

**Title:** The accuracy of automated computer-aided diagnosis for stroke imaging. A critical evaluation of current evidence.

**Authors:** Joanna M Wardlaw,<sup>1</sup> MD, FRCR, FMedSci  
Grant Mair,<sup>1</sup> MD, FRCR  
Rüdiger von Kummer,<sup>2</sup> MD  
Michelle C Williams,<sup>3</sup> PhD, FRCR  
Wenwen Li,<sup>1</sup> PhD  
Amos J Storkey,<sup>4</sup> PhD  
Emanuel Trucco,<sup>5</sup> PhD  
David S. Liebeskind,<sup>6</sup> MD  
Andrew Farrall,<sup>1</sup> PhD, FRC  
Philip M Bath,<sup>7</sup> PhD, FRCP  
Philip White,<sup>8</sup> MD, FRCR

**Affiliations:** 1: Centre for Clinical Brain Sciences, UK Dementia Research Institute Centre at the University of Edinburgh, Chancellor's Building, Little France, Edinburgh, EH16 4SB, UK  
2: Institute of Diagnostic and Interventional Neuroradiology, Universitätsklinikum Carl Gustav Carus, Dresden, Germany  
3: Centre for Cardiovascular Science, University of Edinburgh, Chancellor's Building, 49 Little France, Edinburgh, EH16 4SB, UK  
4: School of Informatics, University of Edinburgh.  
5: VAMPIRE project, Computing, School of Science and Engineering, University of Dundee.  
6: Neurovascular Imaging Res Core, UCLA, Los Angeles, CA  
7: Stroke Trials Unit, Mental Health & Clinical Neuroscience, University of Nottingham, Queen's Medical Centre campus, Derby Road, Nottingham NG7 2UH, UK  
8: Translational and Clinical Research Institute, Newcastle University, Henry Wellcome Building, Framlington Place, Newcastle upon Tyne and Newcastle upon Tyne Hospitals NHS Trust, UK

**Correspondence:** JMW as above

**Running Title:** AI in acute stroke diagnosis

**Key words:** stroke; machine learning; artificial intelligence;

**Disclosures:** GM received speaker and consulting fees from Canon Medical Ltd, Europe. PMB provided consultancy for Diamedica, Phagenesis, Moleac, Sanofi, Nestle. MCW received speaker fees from Canon Medical Systems. DSL consultants for imaging core lab services to Cerenovus, Genentech, Medtronic, Stryker and Rapid Medical. PMW reports institutional unrestricted educational grants from Medtronic, Stryker and Penumbra; Institutional research funding from Microvention and IBEX and consulting fees/core lab services from Microvention and E-VASC. AJF received institutional educational grants from Medtronic for the development of training materials in stroke radiology. AS has received funding from Huawei, ARM and Microsoft Research, Cambridge not related to the topic of the current paper. JMW holds institutional academic grants from the

Fondation Leducq, EU H2020 programme, Stroke Association, Health Data Research UK, British Heart Foundation, Alzheimer's Society, Weston Brain Institute, and Row Fogo Charitable Trust for research unrelated to the

Funding: JMW is supported by the UK Dementia Research Institute, which receives its funding from DRI Ltd, funded by the UK Medical Research Council, Alzheimer's Society and Alzheimer's Research UK. GM is Stroke Association Edith Murphy Foundation Senior Clinical Lecturer (SA L-SMP 18\1000). We acknowledge funding from the Stroke Association (TSA\_CR\_2017/01) and Medical Research Council Proximity to Discovery fund (MC\_PC\_17188) for the RITeS Study. RvK receives royalties as Editor-in-Chief of *Neuroradiology*. WL is funded by Health Data Research UK which receives its funding from the Medical Research Council. MCW is supported by the British Heart Foundation (FS/ICRF/20/26002). PMB is supported by the British Heart Foundation and NIHR Health Technology Programme.

Acknowledgements: PMB is Stroke Association Professor of Stroke Medicine and an Emeritus NIHR Senior Investigator.

## **Abstract**

There is increasing interest in computer applications, [using artificial intelligence \(AI\) methodologies](#), to perform healthcare tasks previously performed by humans, particularly in medical imaging for diagnosis. In stroke, there are now commercial [AI software](#) for use with CT or MR imaging to identify acute ischemic brain tissue pathology, arterial obstruction on CT angiography or as hyperattenuated arteries on CT, brain hemorrhage, or size of perfusion defects. A rapid, accurate diagnosis may aid treatment decisions for individual patients and could [improve](#) outcome if it leads to effective and safe treatment; or conversely, to disaster if a delayed or incorrect diagnosis results in inappropriate treatment. Despite this potential clinical impact, diagnostic tools [including AI methods](#) are not subjected to the same clinical evaluation standards as are mandatory for drugs. Here, we provide an evidence-based review of the pros and cons of automated methods for medical imaging diagnosis, including those based on artificial intelligence (AI), to diagnose acute brain pathology on CT or MRI in patients with stroke.

## **Introduction**

Stroke is common, hospitals are busy, delays=lost brain; diagnosis of the cause should be rapid so that appropriate treatment can start and give the patient the best chance of independent survival. Brain imaging is essential to differentiate ischemic from hemorrhagic stroke and stroke mimics. Furthermore, with advances in treatment options for specific patient subgroups, it is not enough just to identify ischemia or hemorrhage, since the size of the acute lesion, presence of other features (obstructed arteries, mass effect), and pre-stroke changes (leukoaraiosis, old infarcts, brain atrophy) also influence management. Most acute general hospitals assess several patients with suspected stroke each day, meaning that all steps in the process, including diagnosis, should be rapid, timely, efficient and accurate.

It takes many years to train a neuroradiologist, they are scarce in many countries, and serve many disease areas in addition to stroke. Vascular neurologists and stroke physicians learn to recognize early features of ischemic brain on scanning and major contraindications to reperfusion treatment. However, early ischemic changes on non-contrast CT can be subtle and complex, with serious implications hanging on their correct identification, fueling interest in ways to improve their recognition *and quantification*. Methods such as perfusion CT require post-processing to generate a diagnostic image, highlighting abnormalities such as thresholded tissue blood flow. Such 'cleaned' images may appear more user friendly and may facilitate rapid interpretation.

Alongside these longstanding pressures to reduce time and increase diagnostic accuracy, there have been substantial advances in computer vision and artificial intelligence (AI) technologies across all walks of life. At its most general, AI refers to use of computers to solve problems in ways that mimic human behavior; *machine learning* (ML) is of the key technology behind AI where computer algorithms learn from examples ('ground truth') without explicit programming which properties of the data are relevant for a given problem ('feature selection'); and *deep learning* (DL) is a subset of ML that uses biologically inspired neural networks to learn abstract high-order features in any type of data without requiring predetermined inputs.<sup>1</sup>

Medical imaging has been an obvious target for these developments.<sup>2-4</sup> Several commercial CT and MR scan diagnostic software for stroke are now in use in many hospitals. Nonetheless, the major demand for accelerated diagnosis in acute stroke, the fascination with the latest imaging tools, and huge potential financial gains for industry,<sup>5</sup> should not cloud the essential need to demonstrate that AI tools are accurate, are improving not impeding healthcare, and that the benefits outweigh the potential harms for patients.<sup>6</sup>

We assess the current evidence for AI diagnostic imaging tools in stroke, commercially available or in clinical use, the motivations, expectations and work needed to underpin their implementation into clinical practice.

