

# **Automated longitudinal treatment response assessment of brain tumors: a systematic review**

Tangqi Shi<sup>1</sup>, Aaron Kujawa<sup>1</sup>, Christian Linares<sup>1</sup>, Tom Vercauteren<sup>1</sup>, Thomas C Booth<sup>1\*</sup>

<sup>1</sup>School of Biomedical Engineering & Imaging Sciences, King's College London, London, UK.

\*Corresponding author. Thomas C Booth, FRCR PhD, School of Biomedical Engineering & Imaging Sciences, King's College London, Becket House, 1 Lambeth Palace Rd, South Bank, London, UK; E-mail(s): [thomas.booth@kcl.ac.uk](mailto:thomas.booth@kcl.ac.uk); Tel: 02032994865.

Accepted Manuscript

## Abstract

**Background:** Longitudinal assessment of tumor burden using imaging helps to determine whether there has been a response to treatment both in trial and real-world settings. From a patient and clinical trial perspective alike, the time to develop disease progression, or progression-free survival, is an important endpoint. However, manual longitudinal response assessment is time-consuming and subject to interobserver variability. Automated response assessment techniques based on machine learning (ML) promise to enhance accuracy and reduce reliance on manual measurement. This paper evaluates the quality and performance accuracy of recently published studies.

**Methods:** Following PRISMA guidelines and the CLAIM checklist, we searched PUBMED, EMBASE, and Web of Science for articles (January 2010-November 2024). Our PROSPERO-registered study (CRD42024496126) focused on adult brain tumor automated treatment response assessment studies using ML methodologies. We determined the extent of development and validation of the tools and employed QUADAS-2 for study appraisal.

**Results:** Twenty (including seventeen retrospective and three prospective) studies were included. Data extracted included information on the dataset, automated response assessment including pertinent steps within the pipeline (index tests), and reference standards. Only limited conclusions are appropriate given the high bias risk and applicability concerns (particularly regarding reference standards and patient selection), and the low-level evidence. There was insufficient homogenous data for meta-analysis.

**Conclusion:** The study highlights the potential of ML to improve brain tumor longitudinal treatment response assessment. Interpretation is limited due to study bias and limited evidence of generalizability. Prospective studies with external datasets validating the latest neuro-oncology criteria are now required.

**Key Words:** Brain Tumor, Machine Learning, Treatment Response

## **Key Points:**

1. The systematic review emphasizes the role of machine learning in enhancing precision and efficiency in neuro-oncology longitudinal assessments.
2. The review highlights the necessity for further research to address biases and enhance clinical applicability.

## **Importance of the Study**

We present the first systematic review that evaluates machine learning (ML) applications for the longitudinal treatment response assessment of brain tumors. Such technologies have the potential to improve neuro-oncological practice, offering a more precise, consistent, and efficient approach to treatment monitoring in both the clinic and during trials.

We highlight the need for addressing bias risks in the development of automated ML methods. Despite the potential of ML to improve segmentation accuracy and efficiency, systematic errors appear to be common when the enhancing tumor region is measured. From this published work, automated tools do not appear clinic-ready, and further research, especially incorporating external test datasets and prospective datasets, is now needed for more robust validation. Successful demonstration of tool use in the clinic or in clinical trials is also now required to complete clinical validation.

Accepted Manuscript

## 1 Introduction

Brain tumors present significant clinical challenges. For example, due to their infiltrative nature, diffuse gliomas typically have a very poor prognosis with the most common type, glioblastoma, having a median overall survival of only 14.6 months despite standard of care treatment (which typically consists of maximal safe tumor resection, followed by radiotherapy with concomitant and adjuvant temozolomide chemotherapy)<sup>1</sup>. The two-year survival rate is around 30%. Similarly, the presence of brain metastases, which occur in approximately 10% to 20% of adult cancer patients<sup>2</sup>, also represents a challenging clinical scenario due to the blood brain barrier influencing systemic therapeutic delivery. Metastatic invasion therefore complicates treatment decisions and is often associated with a median survival of just a few months. For example, patients with multiple brain metastases treated with whole-brain radiotherapy alone have a median survival of about 3-6 months<sup>3</sup>. To help navigate brain tumor patient management after the initiation of treatment, response assessment using longitudinal imaging has become the clinical standard of care. Regularly scheduled imaging facilitates tracking tumor biology and assessing treatment efficacy, which are important factors influencing decision-making during multidisciplinary team meetings (MDTM or Tumor Boards). The rationale is that disease progression may be identified before clinical symptoms emerge and that may lead to an early intervention - which may plausibly improve therapeutic outcomes and prevent irreversible complications<sup>4,5</sup>.

Longitudinal imaging forms the basis of reference standards for response assessment in clinical trials. In such a research setting, RANO (Response Assessment in Neuro-Oncology) criteria<sup>6-9</sup>, have become essential by providing a standardized approach for assessing the effectiveness of treatments for brain tumors (Box 1). It is important to note that the US Food and Drug Administration (FDA) has endorsed treatment outcomes based on RANO criteria<sup>10</sup>, which ensures that they meet the rigorous standards necessary for regulatory approval in clinical trials. The RANO criteria not only consider changes in tumor size and morphology but also include the patient's clinical presentation and neurological functional status. Standardized clinico-radiological response assessment criteria not only allow comparison of outcomes during trials, but also during routine clinical assessment when applied in an expedient and simplified form to help clinicians to quickly make reliable treatment decisions given the complexities of interpreting MRI data<sup>11</sup>.

**Box 1:** Key terminology relevant for assessing autonomous treatment response assessment studies

Reference Standard: The “reference standard” refers to the best available method to determine the accuracy of diagnostic assessments, establishing a benchmark for evaluating new methods<sup>12</sup>. Here, radiologists’ (1) manual image segmentation and (2) manual tumor assessment — using expert measurement of the Response Assessment in Neuro-Oncology (RANO) criteria — serve as the reference standards.

Response Assessment in Neuro-Oncology (RANO) criteria: The RANO criteria serve as a standardized set of guidelines for evaluating the effectiveness of brain tumor treatments in clinical trials. These criteria were developed to address limitations in previous assessment methods such as the MacDonald criteria<sup>13</sup>. RANO assessment focuses on changes in tumor size (typically using the product of bidimensional perpendicular diameters), measured by T1-weighted post contrast and T2/FLAIR MRI sequences. RANO assessment also incorporates clinical factors (e.g., corticosteroid use and neurological symptoms) alongside imaging. Beyond clinical trials, assessments in routine clinical practice may also use RANO or may largely be based on RANO criteria<sup>11</sup>. First designed for high-grade glioma<sup>6</sup>, updates and extensions of the RANO criteria have been proposed including for specific tumor types (e.g., low-grade gliomas, metastases, meningiomas) and advanced therapies (e.g., immunotherapies). Tumor response can be categorized as progression, stable disease, partial response, or complete response, and the criteria are defined.

Index test: The “index test” refers to the new diagnostic test or assessment method under investigation<sup>12</sup>, which in this review is the automated ML-based assessment, encompassing both segmentation and tumor response evaluation.

Manual longitudinal assessments based on structural MRI protocols can be problematic due to several factors

. High-grade gliomas, for example, exhibit a variety of shapes, and their boundaries can be difficult to precisely define. Moreover, the solid tumor often manifests as a cavity rim, making it challenging to capture the full extent accurately. Indeed, in some cases, large cyst-like high-grade gliomas may not meet the "measurable" criteria unless a solid peripheral nodular component of sufficient size ( $\geq 10$  mm) is present. These complexities highlight the limitations of manual assessments, as they rely on subjective interpretation and can result in inaccuracies in tumor

measurement and monitoring, as well as the need for more standardized and objective approaches. Manual assessments are also resource intensive. In response to these challenges posed by manual assessments in structural MRI protocols there has been a notable advancement in the development of automated assessment tools. These tools, utilizing various machine learning (ML) methodologies (Box 2)<sup>14,15</sup>, aim to automate – and optimize - the longitudinal assessment of treatment response in brain gliomas. In particular, these automated systems are designed to address the limitations of manual assessments by offering more accurate, reproducible and efficient methods for evaluating treatment responses and tumor metrics. In this systematic review, we aimed to analyze and summarize the diagnostic accuracy of current ML algorithms used for longitudinal treatment response assessment. Whilst our primary objective was to examine the overall automated treatment response assessment based on ML, a secondary objective was to investigate the underlying automated tumor segmentation.

**Box 2:** Overview of methods in artificial intelligence.

Artificial Intelligence (AI) encompasses a wide array of computational techniques aimed at enabling machines to mimic human intelligence. Within AI, machine learning (ML) represents a subset of algorithms that learn complex patterns from data without explicit programming for each specific outcome, in order to produce analytical models that can make predictions. Neural networks are a key ML approach inspired by the human brain's structure and use interconnected nodes (like neurons) to process data through sequential layers which perform the pattern recognition process. Deep learning (DL), a subset of ML, uses neural network architectures with multiple layers such as convolutional neural networks (CNNs). CNN architectures like U-Net can be used for specialized tasks like image segmentation. In neuro-oncology, CNN models can be applied to MRI scans to segment tumors or to produce diagnostic, prognostic, predictive or monitoring biomarkers<sup>16</sup>. Other deep learning examples in neuro-oncology might use even more advanced neural network architectures like generative adversarial networks (GANs) that can be used to synthesize specific MRI sequences, such as generating contrast-enhanced T1-weighted images of tumors. In summary, AI techniques can be widely applied to clinical decision-making (using a variety of data including images) and image analysis which can enhance efficiency and accuracy of tasks such as tumor segmentation, treatment response assessment, and disease prognosis.

## 2 Method

This systematic review was registered with PROSPERO (CRD42024496126). The review was organized in line with the Preferred Reporting Items for Systematic Reviews and Meta-Analysis (PRISMA)<sup>17</sup>. Where appropriate, both Quality Assessment of Diagnostic Accuracy Studies 2 (QUADAS-2) methodology<sup>18</sup> alongside ML metrics from the Checklist for Artificial Intelligence in Medical Imaging (CLAIM) were used to assess the risk of bias for each study<sup>19</sup>.

### Search Strategy and Selection Criteria

Search terms were applied to PUBMED, EMBASE, and Web of Science databases to extract original research articles published from January 2010 to November 2024 encompassing both text words and database-specific subject headings (Supplementary Table S1). Specifically, we used Medical Subject Headings (MeSH) for PUBMED; Emtree subject heading terms for EMBASE; and a broader combination of keywords for Web of Science. For example, our PUBMED search strategy was as follows: we used the keywords combination ('automated' or 'automatic' or 'pipeline' or 'AI' or 'artificial intelligence' or 'ML' or 'machine learning' or 'deep learning' or 'radiomics') AND ('brain tumor' or 'brain metastases' or 'glioma' or 'glioblastoma') AND ('longitudinal' or 'follow-up' or 'follow up' or 'treatment response' or 'monitoring biomarker' or 'response assessment' or 'monitoring') AND ('MRI' or 'magnetic resonance imaging' or 'magnetic resonance' or 'MR'), including both English and Chinese language publications, with a date range from January 1, 2010, to November 25, 2024. The 2010 starting point was chosen because it coincided with the publication of the first RANO paper<sup>6</sup>. Preprints, abstracts, reviews, editorials, reports, letters, book chapters, case reports, symposiums, retracted papers or articles without peer review were excluded (Fig. 1). The references of all selected articles were hand-searched to identify any potentially relevant studies missed in the initial database search.

Studies were deemed eligible for inclusion if they assessed the longitudinal treatment response of any type of brain tumor (i.e., both benign and malignant tumors) using ML methodologies. In the context of this study, ML methodologies are those that enable computers to learn from retrospective data, by automatically tuning algorithms

that are not solely composed of explicit instructions, and autonomously make decisions or predictions for prospective use. The definition includes deep learning (DL) which is a type of ML based on artificial neural networks in which multiple layers of processing are used to extract progressively higher-level features from data. The patient cohort was restricted to adult patients ( $\geq 18$  years old) who underwent standardized treatment and subsequent imaging to evaluate treatment outcomes. Ineligible studies included those reporting on only pediatric populations; those without clinical experiments; and those using only animal models. We also excluded all studies without longitudinal analysis, for example prognostic, diagnostic or monitoring biomarkers using a single timepoint to determine treatment response by radiomic analysis. Research that only compared preoperative and early postoperative (i.e.,  $< 72$  hours) imaging timepoints cannot be considered as providing longitudinal assessments of treatment response and were thus excluded.

A meta-analysis could not be performed due to a lack of sufficient homogenous studies identified from the systematic review and marked heterogeneity in the methodology of these included studies.

### **Data Extraction**

A neuroimaging data scientist, T.S., with 2 years of experience in neuroimaging applied to neuro-oncology, independently performed the data extraction and quality assessment. A.K., a neuroimaging data scientist with 8 years of experience in neuroimaging applied to neuro-oncology also independently performed the data extraction. A junior neuro-oncologist (UK specialist trainee grade; US fellow equivalent). C.A.L., with 5 years of experience in neuroimaging applied to neuro-oncology, also independently performed the quality assessment. Discrepancies between the reviewers were considered at research meetings with a senior neuroradiologist (UK consultant; US attending equivalent) T.C.B. with 17 years' experience of neuroimaging applied to neuro-oncology, and T.V., a neuroimaging data scientist with 15 years of experience in neuroimaging applied to neuro-oncology, until a consensus was reached.

