



**HAL**  
open science

## Impact of the reperfusion status for predicting the final stroke infarct using deep learning

Noëlie Debs, Tae-Hee Cho, David Rousseau, Yves Berthezène, Marielle Buisson, Omer Eker, Laura Mechtouff, Norbert Nighoghossian, Michel Ovize, Carole Frindel

### ► To cite this version:

Noëlie Debs, Tae-Hee Cho, David Rousseau, Yves Berthezène, Marielle Buisson, et al.. Impact of the reperfusion status for predicting the final stroke infarct using deep learning. *Neuroimage-Clinical*, 2021, 29, pp.102548. 10.1016/j.nicl.2020.102548 . hal-03113704

**HAL Id: hal-03113704**

**<https://hal.inrae.fr/hal-03113704v1>**

Submitted on 3 Feb 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

# Impact of the reperfusion status for predicting the final stroke infarct using deep learning

Noëlie Debs<sup>a</sup>, Tae-Hee Cho<sup>a,b</sup>, David Rousseau<sup>c</sup>, Yves Berthezène<sup>a,d</sup>, Marielle Buisson<sup>e</sup>, Omer Eker<sup>a,d</sup>, Laura Mechtouff<sup>b,e</sup>, Norbert Nighoghossian<sup>a,b</sup>, Michel Ovize<sup>d</sup>, Carole Frindel<sup>a,\*</sup>

<sup>a</sup>CREATIS, CNRS UMR-5220, INSERM U1206, Université Lyon 1, INSA Lyon, Villeurbanne, France

<sup>b</sup>Department of Vascular Neurology, Hospices Civils de Lyon, Lyon, France

<sup>c</sup>LARIS, UMR IRHS INRA, Université d'Angers, Angers, France

<sup>d</sup>Department of Neuroradiology, Hospices Civils de Lyon, Lyon, France

<sup>e</sup>Department of Cardiology, Clinical Investigation Center, CarMeN INSERM U1060, INRA U1397, INSA Lyon, Université Lyon 1, Hospices Civils de Lyon, Lyon, France

---

## Abstract

### Background

Predictive maps of the final infarct may help therapeutic decisions in acute ischemic stroke patients. Our objectives were to assess whether integrating the reperfusion status into deep learning models would improve their performance, and to compare them to current clinical prediction methods.

### Methods

We trained and tested convolutional neural networks (CNNs) to predict the final infarct in acute ischemic stroke patients treated by thrombectomy in our center. When training the CNNs, non-reperfused patients from a non-thrombectomized cohort were added to the training set to increase the size of this group. Baseline diffusion and perfusion-weighted magnetic resonance imaging (MRI) were used as inputs, and the lesion segmented on day-6 MRI served as the ground truth for the final infarct. The cohort was dichotomized into two subsets, reperfused and non-reperfused patients, from which reperfusion status specific CNNs were developed and compared to one another, and to the clinically-used perfusion-diffusion mismatch model. Evaluation metrics included the Dice similarity coefficient (DSC), precision, recall, volumetric similarity, Hausdorff distance and area-under-the-curve (AUC).

### Results

We analyzed 109 patients, including 35 without reperfusion. The highest DSC were achieved in both reperfused and non-reperfused patients (DSC =  $0.44 \pm 0.25$  and  $0.47 \pm 0.17$ , respectively) when using the corresponding reperfusion status-specific CNN. CNN-based models achieved higher DSC and AUC values compared to those of perfusion-diffusion

mismatch models (reperfused patients:  $AUC = 0.87 \pm 0.13$  vs  $0.79 \pm 0.17$ ,  $P < 0.001$ ; non-reperfused patients:  $AUC = 0.81 \pm 0.13$  vs  $0.73 \pm 0.14$ ,  $P < 0.01$ , in CNN vs perfusion-diffusion mismatch models, respectively).

## Conclusion

The performance of deep learning models improved when the reperfusion status was incorporated in their training. CNN-based models outperformed the clinically-used perfusion-diffusion mismatch model. Comparing the predicted infarct in case of successful *vs* failed reperfusion may help in estimating the treatment effect and guiding therapeutic decisions in selected patients.

*Key words:* Stroke, Prediction, Convolutional neural network, Magnetic resonance imaging, Reperfusion status

---

## 1. Introduction

Early reperfusion, by means of intravenous thrombolysis or thrombectomy, is the main therapeutic goal in acute ischemic stroke (Powers et al., 2019). Acute treatment decisions have increasingly incorporated advanced neuroimaging to estimate patients' prognosis and likelihood of benefiting from revascularization procedures (Albers et al., 2018; Nogueira et al., 2018). Currently, both computed-tomography (CT) and Magnetic Resonance Imaging (MRI) entail threshold-based methods to delineate the still salvageable brain (i.e. ischemic penumbra) from the already lost tissue (infarct core). Specifically in MRI, criteria for the infarct core is based on Apparent Diffusion Coefficient (ADC) extracted from Diffusion-Weighted Imaging (DWI), and criteria for the ischemic penumbra is based on Time to maximum of the residue function ( $T_{max}$ ) extracted from perfusion-weighted imaging. Precisely, infarct core is defined as ADC voxel values  $< 600 \sim 620 \times 10^{-6}$  mm<sup>2</sup>/s, and ischemic penumbra is defined as  $T_{max}$  voxel values  $> 6$  seconds (Kidwell et al., 2013; Olivot et al., 2009). Patients with a large penumbra and limited ischemic core (so-called 'target mismatch' profile) have a high probability of benefiting from reperfusion, even in late time windows (Albers et al., 2018;

---

\*Corresponding Author

*Email addresses:* [noelie.debs@creatis.insa-lyon.fr](mailto:noelie.debs@creatis.insa-lyon.fr) (Noëlie Debs), [tae-hee.cho@chu-lyon.fr](mailto:tae-hee.cho@chu-lyon.fr) (Tae-Hee Cho), [david.rousseau@univ-angers.fr](mailto:david.rousseau@univ-angers.fr) (David Rousseau), [yves.berthezene@chu-lyon.fr](mailto:yves.berthezene@chu-lyon.fr) (Yves Berthezène), [marielle.buisson01@chu-lyon.fr](mailto:marielle.buisson01@chu-lyon.fr) (Marielle Buisson), [omer.eker@chu-lyon.fr](mailto:omer.eker@chu-lyon.fr) (Omer Eker), [laura.mechtouff@chu-lyon.fr](mailto:laura.mechtouff@chu-lyon.fr) (Laura Mechtouff), [norbert.nighoghossian@chu-lyon.fr](mailto:norbert.nighoghossian@chu-lyon.fr) (Norbert Nighoghossian), [michel.ovize@chu-lyon.fr](mailto:michel.ovize@chu-lyon.fr) (Michel Ovize), [carole.frindel@creatis.insa-lyon.fr](mailto:carole.frindel@creatis.insa-lyon.fr) (Carole Frindel)

16 Nogueira et al., 2018). However, these fixed-threshold methods may fail to capture the  
17 significant interindividual heterogeneity observed in stroke progression (Rekik et al., 2012).  
18 While the clinical and imaging characteristics of some patients may clearly indicate urgent  
19 reperfusion therapies, the benefit/risk balance in others can appear more uncertain. Thus,  
20 personalized probability maps of the final infarct would be of high clinical value to guide  
21 acute revascularization decisions and possibly help evaluate novel neuroprotective strategies.

