

Clinical Validation of Segmentation-Based Detection of Glioma Progression

Pablo F. Damasceno^{1,2}, Tyler Gleason^{1,2}, James Hawkins^{1,2}, Tracy Luks^{1,2}, Sharmila Majumdar^{1,2}, Janine M. Lupo^{1,2}, Jason C. Crane^{1,2}, Javier E. Villanueva-Meyer^{1,2}

1 Center for Intelligent Imaging, University of California San Francisco, CA

2 Department of Radiology and Biomedical Imaging, University of California San Francisco, CA

Purpose: To evaluate whether an AI-based method could be used routinely as part of patient care to assist in detecting non-enhancing glioma progression.

Materials and Methods: A 3D U-Net trained (n=481) and validated (n=121) to segment post-surgical lower grade gliomas was used to measure tumor volumes over time and assess progression in a clinical test set. Eight prospective and eight retrospective patients (total 72 exams) who were suspected of progression during their routine outpatient imaging were clinically assessed. Gold standards for progression were derived from clinical reports a posteriori using visual read, and radiologists were blinded to the AI decision at time of reporting.

Results: Progression assessments were presented to radiologists via an easy-to-use, interactive, and interpretable environment in under 10 minutes. Combining prospective and retrospective cases, a final sensitivity of 0.72 and specificity of 0.75 was achieved at progression detection.

Conclusions: Automated detection of glioma progression would provide valuable decision support for routine use.

28 **Introduction.**

29 Radiological assessment of treatment response of glioma is primarily based on change
30 of tumor size on anatomical MRI sequences. Subtle increases in size can be difficult to
31 detect on serial MRI exams, potentially leading to delays in treatment¹⁻³. Although
32 assessment of tumor progression is known to have smaller variance when performed
33 volumetrically⁴, the intensive human labor involved in segmentation limits its routine use
34 in clinical practice⁵.

35 Deep learning methods are particularly suited for medical image evaluation tasks and
36 have been shown to perform well in multiple neuroimaging contexts⁶. However, most
37 reports of AI-informed models focus on performance accuracy, and often do not address
38 the many challenges associated with clinical deployment, such as data quality and
39 variability, usability in clinical workflows, the handling of errors, and transparent
40 presentation of clinically useful and interpretable results⁷.

41 To bridge the gap from research to clinically deployable AI, a lesion segmentation model
42 to assist with the automated detection of tumor progression in lower grade gliomas was
43 developed and deployed for clinical validation. The model was trained on T2-weighted
44 FLAIR images from post-surgical cases, which are of particular importance because
45 resection can lead to brain shifts and image artifacts that make it difficult for automated
46 segmentation algorithms to correctly delineate the tumor⁸.

47 Preliminary results of the clinical validation of our model on accurately detecting glioma
48 progression are presented here to highlight the challenges that must be overcome prior
49 to adoption of a model for routine use in clinical care. The method is based on real-time,
50 neural-network-based lesion segmentation and longitudinal changes in associated
51 volumetrics to predict progression. By integrating automatic image pre-processing, fast
52 and accurate segmentation, and a deployment pipeline for model visualization and
53 evaluation, 3D changes in glioma volume can be tracked over time and presented in an
54 interpretable way to the radiologist.

55 **Project goal and software design.** The primary goal of this project was to evaluate
56 whether an AI-based method could be used routinely as part of patient care to assist in

57 detecting non-enhancing glioma progression. The software package was intended to be:
58 i) easy to use, showing minimal impact on radiologists' workflow; ii) interactive, offering
59 the ability for real-time choice of baseline exam and threshold for progression; and iii)
60 interpretable, showing 3D image segmentations side-by-side with tumor volumes to
61 facilitate interpretation of the model results. For that reason, both the number of times the
62 AI-derived progression agreed with radiologist reports as well as the radiologists'
63 sentiment on software usability were used as success metrics.