Data extracted from each study included the type and grade of brain tumor, and whether classification was based on WHO 2016 or 2021 criteria<sup>20,21</sup>; whether data was obtained from single or multiple institutions; the size of training and testing sets; the types of MRI sequences used; and whether the study was retrospective or prospective. We also analyzed the automated models for longitudinal assessment (assigned as index tests) and any automated sub-components within the pipeline, as well as the reference standard applied (e.g., RANO 2010). Performance metrics for sub-components prior to longitudinal assessment e.g., segmentation, were collected alongside longitudinal assessment performance metrics. Performance metrics extracted were based on index test results compared to reference standard results. Depending on the task, metrics extracted included, for example, Dice coefficient, intraclass correlation coefficient (ICC) or area under the receiver operating characteristic curve (AUC) values.

### **Risk of Bias Assessment**

To evaluate diagnostic accuracy, we applied QUADAS-2 (Quality Assessment of Diagnostic Accuracy Studies 2) methodology<sup>18</sup>, a tool specifically designed for the systematic assessment of the quality of diagnostic accuracy studies. This analytical framework facilitates the appraisal of risk of bias and applicability concerns across four key domains: patient selection, index tests, reference standards, and flow and timing. A three-tiered rating system comprising 'low', 'high', or 'unclear' risk of bias or applicability concerns, was employed.

Each domain was systematically evaluated with carefully prepared criteria. For patient selection, we appraised key aspects such as whether participants were enrolled in the study consecutively or at random, followed a standard treatment protocol, and were subject to any inappropriate exclusions. We also evaluated studies as to whether all appropriate exclusions had been applied. For the index test and reference standard domains, we appraised whether blinding was employed during respective formulation or evaluation, and whether relevant clinical confounding characteristics had been accounted for. We also confirmed whether all participants were included in the analysis, and whether the same reference standard had been applied to them to ensure uniformity. We also determined whether pre-specified thresholds prior to index testing had been fixed a priori.

### 3 Results

In all, 2088 citations fulfilled the search criteria of which the full text of 762 potentially eligible articles were reviewed (Fig. 1). A total of 20 studies were included in the final analysis and are listed in Table 1. <sup>(22-41)</sup>

#### Study characteristics

As shown in Table 1, in terms of dataset utilization, 60.0% (12/20) of studies included grade 2-4 glioma, with 66.7% (8/12) of studies focusing solely on glioblastoma (studies used grade IV glioma WHO 2016 definition). The remainder, 40.0% (8/20) of studies, focused on brain metastases.

Studies differed as to whether they were either validation-only studies, or combined development and validation studies. Of 50.0% (10/20) studies that trained ML models and conducted tests with hold out data, all (100% 10/10) studies utilized internal test data, and 40.0% (4/10) studies additionally employing external test data. The remaining study 5.0% (1/20) only used external test data, and 45.0% (9/20) did not involve model training (i.e., the latter were validation-only studies where the ML model had been developed in a previous study).

There were 55.0% (11/20) studies using multi-institutional data. Few studies incorporated prospective data (15.0%, 3/20). 60.0% (12/20) of studies used post-operative data only for training or testing, while 40.0% (8/20) of studies trained segmentation models using predominantly pre-operative data or combined with a little post-operative data, and then applied the trained model to post-operative testing data. Such a strategy is expedient because pre-operative datasets are arguably more abundant and accessible than post-operative datasets; pre-operative datasets were typically from the Brain Tumor Segmentation (BraTS) Challenge<sup>42</sup>. Despite the benefits of expediency, while current results of test data might be acceptable, the domain gap of such an approach leads to risks of poor generalizability.

Most studies (70.0%, 14/20) employed more than one MRI sequence. Structural imaging sequences (T1-weighted (T1), T1-weighted post contrast (T1 C), T2-weighted (T2), and FLAIR) were predominantly used, with T1 C applied in 90.0% (18/20) of the studies, in keeping with its pivotal role in therapeutic response assessment<sup>6,43</sup>. Some studies 10.0% (2/20) also used diffusion-weighted imaging (DWI) or apparent diffusion coefficient (ADC) maps<sup>44</sup>.

Regarding implementation methodologies, whilst our inclusion criteria required employing some form of ML-based automation, the extent of automation varied with 80.0% (16/20) of studies achieving complete automation from start to finish without manual intervention. The remainder (20.0%, 4/20) used semi-automatic frameworks with manual intervention.

Most studies (85.0%, 17/20) focused on advanced ML algorithms by utilizing deep neural networks for learning and inference, such as employing a 3D U-Net based model<sup>22,23</sup> for segmentation assessment. The remainder, 15.0% (3/20), did not use deep learning.

In longitudinal assessment, most studies (75.0%, 15/20) employed RANO-based criteria, calculating the diameter and/or volume change, as the assessment method for evaluating therapeutic response. Amongst these, 40.0% (6/15) were based on brain metastases RANO (RANO-BM)<sup>7</sup>, 33.3% (5/15) were based on high-grade glioma RANO (criteria from 2010)<sup>6</sup>, and 26.7% (4/15) on modified high-grade glioma RANO (criteria from 2017 where a key difference is that the baseline MRI is the one performed soon after radiotherapy completion)<sup>8</sup>. In terms of measurements, 66.7% (10/15) studies utilized diameter measurement methods, with 30.0% (3/10) employing 'RANO 2010'<sup>6</sup>, 30.0% (3/10) employing 'modified RANO 2017'<sup>8</sup>, and 40.0% (4/10) employing RANO-BM<sup>7</sup> criteria. Following RANO 2010 or modified RANO 2017 criteria<sup>6,8</sup>, a few of these studies additionally used volume measurement methods (20.0%, 3/15) whilst the rest (26.7%, 4/15) only employed volume measurement assessments. In the remaining 25.0% (5/20) of studies where RANO-based criteria were not employed, evaluation metrics were solely based on volume change in 60.0% (3/5), while 40.0% (2/5) of studies incorporated both volume and diameter

changes for assessment. We also note that 40.0% (8/20) of studies were based on newly diagnosed patients and 30.0% (6/20) on recurrent patients, with the remaining 30.0% (6/20) unknown.

### **QUADAS-2 assessments of the included studies**

The results of the QUADAS-2 stratified analysis<sup>18</sup> of both risk of bias and applicability across four domains, (patient selection, index test(s), reference standard, and flow and timing), is presented in Figure 2, Supplementary Table S2 and Supplementary Table S4. In terms of risk of bias and concerns regarding applicability, only 15.0% (3/20), 15.0% (3/20), 20.0% (4/20) and 60.0% (12/20) are considered at 'low' risk in the domains of patient selection, index test(s), reference standard, and flow and timing, respectively. The corollary is that there is either an 'unclear' or 'high' risk of bias and applicability concerns in most studies.

When focusing on particular aspects of bias, a more nuanced picture emerges. In the patient selection domain, we observed that the majority of studies (50.0%, 10/20) had included patients undergoing a clearly stated treatment protocol (e.g., Stupp protocol for glioblastoma)<sup>1,45</sup>; 45.0% (9/20) studies explicitly stated that patients were enrolled as either a consecutive or random sample; and half the studies (60.0%, 12/20) applied appropriate exclusions (e.g., the studies excluded cases without documented original histology or incomplete or poor-quality imaging data). Similarly, several components of good study design reducing bias risk in the index tests and reference standard domains were also evident in many studies. First, researchers had been blinded to the reference standard performance when considering index test performance in most studies (90.0%, 18/20) – and vice versa in (80.0%, 16/20) studies. The blinding ensured the reliability of the outcome measurement and reduced the potential for verification bias. Second, it was clearly stated that trained model parameters had been fixed during the testing process in almost all studies (85.0%, 17/20). Third, in all experiments, at least one senior radiologist was involved in manual annotation.

## 4 Discussion

### Summary of findings

Automated longitudinal treatment response assessment of brain tumors has been achieved for brain metastases, as well as both high- and low-grade gliomas. The prevailing approach remains the development of ML-based methods emulating RANO criteria<sup>6-8</sup>, with glioblastoma being the commonest tumor to be assessed. When ML-based methods were employed as the index test and compared to the reference standard of expert manual assessments, the performance accuracy was generally good. However, there was a high or unclear risk of bias within most studies due to incomplete published information and a lack of rigor in experimental design, which constrains the widespread applicability of the automated systems. Similarly, there was no clear evidence indicating that the automated systems could be applied in clinical settings. Despite the quality assessment findings, and despite the fact that the studies are generally of a low level of evidence<sup>46</sup>, there is value in interrogating these individual studies as they represent the current state of the art and form a baseline for further research.

### Study explanations and relevance from a national and international perspective

This is the first systematic review of automated longitudinal treatment response assessment studies for brain tumors. It has shown that the neuro-oncology imaging research community have leveraged the ability to obtain, process and store digital images, harnessed the improved performance of registration and segmentation tools -and taken together – have built automated treatment response assessment tools. However, whilst analytical validation<sup>5</sup> has been demonstrated to be technically possible for a range of brain tumors, almost all current studies are compromised by bias and are best considered as proof-of-concept studies. One recent multi-reader validation study has largely avoided bias<sup>24</sup>, but the study design does not constitute comprehensive clinical validation<sup>5</sup> (Box 3). Therefore, from published evidence, no tool is definitively ready for clinical use and more research is required to ensure this. Being clinic-ready is important because once performance accuracy is satisfactory in providing treatment response assessment during either a clinical trial investigating therapeutics or during routine patient follow up, there is a high likelihood of benefit for both the patient and healthcare system. The key potential benefit is that a clinically

validated tool will reduce or eradicate interobserver error of treatment response assessment by giving a more reproducible and standardized result, and that the tool will reduce the burden of time and costs spent on clinical trials. Beyond trials, it is also plausible that a clinically validated tool will improve consistency in routine patient follow up in the clinic and therefore allow more rational management decisions.

### **Limitations - studies assessed**

Whilst the studies assessed show numerous strengths they are not without their limitations. First, few studies employed external testing limiting the analytical validation process<sup>5</sup>, so it is unclear whether the tools are generalizable and therefore applicable for further use at other sites. Second, few studies used prospective data, and none employed tools embedded in the clinical or trial workflow, therefore clinical validation<sup>5</sup> steps are still required. Third, given that RANO 2.0 was only published in 2023, no studies using ML-based automation have yet conducted experiments based specifically on the updated criteria<sup>9</sup>. Fourth, there is a clear need for more studies to adopt fully automated measurement techniques (diameters and volume changes calculation) to enhance the accuracy and consistency of RANO assessments; at least a quarter are semi-automated. Fifth, RANO criteria<sup>6-9</sup>, and the Macdonald criteria<sup>13</sup> that preceded them, replaced the previous World Health Organization (WHO) recommendations<sup>47</sup> which considered all brain tumors as solid entities during measurement. The Macdonald criteria allowed assessment of tumor within cavity walls but did not distinguish between the presence of necrosis or a surgical resection in a ‘cyst-like’ cavity. The RANO criteria indicated that any cyst-like cavity should not be measured. However, numerous studies – including automated longitudinal treatment response assessment tools purporting to follow RANO criteria – often demonstrate the inclusion of some or all cyst-like cavities during tumor measurement, e.g.,<sup>22-24,34,36</sup>. It is conceivable that global segmentation competitions where concepts such as ‘entire tumor core’ which have included cyst-like cavities by definition, have disproportionately influenced the current models. In the application of RANO criteria<sup>6-8</sup> as a reference standard, a bidimensional and volumetric measurement systematic error will reduce the accuracy of both the reference standard and the index test<sup>6-8</sup>. Reproducibility is likely to be impacted too as reference standard systematic errors will likely vary between sites. Sixth, segmentation competitions such as BraTS<sup>42</sup> have almost always included pre-operative datasets alone; this poses a challenge for developing generalizable models for the use case of longitudinal treatment response which needs at least some post-

operative datasets. Incorporating more post-operative data in these challenges – as has occurred in BraTS 2024<sup>42,48</sup> is meaningful as it better aligns with the assessment of treatment response in clinical practice and future research. Seventh, few studies explicitly stated that patients were enrolled as either a consecutive or random sample, and in approximately half the studies, there was no evidence to demonstrate that all appropriate exclusions were applied. This lack of clarity suggests a potential bias in patient selection, which could limit the generalizability of the research findings. Eighth, when designing the index test, few studies (15.0%, 3/20) explicitly stated that they considered relevant clinical characteristics important for final longitudinal treatment response assessment. When considering using RANO criteria as a reference standard in glioma, for example, factors such as a change in performance status, a change in the use of corticosteroids, and the start of second-line treatment are essential for final longitudinal treatment response assessment as RANO is a clinico-radiological assessment<sup>49</sup>. Ninth, there was a mismatch between the proposed RANO study scheme and the actual baseline used in the study as the modified RANO baseline is the first scan after radiotherapy, and the RANO 2010 baseline is after surgery but before radiotherapy<sup>6,8</sup>.

