



Automated quantitative tumour response assessment of MRI in neuro-oncology with artificial neural networks: a multicentre, retrospective study

Philipp Kickingereder*, Fabian Isensee*, Irada Tursunova, Jens Petersen, Ulf Neuberger, David Bonekamp, Gianluca Brugnara, Marianne Schell, Tobias Kessler, Martha Foltyn, Inga Harting, Felix Sahn, Marcel Prager, Martha Nowosielski, Antje Wick, Marco Nolden, Alexander Radbruch, Jürgen Debus, Heinz-Peter Schlemmer, Sabine Heiland, Michael Platten, Andreas von Deimling, Martin J van den Bent, Thierry Gorlia, Wolfgang Wick, Martin Bendszus†, Klaus H Maier-Heint†

Summary

Background The Response Assessment in Neuro-Oncology (RANO) criteria and requirements for a uniform protocol have been introduced to standardise assessment of MRI scans in both clinical trials and clinical practice. However, these criteria mainly rely on manual two-dimensional measurements of contrast-enhancing (CE) target lesions and thus restrict both reliability and accurate assessment of tumour burden and treatment response. We aimed to develop a framework relying on artificial neural networks (ANNs) for fully automated quantitative analysis of MRI in neuro-oncology to overcome the inherent limitations of manual assessment of tumour burden.

Methods In this retrospective study, we compiled a single-institution dataset of MRI data from patients with brain tumours being treated at Heidelberg University Hospital (Heidelberg, Germany; Heidelberg training dataset) to develop and train an ANN for automated identification and volumetric segmentation of CE tumours and non-enhancing T2-signal abnormalities (NEs) on MRI. Independent testing and large-scale application of the ANN for tumour segmentation was done in a single-institution longitudinal testing dataset from the Heidelberg University Hospital and in a multi-institutional longitudinal testing dataset from the prospective randomised phase 2 and 3 European Organisation for Research and Treatment of Cancer (EORTC)-26101 trial (NCT01290939), acquired at 38 institutions across Europe. In both longitudinal datasets, spatial and temporal tumour volume dynamics were automatically quantified to calculate time to progression, which was compared with time to progression determined by RANO, both in terms of reliability and as a surrogate endpoint for predicting overall survival. We integrated this approach for fully automated quantitative analysis of MRI in neuro-oncology within an application-ready software infrastructure and applied it in a simulated clinical environment of patients with brain tumours from the Heidelberg University Hospital (Heidelberg simulation dataset).

Findings For training of the ANN, MRI data were collected from 455 patients with brain tumours (one MRI per patient) being treated at Heidelberg hospital between July 29, 2009, and March 17, 2017 (Heidelberg training dataset). For independent testing of the ANN, an independent longitudinal dataset of 40 patients, with data from 239 MRI scans, was collected at Heidelberg University Hospital in parallel with the training dataset (Heidelberg test dataset), and 2034 MRI scans from 532 patients at 34 institutions collected between Oct 26, 2011, and Dec 3, 2015, in the EORTC-26101 study were of sufficient quality to be included in the EORTC-26101 test dataset. The ANN yielded excellent performance for accurate detection and segmentation of CE tumours and NE volumes in both longitudinal test datasets (median DICE coefficient for CE tumours 0·89 [95% CI 0·86–0·90], and for NEs 0·93 [0·92–0·94] in the Heidelberg test dataset; CE tumours 0·91 [0·90–0·92], NEs 0·93 [0·93–0·94] in the EORTC-26101 test dataset). Time to progression from quantitative ANN-based assessment of tumour response was a significantly better surrogate endpoint than central RANO assessment for predicting overall survival in the EORTC-26101 test dataset (hazard ratios ANN 2·59 [95% CI 1·86–3·60] vs central RANO 2·07 [1·46–2·92]; $p < 0·0001$) and also yielded a 36% margin over RANO ($p < 0·0001$) when comparing reliability values (ie, agreement in the quantitative volumetrically defined time to progression [based on radiologist ground truth vs automated assessment with ANN] of 87% [266 of 306 with sufficient data] compared with 51% [155 of 306] with local vs independent central RANO assessment). In the Heidelberg simulation dataset, which comprised 466 patients with brain tumours, with 595 MRI scans obtained between April 27, and Sept 17, 2018, automated on-demand processing of MRI scans and quantitative tumour response assessment within the simulated clinical environment required 10 min of computation time (average per scan).

Interpretation Overall, we found that ANN enabled objective and automated assessment of tumour response in neuro-oncology at high throughput and could ultimately serve as a blueprint for the application of ANN in radiology to improve clinical decision making. Future research should focus on prospective validation within clinical trials and application for automated high-throughput imaging biomarker discovery and extension to other diseases.

Lancet Oncol 2019; 20: 728–740

Published Online

April 2, 2019

[http://dx.doi.org/10.1016/S1470-2045\(19\)30098-1](http://dx.doi.org/10.1016/S1470-2045(19)30098-1)

*Joint first authors

†Joint senior authors

Department of Neuroradiology

(P Kickingereder MD,

I Tursunova MD, J Petersen MSc,

U Neuberger MD,

G Brugnara MD, M Schell MD,

M Foltyn MD, I Harting MD,

M Prager PhD, A Radbruch MD,

Prof S Heiland PhD,

Prof M Bendszus MD),

Neurology Clinic (T Kessler MD,

M Nowosielski MD, A Wick MD,

Prof M Platten MD,

Prof W Wick MD), Department

of Neuropathology, Institute

of Pathology (F Sahn MD,

Prof A von Deimling MD),

and Department of Radiation

Oncology (Prof J Debus MD),

Heidelberg University Hospital,

Heidelberg, Germany;

Medical Image Computing

(F Isensee MSc, J Petersen,

M Nolden PhD,

K H Maier-Hein PhD),

Department of Radiology

(D Bonekamp MD, A Radbruch,

Prof H-P Schlemmer MD),

and German Cancer Consortium

(DKTK) (T Kessler, F Sahn,

Prof M Platten,

Prof A von Deimling,

Prof W Wick) German Cancer

Research Center (DKFZ),

Heidelberg, Germany;

Department of Neurology,

Medical University Innsbruck,

Innsbruck, Austria

(M Nowosielski); Heidelberg

Institute of Radiation

Oncology, Heidelberg,

Germany (Prof J Debus);

Heidelberg Ion-Beam Therapy

Center, Heidelberg, Germany

(Prof J Debus); Department of

Neurology, Mannheim Medical

Center, University of

Heidelberg, Mannheim,

Germany (Prof M Platten);

Funding Medical Faculty Heidelberg Postdoc-Program, Else Kröner-Fresenius Foundation.

Copyright © 2019 Published by Elsevier Ltd. All rights reserved.

Introduction

The development of novel therapies in neuro-oncology requires reliable and accurate endpoints for the assessment of treatment efficacy since both underestimation and overestimation of efficacy restricts trial proficiency. Although overall survival is the most definitive and objective endpoint to assess the efficacy of an investigational treatment, it is sensitive to other interventions (including crossover treatments) that might influence its applicability as an appropriate endpoint.^{1,2} To overcome these limitations, objective responses and progression-free survival are considered endpoints that reliably assess the efficacy of an investigational treatment, specifically in small cohorts and diseases with multiple lines of treatment.³ In neuro-oncology, progression-free survival is assessed according to the Response Assessment in Neuro-Oncology (RANO) working group criteria,^{3,4} which are widely accepted for use in clinical trials⁵ and increasingly used in routine clinical practice to determine treatment response.⁶ These criteria mainly rely on the assessment of treatment response by use of MRI, which can enable both qualitative and quantitative assessment of tumour burden before, during, and after therapy. Underlying the use of RANO is the assumption that tumours grow in spherical

shapes and that the two-dimensional (2D) measurement of a lesion's largest diameter on MRI is a surrogate marker of tumour volume.⁵ However, in clinical practice, this assumption is not always accurate, since brain tumours frequently display very complicated shapes and anisotropic growth, influenced in part by the surrounding anatomic boundaries, host tissue–tumour interface, or treatment-related effects (eg, areas of necrosis and surgical cavities).^{7,8} Consequently, volumetric or three-dimensional (3D) assessment of tumour burden has been of longstanding interest,^{3,4,9} with studies indicating that volumetric measurements might be more reliable and accurate than 2D measurements.^{10,11} Nevertheless, although volumetric assessment might arguably be one of the most quintessential parameters for accurate assessment of tumour burden and response,⁸ it lacks practicability in a clinical setting. Whereas 2D measurements of tumour diameter can be done quickly and without dedicated software, volumetric measurements require sophisticated and time-consuming postprocessing of MRI data with dedicated software.^{12,13}

Here, we present our development of a comprehensive, scalable, and validated approach, relying on artificial neural networks (ANNs), that we implemented in an

Brain Tumor Center at Erasmus MC Cancer Institute, Rotterdam, Netherlands (Prof M J van den Bent MD); and European Organisation for Research and Treatment of Cancer, Brussels, Belgium (T Gorlia PhD)