64 **Materials and Methods.**

65 **Data collection.** The study received institutional review board approval with consent
66 waiver. 605 T2-weighted FLAIR images from a single institution (GE 750 3T scanner, 3D
67 sagittal acquisitions) were split into training/test sets (80%, 20%). Manually or semi-
68 automatically segmented T2-hyperintense lesions were defined on FLAIR images by
69 experienced investigators using open-source software (3D Slicer⁹). Automatic
70 segmentation of lesion volumes was performed using an encoder-decoder architecture¹⁰
71 previously found to perform well on pre-surgical brain tumor cases¹¹. Pre-processing and
72 post-processing steps including data augmentation were performed as described by
73 Myronenko¹⁰, modified for single channel (FLAIR), rather than four channel (FLAIR, T1
74 pre-contrast, T1 post-contrast, T2) input.

75 **Automatic segmentation and clinical integration.** Clara Deploy¹² was used as the
76 inference service for the model and to handle the communication between PACS and an
77 XNAT¹³ server where results were stored and displayed. Selected exams being read in
78 PACS were sent in real time to a virtual machine running Clara Deploy which identified
79 the desired FLAIR series, performed the lesion segmentation, and stored the results in
80 XNAT. Radiologists could then log into a dedicated XNAT webpage where cases were
81 displayed in tabular form together with respective tumor volumes and progression
82 assessment. In-browser visualization of images and AI-derived segmentations was
83 supported by OHIF¹⁴ (Figure 1).

84 **Model interpretability.** Instead of directly modeling progression, image segmentation
85 and volume change heuristics for assessing progression were employed to provide

86 clinicians with a transparent decision process where thresholds, baseline exams, and
87 segmentations are selected and editable at will (Figure 1). This was intended to enhance
88 confidence in results by exposing the visually interpretable segmentations underlying the
89 progression predictions.

90 **Progression detection heuristic.** Tumors were marked as ‘progressed’ for volume
91 changes larger than 40% between baseline and the MRI of interest¹. Because the
92 baseline exam was variable based on treatment, the system allowed the clinician to
93 choose the baseline scan interactively. This allowed, for instance, selection of baseline
94 tumors with volume greater than 1cm³, as suggested by van den Bent et al¹. Once a
95 baseline exam was chosen, all other exams were marked as blue (stable) or red
96 (progressed) with respect to that baseline (Figure 1).

97 **Prospective evaluation.** Eight prospective cases were obtained during a two-week
98 period from patients who were suspected of progression during routine outpatient
99 imaging. The current MRI and two preceding MRIs (24 total exams) were manually
100 pushed from PACS to the AI algorithm for segmentation. Gold standards for progression
101 were derived from clinical reports *a posteriori* using visual read, and radiologists were
102 blinded to the AI decision at the time of their reporting.

103 **Retrospective evaluation.** To evaluate model true and false positive rates, cases with
104 slow progression were identified via text search of radiology reports. In total, 8 patients
105 with 6 visits each for a total of 48 exams were analyzed.

106 **Results.**

107 **Segmentation training.** The model achieved a mean DICE score of 0.87 ± 0.20 for AI-
108 segmented FLAIR lesions on 124 images composing the test set. The 90th percentile
109 DICE score was also 0.87, revealing that most images had an excellent agreement with
110 ground truth.

111 **Prospective evaluation.** Radiology reports revealed that all 24 baseline-to-follow-up
112 pairs were marked as stable, of which 20 were correctly identified as such by the
113 algorithm, resulting in a specificity = 0.83 (see Fig. 2A-B).

114 **Retrospective evaluation.** For the 48 exams analyzed (8 patients, 6 serial exams per
115 patient), 26 out of the total 37 stable cases were correctly identified as such by the model
116 upon correct choice of baseline, for a specificity of 0.70. For the progression cases, 8 out
117 of 11 cases were correctly identified, for a sensitivity of 0.72 (Fig. 2C-D).

118 **Model performance.** The automated volumetrics measurement showed variance
119 between 2-18% which, for all these cases, is below the 9 mL variability observed in
120 manual diameter calculation¹⁵. For all but one patient, AI progression assessment
121 discrepancies with gold standard were due to imaging acquisitions differing significantly
122 from those used during model training (GE 3D Cube). For the remaining case, a small
123 increase in tumor volume from 1.7 mL to 2.7 mL was likely below the threshold for human
124 reader detection^{4,15} and only later marked as progressive.

