

VASARI-auto: equitable, efficient, and economical featurisation of glioma MRI

James K. Ruffle^{1,2}, Samia Mohinta¹, Kelly Pegoretti Baruteau², Rebekah Rajiah¹, Faith Lee¹, Sebastian Brandner³, Parashkev Nachev¹ and Harpreet Hyare^{1,2}

¹Queen Square Institute of Neurology, University College London, London, UK

²Lysholm Department of Neuroradiology, National Hospital for Neurology and Neurosurgery, London, UK

³Division of Neuropathology and Department of Neurodegenerative Disease, Queen Square Institute of Neurology, University College London, London, UK

Correspondence to:

Dr James K Ruffle

Email: j.ruffle@ucl.ac.uk

Address: Institute of Neurology, UCL, London WC1N 3BG, UK

Correspondence may also be addressed to:

Dr Harpreet Hyare

Email: harpreet.hyare@nhs.net

Address: Institute of Neurology, UCL, London WC1N 3BG, UK

Funding

JKR was supported by the Medical Research Council (MR/X00046X/1) and the British Society of Neuroradiology. PN is supported by the Wellcome Trust (213038/Z/18/Z). The UCLH NIHR Biomedical Research Centre supports HH and PN.

Conflict of interest

None to declare.

Manuscript Type

Original article

Authorship

JKR: conceptualisation, methodology, software, validation, formal analysis, resources, writing – original draft, writing - review & editing, visualisation.

SM: methodology, formal analysis, writing - review & editing.

KPG: validation, formal analysis, writing - review & editing.

RR: methodology, writing – original draft, writing - review & editing.

FL: methodology, writing - review & editing.

SB: validation, resources, writing - review & editing.

PN: methodology, resources, writing – original draft, writing - review & editing.

HH: conceptualisation, methodology, validation, formal analysis, resources, writing – original draft, writing - review & editing.

All authors have been involved in the writing of the manuscript and have read and approved the final version.

Acknowledgements

We are most grateful for the support of the British Neuro-Oncology Society, which assisted us in accumulating a list of all UK neuro-oncology centres.

Keywords

Glioma, deep learning, artificial intelligence, VASARI, decision support, radiology, medical imaging.

Abstract

Background and Purpose: The VASARI (Visually AcceSAble Rembrandt Images) MRI feature set is a quantitative system designed to standardise glioma imaging descriptions. Though effective, deriving VASARI is time-consuming and seldom used in clinical practice. This is a problem that machine learning could plausibly automate.

Materials and Methods: Using glioma MRI data of 1172 patients, we developed VASARI-auto, an automated labelling software applied to open-source lesion masks and our openly available tumour segmentation model. In parallel, two consultant neuroradiologists independently quantified VASARI features in a randomly chosen subsample of 100 glioblastoma cases. We quantified: 1) agreement across neuroradiologists and VASARI-auto; 2) calibration of performance equity; 3) an economic workforce analysis; and 4) fidelity in predicting patient overall survival. Neuroradiologists were blinded to VASARI-auto software development and evaluations, and software developers were blinded to neuroradiologist labelling.

Results: Tumour segmentation was compatible with the current state of the art (mean segmentation Dice coefficient 0.95) and equally performant regardless of age or sex. A modest inter-rater variability between in-house neuroradiologists (Cohen's Kappa 0.49) was comparable to between neuroradiologists and VASARI-auto (Cohen's Kappa 0.41), with far higher agreement between VASARI-auto methods (Cohen's Kappa 0.94). The time for neuroradiologists to derive VASARI was substantially higher than VASARI-auto (mean time per case 317 vs 3 seconds, $p < 0.0001$). A UK hospital workforce analysis forecast that three years of VASARI featurisation would demand 29,777 consultant neuroradiologist workforce hours and >£1.5 million (\$1.9 million), reducible to 331.95 hours of computing time (and £146 of power) with VASARI-auto. The best-performing survival model utilised VASARI-auto features (R^2 0.25) instead of those derived by neuroradiologists (R^2 0.21).

Conclusion: VASARI-auto is a highly efficient automated labelling system with equitable performance across patient age or sex, a favourable economic profile if used as a decision support tool, and with non-inferior fidelity in downstream patient survival prediction. Future work should iterate upon and integrate such tools to enhance patient care.

Summary Statement

We develop and release VASARI-auto, a performant, equitable, efficient, economically favourable, and survival-predictive automated featurisation software for glioma MRI.

Key results

- Lesion segmentation performance was compatible with the current state of the art.
- VASARI featurisation demonstrated similar agreement between experienced consultant neuroradiologists but was derived considerably faster, with greater quantitative stability across patients and higher fidelity in downstream survival prediction.
- A UK-wide workforce analysis of all neuro-oncology centres forecast that this tool could save over £1.5 million (\$1.9 million) in workforce costs within three years of routine clinical care.

Introduction

The VASARI (Visually AcceSAbLe Rembrandt Images) MRI feature set is a quantitative scoring system designed to facilitate accurate and reliable imaging descriptions of adult gliomas(1), initially developed in 2010 as part of The Cancer Genome Atlas (TCGA) initiative from the Repository for Molecular BRAIn Neoplasia DaTA (REMBRANDT) study(2). It uses controlled and predefined terminology to define hallmark characteristics of glioma – including location, proportions of constituent components (such as oedema, enhancing and nonenhancing tumour), and other associated features such as cortical, ependymal, or deep white matter involvement.

VASARI's inception was intended to yield more consistent imaging interpretations, irrespective of its rater, centre or imaging approach(1). Indeed, it has shown promise towards better standardisation of adult glioblastoma, with multiple studies consistently demonstrating reasonable inter-observer agreement across constituent VASARI features beyond what could be expected from a conventional means of reporting(3-5). It has also been used with clinical and genomic data to effectively predict tumour histological grade, progression, mutation status, risk of recurrence and overall patient survival, implying a broader clinical utility(5-11). Though initially developed for adult glioblastoma, the VASARI feature set has been trialled in several novel clinical contexts, including the diagnosis of paediatric brain tumours(12) and rarer neuroepithelial malignancies(13), where it has shown potential as a clinical aid.

However, despite good evidence to support implementing the VASARI feature set as a clinical tool, it can be prohibitively time-consuming. Some studies report manual segmentation times of 20-40 minutes per case(14, 15). In an inevitably resource-limited and overstretched healthcare system(16), such a time constraint inevitably obstructs translation into real-world care.

Though the task is complex, it is theoretically deliverable by machine vision. Over the last few decades, lesion segmentation has formed a cornerstone of innovation across neuro-oncology(17-20), medical imaging(21, 22), biomedical engineering(23), machine and deep learning(24). The ability to segment an anatomical or pathological lesion in three dimensions confers the ability to evaluate it quantitatively – moving beyond visual qualitative assessment – with greater richness and fidelity than conventional two-dimensional measurements repeatedly shown to be often spurious and inconsistent between radiologists(25-27), and with greater sensitivity to the heterogeneity of the underlying pathological patterns(28). Enabling radiological image segmentation opens many possibilities for downstream innovation in neuro-oncological healthcare and research, ranging from clinical stratification, outcome prediction, response assessment, treatment allocation and risk quantification, many of which have already shown great promise. The underlying goal is to enhance the individual fidelity of data-driven decision-making, facilitating better patient-centred care(29-31), a remit especially warranted in neuro-oncology(32, 33).

Given this, we developed VASARI-auto, an automated VASARI feature set labelling tool (Figure 1). With a required input of patient lesion segmentations only – engineered by design to maximise patient confidentiality - we herein illustrate its high performance, efficiency, equitably, and downstream survival predictive utility in a multi-site patient cohort large-in-kind, and with real-world healthcare provider simulations illustrating tangible added value that can enhance clinical neuro-oncology workflows.