### **Limitations - review process**

This paper is the first systematic review of automated pipelines that assess brain tumor treatment response using ML. However, the review still has some limitations. First, our review did not include those ML studies that only provide theoretical foundations or preliminary data without experiments<sup>50</sup>. While these studies do not involve direct clinical experiments, it is plausible that they are potentially valuable for the development of longitudinal treatment response assessment tools. Second, we also excluded those automated pipelines without ML even if they were able to assess automated longitudinal treatment response, as they were beyond the remit of the systematic review. Nonetheless for comparison, some important studies that utilize automatic algorithms, such as those using the region-growing algorithm to achieve automation<sup>51-53</sup>, are shown in Supplementary Table S3. Third, studies were excluded which focused on the development of prognostic biomarkers for overall survival (OS) based on radiographic feature changes<sup>54</sup>. However, there is some overlap in longitudinal research methodology which might be relevant to treatment response assessment. Fourth, publication bias may also have affected the range of automated pipelines of treatment response included in this systematic review. Related to this, the exclusion of pre-

prints and non-peer reviewed material may exacerbate publication bias. In particular, given that some in the data science community may not submit their work in peer-reviewed journals, as peer review is relatively slow compared to the speed at which data science develops, it is plausible that publication bias relates to the make-up of the researcher team<sup>12</sup>. For example, more clinically orientated teams may be more inclined to publish in a peer reviewed journal compared to more data science-orientated teams who sometimes use pre-prints alone or full-length conference proceedings<sup>12</sup>. Fifth, our search strategy may not have captured all relevant studies.

### **Current Evidence in the Field**

This is the first systematic review of automated longitudinal treatment response assessment studies for brain tumors. The response evaluation criteria in solid tumors (RECIST) is a widely-used reference standard for evaluating efficacy of therapies in patients with solid tumors which are included in clinical trials and it is widely used and accepted by regulatory agencies<sup>55</sup>. Similar to brain tumor response assessment using automated RANO assessments, automation of RECIST is desirable as it can potentially streamline the process and potentially reduce the variability of results within the RECIST<sup>55</sup>. However, there remain technical challenges which must be overcome to ensure reproducibility, and currently there are no clinic-ready automated RECIST studies to the best of our knowledge. The current evidence suggests that RANO<sup>6-8</sup> appears more likely to achieve full automation with fewer remaining challenges compared to RECIST<sup>55</sup>.

### **Implications for clinical practice and future research**

These automated tools can enhance the overall efficiency of tumor treatment response assessment and reduce interobserver variability. Whilst the main intention for RANO assessments - including when using automated tools - is to produce measurable, standardized, and meaningful outcomes for clinical trials, approximating RANO assessment into routine real-world follow up assessment may also help clinicians make more reproducible treatment decisions given the complexities of interpreting MRI data<sup>11</sup>. If further developed and validated, it is plausible that automation might overcome the time-consuming process preventing RANO-like assessments being routinely

available in the clinic. Considering the directions for future developments, attention should be focused on the main aspects detailed hereafter.

First, new automated tool developments incorporating the requirements of criteria like RANO 2.0 (which incorporate bidimensional diameters and volumes, use different contrasts, and allow assessment of both high-grade gliomas and low-grade gliomas<sup>9</sup>), are likely to support future clinical trials and be more translatable to the clinic. Second, adherence to the incorporation of tumor tissue demarcated by T1 C, as opposed to including voluminous cyst-like regions, is needed in ML models to keep RANO assessments as they were intended during inception over a decade ago. Third, to enhance utility, given that treatment response assessment of brain tumors is a clinico-radiological assessment, tool interfaces would benefit from having the option to integrate clinical information such as changes in performance status or steroid use and/or contain user warnings to not confound assessment by not considering confounders such as early second-line treatment<sup>12</sup> which is a common concern for RANO assessments<sup>8,9</sup>. Fourth, consideration should also be given to distinguishing target lesions, measurable lesions, and whether lesions are inside or outside the radiotherapy field – all of which are requisites for the assessment of longitudinal treatment responses in RANO<sup>9</sup>. Fifth, it should be noted that there is ongoing controversy regarding the optimal approach to measuring brain tumors over time. For glioblastomas, it remains unclear whether measurements should focus solely on the enhancing portion or also include non-enhancing FLAIR signal abnormalities. Similarly, for IDH-mutant gliomas, standardized measurement strategies are not yet well-defined. Additionally, as highlighted by RANO 2.0 criteria<sup>9</sup>, volumetric measurements have not been shown to be unequivocally superior to orthogonal diameters. Furthermore, clinical trials seeking to modify the Stupp protocol may influence the selection of a "post-treatment baseline," but the implications of these modifications remain uncertain. In summary, the RANO criteria, as articulated from inception in 2010, are best considered as works in progress and will continue to evolve.<sup>6</sup>

In terms of evidence generation, there is a need for more validation studies containing external test datasets and prospective datasets to demonstrate that automated tools are ready for downstream clinical requirements. When selecting patient cohorts, it is advisable to consider using consecutive or random samples, follow a standard

treatment protocol, avoid inappropriate exclusions (i.e., cherry-picking), and apply all necessary exclusions (e.g., lack of histology, incomplete or poor-quality imaging data).

It is acknowledged that a focus on RANO assessment, especially for high-grade gliomas, fails to consider approaches other than fixed interval imaging where utility is unclear<sup>4,11</sup>. It is also acknowledged that a focus on RANO assessment, especially for high-grade gliomas, and the use of T1C images alone in many studies, fails to use all the data produced during an MRI scan. Nonetheless, some studies included in the current review were not constrained to T1C only, potentially allowing for a more nuanced understanding of tumor behavior and response to treatment across different imaging and clinical scenarios. For example, in high-grade gliomas DWI/ADC is a surrogate marker of tumor cell density and has been used to assess aggressiveness, while T2 contrast is effective in displaying broader tumor-host behavior<sup>56-59</sup>. Future research may go beyond the limitations of RANO assessment and even incorporate advanced MRI as well as multi-modal techniques such as a combination of MRI and positron emission tomography<sup>60</sup>.

Similarly, treatment response assessment may be improved through another key area of research in neuro-oncology imaging which is the development of advanced MRI and radiomic prognostic, predictive and monitoring biomarkers which have a role in treatment response assessment<sup>56</sup>. Integrating radiomic biomarkers with automated longitudinal treatment response assessment frameworks offers a promising avenue for enhancing precision in tracking tumor responses across various brain tumor types<sup>54,61</sup>.

#### Box 3: Comprehensive clinical validation

Comprehensive clinical validation would require not only embedding tools within the clinical or trial workflow but also demonstrating robust real-world utility across diverse patient populations and clinical settings. This also includes validation datasets from multiple institutions to ensure generalizability; testing with varied imaging protocols and scanner vendors; and prospective trials that assess clinical outcomes when using these tools in

decision-making.

## 5 Conclusion

The systematic review demonstrates the potential of automated tools to enhance the accuracy and reliability of treatment response assessments in brain tumors. Studies achieving complete automation from start to finish without manual intervention will contribute to the consistency and efficiency of data processing, likely minimizing the potential for human error and workload. However, automated tools are not clinic-ready and further research, especially incorporating external test datasets and prospective datasets, is now needed for more robust validation. Successful demonstration of tool use in the clinic or in clinical trials is also now required to complete clinical validation.

## Ethics

This study did not involve experimental investigations on human or animal subjects.

## Funding

T.C.B. is supported by the UK Medical Research Council (MR/W021684/1). This work was also supported by the Wellcome EPSRC Centre for Medical Engineering at King's College London (203148/Z/16/Z) (including authors T.V. and T.C.B.). T.S. is supported by a K-CSC scholarship.

## Conflict of Interest

There is no conflict of interest for all authors. T.S. - none declared. A.K. - none declared. C.A.L. - none declared. T.C.B. - consultancy: Microvention; payment/honoraria for education lectures: Siemens Healthineers Speakers Bureau, Medtronic Speakers Bureau; support for attending meetings/ travel: Balt, whose activity is not related to the

present article. T.V. - co-founder and shareholder of Hypervision Surgical whose activity is not related to the present article.

## **Authorship**

Conception and design: T.S., T.C.B., T.V. Data acquisition and preparation, data analysis, manuscript drafting, data interpretation, critical review of the work and manuscript, final approval of manuscript, accountability for all aspects of the work: all authors.

## **Data Availability**

Data generated or analyzed during the study are available from the corresponding author by request.

Accepted Manuscript

## References

1. Stupp R, Mason WP, van den Bent MJ, et al. Radiotherapy plus concomitant and adjuvant temozolomide for glioblastoma. *N Engl J Med*. 2005; 352(10):987-996.
2. Lin X, DeAngelis LM. Treatment of Brain Metastases. *J Clin Oncol*. 2015; 33(30):3475-3484.
3. Soffiatti R, Abacioglu U, Baumert B, et al. Diagnosis and treatment of brain metastases from solid tumors: guidelines from the European Association of Neuro-Oncology (EANO). *Neuro Oncol*. 2017; 19(2):162-174.
4. Booth TC, Thompson G, Bulbeck H, et al. A Position Statement on the Utility of Interval Imaging in Standard of Care Brain Tumour Management: Defining the Evidence Gap and Opportunities for Future Research. *Front Oncol*. 2021; 11:620070.
5. Cagney DN, Sul J, Huang RY, Ligon KL, Wen PY, Alexander BM. The FDA NIH Biomarkers, EndpointS, and other Tools (BEST) resource in neuro-oncology. *Neuro Oncol*. 2018; 20(9):1162-1172.
6. Wen PY, Macdonald DR, Reardon DA, et al. Updated response assessment criteria for high-grade gliomas: response assessment in neuro-oncology working group. *J Clin Oncol*. 2010; 28(11):1963-1972.
7. Lin NU, Lee EQ, Aoyama H, et al. Response assessment criteria for brain metastases: proposal from the RANO group. *Lancet Oncol*. 2015; 16(6):e270-278.
8. Ellingson BM, Wen PY, Cloughesy TF. Modified Criteria for Radiographic Response Assessment in Glioblastoma Clinical Trials. *Neurotherapeutics*. 2017; 14(2):307-320.
9. Wen PY, van den Bent M, Youssef G, et al. RANO 2.0: Update to the Response Assessment in Neuro-Oncology Criteria for High- and Low-Grade Gliomas in Adults. *J Clin Oncol*. 2023; 41(33):5187-5199.
10. Wen PY, Cloughesy TF, Ellingson BM, et al. Report of the Jumpstarting Brain Tumor Drug Development Coalition and FDA clinical trials neuroimaging endpoint workshop (January 30, 2014, Bethesda MD). *Neuro Oncol*. 2014; 16 Suppl 7(Suppl 7):vii36-47.
11. Booth TC, Luis A, Brazil L, et al. Glioblastoma post-operative imaging in neuro-oncology: current UK practice (GIN CUP study). *Eur Radiol*. 2021; 31(5):2933-2943.
12. Booth TC, Grzeda M, Chelliah A, et al. Imaging Biomarkers of Glioblastoma Treatment Response: A Systematic Review and Meta-Analysis of Recent Machine Learning Studies. *Front Oncol*. 2022; 12:799662.
13. Macdonald DR, Cascino TL, Schold SC, Jr., Cairncross JG. Response criteria for phase II studies of supratentorial malignant glioma. *J Clin Oncol*. 1990; 8(7):1277-1280.

14. Xu Y, Liu X, Cao X, et al. Artificial intelligence: A powerful paradigm for scientific research. *Innovation (Camb)*. 2021; 2(4):100179.
15. Rudie JD, Rauschecker AM, Bryan RN, Davatzikos C, Mohan S. Emerging Applications of Artificial Intelligence in Neuro-Oncology. *Radiology*. 2019; 290(3):607-618.
16. Booth TC, Williams M, Luis A, Cardoso J, Ashkan K, Shuaib H. Machine learning and glioma imaging biomarkers. *Clin Radiol*. 2020; 75(1):20-32.
17. Page MJ, McKenzie JE, Bossuyt PM, et al. The PRISMA 2020 statement: An updated guideline for reporting systematic reviews. *Int J Surg*. 2021; 88:105906.
18. Whiting PF, Rutjes AW, Westwood ME, et al. QUADAS-2: a revised tool for the quality assessment of diagnostic accuracy studies. *Ann Intern Med*. 2011; 155(8):529-536.
19. Mongan J, Moy L, Kahn CE, Jr. Checklist for Artificial Intelligence in Medical Imaging (CLAIM): A Guide for Authors and Reviewers. *Radiol Artif Intell*. 2020; 2(2):e200029.
20. Louis DN, Perry A, Wesseling P, et al. The 2021 WHO Classification of Tumors of the Central Nervous System: a summary. *Neuro Oncol*. 2021; 23(8):1231-1251.
21. Louis DN OH, Wiestler OD, Cavenee WK. WHO Classification of Tumours of the Central Nervous System: IARC Publication; 2016.
22. Chang K, Beers AL, Bai HX, et al. Automatic assessment of glioma burden: a deep learning algorithm for fully automated volumetric and bidimensional measurement. *Neuro Oncol*. 2019; 21(11):1412-1422.
23. Nalepa J, Kotowski K, Machura B, et al. Deep learning automates bidimensional and volumetric tumor burden measurement from MRI in pre- and post-operative glioblastoma patients. *Comput Biol Med*. 2023; 154:106603.
24. Vollmuth P, Foltyn M, Huang RY, et al. Artificial intelligence (AI)-based decision support improves reproducibility of tumor response assessment in neuro-oncology: An international multi-reader study. *Neuro Oncol*. 2023; 25(3):533-543.
25. Rudie JD, Calabrese E, Saluja R, et al. Longitudinal Assessment of Posttreatment Diffuse Glioma Tissue Volumes with Three-dimensional Convolutional Neural Networks. *Radiol Artif Intell*. 2022; 4(5):e210243.
26. Strack C, Pomykala KL, Schlemmer HP, Egger J, Kleesiek J. "A net for everyone": fully personalized and unsupervised neural networks trained with longitudinal data from a single patient. *BMC Med Imaging*. 2023;