125 **System Usability.** All prospective results were available for the radiologist within 10
126 minutes of DICOM transmission (90th percentile: 4 minutes and 40 seconds). 72 out of
127 73 cases were processed correctly (success rate: 99%). Two radiologists reported that
128 the results and presentation were interpretable and would provide valuable decision
129 support for routine use if data selection could be automated.

130 **Discussion and Conclusions.**

131 This work presents the first steps toward clinical deployment of an AI-based solution for
132 real-time detection of glioma progression. Combining prospective and retrospective
133 cases, a final sensitivity of 0.72 and specificity of 0.75 was achieved. 90% of the
134 successfully processed cases were available in under 5 minutes and only one case in 73
135 failed to be processed. Results were presented via an easy-to-use, interactive, and
136 explainable environment. Radiologists could manually choose baseline cases, adjust
137 threshold for progression, visually inspect segmentations associated with the exam, and
138 graphically check how volume varied over time.

139 Discrepancies between model and ground truth were mostly due to imaging acquisitions
140 differing from those used during model training (Fig. 3). This should be mitigated through
141 additional training iterations on a more representative clinical cohort, with varied
142 acquisition parameters. Additionally, if inadequate data does arrive at the AI model, the
143 output should explicitly note that the model confidence is low as opposed to simply giving
144 the wrong answer.

145 A potential solution for both issues would involve a “human in the loop” scenario, where
146 images would be first sent to a dedicated QC team responsible for checking (input /
147 output) data quality prior to notifying the doctor about the results. This would yield a
148 database of common ‘real-life’ failures that future models could be trained on.
149 Additionally, QC team triage would mitigate the risk of unintended clinical consequences⁷
150 and increase trust from clinicians and patients.

151 The aforementioned performance nevertheless highlights the benefit of the automated
152 method to reduce variability in FLAIR lesion measurement and hence increase accuracy
153 in disease progression detection. To maximize applicability at all treatment timepoints,
154 the ability to change the baseline provides a needed function particularly in the context of
155 changing therapies or clinical trials.