Correspondence to: Dr Philipp Kickingereder, Department of Neuroradiology, Heidelberg University Hospital, Heidelberg 69120, Germany philipp.kickingereder@med.uni-heidelberg.de

Research in context

Evidence before this study

MRI is a key method for detection, staging, and evaluation of response to treatment in neuro-oncology. The Response Assessment in Neuro-Oncology (RANO) criteria have been introduced to standardise assessment of patients with neuro-oncological tumours in both clinical trials and daily clinical practice. However, these criteria primarily rely on manual two-dimensional measurements of target lesions, which restricts both the reliability and accuracy of assessment of tumour burden and treatment response. Consequently, longstanding interest has existed in the use of volumetric assessment of tumour burden. We searched PubMed on Oct 31, 2018, with no date restrictions for publications in English using the terms (“neuro-oncology” OR “brain tumor” OR “brain tumour” OR “glioma” OR “glioblastoma”) AND (“volumetry” OR “volumetric”). Our search did not identify any articles that assessed the use of automated quantitative assessment of tumour response in comparison with RANO assessment. Previous studies have indicated that volumetric measurements might be more reliable and accurate than manual two-dimensional measurements. However, non-automated volumetric measurements are not practical in a clinical setting and have previously been cited as a labour-intensive, time-consuming, and complex task that prevents clinical translation. We aimed to develop a

comprehensive and scalable approach, relying on artificial neural networks (ANNs), that enables fully automated quantitative assessment of tumour burden by use of MRI in neuro-oncology.

Added value of this study

We showed that automated volumetric quantification of tumour burden is highly accurate on independent large-scale datasets. Application of this method for automated quantitative identification of disease progression in a multicentre clinical trial dataset outperformed RANO both in terms of reliability and as a surrogate endpoint for predicting overall survival. To facilitate clinical translation, we integrated this approach for fully automated quantitative analysis of MRI in neuro-oncology into an application-ready software infrastructure.

Implications of all the available evidence

Our results show that automated quantitative analysis of MRI using a comprehensive deep-learning approach with ANN could allow radiologists and clinicians to overcome the inherent limitations of manual assessment of tumour burden. This approach could greatly improve and standardise assessment of tumour response in routine clinical practice and clinical trials and might become a valuable asset for decision making in neuro-oncology.

application-ready software infrastructure that enables fully automated quantitative analysis of MRI in neuro-oncology. Specifically, we aimed to investigate the potential of this approach for automated quantitative assessment of tumour response to overcome the inherent limitations of manual assessment of tumour burden.

Methods

Study design and participants

In this multicentre, retrospective study, we analysed MRI data from patients with brain tumours that were acquired at the Heidelberg University Hospital (Heidelberg, Germany) or as part of the European Organisation for Research and Treatment of Cancer (EORTC)-26101 trial, which was run at 38 institutions in Europe, to develop, train, and test an ANN for automated interpretation of MRI data in clinical settings and prediction of time to progression in these patients.

First, we created a training dataset for the ANN from a non-consecutive cohort of adult patients (aged ≥ 18 years) with histologically confirmed glioblastoma or lower-grade glioma (including diffuse astrocytic and oligodendroglial WHO grade II and III tumours as defined by the 2016 WHO Classification of Tumors of the Central Nervous System¹⁴) being treated at the Heidelberg University Hospital (Heidelberg training dataset). Specifically, the Heidelberg training dataset consisted of MRI data from a single timepoint (one MRI scan per patient) either preoperatively from initial diagnosis, or early postoperatively (< 72 h after surgery) or at follow-up, and was specifically ensembled to represent the broad phenotypic appearance of brain tumours on MRI during disease evolution. Appropriate MRI scans were manually identified to enrich the dataset with comparatively uncommon and difficult cases on the basis of the judgment of the neuroradiologist (PK, DB)—eg, cases with complex resection cavities, extensive alterations after treatment, different CE pattern (faint, dot-like, or multifocal pattern), or presence of non-tumoral contrast enhancement (intratumoral or peritumoral blood vessels). We had no further inclusion or exclusion criteria. Second, for independent testing of the ANN once it has been developed, we selected a non-consecutive cohort of adult patients (aged ≥ 18 years) with histologically confirmed glioblastoma or lower-grade glioma from the Heidelberg University Hospital, with individual MRI scans from multiple timepoints for each patient (Heidelberg test dataset). We had no further inclusion or exclusion criteria. This cohort was selected in parallel with and independent of the training dataset. This dataset was a longitudinal dataset with preoperative and consecutive follow-up scans. Finally, another cohort of adult patients (aged ≥ 18 years) with brain tumours undergoing routine MRI scans at the Heidelberg University Hospital was selected for testing of the developed infrastructure for automated tumour segmentation and quantitative assessment of tumour

response in a simulated clinical environment (Heidelberg simulation dataset). MRI scans from all three Heidelberg datasets were acquired according to an established protocol as described previously,^{12,13,15} and included T1-weighted images before (T1-w) and after (cT1-w) administration of contrast agent, and fluid-attenuated inversion recovery (FLAIR) and T2-weighted (T2-w) images (requiring either 3D or 2D with axial orientation). Retrospective evaluation of imaging data from Heidelberg University Hospital was approved by the local ethics committee of the University of Heidelberg and informed consent was waived.

For independent testing of the ANN, we collected data from the EORTC-26101 study, which was a prospective randomised phase 2/3 trial in patients with first progression of a glioblastoma after standard chemoradiotherapy. In brief, the phase 2 trial¹⁶ assessed the optimal treatment sequence of bevacizumab and lomustine (four treatment groups with single drug *vs* sequential *vs* combination), whereas the subsequent phase 3 trial⁷ (two treatment groups) compared patients treated with lomustine alone with those receiving a combination of lomustine and bevacizumab. Overall, the EORTC-26101 study included 596 consecutively recruited patients ($n=159$ in phase 2, $n=437$ in phase 3) with 2593 individual MRI scans acquired at 38 institutions in Europe. The MRI scans were acquired with a uniform imaging protocol at baseline and every 6 weeks until week 24, and thereafter once every 3 months until last follow-up. Data we collected were from T1-w, cT1-w, FLAIR, and T2-w images. The study was conducted in accordance with the Declaration of Helsinki and the protocol was approved by local ethics committees and patients provided written informed consent (EudraCT number 2010-023218-30, ClinicalTrials.gov number NCT01290939). Full study design and outcomes have been published previously.^{7,16} Access to the EORTC-26101 trial data for the present study was granted through an EORTC external research project. Of the data collected, we excluded MRI scans from the EORTC-26101 test dataset if data were corrupted; in the case of incomplete availability of T1-w, cT1-w, FLAIR, or T2-w sequences; or if there were heavy motion artifacts that also precluded central RANO assessment.

Procedures

Figure 1 depicts the analysis workflow. MRI data from the Heidelberg training dataset, Heidelberg test dataset, and EORTC-26101 test dataset were preprocessed identically (full description in appendix p 2). Briefly, preprocessing included brain extraction (ie, removing the skull and extracranial tissue from images), followed by image registration, calculation of T1-subtraction volumes, and preparation of volumetric tumour segmentation masks for each MRI scan. The preparation of volumetric tumour segmentation masks included volumetric delineation of both contrast-enhancing (CE) tumours and