## **Review of Evidence**

### **1. What could automated AI imaging technologies achieve in acute stroke diagnosis?**

AI technologies could improve accuracy, speed and standardization of stroke diagnosis,<sup>7-11</sup> particularly in low throughput Centres,<sup>12</sup> and improve prognostication using quantitative measures of acute and chronic brain injury. AI tools could also improve workflows and communication along thrombectomy referral pathways, reduce time to treatment, improve clinical outcomes and save clinicians and radiologists time.<sup>2</sup>

AI software could be most useful to assess the cause of stroke and its likely pathology:

1. Brain focal ischemia vs. intracranial hemorrhage vs. “stroke mimics” (migraine, seizure), the latter two requiring different management to ischemic stroke. Intracranial hemorrhages (brain hemorrhage, subarachnoid hemorrhage (SAH), subdural hematoma) can be visualized on CT and MRI with high accuracy. “Stroke mimic” is a clinical diagnosis based on the exclusion of brain ischemia and hemorrhage, noting that in patients with ischemic stroke symptoms, CT and MRI may appear normal, or show various degrees of reduced attenuation/altered signal intensity and/or swelling in the affected tissue.
2. Arterial pathology: site of acute embolic or thrombotic arterial occlusion, secondary features (e.g. collateral blood flow), and underlying pathologies (atherosclerosis, dissection, inflammation, vasospasm). This is, perhaps, the most important determinant of acute therapy decisions and related prognosis assessment. Thin slice non-enhanced CT, angiography (CTA, MRA, DSA) and arterial wall imaging may differentiate between embolic or local arterial disease and have an impact on secondary stroke prophylaxis.
3. Early ischemic brain tissue alterations: isolated tissue swelling due to autoregulatory vasodilation,<sup>13</sup> tissue swelling with reduced attenuation (indicating early focal net water uptake or ionic edema triggered by critically low CBF <15 ml/100g/min, which cannot be tolerated by brain tissue for more than about 30 min), ion pump failure triggered by a CBF < 30 ml/100g/min with consequent neuronal dysfunction and cellular edema indicated by ADC decrease/increase of DWI signal intensity. The still limited understanding of ischemia-effected brain tissue alterations makes it a potential mistake that algorithms are trained to detect “brain infarction” on non-enhanced CT or DWI within 6 hours of symptom onset,<sup>14,15</sup> since focal brain ischemia of up to 6-hour duration does not induce coagulation necrosis.<sup>16,17</sup>
4. Pre-stroke changes: leukoaraiosis (or white matter hyperintensities), prior infarcts or hemorrhages, and brain atrophy are all associated with worse short-term (haemorrhagic transformation, dependency, death)<sup>18,19</sup> and long term (dependency, recurrent stroke, death, cognitive impairment, dementia) outcomes.<sup>19</sup>

Commonly used approaches for developing AI tools in stroke start with identifying attenuation or signal change that indicates ischemic tissue and its extent, adapted from visual rating tools. ~~Visual rating tools, such as the Alberta Stroke Programme Early CT Score (ASPECTS) is a widely used visual rating tool;~~ used to help diagnose acute ischemic stroke and select patients for thrombolysis and thrombectomy. Table 1 lists eighteen-19 vendors that currently provide 31-32 different commercially available AI software packages for ischemic stroke to assess a

range of features including: ischemic tissue change (CT or MRI); hyperattenuated arteries (a surrogate for arterial thrombus); hemorrhage on CT; large artery occlusion on CTA; and 'salvageable tissue' from CT or MR perfusion imaging. ~~Despite several drawbacks of the ASPECTS (not an independent outcome predictor,<sup>18</sup> variable ASPECTS cut points, validity,<sup>20</sup> etc),~~ ASPECTS is used by seven of eighteen companies currently offering commercial stroke AI software.<sup>7, 21</sup> ~~despite several drawbacks of ASPECTS (not an independent outcome predictor,<sup>18</sup> variable cut points, validity<sup>20</sup>).~~ Four companies provide comprehensive packages for handling non-enhanced, angiographic and perfusion imaging in one workflow (Brainomix, NICO.Lab, RapidAI and Viz.AI), while three companies combine assessment of non-enhanced and angiographic imaging only (Aidoc, Avicenna, Circle Neurovascular Imaging). Other companies assess different combinations or individual components of ischemic stroke or haemorrhagic stroke (Table 1).

## **2. What drives the increasing use of automated analysis technologies?**

There is undoubtedly a real need to aid busy doctors at the hospital front door, and radiologists are also increasingly pressured. However, two decades after early thrombolysis trials, clinical awareness of acute stroke imaging features is much more established. Therefore, any AI diagnostic tools have to be exceptionally fast, easy to use, and accurate on 'real world' data (see Section 6) in order to add value. AI tools also require users to be trained in their proper use and interpretation.

Commercially, data and AI offer massive financial gains for successful products.<sup>5</sup> Worldwide spending on AI was estimated at US\$38 billion in 2019 and is predicted to rise to \$98 billion by 2023.<sup>22</sup> Investments in AI-based medical imaging companies in the USA reached US\$1.17 billion between January 2014 and January 2019, doubling since 2012-2017, while the number of companies in the AI market had trebled,<sup>5</sup> including new industries devoted to image classification (section 7). Major medical imaging manufacturers are incorporating AI tools [into consoles](#) to retain a marketing edge. AI requires data storage capacity and computing power: between Jan 2014 and Jan 2019, there were over US\$435 billion-worth of cloud-based medical imaging transactions, indicating massive investments in these areas.<sup>5</sup>

## **3. What clinically relevant testing should automated analysis undergo in acute stroke?**

Currently, AI software for radiology can only be marketed in the European Economic Area (EEA) after achieving a CE mark (Conformité Européenne), indicating that the technology conforms with European health, safety, and environmental protection standards, and in the USA with FDA (Food and Drug Administration) approval. However, the clinical standards for achieving these certifications are low as compared with licensing a new drug. Both CE and FDA systems have different classes of approval depending on the perceived risk to the patient and to software users. Since radiology software have so far been designed to *support* rather than *replace* physicians,

they require usually only Class I or II approval denoting low to medium risk. The 'decision support' labelling raises important questions about what happens when the clinician disagrees with the AI diagnoses, and who is to blame when one or other gets it wrong. Nonetheless, in both territories, Class I approval is awarded without external scrutiny and companies self-declare these products as compliant. While Class II approval usually includes the submission of evidence assessed by an independent body, companies can provide this evidence from their own internal testing without peer-review or publication. Indeed, a recent independent review of all CE-marked AI software for radiology in Europe found that 64 of 100 products had no published peer-reviewed evidence of efficacy.<sup>3</sup>

It is beyond the scope of this article to assess the depth of peer-reviewed evidence for every available software, so we focus on commercial products with the most citations according to recent reviews (Brainomix, RapidAI, and Viz.ai).<sup>8-11</sup>

Reporting guidelines and methodological standards for developing AI in medical imaging are available<sup>23</sup> including SPIRIT-AI,<sup>24</sup> CONSORT-AI,<sup>25</sup> and checklists.<sup>26</sup> Several societies have released their own guidelines and position statements.<sup>26,27</sup> The FDA <https://www.fda.gov/media/145022/download>, and British Standards Institute (BSI) and Medicines and Healthcare Regulatory Agency (MHRA) <https://standardsdevelopment.bsigroup.com/committees/50299208> (pending public consultation) provide standards.