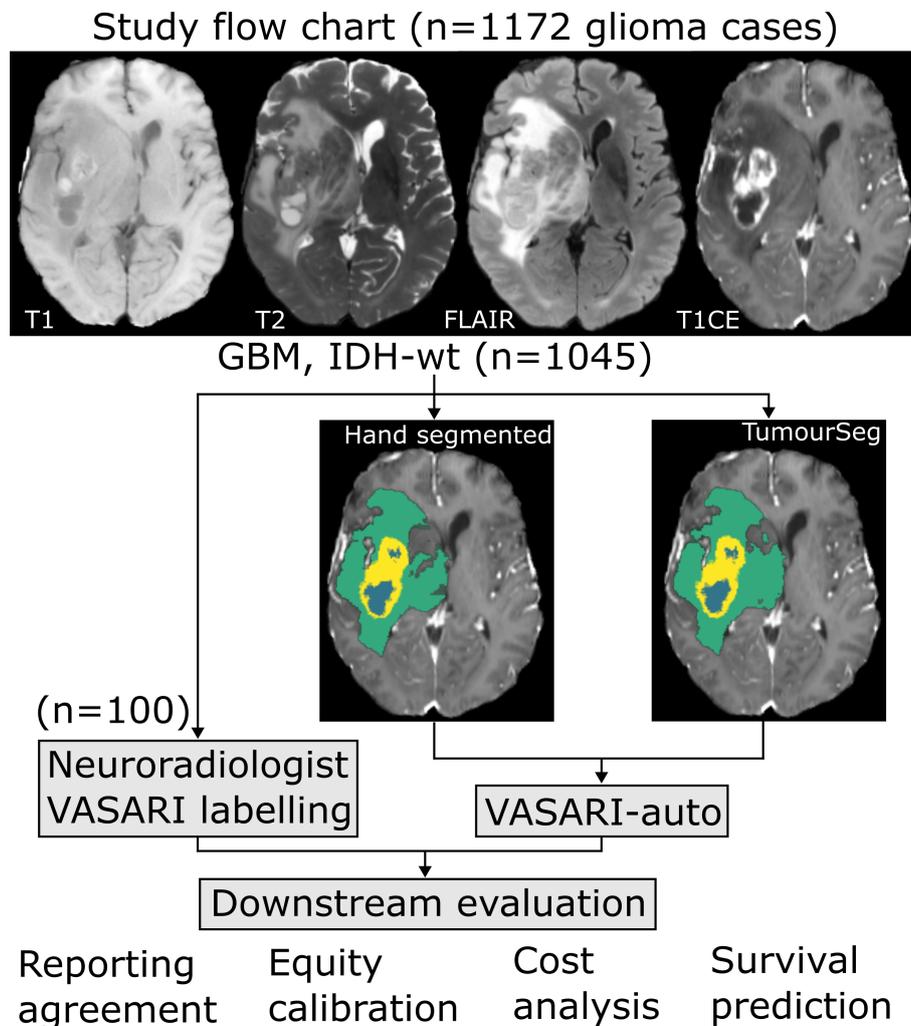


Figure 1. Study flow chart. Volumetric T1, T2, FLAIR, and post-contrast T1-weighted imaging were acquired for all participants. By random assignment of 100 glioblastoma, IDH-wt cases, two experienced consultant neuroradiologists reviewed structural imaging and recorded VASARI features and were timed doing so. In parallel, we developed VASARI-auto, an automated software to determine VASARI features. We derived these features using VASARI-auto from both semi-supervised hand-annotated lesion masks from a separate group of neuroradiologists, and using our previously published and openly available tumour segmentation model ‘TumourSeg’[20, 34]. We subsequently undertook multiple downstream evaluations of both neuroradiologist VASARI labelling and that from VASARI-auto, evaluating: 1) agreement both between neuroradiologists, between software, and between neuroradiologist and software; 2) equity calibration to determine if neuroradiologist and software labelling were equitably performant for all ages and sexes; 3) a simulated economic analysis determining the cost to undertake labelling with neuroradiologists or VASARI-auto based on real-world clinical workloads; and 4) in using these data to predict patient overall survival. Neuroradiologists were blinded to all software development and evaluations from VASARI-auto, and likewise, software developers were blinded to all neuroradiologist labelling until the final downstream evaluation stage.

Materials and Methods

Data

We utilised open-source neuro-oncology patient imaging data (n=1172) provided by The University of California San Francisco Preoperative Diffuse Glioma MRI (UCSF-PDGM) (n=501)[35] and the University of Pennsylvania Glioblastoma Imaging, Genomics, and Radiomics (UPenn-GBM) (n=671)[36] repositories. We firstly contacted the corresponding authors of both datasets to clarify which participant imaging was part of The Brain Tumour Segmentation Challenge (BraTS)[37] since we used BraTs to train our in-house tumour segmentation model (herein referred to as 'TumourSeg')[20, 34], and as such needed to prevent any possibility of an information leak. We excluded any such cases from the data pool. Further, we subsampled to study only patients with a confirmed molecular diagnosis of glioblastoma, IDH-wt[33], for which VASARI featurisation was initially developed. Each patient dataset included volumetric and brain-extracted T1, T2, FLAIR, and post-contrast T1-weighted MRI sequences (Table 1).

In all patients, age, sex, overall survival (in days) and a lesion segmentation were available. Separately undertaken by the original UCSF-PDGM and UPenn-GBM repository authors, each patient neuroimaging set first underwent automated segmentation using an ensemble model consisting of the prior top-scoring BraTS challenge algorithms, which was then manually corrected by a group of annotators with varying experience and approved by one of two neuroradiologists with more than 15 years of attending experience each[35, 36].

Ethical approval

UCSF-PDGM data collection followed relevant guidelines and regulations and was approved by the UCSF institutional review board with a waiver for consent[36]. For UPenn-GBM, collection, analysis, and release of the UPenn-GBM data was approved by the Institutional Review Board at the University of Pennsylvania Health System (UPHS), and informed consent was obtained from all participants[36].

Parameter	All patients	Glioblastoma, IDH-wt cohort	VASARI labelled by consultant neuroradiologists
Number of patients	1172	1045	100
Age (years) \pm SD	60 \pm 13.84	62 \pm 12.24	61 \pm 13.49
Sex	704 male, 468 female	627 male, 418 female	56 male, 44 female
Molecular diagnosis	Glioblastoma, IDH-wt (89.16%), Astrocytoma, IDH-mut (9.73%), Oligodendroglioma, IDH-mut and 1p/19q co-deleted (1.11%).	Glioblastoma, IDH-wt (100%)	Glioblastoma, IDH-wt (100%)
Overall survival (days)	503 \pm 453.90	441 \pm 370.29	436 \pm 462.21

Neuroimaging

Neuroradiologist VASARI-featurisation

From this glioblastoma, IDH-wt cohort, we drew a random sample of 100 unique patients. All were randomly assigned to one of two consultant neuroradiologists (HH, KPB) with more than 15 years of experience in neuro-oncology, who quantified VASARI features from structural neuroimaging. This is except for features that required diffusion or non-brain-extracted sequences. The time taken to derive VASARI features in each patient case was recorded. From a random number generator, we drew a random integer between 10-15 (which drew 13), for which we randomly allocated 13 duplicate cases to both neuroradiologists to ascertain inter-rater agreement. Neuroradiologists were blinded to all software and model development.

Tumour segmentation

We deployed our in-house tumour segmentation model (TumourSeg) to all cases, described in significant detail elsewhere(20), and made openly available online(34). In brief, this model is a high-resolution convolutional 3D U-Net implemented with nnU-Net(38), a pipeline with proven high performance in semantic segmentation across a range of micro and macroscopic tasks(38-40). Our model was trained on the BraTS training dataset of 1251 participants with 5-fold cross-validation, with additional evaluation on cases acquired at our local tertiary neuro-oncology centre, the National Hospital for Neurology and Neurosurgery(20).

Nonlinear registration with enantiomorphic normalisation

Having segmented lesions in native space, structural imaging and lesion segmentation masks were nonlinearly registered to 1mm MNI space with Statistical Parametric Mapping (SPM) using enantiomorphic correction(41, 42). The advantage of enantiomorphic correction is that the risks of registration errors secondary to a lesion are minimised by leveraging a given patient's normal structural neuroanatomy on the unaffected contralesional hemisphere(41). A neuroradiologist (JKR) manually reviewed all imaging data at multiple stages of the data pre-processing.

VASARI-auto software development

We developed a fully automated pipeline – 'VASARI-auto' – to derive VASARI features from lesion masks. Lesion masks could be of any source, whether manually traced, from our internal tumour segmentation model, or other lesion segmentation tools. VASARI-auto required data to be held in MNI registered space (prototyped in a 1mm³ volumetric resolution, but deployable in any). We pooled neuroanatomical atlases for all brain lobes, as well as the brainstem, insula, thalamus, corpus callosum, internal capsule, ventricles, and cortex, for the derivation of locational-based features. For each case, VASARI-auto loaded the multi-channel tumour segmentation (with separate labels for enhancing tumour, nonenhancing tumour, and perilesional signal change) and – following pre-existing VASARI reporting standards(1) –

derived the following: F1 - tumour location; F2 – side of tumour epicentre; F4 – enhancement quality; F5 – proportion enhancing; F6 – proportion of nonenhancing tumour; F7 – proportion of necrosis; F9 – multifocal/multicentric lesional status; F11 – the thickness of the enhancing margin; F14 – the proportion of oedema; F19 – ependymal invasion; F20 – cortical involvement; F21 – deep white matter invasion; F22 whether nonenhancing tumour crossed the midline; F23 – whether enhancing tumour crossed the midline; and F24 – the presence of satellite lesions.

Our code did not quantify a few VASARI features that require either non-brain extracted data (F25 – calvarial modelling) or the original structural neuroimaging data (F10 T1/FLAIR ratio; F12-13 – definition of the enhancing and nonenhancing margin, F18 – pial invasion, and F16 – haemorrhage), the reason being was that we wished to develop an automated tool immediately usable with irrevocably anonymised lesion segmentation data without the requirement for raw volumetric neuroimaging. We similarly did not quantify F17 - diffusion changes since DWI was not available for many cases in the external data, beyond our control. We also did not quantify the presence of F8 – cysts since most brain tumour segmentation models (including ours) rely on BraTS lesion labels of enhancing tumour, nonenhancing tumour, and perilesional signal change, but with no distinction for cysts. Therefore, we felt any attempts to model cyst presence would be liable to confabulation. We similarly did not quantify F3 – eloquence, for lack of appropriate brain masks to robustly model it; moreover, we did not wish to detract from a gold standard of a neurosurgeon’s electrical stimulation assessment for eloquent-sparing resections[43].

The requirements to run VASARI-auto are given below in the software subsection. We also recorded time to quantify VASARI features with VASARI-auto, both already pre-generated lesion masks, and when paired with TumourSeg[20].

Downstream evaluation

Reporting agreement

We compared agreement in all VASARI featurisation between 1) consultant neuroradiologists, 2) consultant neuroradiologists and VASARI-auto, and 3) between VASARI-auto when using either the source semi-supervised and neuroradiology-reviewed segmentations to VASARI-auto using TumourSeg[20, 34]. Agreement was quantified by Cohen’s Kappa[44], which furthermore was appropriately linearly weighted for non-Boolean VASARI features.