See Online for appendix

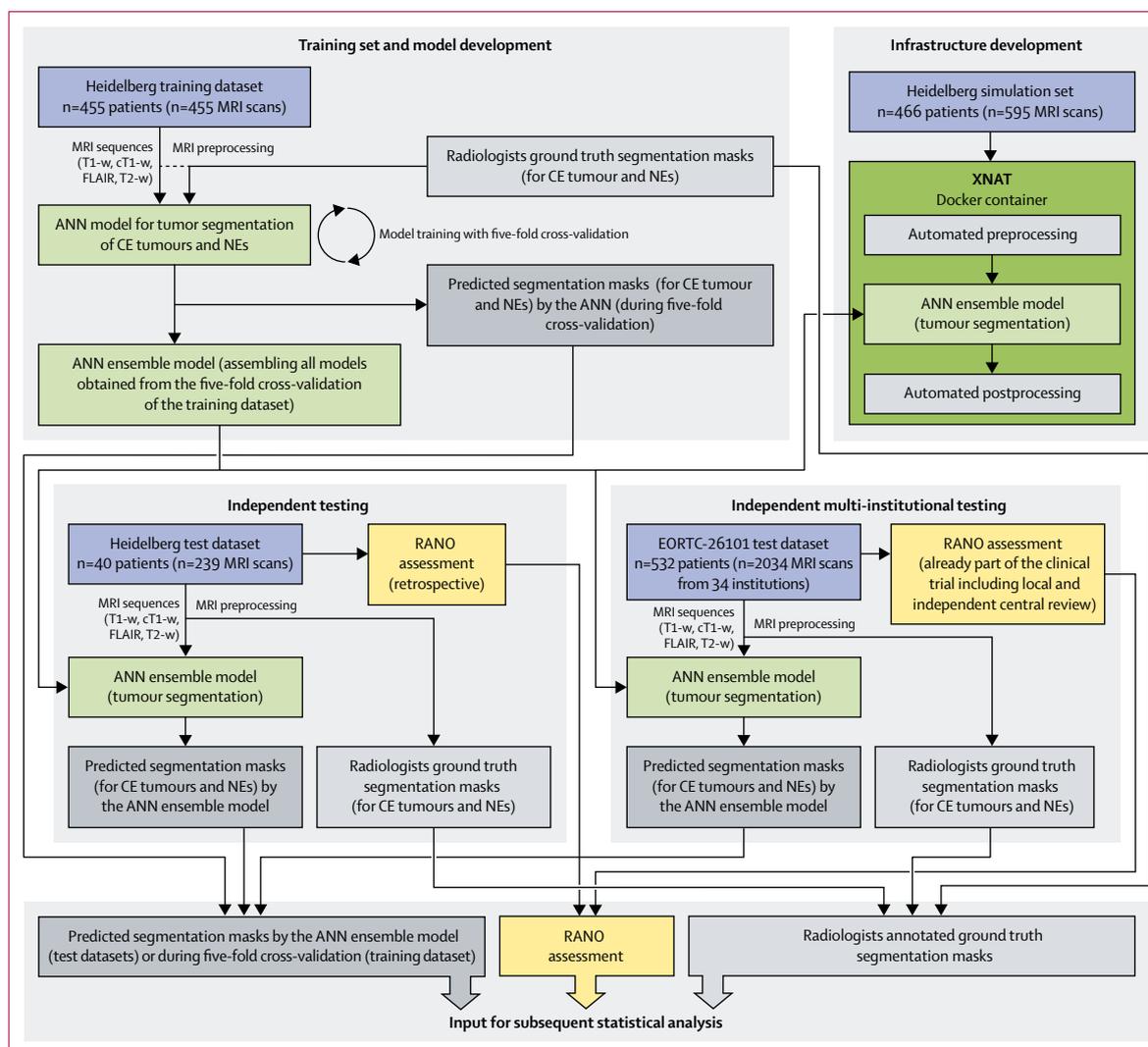


Figure 1: Flowchart of the procedures for training and model development, testing, statistical analysis, and infrastructure development for automated brain tumour segmentation

ANN=artificial neural network. CE=contrast enhancing. cT1-w=T1-weighted images after contrast agent. EORTC=European Organisation for Research and Treatment of Cancer. FLAIR=fluid-attenuated inversion recovery. NEs=non-enhancing T2-signal abnormalities. RANO=Response Assessment in Neuro-Oncology. T1-w=T1-weighted images before contrast agent. T2-w=T2-weighted images.

non-enhancing T2-signal abnormalities (NEs; defined as T2-FLAIR hyperintense abnormality excluding the contrast-enhancing and necrotic portion of the tumour, resection cavity, and obvious leukoaraiosis) by experienced neuroradiologists (for the Heidelberg training and test dataset: by PK, a radiology resident with 6 years of experience, and subsequently checked by DB, a board-certified radiologist and neuroradiologist with 17 years of experience in image processing; for the EORTC-26101 test dataset: by IT, a radiology resident with 3 years of experience and subsequently checked by PK. Any discrepancies were resolved through a consensus discussion) using a semi-automated approach as described previously.^{12,13} All MRI sequences were normalised independently by subtracting the mean value from each voxel

and dividing by its SD. Voxels outside the brain mask were set to zero.

The architecture of our developed ANN was inspired by our work¹⁷ in the Brain Tumor Segmentation (BraTS) challenge¹⁸ and is based on U-Net architecture.¹⁹ The U-Net consists of an encoder and a decoder network that are interconnected with skip connections. Conceptually, the encoder network is used to aggregate semantic information at the cost of decreased spatial information. The decoder is the counterpart to the encoder and reconstructs the spatial information while accounting for the semantic information extracted by the encoder. Skip connections are used to transfer feature maps from the encoder to the decoder to allow for even more precise localisation of the tumour. Our adaptation of the U-Net

(full description of the network architecture, and applied training and testing procedures are in the appendix [pp 2–4, 11]) makes use of residual connections in the encoder²⁰ while keeping the decoder relatively lightweight. During training of the ANN, it processes large input patches (128×128×128 voxels) to effectively capture as much contextual information as possible. To encourage the training of the bottleneck layers, we made use of auxiliary loss layers deep in the network. We used the Heidelberg training dataset to train and validate the ANN (with five-fold cross-validation). For this training, we provided the ANN with the four different MRI sequences (T1-w, cT1-w, FLAIR, and T2-w sequences) for each MRI scan and the corresponding tumour segmentation masks generated by the radiologists as input. These segmentation masks were used as so-called ground truth measures and allowed the ANN to learn the phenotypic appearance of brain tumours on MRI, and consequently enabled automated identification and volumetric delineation of CE tumours and NEs on MRI. In the Heidelberg training dataset, the tumour segmentation masks were predicted during the five-fold cross-validation procedure. Thereby, the Heidelberg training dataset was randomly partitioned into five equal size subsamples (20% of patients). Of the five subsamples, a single subsample (20% of patients) was retained as the validation data for testing the model, and the remaining four subsamples (80% of patients) were used as training data. The cross-validation process was then repeated five times (the folds), with each of the five subsamples used once as the validation data.

We used the longitudinal Heidelberg test dataset and the longitudinal EORTC-26101 test dataset to independently do large-scale testing of the performance of the ANN. Specifically, to predict the segmentation masks with the CE tumours and NEs on MRI in both test datasets, we used the four different MRI sequences (T1-w, cT1-w, FLAIR, and T2-w sequences) from each MRI scan as input into an ANN ensemble model consisting of the five ANN models obtained during cross-validation of the Heidelberg training dataset. The predicted tumour segmentation masks generated by the ANN within the Heidelberg training dataset, Heidelberg test dataset, and EORTC-26101 test dataset served as a fundamental input for all subsequent analyses.

For both longitudinal test datasets, we quantitatively assessed the volumetric tumour response and calculated the time to progression separately, once on the basis of the ground truth tumour segmentation masks generated by the radiologist and once on the basis of the automatically generated ANN-based tumour segmentation masks. We defined tumour progression as an increase in tumour volume (compared with baseline or best response) in either CE tumour or NEs, beyond a minimal tumour volume of 1 cm³; or occurrence of a new CE lesion outside of the CE tumour volume from the previous MRI scan (identified automatically using

dedicated algorithms with the respective segmentation masks over time as input; extended description is on appendix p 5). For the increase in volume, we applied a uniform threshold of 40% to qualify for progression on the basis of CE tumour, whereas for NE we applied a threshold of 40% for patients with lower-grade glioma and 100% for patients with glioblastoma. These volumetric thresholds are justified by the equivalent 2D threshold of use in the RANO criteria, except for the 100% increase in volume in NEs for patients with glioblastoma, for which the RANO working group has not yet defined a 2D threshold (appendix p 5).^{3,4}

Furthermore, in both longitudinal test datasets, we did conventional assessment of tumour response according to RANO criteria.^{3,4} In the Heidelberg test dataset, we retrospectively did RANO assessments (PK and DB), with discrepancies resolved through consensus discussion. For the EORTC-26101 test dataset both local assessment and independent central RANO review had already been done as part of the clinical trial, and so we extracted these data from the trial database. We considered the central RANO review to be an unbiased reference standard because two independent expert radiologists did the review, with discrepancies resolved through consensus discussion, and, by contrast with local RANO assessment, they were masked to the type of treatment, neurological status, steroid doses, and the local RANO investigator's assessment. The RANO assessment in both test datasets was only based on imaging criteria and no additional clinical criteria to allow precise comparison with our quantitative assessment method for tumour response.

We developed an application-ready software infrastructure (appendix pp 5, 6) using the XNAT open-source imaging informatics software platform components. We aimed to enable translation and application of our ANN for automated tumour segmentation and quantitative assessment of tumour response in daily clinical practice and clinical trials, with a specific focus on optimising the processing pipeline—ie, that processing is completed in a clinically acceptable timeframe. In routine clinical practice, automated on-demand processing of an MRI scan is triggered after the images have been acquired on the MRI scanner (or alternatively, for example, within clinical trials uploaded to the XNAT server). Processing of the MRI scans is fully automated and does not require any additional manual intervention. The processed results (superimposed tumour segmentation mask on individual MRI sequences and chart depicting longitudinal tumour volume dynamics) are automatically pushed back to the picture archiving and communication system, where they are available for interpretation. We extensively tested the developed infrastructure in a simulated clinical environment (the Heidelberg simulation dataset) with automated (retrospective) processing of all MRI scans.