~~There is little evidence that~~ Few if any evaluations of AI software are meeting these standards, generally or in stroke.<sup>6,8,9,23,24,28,29</sup> A recent systematic review of reporting quality of studies of ML in medical diagnosis found 28 includable studies but none mentioned a reporting guideline, few mentioned the distribution of disease severity or alternative diagnoses, most studies had a long delay between the reference standard and ML diagnoses, and in half of studies the population was of uncertain relevance to the clinical setting.<sup>28</sup> Five of the 28 included studies concerned brain imaging (~~addressing multiple sclerosis, attention deficit hyperactivity disorder, distinguishing minimal consciousness from unconsciousness, and Alzheimer's disease~~), but none of which addressed stroke.<sup>28</sup>

#### **4. Accuracy of AI in stroke:**

*i) How accurate is AI software for differentiating acute ischemic from haemorrhagic stroke and stroke mimics, and on which imaging modalities?*

Recent systematic reviews,<sup>8,9</sup> narrative reviews,<sup>7,10-12,14</sup> a review protocol<sup>30</sup> and a pending combined analysis of nine large stroke trials<sup>31</sup> show that most studies of AI have focused on ischemic not haemorrhagic stroke, and CT not MRI. The systematic reviews identified 20 (tissue and arterial changes<sup>8</sup>) and 68 (non-contrast CT only<sup>9</sup>) includable studies, but most studies were small, provided little documentation of the patient characteristics, recruitment or CT characteristics, and focused on comparing AI against feature detection by humans, not on

clinical outcomes. Rates of failed scan processing were often omitted. Many methodological differences between studies precluded formal meta-analyses. All published studies and reviews focused on AI detection of ischemic stroke features including arterial obstruction,<sup>7-9,12</sup> with only one pending study<sup>31</sup> assessing if an AI software can differentiate ischemic from haemorrhagic stroke or mimics, ~~which are just as important in clinical practice.~~

The extent to which AI software may be affected by patient-related (background brain changes or alternative brain pathology, movement and position during scanning, heart failure, metallic artefacts<sup>32-35</sup>) and imaging-related (scanner manufacturer, acquisition parameters such as slice thickness, or CT gantry position<sup>34, 36</sup>) factors is beginning to emerge. These affect the likelihood of successful image processing, agreement with a reference standard, and strength of associations with clinical outcome. However, most published evidence excluded difficult cases before analysis, while results of studies that retained difficult cases were often less positive.<sup>32,37,38</sup>

*ii) How good are these technologies for identifying key features of acute ischemic stroke that are of prognostic importance?*

Few studies assessed AI software-based diagnoses and clinical outcomes. Amongst commercially available AI tools, we reviewed the published literature for the three most established providers offering comprehensive imaging packages: Brainomix, RapidAI, and Viz.ai, (Table 2). We used published reviews<sup>7-9,21</sup> updated by searching Pubmed for company and software names, and the vendor's websites. We focus on studies with the largest test datasets (ideally >100 patients), and report diagnostic accuracy statistics for stroke feature detection.

Three studies assessed detection of tissue ischemia (all Brainomix, all retrospective, two independent of the company,<sup>32,39</sup> total patients n=367),<sup>32,39,40</sup> with sensitivities of 44-83% and specificities of 57-93% (Table 2). Six studies assessed detection of large vessel occlusion (LVO) (three vendors, all retrospective, one independent of the vendor,<sup>41</sup> total n=2635)<sup>38,41-45</sup> with sensitivities of 80-96% and specificities of 90-98%. Only one study each assessed hemorrhage detection<sup>46</sup> and MRI diffusion or perfusion imaging.<sup>47</sup> Compared with optimal circumstances, the agreement between each of these software and experts was poorer in patients with leukoaraiosis, old infarcts, or other parenchymal defects.<sup>7</sup> This underscores the importance of evaluating AI tools in realistic and common clinical settings where patients are often older, have multiple conditions or delayed presentations,<sup>48</sup> and not relying on results when tested in the simplified training scenarios common in public datasets.<sup>5,49</sup>

*iii). On a population level, how many false positives or negatives might arise per 100 typical suspected strokes and what are the implications for patient outcomes?*

Given the range of published sensitivity and specificity results for stroke feature detection by AI software above, we translate the results as follows, Figure 1. For every 100 patients assessed using these software:



- With ischemic stroke, ischemia will be correctly detected in 44-83 but missed in 17-56.
- Without ischemic stroke, ischemia will be incorrectly detected in 7-43.
- With LVO, occlusion will be correctly detected in 80-96 but missed in 4-20.
- Without LVO, occlusion will be incorrectly detected in 2-10.

*iv). Have (any) automated, including AI, technologies that are proposed for use in stroke, undergone proper prospective randomised blinded outcome clinical trial assessment to determine impact on clinical outcomes or health economics?*

Randomised controlled trials of AI technologies are scarce and mostly ongoing: a recent survey of trials' registries and the literature identified only one RCT comparing DL with clinicians in medical imaging (on breast ultrasound).<sup>6,29</sup> We identified one ongoing multicenter RCT testing the impact of Viz LVO on stroke workflow and 90 day clinical outcome in 500 participants admitted with stroke and suspected LVO in the USA (Automated Detection and Triage of Large Vessel Occlusions Using Artificial Intelligence for Early and Rapid Treatment, ALERT, NCT04142879). Diagnostic accuracy is not a primary or specified secondary outcome in ALERT. Another ongoing multicenter RCT in China (GOLDEN BRIDGE II, NCT04524624) is testing AI identification of stroke on diffusion imaging plus decision support versus usual care in 21689 patients requiring secondary stroke prevention.

Three studies compared times to thrombectomy before and after introduction of RAPID<sup>41,50</sup> or Viz LVO,<sup>51</sup> reporting average reductions of 30 minutes to groin puncture; however, all were retrospective, two only report [the small](#) numbers of patients who all received thrombectomy, and before-after comparisons are unable to address many sources of bias.

***5. How do stroke AI technologies compare with other medical AI technologies, particularly medical imaging AI, in terms of stage of development and quality and thoroughness of assessment?***

Stroke AI is similar to other medical imaging AI – great hopes but important challenges for delivery into clinical practice. These challenges reflect data curation, model development, relevance to clinical practice, potential to introduce and amplify biases, the AI tool's transparency, and evidence of accuracy, impact on outcomes and cost effectiveness meeting RCT evidence standards.

The quality, quantity, diversity, and provenance of the data used to train a ML model are critical to its utility in clinical practice. Many current papers describe ML models trained on one small dataset from one hospital,<sup>52,53</sup> [insufficient](#) to be useful on the variety seen in clinical practice.<sup>53</sup> Commercial ML models have similar issues, particularly “black box” models where the training is not described.

Datasets created as part of international challenges, e.g. the RSNA 2019 Brain CT Hemorrhage Challenge<sup>54</sup> (874035 images, multiple institutions) [have limitations, including of are limited in](#) diversity, accuracy and

reproducibility.<sup>55</sup> ~~Over-Over-tuning of the software to particular datasets (overfitting) results in poor impedes~~ generalization. ~~There is a risk that~~ Sometimes AI models identify confounders rather than target disease, e.g. hip fracture detection, where scanner model (emergency department) and “priority request” marker predicted fracture better than the imaging findings ~~themselves~~.<sup>56</sup> ‘Explainable AI models’ may help to show the underlying features identified by a DL model to avoid ‘black box’ problems. ~~Even so, it is an open question as to how much such explanations are valid, interpretable by a clinician, or even what ‘explainability’ should mean in a clinical setting.~~

Bias emerging during the development of any automated system is common in AI studies,<sup>57</sup> and can reflect issues with the underlying datasets or model development techniques. Numerous examples of biases have emerged, including those related to sex,<sup>58</sup> race,<sup>59</sup> and geography,<sup>60</sup> which could inadvertently exacerbate underlying healthcare inequalities.<sup>61</sup>

There are increasing suggestions that AI can ‘perform better than humans’ in medical tasks. A recent systematic review found 10 RCTs testing DL versus clinicians (2 completed, 8 ongoing), only one of which was in medical imaging (breast ultrasound, ongoing).<sup>6</sup> In contrast, they found 81 non-randomised comparisons (nine prospective, six relevant clinical environment). The AI was said to be better than or comparable to the clinician in 47% of the 81 studies. However, development and testing often used the same dataset, had small numbers of human comparators (e.g. five), most studies had high risk of bias, and few adhered to reporting standards.<sup>6</sup>

Another comparison of the diagnostic accuracy of DL versus clinicians identified 82 studies in which it was possible to calculate accuracy in 69. Mean sensitivity was 79.1% (range 9.7–100%), specificity was 88.3% (range 38.9–100%),<sup>29</sup> but many studies did not compare the DL and clinicians *on the same data*, did not externally validate their results, or report their methods adequately. Amongst the 14 studies with external validation that tested DL and clinicians *on the same sample*, the pooled sensitivity was 87.0% (95% CI 83.0–90.2%) versus 86.4% (79.9–91.0%), and pooled specificity was 92.5% (85.1–96.4%) versus 90.5% (80.6–95.7%) for DL versus clinicians respectively.<sup>29</sup> Of note, there were no studies of AI in stroke amongst the 82 studies, and only two studies concerned brain imaging (one MRI in dementia, one CT in head trauma).