- 23(1):174.
27. Jalalifar SA, Soliman H, Sahgal A, Sadeghi-Naini A. Automatic Assessment of Stereotactic Radiation Therapy Outcome in Brain Metastasis Using Longitudinal Segmentation on Serial MRI. *IEEE J Biomed Health Inform.* 2023; 27(6):2681-2692.
  28. Kickingereder P, Isensee F, Tursunova I, et al. Automated quantitative tumour response assessment of MRI in neuro-oncology with artificial neural networks: a multicentre, retrospective study. *Lancet Oncol.* 2019; 20(5):728-740.
  29. Chen J, Meng L, Bu C, Zhang C, Wu P. Feature pyramid network-based computer-aided detection and monitoring treatment response of brain metastases on contrast-enhanced MRI. *Clin Radiol.* 2023; 78(11):e808-e814.
  30. Meier R, Knecht U, Loosli T, et al. Clinical Evaluation of a Fully-automatic Segmentation Method for Longitudinal Brain Tumor Volumetry. *Sci Rep.* 2016; 6:23376.
  31. Jayachandran Preetha C, Meredig H, Brugnara G, et al. Deep-learning-based synthesis of post-contrast T1-weighted MRI for tumour response assessment in neuro-oncology: a multicentre, retrospective cohort study. *Lancet Digit Health.* 2021; 3(12):e784-e794.
  32. Cho J, Kim YJ, Sunwoo L, et al. Deep Learning-Based Computer-Aided Detection System for Automated Treatment Response Assessment of Brain Metastases on 3D MRI. *Front Oncol.* 2021; 11:739639.
  33. Hsu DG, Ballangrud A, Prezelski K, et al. Automatically tracking brain metastases after stereotactic radiosurgery. *Phys Imaging Radiat Oncol.* 2023; 27:100452.
  34. Kleesiek J, Petersen J, Doring M, et al. Virtual Raters for Reproducible and Objective Assessments in Radiology. *Sci Rep.* 2016; 6:25007.
  35. Ozkara BB, Federau C, Dagher SA, et al. Correlating volumetric and linear measurements of brain metastases on MRI scans using intelligent automation software: a preliminary study. *J Neurooncol.* 2023; 162(2):363-371.
  36. Suter Y, Notter M, Meier R, et al. Evaluating automated longitudinal tumor measurements for glioblastoma response assessment. *Front Radiol.* 2023; 3:1211859.
  37. Zhang J, LaBella D, Zhang D, et al. Development and Evaluation of Automated Artificial Intelligence-Based Brain Tumor Response Assessment in Patients with Glioblastoma. *AJNR Am J Neuroradiol.* 2024.

38. Prezelski K, Hsu DG, del Balzo L, et al. Artificial-intelligence-driven measurements of brain metastases' response to SRS compare favorably with current manual standards of assessment. *Neuro-Oncol Adv.* 2024; 6(1).
39. Son S, Joo B, Park M, et al. Development of RLK-Unet: a clinically favorable deep learning algorithm for brain metastasis detection and treatment response assessment. *Front Oncol.* 2023; 13:1273013.
40. Kotowski K, Machura B, Nalepa J. Robustifying Automatic Assessment of Brain Tumor Progression from MRI. *Lect Notes Comput Sc.* 2023; 13769:90-101.
41. Hammer Y, Najjar W, Kahanov L, Joscowicz L, Shoshan Y. Two is better than one: longitudinal detection and volumetric evaluation of brain metastases after Stereotactic Radiosurgery with a deep learning pipeline. *J Neuro-Oncol.* 2024; 166(3):547-555.
42. Menze BH, Jakab A, Bauer S, et al. The Multimodal Brain Tumor Image Segmentation Benchmark (BRATS). *IEEE Trans Med Imaging.* 2015; 34(10):1993-2024.
43. Wen PY, Chang SM, Van den Bent MJ, Vogelbaum MA, Macdonald DR, Lee EQ. Response Assessment in Neuro-Oncology Clinical Trials. *J Clin Oncol.* 2017; 35(21):2439-2449.
44. Isensee F, Jager PF, Kohl SAA, Petersen J, Maier-Hein K. Automated Design of Deep Learning Methods for Biomedical Image Segmentation. *arXiv: Computer Vision and Pattern Recognition.* 2019.
45. Li M, Song X, Zhu J, Fu A, Li J, Chen T. The interventional effect of new drugs combined with the Stupp protocol on glioblastoma: A network meta-analysis. *Clin Neurol Neurosurg.* 2017; 159:6-12.
46. Bob Phillips CB, Dave Sackett, Doug Badenoch, Sharon Straus, Brian Haynes, Martin Dawes. Oxford Centre for Evidence-Based Medicine: Levels of Evidence. 2009; <https://www.cebm.ox.ac.uk/resources/levels-of-evidence/oxford-centre-for-evidence-based-medicine-levels-of-evidence-march-2009>.
47. Miller AB, Hoogstraten B, Staquet M, Winkler A. Reporting results of cancer treatment. *Cancer.* 1981; 47(1):207-214.
48. Verdier MCd, Saluja R, Gagnon L, et al. The 2024 Brain Tumor Segmentation (BraTS) Challenge: Glioma Segmentation on Post-treatment MRI. *ArXiv.* 2024; abs/2405.18368.
49. Chelliah A, Wood DA, Canas LS, et al. Glioblastoma and radiotherapy: A multicenter AI study for Survival Predictions from MRI (GRASP study). *Neuro Oncol.* 2024; 26(6):1138-1151.
50. Suter Y, Knecht U, Valenzuela W, et al. The LUMIERE dataset: Longitudinal Glioblastoma MRI with expert

- RANO evaluation. *Sci Data*. 2022; 9(1):768.
51. Huber T, Alber G, Bette S, et al. Progressive disease in glioblastoma: Benefits and limitations of semi-automated volumetry. *PLoS One*. 2017; 12(2):e0173112.
  52. Bauknecht HC, Klingebiel R, Hein P, et al. Effect of MRI-based semiautomatic size-assessment in cerebral metastases on the RANO-BM classification. *Clin Neuroradiol*. 2020; 30(2):263-270.
  53. Tan J, Liu C, Li Y, et al. Assessment of immunotherapy response in intracranial malignancy using semi-automatic segmentation on magnetic resonance images. *Front Immunol*. 2022; 13:1029656.
  54. Wang C, Sun W, Kirkpatrick J, Chang Z, Yin FF. Assessment of concurrent stereotactic radiosurgery and bevacizumab treatment of recurrent malignant gliomas using multi-modality MRI imaging and radiomics analysis. *J Radiosurg SBRT*. 2018; 5(3):171-181.
  55. Fournier L, de Geus-Oei LF, Regge D, et al. Twenty Years On: RECIST as a Biomarker of Response in Solid Tumours an EORTC Imaging Group - ESOI Joint Paper. *Front Oncol*. 2021; 11:800547.
  56. Surov A, Gottschling S, Mawrin C, et al. Diffusion-Weighted Imaging in Meningioma: Prediction of Tumor Grade and Association with Histopathological Parameters. *Transl Oncol*. 2015; 8(6):517-523.
  57. Young GS. Advanced MRI of adult brain tumors. *Neurol Clin*. 2007; 25(4):947-973, viii.
  58. Charles-Edwards EM, deSouza NM. Diffusion-weighted magnetic resonance imaging and its application to cancer. *Cancer Imaging*. 2006; 6(1):135-143.
  59. Radbruch A, Lutz K, Wiestler B, et al. Relevance of T2 signal changes in the assessment of progression of glioblastoma according to the Response Assessment in Neurooncology criteria. *Neuro Oncol*. 2012; 14(2):222-229.
  60. Booth TC, Wiegers EC, Warnert EAH, et al. High-Grade Glioma Treatment Response Monitoring Biomarkers: A Position Statement on the Evidence Supporting the Use of Advanced MRI Techniques in the Clinic, and the Latest Bench-to-Bedside Developments. Part 2: Spectroscopy, Chemical Exchange Saturation, Multiparametric Imaging, and Radiomics. *Front Oncol*. 2021; 11:811425.
  61. Ghimire P, Kinnersley B, Karami G, et al. Radiogenomic biomarkers for immunotherapy in glioblastoma: A systematic review of magnetic resonance imaging studies. *Neurooncol Adv*. 2024; 6(1):vdae055.

**Figure 1.** Flow diagram of search strategy. The flowchart depicts the systematic review search strategy and selection process. Initially, 2088 records were identified from PUBMED, EMBASE, and Web of Science databases of which 1200 records were examined further. Of these, 918 studies were retrieved for detailed evaluation, including 4 additional studies extracted from citation searches. Following abstract analysis, 762 full-text articles were assessed for eligibility. This process culminated in the inclusion of 20 studies in the final analysis.

**Figure 2.** Summary of the QUADAS-2 assessments of the included studies. Graphical representation of included studies (in percentages) in each key domain in terms of the risk of bias and the concerns regarding applicability. each bar signifies the assessed risk levels, with blue indicating 'low' risk/concern, orange signifying 'high' risk/concern, and gray denoting 'unclear' risk/concern.

Accepted Manuscript

**Table 1.** Studies applying machine learning to automated longitudinal treatment response assessment for brain tumors. This table synthesizes 20 studies highlighting the datasets used, research designs, and machine learning models employed. It summarizes the automated processes implemented in the pipeline, such as brain extraction, tumor segmentation, and volume measurement. Performance metrics are presented based on published information or calculated from available data – the metrics vary according to task.

Study	Dataset <ol style="list-style-type: none"> <li>1. Tumor type</li> <li>2. Retrospective/ Prospective</li> <li>3. Single/Multiple sites</li> <li>4. Sequences</li> <li>5. Numbers</li> </ol>	Pipeline steps	Automated model(s)  (index test(s), in comparison with reference standard(s))	Performance metrics for tasks prior to longitudinal assessment e.g., segmentation	Longitudinal performance metrics
Accepted Manuscript					

<p>Chang, et al.<sup>23</sup></p>	<p>Dataset 1 Preoperative cohort:</p> <ol style="list-style-type: none"> <li>Grade II-IV glioma (pre 2021 WHO classification)</li> <li>Retrospective</li> <li>Multisite (4 sites)</li> <li>T1 C, FLAIR</li> <li>843 patients with 843 MRIs (3 site cohorts split 4:1 into train/test sets, one site as external test set with 157 patients)</li> </ol> <p>Dataset 2 Post-operative cohort:</p> <ol style="list-style-type: none"> <li>Grade IV glioma (Glioblastoma) (pre 2021 WHO classification)</li> <li>Retrospective</li> <li>Single-site</li> <li>T1, T1 C, FLAIR</li> <li>54 patients with 713 MRIs (train: test = 4: 1)</li> </ol> <p>Dataset 3: A randomly selected cohort from Dataset 1 and 2: 42 patients (30 train; 12 test)</p>	<ol style="list-style-type: none"> <li>Preprocessing <ol style="list-style-type: none"> <li>Resampling</li> <li>N4 bias correction</li> <li>Registration</li> <li>Intensity normalization</li> <li>Skull-stripping (Brain extraction)</li> </ol> </li> <li>Segmentation</li> <li>Modified RANO<sup>8</sup>-based bi-dimensional product calculation</li> </ol>	<p>3D U-Net compared to expert manual extraction and segmentation</p> <p>Automated modified RANO<sup>8</sup> calculation compared to manual calculation</p> <p>Automated volume assessment compared to volume assessment based on manual segmentation</p>	<ol style="list-style-type: none"> <li>Brain extraction: DSC (95% CI) = 0.94 (0.92 – 0.95) (Dataset 3 test)</li> <li>FLAIR segmentation: DSC = 0.80 (0.75 – 0.80) &amp; volume ICC = 0.92 (P &lt; 0.001) (Dataset 1 test); DSC = 0.82 (0.79 - 0.84) &amp; volume ICC = 0.92 (P &lt; 0.001) (Dataset 1 external test); DSC = 0.70 (0.67 – 0.73) &amp; volume ICC = 0.92 (P &lt; 0.001) (Dataset 2 test)</li> <li>ET segmentation: DSC = 0.70 (0.66 – 0.73) &amp; volume ICC = 0.97 (P &lt; 0.001) (Dataset 2 test)</li> </ol>	<ol style="list-style-type: none"> <li>FLAIR volume longitudinal change ICC = 0.92 (P &lt; 0.001) (Dataset 2 test)</li> <li>ET volume longitudinal change ICC = 0.97 (P &lt; 0.001) (Dataset 2 test)</li> <li>RANO bi-dimensional product (ET) ICC range 0.50-0.77 (P &lt; 0.001) (individual rater differences) (Dataset 2 test)</li> <li>RANO bi-dimensional product (ET) longitudinal change ICC = 0.85 (P &lt; 0.001) (combined raters) (Dataset 2 test)</li> </ol>
-----------------------------------	---	--	---	--	---