Outcomes

We had four main objectives. The first objective was to evaluate the accuracy of automated volumetric tumour segmentation by use of the ANN in comparison with radiologist ground truth tumour segmentation. This outcome was to be assessed in the Heidelberg training dataset, Heidelberg test dataset, and the EORTC-26101 test dataset. Our second objective was to assess within both longitudinal test datasets (the Heidelberg test dataset and the EORTC-26101 test dataset) the spatial and temporal tumour volume dynamics in each patient to calculate a quantitative time to progression that is volumetrically defined and investigate the extent of agreement (reliability) of this quantitative volumetrically defined time to progression (ie, time to progression calculated from ANN-based automated tumour segmentation masks *vs* time to progression calculated from radiologist ground truth tumour segmentation masks) and to compare this reliability (in the EORTC-26101 test dataset) with the extent of agreement (reliability) for the time to progression from RANO (ie, time to progression calculated from RANO local assessment *vs* time to progression calculated from RANO central review). Time to progression was calculated from the date of baseline MRI after surgery in the Heidelberg test dataset and from the date of randomisation in the EORTC-26101 test dataset (censored at the date of last MRI if no progression occurred during follow-up). Our third objective was to compare the performance of quantitative volumetrically defined time to progression versus time to progression calculated from RANO central review (unbiased reference standard) as surrogate endpoints for predicting overall survival within the EORTC-26101 test dataset (overall survival and RANO information taken from the EORTC-26101 trial database; and overall survival calculated from the date of randomisation until death or last follow-up).²¹ Our final objective was to implement the ANN for automated tumour segmentation and quantitative assessment of tumour response in an application-ready software infrastructure and apply it to the Heidelberg simulation dataset.

Statistical analysis

Cohort size for each of the included datasets were determined by the availability of MRI data and not derived from a power calculation. Detailed information on statistical analyses are in the appendix (pp 4, 5). Briefly, we assessed the accuracy of the automatically generated ANN-based tumour segmentation masks for delineating CE tumours and NEs in comparison with the reference—ie, the ground truth segmentation masks generated by a radiologist—in the Heidelberg training dataset, Heidelberg test dataset, and EORTC-26101 test dataset using DICE similarity coefficients for segmentation agreement and Bland-Altman plots and concordance correlation coefficients for volume agreement. The DICE similarity coefficient is a standard measure to report the

performance of a segmentation¹⁸ and measures the extent of spatial overlap between two binary segmentation masks. The DICE similarity coefficient is defined as twice the size of the intersection between the two masks (ground truth [GT] and predicted segmentation mask [PM]), normalised by the sum of their volumes

$$\text{DICE} = \frac{2 |\text{GT} \cap \text{PM}|}{|\text{GT}| + |\text{PM}|}$$

The DICE coefficient can range between 0 (no overlap) and 1 (perfect agreement). The reported 95% CIs for the median DICE coefficients were calculated using bootstrapping with 1000 iterations.

We calculated the relative and absolute agreement in the time to progression between the different methods (quantitative volumetric assessment [based on the spatial and temporal changes in the radiologist's ground truth *vs* automatically generated ANN tumour volumes] *vs* RANO [based on local assessment *vs* central review]) using the Heidelberg test dataset and EORTC-26101 test dataset (the EORTC-26101 test dataset was limited to a subset of patients with complete data; see appendix p 19). We used a one-tailed two-sample test for equality of proportions to assess whether the agreement in time to progression within the EORTC-26101 test dataset was higher for quantitative volumetric assessment (ie, time to progression from ANN-based automated tumour segmentation *vs* time to progression from radiologist ground truth tumour segmentation) than the RANO assessment (ie, time to progression from RANO local assessment *vs* from central review). We generated Kaplan-Meier plots and log-rank tests to assess whether disagreement in the time to progression on a dataset level was higher for RANO than quantitative volumetric assessment.

In the EORTC-26101 test dataset, we assessed the performance of time to progression from RANO (using central RANO review as an unbiased reference standard) and the time to progression from automated ANN-based assessment as surrogate endpoints for predicting overall survival. Specifically, we generated Cox proportional hazards regression models for overall survival, with the time to progression (either from central RANO review or from automated ANN assessment) included as a time-dependent covariate. Hazard ratios (HRs) indicated risk of death at any time during the study period,²² and Z values correspond to the ratio of each regression coefficient to its SE. To control for confounding treatment-specific effects in the EORTC-26101 trial, we included the treatment group as a binary covariate (initial treatment containing bevacizumab *vs* no bevacizumab). We assessed the performance of each Cox proportional hazards regression model with Harrell's concordance index (c-index, with 95% CIs calculated using bootstrapping with 1000 iterations), which is a standard performance measure for model assessment in survival analysis.²³ High c-index values indicate better

performance (ie, better discriminative ability) of the model for predicting overall survival.²³ We did an analysis of deviance that allowed us to compare the reduction in the log-likelihood between different Cox proportional hazards regression models using χ^2 test statistics. Specifically, we compared Cox proportional hazards regression models with and without inclusion of time to progression from ANN as additional time-dependent covariate (to a Cox proportional hazards regression model with time to progression from central RANO review as time-dependent covariate and treatment group as a binary covariate).

Finally, we assessed the prognostic relevance of baseline CE tumour and NE volumes and their early changes between baseline and first follow-up in the EORTC-26101 test dataset for predicting overall survival in a multivariable context including both clinical and molecular parameters. Importance of individual covariates in the Cox proportional hazards regression model was assessed by computing the Wald χ^2 statistic and the proportion of the overall model χ^2 that is due to each covariate (full description of analysis is in the appendix p 5).

p values of less than 0.05 were considered significant. We did all statistical analyses using R version 3.5.1.

Role of the funding source

The funders of the study had no role in study design, data collection, data analysis, data interpretation, or writing of the report. PK, FI, MB, and KHM-H had full access to all the data in the study and had final responsibility for the decision to submit for publication.

Results

The Heidelberg training dataset comprised 455 non-consecutive patients with brain tumours being treated in the Heidelberg University Hospital between July 29, 2009, and March 17, 2017. Of 455 patients, 364 (80%) had histologically confirmed glioblastoma and 91 (20%) lower-grade glioma. Baseline characteristics of these patients are shown in the appendix (p 18). We collected data for one MRI per patient, of which 319 (70%) were preoperative from initial diagnosis and 136 (30%) were early postoperative (<72 h after surgery) or from follow-up scans. These MRI scans, in conjunction with individual radiologist ground truth tumour segmentation masks, were used to train the ANN. Within the Heidelberg training dataset, the ANN acquired the relevant knowledge to accurately delineate the clinically relevant CE tumour and NE compartments (obtained metrics are shown in the appendix [p 9]).

For independent testing and large-scale application of the ANN, two longitudinal testing datasets were compiled: the single-institution Heidelberg test dataset and the multi-institutional EORTC-26101 test dataset. The Heidelberg test dataset comprised 40 non-consecutive patients who were being treated in the Heidelberg

University Hospital between July 29, 2009, and March 17, 2017, with histologically confirmed glioblastoma (n=25 [63%]), or lower grade glioma (n=15 [38%]). These 40 patients had MRI data from 239 scans, with data from multiple timepoints (median of five scans per patient [IQR four to six]). The EORTC-26101 study included 596 patients (n=159 in phase 2, n=437 in phase 3) with 2593 individual MRI scans acquired at 38 institutions in Europe. We excluded 559 MRI scans because the data were corrupted after conversion of file formats from DICOM to *NIfTI* (due to non-standardised centre-specific anonymisation of DICOM files or corrupt or incomplete DICOM files; n=178); incomplete availability of T1-w, cT1-w, FLAIR, and T2-w sequences (requiring either 3D acquisitions or 2D with axial orientation; n=374); or heavy motion artifacts (also precluding central RANO assessment; n=7). Therefore, the multicentre EORTC-26101 trial dataset comprised 532 patients, all with histologically confirmed glioblastomas, from 34 institutions with 2034 MRI scans (median of four scans per patient [IQR three to five]) done on 16 different MRI scanners from four manufacturers (information not available for 102 MRI scans) between Oct 26, 2011, and Dec 3, 2015 (extended description and full list of MRI scanners in appendix [p 18]).

Independent testing in the Heidelberg test dataset yielded median DICE coefficients of 0.89 (95% CI 0.86–0.90) for CE tumours and 0.93 (0.92–0.94) for NEs and in the EORTC-26101 test dataset of 0.91 (0.90–0.92) for CE tumours and 0.93 (0.93–0.94) for NEs (figure 2A, table 1). The performance of the ANN remained stable after application to the broad multicentre setting of the EORTC-26101 test dataset, which was coupled with high agreement between the radiologist ground truth tumour volumes and those automatically predicted by the ANN across both test datasets (concordance correlation coefficients for CE and NE each ≥ 0.98 ; figure 2B and 2C, table 1). The performance of the ANN for the segmentation of CE tumours in the EORTC-26101 test dataset was significantly improved when using 3D T1 and cT1 sequences compared with corresponding 2D sequences (p<0.0001; appendix p 9). These results with an integrative discussion of the individual performance metrics obtained within each dataset are in the appendix (pp 9, 10, 13–16).