#### **6. What evidence is there that AI technologies will meet the needs of users, including community practitioner, neurologist, radiologist, and the patient?**

Do AI tools reduce ‘door to needle’ time? Or ~~do might~~ AI tools ~~increase-worsen ‘analysis-paralysis’~~, treatment delays, or deny some patients effective treatments?<sup>29</sup> There are no completed prospective randomised trials of the impact of stroke AI tools on workflows or clinical outcomes, only before-after evaluations (see 4iv above).<sup>41,50,51</sup> It Unfortunately, it is common in hospitals for new digital systems to slow, not accelerate.

workflows.<sup>62</sup> ~~it should not be assumed that AI imaging tools will accelerate treatment.~~ Attention of AI should focus on improving routine medical workflows including image management and electronic case records to reduce time wasted, improve information content and diagnostic utility.<sup>2,63</sup> Some AI tools perform tasks that are not that helpful,<sup>64</sup> provide clinically irrelevant measures, or operate slower than a seasoned user of existing medical computing systems, delaying uptake of AI into general radiology.<sup>5</sup>

Different algorithms may give different results. Comparison of three AI tools for ASPECTS scoring<sup>7,21</sup> showed the highest correlation between the expert read and Brainomix (ICC=0.871 (0.818, 0.909),  $p < 0.001$ ) but comparable area-under-curve between the AI applications and expert consensus (Brainomix: AUC 0.759 (0.670–0.848),  $p < 0.001$ ; Frontier V2: AUC 0.752 (0.660–0.843),  $p < 0.001$ ; RAPID: AUC 0.734 (0.634–0.831),  $p < 0.001$ ). AI software may help less experienced doctors interpret acute stroke CT scans, e.g. use of Syngo.via Frontier ASPECT Score Prototype improved the correlation between junior radiologists and the reference standard from good ( $r=0.680$ ) to excellent ( $r =0.852$ ) in one small study.<sup>65</sup>

Apparently good performance at group level may mask important variation at the individual level. Amongst 12 ML and seven statistical models to predict cardiovascular disease risk using data from 3.6 million patients, the models had similar population level performance (C statistics of about 0.87), but varied widely in their prediction of individual risks particularly at higher risks,<sup>49</sup> and compared with the risk predicted by a reference model, about 60% of patients would have been classed as lower risk by another model.<sup>49</sup>

AI tool evaluation is usually restricted to single technical measures in controlled settings that only indirectly relate to the intended clinical tool use. However, the clinical need is rarely 'raw classification' but rather diagnostic or therapeutic decision support.<sup>66</sup> DL tools are notoriously sensitive to changes in input characteristics, and not customarily stress-tested across different scanners, parameters, clinics, patient groups, pre-processing tools, etc., further hindered by the black-box nature of DL methods and commercial confidentiality.<sup>67</sup> The benefits and challenges of introducing computational innovations into existing clinical ecosystems<sup>62</sup> remains sparsely assessed; AI software may require specific tailoring to suit different settings and institutions.<sup>68</sup>

What about patients? There is little participant involvement in development of AI software<sup>69</sup> despite concerns.<sup>70</sup> Few people want to receive a terminal diagnosis from a robot,<sup>71</sup> or want major treatment decisions for a life-threatening disease (like stroke) to be based primarily on AI-determined green, yellow or red areas on a scan characteristics, particularly when these differ between AI software.<sup>7,72</sup> Increasing availability of AI diagnostic tools and their potential including front line use by less experienced doctors, risks subverting medical judgement through inflexible application of easy-to-derive threshold values - e.g. treat if  $\leq 70$ ml core, not if  $>70$  ml. What about 75 ml, or 85 ml? Treatment decisions with the very powerful reperfusion therapies now available for stroke must *consider the whole patient* and not place inappropriate weight on perfusion maps that only represent brief

snapshots in time of a highly dynamic disorder, particularly when use of a different software is very likely to give different results.<sup>7</sup>

### **7. What could improve translation of AI technologies into benefits for patients and health care systems?**

There is a gulf between AI software developers and the intended clinical uses of the software, and a need for consensus-based, interdisciplinary and comprehensive scoping of user needs and constraints to guide effective development of AI tools.<sup>24,25,63</sup> The accuracy of an AI tool depends on its data input. AI developers highlight the need for ‘ground truth’ for training the software, meaning images where relevant features have been demarcated, often by hand, and in large datasets.<sup>54,55</sup> While large collections of medical imaging data are increasingly available, manual annotation is a massive, time consuming task and must be done by humans, therefore not many large annotated datasets exist. Some companies are outsourcing the work to low paid workers in low income countries.<sup>22</sup> The market for data-labelling services may triple to US\$5billion by 2023, with companies like ‘Mechanical Turk’ (owned by Amazon) providing freelancers ready to perform ‘micro-tasks’ like tagging images, or ‘Hive’ which runs online data labelling ‘games’ where operators earn money for labelling features,<sup>22</sup> which questions the reliability of the ground truth thereby derived.

Can we accelerate reliable, representative and diverse dataset availability, and, is ‘ground truth’ really essential, or could correlated variables like clinical outcomes be used instead? AI tools could be even more valuable if they could discover novel features or markers of severity, or predict clinically-relevant outcomes and treatment response, and thus improve clinical management. Table 3 lists important factors, including more accessible large-scale data, standardization of pre-processing pipelines, sharing open source codes, adoption of guidelines for reporting of AI development, closer working between ML/DL developers and clinicians, and better standards for evaluating AI tools against relevant clinical outcomes, control of confounders, correlates and colliders that impede AI performance

### **Discussion and Conclusions**

While AI tools hold great promise in stroke, a very large amount of much more work is required to demonstrate their clinical value and cost-effectiveness to patients, doctors, and health-care providers. AI development requires more focus on multidisciplinary teams including AI experts, IT experts, radiologists and strokologists<sup>63</sup> that listening to each other carefully, to gain from the undoubted huge potential of AI. Currently, without a more cohesive multidisciplinary and less commercially motivated effort, it seems that the stroke-AI is at risk of *Verschlimmbesserung*, meaning an ‘attempted advance without improvement, or even with worsening’. We should not treat green, yellow and magenta regions on scans, but the specific pathology within individual

Formatted: Font: Italic

patients. Imaging and AI image analysis should  ~~demonstrate clearly through objective, relevant and reliable evidence that improve help~~ patient management and  ~~improve~~ clinical outcomes,  ~~not hinder care or worsen outcomes. Tools should work for, not control, us.~~

