: 30553900]
31. Food US and Administration Drug. Proposed Regulatory Framework for Modifications for Artificial Intelligence/Machine Learning - Based Software as a Medical Device (SaMD). URL: <https://www.fda.gov/media/122535/download> [Accessed on 17th August, 2019]. 2019.
32. Quigley HA. Glaucoma. *Lancet*. 2011;377(9774):1367–1377. [PubMed: 21453963]
33. Harding S, Greenwood R, Aldington S, et al. Grading and disease management in national screening for diabetic retinopathy in England and Wales. *Diabet Med*. 2003;20(12):965–971. [PubMed: 14632697]
34. Wilkinson CP, Ferris FL 3rd, Klein RE, et al. Proposed international clinical diabetic retinopathy and diabetic macular edema disease severity scales. *Ophthalmology*. 2003;110(9):1677–1682. [PubMed: 13129861]
35. Davis JL MG The relationship between Precision-Recall and ROC curves. Proceedings of the 23rd international conference on Machine learning. 2006:233–240.
36. Bossuyt PM, Reitsma JB, Bruns DE, et al. Towards complete and accurate reporting of studies of diagnostic accuracy: The STARD Initiative. *Ann Intern Med*. 2003;138(1):40–44. [PubMed: 12513043]
37. Collins GS, Reitsma JB, Altman DG, Moons KG. Transparent Reporting of a multivariable prediction model for Individual Prognosis or Diagnosis (TRIPOD): the TRIPOD statement. *Ann Intern Med*. 2015;162(1):55–63. [PubMed: 25560714]

Table 1:

The Readers' Guide to Browse Through an Artificial Intelligence Article in Ophthalmology

Introduction	<p>Validity of the research question</p> <p>Comprehensive literature search of similar technologies related to the specific disease</p> <p>Clinical unmet need</p> <p>“Value-add” of the proposed AI system</p>
Methods	
Core components	Clinical datasets and technical network
Clinical datasets	Division of training, validation and testing datasets
Dataset descriptions	<ol style="list-style-type: none"> 1. Number of images, eyes and patients 2. Inclusion and exclusion criteria for these patients 3. Study design (prospective vs retrospective), patients demographics (optional) 4. Recruitment methods (consecutive, randomised and etc) and sites 5. Prevalence of positive vs control cases 6. Types input data - clinical data, imaging test or others
Technical methodology	<ol style="list-style-type: none"> 1. Technical approach (deep learning, machine learning or statistical approach) 2. Types of neural network 3. Operational flow of an AI system
Assessment of the diagnostic performance	<ol style="list-style-type: none"> 1. Power calculation of the testing datasets 2. Receivers' operating curve (AUC), sensitivity and specificity (with 95% confidence interval) 3. Accuracy, positive predictive or negative predictive value 4. Cohen's kappa
Reference Standard	<p>Numbers and experience of graders</p> <p>(e.g. Graders from reading centers, retinal specialists and etc)</p> <p>Disease Classification System</p>
Statistical analysis and results (all with 95% CI)	<ol style="list-style-type: none"> 1. Area under receivers' operating curve (AUC) 2. Sensitivity and specificity 3. Accuracy, positive predictive value and negative predictive value 4. Cohen Kappa 5. Dice coefficients (for segmentation tasks)
Discussion	
Clinical translational value	<ol style="list-style-type: none"> 1. Clinical application of the AI solution 2. Limitation of the AI systems 3. Potential deployment methods