By applying the outlined criteria for quantitative identification of disease progression, patients most frequently qualified for progression because of an increase in the CE tumour volume: 19 (48%) of 40 patients in the Heidelberg test dataset and 189 (62%) of 306 patients with complete data in the EORTC-26101 test dataset (appendix p 20). However, beyond disease progression determined via pure volumetric thresholds, a notable number of patients (four [10%] of 40 patients in the Heidelberg test dataset and 26 [8%] of 306 in the EORTC-26101 test dataset) only developed new anatomically distinct CE tumour lesions during follow-up that

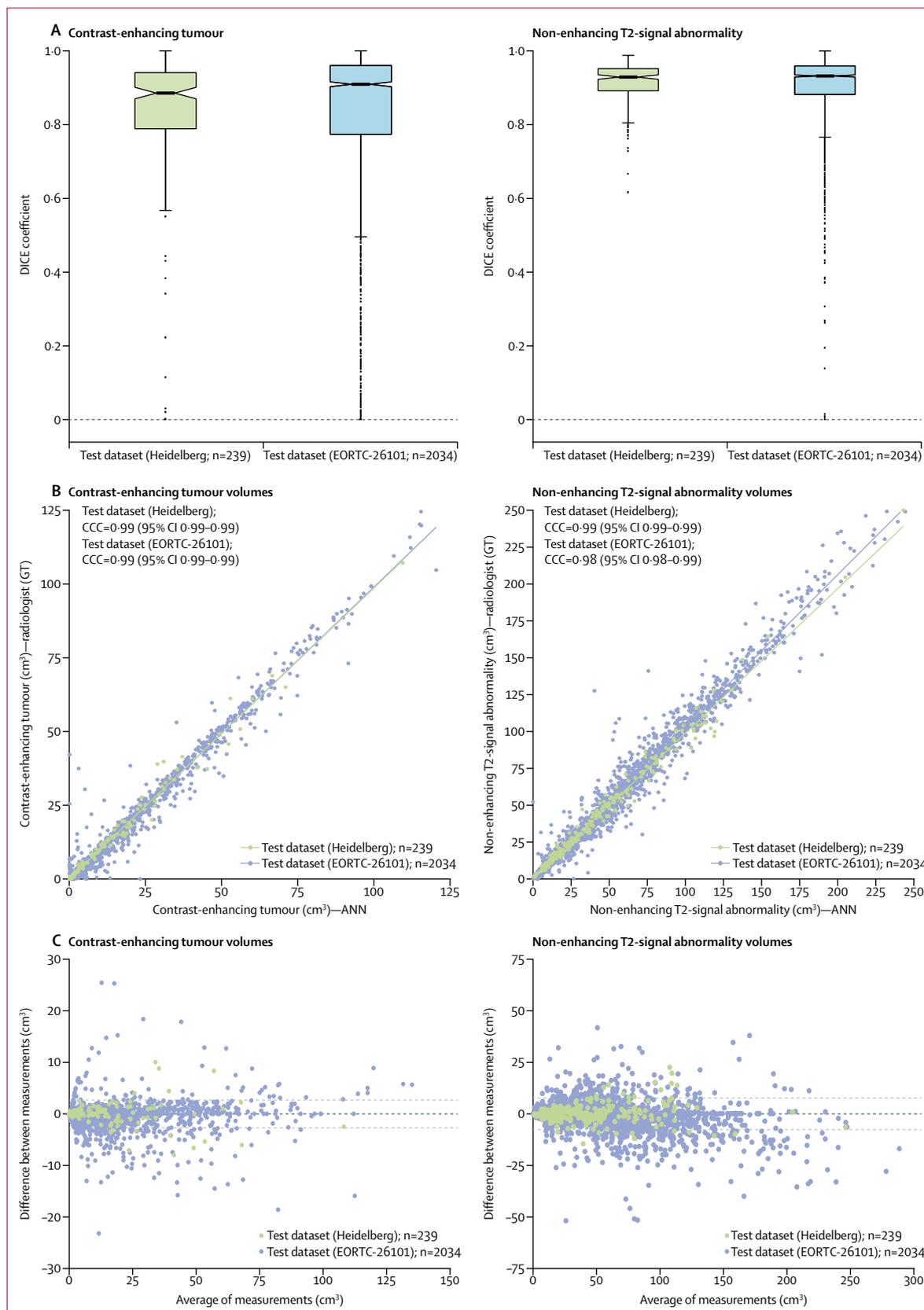


Figure 2: Agreement between the automated volumetric segmentation with the ANN and the radiologist-generated ground truth for tumour segmentation (A) and tumour volume (B, C) (A) Data are median DICE coefficients for tumour segmentation. Blocks show IQR of datapoints, with the horizontal central line showing the median. The sides of blocks are indented and indicate the 95% CI of the median. Whiskers adjacent to the boxes represent 1.5 times the IQR. Dots are outliers. Outliers with DICE coefficients of 0 primarily reflect the uncertainty of accurate tumour segmentation in the post-treatment setting (ie, differentiating true contrast-enhancing tumour from reactive gliosis; more details in appendix pp 9, 10). (B) Data are concordance correlation coefficients (CCCs). (C) Bland-Altman plot. EORTC=European Organisation for Research and Treatment of Cancer. GT=ground truth. ANN=artificial neural network.

would not have qualified for disease progression with volumetric thresholds alone, thereby depicting the relevance of additionally integrating automated screening for newly appearing tumour lesions as a distinct criterion (appendix p 12).

Agreement in quantitative volumetrically defined time to progression (based on radiologist ground truth vs automated assessment with ANN) was 90% (36 of 40 patients) in the Heidelberg test dataset and 87% (266 of 306 patients) in the EORTC-26101 test dataset (table 2). With agreement in only 51% (155 of 306) of patients in the EORTC-26101 test dataset, the reference benchmark (ie, agreement in time to progression between local and central RANO assessment) was significantly lower than the agreement in the quantitative volumetric data between the ANN and radiologist ground truth ($p < 0.0001$; table 2). The higher reliability of the quantitative volumetrically defined time to progression than the reference benchmark was also reflected by the

corresponding Kaplan-Meier plots for time to progression, which showed no significant difference for the quantitative volumetrically defined time to progression based on ground truth versus ANN on a dataset level ($p = 0.94$ for the Heidelberg test dataset, $p = 0.77$ for the EORTC-26101 test dataset); however, a significant difference in the time to progression based on local versus independent central RANO assessment was seen in the EORTC-26101 test dataset ($p < 0.0001$; figure 3).

We compared the performance of calculation of time to progression determined via quantitative volumetric assessment (using ANN) with those determined via RANO (using central assessment as an unbiased reference standard) as surrogate endpoints for predicting overall survival in the EORTC-26101 test dataset. The Cox regression model for overall survival with the time to progression from central RANO as a time-dependent covariate yielded an HR of 2.07 (95% CI 1.46–2.92) with a Z value of 4.12 and a c-index of 0.57 (95% CI 0.54–0.61; $p < 0.0001$; table 3). By contrast, the Cox regression model for overall survival with time to progression from ANN as a time-dependent covariate yielded an HR of 2.59 (95% CI 1.86–3.60) with a Z value of 5.64 and a c-index of 0.62 (95% CI 0.59–0.66; $p < 0.0001$; table 3). The treatment regimen in the EORTC-26101 trial had no confounding effect in either model ($p = 0.34$ for both). The inclusion of the time to progression from ANN as an additional time-dependent covariate yielded a significantly improved model fit over a Cox proportional hazards regression model that only included time to progression from central RANO review as time-dependent covariate and the treatment group as a binary covariate ($\chi^2 = 21.95$; $p < 0.0001$).

The quantitative volumetrically defined criteria for disease progression (40% volume increase in CE tumour for glioblastoma, 40% volume increase in NE for lower-grade glioma, or appearance of a new anatomically distinct lesion) reflect the equivalent (2D) thresholds

	Heidelberg test dataset	EORTC-26101 test dataset
Tumour segmentation agreement		
Contrast-enhancing tumour	0.89 (0.86–0.90)	0.91 (0.90–0.92)
Non-enhancing T2-signal abnormality	0.93 (0.92–0.94)	0.93 (0.93–0.94)
Tumour volume agreement		
Contrast enhancing tumour	0.99 (0.99–1.00)	0.99 (0.99–0.99)
Non-enhancing T2 signal abnormality	0.99 (0.99–0.99)	0.98 (0.98–0.99)

Data are median DICE coefficient for tumour segmentation agreement, and concordance correlation coefficient for tumour volume agreement, with 95% CIs in parentheses. EORTC=European Organisation for Research and Treatment of Cancer. ANN=artificial neural network.

Table 1: Agreement between tumour segmentation masks and tumour volumes predicted by the ANN and those generated by the radiologist (ground truth)

	Absolute agreement (both time and event)		Details on disagreement			
	Yes	No	Earlier progressive disease with alternative	Later progressive disease with alternative	No progressive disease with alternative (but with reference)	Progressive disease with alternative (but not with reference)
Heidelberg test dataset (n=40); reference method vs alternative method						
Quantitative (radiologist) vs quantitative (ANN)	36 (90%)	4 (10%)	2 (5%)	0	2 (5%)	0
Quantitative (radiologist) vs RANO*	29 (73%)	11 (28%)	3 (8%)	6 (15%)	2 (5%)	0
EORTC-26101 test dataset (n=306); reference method vs alternative method						
Quantitative (radiologist) vs quantitative (ANN)	266 (87%)	40 (13%)	23 (8%)	13 (4%)	4 (1%)	0
Quantitative (radiologist) vs RANO (central)	156 (51%)	150 (49%)	110 (36%)	17 (6%)	23 (8%)	0
Quantitative (radiologist) vs RANO (local)	181 (59%)	125 (41%)	36 (12%)	51 (17%)	28 (9%)	10 (3%)
RANO (central) vs RANO (local)	155 (51%)	151 (49%)	15 (5%)	108 (35%)	24 (8%)	4 (1%)

ANN=artificial neural network. RANO=Response Assessment in Neuro-Oncology. EORTC=European Organisation for Research and Treatment of Cancer. *Disagreement between the two RANO readers in five (13%) of 40 patients, resolved through consensus discussion.