'Digital' methods complement but cannot replace the *human touch* in medicine.<sup>70</sup> Patients should be actively  
'Digital' methods complement but cannot replace the *human touch* in medicine.<sup>70</sup> Patients should be actively  
'Digital' methods complement but cannot replace the *human touch* in medicine.<sup>70</sup> Patients should be actively  
'Digital' methods complement but cannot replace the *human touch* in medicine.<sup>70</sup> Patients should be actively  
'Digital' methods complement but cannot replace the *human touch* in medicine.<sup>70</sup> Patients should be actively  
involved in the design and evaluation of AI tools usage, since  ~~patient groups are clearly very wary of nobody~~  
 ~~wants to have too much unregulated use of computers in clinical decision-making their fate decided by a robot~~  
 ~~using an arbitrary, variable threshold.~~<sup>71,72</sup>

Not all measurement has value. Rather, consider what features would be treatment relevant to detect  
and quantify. Imaging findings in stroke patients help us chose the most effective treatment.  
Spontaneous brain hemorrhage requires angiography to find and treat a vascular malformation or aneurysm.  
Exclusion of hemorrhage but thrombus within the MCA requires immediate thrombectomy. Stroke imaging is  
complex, not a 'single feature' process, and perhaps a more difficult place to initiate AI tool development than it  
might seem superficially. Low ASPECTS is an independent predictor of poor outcome, but patients may still  
benefit from treatment, therefore reducing information to a single binary variable such as  
ASPECTS score would seem to be a retrograde step.

Can costs of AI tools be realistic? Currently one typical commercial AI software costs around  
US\$47,868 for one hospital for one year in the UK, equivalent to about a third of a hospital consultant's salary,  
and seems an unreasonable amount of money for something which only identifies a few features in one  
disease, should only be used by an experienced medic, and thus does not replace anything or  
anyone, and has a limited evidence base.

The essential next steps are first, to be aware of the limitations where commercial AI tools are in use, and second,  
to obtain provide reliable evidence of benefits versus harms of imaging AI tools' performance. This would best be  
tested in large scale randomised trials, to minimize bias, and in the clinical settings in which they will be used to  
ensure applicability to real world clinical practice. It is no longer appropriate to show only diagnostic accuracy on  
selected retrospective datasets, without reporting the failures. We need to see how AI tools work in clinical

practice, how they integrate into patient care, and if and how much they are beneficial, through providing evidence that they ~~change~~ing management for the better, improving outcomes, and ~~being~~are cost effective.

## References

1. Chen D, Liu S, Kingsbury P, Sohn S, Storlie CB, Habermann EB, et al. Deep learning and alternative learning strategies for retrospective real-world clinical data. *NPJ Digit Med*. 2019;2:43
2. Schwalbe N, Wahl B. Artificial intelligence and the future of global health. *Lancet*. 2020;395:1579-1586
3. van Leeuwen KG, Schalekamp S, Rutten MJCM, van Ginneken B, de Rooij M. Artificial intelligence in radiology: 100 commercially available products and their scientific evidence. *European Radiology*. 2021;31:3797-3804
4. Yao AD, Cheng DL, Pan I, Kitamura F. Deep learning in neuroradiology: A systematic review of current algorithms and approaches for the new wave of imaging technology. *Radiology: Artificial Intelligence*. 2020;2:e190026
5. Alexander A, Jiang A, Ferreira C, Zurkiya D. An intelligent future for medical imaging: A market outlook on artificial intelligence for medical imaging. *J Am Coll Radiol*. 2020;17:165-170
6. Nagendran M, Chen Y, Lovejoy CA, Gordon AC, Komorowski M, Harvey H, et al. Artificial intelligence versus clinicians: Systematic review of design, reporting standards, and claims of deep learning studies. *BMJ*. 2020;368:m689
7. Soun JE, Chow DS, Nagamine M, Takhtawala RS, Filippi CG, Yu W, et al. Artificial intelligence and acute stroke imaging. *Am J Neuroradiol*. 2021;42:2-11
8. Murray NM, Unberath M, Hager GD, Hui FK. Artificial intelligence to diagnose ischemic stroke and identify large vessel occlusions: A systematic review. *J Neurointerv Surg*. 2020;12:156-164
9. Mikhail, Le MGD, Mair G. Computational image analysis of nonenhanced computed tomography for acute ischemic stroke: A systematic review. *J Stroke Cerebrovasc Dis*. 2020;29:104715
10. Zhu G, Jiang B, Chen H, Tong E, Xie Y, Faizy TD, et al. Artificial intelligence and stroke imaging: A west coast perspective. *Neuroimaging Clin N Am*. 2020;30:479-492
11. Yeo M, Tahayori B, Kok HK, Maingard J, Kutaiba N, Russell J, et al. Review of deep learning algorithms for the automatic detection of intracranial hemorrhages on computed tomography head imaging. *J Neurointerv Surg*. 2021;13:369-378
12. Bivard A, Churilov L, Parsons M. Artificial intelligence for decision support in acute stroke — current roles and potential. *Nature Reviews Neurology*. 2020;16:575-585
13. Na DG, Kim EY, Ryoo JW, Lee KH, Roh HG, Kim SS, et al. Ct sign of brain swelling without concomitant parenchymal hypoattenuation: Comparison with diffusion- and perfusion-weighted mr imaging. *Radiology*. 2005;235:992-998
14. Mouridsen K, Thurner P, Zaharchuk G. Artificial intelligence applications in stroke. *Stroke*. 2020;51:2573-2579
15. Qiu W, Kuang H, Teleg E, Ospel JM, Sohn SI, Almekhlafi M, et al. Machine learning for detecting early infarction in acute stroke with non-contrast-enhanced ct. *Radiology*. 2020;294:638-644
16. Garcia JH, Yoshida Y, Chen H, Li Y, Zhang ZG, Lian J, et al. Progression from ischemic injury to infarct following middle cerebral artery occlusion in the rat. *Am J Pathol*. 1993;142:623-635

17. Wu S, Mair G, Cohen G, Morris Z, von Heijne A, Bradey N, et al. Hyperdense artery sign, symptomatic infarct swelling and effect of alteplase in acute ischemic stroke. . *Stroke & Vascular Neurology*. 2020;0
18. The IST-3 Collaborative Group. Association between brain imaging signs, early and late outcomes, and response to intravenous alteplase after acute ischemic stroke in the third international stroke trial (ist-3): Secondary analysis of a randomised controlled trial. *Lancet Neurol*. 2015;14:485-496
19. Georgakis MK, Duering M, Wardlaw JM, Dichgans M. WMH and long-term outcomes in ischemic stroke: A systematic review and meta-analysis. *Neurology*. 2019;92:e1298-e1308
20. Suss RA. ASPECTS, the mismeasure of stroke: A metrological investigation. *OSF Preprints*. 2019;31st Dec
21. Hoelter P, Muehlen I, Goelitz P, Beuscher V, Schwab S, Doerfler A. Automated aspect scoring in acute ischemic stroke: Comparison of three software tools. *Neuroradiology*. 2020;62:1231-1238
22. Editorial. Artificial intelligence: Human-machine interface. *The Economist*. 2019;19th October:70
23. Luo W, Phung D, Tran T, Gupta S, Rana S, Karmakar C, et al. Guidelines for developing and reporting machine learning predictive models in biomedical research: A multidisciplinary view. *J Med Internet Res*. 2016;18:e323
24. Rivera SC, Liu X, Chan AW, Denniston AK, Calvert MJ. Guidelines for clinical trial protocols for interventions involving artificial intelligence: The spirit-ai extension. *BMJ* 2020;370:m3210
25. Liu X, Rivera SC, Moher D, Calvert MJ, Denniston AK. Reporting guidelines for clinical trial reports for interventions involving artificial intelligence: The consort-ai extension. *BMJ*. 2020;370:m3164
26. Mongan J, Moy L, Kahn CEJ. Checklist for artificial intelligence in medical imaging (claim): A guide for authors and reviewers. *Radiology: Artificial Intelligence*. 2020;2:e200029
27. Weikert T, Francone M, Abbara S, Baessler B, Choi BW, Gutberlet M, et al. Machine learning in cardiovascular radiology: Escr position statement on design requirements, quality assessment, current applications, opportunities, and challenges. *European Radiology*. 2021;31:3909-3922
28. Yusuf M, Atal I, Li J, Smith P, Ravaud P, Fergie M, et al. Reporting quality of studies using machine learning models for medical diagnosis: A systematic review. *BMJ Open*. 2020;10:e034568
29. Liu X, Faes L, Kale AU, Wagner SK, Fu DJ, Bruynseels A, et al. A comparison of deep learning performance against health-care professionals in detecting diseases from medical imaging: A systematic review and meta-analysis. *Lancet Digital Health* 2019;1:e271-e297
30. Kundeti SR, Vaidyanathan MK, Shivashankar B, Gorthi SP. Systematic review protocol to assess artificial intelligence diagnostic accuracy performance in detecting acute ischemic stroke and large-vessel occlusions on ct and mr medical imaging. *BMJ Open*. 2021;11:e043665
31. Mair G, Chappell F, Martin C, Dye D, Bath PM, Muir KW, et al. Real-world independent testing of e-aspects software (rites): Statistical analysis plan. *AMRC Open Research*. 2020;2
32. Guberina N, Dietrich U, Radbruch A, Goebel J, Deuschl C, Ringelstein A, et al. Detection of early infarction signs with machine learning-based diagnosis by means of the alberta stroke program early ct score (aspects) in the clinical routine. *Neuroradiology*. 2018;60:889-901