22 Convolutional neural networks (CNNs), a subtype of machine learning, are flexible, data-  
23 driven methods capable of automatic non-linear feature extraction, with promising results in  
24 stroke lesion segmentation (Qiu et al., 2020). A well-acknowledged limitation of CNNs is the  
25 large quantity of data required for their training and validation. Only a limited number of  
26 studies, with heterogeneous treatment paradigms and evaluations metrics, have evaluated  
27 CNNs for the prediction of the final stroke lesion from baseline MRI (Nielsen et al., 2018;  
28 Pinto et al., 2018; Winzeck et al., 2018; Yu et al., 2020) or CT (Robben et al., 2020). Sample  
29 size and performance were modest ( $\sim 50$  to  $\sim 200$  patients, Dice similarity coefficient  $\sim 0.50$  or  
30 lower), illustrating both the inherent difficulty of prediction tasks and scarcity of high-quality  
31 data, compared to simpler image segmentation tasks.

32 In the present work, we evaluated the impact of integrating the reperfusion status on the  
33 performance of CNNs for predicting the final infarct in patients with proximal intracranial  
34 occlusions treated by thrombectomy. Reperfusion is the single most important clinical  
35 metadata known to influence the progression of ischemic lesions from the baseline imaging  
36 (used as inputs to CNN) to the final infarct (Tsai and Albers, 2015). Previous studies  
37 have investigated direct integration of the reperfusion status during the learning process of  
38 CNN-based methods (Pinto et al., 2018; Robben et al., 2020). Another dichotomized the  
39 training set according to the reperfusion status with random forest-based methods (McKinley  
40 et al., 2017), but has not been evaluated with CNNs. We hypothesized that training CNNs  
41 from reperfusion status-specific subcohorts could improve their performance. Our objectives  
42 were: (1) to assess the impact of the reperfusion status on CNN-based predictive models; (2)  
43 to compare the predictive value of these CNNs against the threshold-based perfusion-diffusion  
44 mismatch models. An ancillary objective was to assess the relative predictive importance of  
45 the MRI inputs with an ablation study.

## 2. Material and methods

### 2.1. Data

We describe the HIBISCUS-STROKE and I-KNOW cohorts, from which the final stroke lesion was assessed. This section details the MRI protocol, patient inclusion criteria and image post-processing steps (upsampling, registration, normalization).

#### 2.1.1. Patients and imaging protocol

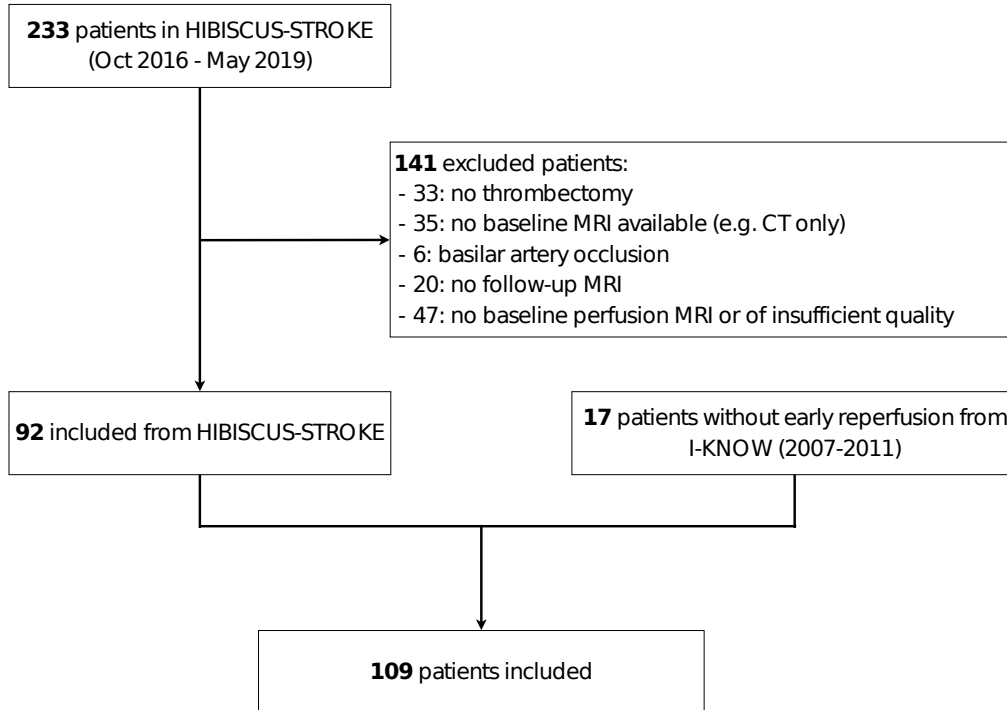
Patients were included from the HIBISCUS-STROKE and I-KNOW cohorts. HIBISCUS-STROKE is an ongoing monocentric observational cohort enrolling patients with a large intracranial artery occlusion treated by thrombectomy, following a baseline diffusion-perfusion MRI. I-KNOW (2007-2011) was a prospective multicenter observational study of stroke patients with both admission and several follow-up MRI. A subset of these patients underwent an acute follow-up perfusion MRI ( $\sim 3$  hours from the baseline MRI) to assess early reperfusion (Cho et al., 2015). In total, 109 patients were analyzed as shown in Figure 1. Early reperfusion was observed in 74 patients, while 35 had no reperfusion (17 from I-KNOW and 18 from HIBISCUS-STROKE). Baseline patients' characteristics are summarized in Appendix A.2. The inclusion and exclusion criteria for both cohorts are detailed in Appendix A.1. All patients from both cohorts gave their informed consent and the imaging protocol was approved by the regional ethics committee.

In both cohorts, all patients underwent the following MRI protocol on admission: diffusion-weighted-imaging (DWI), T2-weighted fluid-attenuated-inversion-recovery (FLAIR), T2-gradient echo, MR-angiography and dynamic susceptibility-contrast perfusion imaging (DSC-PWI). A follow-up FLAIR was performed several days after admission (specifically, 6 and 30 days in HIBISCUS-STROKE and I-KNOW, respectively). MRI acquisition parameters are described in Appendix A.3.

#### 2.1.2. Image post-processing

Parametric maps were extracted from the DSC-PWI by circular singular value decomposition of the tissue concentration curves (Olea Sphere, Olea Medical, La Ciotat, France): cerebral blood flow (CBF), cerebral blood volume (CBV), mean transit time (MTT), time to maximum ( $T_{max}$ ) and time to peak (TTP). Lesions on the baseline DWI and final FLAIR were segmented by an expert (THC) blinded to the clinical data with a semi-automated method (3D Slicer, <https://www.slicer.org/>). Specifically, a region-of-interest-controlled thresholding was used with manual corrections when required (for the DWI lesion, an ADC upper threshold of  $620 \times 10^{-6} \text{ mm}^2/\text{s}$  was used).