Table 2: Agreement in the time to progression on a patient level between the different methods

mandated by the RANO criteria.^{3,4} We chose the applied threshold of 100% volume increase in NE for glioblastoma independently, because the RANO criteria do not provide an equivalent 2D threshold and only suggest that a significant increase in NE qualifies for disease progression. The basis for our conservative threshold of 100% is the theory that dynamics in NE volume are less specific than for CE tumour volumes to determine tumour burden in patients with glioblastoma. This assumption is supported by our findings in the EORTC-26101 test dataset that both baseline CE tumour volume, and early changes in this volume, were the covariates in the multivariable Cox model for overall survival that showed the greatest contribution to the overall model χ^2 compared with other clinical and molecular parameters (table 4). Specifically, baseline CE tumour volume (HR of 1.02 per 1 cm³, 95% CI 1.01–1.03; $p < 0.0001$) and early change in CE tumour volume (HR of 1.04 per 100% increase, 95% CI 1.02–1.06; $p < 0.0001$), showed the highest χ^2 values (18.87 and 19.88) and contributed 25% and 26% to the overall model χ^2 value of 76.97. The next highest χ^2 values were O⁶-methylguanine-DNA methyltransferase (*MGMT*) promoter methylation status, with an χ^2 value of 11.42 (HR 0.61, 95% CI 0.46–0.81; $p = 0.00073$) and glucocorticoid intake, with an χ^2 value of 6.64 (HR 1.52, 95% CI 1.11–2.09; $p = 0.0099$); thus contributing 15% and 9% to the overall model χ^2 value (table 4). Moreover, the baseline NE volume, and early changes in this volume, did not show independent significance within this multivariable model (table 4), thereby supporting our chosen conservative threshold of 100% increase in NE tumour volume in patients with glioblastoma.

To facilitate adoption of our approach for automated tumour segmentation and quantitative volumetric assessment of tumour response in clinical trials and routine clinical practice, we developed a fully automated application-ready processing pipeline for MRI scans (schematic illustration of the workflow is in appendix p 17). This approach enables seamless manufacturer neutral integration into existing clinical infrastructures. We applied the approach within a simulated clinical environment and did fully automated processing (including quantitative tumour response assessment) of a simulation dataset drawn from the Heidelberg University Hospital. This simulation dataset comprised 466 patients with primary intra-axial brain tumours undergoing routine MRI (n=241 [52%] glioblastoma, n=177 [38%] lower-grade glioma [diffuse astrocytic and oligodendroglial WHO grade II and III tumours], and n=48 [10%] other histological entities [pilocytic astrocytoma, pleomorphic xanthoastrocytoma, hemangiopericytoma, dysembryoplastic neuroepithelial tumor, ganglioglioma, medulloblastoma, central neurocytoma, and primary central nervous system lymphoma]), with MRI data from 595 scans collected between April 27, and Sept 17, 2018

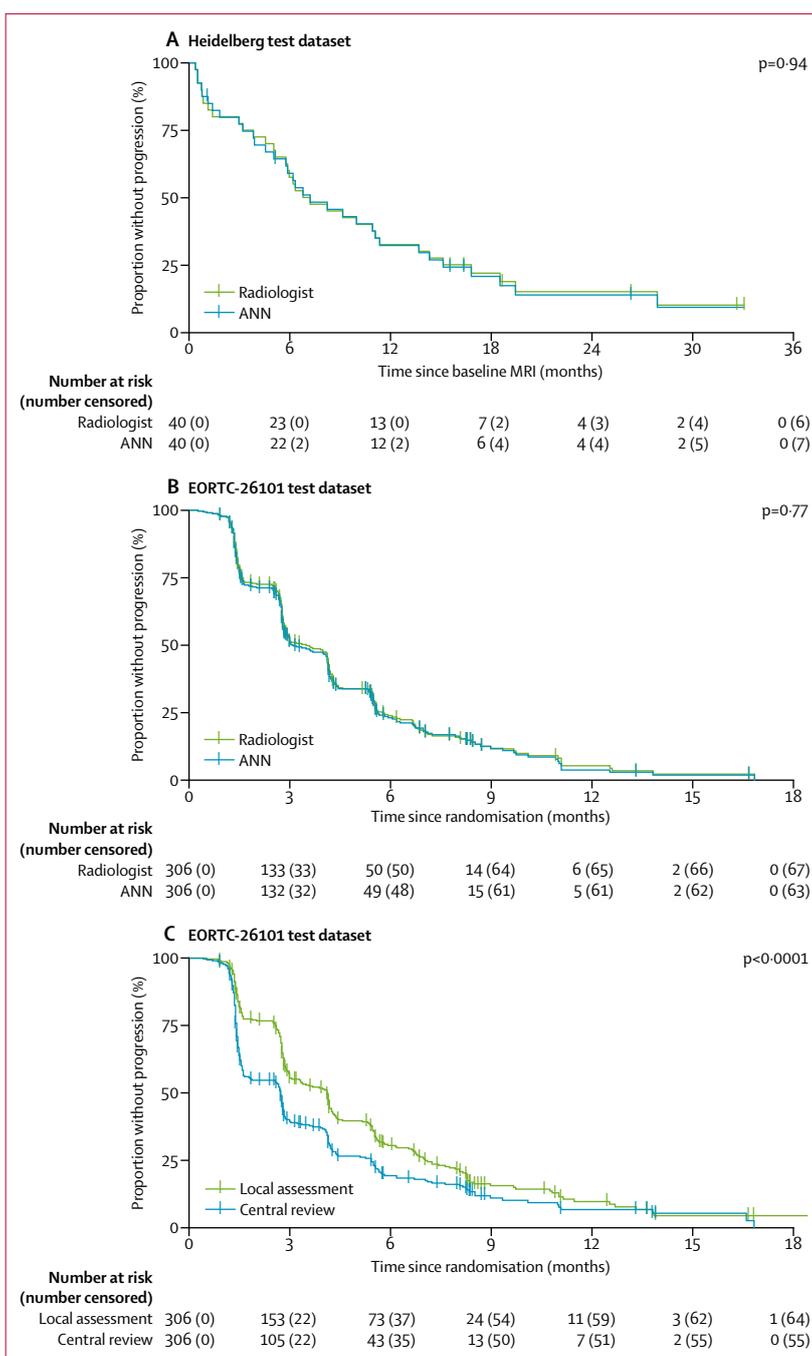


Figure 3: Quantitative volumetrically defined time to progression in the Heidelberg test dataset (A) and EORTC-26101 test dataset (B), and RANO-defined time to progression in the EORTC-26101 test dataset (C) ANN=artificial neural network. EORTC=European Organisation for Research and Treatment of Cancer. RANO=Response Assessment in Neuro-Oncology.

(Heidelberg simulation dataset). We yielded an average computational processing time of 10 min 14 s per MRI exam, thus staying well within a clinically acceptable range (individual data are not shown). The processing pipeline can already accommodate three routines running simultaneously and can be scaled up linearly by adding

	Point estimate	Z value	p value
Quantitative ANN assessment			
Time to progression*	HR 2.59 (1.86–3.60)	5.64	<0.0001
Treatment regimen†	HR 1.14 (0.87–1.47)	0.95	0.34
c-index of the model	0.62 (0.59–0.66)
Central RANO assessment			
Time to progression*	HR 2.07 (1.46–2.92)	4.12	<0.0001
Treatment regimen†	HR 1.14 (0.87–1.47)	0.95	0.34
c-index of the model	0.57 (0.54–0.61)

ANN=artificial neural network. HR=hazard ratio. RANO=Response Assessment in Neuro-Oncology. Z value is the ratio of each regression coefficient to its SE. 95% CIs are shown in parentheses where appropriate. *Time to progression is included as a time-dependent covariate. †Included as binary covariate (initial treatment containing bevacizumab vs no bevacizumab).