Accepted Manuscript

<p>Nalepa, et al.<sup>23</sup></p>	<p>Dataset 1 BraTS 2020 preoperative dataset:</p> <ol style="list-style-type: none"> <li>Grade 4 glioma (Glioblastoma) (WHO, 2021)</li> <li>Retrospective</li> <li>Multisite</li> <li>T1, T1 C, T2, FLAIR</li> <li>660 patients with 660 MRIs (369 train; 125 validate; 166 test)</li> </ol> <p>Dataset 2 Phase 3 cohort:</p> <ol style="list-style-type: none"> <li>Grade 4 glioma (Glioblastoma) (WHO, 2021)</li> <li>Retrospective</li> <li>Multisite</li> <li>T1, T1 C, T2, FLAIR</li> <li>100 patients with 100 preoperative MRIs (100 train)</li> <li>504 patients with 504 postoperative MRIs (464 train; 40 test)</li> </ol>	<ol style="list-style-type: none"> <li>Preprocessing <ol style="list-style-type: none"> <li>Co-registration</li> <li>Skull-stripping</li> <li>Resampling</li> </ol> </li> <li>Segmentation</li> <li>RANO<sup>6</sup>-based bi-dimensional product calculation</li> </ol>	<p>Confidence-aware nnU-Net compared to manual segmentation</p> <p>Automated RANO<sup>6</sup> calculation compared to manual calculation</p> <p>Automated volume assessment compared to volume assessment based on manual segmentation</p>	<ol style="list-style-type: none"> <li>Segmentation performance for ET (Dataset 1 validate and test) mean DSC = 0.744 (95%CI = 0.690-0.799); mean H95 = 39.624 mm</li> <li>Segmentation performance (Dataset 2 test) mean DSC (95% CI) = 0.692 (0.628–0.757), 0.677 (0.631–0.724) and 0.691 (0.604–0.778) for ET, ED, and surgical cavity; mean H95 (25p–75p) = 9.221 mm (6.437–12.000 mm), 9.455 mm (7.176–11.730 mm) and 7.956 mm (5.938–9.975 mm) for ET, ED, and surgical cavity</li> <li>Automatic segmentation volumetric measurements agreement with GT (Dataset 2) ICC (ET): 0.959 mm<sup>3</sup> (p &lt; 0.001) ICC (cavity): 0.960 mm<sup>3</sup> (p &lt; 0.001) ICC (ED): 0.703 mm<sup>3</sup> (p &lt; 0.703)</li> </ol>	<ol style="list-style-type: none"> <li>Inter-rater agreement for RANO bidimensional measurements (Dataset 2 test) manual RANO compared to Automated RANO (Diameters) (ET) ICC = 0.299–0.866 (p &lt; 0.001) manual RANO compared to Automated RANO (Product) (ET) ICC = 0.292–0.858 (p &lt; 0.001) maximum manual RANO compared to Automated RANO (Diameters) ICC: 0.915 (p &lt; 0.001) maximum manual RANO compared to Automated RANO (Product) ICC: 0.919 (p &lt; 0.001)</li> </ol>
<p>Vollmuth, et al.<sup>24</sup></p>	<p>Dataset 1 Heidelberg Cohort:</p> <ol style="list-style-type: none"> <li>Grade II-IV glioma (pre 2021 WHO classification)</li> <li>Retrospective</li> <li>Single-site</li> <li>T1, T1 C, T2, FLAIR, DWI, ADC</li> <li>30 patients with 450 pairs assessment results</li> </ol>	<ol style="list-style-type: none"> <li>Preprocessing <ol style="list-style-type: none"> <li>Skull-stripping</li> <li>Co-registration</li> <li>T1 subtraction</li> </ol> </li> <li>Segmentation</li> <li>Calculation of TTP</li> </ol>	<p>Automated nnU-Net based segmentation and modified RANO<sup>8</sup> assessment compared to manual segmentation and assessment</p> <p>AI-based TTP assessment compared to Manual TTP assessment</p>	<p>N/A</p>	<ol style="list-style-type: none"> <li>TTP Assessment comparison between investigators using AI assistance (95% CI): CCC = 0.91 (0.82-0.95) (p = 0.005) (Dataset 1)</li> <li>LGG TTP (95% CI): CCC = 0.90 (0.76-0.95) (p = 0.008) (Dataset 1)</li> <li>Glioblastomas TTP (95% CI): CCC = 0.83 (0.75–0.92) (p = 0.016) (Dataset 1)</li> <li>SD TTP Measurements (95% CI): 4.8 months (3.7-6.2 months) (p = 0.004) (Dataset 1)</li> <li>SD LGG TTP (95% CI) = -1.7 months (-4.2 to -1.1 months) (Dataset 1)</li> <li>SD Glioblastoma TTP (95% CI) = -0.1 months (-0.5 to 0.0 months) (p &lt; 0.001) (Dataset 1)</li> </ol>

Rudie, et al. <sup>25</sup>	<p>Dataset 1 BraTS 2020 preoperative dataset:</p> <ol style="list-style-type: none"> <li>Grade 4 glioma (Glioblastoma) (WHO, 2021)</li> <li>Retrospective</li> <li>Multisite</li> <li>T1, T1 C, T2, FLAIR</li> <li>369 patients with 369 MRIs (All for training an initial segmentation network)</li> </ol> <p>Dataset 2 Retrospective posttreatment cohort:</p> <ol style="list-style-type: none"> <li>Grade 2-4 glioma (WHO, 2021)</li> <li>Retrospective</li> <li>Single-site</li> <li>T1, T1 C, T2, FLAIR</li> <li>298 patients (198 train; 100 test) with 596 MRIs</li> </ol>	<ol style="list-style-type: none"> <li>Preprocessing       <ol style="list-style-type: none"> <li>DICOM to NiftI conversion</li> <li>Registration</li> <li><math>1 \times 1 \times 1</math> interpolation</li> <li>Skull-stripping</li> <li>Bias correction</li> </ol> </li> <li>Additional preprocessing between two timepoints for longitudinal change networks       <ol style="list-style-type: none"> <li>Registration</li> <li>Subtraction</li> </ol> </li> <li>Segmentation</li> <li>Longitudinal volumetric change classification</li> </ol>	<p>nnU-Net segmentation network compared to 3D U-Net trained only on the dataset 1</p> <p>nnU-Net longitudinal change network compared to attending neuroradiologists manual longitudinal volumetric change classification</p>	<ol style="list-style-type: none"> <li>Segmentation network (Dataset 2 test) [Mean <math>\pm</math> SD; Median with 25%–75% IQRs]       <p>WT: DSC = [0.86 <math>\pm</math> 0.10; 0.89 (0.84–0.93)]; Volume Similarity = [0.94 <math>\pm</math> 0.10; 0.96 (0.92–0.98)]; HD95 (mm) = [6.9 <math>\pm</math> 10.0; 3.3 (1.7–7.1.0)]</p> <p>ED: DSC = [0.85 <math>\pm</math> 0.11; 0.8 (0.83–0.92)]; Volume Similarity = [0.94 <math>\pm</math> 0.09; 0.96 (0.92–0.99)]; HD95 (mm) = [6.6 <math>\pm</math> 10.1; 3.0 (1.4–6.7)]</p> <p>TC: DSC = [0.71 <math>\pm</math> 0.27; 0.82 (0.55–0.92)]; Volume Similarity = [0.82 <math>\pm</math> 0.25; 0.95 (0.74–0.98)]; HD95 (mm) = [8.6 <math>\pm</math> 14.6; 10.4 (1.4–8.3)]</p> <p>ET: DSC = [0.71 <math>\pm</math> 0.26; 0.82 (0.55–0.92)]; Volume Similarity = [0.83 <math>\pm</math> 0.25; 0.96 (0.80–0.99)]; HD95 (mm) = [8.2 <math>\pm</math> 14.7; 10.4 (1.0–7.9)]</p> <p>NCR: DSC = [0.65 <math>\pm</math> 0.29; 0.72 (0.49–0.88)]; Volume Similarity = [0.80 <math>\pm</math> 0.26; 0.90 (0.74–0.97)]; HD95 (mm) = [5.9 <math>\pm</math> 8.1; 10.4 (1.4–6.0)]</p> </li> </ol>	<ol style="list-style-type: none"> <li>Longitudinal change network (Dataset 2 test)       <p>[Mean <math>\pm</math> SD; Median with 25%–75% IQRs]</p> <p>ED change: DSC = [0.73<math>\pm</math>0.25; 0.83 (0.64–0.88)]; Volume Similarity = [0.84<math>\pm</math>0.27; 0.94 (0.85–0.98)]; HD95 (mm) = [10.3<math>\pm</math>11.6; 5.7 (2.0–15.1)]</p> <p>ET change: DSC = [0.60<math>\pm</math>0.26; 0.67 (0.45–0.81)]; Volume Similarity = [0.73<math>\pm</math>0.27; 0.86 (0.68–0.92)]; HD95 (mm) = [14.2<math>\pm</math>16.9; 5.4 (2.5–19.4)]</p> </li> <li>Longitudinal classification performance for three classes (Dataset 2 test)       <p>ED Longitudinal change network: Sensitivity = 0.91; Specificity = 0.91; PPV = 0.89; NPV = 0.93; F1 = 0.85; Accuracy = 0.91; P = 0.84</p> <p>ET Longitudinal change network: Sensitivity = 0.88; Specificity = 0.92; PPV = 0.88; NPV = 0.92; F1 = 0.88; Accuracy = 0.90; P = 0.61</p> </li> <li>Longitudinal classification performance for two classes (Dataset 2 test)       <p>ED Longitudinal change network: Sensitivity = 0.94; Specificity = 0.93; PPV = 0.86; NPV = 0.97; F1 = 0.90; Accuracy = 0.93; P = 0.81</p> <p>ET Longitudinal change network: Sensitivity = 0.87; Specificity = 0.94; PPV = 0.87; NPV = 0.94; F1 = 0.87; Accuracy = 0.92; P = 0.48</p> </li> </ol>
-----------------------------	--	--	--	---	---

Accepted Manuscript

Strack, et al. <sup>26</sup>	<p>Dataset 1 Local dataset:</p> <ol style="list-style-type: none"> <li>1. Grade IV glioma (glioblastoma) (pre 2021 WHO classification)</li> <li>2. Retrospective</li> <li>3. Single-site</li> <li>4. T1 C</li> <li>5. 15 patients</li> </ol> <p>Dataset 2 TCIA dataset:</p> <ol style="list-style-type: none"> <li>1. Grade IV glioma (glioblastoma) (pre 2021 WHO classification)</li> <li>2. Retrospective</li> <li>3. Multisite</li> <li>4. T1 C</li> <li>5. 20 patients with 40 MRIs</li> </ol>	<ol style="list-style-type: none"> <li>1. Preprocessing <ol style="list-style-type: none"> <li>a. Resampling</li> <li>b. Histogram matching</li> <li>c. Normalization</li> <li>d. Brain centering</li> <li>e. Skull-stripping</li> </ol> </li> <li>2. Augmentation <ol style="list-style-type: none"> <li>a. Shifting</li> <li>b. Rotation</li> <li>c. Gaussian noise</li> </ol> </li> <li>3. Segmentation by BraTS model</li> <li>4. Wasserstein GANs learning changes between time 1 and time 2 images</li> <li>5. Modified RANO-based classification according to ET volume change</li> </ol>	Automated volume change assessment and classification based on Wasserstein GANs compared to manual volume assessment and modified RANO <sup>8</sup> classification	N/A	<p>ROC analysis of tumor change</p> <p>micro-average AUC = 0.87 (Dataset 1 &amp; 2); total tumor growth AUC = 0.87 (Dataset 1); total AUC = 0.86 (Dataset 2); tumor growth AUC = 0.72 (Dataset 1); tumor reduction AUC = 0.75 (Dataset 1); tumor growth AUC = 0.94 (Dataset 2); tumor reduction AUC = 0.94 (Dataset 2)</p> <p>b. RANO classification</p> <p>overall sensitivity = 0.66 (Dataset 1 &amp; 2); overall specificity = 0.83 (Dataset 1 &amp; 2); total accuracy = 0.66 (Dataset 1 &amp; 2); sensitivity = 0.65 (Dataset 1); specificity = 0.82 (Dataset 1); sensitivity = 0.64 (Dataset 2); specificity = 0.82 (Dataset 2)</p>
------------------------------	---	--	--	-----	---