Figure 1: Patient inclusion flowchart.



79 Images from HIBISCUS-STROKE were coregistered within subjects to the baseline DWI  
80 MRI using non-linear registration with Ants (Avants et al., 2011). Images from I-KNOW were  
81 coregistered within subjects to the PWI-DSC MRI (matrix 128x128) using affine registration  
82 with Statistical Parametric Mapping 8. Once co-registration was performed, HIBISCUS-  
83 STROKE patients had all MRI slices of size 192x192 compared to 128x128 for I-KNOW  
84 patients. As I-KNOW patients were largely in the minority (17 patients out of the 109  
85 total patients), we up-sampled the images of I-KNOW patients to 192x192. The skull from  
86 all patients was removed using FSL (Smith et al., 2001). Finally, images were normalized  
87 between 0 and 1 to ensure inter-patient standardization.

## 88 2.2. Early reperfusion and training sets

89 We describe reperfusion criteria and we define the training sets.

### 90 2.2.1. Assessment of early reperfusion

91 In HIBISCUS-STROKE, early reperfusion was assessed at the end of the endovascular  
92 procedure with the modified Thrombolysis in Cerebral Infarction (mTICI) score (grade 0:

93 no reperfusion; grade 1: anterograde reperfusion past the initial occlusion, but limited distal  
94 branch filling with little or slow distal reperfusion; grade 2a: anterograde reperfusion of less  
95 than half of the occluded target artery previously ischemic territory; grade 2b: anterograde  
96 reperfusion of more than half of the previously occluded target artery ischemic territory;  
97 grade 2c: near complete reperfusion, i.e.  $>90\%$  but less than mTICI 3; grade 3: complete  
98 anterograde reperfusion) (Zaidat et al., 2013). Angiographic reperfusion was defined by  
99 mTICI scores of 2b-3, while patients without reperfusion had mTICI scores of 0-2a.

100 In I-KNOW, no patient was treated by endovascular procedures. Early reperfusion was  
101 assessed 3 hours after the first MRI (H3) and was defined as voxels with  $T_{max} \geq 6$  s at  
102 admission (H0) and  $T_{max} < 6$  s at H3. Acute reperfusion was defined by a reperfusion ratio  
103 (volume of reperfused voxels at H3/perfusion lesion volume at H0) of  $\geq 50\%$ .

### 104 2.2.2. Training sets

105 Three distinct training sets and corresponding models were built to assess the impact  
106 of reperfusion on the accuracy of final infarct prediction: a ‘general’ model, trained on the  
107 entire cohort irrespective of the reperfusion status (*all* training set); a ‘reperfused’ model,  
108 trained only with reperfused patients (*reperfused* training set); a ‘non-reperfused’ model,  
109 trained only with non-reperfused patients (*non reperfused* training set). Given the high  
110 rate of angiographic success in patients treated by thrombectomy (mTICI score of 2b-3 in  
111  $>70\%$  of patients) (Goyal et al., 2016), we expected a limited proportion of non-reperfused  
112 patients from HIBISCUS-STROKE. We thus included patients without early reperfusion  
113 from I-KNOW (identified by the H3 perfusion MRI follow-up) in order to improve this  
114 imbalance. I-KNOW patients were only included in the training set of the general and the  
115 non-reperfused models, but were not included in any testing set.

### 116 2.3. Proposed CNN architecture

117 We used a U-Net architecture, a multi-scale network that has already shown its potential  
118 for infarct prediction tasks (Winzeck et al., 2018; Yu et al., 2020). Perfusion and diffusion  
119 MRI were used as inputs, as both modalities are complementary to evaluate the risk of  
120 infarction (Barber et al., 1998). More precisely, a total of five inputs were used : DWI and  
121 ADC for diffusion MRI, as well as  $T_{max}$ , CBF and CBV for perfusion MRI. Previous studies  
122 in other medical applications have evaluated methods for combining the input data into  
123 CNNs, showing the merit of late fusion strategies (Aygün et al., 2018; Dolz et al., 2018a,b; Nie  
124 et al., 2016). Late fusion incorporates each input independently into distinct convolutional  
125 branches, subsequently merging features at a higher level. This strategy was chosen for its

126 potential to better integrate each MRI input and the impact of reperfusion status. The  
 127 comparison of the early and late fusion strategies is presented in Appendix C.

128 The five inputs (DWI, ADC,  $T_{max}$ , CBV, CBF) were fed into our late fusion network of  
 129 5 distinct convolution branches. The proposed architecture is depicted in Figure 2, and its  
 130 encoding layers are detailed in Table 1. Each input consisted of whole 2D images (192x192).  
 131 No patches were used in order to secure a large spatial context for lesion prediction. The  
 132 network produced probability maps with 3 classes: lesion, healthy tissue, background. The  
 133 lesion probability map was thresholded at 0.5 to define the final infarct. Training and  
 134 configuration of the network are detailed in Appendix B.

Table 1: Encoding layers of the proposed late fusion U-net. The encoder is composed of 5 convolution blocks (Conv Block), maxpooling operations (2D MaxPooling) and dropout. The Conv Block is made of: 2D convolution (3\*3)+ batch normalization + 2D convolution (3\*3)+ batch normalization.

Layer (type)	Output shape
Conv Block 1	192*192*8
2D MaxPooling	96*96*8
Conv Block 2	96*96*16
2D MaxPooling	48*48*16
Conv Block 3	48*48*32
2D MaxPooling	24*24*32
Conv Block 4	24*24*64
Dropout + 2D Maxpooling	12*12*64
Conv Block 5 + Dropout	12*12*128
Concatenation	12*12*640