Table 3: Cox proportional hazards regression models for overall survival with time to progression in the EORTC-26101 test dataset by assessment method

	Hazard ratio	Wald χ^2	p value
Baseline CE tumour volume, cm ^{3*}	1.02 (1.01–1.03)	18.87	<0.0001
Early change in CE tumour volume, %*	1.04 (1.02–1.06)	19.88	<0.0001
Baseline NE volume, cm ^{3†}	1.00 (1.00–1.00)	0.03	0.87
Early change in NE volume, %†	1.14 (0.92–1.41)	1.47	0.22
Age, years‡	1.01 (0.99–1.02)	1.34	0.25
Sex (female vs male)	0.92 (0.69–1.23)	0.30	0.59
WHO performance status (>0 vs 0)	1.17 (0.85–1.63)	0.92	0.34
MGMT promotor methylation status (methylated vs unmethylated)	0.61 (0.46–0.81)	11.42	0.00073
Glucocorticoids intake (yes vs no)	1.52 (1.11–2.09)	6.64	0.0099

The Cox model included tumour volumes automatically predicted by the artificial neural network (from baseline MRI and the early change in those volumes between baseline and first follow up MRI as covariates). CE=contrast enhancing. NE=non-enhancing T2 signal abnormality. MGMT=O⁶-methylguanine-DNA methyltransferase. *Included as continuous variable (hazard ratios correspond to an increase of 1 cm³). †Included as continuous variable (hazard ratios correspond to an increase of 100%). ‡Included as continuous variable (hazard ratios correspond to an increase of 1 year).

Table 4: Multivariable Cox proportional hazards regression model for overall survival in the EORTC-26101 test dataset

additional processing nodes to the cluster without any need to interrupt the existing workflow (appendix p 6).

Discussion

We showed that automated quantitative analysis of MRI using a comprehensive deep-learning approach with ANN could be a valuable tool for clinical decision making in neuro-oncology. Specifically, the standardisation of our approach has great promise to decrease inter-observer variability of assessment of tumour response that often occurs with RANO criteria. We showed robust performance and generalisability of our ANN in the EORTC-26101 trial dataset, which was across 34 institutions including all major MRI manufacturers, with a broad variety of scanner types and field strengths. Moreover, our results suggest superiority of quantitative volumetric assessment of tumour response, both in terms of reliability and performance as a surrogate

endpoint for predicting overall survival. These findings point out the inherent limitations of the 2D RANO criteria, which only serve as an imperfect surrogate parameter for the assessment of brain tumours that frequently display complicated shapes and anisotropic growth. Additionally, our implementation into a fully automated application-ready processing pipeline for MRI scans in the open-source XNAT framework holds great promise for standardisation of tumour response assessment in neuro-oncology across institutions and clinical trials. Specifically, this processing pipeline not only allows seamless integration that is manufacturer neutral into routine clinical practice independent of pre-existing infrastructures, but also enables investigators to make use of existing XNAT capabilities to manage and coordinate the analysis of MRI data in large multisite clinical trials.

Although quantitative volumetric assessment of tumour response might arguably be one of the most quintessential parameters for accurate assessment of tumour burden and response,^{8,24} it has previously been cited as a labour-intensive, time-consuming, and complex task—even in the case of semi-automated techniques²⁵—which ultimately prevents clinical adoption.^{3,5,26} Our integration of the robust ANN-based tumour segmentation algorithm into a fully automated application-ready processing pipeline for MRI scans allows investigators to overcome this bottleneck that has previously restricted automated and quantitative analysis of MRIs in neuro-oncology. Although objective and automated assessment of tumour response such as we have presented here is the most evident application of this technology, this technology could also be extended to a broad variety of other applications, including automated high-throughput imaging biomarker discovery (eg, volumetric quantification of advanced MRI parameters, such as apparent diffusion coefficients or relative cerebral blood volumes, and radiomics analysis²⁷) or automated contouring of target volumes for radiotherapy treatment planning, all of which require tumour segmentation masks as a fundamental input. For example, in terms of imaging biomarker discovery, we not only confirm the prognostic importance of baseline CE tumour volume,²⁸ but also show that this parameter outperforms well known molecular (eg, MGMT promoter methylation status) or established clinical characteristics within the EORTC-26101 dataset, and consequently provide further rationale to include imaging parameters into clinical trial design.²⁹

Extensive investigation into the comparison of 2D measurements (including RANO) with volumetric measurements of tumour burden has been done,^{10,11,26,30} and a consensus exists that volumetric measurements are more reliable and accurate than 2D measurements,^{4,8,10,11,24} which is also supported by our findings in the EORTC-26101 test dataset (with an increase in reliability from 51% for local vs central RANO assessment to 87% for

quantitative volumetric assessment based on ANN *vs* radiologist ground truth). However, we acknowledge that the added value of quantitative volumetric assessment might be less pronounced when comparing two neuroradiologists with extensive RANO experience than when comparing the readings of less experienced RANO readers with those of highly experienced RANO readers, such as in the EORTC-26101 test dataset (with potentially less experienced local RANO readers *vs* highly experienced central RANO readers). Moreover, whether the higher reliability and accuracy of quantitative volumetric measurements than 2D measurements would translate into clinical relevance has been uncertain until now. Indeed, the few studies^{26,30} that have compared non-automated volumetric assessment of tumour response with RANO had divergent results, with some studies suggesting added value of volumetric assessment of tumours,³⁰ whereas others did not find evidence to favour volumetric assessment over RANO as a surrogate endpoint for predicting overall survival.²⁶ However, all these previous comparisons did not consider that patients might only develop new, anatomically distinct lesions, which would immediately qualify for progression with RANO but might not qualify for progression on the basis of volumetric thresholds because the overall tumour volume could remain below the prespecified threshold. We overcame this limitation and introduced an algorithm to automatically identify the occurrence of new tumour lesions during follow-up. By integrating this algorithm into our automated processing pipeline, we identified that up to 10% of patients fell into this category and would otherwise not have qualified for tumour progression at this timepoint. This finding was of substantial importance for unbiased comparison of the performance of quantitative volumetric assessment of tumour response versus RANO in our study. Consequently, by use of time-dependent Cox regression modelling, our results suggest superiority of time to progression calculated by the automated quantitative ANN-based assessment of tumour response over central RANO assessment as a surrogate endpoint for predicting overall survival.

Our study had some limitations. First, we acknowledge the retrospective nature of the study and the relatively small, single-centre dataset used for training of the ANN. Although we specifically enriched the Heidelberg training dataset with comparatively uncommon and difficult cases, a larger dataset might allow further improvement of the accuracy of the ANN. Second, given the short follow-up period in the Heidelberg simulation dataset, we were unable to investigate the accuracy of automated quantitative volumetric assessment of tumour response in comparison with RANO in this dataset. Third, the suggested added value of automated quantitative volumetric assessment of tumour response compared with RANO, both in terms of reliability and performance, as a surrogate endpoint for predicting overall survival in the EORTC-26101 test dataset requires further

validation in a prospective setting. This investigation is currently ongoing via application of the fully automated MRI-processing pipeline within the XNAT infrastructure as part of central neuroradiological assessment for the N²M² umbrella multicentre trial (NCT03158389) in newly diagnosed patients with non-MGMT hypermethylated glioblastoma.³¹ Moreover, refinement from a methodical perspective will focus on further improving the segmentation performance of the ANN. Specifically, although the ANN and its accompanying training scheme were heavily inspired by our contribution¹⁷ to the BraTS 2017 challenge, we acknowledge that a complete understanding of all design choices and their relative contribution to segmentation performance could point us towards potential further improvements and thus constitutes a valuable topic of research for further projects. From a clinical perspective, we will focus on also including advanced MRI parameters (eg, apparent diffusion coefficients or relative cerebral blood volumes) into the automated analysis workflow, which is of specific importance in the era of immunotherapy—eg, for early separation of pseudoprogression from true progression.^{32,33} Finally, the scalability and flexibility of our approach will enable further extension to other disease entities (eg, quantification of lesion load in multiple sclerosis).³⁴

Overall, our results suggest that ANN can enable objective and automated assessment of tumour response and imaging biomarker discovery in neuro-oncology at high throughput, and could ultimately serve as a blueprint for the application of ANN in radiology to improve clinical decision making.

Contributors

PK, MB, WW, FI, and KHM-H designed the study. PK, IT, MNow, UN, DB, GB, MS, and MF did quality control of MRI data. PK, IT, DB, UN, GB, MS, and MF preprocessed the MRI data. FI developed, trained, and applied the artificial neural network. PK and FI postprocessed the data generated by the artificial neural network. JP, MNol, PK, MPr, FI, and KHM-H did clinical translation (ie, design and development of the automated postprocessing workflow, integration into the XNAT infrastructure, and application in the simulated clinical environment). WW, MB, IH, MJvdB, TG, FS, AvD, and MPI critically contributed to the primary analysis of the relevant data from the EORTC-26101 trial that were used within this study. FS, TK, and AvD analysed the methylation array data in the EORTC-26101 dataset. PK and TG did statistical analyses. PK, FI, MB, WW, and KHM-H interpreted the findings with essential input from all coauthors. PK, FI, and JP prepared the first draft of the manuscript. All authors critically revised the manuscript for important intellectual content and approved the final version.