Accepted Manuscript

<p>Jalalifar, et al.<sup>27</sup></p>	<p>Dataset 1:</p> <ol style="list-style-type: none"> <li>1. Brain metastasis</li> <li>2. Retrospective</li> <li>3. Single-site</li> <li>4. T1 C, FLAIR</li> <li>5. 116 patients with 152 tumors (train: 96 patients with 130 tumors; independent test: 20 patients with 22 tumors)</li> </ol>	<ol style="list-style-type: none"> <li>1. Preprocessing <ol style="list-style-type: none"> <li>a. Resampling</li> <li>b. Voxel intensity normalization</li> <li>c. Co-registration</li> </ol> </li> <li>2. Segmentation</li> <li>3. Calculating the changes of longest diameters and volume</li> <li>4. Treatment response classification as shrinkage/steady/enlargement</li> <li>5. Automatic detection of LC/LF and ARE outcome</li> </ol>	<p>A proposed combination model of 2D U-Nets, 3D U-Net and MSGA for segmentation compared to manual segmentation by expert oncologists</p> <p>Automated longest diameters and volume assessment compared to manual longest diameters assessment based on RANO-BM<sup>7</sup> and volumetric assessment criteria</p> <p>Automatic detection of LC/LF and ARE outcome compared to manual assessments by expert oncologists.</p>	<p>a. Tumor segmentation of baseline and follow-up scans (Dataset 1 independent test)</p> <p>DSC = [0.84 ± 0.07, 0.92 ± 0.04]</p> <p>HD95 (mm) = [2.1 ± 0.6, 3 ± 0.6]</p> <p>VEE (cc) = [0.44 ± 0.4, 0.62 ± 0.6]</p> <p>VEE = [0.10 ± 0.05, 0.20 ± 0.09]</p>	<p>a. Tumor size status detecting (Dataset 1 independent test)</p> <p>Accuracy = 0.86; Precision (Increase) = 0.90; Precision (Stable) = 0.75; Precision (Decrease) = 1.00; Recall (Increase) = 0.90; Recall (Stable) = 0.91; Recall (Decrease) = 0.76</p> <p>b. Tumor response assessments by Longest Diameter of Tumor (Dataset 1 independent test)</p> <p>Accuracy = 0.84; (Enlargement, PD) Precision = 0.78; (Steady, SD) Precision = 0.92; (Shrinkage, PR) Precision = 0.82; (Enlargement, PD) Recall = 0.90; (Steady, SD) Recall = 0.82; (Shrinkage, PR) Recall = 0.82</p> <p>c. Tumor response assessments by Tumor Volume (Dataset 1 independent test)</p> <p>Accuracy = 0.81; (Enlargement, PD) Precision = 0.76; (Steady, SD) Precision = 0.86; (Shrinkage, PR) Precision = 0.80; (Enlargement, PD) Recall = 0.80; (Steady, SD) Recall = 0.89; (Shrinkage, PR) Recall = 0.71</p> <p>d. Detecting LC/LF and ARE outcomes by RANO-BM (Dataset 1 independent test)</p> <p>Accuracy (LC/LF) = 0.91; Sensitivity (LC/LF) = 0.89; Specificity (LC/LF) = 0.92; Accuracy (ARE) = 0.91; Sensitivity (ARE) = 1.00; Specificity (ARE) = 0.89</p>
---------------------------------------	---	---	---	--	--

Accepted Manuscript

<p>Kickingereder, et al.<sup>28</sup></p>	<p>Dataset 1 Heidelberg training dataset:</p> <ol style="list-style-type: none"> <li>Grade II-IV glioma (pre 2021 WHO classification)</li> <li>Retrospective</li> <li>Single-site</li> <li>T1, T2, T1 C, FLAIR</li> <li>455 patients with 455 MRIs (five-fold)</li> </ol> <p>Dataset 2 Heidelberg independent test dataset:</p> <ol style="list-style-type: none"> <li>Grade II-IV glioma (pre 2021 WHO classification)</li> <li>Retrospective</li> <li>Single-site</li> <li>T1, T2, T1 C, FLAIR</li> <li>40 patients with 239 MRIs</li> </ol> <p>Dataset 3 Heidelberg simulation dataset:</p> <ol style="list-style-type: none"> <li>Grade II-IV glioma (pre 2021 WHO classification)</li> <li>Retrospective</li> <li>Single-site</li> <li>T1, T2, T1 C, FLAIR</li> <li>466 patients with 595 MRIs</li> </ol> <p>Dataset 4 EORTC-26101 external testing dataset:</p> <ol style="list-style-type: none"> <li>Grade IV glioma (glioblastoma) (pre 2021 WHO classification)</li> <li>Prospective</li> <li>Multisite</li> <li>T1, T2, T1 C, FLAIR</li> </ol>	<ol style="list-style-type: none"> <li>Preprocessing <ol style="list-style-type: none"> <li>DICOM to NiftI conversion</li> <li>Reorientation</li> <li>Skull-stripping</li> <li>Registration</li> <li>T1 subtraction</li> </ol> </li> <li>Segmentation</li> <li>Tumor response classification and TTP calculation</li> </ol>	<p>U-Net-based model for segmentation compared to manual segmentation</p> <p>Automated volume assessment compared to manual volume assessment based on RANO in 2010<sup>6</sup></p> <p>Automated TTP calculation compared to manual TTP assessment</p>	<ol style="list-style-type: none"> <li>CE segmentation agreement DSC (95% CI) = 0.89 (0.86-0.90) (Dataset 2); DSC (95% CI) = 0.91 (0.90-0.92) (Dataset 4)</li> <li>CE volume agreement DSC (95% CI) = 0.99 (0.99-1.00) (Dataset 2); DSC (95% CI) = 0.99 (0.99-0.99) (Dataset 4)</li> <li>NE segmentation agreement DSC (95% CI) = 0.93 (0.92-0.94) (Dataset 2); DSC (95% CI) = 0.93 (0.93-0.94) (Dataset 4)</li> <li>NE volume agreement DSC (95% CI) = 0.99 (0.99-0.99) (Dataset 2); DSC (95% CI) = 0.98 (0.98-0.99) (Dataset 4)</li> <li>Concordance Correlation Coefficients <math>\geq 0.98</math> (Dataset 2 &amp; 4)</li> </ol>	<p>a. Agreement in quantitative volumetrically defined TTP is 0.90, P = 0.94 (Dataset 2) and 0.87, P = 0.77 (Dataset 4)</p> <p>Note: no RANO assessment evaluation</p>
---	---	---	--	---	--

	5. 532 patients with 2034 MRIs				
Chen, et al. <sup>29</sup>	<p>Dataset 1 longitudinal dataset:</p> <ol style="list-style-type: none"> <li>1. Brain metastases</li> <li>2. Retrospective</li> <li>3. Single-site</li> <li>4. T1</li> <li>5. 85 patients with 170 MRIs</li> </ol>	<ol style="list-style-type: none"> <li>1. BMs detection</li> <li>2. Calculating changes in volume and number of BM lesions</li> </ol>	<p>FPN-based CAD of United Imaging Intelligence (uAI) Discover-BMs software compared to manual detection and volume change measurement</p> <p>Automated volume and number of BMs lesion change measurement compared to manual measurement based on RANO-BM<sup>7</sup></p>	<ol style="list-style-type: none"> <li>a. Metastasis lesions detection (Dataset 1) Sensitivity = 0.99; FNs (per scan) = 0.06; FPs (per scan) = 0.53; X<sup>2</sup> = 31.15, p &lt; 0.05</li> <li>b. Follow-up metastasis lesions detection (Dataset 1) Sensitivity = 0.98; FNs (per scan) = 0.08; FPs (per scan) = 0.39; X<sup>2</sup> = 21.09, p &lt; 0.05</li> </ol>	<ol style="list-style-type: none"> <li>a. Agreement of treatment response between automated and manual assessment (Dataset 1): kappa = 0.941, p &lt; 0.05</li> </ol>
Meier, et al. <sup>30</sup>	<p>Dataset 1 longitudinal dataset:</p> <ol style="list-style-type: none"> <li>1. Grade IV glioma (glioblastoma) (pre 2021 WHO classification)</li> <li>2. Prospective</li> <li>3. Single-site</li> <li>4. T1, T1 C, T2, FLAIR</li> <li>5. 14 patients with 64 MRIs</li> </ol>	<ol style="list-style-type: none"> <li>1. Preprocessing by BraTumIA <ol style="list-style-type: none"> <li>a. Skull-stripping</li> <li>b. Intermodality registration</li> <li>c. Bias field correction</li> </ol> </li> <li>2. Voxel-wise segmentation for NCE-T2 and ET-T1 C by BraTumIA</li> <li>3. Volume change measurement</li> </ol>	<p>Machine learning based BraTumIA software segmentation compared to manual segmentation</p> <p>Automated volume changes assessment by BraTumIA compared to manual volume changes assessment</p>	<ol style="list-style-type: none"> <li>a. Volume correlations between BraTumIA and raters (Dataset 1) r-values 0.95 to 0.96, p &lt; 0.001</li> <li>b. Relative over- or underestimation of the volumes (Dataset 1) B compared to. R1 0.52 to 9.9</li> </ol>	<ol style="list-style-type: none"> <li>a. Volume change correlations between BraTumIA and raters (Dataset 1) r-values 0.83 to 0.96, p &lt; 0.001</li> </ol>

Accepted Manuscript

<p>Preetha, et al.<sup>31</sup></p>	<p>Dataset 1 Heidelberg cohort:</p> <ol style="list-style-type: none"> <li>Grade IV glioma (glioblastoma) (pre 2021 WHO classification)</li> <li>Retrospective</li> <li>Single-site</li> <li>T1, T1 C, T2, FLAIR, ADC</li> <li>775 patients with 775 MRIs</li> </ol> <p>Dataset 2 CORE longitudinal cohort:</p> <ol style="list-style-type: none"> <li>Grade IV glioma (glioblastoma) (pre 2021 WHO classification)</li> <li>Prospective</li> <li>Multisite</li> <li>T1, T1 C, T2, FLAIR, ADC</li> <li>260 patients with 1083 MRIs</li> </ol> <p>Dataset 3 CENTRIC longitudinal cohort:</p> <ol style="list-style-type: none"> <li>Grade IV glioma (glioblastoma) (pre 2021 WHO classification)</li> <li>Prospective</li> <li>Multisite</li> <li>T1, T1 C, T2, FLAIR, ADC</li> <li>505 patients with 3147 MRIs</li> </ol> <p>Dataset 4 EORTC-26101 longitudinal cohort:</p> <ol style="list-style-type: none"> <li>Grade IV glioma (glioblastoma) (pre 2021 WHO classification)</li> <li>Prospective</li> <li>Multisite</li> <li>T1, T1 C, T2, FLAIR, ADC</li> <li>521 patients with 1924 MRIs</li> </ol>	<ol style="list-style-type: none"> <li>Preprocessing <ol style="list-style-type: none"> <li>DICOM to NIFTI conversion</li> <li>Reorientation</li> <li>Skull-stripping</li> <li>Registration</li> <li>Resampling</li> <li>Normalization</li> <li>T1 Subtraction</li> </ol> </li> <li>Synthetic T1 C imaging generation</li> <li>ET segmentation map generation</li> <li>ET volume change calculation</li> </ol>	<p>The combination model of U-Net and CGAN for generation, segmentation, calculation based on RANO in 2010<sup>6</sup> compared to manual corresponding assessment</p>	<ol style="list-style-type: none"> <li>Comparison on automated and manual T1 subtraction generation (Dataset 4) <p>CGAN-SSIM (95% CI) = 0.818 (0.817-0.820), <math>p &lt; 0.0001</math></p> <p>U-Net-SSIM (95% CI) = 0.809 (0.807-0.810), <math>p &lt; 0.0001</math></p> </li> <li>Agreement in CE segmentations and volumes between automatic and manual assessment (Dataset 4) <p>CCC (95% CI) = 0.782 (0.751-0.807), <math>p &lt; 0.0001</math></p> <p>Spatial agreement Sørensen-DSC: <math>r</math> (95% CI) = 0.438 (0.401-0.475), <math>p &lt; 0.0001</math></p> </li> </ol>	<ol style="list-style-type: none"> <li>Agreement in TTP: comparison based on automated and manual T1 subtraction generation assessments (Dataset 4) <p>automated 4.2 months (95%CI 4.1-5.2)</p> <p>manual 4.3 months (95%CI 4.1-5.5)</p> <p><math>p = 0.33</math></p> </li> </ol>
-------------------------------------	---	--	--	---	---

<p>Cho, et al.<sup>32</sup></p>	<p>Dataset 1 SNUBH training data:</p> <ol style="list-style-type: none"> <li>Brain metastases</li> <li>Retrospective</li> <li>Single-site</li> <li>T1, T1 C, T2, FLAIR</li> <li>174 patients with 127 MRIs</li> </ol> <p>Dataset 2 SNUBH temporal test set #1:</p> <ol style="list-style-type: none"> <li>Brain metastases</li> <li>Retrospective</li> <li>Single-site</li> <li>T1, T1 C, T2, FLAIR</li> <li>40 patients with 20 MRIs</li> </ol> <p>Dataset 3 SNUBH temporal test set #2:</p> <ol style="list-style-type: none"> <li>Brain metastases</li> <li>Retrospective</li> <li>Single-site</li> <li>T1, T1 C, T2, FLAIR</li> <li>12 MRIs</li> </ol> <p>Dataset 4 SNUH external geographic test:</p> <ol style="list-style-type: none"> <li>Brain metastases</li> <li>Retrospective</li> <li>Multisite</li> <li>T1, T1 C, T2, FLAIR</li> <li>24 patients with BM and 11 patients without BM</li> </ol>	<ol style="list-style-type: none"> <li>Preprocessing <ol style="list-style-type: none"> <li>Normalization</li> <li>Isotropic reconstruction</li> </ol> </li> <li>Brain segmentation</li> <li>Brain parenchyma extraction</li> <li>BM detection using 3D U-Net</li> <li>BM segmentation using 2D U-Net (DenseNet 201)</li> <li>3D rigid registration</li> <li>Volumetric changes calculation</li> </ol>	<p>DL-CAD compared to MD for segmentation, detection</p> <p>Automated compared to manual RANO-BM<sup>7</sup> (the changes of the sum in longest diameters) and volumetric response criteria</p>	<p>a. BM detection</p> <p>Sensitivity (95% CI) = 0.58 (0.53-0.63); DSC = 0.67 ± 0.23; FP/scan = 2.50 (Dataset 2)</p> <p>Sensitivity (95% CI) = 0.80 (0.61-0.92); DSC = 0.76 ± 0.26; FP/scan = 2.20 (Dataset 3)</p> <p>Sensitivity (95% CI) = 0.76 (0.66-0.84); DSC = 0.66 ± 0.22; FP/scan = 7.60 (Dataset 4)</p> <p>b. BM measuring &gt;= 5mm</p> <p>Sensitivity (95% CI) = 0.75 (0.70-0.80); DSC = 0.69 ± 0.22; FP/scan = 0.80 (Dataset 2)</p> <p>Sensitivity (95% CI) = 0.95 (0.74-1.00); DSC = 0.82 ± 0.20; FP/scan = 0.50 (Dataset 3)</p> <p>Sensitivity (95% CI) = 0.88 (0.77-0.95); DSC = 0.68 ± 0.20; FP/scan = 1.90 (Dataset 4)</p>	<p>a. Agreement of the response assessment in RANO-BM and volumetry (Dataset 2 &amp; 3 &amp; 4)</p> <p>RANO-BM measurement: k (95% CI) = 0.52 (0.26-0.79); volumetric measurement: k (95% CI) = 0.68 (0.41-0.94)</p>
---------------------------------	--	--	---	---	--