156

157

158 References.

- 159 (1) van den Bent, M. J.; Wefel, J. S.; Schiff, D.; Taphoorn, M. J. B.; Jaeckle, K.; Junck, L.; Armstrong,
160 T.; Choucair, A.; Waldman, A. D.; Gorlia, T.; Chamberlain, M.; Baumert, B. G.; Vogelbaum, M. A.;
161 Macdonald, D. R.; Reardon, D. A.; Wen, P. Y.; Chang, S. M.; Jacobs, A. H. Response Assessment
162 in Neuro-Oncology (a Report of the RANO Group): Assessment of Outcome in Trials of Diffuse Low-
163 Grade Gliomas. *Lancet Oncol.* **2011**, *12* (6), 583–593. [https://doi.org/10.1016/S1470-](https://doi.org/10.1016/S1470-2045(11)70057-2)
164 [2045\(11\)70057-2](https://doi.org/10.1016/S1470-2045(11)70057-2).
- 165 (2) Villanueva-Meyer, J. E.; Mabray, M. C.; Cha, S. Current Clinical Brain Tumor Imaging. *Neurosurgery*
166 **2017**, *81* (3), 397–415. <https://doi.org/10.1093/neuros/nyx103>.
- 167 (3) Upadhyay, N.; Waldman, A. D. Conventional MRI Evaluation of Gliomas. *Br. J. Radiol.* **2011**, *84*
168 (special_issue_2), S107–S111. <https://doi.org/10.1259/bjr/65711810>.
- 169 (4) Huang, R. Y.; Young, R. J.; Ellingson, B. M.; Veeraraghavan, H.; Wang, W.; Tixier, F.; Um, H.;
170 Nawaz, R.; Luks, T.; Kim, J.; Gerstner, E. R.; Schiff, D.; Peters, K. B.; Mellinghoff, I. K.; Chang, S.
171 M.; Cloughesy, T. F.; Wen, P. Y. Volumetric Analysis of IDH-Mutant Lower-Grade Glioma: A Natural
172 History Study of Tumor Growth Rates before and after Treatment. *Neuro-Oncol.* **2020**, *22* (12),
173 1822–1830. <https://doi.org/10.1093/neuonc/noaa105>.
- 174 (5) Vaziri, S.; Lafontaine, M.; Olson, B.; Crane, J. C.; Chang, S.; Lupo, J.; Nelson, S. J. Rapid
175 Assessment of Lesion Volumes for Patients with Glioma Using the Smartbrush Software Package.
176 *Neuro-Oncol.* **2014**, *16* (suppl 5), v156–v156. <https://doi.org/10.1093/neuonc/nou264.77>.
- 177 (6) Zaharchuk, G.; Gong, E.; Wintermark, M.; Rubin, D.; Langlotz, C. P. Deep Learning in
178 Neuroradiology. *Am. J. Neuroradiol.* **2018**, *39* (10), 1776–1784. <https://doi.org/10.3174/ajnr.A5543>.
- 179 (7) Mongan, J.; Kohli, M. Artificial Intelligence and Human Life: Five Lessons for Radiology from the 737
180 MAX Disasters. *Radiol. Artif. Intell.* **2020**, *2* (2), e190111. <https://doi.org/10.1148/ryai.2020190111>.
- 181 (8) Zeng, K.; Bakas, S.; Sotiras, A.; Akbari, H.; Rozycki, M.; Rathore, S.; Pati, S.; Davatzikos, C.
182 Segmentation of Gliomas in Pre-Operative and Post-Operative Multimodal Magnetic Resonance
183 Imaging Volumes Based on a Hybrid Generative-Discriminative Framework. In *Brainlesion: Glioma,*
184 *Multiple Sclerosis, Stroke and Traumatic Brain Injuries*; Crimi, A., Menze, B., Maier, O., Reyes, M.,
185 Winzeck, S., Handels, H., Eds.; Lecture Notes in Computer Science; Springer International
186 Publishing: Cham, 2016; Vol. 10154, pp 184–194. https://doi.org/10.1007/978-3-319-55524-9_18.
- 187 (9) Fedorov, A.; Beichel, R.; Kalpathy-Cramer, J.; Finet, J.; Fillion-Robin, J.-C.; Pujol, S.; Bauer, C.;
188 Jennings, D.; Fennessy, F.; Sonka, M.; Buatti, J.; Aylward, S.; Miller, J. V.; Pieper, S.; Kikinis, R. 3D
189 Slicer as an Image Computing Platform for the Quantitative Imaging Network. *Magn. Reson. Imaging*
190 **2012**, *30* (9), 1323–1341. <https://doi.org/10.1016/j.mri.2012.05.001>.
- 191 (10) Myronenko, A. 3D MRI Brain Tumor Segmentation Using Autoencoder Regularization; Springer,
192 2018; pp 311–320.
- 193 (11) Villanueva-Meyer, J.; Damasceno, P. F.; LaFontaine, M.; Hawkins, J.; Luks, T.; Crane, J.; Lupo, J.
194 Integrating Automated Lesion Segmentations from Single-Images into Routine Clinical Workflow for
195 Volumetric Response Assessment. *Neuro-Oncol.* **2020**, *22* (Supplement_2), ii157–ii157.
196 <https://doi.org/10.1093/neuonc/noaa215.657>.
- 197 (12) *Clara Deploy*; NVIDIA.
- 198 (13) Marcus, D. S.; Olsen, T.; Ramaratnam, M.; Buckner, R. L. XNAT: A Software Framework for
199 Managing Neuroimaging Laboratory Data. In *Proceedings of the 12th annual meeting of the*
200 *organization for human brain mapping, Florence*; 2006.
- 201 (14) Urban, T.; Ziegler, E.; Lewis, R.; Hafey, C.; Sadow, C.; Abbeele, A. D. V. den; Harris, G. J.
202 LesionTracker: Extensible Open-Source Zero-Footprint Web Viewer for Cancer Imaging Research
203 and Clinical Trials. *Cancer Res.* **2017**, *77* (21), e119–e122. [https://doi.org/10.1158/0008-5472.CAN-](https://doi.org/10.1158/0008-5472.CAN-17-0334)
204 [17-0334](https://doi.org/10.1158/0008-5472.CAN-17-0334).
- 205 (15) Sorensen, A. G.; Patel, S.; Harmath, C.; Bridges, S.; Synnott, J.; Sievers, A.; Yoon, Y.-H.; Lee, E.
206 J.; Yang, M. C.; Lewis, R. F.; Harris, G. J.; Lev, M.; Schaefer, P. W.; Buchbinder, B. R.; Barest, G.;
207 Yamada, K.; Ponzio, J.; Kwon, H. Y.; Gemmete, J.; Farkas, J.; Tievsky, A. L.; Ziegler, R. B.; Salhus,
208 M. R. C.; Weisskoff, R. Comparison of Diameter and Perimeter Methods for Tumor Volume
209 Calculation. *J. Clin. Oncol.* **2001**, *19* (2), 551–557. <https://doi.org/10.1200/JCO.2001.19.2.551>.
- 210 (16) Flores, M.; Dayan, I.; Roth, H.; Zhong, A.; Harouni, A.; Gentili, A.; Abidin, A.; Liu, A.; Costa, A.;
211 Wood, B.; Tsai, C.-S.; Wang, C.-H.; Hsu, C.-N.; Lee, C.; Ruan, C.; Xu, D.; Wu, D.; Huang, E.;