Declaration of interests

DB reports personal fees from Profound Medical outside of the submitted work. FS reports personal fees from Agilent and Illumina outside of the submitted work. AR reports grants, personal fees, and other support from Bayer and Guerbet, and personal fees and other support from GE-Healthcare, Bracco, Siemens, Sanofi, and Medscape outside of the submitted work. JD reports grants from ViewRay, CRI—The Clinical Research Institute, Accuray, RaySearch Laboratories, Vision RT, Merck Serono, Astellas Pharma, AstraZeneca, Siemens Healthcare, Solution Akademie, Egomed, Quintiles, Pharmaceutical Research Association, Boehringer Ingelheim, PTW-Freiburg, and Nanobiotix outside of the submitted work. SH reports grants from German Research Council and Dietmar-Hopp Foundation outside of the submitted work. MPI reports non-financial support from Pfizer, and grants and personal fees from Bayer outside of the submitted work. MPI has a patent IDH1 vaccines

licensed, a patent H3 vaccine pending, and a patent AHR inhibitor with royalties paid to Bayer. MJvdB reports personal fees from Roche, Cellgene, Bristol-Myers Squibb, AGIOS, Merck Sharpe & Dohme, and Boehringer Ingelheim; and grants and personal fees from AbbVie outside of the submitted work. WW reports grants from Apogenix, Boehringer Ingelheim, Pfizer; grants and personal fees from Merck Sharp and Dohme and Roche; and personal fees from Bristol-Myers Squibb and Celldex outside of the submitted work. MB reports personal fees from Boehringer Ingelheim, Merck, Bayer, Teva, B Braun, Springer, and Vascular Dynamics; grants and personal fees from Novartis, Codman, and Guerbet; and grants from Siemens, Hopp Foundation, German Research Council, the European Union, Stryker, and Medtronic outside of the submitted work. All other authors declare no competing interests.

Acknowledgments

PK was supported by the Medical Faculty Heidelberg Postdoc-Program and the Else Kröner-Fresenius Foundation (Else-Kröner Memorial Scholarship). FS was supported by the Else Kröner-Fresenius Foundation (EKFS Excellence Scholarship).

References

- Chinot OL, Wick W, Mason W, et al. Bevacizumab plus radiotherapy-temozolomide for newly diagnosed glioblastoma. *N Engl J Med* 2014; **370**: 709–22.
- Gilbert MR, Dignam JJ, Armstrong TS, et al. A randomized trial of bevacizumab for newly diagnosed glioblastoma. *N Engl J Med* 2014; **370**: 699–708.
- van den Bent MJ, Wefel JS, Schiff D, et al. Response assessment in neuro-oncology (a report of the RANO group): assessment of outcome in trials of diffuse low-grade gliomas. *Lancet Oncol* 2011; **12**: 583–93.
- 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**: 1963–72.
- 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**: 2439–49.
- Thust SC, Heiland S, Falini A, et al. Glioma imaging in Europe: a survey of 220 centres and recommendations for best clinical practice. *Eur Radiol* 2018; **28**: 3306–17.
- Wick W, Gorlia T, Bendszus M, et al. Lomustine and bevacizumab in progressive glioblastoma. *N Engl J Med* 2017; **377**: 1954–63.
- Korn RL, Crowley JJ. Overview: progression-free survival as an endpoint in clinical trials with solid tumors. *Clin Cancer Res* 2013; **19**: 2607–12.
- Yang D. Standardized MRI assessment of high-grade glioma response: a review of the essential elements and pitfalls of the RANO criteria. *Neurooncol Pract* 2016; **3**: 59–67.
- Chow DS, Qi J, Guo X, et al. Semiautomated volumetric measurement on postcontrast MR imaging for analysis of recurrent and residual disease in glioblastoma multiforme. *AJNR Am J Neuroradiol* 2014; **35**: 498–503.
- Sorensen AG, Patel S, Harmath C, et al. Comparison of diameter and perimeter methods for tumor volume calculation. *J Clin Oncol* 2001; **19**: 551–57.
- Kickingereder P, Götz M, Muschelli J, et al. Large-scale radiomic profiling of recurrent glioblastoma identifies an imaging predictor for stratifying anti-angiogenic treatment response. *Clin Cancer Res* 2016; **22**: 5765–71.
- Kickingereder P, Neuberger U, Bonekamp D, et al. Radiomic subtyping improves disease stratification beyond key molecular, clinical and standard imaging characteristics in patients with glioblastoma. *Neuro Oncol* 2018; **20**: 848–57.
- Louis DN, Perry A, Reifenberger G, et al. The 2016 World Health Organization classification of tumors of the central nervous system: a summary. *Acta Neuropathol* 2016; **131**: 803–20.
- Ellingson BM, Bendszus M, Boxerman J, et al. Consensus recommendations for a standardized brain tumor imaging protocol in clinical trials. *Neuro Oncol* 2015; **17**: 1188–98.
- Wick W, Stupp R, Gorlia T, et al. Phase II part of EORTC study 26101: the sequence of bevacizumab and lomustine in patients with first recurrence of a glioblastoma. *J Clin Oncol* 2016; **34** (suppl 15): 2019 (abstr).
- Isensee F, Kickingereder P, Wick W, Bendszus M, Maier-Hein KH. Brain tumor segmentation and radiomics survival prediction: contribution to the BRATS 2017 challenge. In: Crimi A, Bakas S, Kuijf H, Menze B, Reyes M, eds. Brainlesion: glioma, multiple sclerosis, stroke and traumatic brain injuries. *BrainLes* 2017. Lecture notes in computer science, vol 10670. Springer International Publishing, 2018: 287–97.
- Menze BH, Jakab A, Bauer S, et al. The multimodal brain tumor image segmentation benchmark (BRATS). *IEEE Trans Med Imaging* 2015; **34**: 1993–2024.
- Ronneberger O, Fischer P, Brox T. U-net: convolutional networks for biomedical image segmentation. In: Navab N, Hornegger J, Wells W, Frangi A, eds. Medical image computing and computer-assisted intervention—MICCAI 2015. MICCAI 2015. Lecture notes in computer science, vol 9351. Springer International Publishing, 2015: 234–41.
- He K, Zhang X, Ren S, Sun J. Identity mappings in deep residual networks. In: Leibe B, Matas J, Sebe N, Welling M, eds. Computer vision—ECCV 2016. ECCV 2016. Lecture notes in computer science, vol 9908. Springer International Publishing, 2016: 630–45.
- Han K, Ren M, Wick W, et al. Progression-free survival as a surrogate endpoint for overall survival in glioblastoma: a literature-based meta-analysis from 91 trials. *Neuro Oncol* 2014; **16**: 696–706.
- Sedgwick P, Joeke K. Interpreting hazard ratios. *BMJ* 2015; **351**: h4631.
- Tripepi G, Jager KJ, Dekker FW, Zoccali C. Statistical methods for the assessment of prognostic biomarkers (part I): discrimination. *Nephrol Dial Transplant* 2010; **25**: 1399–401.
- Sorensen AG, Batchelor TT, Wen PY, Zhang WT, Jain RK. Response criteria for glioma. *Nat Clin Pract Oncol* 2008; **5**: 634–44.
- Yushkevich PA, Piven J, Hazlett HC, et al. User-guided 3D active contour segmentation of anatomical structures: significantly improved efficiency and reliability. *Neuroimage* 2006; **31**: 1116–28.
- Gahrman R, van den Bent M, van der Holt B, et al. Comparison of 2D (RANO) and volumetric methods for assessment of recurrent glioblastoma treated with bevacizumab—a report from the BELOB trial. *Neuro Oncol* 2017; **19**: 853–61.
- Kickingereder P, Andronesi OC. Radiomics, metabolic, and molecular MRI for brain tumors. *Semin Neurol* 2018; **38**: 32–40.
- Ellingson BM, Abrey LE, Nelson SJ, et al. Validation of postoperative residual contrast-enhancing tumor volume as an independent prognostic factor for overall survival in newly diagnosed glioblastoma. *Neuro Oncol* 2018; **20**: 1240–50.
- Erickson BJ, Galanis E. Where size matters: imaging-based biomarkers for patient stratification. *Neuro Oncol* 2017; **19**: 7–8.
- Boxerman JL, Zhang Z, Safriel Y, et al. Early post-bevacizumab progression on contrast-enhanced MRI as a prognostic marker for overall survival in recurrent glioblastoma: results from the ACRIN 6677/RTOG 0625 Central Reader Study. *Neuro Oncol* 2013; **15**: 945–54.
- Wick W, Dettmer S, Berberich A, et al. N2M2 (NOA-20) phase I/II trial of molecularly matched targeted therapies plus radiotherapy in patients with newly diagnosed non-MGMT hypermethylated glioblastoma. *Neuro Oncol* 2019; **21**: 95–105.
- Antonios JP, Soto H, Everson RG, et al. Detection of immune responses after immunotherapy in glioblastoma using PET and MRI. *Proc Natl Acad Sci USA* 2017; **114**: 10220–25.
- Okada H, Weller M, Huang R, et al. Immunotherapy response assessment in neuro-oncology: a report of the RANO working group. *Lancet Oncol* 2015; **16**: e534–42.
- Brugnara G, Isensee F, Wick W, Bendszus M, Maier-Hein KH, Kickingereder P. Automated volumetric assessment of multiple sclerosis disease burden and activity with artificial neural networks. *Insights Imaging* 2019; **10** (suppl 1): 22 (abstr).