We also quantified the balanced accuracy in VASARI featurisation between consultant neuroradiologists (the ground truth) and VASARI-auto (the prediction), as well as the balanced accuracy between VASARI-auto using the source semi-supervised and neuroradiology-reviewed segmentations (the ground truth) and VASARI-auto using TumourSeg (the prediction)[20, 34].

Equity calibration

We quantified software and reporting patient equity(45, 46) for all analysis steps. For tumour segmentation, we compared model performance by the Dice coefficient across all lesional compartments (whole tumour, enhancing tumour, nonenhancing tumour, perilesional signal change) for male and female sex and for all decades of age included in the cohort (20-90 years). We similarly compared Cohen's Kappa agreement metrics across male and female sex and all decades of age.

Efficiency, economic and workforce analysis

We statistically compared the time required to record VASARI features between 1) consultant neuroradiologists, 2) VASARI-auto with tumour segmentations already supplied, and 3) VASARI-auto with TumourSeg.

Next, we undertook an economic and workforce analysis, simulating neuro-oncology workload across the UK, and informed by that within our centre (the National Hospital for Neurology and Neurosurgery). Every week, a neuro-oncology multidisciplinary team (MDT) meeting is held to discuss all referrals, ongoing cases, and management plans, which includes a neuroradiological review of all cases. We reviewed the last three years of neuro-oncology MDT lists (2020-2023) and quantified the minimum and maximum number of cases to be discussed each week, which was of range 30-75. We determined the minimum and maximum pay scales for consultants in the National Health Service workforce as of March 2024, which vary depending on years of service(47). We similarly quantified power consumption costs to run a reasonably powerful computer (1200 kilowatt Hour (kWh), based on UK energy tariffs as of March 2024(48). We curated a list of all UK neuro-oncology centres (n=40), kindly provided by the British Society of Neuro-Oncology, to simulate UK-wide neuro-oncology workloads.

Having derived this data, we simulated the next three years (2024-2027) of MDT clinical workload at each centre. A random number of MDT cases was simulated weekly using the previous minimum-maximum caseload through 2020-2023. We then simulated a random choice of neuroradiology consultants who would be allocated to present a given week's neuro-oncology MDT, with their salary drawn randomly from the NHS consultant pay scales. We then randomly simulated the time taken to quantify VASARI features across all cases, where time per case was drawn from a random uniform distribution informed by the time taken for neuroradiologists to quantify all 100 cases in our earlier analysis. From this, we quantified the workload and financial cost if each patient had undergone VASARI featurisation by a neuroradiologist. We similarly quantified the time and expense of power if VASARI-auto and VASARI-auto with TumourSeg had undertaken featurisation. We undertook this process with five iterations to ensure model stability/robustness to outliers.

Survival prediction

Lastly, we fitted linear regression models seeking to predict patient overall survival (OS) (in days) from VASARI features. These were in the formulation $OS \sim 1 + f_1 + f_2 + f_n$, where f_n denotes each VASARI feature. We fitted separate models using VASARI features quantified from 1) consultant neuroradiologists, 2) VASARI-auto using the source semi-supervised and neuroradiology-reviewed segmentations, and 3) VASARI-auto using TumourSeg(20, 34), from which we compared the quality of fit. Although relatively large in kind (n=100), with the relatively small sample used here, we deliberately chose not to model with nonlinear or machine learning models nor partition data into train or test datasets, which would otherwise be highly liable to overfit in such an instance.

Analytic compliance

All analyses were performed and reported following international TRIPOD and PROBAST-AI guidelines(49).

Code, model, and data availability

The software for VASARI-auto shall be openly available upon publication at <https://github.com/jamesruffle/vasari-auto>. Our tumour segmentation model and code are openly and freely available at <https://github.com/high-dimensional/tumour-seg>. All patient data utilised in this article is freely and openly available(35, 36).

Software

The following software and models were used for analyses: Matplotlib(50), MONAI(51), Nibabel(52), Nilearn(53), NumPy(54), pandas(55), PyTorch(56), seaborn(57), scikit-learn(44), and TumourSeg(20, 34), all within a Python environment.

Compute

All experiments were performed on a 64-core Linux workstation with 256Gb of RAM and an NVIDIA 3090Ti GPU.

Results

Cohort

The brain tumour patient cohort included 56 male and 44 female participants, mean age \pm standard deviation age 61 years \pm 13.49. Mean overall survival (in days) was 436 ± 462.21 days. Seventy-four participants were included from UPenn-GBM, and 26 were included in UCSF-PDGM. There were no significant differences in age, sex, or survival between participants at either site, indicating a well-standardised and representative multi-site sample.

Segmentation

A comparison of tumour segmentation using our in-house developed model TumourSeg(20, 34) to the externally curated semi-supervised labels showed a mean Dice segmentation performance of 0.95 ± 0.05 for the whole tumour, 0.89 ± 0.07 for the enhancing tumour, 0.86 ± 0.11 for the nonenhancing tumour, and 0.91 ± 0.06 for the perilesional signal change. A visual overlay of lesion segmentations to the brain showed no spatial discrepancy (Figure 2). There was no significant difference in segmentation performance between male and female sex, and across all decades of age, indicating an equitably performant tumour segmentation model.

Tumour segmentation equitable calibration

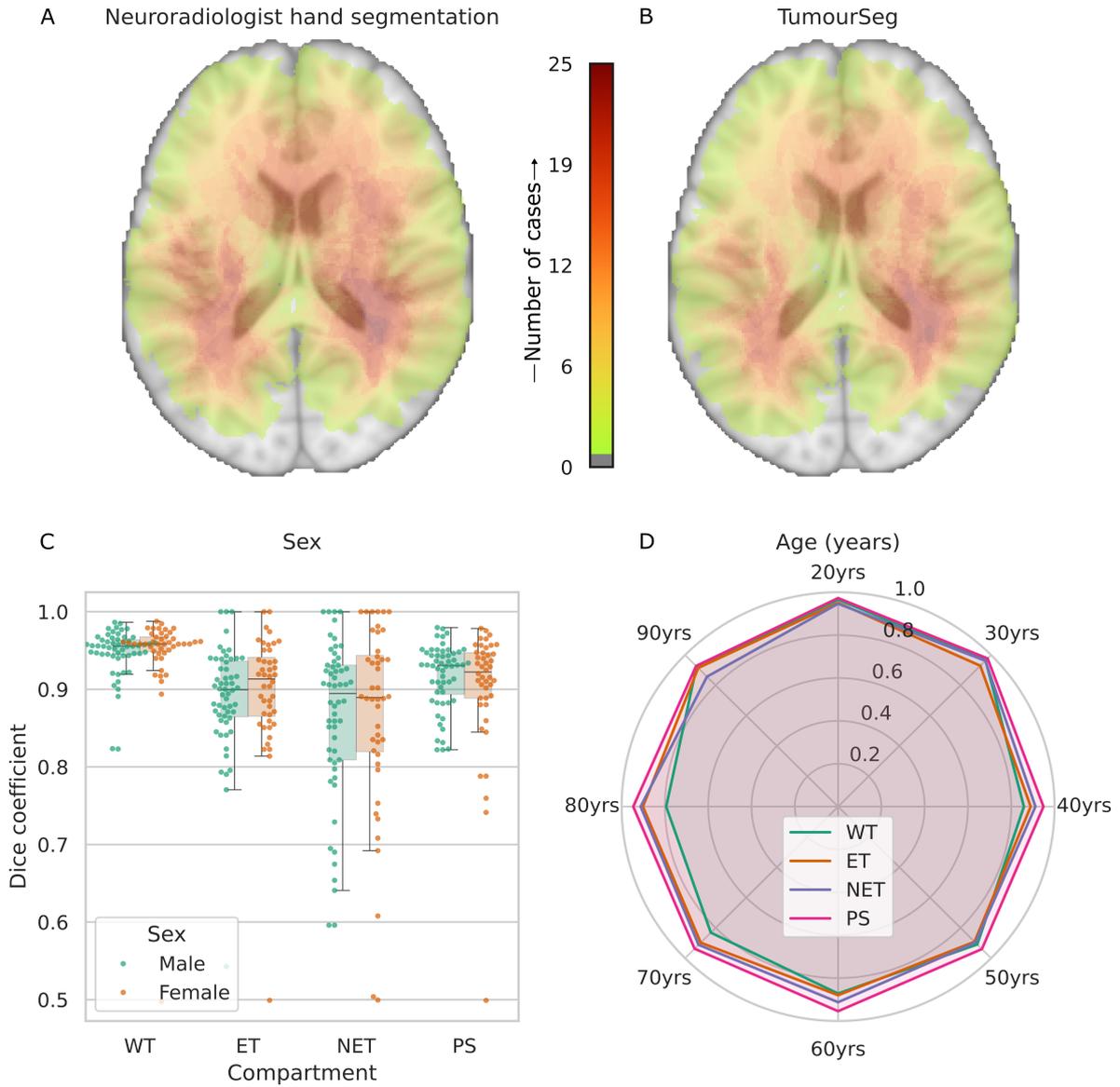


Figure 2. Tumour segmentation equitable calibration. A-B) Heatmap of tumour location derived from semi-supervised external neuroradiologist hand segmentations (A) and from our internal model TumourSeg (B) shows the two to be highly similar indicative of spatially equitable intracranial performance. C-D) Box and whisker (C) and radar (D) plots depict tumour segmentation model performance by Dice coefficient across whole tumour (WT), enhancing tumour (ET), nonenhancing tumour (NET), and perilesional signal change (PS), illustrating that tumour segmentation is equally performant across both male (green) and female (orange) patients (C), and across all decades of life (D).