33. Bulwa Z, Dasenbrock H, Osteraas N, Cherian L, Crowley RW, Chen M. Incidence of unreliable automated computed tomography perfusion maps. *J Stroke Cerebrovasc Dis.* 2019;28:104471
34. Purrucker JC, Mattern N, Herweh C, Möhlenbruch M, Ringleb PA, Nagel S, et al. Electronic alberta stroke program early ct score change and functional outcome in a drip-and-ship stroke service. *J Neurointerv Surg.* 2020;12:252-255
35. Potreck A, Seker F, Mutke MA, Weyland CS, Herweh C, Heiland S, et al. What is the impact of head movement on automated ct perfusion mismatch evaluation in acute ischemic stroke? *J Neurointerv Surg.* 2021
36. Neuberger U, Nagel S, Pfaff J, Ringleb PA, Herweh C, Bendszus M, et al. Impact of slice thickness on clinical utility of automated alberta stroke program early computed tomography scores. *Eur Radiol.* 2020;30:3137-3145
37. Kral J, Cabal M, Kasickova L, Havelka J, Jonszta T, Volny O, et al. Machine learning volumetry of ischemic brain lesions on ct after thrombectomy-prospective diagnostic accuracy study in ischemic stroke patients. *Neuroradiology.* 2020;62:1239-1245
38. Yahav-Dovrat A, Saban M, Merhav G, Lankri I, Abergel E, Eran A, et al. Evaluation of artificial intelligence-powered identification of large-vessel occlusions in a comprehensive stroke center. *AJNR* 2021;42:247-254
39. Ferreti LA, Leitao CA, Teixeira BCA, Lopes Neto FDN, ZÉtola VF, Lange MC. The use of e-ASPECTS in acute stroke care: Validation of method performance compared to the performance of specialists. *Arq Neuropsiquiatr.* 2020;78:757-761
40. Nagel S, Sinha D, Day D, Reith W, Chapot R, Papanagiotou P, et al. E-aspects software is non-inferior to neuroradiologists in applying the ASPECTS score to computed tomography scans of acute ischemic stroke patients. *Int J Stroke.* 2017;12:615-622
41. Adhya J, Li C, Eisenmenger L, Cerejo R, Tayal A, Goldberg M, et al. Positive predictive value and stroke workflow outcomes using automated vessel density (RAPID-CTA) in stroke patients: One year experience. *Neuroradiol J.* 2021:19714009211012353
42. Seker F, Pfaff JAR, Mokli Y, Berberich A, Namias R, Gerry S, et al. Diagnostic accuracy of automated occlusion detection in CT angiography using e-cta. *Int J Stroke.* 2021:1747493021992592
43. Dehkharghani S, Lansberg M, Venkatsubramanian C, Cereda C, Lima F, Coelho H, et al. High-performance automated anterior circulation ct angiographic clot detection in acute stroke: A multireader comparison. *Radiology.* 2021;298:665-670
44. Amukotuwa SA, Straka M, Smith H, Chandra RV, Dehkharghani S, Fischbein NJ, et al. Automated detection of intracranial large vessel occlusions on computed tomography angiography: A single center experience. *Stroke.* 2019;50:2790-2798
45. Golan D, Shalitin O, Sudry N, J. M. Ai-powered stroke triage system performance in the wild. *J Exper Stroke Trans Med.* 2020;12:01-04
46. Heit JJ, Coelho H, Lima FO, Granja M, Aghaebrahim A, Hanel R, et al. Automated cerebral hemorrhage detection using rapid. *AJNR.* 2021;42:273-278

47. Straka M, Albers GW, Bammer R. Real-time diffusion-perfusion mismatch analysis in acute stroke. *J Magn Reson Imaging*. 2010;32:1024-1037
48. Aggarwal R, Sounderajah V, Martin G, Ting DSW, Karthikesalingam A, King D, et al. Diagnostic accuracy of deep learning in medical imaging: A systematic review and meta-analysis. *NPJ Digit Med*. 2021;4:65
49. Li Y, Sperrin M, Ashcroft DM, van Staa TP. Consistency of variety of machine learning and statistical models in predicting clinical risks of individual patients: Longitudinal cohort study using cardiovascular disease as exemplar. *BMJ*. 2020;371:m3919
50. Al-Kawaz M, Primiani C, Urrutia V, Hui F. Impact of rapidai mobile application on treatment times in patients with large vessel occlusion. *J NeuroInterven Surg*. 2021:neurintsurg-2021-017365
51. Morey JR, Zhang X, Yaeger KA, Fiano E, Marayati NF, Kellner CP, et al. Real-world experience with artificial intelligence-based triage in transferred large vessel occlusion stroke patients. *Cerebrovasc Dis*. 2021;50:450-455
52. Domingues I, Pereira G, Martins P, Duarte H, Santos J, Abreu PH. Using deep learning techniques in medical imaging: A systematic review of applications on ct and pet. *Artificial Intelligence Review*. 2020;53:4093-4160
53. Roberts M, Driggs D, Thorpe M, Gilbey J, Yeung M, Ursprung S, et al. Common pitfalls and recommendations for using machine learning to detect and prognosticate for covid-19 using chest radiographs and ct scans. *Nature Machine Intelligence*. 2021;3:199-217
54. Flanders A, Prevedello L, Shih G, Halabi SS, Kalpathy-Cramer J, Ball R, et al. Construction of a machine learning dataset through collaboration: The rsna 2019 brain ct hemorrhage challenge. *Radiology: Artificial Intelligence*. 2020;2:e190211
55. Prevedello LM, Halabi SS, Shih G, Wu CC, Kohli MD, Chokshi FH, et al. Challenges related to artificial intelligence research in medical imaging and the importance of image analysis competitions. *Radiology: Artificial Intelligence*. 2019;1:e180031
56. Badgeley MA, Zech JR, Oakden-Rayner L, Glicksberg BS, Liu M, Gale W, et al. Deep learning predicts hip fracture using confounding patient and healthcare variables. *NPJ Digit Med*. 2019;2:31
57. Andaur Navarro CL, Damen JAA, Takada T, Nijman SWJ, Dhiman P, Ma J, et al. Risk of bias in studies on prediction models developed using supervised machine learning techniques: Systematic review. *BMJ*. 2021;375:n2281
58. Larrazabal AJ, Nieto N, Peterson V, Milone DH, Ferrante E. Gender imbalance in medical imaging datasets produces biased classifiers for computer-aided diagnosis. *PNAS*. 2020;117:12592-12594
59. Obermeyer Z, Powers B, Vogeli C, Mullainathan S. Dissecting racial bias in an algorithm used to manage the health of populations. *Science*. 2019;366:447-453
60. Kaushal A, Altman R, Langlotz C. Geographic distribution of us cohorts used to train deep learning algorithms. *JAMA*. 2020;324:1212-1213
61. Leslie D, Mazumder A, Peppin A, Wolters MK, Hagerty A. Does “ai” stand for augmenting inequality in the era of covid-19 healthcare? *BMJ*. 2021;372:n304