<p>Hsu, et al.<sup>33</sup></p>	<p>Dataset 1:</p> <ol style="list-style-type: none"> <li>Brain metastases</li> <li>Retrospective</li> <li>Single-site</li> <li>T1 C</li> <li>20 patients</li> </ol> <p>Dataset 2:</p> <ol style="list-style-type: none"> <li>Brain metastases</li> <li>Retrospective</li> <li>Single-site</li> <li>T1 C</li> <li>train 409 patients with 1345 BMs; test 102 patients with 367 BMs</li> </ol> <p>Dataset 3:</p> <ol style="list-style-type: none"> <li>Brain metastases</li> <li>Retrospective</li> <li>Single-site</li> <li>T1 C</li> <li>32 patients with 123 BMs</li> </ol>	<ol style="list-style-type: none"> <li>BM registration</li> <li>BM segmentation</li> <li>BM tracking</li> <li>Calculation of GT Compared to change and the percent changes of and 3LD and ESD measurement</li> </ol>	<p>Metastasis Tracking with Repeated Observations 3D CNN-based software compared to manual assessments for registration, segmentation, the changes of volume and diameters measurements</p>	<ol style="list-style-type: none"> <li>The average shift across all points of registration <math>1.5 \pm 0.2\text{mm}</math> (95% CI) (Dataset 1)</li> <li>Segmentation performance (Dataset 2 test)</li> </ol> <p>Sensitivity (95% CI) = <math>95\% \pm 3\%</math>; False positive rate (95% CI) = <math>2.4 \pm 0.5</math> per patient; DSC (95% CI) = <math>0.76 \pm 0.03</math></p>	<ol style="list-style-type: none"> <li>Detection rate of new or unirradiated BMs 72% (Dataset 3)</li> <li>Correlation of size responses <math>R^2 = 0.80</math> (Dataset 3)</li> <li>Pearson correlation coefficient for size changes of non-disappeared lesions was 0.88 for 3LD (3-dimensional longest diameter) and 0.86 for ESD (equivalent spherical diameter), with <math>p &lt; 0.001</math> (Dataset 3)</li> </ol>
---------------------------------	---	--	---	---	--

Accepted Manuscript

<p>Kleesiek, et al.<sup>34</sup></p>	<p>Dataset 1 longitudinal data:</p> <ol style="list-style-type: none"> <li>1. Grade IV glioma (glioblastoma) (pre 2021 WHO classification)</li> <li>2. Retrospective</li> <li>3. Single-site</li> <li>4. T1, T1 C, T2, FLAIR</li> <li>5. 15 patients with 71 MRIs</li> </ol> <p>Dataset 2 Brain Tumor Segmentation Challenge 2013:</p> <ol style="list-style-type: none"> <li>1. Grade II–IV gliomas (pre 2021 WHO classification)</li> <li>2. Retrospective</li> <li>3. Multisite</li> <li>4. T1, T1 C, T2, FLAIR</li> <li>5. 30 MRIs</li> </ol>	<ol style="list-style-type: none"> <li>1. Preprocessing for dataset 1 <ol style="list-style-type: none"> <li>a. N3 bias field correction</li> <li>b. Resampling</li> <li>c. Longitudinal registration</li> <li>d. Skull-stripping</li> <li>e. Brain mask generation</li> <li>f. Intra-individual registration</li> <li>g. Brain mask application</li> </ol> </li> <li>2. Preprocessing for dataset 2 <ol style="list-style-type: none"> <li>a. N3 bias field correction</li> <li>b. Normalization</li> </ol> </li> <li>3. Preprocessing for both datasets <ol style="list-style-type: none"> <li>a. T1 Subtraction</li> <li>b. Feature Extraction</li> </ol> </li> <li>4. Volumetry segmentation</li> <li>5. Volumetric assessment classification</li> </ol>	<p>Random forest-based segmentation and volumetric measurements compared to manual volumetric and RANO in 2010<sup>6</sup> measurements</p>	<ol style="list-style-type: none"> <li>a. GTV segmentation <p>DSC = 0.636 (Dataset 1)</p> <p>DSC = 0.963 (Dataset 2)</p> </li> </ol>	<ol style="list-style-type: none"> <li>a. The change of GTV with the virtual raters correlation <math>r = 0.995</math>, <math>p &lt; 0.0001</math> (Dataset 1)</li> </ol>
<p>Ozkara, et al.<sup>35</sup></p>	<p>Dataset 1 Longitudinal Dataset:</p> <ol style="list-style-type: none"> <li>1. Brain metastases</li> <li>2. Retrospective</li> <li>3. Single-site</li> <li>4. T1 C</li> <li>5. 180 patients</li> </ol>	<ol style="list-style-type: none"> <li>1. Tumor segmentation by DL-based algorithm</li> <li>2. Calculation of volume and longest diameters change measurement by thresholding functions</li> </ol>	<p>Automation ML-based software Jazz for segmentation and volumetric assessment compared to manual measurement based on RANO-BM<sup>7</sup> (the change of longest diameter and volume)</p>	<p>N/A</p>	<ol style="list-style-type: none"> <li>a. The agreement of volume changes measurement (Dataset 1)</li> </ol> <p>ICC = 0.98 (95% CI, 0.97-0.98)</p>

Accepted Manuscript

<p>Suter, et al. 2023<sup>36</sup></p>	<p>Dataset 1 LUMERE post-operative Dataset:</p> <ol style="list-style-type: none"> <li>1. Grade IV glioma (glioblastoma) (pre 2021 WHO classification)</li> <li>2. Retrospective</li> <li>3. Single-site</li> <li>4. T1, T1 C, T2, FLAIR</li> <li>5. 80 patients with 502 MRIs</li> </ol> <p>Dataset 2 Scans identified from Dataset 1 containing target lesions:</p> <p>129 MRIs</p>	<ol style="list-style-type: none"> <li>1. Preprocessing <ol style="list-style-type: none"> <li>a. Resampling</li> <li>b. Skull-stripping</li> </ol> </li> <li>2. Segmentation</li> <li>3. Automated 2D measurement (the product of longest perpendicular diameters in the axial space)</li> <li>4. Automated volumetry (quantifying the contrast enhancement volume by counting the voxels of the segmentation label)</li> <li>5. Automated 2.5D measurement (the product of the longest diameters in the tumor 3D space.)</li> <li>6. Classification of treatment response</li> <li>7. Calculation of TTP</li> </ol>	<p>DL-based BraTumIA software and HD-GLIO compared to manual assessments for segmentation, modified RANO<sup>8</sup> measurements, classification and TTP calculation</p>	<p>N/A</p>	<p>a. Agreement of 2D measurements with manual measurements (Dataset 2)</p> <p>HD-GLIO: 0.81, BraTumIA: 0.80</p>
<p>Zhang, et al.<sup>37</sup></p>	<p>Dataset 1 Longitudinal dataset:</p> <ol style="list-style-type: none"> <li>1. Glioblastoma (pre 2021 WHO classification)</li> <li>2. Retrospective</li> <li>3. Single-site</li> <li>4. T1, T1 C, T2, FLAIR</li> <li>5. 634 patients with 3403 MRIs</li> </ol>	<ol style="list-style-type: none"> <li>1. Preprocessing <ol style="list-style-type: none"> <li>a. Field-of-View Standardization</li> <li>b. Skull-stripping</li> </ol> </li> <li>2. Segmentation</li> <li>3. Volumetric measurements</li> <li>4. BT-RADS classification</li> </ol>	<p>nnU-Net based segmentation to manual segmentation and volumetric assessments compared to volumetric assessments based on BT-RADS (volume)</p>	<p>a. Mean <math>\pm</math> SD = 0.8861 <math>\pm</math> 0.2476 for enhancing tumor and 0.9833 <math>\pm</math> 0.0372 for surrounding non-enhancing FLAIR signal abnormality (Dataset 1 internal validation test)</p>	<p>a. The agreement across BT-RADS (F1: 0.587-0.755) (Dataset 1 internal test)</p> <p>b. Kaplan-Meier Survival Analysis: Worse survival for human-assessed progression vs. AI (Log-rank P=0.007) (Dataset 1 internal test)</p> <p>c. Cox Proportional Hazard Model Analysis: AI assessments less accurate for survival prediction (P=0.012) (Dataset 1 internal test)</p>

<p>Prezelski, et al.<sup>38</sup></p>	<p>Dataset 1 Longitudinal dataset:</p> <ol style="list-style-type: none"> <li>Brain metastases</li> <li>Retrospective</li> <li>Single-site</li> <li>T1</li> <li>71 patients with 176 BMs, 629 MRIs</li> </ol>	<ol style="list-style-type: none"> <li>BM detection</li> <li>Segmentation</li> <li>Rigid registration</li> <li>BM volume and longest 3D diameter changes calculation</li> <li>BM classification</li> </ol>	<p>Metastasis Tracking with Repeated Observations 3D CNN-based software compared to manual assessments for detection, segmentation, the calculation of volume and longest 3D diameter changes and classification based on RANO-BM<sup>7</sup></p>	<p>a. BM Detection (Dataset 1)</p> <p>Sensitivity: Nearly 100% for larger lesions, drops below 90% for lesions smaller than 5 mm</p>	<p>a. BM Volume and Longest 3D Diameter Changes Calculation (Dataset 1)</p> <p>Correlation Coefficient (<math>R^2</math>): 0.76 (<math>P = .0001</math>)</p> <p>Comparison between Manual and METRO Measurements: METRO's longest 3D diameter is generally longer than the manual axial diameter</p> <p>b. BM Classification (Dataset 1)</p> <p>Sensitivity: 0.72</p> <p>Specificity: 0.95</p> <p>Precision: 0.81 (Increasing), 0.32 (Stable), 0.36 (Decreasing), 0.66 (Unappreciable)</p> <p>Recall: 0.72 (Increasing), 0.55 (Stable), 0.22 (Decreasing), 0.72 (Unappreciable)</p> <p>Specificity: 0.95 (Increasing), 0.82 (Stable), 0.85 (Decreasing), 0.77 (Unappreciable)</p> <p>F1-score: 0.76 (Increasing), 0.40 (Stable), 0.27 (Decreasing), 0.69 (Unappreciable)</p>
<p>Son, et al.<sup>39</sup></p>	<p>Dataset 1 Segmentation dataset:</p> <ol style="list-style-type: none"> <li>Brain metastases</li> <li>Retrospective</li> <li>Single-site</li> <li>T1, T1 C, T2, FLAIR, BB T1</li> <li>128 patients with 1339 BMs</li> </ol> <p>Dataset 2 Treatment response dataset:</p> <ol style="list-style-type: none"> <li>Brain metastases</li> <li>Retrospective</li> <li>Single-site</li> <li>T1, T1 C, T2, FLAIR, BB T1</li> <li>58 patients with 629 BMs</li> </ol>	<ol style="list-style-type: none"> <li>Skull-stripping</li> <li>BM detection</li> <li>BM segmentation</li> <li>Volumetric changes calculation</li> <li>RANO-BM<sup>7</sup> classification</li> </ol>	<p>RLK-UNet based architecture compared to manual assessments for detection, segmentation, volume changes and classification based on RANO-BM<sup>7</sup></p>	<p>a. Detection Performance: (Dataset 1)</p> <p>Sensitivity: 86.9%</p> <p>Precision: 79.6%</p> <p>False Positives per Scan: 1.76</p> <p>b. Segmentation Performance: (Dataset 2)</p> <p>All BMs (DSC): 0.663</p> <p>Large BMs (DSC): 0.851</p> <p>Small BMs (DSC): 0.535</p> <p>Pearson Correlation Coefficient: 0.96</p> <p>Bland-Altman Analysis: Mean difference of 0.01 cm<sup>3</sup></p>	<p>a. Agreement on treatment response assessment: (Dataset 2)</p> <p>ICC: 0.84 (95% CI: 0.75-0.91)</p> <p>b. Agreement in Response Assessment: 87.9% (51/58 patients)</p> <p>Overestimation of Treatment Response: 6.8% (4/58 patients)</p> <p>Underestimation of Treatment Response: 5.1% (3/58 patients)</p>