## 135 2.4. Evaluation

### 136 2.4.1. Ground truth

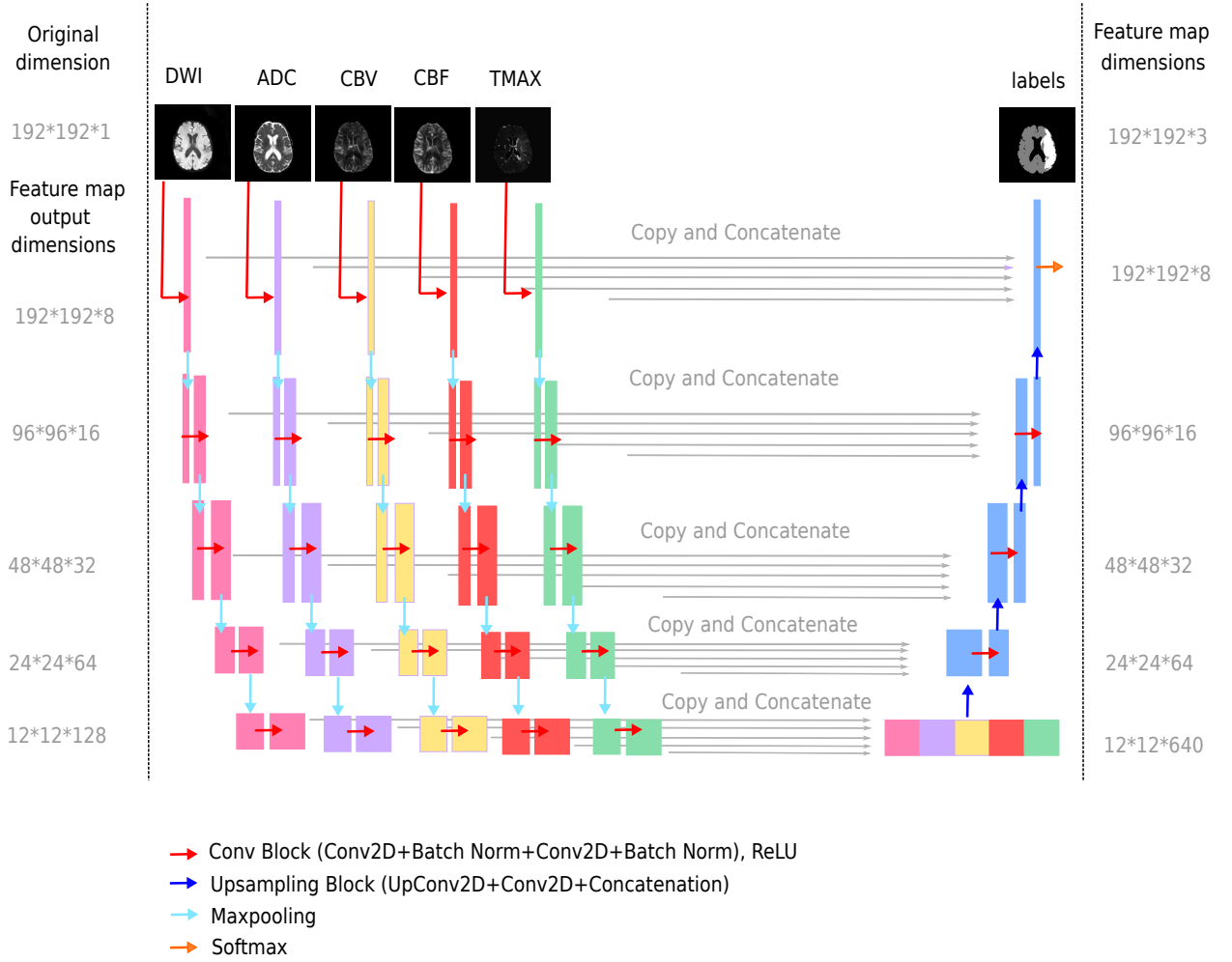
137 The final lesion is given by the FLAIR MRI, which was performed several days after  
 138 admission (specifically, 6 and 30 days in HIBISCUS-STROKE and I-KNOW, respectively).  
 139 The brain mask and the final lesion on the FLAIR MRI were segmented by experts using  
 140 semi-automatic intensity-based thresholding. The ground truth for each patient was therefore  
 141 a 3D mask with 3 classes : one class for background, one class for healthy tissues and one  
 142 class for the lesion.

### 143 2.4.2. Metrics

144 Standard metrics for assessing image segmentation/prediction tasks were used: the Dice  
 145 similarity coefficient (DSC), precision, recall, volumetric similarity (VS), and Hausdorff



Figure 2: Overview of the proposed deep learning architecture. Top left: The network takes five MRI images (2D slices from DWI, ADC, CBV, CBF,  $T_{max}$  volumes) as input. Below: Each input image is processed independently on 5 separate branches. Pink, purple, yellow, red and green feature maps result from 2D-convolutions and maxpooling. The output of the 5 branches are then concatenated, and upsampled through 2D-deconvolution layers. The network produces an output map with 3 classes (lesion, healthy tissue and background). Top Right : The predicted lesion has to be compared to the true lesion from the final FLAIR.



146 distance (HD) (Taha and Hanbury, 2015). The DSC measures the relative overlap of the  
 147 prediction with the ground truth ( $TP$ ,  $FN$  and  $FP$  are respectively the true positive, false  
 148 negative and false positive voxels):

$$DSC = \frac{2 \cdot TP}{FN + FP + 2 \cdot TP}. \quad (1)$$

Precision (also know as positive predictive value) measures the percentage of voxels identified as lesion that have been classified correctly, while recall (also know as sensitivity) measures the percentage of actual lesion voxels that have been classified correctly:

$$Precision = \frac{TP}{TP + FP}, \quad (2)$$

$$Recall = \frac{TP}{TP + FN}. \quad (3)$$

149 The VS gives a relative ratio between the prediction and the ground truth volumes, without  
150 considering any overlap of the two volumes:

$$VS = 1 - \frac{|FN - FP|}{2 \cdot TP + FP + FN}, \quad (4)$$

151 The HD is a measure of the distance of the largest error between the prediction ( $A$ ) and  
152 ground truth ( $B$ ):

$$HD(A, B) = \max(h(A, B), h(B, A)) \quad \text{where} \quad h(A, B) = \max_{a \in A} \min_{b \in B} \|a - b\|. \quad (5)$$

153 The area-under-the-curve (AUC) is widely used in medical evaluation. Based on the ROC  
154 curve (Hajian-Tilaki, 2013), it provides an aggregated performance measure of an image  
155 modality or parametric map across all possible threshold values. However, the overwhelming  
156 number of non-infarcted voxels relative to infarcted ones can drive high AUC values while the  
157 extent and location of the infarct is poorly predicted (Jonsdottir et al., 2009). Several studies  
158 thus favored the DSC, which is more specific for lesion prediction (Winder et al., 2019; Yu  
159 et al., 2020). We presented AUC values in order to facilitate comparisons with some previous  
160 studies, notably when comparing CNN-based models and the clinical perfusion-diffusion  
161 mismatch model (Nielsen et al., 2018; Yu et al., 2020).

### 162 2.4.3. Perfusion-diffusion mismatch model

163 Our CNN-based predictive models were compared with the current reference method used  
164 in clinical practice. According to the perfusion-diffusion mismatch model, the projected final  
165 infarct can be defined as follows: (1) in reperfused patients, the final infarct is represented  
166 by the baseline diffusion lesion; (2) in non-reperfused patients, the final infarct is defined as  
167 the union of the acute diffusion lesion and the ischemic penumbra (voxels with a  $T_{max} > 6$   
168 seconds and normal DWI)(Olivot et al., 2009). The AUC of the perfusion-diffusion mismatch  
169 model to predict the final infarct was assessed in patients with and without reperfusion.

170 Non-infarcted voxels were those not included in the diffusion lesion in reperfused patients,  
171 and those not included in the diffusion  $\cup$  penumbra in non-reperfused patients. Infarcted  
172 voxels were the complementary voxels. The AUC was computed as in Jonsdottir et al. (2009).

#### 173 2.4.4. Statistical analyses

174 A two-sided Wilcoxon signed-rank test was performed in order to compare the perfor-  
175 mances of: (1) reperfused *vs* general, non-reperfused *vs* general and reperfused *vs* non-  
176 reperfused models; (2) models with all MRI inputs *vs* models with ablation of one or more  
177 MRI inputs; (3) reperfused model *vs* diffusion lesion model; (4) non-reperfused model *vs*  
178 diffusion  $\cup$  penumbra lesion model. Statistical analyses were performed using R version 3.5.1.