62. Overhage JM, McCallie Jr D. Physician time spent using the electronic health record during outpatient encounters: A descriptive study. *Ann Int Med.* 2020;172:169-174
63. Gilbert F, Smye S, Schönlieb C-B. Artificial intelligence in clinical imaging: A health system approach. *Clinical Radiology.* 2020;75:3-6
64. Koçak B, Durmaz E, Ateş E, Kılıçkesmez Ö. Radiomics with artificial intelligence: A practical guide for beginners. *Diagn Interv Radiol.* 2019;25:485-495
65. Li L, Chen Y, Bao Y, Jia X, Wang Y, Zuo T, et al. Comparison of the performance between frontier aspects software and different levels of radiologists on assessing ct examinations of acute ischemic stroke patients. *Clin Radiol.* 2020;75:358-365
66. Montagnon E, Cerny M, Cadrin-Chênevert A, Hamilton V, Derennes T, Ilinca A, et al. Deep learning workflow in radiology: A primer. *Insights Imaging.* 2020;11:22
67. Baselli G, Codari M, Sardanelli F. Opening the black box of machine learning in radiology: Can the proximity of annotated cases be a way? *Eur Radiol Exp.* 2020;4:30
68. Hickman SE, Baxter GC, Gilbert FJ. Adoption of artificial intelligence in breast imaging: Evaluation, ethical constraints and limitations. *Br J Cancer.* 2021;125:15-22
69. Gibney E. The battle for ethical ai at the world's biggest machine-learning conference. *Nature.* 2020;577:609-610
70. Richardson JP, Smith C, Curtis S, Watson S, Zhu X, Barry B, et al. Patient apprehensions about the use of artificial intelligence in healthcare. *NPJ Digit Med.* 2021;4:140
71. Mittelman M, Markham S, Taylor M. Patient commentary: Stop hyping artificial intelligence—patients will always need human doctors. *BMJ.* 2018;363:k4669
72. Longoni C, Bonezzi A, Morewedge CK. Resistance to medical artificial intelligence. *J Consum Res.* 2019;46:629-650
73. Winzeck S, Hakim A, McKinley R, Pinto J, Alves V, Silva C, et al. Isles 2016 and 2017-benchmarking ischemic stroke lesion outcome prediction based on multispectral mri. *Front Neurol.* 2018;9:679
74. Muschelli J. Recommendations for processing head ct data. *Front Neuroinform.* 2019;13:61
75. Maier-Hein L, Eisenmann M, Reinke A, Onogur S, Stankovic M, Scholz P, et al. Why rankings of biomedical image analysis competitions should be interpreted with care. *Nat Commun.* 2018;9:5217
76. Kelly CJ, Karthikesalingam A, Suleyman M, Corrado G, King D. Key challenges for delivering clinical impact with artificial intelligence. *BMC Med.* 2019;17:195

Table 1. Current commercially available AI-based software for stroke.

Company (Country)	Product	Automated functions (imaging modality)	CE Approval (Class) for EU Marketing	FDA Approval (Class) for US Marketing
Aidoc (Israel)	Briefcase	Detect hemorrhage (CT)	Yes (I)	Yes (II)
	Aidoc LVO	Detect LVO (CTA)	Yes (I)	Yes (II)
Avicenna.AI (France)	Cina-ASPECTS	Provide ASPECTS (CT)	Yes (I)	No
	Cina-ICH	Detect hemorrhage (CT)	Yes (I)	Yes (II)
	Cina-LVO	Detect LVO (CTA)	Yes (I)	Yes (II)
Brainomix (UK)	e-ASPECTS	Provide ASPECTS, detect dense MCA, detect hemorrhage (CT)	Yes (IIa)	No
	e-CTA	Detect LVO, provide collateral scoring (CTA)	Yes (IIa)	Yes (II)
	e-CTP	Process perfusion data (CTP)	Yes (IIa)	No
Cercare Medical (Denmark)	Cercare Stroke	Process perfusion data (CTP, MRP), detect ischemic lesions (MRI)	Yes (IIa)	No
Circle Neurovascular Imaging (Canada)	StrokeSENS	Provide ASPECTS (CT), provide collateral assessment, detect LVO (CTA)	Yes (II), does not include collaterals	Yes (I)*
Deep01 (Taiwan)	DeepCT	Detect hemorrhage (CT)	Yes (I)	Yes (II)
General Electric (USA)	Stroke VCAR	Detect hemorrhage (CT)	Yes (not declared)	Yes (II)
	CT Perfusion 4D	Process perfusion data (CTP)	Yes (not declared)	Yes (II)
Icometrix (Belgium)	Icobrain CVA	Process perfusion data (CTP)	Yes (I)	Yes (II)
InferVision Med Tech (China)	InferRead CT Stroke.AI	Provide ASPECTS, detect hemorrhage (CT)	Yes (IIa)	Yes (II)
JLK Inc. (South Korea)	JBS-01K	Detect ischemic lesions (MRI), provide stroke classification	Yes (I)	No
	JBS-04K	Detect hemorrhage (CT)	Yes (I)	No
Keya Medical (China)	CuraRad-ICH	Detect hemorrhage (CT)	No	Yes (II)

MaxQ.ai (Israel)	Accipio IX	Detect hemorrhage (CT)	Yes (not declared)	Yes (II)
NICO.Lab (Netherlands)	StrokeViewer	Provide ASPECTS, detect hemorrhage (CT), LVO detection, provide collateral assessment (CTA), process perfusion data (CTP)	Yes (I)	Yes (II)
<u>Olea Medical (France)</u>	<u>Olea Sphere</u>	<u>Process perfusion data (CTP)</u>	<u>Yes (IIa)</u>	<u>Yes (II)</u>
Qure.ai (India)	qER	Detect hemorrhage (CT)	Yes (IIa)	Yes (II)
RapidAI (USA)	Rapid	Process perfusion data and provide flow dynamics (CTA, CTP), detect DWI lesions (MRI)	Yes (I)	Yes (II)
	Rapid ASPECTS	Provide ASPECTS (CT)	Yes (I)	Yes (II)
	Rapid ICH	Detect hemorrhage (CT)	Yes (I)	Yes (II)
	Rapid LVO	Detect LVO (CTA)	Yes (I)	Yes (II)
Siemens (Germany)	syngo.CT ASPECTS	Provide ASPECTS (CT)	Yes (not declared)	No
	syngo.CT Neuro Perfusion	Process perfusion data (CTP)	Yes (not declared)	Yes (II)
Viz.ai (USA)	Viz ICH	Detect hemorrhage (CT)	Yes (not declared)	Yes (II)
	Viz LVO (ContaCT)	Detect LVO (CTA)	Yes (not declared)	Yes (II)
	Viz CTP	Process perfusion data (CTP)	Yes (not declared)	Yes (II)
Zebra Medical Vision (Israel)	HealthICH	Detect hemorrhage (CT)	Yes (not declared)	Yes (II)

**Notes:** Both CE and FDA classification use 3 classes depending on risk to the patient and/or user according to the *intended* use of the device: I=low risk, II=medium risk, III=high risk. Both classification systems have a more stringent process for classifying higher risk devices. For CE, I can be self-certified by manufacturer, II&III require audit of validation results by a notified body. For FDA, I&II require 501k (prove equivalence to device already approved for marketing, or de novo assessment if novel) while III requires premarket approval (PMA) including software validation results.