<p>Kotowski, et al.<sup>40</sup></p>	<p>Dataset 1 BraTS2021 training dataset:</p> <ol style="list-style-type: none"> <li>1. Glioblastoma (pre 2021 WHO classification)</li> <li>2. Retrospective</li> <li>3. Multisite</li> <li>4. T1, T1 C, T2, FLAIR</li> <li>5. 1251 patients with 1251 MRIs</li> </ol> <p>Dataset 2 Brain Tumor Progression dataset:</p> <ol style="list-style-type: none"> <li>1. Glioblastoma (pre 2021 WHO classification)</li> <li>2. Retrospective</li> <li>3. Multisite</li> <li>4. T1, T1 C, T2, FLAIR</li> <li>5. 20 patients with 40 MRIs</li> </ol>	<ol style="list-style-type: none"> <li>1. Preprocessing by CaPTK <ol style="list-style-type: none"> <li>a. Reorientation</li> <li>b. Resampling</li> <li>c. Denoising</li> <li>d. Bias correction</li> <li>e. Co-registration</li> </ol> </li> <li>2. Skull-stripping</li> <li>3. Segmentation</li> <li>4. Bidimensional and volumetric measurements</li> </ol>	<p>HD-BET 3D U-Net based brain extraction,</p> <p>nnU-Net segmentation, AutoRANO and volumetric measurement compared to manual measurement based on RANO 2010<sup>6</sup></p>	<p>a. Segmentation DSC: (Dataset 2)</p> <p>DeepMedic: mean 0.72, median 0.77 (95% CI: 0.66-0.79)</p> <p>HD-BET: mean 0.73, median 0.79 (95% CI: 0.67-0.80)</p>	<p>a. Bidimensional measurements spearman's correlation coefficient: (Dataset 2)</p> <p>DeepMedic: 0.58</p> <p>HD-BET: 0.68</p> <p>b. Measurable ET Volume Correlation Coefficient: (Dataset 2)</p> <p>DeepMedic: 0.90</p> <p>HD-BET: 0.93</p> <p>c. Full ET Volume Correlation Coefficient: (Dataset 2)</p> <p>DeepMedic: 0.89</p> <p>HD-BET: 0.93</p>
--------------------------------------	--	---	---	--	---

Accepted Manuscript

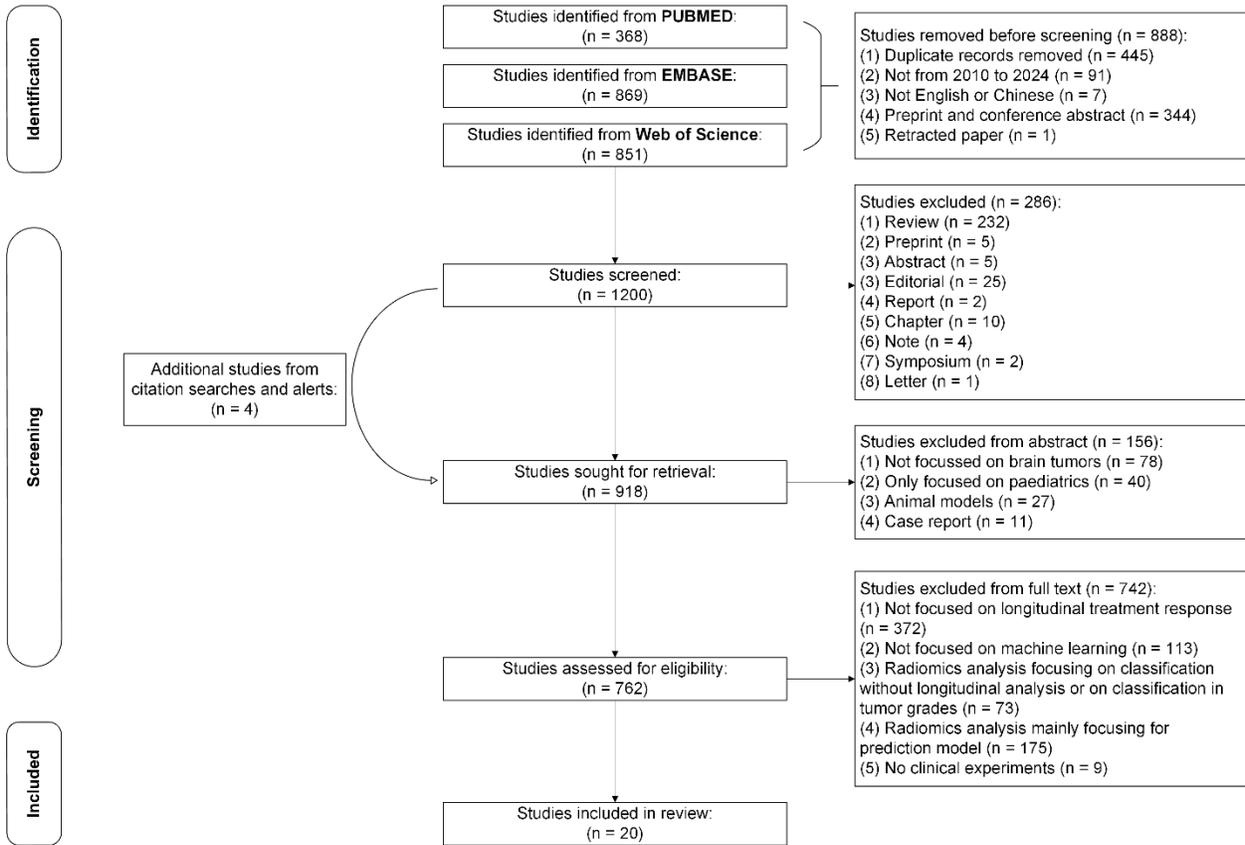
<p>Hammer, et al.<sup>41</sup></p>	<p>Dataset 1 D-STUDIES dataset:</p> <ol style="list-style-type: none"> <li>Brain metastases</li> <li>Retrospective</li> <li>Single-site</li> <li>T1 C</li> <li>226 patients</li> </ol> <p>Dataset 2: D-SCANS dataset (created from Dataset 1 for time-sequenced analysis and pairing creation)</p> <ol style="list-style-type: none"> <li>Brain metastases</li> <li>Retrospective</li> <li>Single-site</li> <li>T1 C</li> <li>226 patients with 500 MRIs</li> </ol> <p>Dataset 3: D-SCAN-PAIRS dataset (created by pairing pre-SRS with post-SRS scans)</p> <ol style="list-style-type: none"> <li>Brain metastases</li> <li>Retrospective</li> <li>Single-site</li> <li>T1 C</li> <li>271 pairs of time-ordered MRI scans (train/validate 205 pairs from 169 patients; test: 66 pairs from 57 patients)</li> </ol> <p>Dataset 4: D-LESIONS-GT dataset (Dataset 2 annotations)</p> <ol style="list-style-type: none"> <li>Brain metastases</li> <li>Retrospective</li> <li>Single-site</li> <li>1,889 lesions annotated (From 439 scans; 1,571 lesions manually annotated)</li> </ol>	<ol style="list-style-type: none"> <li>Brain segmentation</li> <li>Registration</li> <li>Simultaneous lesion detection and segmentation</li> <li>Detection and classification of lesion changes</li> <li>Quantification</li> </ol>	<p>SimU-Net based detection, segmentation, classification compared to manual measurements</p>	<p>a. Lesion Detection (Dataset 3 test): (Lesions &gt; 10 mm) Precision: <math>1.00 \pm 0.00</math> Recall: <math>1.00 \pm 0.00</math></p> <p>(Lesions &gt; 5 mm) Recall (all scenarios): <math>0.95-0.96 \pm 0.13-0.14</math> Simultaneous without prior &amp; Standalone pairs (Precision): <math>0.92-0.93 \pm 0.18-0.19</math> Other scenarios (Precision): <math>0.86-0.89 \pm 0.24-0.28</math></p> <p>(Lesions of All Sizes) Standalone scenarios (Recall): <math>0.82-0.83 \pm 0.28-0.29</math> Other scenarios (Recall): <math>-0.80 \pm 0.28-0.31</math> Simultaneous with prior (Precision): <math>0.83 \pm 0.24</math> Other scenarios (Precision): <math>0.75-0.78 \pm 0.26-0.28</math></p> <p>b. Lesion Segmentation (Dataset 3 test) DSC: <math>0.80-0.90 \pm 0.10-0.21</math> ASSD: <math>0.27-0.62 \pm 0.35-1.27</math> mm Simultaneous without prior: DSC <math>0.83-0.90 \pm 0.10-0.22</math></p>	<p>a. Lesion Matching (Dataset 2) Precision and Recall: <math>1.00 \pm 0.00</math> b. Lesion Change Classification: (Dataset 2) Precision and Recall: <math>1.00 \pm 0.00</math></p>
------------------------------------	---	--	---	---	--

	<p>by experts; From 61 scans: 318 lesions refined from SimU-Net predictions)</p> <p>Dataset 5: D-LESION-PAIRS-GT dataset (Dataset 3 annotations)</p> <ol style="list-style-type: none"> <li>1. Brain metastases</li> <li>2. Retrospective</li> <li>3. Single-site</li> <li>4. 2,055 lesions annotated: (Includes lesions from pre-SRS scan repetitions)</li> </ol>				
--	--	--	--	--	--

WHO, 2021 (World Health Organization Classification of Tumors of the Central Nervous System in 2021)<sup>20</sup>. T1 = T1-weighted; T2 = T2-weighted; T1 C = Post Contrast T1-weighted; FLAIR = Fluid-Attenuated Inversion Recovery; DWI = Diffusion-Weighted Imaging; ADC = Apparent Diffusion Coefficient; RANO = Response Assessment in Neuro-Oncology; ICC = Intraclass Correlation Coefficient; ET = Enhancing tumor; ED = peritumoral edematous, infiltrated, or treatment-changed tissue; NCR = necrotic core; TC = tumor core (AT+NCR); WT = whole tumor (ED+AT+NCR); HD95 = Hausdorff 95th Percentile Distance; TTP = Time to Progression; CCC = Concordance Correlation Coefficient; LGG = Low-Grade Glioma; SD = Standard Deviation; AUC = Area-under-the-curve; 25p-75p = 25% percentile-75% percentile; GT = ground truth; AUC = area under the curve; LF = Local Failure; LC = Local Control; SRS = Stereotactic Radiotherapy; ARE = Adverse Radiation Effect; PD = Progressive Disease; PR = Partial Response; SD = Stable Disease; RANO = Response Assessment in Neuro-oncology; RANO-BM = Response Assessment in Neuro-oncology-Brain Metastases; P = P-value; BMs = Brain metastases; CAD = Computer-Aided Detection; MD = Manual Detection; SSIM = Structural Similarity Index Measure; DL = Deep Learning; BraTS = Brain Tumor Segmentation Challenge; ANN = Artificial Neural Network; 3LD = 3D longest diameter; ESD = Equivalent Sphere Diameter; GTVs = Gross Tumor Volume; CI = Confidence Interval; DSC = Dice Similarity Coefficient; CaPTK = Cancer Imaging Phenomics Toolkit; ASSD = Average Symmetric Surface Distance

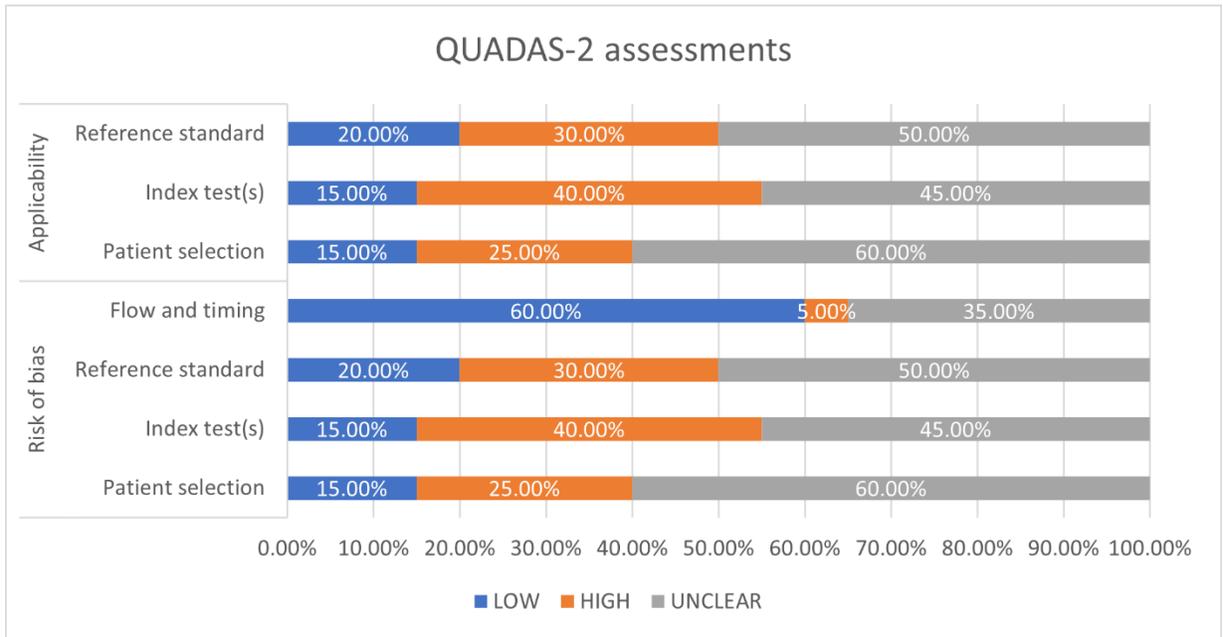
Accepted Manuscript

Figure 1



Accepted

Figure 2



Accepted Manuscript