Agreement and accuracy evaluations

There was a modest inter-rater agreement between consultant neuroradiologists, with a mean Cohen's Kappa of 0.49 ± 0.32 (Figure 3). The features of highest agreement between neuroradiologists were: F9 – whether a lesion was multifocal/multicentric (Cohen's Kappa 1.00); F20 – cortical involvement (Cohen's Kappa 0.71); and F21 – deep white matter invasion (Cohen's Kappa 0.71). The features of least agreement between neuroradiologists were F24 – presence of satellite lesions (Cohen's Kappa -0.23), F11 – thickness of an enhancing margin (Cohen's Kappa -0.03); and F19 – ependymal invasion (Cohen's Kappa 0.33). Agreement between neuroradiologists and VASARI-auto was relatively similar, with a mean Cohen's Kappa of 0.42 ± 0.34 . The features of highest agreement between neuroradiologists and VASARI-auto were: F9 – multifocal/multicentric (Cohen's Kappa 1.00); F2 – side of the epicentre (Cohen's Kappa 0.93), and F1 – tumour location (Cohen's Kappa 0.75). The features of least agreement between neuroradiologists and VASARI-auto were F6 – the proportion of nonenhancing tumour (Cohen's Kappa -0.16), F14 – the proportion of oedema (Cohen's Kappa 0.02), and F20 – cortical involvement (Cohen's Kappa 0.11). Agreement between VASARI-auto featurisations, whether using externally curated lesion segmentations or TumourSeg, was substantially higher, with a mean Cohen's Kappa of 0.94 ± 0.10 . For this comparison, all feature agreement was 0.88 or higher, apart from F24 – the presence of satellite lesions (Cohen's Kappa 0.59).

Treating neuroradiologist labels as a ground truth, VASARI-auto achieved a mean accuracy of $66\% \pm 21\%$ (Figure 3). The greatest accuracy was in F2 – side of tumour epicentre (96.55%), F24 – the presence of satellite lesions (86.20%), and F1 – tumour location (80.46%). The lowest accuracy with VASARI-auto was in F6 – the proportion of nonenhancing tumour (20.69%), F14 – the proportion of oedema (32.18%), and F5 – the proportion of enhancing tumour (49.42%). In contrast, when treating VASARI-auto when derived from the external semi-supervised lesion labels as a ground truth, VASARI-auto accuracy using the in-house tumour segmentation model was much more stable, with a mean accuracy of $97.40 \pm 0.03\%$.

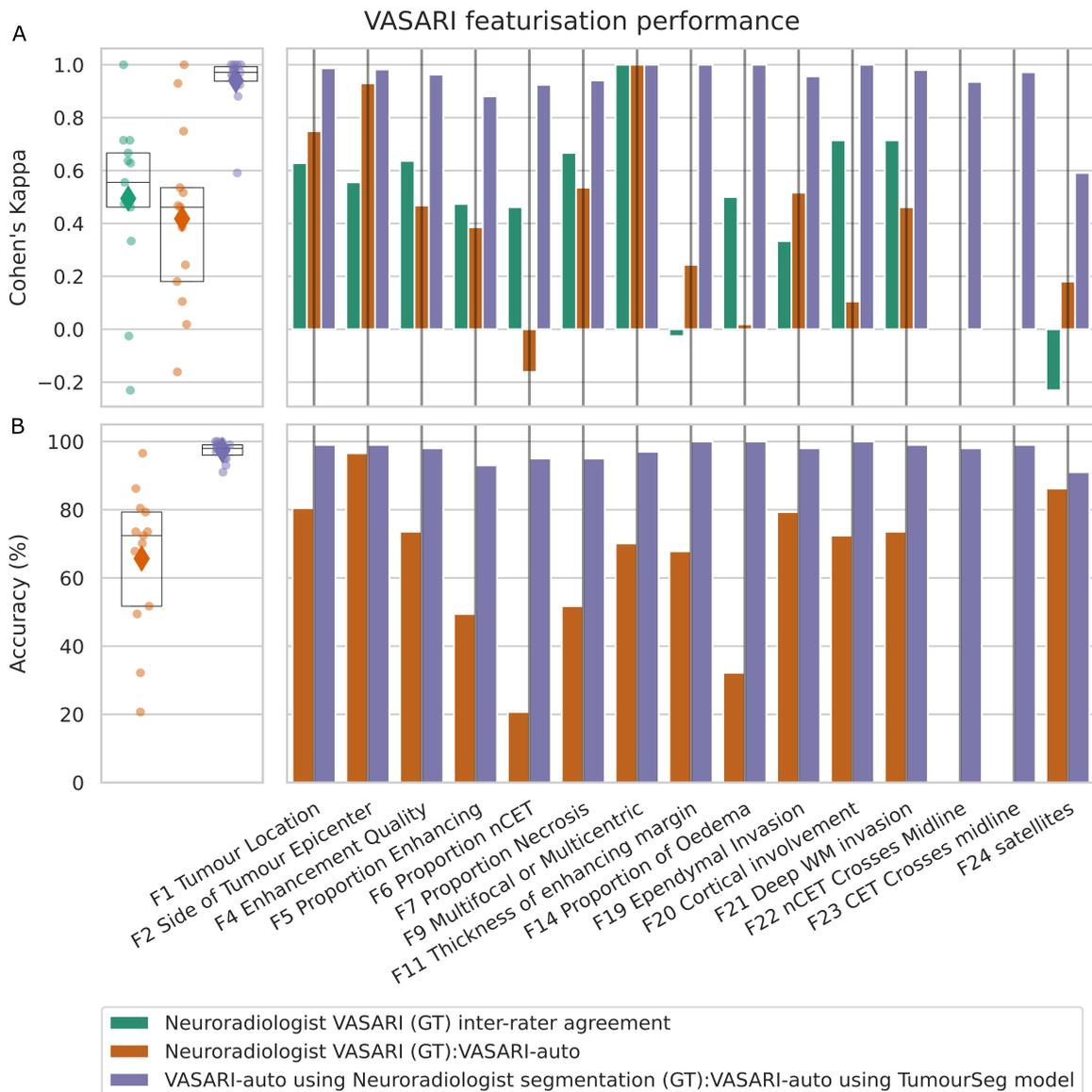


Figure 3. VASARI featurisation performance. A) Cohen's Kappa agreement between neuroradiologist reporters (green), between neuroradiologists and VASARI-auto (orange), and between VASARI-auto with and without using TumourSeg (purple) shows inter-rater variability between neuroradiologists but quantitatively higher agreement and stability between both VASARI-auto methods. B) Accuracy between neuroradiologist VASARI reporting (the ground truth, GT) and VASARI-auto (orange), and between VASARI-auto with and without using TumourSeg (purple). VASARI-auto is generally performant compared to neuroradiologists, although some discrepancies are evident due to diverging definitions of what is referred to as a nonenhancing tumour and what is oedema. VASARI-auto cross-model comparison is highly accurate. Abbreviations: nCET, non-contrast-enhancing tumour; WM, white matter.

Efficiency

The use of VASARI-auto in VASARI featurisation – regardless of whether used in isolation or when paired with our tumour segmentation model – was significantly faster per case compared to consultant neuroradiologists ($p < 0.0001$) (Figure 4). The mean time to quantify was 3.03 ± 0.59 seconds with VASARI-auto, which was significantly higher but notably still efficient at 15.47 ± 1.56 [95%CI] using VASARI-auto with TumourSeg ($p < 0.0001$). In comparison, the mean time to quantify was 317.46 (i.e., 5.28 minutes) ± 96.89 seconds with consultant neuroradiologists.

Simulated workforce analysis

The simulated workforce analysis forecast that, over 2024-2027, a total cumulative 8150 ± 168 cases would require discussion at each weekly neuro-oncology MDT (Figure 4). For VASARI featurisation to be undertaken in all cases, this would demand 744.43 ± 15.54 consultant neuroradiologist workforce hours, equating to $\pounds 39,373.37 \pm \pounds 864.22$ in salary remuneration for hours worked. In contrast, quantifying VASARI features with VASARI-auto for all cases over three years would require 8.30 ± 0.15 hours of computing time (time comparison $p < 0.0001$), equating to approximately $\pounds 3.65 \pm 0.12$ for power costs (cost comparison $p < 0.0001$). If combined with tumour segmentation, this time and expense would rise slightly to 34.89 ± 0.70 hours of computing time and $\pounds 15.17 \pm 0.55$ for power costs (both of which remained significantly less than with neuroradiologist labelling). Time taken and costs remained significantly greater for featurisation by neuroradiologists compared to VASARI-auto ($p < 0.0001$).

We scaled this up to all 40 neuro-oncology centres across the UK. For VASARI featurisation to be undertaken in all UK cases, this would demand 29,777.39 consultant neuroradiologist workforce hours, equating to $\pounds 1,574,935$ in salary remuneration for hours worked. In contrast, quantifying VASARI features with VASARI-auto for all cases over three years would require 331.95 hours of computing time, equating to approximately $\pounds 145.85$ for power costs. If combined with tumour segmentation, this time and expense would rise slightly to 1394.42 hours of computing time and $\pounds 606.75$ for power costs. Both time taken and cost were significantly greater for featurisation by neuroradiologists compared to VASARI-auto with or without the addition of TumourSeg (all $p < 0.0001$).