212 Kitamura, F.; Lacey, G.; Corradi, G. C. de A.; Shin, H.-H.; Obinata, H.; Ren, H.; Crane, J.; Tetreault,
213 J.; Guan, J.; Garrett, J.; Park, J. G.; Dreyer, K.; Juluru, K.; Kersten, K.; Rockenbach, M. A. B. C.;
214 Linguraru, M.; Haider, M.; AbdelMaseeh, M.; Rieke, N.; Damasceno, P.; Silva, P. M. C. e; Wang, P.;
215 Xu, S.; Kawano, S.; Sriswa, S.; Park, S. Y.; Grist, T.; Buch, V.; Jantarabekul, W.; Wang, W.; Tak,
216 W. Y.; Li, X.; Lin, X.; Kwon, F.; Gilbert, F.; Kaggie, J.; Li, Q.; Quraini, A.; Feng, A.; Priest, A.; Turkbey,
217 B.; Glicksberg, B.; Bizzo, B.; Kim, B. S.; Tor-Diez, C.; Lee, C.-C.; Hsu, C.-J.; Lin, C.; Lai, C.-L.; Hess,
218 C.; Compas, C.; Bhatia, D.; Oermann, E.; Leibovitz, E.; Sasaki, H.; Mori, H.; Yang, I.; Sohn, J. H.;
219 Murthy, K. N. K.; Fu, L.-C.; Mendonça, M. R. F. de; Fralick, M.; Kang, M. K.; Adil, M.; Gangai, N.;
220 Vateekul, P.; Elnajjar, P.; Hickman, S.; Majumdar, S.; McLeod, S.; Reed, S.; Graf, S.; Harmon, S.;
221 Kodama, T.; Puthanakit, T.; Mazzulli, T.; Lavor, V. de L.; Rakvongthai, Y.; Lee, Y. R.; Wen, Y.
222 *Federated Learning Used for Predicting Outcomes in SARS-COV-2 Patients*; preprint; In Review,
223 2021. <https://doi.org/10.21203/rs.3.rs-126892/v1>.
224