### 179 3. Results

#### 180 3.1. Performance of the general, reperfused and non-reperfused CNNs

181 The performances and comparisons of the general, reperfused and non-reperfused models  
182 tested in reperfused and non-reperfused patients are presented in Table 2.

183 Among reperfused patients, the non-reperfused model was inferior to either the reperfused  
184 or general models for all metrics except for precision (Tables 2-a and 2-b). The model seems to  
185 predict many false negative voxels (low recall), many outlier voxels (high hausdorff distance),  
186 and a different volume than expected (low VS). Conversely, no clear-cut performance difference  
187 was found between the reperfused and general models.

188 Among non-reperfused patients, the non-reperfused model had better or similar perfor-  
189 mance than the reperfused model for all metrics except for recall (Tables 2-c and 2-d). The  
190 model seems to predict the lesion well in terms of volume and localisation (high VS and high  
191 DSC), with few false positive voxels (high precision) but some false negative voxels (medium  
192 recall). No clear overall difference was observed between the non-reperfused and general  
193 models, or between the reperfused and general models.

194 The predicted infarct volumes were significantly larger with the non-reperfused compared  
195 to the reperfused model (39.7 mL (61.3-20) vs 17.5 mL (28-5.1),  $p = 4.5e - 16$  for the non-  
196 reperfused and reperfused models, respectively; median with interquartile range). Accordingly,  
197 significant differences of VS between these two models were observed (Tables 2-b and -d).  
198 Figure 3 illustrates and compares the output of the two CNNs (reperfused and non-reperfused)  
199 for two patients with distinct reperfusion status.

Table 2: Performance metrics of the general, reperfused and non-reperfused models among (a) **reperfused** and (c) **non-reperfused** patients (average values  $\pm$  standard deviation). Bold values correspond to the best value of the respective evaluation metric (column-wise). P-values from two-sided wilcoxon signed-rank tests comparing the general, reperfused and non-reperfused models among (b) **reperfused** and (d) **non-reperfused** patients. Bold values correspond to significant differences, with (\*) indicating  $P < 0.05$ , (\*\*) indicating  $P < 0.01$  and (\*\*\*) indicating  $P < 0.001$ . Note that tests were not corrected for multiple comparisons, and correspond to independent two-by-two comparisons

(a) *Performance metrics among reperfused patients*

Model	DSC	VS	Precision	Recall	HD
General	0.43 $\pm$ 0.24	0.69 $\pm$ 0.27	0.55 $\pm$ 0.28	0.43 $\pm$ 0.25	<b>33.23 <math>\pm</math> 15.6</b>
Reperfused	<b>0.44 <math>\pm</math> 0.25</b>	<b>0.70 <math>\pm</math> 0.27</b>	0.50 $\pm$ 0.27	<b>0.50 <math>\pm</math> 0.26</b>	38.58 $\pm$ 18.1
Non-reperfused	0.35 $\pm$ 0.21	0.57 $\pm$ 0.28	<b>0.60 <math>\pm</math> 0.25</b>	0.31 $\pm$ 0.24	40.05 $\pm$ 15.6

(b) *Model comparisons among reperfused patients*

Two-sided Test	DSC P-value	VS P-value	Precision P-value	Recall P-value	HD P-value
General vs Reperfused	0.43	0.53	<b>3.7e-6 (***)</b>	<b>1.4e-6 (***)</b>	<b>0.048 (*)</b>
General vs Non-Reperfused	<b>1.4e-8 (***)</b>	<b>4.3e-6 (***)</b>	<b>0.0069 (**)</b>	<b>1.0e-10 (***)</b>	<b>0.0041 (**)</b>
Reperfused vs Non-Reperfused	<b>2.3e-8 (***)</b>	<b>1.6e-5 (***)</b>	<b>2.9e-7 (***)</b>	<b>2.7e-11 (***)</b>	0.65

(c) *Model performance among non-reperfused patients*

Model	DSC	VS	Precision	Recall	HD
General	0.44 $\pm$ 0.21	0.66 $\pm$ 0.26	0.39 $\pm$ 0.25	0.63 $\pm$ 0.21	<b>30.61 <math>\pm</math> 16.1</b>
Reperfused	0.44 $\pm$ 0.22	0.63 $\pm$ 0.25	0.36 $\pm$ 0.23	<b>0.69 <math>\pm</math> 0.22</b>	44.53 $\pm$ 16.7
Non-reperfused	<b>0.47 <math>\pm</math> 0.17</b>	<b>0.74 <math>\pm</math> 0.13</b>	<b>0.49 <math>\pm</math> 0.22</b>	0.52 $\pm$ 0.21	37.70 $\pm$ 17.7

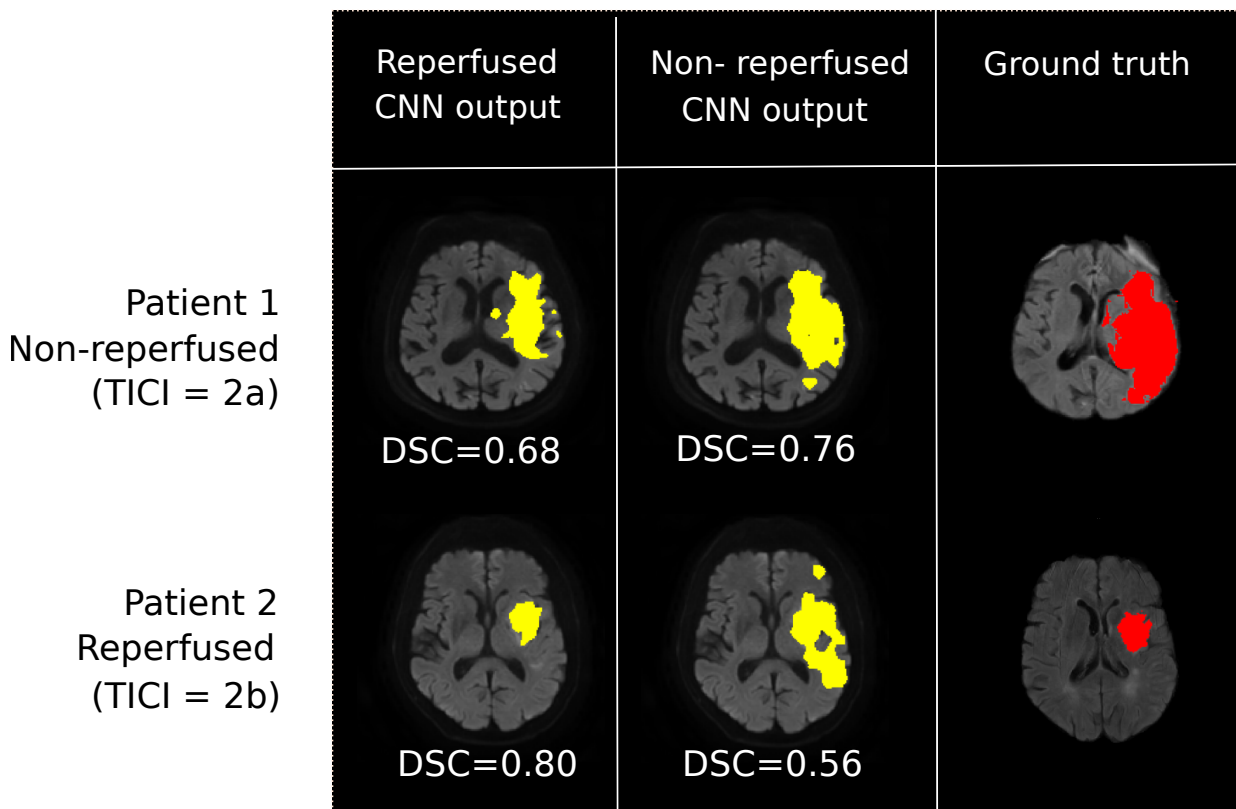
(d) *Model comparisons among non-reperfused patients*