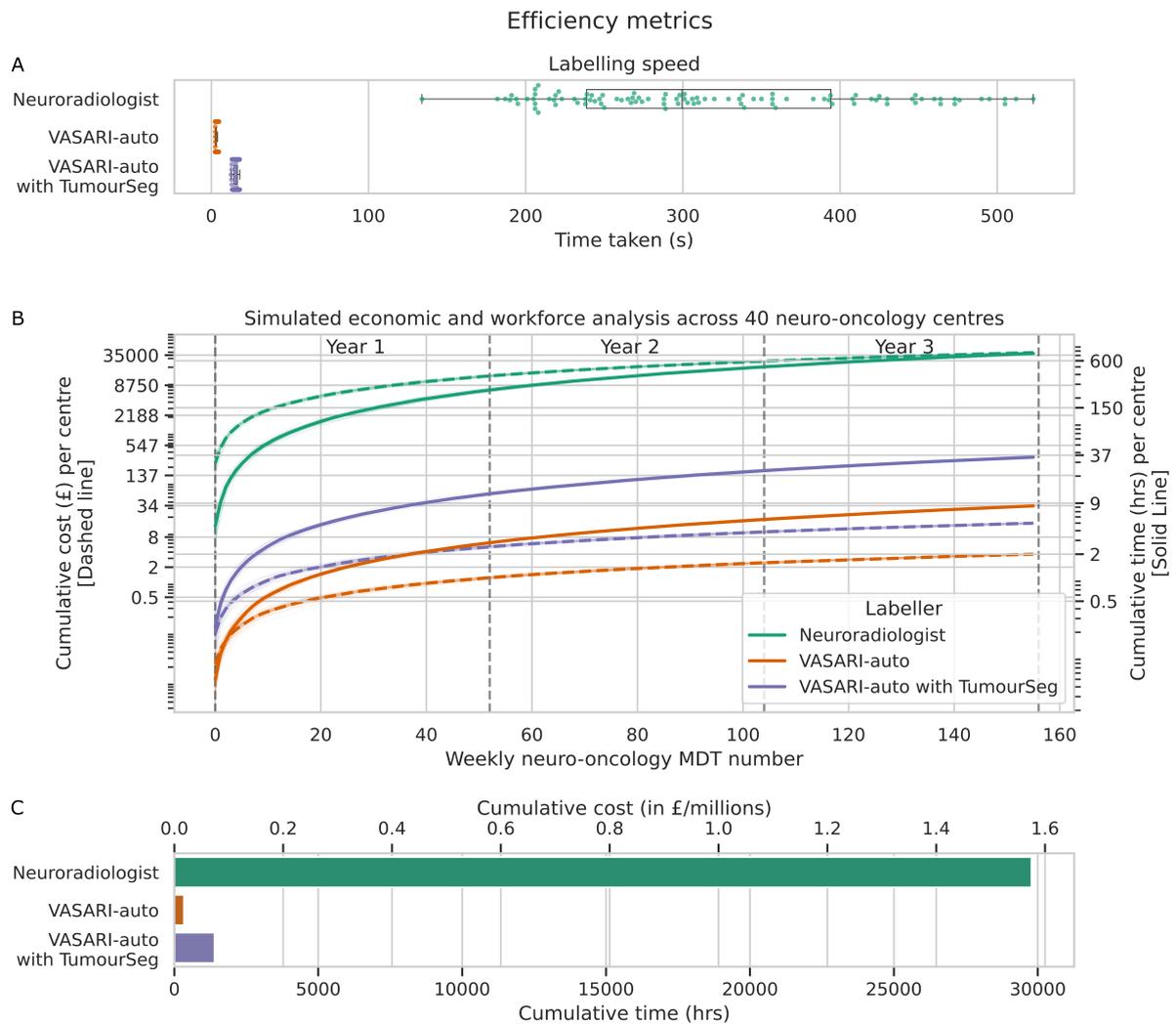


Figure 4. Efficiency, economic and workforce planning analysis. A) The time taken for a neuroradiologist to derive a VASARI feature set for a single patient (green) is substantially higher than with either VASARI-auto (orange) or using VASARI-auto paired with our tumour segmentation model, TumourSeg (purple). B) Simulated economic and workforce analysis, where weekly neuro-oncology multidisciplinary team (MDT) workload is drawn from a random uniform distribution based upon the last three years of workload at our centre. Thin individual lines represent different simulation runs to emulate the forty different UK neuro-oncology centres, with thicker lines representing the epoch mean. Cumulative financial cost (dashed line) and time taken (solid line) for a neuroradiologist (green), VASARI-auto (orange), and VASARI-auto with TumourSeg model (purple) to derive VASARI features for each patient for each week of the neuro-oncology MDT. The financial cost for consultant neuroradiologists is drawn from a uniform distribution of current National Health Service consultant pay scales provided by the British Medical Association. The power cost for VASARI-auto is based upon current energy prices for a modest 1200 Kilowatt-hour (kWh) GPU-supported computer. Note both y-axes are logged. C) Total UK-wide cumulative cost in time and resource for a neuroradiologist (green), VASARI-auto (orange), and VASARI-auto with TumourSeg (purple) to derive VASARI features for each patient.

Performance equity

A critical performance measure of any automated tool in healthcare is invariance across patient background characteristics(46). We compared reporting agreement between 1) neuroradiologists, 2) neuroradiologists and VASARI-auto, and 3) VASARI-auto when applied to the externally curated tumour segmentations or with the in-house developed TumourSeg(20), with respect to patient age and sex, using Cohen's Kappa (Figure 5). There was no evidence of reporting inequity between neuroradiologists and between neuroradiologists and VASARI-auto (allowing for the more limited distribution of demographics for those patients double-reported by neuroradiologists). Similarly, agreement between VASARI-auto using external or internal lesion segmentations was equally performant across patient age and sex, all of which indicate equitability of VASARI-auto and tumour segmentation models.

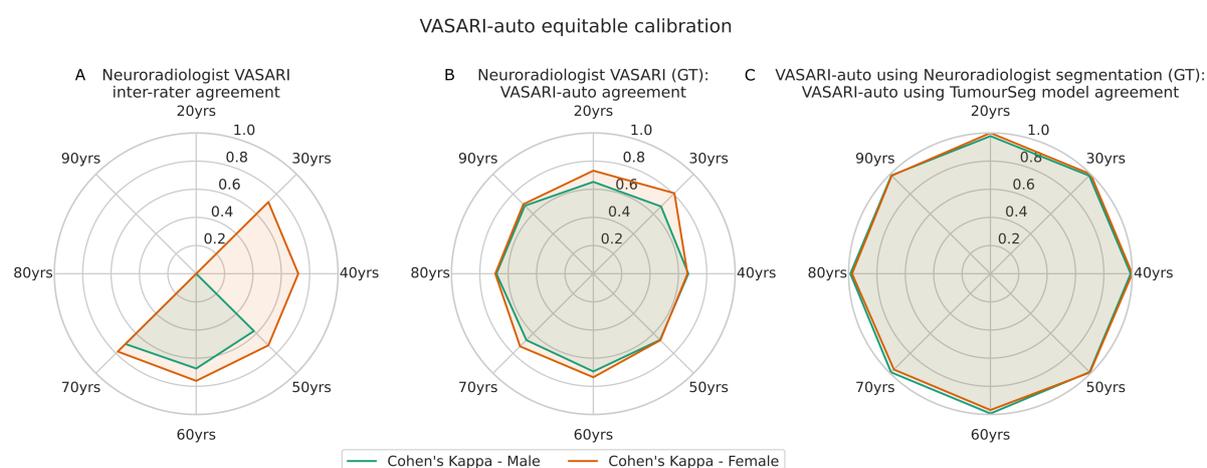


Figure 5. VASARI featurisation equitable calibration. A-C) radar plots showing Cohen's Kappa agreement aligned to male (green) or female (orange) patient sex across all decades of life. A) inter-rater agreement between neuroradiologists, B) between neuroradiologists and VASARI-auto, and C) between VASARI-auto using external or internal (TumourSeg) curated lesion segmentations.

Survival prediction

The clinical utility of any feature is ultimately determined by its downstream predictive, prescriptive, or inferential power. Fidelity in overall survival prediction using VASARI features was qualitatively similar whether using feature sets derived by consultant neuroradiologists, from VASARI-auto applied to the semi-supervised and neuroradiology reviewed segmentations, or VASARI-auto paired with TumourSeg (Figure 6). Quantitatively, the best linear regression fit was achieved with VASARI-auto using the semi-supervised and neuroradiology-reviewed segmentations (R^2 0.245), followed closely by VASARI-auto using TumourSeg (R^2 0.227), with slightly weaker performance when using the consultant neuroradiologist-labelled VASARI features (R^2 0.205).

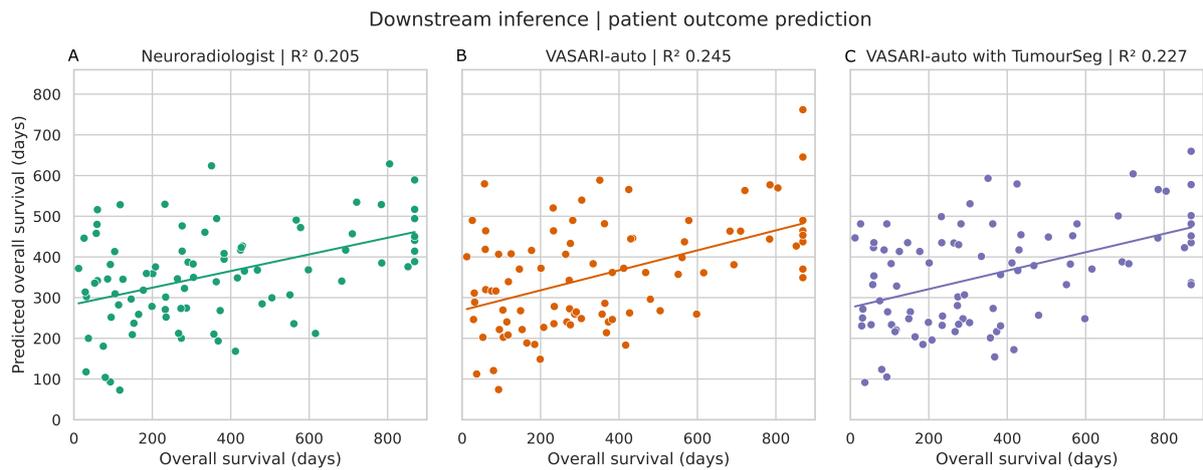


Figure 6. Downstream inference with patient outcome prediction. Results of linear regression predicting overall patient survival in days using VASARI-features derived by A) neuroradiologists (green), B) VASARI-auto from the semi-supervised external segmentations (orange), or C) VASARI-auto combined with our in-house developed model, TumourSeg (purple). X-axes illustrate the actual survival, whereas y-axes illustrate predicted survival. There is highly similar qualitative performance in survival prediction regardless of whether a neuroradiologist or VASARI-auto labels it, although quantitatively, the R^2 is higher with both VASARI-auto assessments.