225

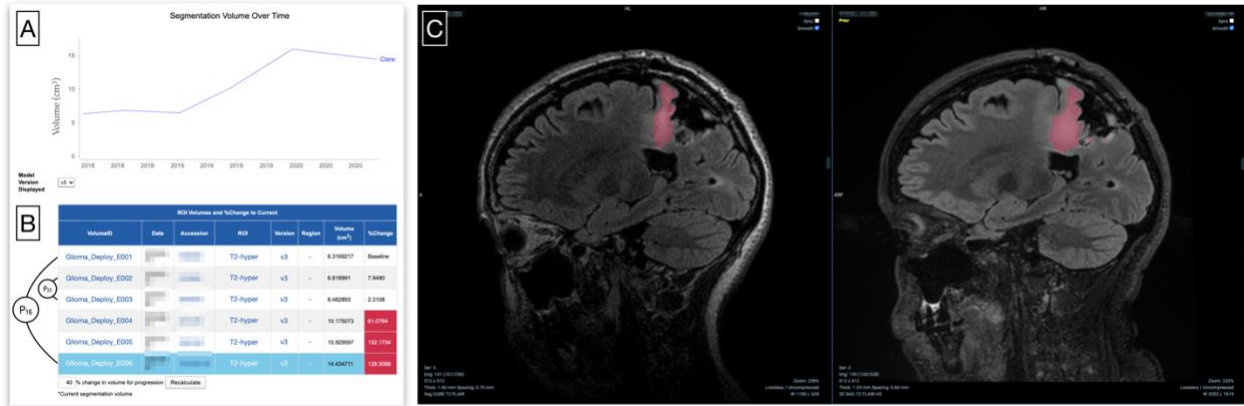
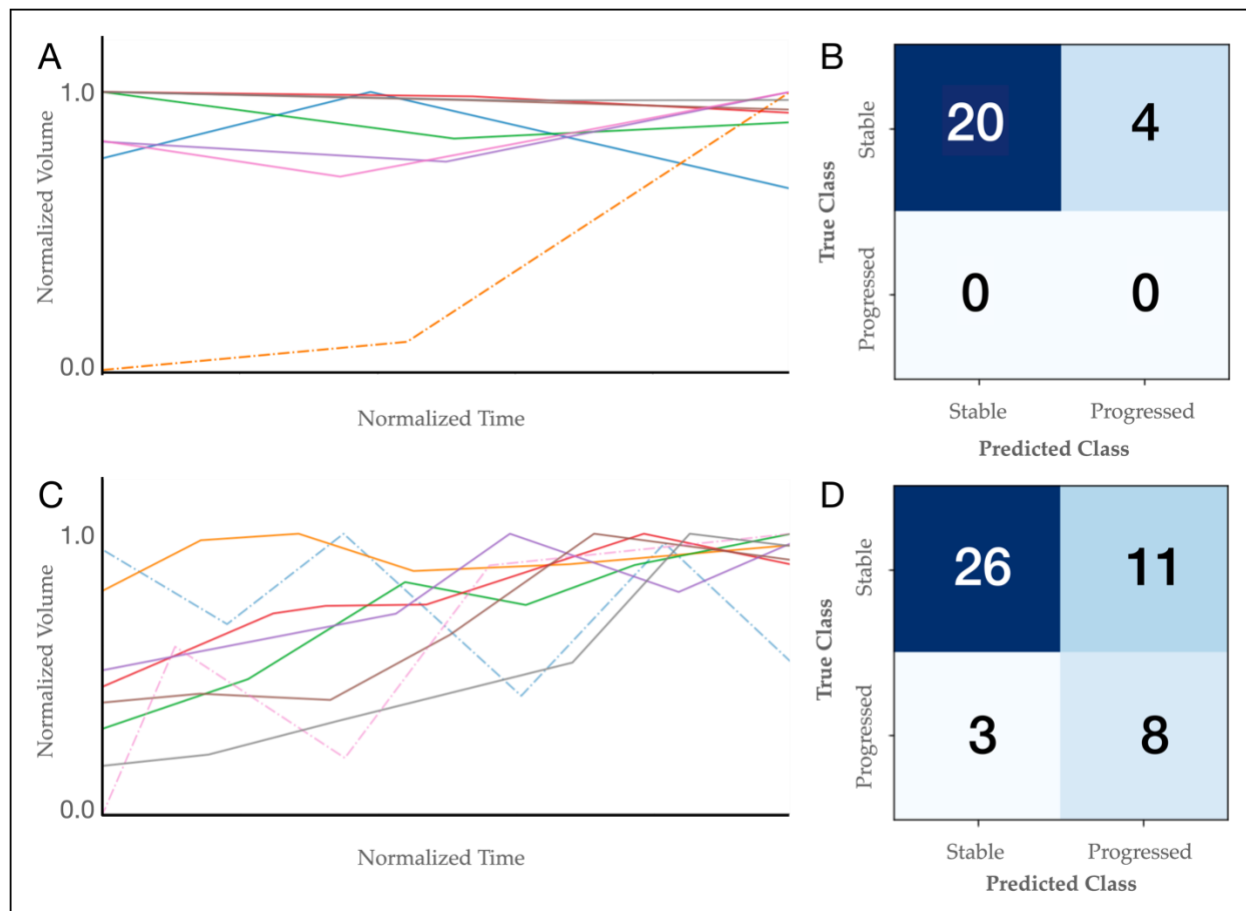