Two-sided Test	DSC P-value	VS P-value	Precision P-value	Recall P-value	HD P-value
General vs Reperfused	0.93	0.55	0.13	<b>0.021 (*)</b>	<b>0.0023 (**)</b>
General vs Non-Reperfused	0.17	0.21	<b>0.0016 (**)</b>	<b>0.00084 (***)</b>	0.11
Reperfused vs Non-Reperfused	0.13	<b>0.034 (*)</b>	<b>0.00067 (***)</b>	<b>5.3e-5 (***)</b>	0.12

### 3.2. Comparison of CNN-based models and the perfusion-diffusion mismatch model

In both reperfused and non-reperfused patients, the DSC, VS and recall of CNN-based models were superior to those of the perfusion-diffusion mismatch models (Table 3). Final lesion predicted by CNNs are therefore more spatially and volumetrically coherent (high DSC and VS), and have fewer false negative voxels than the mismatch model. At the patient level, higher DSC values were achieved with CNN-based models in 68% and 89% of the reperfused and non-reperfused patients, respectively. Conversely, the precision of mismatch models was higher than that of CNN, suggesting more false positive voxels with the latter methods.

Figure 3: CNN-based predictions of the final infarct using the reperfused and non-reperfused models, applied in: patient 1 (no reperfusion, TICI=2a); patient 2 (reperfused, TICI=2b).



208 CNN-based models achieved higher AUC values compared to those of perfusion-diffusion  
 209 mismatch models (reperfused patients:  $0.87 \pm 0.13$  vs  $0.79 \pm 0.17$ ,  $P < 0.001$ ; non-reperfused  
 210 patients:  $0.81 \pm 0.13$  vs  $0.73 \pm 0.14$ ,  $P < 0.01$ , in CNN vs perfusion-diffusion mismatch models,  
 211 respectively). Cases illustrating successful or suboptimal outputs from CNN and mismatch  
 212 models are presented in Figure 4.

213 The comparison of CNNs and perfusion-diffusion mismatch model was included as the  
 214 latter remains the reference method in clinical practice. The mismatch model only provides  
 215 a crude threshold-based segmentation of baseline images, and may not match the feature  
 216 extraction potential of CNNs. Also, the mismatch model is only based on ADC and Tmax in  
 217 order to predict the final lesion outcome, whereas our model is based on more inputs (DWI,  
 218 ADC, Tmax, CBV, CBF).

Table 3: Comparison of CNN-based and perfusion-diffusion mismatch models. Among reperfused patients (upper rows), the CNN-based reperfused model was compared to the threshold-based diffusion lesion. Among non-reperfused patients (lower rows), the CNN-based non-reperfused model was compared to the threshold-based diffusion  $\cup$  penumbra lesion. Bold values correspond to the best value of the respective evaluation metric (column-wise). A two-sided wilcoxon signed-rank test was performed between the proposed models and the clinical models, with (.) indicating  $P < 0.10$ , (\*) indicating  $P < 0.05$ , (\*\*) indicating  $P < 0.01$  and (\*\*\*) indicating  $P < 0.001$ .

<i>Reperfused patients</i>					
Model	DSC	VS	Precision	Recall	HD
CNN	<b>0.44 <math>\pm</math> 0.21</b> (*)	<b>0.66 <math>\pm</math> 0.26</b> (***)	0.39 $\pm$ 0.25	<b>0.63 <math>\pm</math> 0.21</b> (***)	30.61 $\pm$ 16.1
Perfusion-diffusion mismatch	0.41 $\pm$ 0.23	0.56 $\pm$ 0.27	<b>0.71 <math>\pm</math> 0.31</b> (***)	0.33 $\pm$ 0.20	<b>19.34 <math>\pm</math> 10.3</b> (***)
<i>Non-reperfused patients</i>					
Model	DSC	VS	Precision	Recall	HD
CNN	<b>0.47 <math>\pm</math> 0.17</b> (***)	<b>0.74 <math>\pm</math> 0.13</b> (***)	0.49 $\pm$ 0.22	<b>0.52 <math>\pm</math> 0.21</b> (***)	<b>37.70 <math>\pm</math> 17.7</b> (***)
Perfusion-diffusion mismatch	0.26 $\pm$ 0.17	0.31 $\pm$ 0.21	<b>0.84 <math>\pm</math> 0.16</b> (***)	0.17 $\pm$ 0.13	69.15 $\pm$ 7.7

### 219 3.3. Value of the MRI inputs for predicting the final infarct

220 An ablation study was performed with the reperfused and non-reperfused models (tested  
 221 only in reperfused and non-reperfused patients, respectively) in order to evaluate the relative  
 222 importance of the different MRI inputs for predicting the final infarct. In both reperfused  
 223 and non-reperfused patients, the full CNN models (i.e. including DWI, ADC,  $T_{max}$ , CBF  
 224 and CBV) had similar performances compared to models without CBF and CBV, suggesting  
 225 these latter inputs had limited predictive value (lines 1 and 2 from Tables 4-a and 4-b).  
 226 Conversely, adding the diffusion data (DWI and ADC) to  $T_{max}$  maps significantly increased  
 227 the DSC of these CNNs. This performance increase was more pronounced among reperfused  
 228 patients compared to those without reperfusion.

## 229 4. Discussion

### 230 4.1. Impact of the reperfusion status on CNN performance

231 Our study showed that the performance of CNN-based models improved when trained  
 232 from reperfusion status-specific subgroups. The predicted lesion had better overlap (i.e.  
 233 higher DSC) with the final infarct in both reperfused and non-reperfused patients, when  
 234 using the corresponding reperfusion status-specific CNN.

235 Baseline imaging features do have significant predictive value, and CNNs trained without  
 236 data on reperfusion can successfully predict the final lesion in some patients (Yu et al.,  
 237 2020). This may in part reflect the mostly homogenous profile of patients currently treated  
 238 by thrombectomy (i.e. limited cerebral damage at baseline and successful reperfusion).  
 239 Indeed, the training set for our general CNN consisted of  $\sim 70\%$  of reperfused patients, and

Figure 4: Output predictions from CNN models compared with the PWI-DWI mismatch model. Five tested patients are shown: two successful cases when CNN models outperform PWI-DWI mismatch in reperfused and non-reperfused patients (patient A with  $TICI=2a$  and patient B with  $TICI=3$ ) and three difficult patients to predict, for both CNN and PWI-DWI mismatch models (patient C with  $TICI=2a$ , patient D with  $TICI=3$  and patient E with  $TICI=2b$ ). For each prediction model, patient-wide DSC is specified.