Discussion

We present VASARI-auto, an automated system for deriving VASARI features from glioma imaging using tumour segmentations alone. Our evaluation shows high accuracy, greater consistency than inter-agreement between neuroradiologists, and equitable performance across age and sex. We show VASARI-auto could save time and resources at each centre, equating over three years to 771 hours of consultant neuroradiologist time or ~£40,000 (>\$50,000) in NHS finance terms, given the workload of a neuro-oncology centre such as ours. Scaled across the UK, the saving is anticipated to be more than £1.5 million (\$1.9 million). Framed differently, such software would enable these workforce hours to be reallocated to other areas of unmet clinical need. We furthermore show that patient survival forecasting is non-inferior when using these automated models, demonstrating the preservation of feature fidelity.

Adding value with AI-assisted practice

Despite being well-validated in research to provide well-structured information on the imaging appearances of glioma, the VASARI feature set is seldom used in clinical practice. The causes for this are multifactorial but are likely a combination of high clinical workload – it is time-consuming to record – and lack of sufficient level of neuroradiology training and experience. Our software substantially lowers the barriers to adopting VASARI scoring while maintaining fidelity and assuring patient equity. Its introduction provides a means of extracting more detailed patient-personalized information, aiming to improve clinical care at a very modest cost in either time or financial terms. Particularly pertinent in the UK, where the number of

radiologists per 100,000 population is one of the lowest in Western Europe(58) – only 7 – such decision support tools add high value or even free up an already overstretched workforce to allow work in other clinical areas.

A critical measure of the value of any feature, automated or manual, is downstream utility, such as survival prediction. Our analysis shows non-inferior – rather, quantitatively higher – predictive fidelity in using VASARI-auto features over those curated by consultant neuroradiologists. Demonstrating non-inferiority in software that is resource-cheap, contrasted with a time-consuming process for experienced neuroradiologists, is essential for software that provides inferior care than the current clinical standard, regardless of any efficiency or cost saving, adds little value.

Maximising reporting consistency

Clinicians' opinions—whether radiologists or others—often differ. This is to be expected: diseases are typically heterogeneous. Patients, too, are heterogeneous: a successful treatment approach for one might not be suitable for another(30, 59). However, a model capable of absorbing heterogeneity can yield a quantitative description that exhibits consistency across the population concerning a critical decision. From follow-up monitoring of tumours, we know that conventional two-dimensional measurements can be highly inconsistent between radiologists(25-27), motivating the pursuit of alternative approaches.

What are the characteristics of an optimal approach? The ideal would be to absorb all variation irrelevant to the task. For example, where measurements are undertaken by manual annotation – which one should note is the currently adopted clinical practice globally, despite their empirically observed limitations(25-27) – this is trivial to ameliorate using automated software and relatively simple mathematics. We exemplify this here, showing only modest inter-rater agreement between highly experienced neuroradiologists that can be stabilised and standardised with automated methods. Particularly pertinent examples are in deriving VASARI features (or, for that matter, any other radiological feature outside the scope of this article) that are ultimately quantitative. Where the quality of a lesion segmentation is validated, such as we show in our comparison between source segmentations and those using our in-house model, then the mathematical derivation of precise proportions of lesional compartments, such as enhancing tumour, nonenhancing tumour, and perilesional signal change, is a simple mathematical operation of compartmental ratios. This is especially true for quantitative features that are harder to quantify intuitively – such as the thickness of enhancing tumour. Gliomas – particularly glioblastoma – are highly variable in their appearance. One part of an enhancing margin (if any) might be considerably thicker than another: how do we measure this? We would argue that the wrong answer (although commonplace in clinical reporting) would be to hedge an approximation between the lower and upper limits. Instead, a more robust solution is a simple mathematical derivation operating on a lesion segmentation(20). However, the difficulty one faces, as is evidenced in these works, is where ambiguity in the ground truth – namely, what is a nonenhancing tumour and what is oedema – compounds an

assessment regardless of whether derived by an experienced neuroradiologist or by software. An answer to this problem is unlikely to be solved by clinical experience, status quo imaging techniques or software, but rather by innovation across all three.

Maximising performance equity

Healthcare should be equitable, which extends to any such tool at our clinical disposal(45, 46). Artificial intelligence is one of the domains seeing the quickest growth in all research, industry, and society, with many purported applications across medicine(31). Yet equitable calibration to ensure that software brings benefits to all is relatively rarely quantified. For these reasons, we assess performance equity, not only of VASARI-auto but also of our tumour segmentation model. Though confined to age and sex – both notably legible from brain imaging(60) – the approach can be scaled through representation learning to encompass any characteristic(46).

Limitations

Our study has limitations. First, although drawn from a larger cohort of 1172, we utilise a sample of 100 patients with glioma who have undergone comprehensive clinical VASARI featurisation by experienced neuroradiologists. Although large for the domain, this is statistically a relatively small sample, and further validation should be undertaken at a larger scale. This sample, however, is carefully curated and includes imaging from two major US medical centres, for which we could evaluate both the performance of our tumour segmentation model and VASARI-auto.

Second, we could not incorporate all features of the VASARI set in the software, specifically those that required structural neuroimaging, additional sequences, or where variability/confabulation could occur. This decision was deliberate, for we wished to develop software that did not use patient-identifiable data. Instead, only a lesion segmentation is required. Future work should expand upon this to include these remaining features.

Thirdly, the extent of the VASARI-auto featurisation pipeline was gated by the availability of widely used lesion compartment labels, namely enhancing tumour, nonenhancing tumour, and oedema(37). Therefore, our software could not provide VASARI data on haemorrhagic change because no label exists in the source data(35, 36): it cannot learn what it has not been taught(20). There are evolving opinions across neuro-oncology as to what may be oedema and what is nonenhancing tumour: it is for this reason we use the terminology of ‘perilesional signal change’ in discussing the segmentation pipeline. Moreover, it is the reason for lower agreement between VASARI-auto and neuroradiologist reporting for the lesion proportion features, since the software is guided by the status quo where such perilesional signal is referred to as oedema, though some radiologists might instead label as nonenhancing tumour. In any case, it should be stressed that the values themselves do not actually matter here, but of far greater importance is that from lesion segmentation, a more robust and standardised

assessment across a cohort of patients is yielded, likely the reason for stronger performance in downstream survival prediction.

Fourthly, our economic cost analysis assumes a VASARI feature set is undertaken in all cases assigned to the neuro-oncology MDT. This is an upper bound: VASARI is seldom used for the time and level of professional training it demands. Furthermore, given the sharp rise in the volume of medical imaging undertaken for patients globally(58, 61), it is likely that the hours and cost incurred for radiologists to featurise these cases are a significant under-representation. However, precisely to that point, one should consider the economic analysis to highlight a gain in healthcare value at negligible time or financial cost.

Conclusions

VASARI-auto can characterise glioma efficiently, effectively, and equitably. The use of VASARI-auto-derived features in predicting patient survival is non-inferior to the use of those manually curated by experienced consultant neuroradiologists. Translation to the clinical frontline with an automated derivation of these features may enhance existing clinical practice with negligible cost to a centre, serving as a decision support tool to provide healthcare providers with more information to facilitate standardised, equitable, and more personalised patient care.