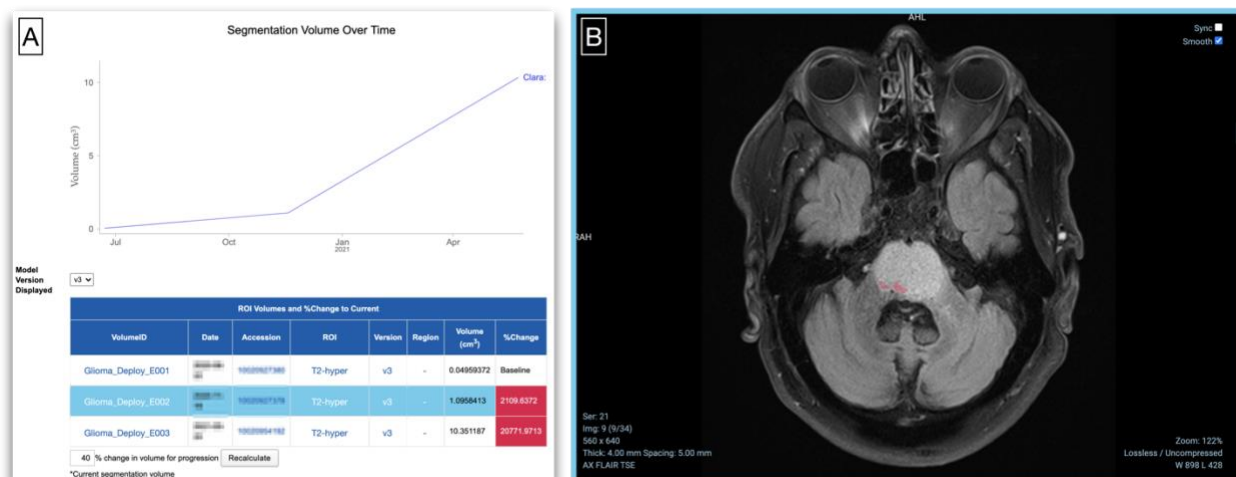
Figure 1. Deployed tumor progression detection pipeline in action. Visual representation of glioma volume over time (A) followed by a table of individual exams (B) containing tumor volumes, and percentage change with respect to baseline (highlighted in blue). Arrow connections exemplify two *baseline-to-follow-up* pairs. Exams whose tumor volume exceeds a chosen threshold (defaulted to 40%) are marked in red to indicate progression. Interactive options include: i) selection of AI model version; ii) selection of baseline exam; iii) manipulation of threshold percentage value; and iv) ability to choose which exam to be visualized. (C) in-browser visualization of two exams – baseline (left) and latest (right) – as well as AI-generated glioma segmentation (pink). The volume increase is clear and facilitates the visual confirmation of progression or identification of model mistakes.



226 **Figure 2. Glioma segmentation volumes and progression detection.** Glioma volumes
227 as a function of time for prospective (A) and retrospective (C) cases. Dashed lines show
228 cases where the model made significant mistakes due to acquisition protocol issues.
229 Volumes were normalized by the maximum volume observed to facilitate collected
230 visualization.

231

232



233

234 **Figure 3. An example of an AI tumor progression misclassification.** (A) Glioma
235 volume over time for a patient whose prior visits came from another institution. Because
236 the prior data was a 2D acquisition – and therefore different from what the model was
237 trained on – the AI model was not able to segment the glioma clearly present in the image
238 (B).