



**Supplementary Figure S1.** Segmentation performance of nnU-Net models by tumor size across different datasets. (A) Tumor segmentation performance of nnU-Net models on baseline (BL), (B,E) after two cycles (C2), (C,F) and after four cycles (C4) test sets across different tumor sizes. Horizontal bar indicates significant difference in paired Wilcoxon signed rank test ( $p < 0.05$ ). (D), the detailed quantitative results used for the boxplots in (A). (G), the detailed quantitative results used for the boxplots in (B) and (E). (H), the detailed quantitative results used for the boxplots in (C) and (F).

- nnU-Net model training configuration and procedure

The source code of nnU-Net was downloaded from <https://github.com/MIC-DKFZ/nnUNet>. All trainings used the default set-up<sup>36</sup>, and important parameters included the following: data augmentation was implemented using the batchgenerators framework that was included in nnU-Net (a series of augmentations involving spatial, color, noise, and crop processes); the training length was fixed at 1000 epochs, with 250 training iterations per epoch; and the learning rate was 0.01 initially and then was gradually reduced using the “polyLR” schedule<sup>55</sup>. All training used the 2D nnU-Net framework and full-resolution 3D nnU-Net framework and then ensembled these frameworks for inferences. A five-fold cross-validation was performed on the development dataset (training + validation, randomly split at 4:1).