References

1. TCIA. VASARI Research Project. <https://wiki.cancerimagingarchive.net/display/Public/VASARI+Research+Project>. Published 2020. Accessed 2024 07/03/2024.
2. Gusev Y, Bhuvaneshwar K, Song L, Zenklusen JC, Fine H, Madhavan S. The REMBRANDT study, a large collection of genomic data from brain cancer patients. *Sci Data* 2018;5:180158. doi: 10.1038/sdata.2018.158
3. Gemini L, Tortora M, Giordano P, Prudente ME, Villa A, Vargas O, Giugliano MF, Somma F, Marchello G, Chiaramonte C, Gaetano M, Frio F, Di Giorgio E, D'Avino A, Tortora F, D'Agostino V, Negro A. Vasari Scoring System in Discerning between Different Degrees of Glioma and IDH Status Prediction: A Possible Machine Learning Application? *J Imaging* 2023;9(4). doi: 10.3390/jimaging9040075
4. Park CJ, Han K, Kim H, Ahn SS, Choi D, Park YW, Chang JH, Kim SH, Cha S, Lee SK. MRI Features May Predict Molecular Features of Glioblastoma in Isocitrate Dehydrogenase Wild-Type Lower-Grade Gliomas. *AJNR Am J Neuroradiol* 2021;42(3):448-456. doi: 10.3174/ajnr.A6983
5. Setyawan NH, Choridah L, Nugroho HA, Malueka RG, Dwianingsih EK. Beyond invasive biopsies: using VASARI MRI features to predict grade and molecular parameters in gliomas. *Cancer Imaging* 2024;24(1):3. doi: 10.1186/s40644-023-00638-8
6. Peeken JC, Hesse J, Haller B, Kessel KA, Nusslin F, Combs SE. Semantic imaging features predict disease progression and survival in glioblastoma multiforme patients. *Strahlenther Onkol* 2018;194(6):580-590. doi: 10.1007/s00066-018-1276-4
7. Peeken JC, Goldberg T, Pyka T, Bernhofer M, Wiestler B, Kessel KA, Tafti PD, Nusslin F, Braun AE, Zimmer C, Rost B, Combs SE. Combining multimodal imaging and treatment features improves machine learning-based prognostic assessment in patients with glioblastoma multiforme. *Cancer Med* 2019;8(1):128-136. doi: 10.1002/cam4.1908
8. Nicolasjilwan M, Hu Y, Yan C, Meerzaman D, Holder CA, Gutman D, Jain R, Colen R, Rubin DL, Zinn PO, Hwang SN, Raghavan P, Hammoud DA, Scarpance LM, Mikkelsen T, Chen J, Gevaert O, Buetow K, Freymann J, Kirby J, Flanders AE, Wintermark M, Group TGPR. Addition of MR imaging features and genetic biomarkers strengthens glioblastoma survival prediction in TCGA patients. *Journal of neuroradiology Journal de neuroradiologie* 2015;42(4):212-221. doi: 10.1016/j.neurad.2014.02.006
9. Jain R, Poisson LM, Gutman D, Scarpance L, Hwang SN, Holder CA, Wintermark M, Rao A, Colen RR, Kirby J, Freymann J, Jaffe CC, Mikkelsen T, Flanders A. Outcome prediction in patients with glioblastoma by using imaging, clinical, and genomic biomarkers: focus on the nonenhancing component of the tumor. *Radiology* 2014;272(2):484-493. doi: 10.1148/radiol.14131691
10. Zhou H, Vallieres M, Bai HX, Su C, Tang H, Oldridge D, Zhang Z, Xiao B, Liao W, Tao Y, Zhou J, Zhang P, Yang L. MRI features predict survival and molecular markers in diffuse lower-grade gliomas. *Neuro Oncol* 2017;19(6):862-870. doi: 10.1093/neuonc/now256

11. Wang J, Yi X, Fu Y, Pang P, Deng H, Tang H, Han Z, Li H, Nie J, Gong G, Hu Z, Tan Z, Chen BT. Preoperative Magnetic Resonance Imaging Radiomics for Predicting Early Recurrence of Glioblastoma. *Front Oncol* 2021;11:769188. doi: 10.3389/fonc.2021.769188
12. Biswas A, Amirabadi A, Wagner MW, Ertl-Wagner BB. Features of Visually Accessible Rembrandt Images: Interrater Reliability in Pediatric Brain Tumors. *AJNR Am J Neuroradiol* 2022;43(2):304-308. doi: 10.3174/ajnr.A7399
13. Li L, Fu Y, Zhang Y, Mao Y, Huang D, Yi X, Wang J, Tan Z, Jiang M, Chen BT. Magnetic resonance imaging findings of intracranial extraventricular ependymoma: A retrospective multi-center cohort study of 114 cases. *Cancer Med* 2023;12(15):16195-16206. doi: 10.1002/cam4.6279
14. Deeley MA, Chen A, Datteri R, Noble JH, Cmelak AJ, Donnelly EF, Malcolm AW, Moretti L, Jaboin J, Niermann K, Yang ES, Yu DS, Yei F, Koyama T, Ding GX, Dawant BM. Comparison of manual and automatic segmentation methods for brain structures in the presence of space-occupying lesions: a multi-expert study. *Phys Med Biol* 2011;56(14):4557-4577. doi: 10.1088/0031-9155/56/14/021
15. Wan Y, Rahmat R, Price SJ. Deep learning for glioblastoma segmentation using preoperative magnetic resonance imaging identifies volumetric features associated with survival. *Acta Neurochir (Wien)* 2020;162(12):3067-3080. doi: 10.1007/s00701-020-04483-7
16. NHS. Guide to NHS waiting times in England. <https://www.nhs.uk/nhs-services/hospitals/guide-to-nhs-waiting-times-in-england/>. Published 2024. Accessed 2024 08/03/2024.
17. Peng J, Kim DD, Patel JB, Zeng X, Huang J, Chang K, Xun X, Zhang C, Sollee J, Wu J, Dalal DJ, Feng X, Zhou H, Zhu C, Zou B, Jin K, Wen PY, Boxerman JL, Warren KE, Poussaint TY, States LJ, Kalpathy-Cramer J, Yang L, Huang RY, Bai HX. Corrigendum to: Deep learning-based automatic tumor burden assessment of pediatric high-grade gliomas, medulloblastomas, and other leptomeningeal seeding tumors. *Neuro-Oncology* 2021;23(12):2124-2124. doi: 10.1093/neuonc/noab226
18. Xue J, Wang B, Ming Y, Liu X, Jiang Z, Wang C, Liu X, Chen L, Qu J, Xu S, Tang X, Mao Y, Liu Y, Li D. Deep learning-based detection and segmentation-assisted management of brain metastases. *Neuro-Oncology* 2020;22(4):505-514. doi: 10.1093/neuonc/noz234
19. Lu S-L, Xiao F-R, Cheng JC-H, Yang W-C, Cheng Y-H, Chang Y-C, Lin J-Y, Liang C-H, Lu J-T, Chen Y-F, Hsu F-M. Randomized multi-reader evaluation of automated detection and segmentation of brain tumors in stereotactic radiosurgery with deep neural networks. *Neuro-Oncology* 2021;23(9):1560-1568. doi: 10.1093/neuonc/noab071
20. Ruffle J, K., Mohinta S, Gray R, Hyare H, Nachev P. Brain tumour segmentation with incomplete imaging data. *Brain Commun* 2023. doi: 10.1093/braincomms/fcad118
21. Lenchik L, Heacock L, Weaver AA, Boutin RD, Cook TS, Itri J, Filippi CG, Gullapalli RP, Lee J, Zagurovskaya M, Retson T, Godwin K, Nicholson J, Narayana PA. Automated Segmentation of Tissues Using CT and MRI: A Systematic Review. *Acad Radiol* 2019;26(12):1695-1706. doi: 10.1016/j.acra.2019.07.006

22. Suetens P, Bellon E, Vandermeulen D, Smet M, Marchal G, Nuyts J, Mortelmans L. Image segmentation: methods and applications in diagnostic radiology and nuclear medicine. *Eur J Radiol* 1993;17(1):14-21. doi: 10.1016/0720-048x(93)90023-g
23. Ashburner J, Friston KJ. Unified segmentation. *Neuroimage* 2005;26(3):839-851. doi: 10.1016/j.neuroimage.2005.02.018
24. Menze BH, Jakab A, Bauer S, Kalpathy-Cramer J, Farahani K, Kirby J, Burren Y, Porz N, Slotboom J, Wiest R, Lanczi L, Gerstner E, Weber MA, Arbel T, Avants BB, Ayache N, Buendia P, Collins DL, Cordier N, Corso JJ, Criminisi A, Das T, Delingette H, Demiralp C, Durst CR, Dojat M, Doyle S, Festa J, Forbes F, Geremia E, Glocker B, Golland P, Guo X, Hamamci A, Iftekharuddin KM, Jena R, John NM, Konukoglu E, Lashkari D, Mariz JA, Meier R, Pereira S, Precup D, Price SJ, Raviv TR, Reza SM, Ryan M, Sarikaya D, Schwartz L, Shin HC, Shotton J, Silva CA, Sousa N, Subbanna NK, Szekely G, Taylor TJ, Thomas OM, Tustison NJ, Unal G, Vasseur F, Wintermark M, Ye DH, Zhao L, Zhao B, Zikic D, Prastawa M, Reyes M, Van Leemput K. The Multimodal Brain Tumor Image Segmentation Benchmark (BRATS). *IEEE Trans Med Imaging* 2015;34(10):1993-2024. doi: 10.1109/TMI.2014.2377694
25. Zhao B, James LP, Moskowitz CS, Guo P, Ginsberg MS, Lefkowitz RA, Qin Y, Riely GJ, Kris MG, Schwartz LH. Evaluating variability in tumor measurements from same-day repeat CT scans of patients with non-small cell lung cancer. *Radiology* 2009;252(1):263-272. doi: 10.1148/radiol.2522081593
26. McNitt-Gray MF, Kim GH, Zhao B, Schwartz LH, Clunie D, Cohen K, Petrick N, Fenimore C, Lu ZQ, Buckler AJ. Determining the Variability of Lesion Size Measurements from CT Patient Data Sets Acquired under "No Change" Conditions. *Transl Oncol* 2015;8(1):55-64. doi: 10.1016/j.tranon.2015.01.001
27. Dempsey MF, Condon BR, Hadley DM. Measurement of tumor "size" in recurrent malignant glioma: 1D, 2D, or 3D? *AJNR Am J Neuroradiol* 2005;26(4):770-776.
28. Mandal AS, Romero-Garcia R, Hart MG, Suckling J. Genetic, cellular, and connectomic characterization of the brain regions commonly plagued by glioma. *Brain* 2020;143(11):3294-3307. doi: 10.1093/brain/awaa277
29. Topol E. The Topol Review: Preparing the healthcare workforce to deliver the digital future. In: NHS, ed.2019.
30. Rajpurkar P, Chen E, Banerjee O, Topol EJ. AI in health and medicine. *Nature Medicine* 2022. doi: 10.1038/s41591-021-01614-0
31. Ruffle JK, Farmer AD, Aziz Q. Artificial Intelligence Assisted Gastroenterology - Promises and Pitfalls. *Am J Gastroenterol* 2019; 114(3):422-428. doi: 10.1038/s41395-018-0268-4.
32. Ruffle JK, Mohinta S, Pombo G, Gray R, Kopanitsa V, Lee F, Brandner S, Hyare H, Nachev P. Brain tumour genetic network signatures of survival. *Brain* 2023.
33. Louis DN, Perry A, Wesseling P, Brat DJ, Cree IA, Figarella-Branger D, Hawkins C, Ng HK, Pfister SM, Reifenberger G, Soffietti R, von Deimling A, Ellison DW. The 2021 WHO

Classification of Tumors of the Central Nervous System: a summary. *Neuro Oncol* 2021;23(8):1231-1251. doi: 10.1093/neuonc/noab106

34. Ruffle JK. Brain tumour segmentation with incomplete imaging data. GitHub. <https://github.com/high-dimensional/tumour-seg>. Published 2023.