\* Approved for data transfer only, not automated processing.

Details extracted from publicly available [EU](#) and [FDA](#) data and vendor websites, correct to 31<sup>st</sup> Aug 2021.

Commented [WJ1]: Now to 31 Jan 2022?

Table 2. Accuracy of three commercially available AI-based software in stroke.

Company	Software	Detection of ischemic brain injury		Detection of LVO	Detection of hemorrhage
		CT	MRI/CTP		
Brainomix	e-ASPECTS	Retrospective, <b>132</b> , 44-45%, 91-93%, follow-up CT <sup>40</sup> Retrospective*, <b>119</b> , 83%, 57%, follow-up CT <sup>32</sup> Retrospective*, <b>116</b> , 75%, 73%, experts with all data including follow-up CT or MRI. <sup>39</sup>			
	e-CTA			Retrospective, <b>301</b> , 84%, 96%, experts with all data including follow-up imaging. <sup>42</sup>	
	e-CTP				
Rapid	Rapid		Retrospective, <b>63</b> , 100%, 91%, experts <sup>47</sup>		
	Rapid ASPECTS				
	Rapid ICH				Retrospective, <b>308</b> , 96%, 95%, expert consensus <sup>46</sup>
	Rapid LVO			Retrospective*, <b>310</b> , 80%, NS, experts <sup>41</sup> Retrospective, <b>217</b> , 96%, 98%, experts <sup>43</sup> Retrospective, <b>477</b> , 92-94%, 97-98%, experts <sup>44</sup>	
Viz.ai	Viz ICH				
	Viz LVO (ContaCT)			Retrospective, <b>1167</b> , 81-82%, 90-96%, experts <sup>38</sup> Retrospective, <b>163</b> , 96%, 94%, NS <sup>45</sup>	
	Viz CTP				
<b>SUMMARY</b>		<b>367 patients, 1 software</b> <b>Sensitivity = 44-83%</b> <b>Specificity = 57-93%</b>		<b>2635 patients, 3 different software</b> <b>Sensitivity = 80-96%</b> <b>Specificity = 90-98%</b>	

**Note:** Results are [study design, n, sensitivity, specificity, reference standard], unless otherwise stated. \* Indicates study conducted independent of company. Blank boxes indicate no suitable papers were identified. Shaded boxes indicate no papers expected. NS = Not specified.

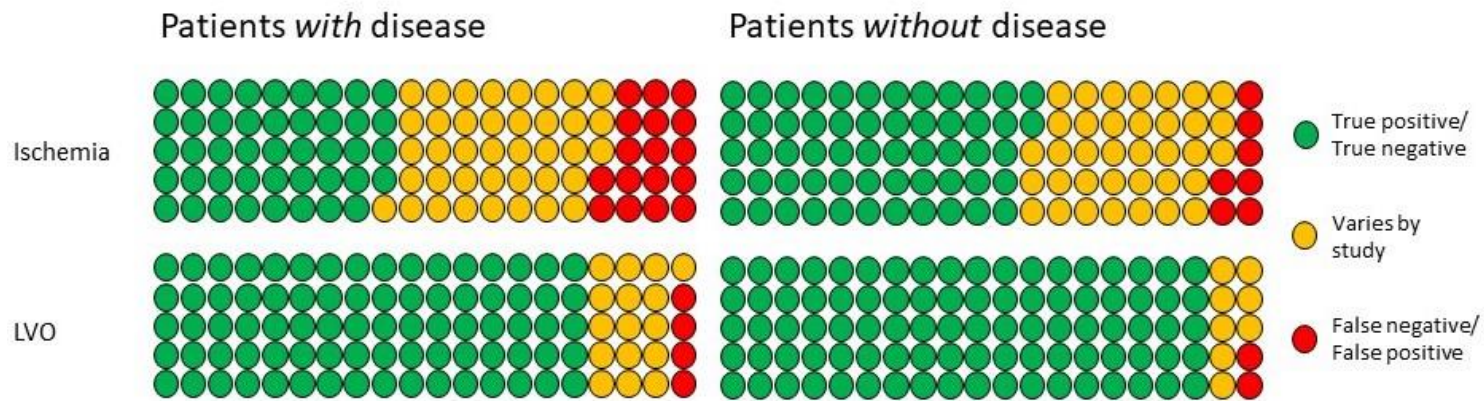
**Table 3.** Steps needed to accelerate AI software development

Need	Comment
More accessible large-scale data repositories:	RSNA (Radiology Society of North America) Head CT Challenge: <sup>54</sup> ICH detection (>800,000 scans); ASFNR (American Society of Functional Neuroradiology) Head CT Challenge: CT pathology at different ages ( <a href="https://aichallenge.asfnr.org/">https://aichallenge.asfnr.org/</a> ); ISLES (Ischemic Stroke Lesion Segmentation) Challenge (n=75); <sup>73</sup> large trials e.g. Third International Stroke Trial (IST-3) (3035 pts, >10,000 scans) UK Biobank ( <a href="http://ukbiobank.ac.uk">ukbiobank.ac.uk</a> ); HDR UK (Health Data Research)
Guidelines to standardize pre-processing pipelines <sup>74</sup>	Preprocessing steps: reading DICOM data, converting DICOM to other formats, brain extraction from skull, defacing, registration, etc RSNA 2019 ICH detection challenge <sup>54</sup> shows many common problems that can occur e.g., under-labelling data, human errors, imbalanced classes, inappropriate de-identification and anonymization across data sources.
Guidelines and standards for reporting of machine learning models <sup>6, 63</sup>	Transparent Reporting of a Multivariable Prediction Model for Individual Prognosis or Diagnosis - statement specific to Machine Learning (TRIPOD-ML), Standard Protocol Items: Recommendations for Interventional Trials (SPIRIT)-AI and Consolidated Standards of Reporting Trials (CONSORT)-AI. <sup>25</sup>
Collaboration between ML practitioners and all relevant disciplines	Physicians and radiologists benefit from becoming familiar with basic concepts in machine learning. <sup>63</sup> Machine Learning methods must be developed by those familiar with details and context of the real clinical settings including decision triage, practicalities of the specific patient population, typical medical imaging, image interpretation, definition of treatment relevant imaging finding, and data preparation.
Testing and validation protocols	Should move beyond technical performance evaluation (e.g., ROC, accuracy, sensitivity /specificity) to determine the clinical value of a system, given limited annotator availability, experience, accuracy and interest. The limits of testing criteria used in international challenges are evident. <sup>75</sup>
Attuned algorithmic development	Move beyond classifiers: methods must control for confounders, correlates and colliders that introduce bias and produce non-robust methods that collapse with potentially dangerous consequences when deployed in different real settings <sup>76</sup>

Open-source AI code and data	Proprietary code impedes replication, reproducibility, clinical validation: it is difficult and costly. <sup>63</sup> Commercial development based on ill-established methods that the community cannot verify in clinical settings risks reputational damage.
------------------------------	--



**Figure 1.** Potential clinical implications of AI software use for stroke feature detection per 100 patients assessed, derived from data in Table 2.



**Note:** Orange circles indicate the overlapping range of values provided by different studies.