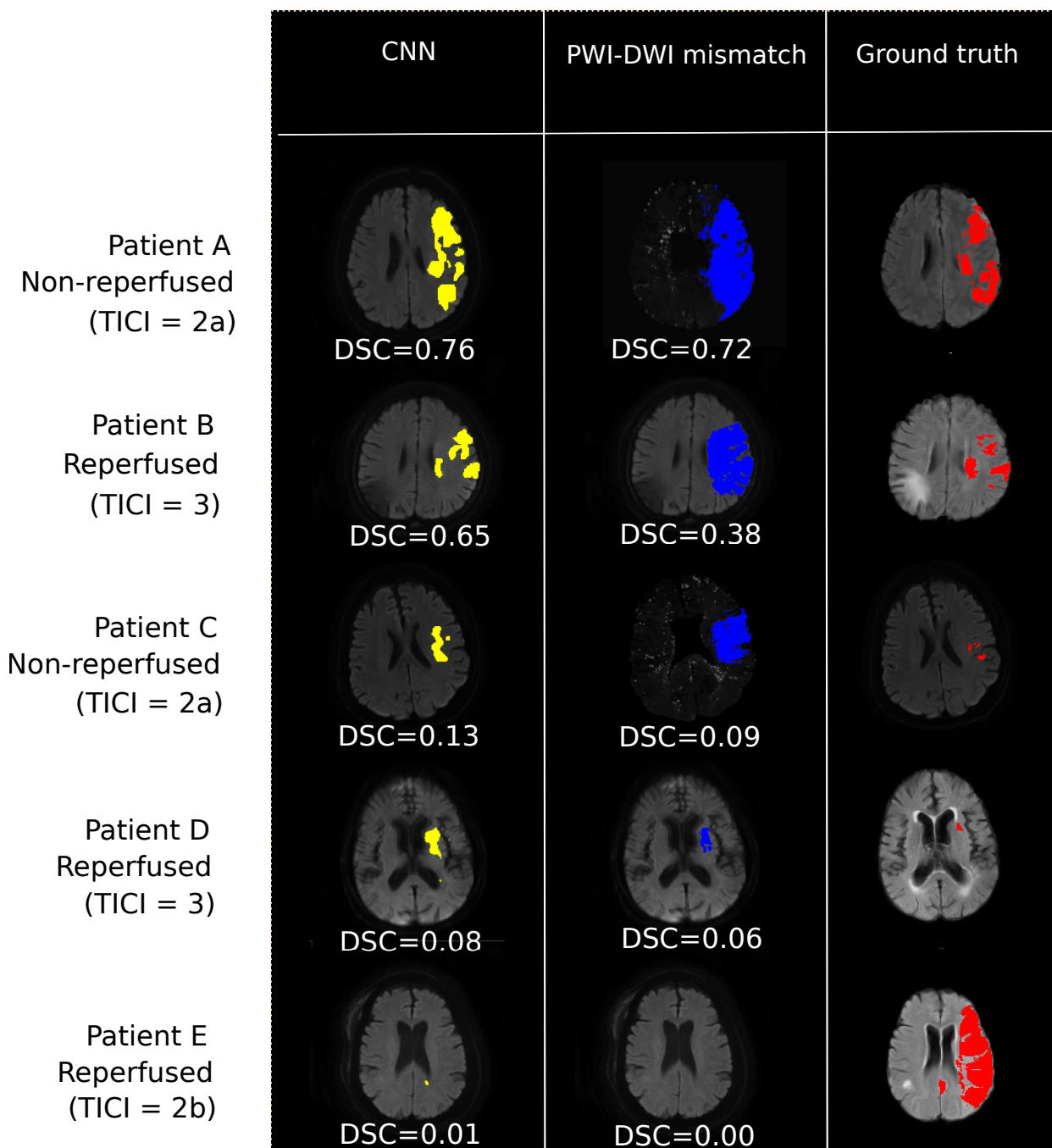


Table 4: Evaluation metrics of the reperfused and non-reperfused models after successive ablation of the MRI inputs, tested among **(a) reperfused** and **(b) non-reperfused** patients, respectively (average values  $\pm$  standard deviation). Bold values correspond to the best value of the respective evaluation metric (column-wise). A two-sided wilcoxon signed-rank test was performed between the full models with all 5 MRI inputs and the ablated ones, with (.) indicating  $P < 0.10$ , (\*) indicating  $P < 0.05$ , (\*\*) indicating  $P < 0.01$  and (\*\*\*) indicating  $P < 0.001$ .

<i>(a) Reperfused model: ablation study among reperfused patients</i>					
Input MRI	DSC	VS	Precision	Recall	HD
DWI+ADC+T <sub>max</sub> +CBF+CBV	<b>0.44 <math>\pm</math> 0.21</b>	0.66 $\pm$ 0.26	0.39 $\pm$ 0.25	<b>0.63 <math>\pm</math> 0.21</b>	30.61 $\pm$ 16.1
DWI+ADC+T <sub>max</sub>	0.44 $\pm$ 0.25	<b>0.70 <math>\pm</math> 0.26</b>	<b>0.54 <math>\pm</math> 0.28</b> (***)	0.46 $\pm$ 0.27 (**)	35.13 $\pm$ 15.6 (.)
DWI	0.42 $\pm$ 0.24 (*)	0.70 $\pm$ 0.26	0.51 $\pm$ 0.28	0.44 $\pm$ 0.27 (***)	31.28 $\pm$ 16.1 (**)
ADC	0.40 $\pm$ 0.24 (***)	0.67 $\pm$ 0.28 (.)	0.47 $\pm$ 0.27 (.)	0.43 $\pm$ 0.27 (***)	34.35 $\pm$ 20.4 (*)
T <sub>max</sub>	0.32 $\pm$ 0.20 (***)	0.63 $\pm$ 0.30 (*)	0.44 $\pm$ 0.25 (*)	0.35 $\pm$ 0.25 (***)	<b>29.99 <math>\pm</math> 13.7</b> (**)

<i>(b) Non-reperfused model: ablation study among non-reperfused patients</i>					
Input MRI	DSC	VS	Precision	Recall	HD
DWI+ADC+T <sub>max</sub> +CBF+CBV	<b>0.47 <math>\pm</math> 0.17</b>	<b>0.74 <math>\pm</math> 0.13</b>	0.49 $\pm$ 0.22	<b>0.52 <math>\pm</math> 0.21</b>	37.70 $\pm$ 17.7
DWI+ADC+T <sub>max</sub>	0.47 $\pm$ 0.18	0.74 $\pm$ 0.16	<b>0.52 <math>\pm</math> 0.22</b>	0.50 $\pm$ 0.22	35.77 $\pm$ 20.2
DWI	0.45 $\pm$ 0.17	0.71 $\pm$ 0.17	0.50 $\pm$ 0.22	0.50 $\pm$ 0.25	33.20 $\pm$ 17.2
ADC	0.42 $\pm$ 0.15 (.)	0.73 $\pm$ 0.23	0.47 $\pm$ 0.18	0.46 $\pm$ 0.21 (.)	28.35 $\pm$ 12.9 (**)
T <sub>max</sub>	0.40 $\pm$ 0.19 (*)	0.65 $\pm$ 0.21	0.50 $\pm$ 0.29	0.46 $\pm$ 0.24 (***)	<b>26.86 <math>\pm</math> 13.3</b> (.)