35. Calabrese E, Villanueva-Meyer JE, Rudie JD, Rauschecker AM, Baid U, Bakas S, Cha S, Mongan JT, Hess CP. The University of California San Francisco Preoperative Diffuse Glioma MRI Dataset. *Radiol Artif Intell* 2022;4(6):e220058. doi: 10.1148/ryai.220058

36. Bakas S, Sako C, Akbari H, Bilello M, Sotiras A, Shukla G, Rudie JD, Santamaría NF, Kazerooni AF, Pati S, Rathore S, Mamourian E, Ha SM, Parker W, Doshi J, Baid U, Bergman M, Binder ZA, Verma R, Lustig RA, Desai AS, Bagley SJ, Mourelatos Z, Morrissette J, Watt CD, Brem S, Wolf RL, Melhem ER, Nasrallah MP, Mohan S, O'Rourke DM, Davatzikos C. The University of Pennsylvania glioblastoma (UPenn-GBM) cohort: advanced MRI, clinical, genomics, & radiomics. *Scientific Data* 2022;9(1):453. doi: 10.1038/s41597-022-01560-7

37. Baid U, Ghodasara S, Bilello M, Mohan S, Calabrese E, Colak E, Farahani K, Kalpathy-Cramer J, Kitamura F, Pati S, Prevedello L, Rudie J, Sako C, Shinohara R, Bergquist T, Chai R, Eddy J, Elliott J, Reade W, Bakas S. The RSNA-ASNR-MICCAI BraTS 2021 Benchmark on Brain Tumor Segmentation and Radiogenomic Classification, 2021.

38. Isensee F, Jaeger PF, Kohl SAA, Petersen J, Maier-Hein KH. nnU-Net: a self-configuring method for deep learning-based biomedical image segmentation. *Nat Methods* 2021;18(2):203-211. doi: 10.1038/s41592-020-01008-z

39. Isensee F, Jaeger PF, Full PM, Vollmuth P, Maier-Hein K. nnU-Net for Brain Tumor Segmentation. *BrainLes@MICCAI2020*.

40. Antonelli M, Reinke A, Bakas S, Farahani K, Kopp-Schneider A, Landman BA, Litjens G, Menze B, Ronneberger O, Summers RM, van Ginneken B, Bilello M, Bilic P, Christ PF, Do RKG, Gollub MJ, Heckers SH, Huisman H, Jarnagin WR, McHugo MK, Napel S, Pernicka JSG, Rhode K, Tobon-Gomez C, Vorontsov E, Meakin JA, Ourselin S, Wiesenfarth M, Arbelaez P, Bae B, Chen S, Daza L, Feng J, He B, Isensee F, Ji Y, Jia F, Kim I, Maier-Hein K, Merhof D, Pai A, Park B, Perslev M, Rezaiifar R, Rippel O, Sarasua I, Shen W, Son J, Wachinger C, Wang L, Wang Y, Xia Y, Xu D, Xu Z, Zheng Y, Simpson AL, Maier-Hein L, Cardoso MJ. The Medical Segmentation Decathlon. *Nat Commun* 2022;13(1):4128. doi: 10.1038/s41467-022-30695-9

41. Nachev P, Coulthard E, Jager HR, Kennard C, Husain M. Enantiomorphic normalization of focally lesioned brains. *Neuroimage* 2008;39(3):1215-1226. doi: 10.1016/j.neuroimage.2007.10.002

42. Ashburner J, Friston KJ. Nonlinear spatial normalization using basis functions. *Hum Brain Mapp* 1999;7(4):254-266. doi: 10.1002/(SICI)1097-0193(1999)7:4<#x0003c;254::AID-HBM4#x0003e;3.0.CO;2-G

43. Ritaccio AL, Brunner P, Schalk G. Electrical Stimulation Mapping of the Brain: Basic Principles and Emerging Alternatives. *J Clin Neurophysiol* 2018;35(2):86-97. doi: 10.1097/WNP.0000000000000440

44. Pedregosa F, Varoquaux G, Gramfort A. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* 2011;12:2825-2830.
45. Abramoff MD, Tarver ME, Loyo-Berrios N, Trujillo S, Char D, Obermeyer Z, Eydelman MB, Foundational Principles of Ophthalmic I, Algorithmic Interpretation Working Group of the Collaborative Community for Ophthalmic Imaging Foundation WDC, Maisel WH. Considerations for addressing bias in artificial intelligence for health equity. *NPJ Digit Med* 2023;6(1):170. doi: 10.1038/s41746-023-00913-9
46. Carruthers R, Straw I, Ruffle JK, Herron D, Nelson A, Bzdok D, Fernandez-Reyes D, Rees G, Nachev P. Representational ethical model calibration. *NPJ Digit Med* 2022;5(1):170. doi: 10.1038/s41746-022-00716-4
47. Employers N. Pay and Conditions Circular (M&D) 4/2023. https://www.nhsemployers.org/system/files/2023-08/Pay%20and%20Conditions%20Circular%20%28MD%29%204-2023%20FINAL_0.pdf. Published 2023. Accessed 2024 05/03/24.
48. sust-it. Energy Cost Calculator - UK: Price Cap (Jan 2024). <https://www.sust-it.net/energy-calculator.php>. Published 2024. Accessed 2024 05/03/2024.
49. Collins GS, Dhiman P, Andaur Navarro CL, Ma J, Hooft L, Reitsma JB, Logullo P, Beam AL, Peng L, Van Calster B, van Smeden M, Riley RD, Moons KG. Protocol for development of a reporting guideline (TRIPOD-AI) and risk of bias tool (PROBAST-AI) for diagnostic and prognostic prediction model studies based on artificial intelligence. *BMJ Open* 2021;11(7):e048008. doi: 10.1136/bmjopen-2020-048008
50. Hunter JD. Matplotlib: A 2D Graphics Environment. *Computing in Science & Engineering* 2007;9(3):90-95. doi: 10.1109/MCSE.2007.55
51. Consortium TM. Project MONAI. Zenodo 2020. doi: <https://doi.org/10.5281/zenodo.4323059>
52. Brett, Matthew, Markiewicz, Hanke C. nipy/nibabel: 3.2.1 (Version 3.2.1). Zenodo 2020. doi: <http://doi.org/10.5281/zenodo.4295521>
53. contributors N. nilearn. <https://github.com/nilearn/nilearn>. Published 2024.
54. Harris CR, Millman KJ, van der Walt SJ, Gommers R, Virtanen P, Cournapeau D, Wieser E, Taylor J, Berg S, Smith NJ, Kern R, Picus M, Hoyer S, van Kerkwijk MH, Brett M, Haldane A, del Río JF, Wiebe M, Peterson P, Gérard-Marchant P, Sheppard K, Reddy T, Weckesser W, Abbasi H, Gohlke C, Oliphant TE. Array programming with NumPy. *Nature* 2020;585(7825):357-362. doi: 10.1038/s41586-020-2649-2
55. Reback J, McKinney W, jbrockmendel. pandas-dev/pandas: Pandas 1.0.3 (Version v1.0.3). Zenodo 2020. doi: <http://doi.org/10.5281/zenodo.3715232>
56. Paszke A, Gross S, Massa F, Lerer A, Bradbury J, Chanan G, Killeen T, Lin Z, Gimelshein N, Antiga L, Demsmaison A, Köpf A, Yang E, DeVito Z, Raison M, Tejani A, Chilamkurthy S, Steiner B, Fang L, Bai J, Chintala S. PyTorch: An Imperative Style, High-Performance Deep Learning Library. *NeurIPS* 2019.

57. Waskom M, Seaborn_Development_Team. seaborn. Zenodo 2020. doi: 10.5281/zenodo.592845
58. Piorkowska M, Goh V, Booth TC. Post Brexit: challenges and opportunities for radiology beyond the European Union. *Br J Radiol* 2017;90(1072):20160852. doi: 10.1259/bjr.20160852
59. Ruffle JK, Farmer AD, Aziz Q. Artificial Intelligence-Assisted Gastroenterology—Promises and Pitfalls. *The American journal of gastroenterology* 2018. doi: 10.1038/s41395-018-0268-4
60. Ruffle JK, Gray RJ, Mohinta S, Pombo G, Kaul C, Hyare H, Rees G, Nachev P. Computational limits to the legibility of the imaged human brain. *NeuroImage* 2024:120600. doi: <https://doi.org/10.1016/j.neuroimage.2024.120600>
61. Smith-Bindman R, Miglioretti DL, Larson EB. Rising use of diagnostic medical imaging in a large integrated health system. *Health Aff (Millwood)* 2008;27(6):1491-1502. doi: 10.1377/hlthaff.27.6.1491