240 this case-mix likely accounts for the lack of significant difference between the general and  
 241 reperfused models.

242 Still, the pathophysiological rationale for integrating the reperfusion status in predictive  
 243 models is strong. Timely reperfusion is closely associated with increased penumbra salvage  
 244 and reduced final infarct size (Cho et al., 2015). **We propose that a new patient’s eligibility**  
 245 **to treatment could be assessed by using both CNNs (the one trained from reperfused**  
 246 **and the other from non-reperfused patients). The clinician would thus have a dual set of**  
 247 **predictive maps allowing a comparison of the projected infarct with and without reperfusion,**  
 248 **and an estimation of the treatment effect.** A mismatch between these two models (i.e. a  
 249 smaller infarct in case of a successful thrombectomy that achieved reperfusion, than in the  
 250 no-reperfusion model) would indicate that this patient is likely to benefit from therapy  
 251 (responder). Conversely, a similar output from the reperfused and non-reperfused models  
 252 would suggest a limited effect of therapy (non-responder). In our selected dataset, the final  
 253 predicted infarct was substantially larger with the non-reperfused CNN in 53 ( $\sim$ 50%) patients  
 254 when considering the following criteria: DSC between the two CNNs  $< 0.5$  and non-reperfused  
 255 CNN lesion volume  $\geq 20\%$  larger than the output of the reperfused CNN. Conversely, the  
 256 absence of a clear difference between the two models would suggest limited benefit from  
 257 reperfusion therapies. Reliable predictions of the final infarct may also help in evaluating



258 novel neuroprotection strategies, by comparing the projected *vs* observed infarct size in  
259 patients with ischemia-reperfusion (Hougaard et al., 2013). This approach may facilitate  
260 the screening of a larger number of putative neuroprotectants at lesser cost than full-sized  
261 controlled trials.

262 **Our results indicate that CNN can successfully take into account reperfusion by condi-**  
263 **tioning the training dataset according to this clinical status, in order to achieve more robust**  
264 **predictions.** The full validation of this approach will require a multicentric collaboration in  
265 order to collect high quality longitudinal data, including cases without reperfusion.

#### 266 4.2. Comparison to current clinical prediction methods

267 Our CNN models achieved higher AUC and DSC than the perfusion-diffusion mismatch  
268 models currently used in clinical practice (patient A and B in Figure 4 are illustrative cases).  
269 Our results were in the same range as those of recently reported CNNs: the best model of  
270 the ISLES challenge achieved a DSC of 0.38 (Winzeck et al., 2018); Nielsen et al. (2018)  
271 reported a mean AUC of 0.88, while Yu et al. (2020) reported a mean DSC and AUC of 0.53  
272 and 0.89, respectively. However, a strict comparison is not possible as the cited studies were  
273 all performed on different datasets, and in the light of different time-windows of prediction.

274 We also confirmed that predicting the final infarct remains a challenging task. Mean  
275 DSC were modest (0.44 and 0.47 for the reperfused and non-reperfused model, respectively),  
276 corresponding to an assortment of highly accurate predictions (DSC>0.7) and failure of both  
277 CNNs and perfusion-diffusion mismatch models in other cases (e.g. patient C, D and E in  
278 Figure 4). Partial and sometimes extensive reversal of the diffusion lesion can be observed  
279 (patients C and D in Figure 4), especially in the event of early reperfusion (Yoo et al., 2019).  
280 This phenomenon may particularly affect patients with small baseline DWI lesion, in whom  
281 even limited discrepancies between the predicted and observed infarct may result in very  
282 low DSC values. Still, no significant correlation was found between the DSC and baseline  
283 DWI lesion volume ( $r=0.038$ ,  $p=0.72$ ). Also, baseline imaging cannot account for subsequent  
284 events that may alter the progression of ischemic lesions (e.g. patient E in Figure 4: a  
285 possible case of reocclusion after a successful reperfusion). These patients illustrate the  
286 heterogeneity and complexity of stroke lesion progression. Reinforcement learning could help  
287 improve the performance of CNNs by training more specifically on these underrepresented  
288 patients (Arulkumaran et al., 2017).

### 289 4.3. Predictive value of the MRI inputs

290 The ablation study showed that CBF and CBV had limited impact on the performance of  
291 our CNN. This result is in line with the common qualitative observation that the perfusion  
292 lesion is less conspicuous on CBF or CBV maps compared to  $T_{max}$  maps. A previous voxel  
293 and threshold-based study had also observed that these parameters were poor predictors of  
294 the final infarct (Christensen et al., 2009).

295 Thus, ADC, DWI and  $T_{max}$  could constitute the main inputs for the network predicting  
296 the final infarct. Similarly, Livne et al. have shown that both perfusion parameters and  
297 DWI made significant predictive contributions, albeit with a different method (extreme  
298 gradient tree boosting) and among patients who were not treated by thrombectomy and thus  
299 had a significantly lower rate of reperfusion (Livne et al., 2018). Our study was conducted  
300 among thrombectomy-treated patients with a reperfusion rate of 80%, in whom the baseline  
301 DWI lesion is known to have a strong correlation with the final infarct. Our results further  
302 suggest that  $T_{max}$  maps may have a greater predictive value among non-reperused patients,  
303 which would be consistent with previously available data. Wheeler et al. (2013) had shown a  
304 strong correlation between the baseline diffusion lesion and final infarct volume in reperused  
305 patients, and a high correlation between the  $T_{max} > 6$  seconds lesion and final infarct volume  
306 for non-reperused patients.

307 These observations support our chosen deep learning architecture. The late fusion  
308 configuration allows for better integration of the distinct information contained in perfusion  
309 and diffusion imaging. Training reperfusion status-specific models entail assigning distinct  
310 weights to each MRI input. The performance of CNNs built with an early fusion configuration  
311 are presented in Appendix C. Early fusion had overall worse performance than late fusion.  
312 Fewer performance differences were also observed between the general, reperused and  
313 non-reperused models, suggesting that early fusion may overlook the reperfusion status.

### 314 4.4. Limitations

315 Our study presents several limitations. Patients were included from two cohorts with dif-  
316 ferent treatment protocols: HIBISCUS-STROKE involved patients treated by thrombectomy,  
317 whereas I-KNOW was a multicentric observational study of patients managed conservatively  
318 or with intravenous thrombolysis without any endovascular procedure. However, I-KNOW  
319 only contributed patients with proximal occlusions without reperfusion, who likely have a  
320 very similar course to failed thrombectomy cases. Methods for assessing early reperfusion  
321 differed between these two cohorts. Nevertheless, as proposed in a previous study, MRI and  
322 angiographic data can be pooled when evaluating reperfusion (Marks et al., 2014). Several

323 precautions were observed to limit potential biases: (i) TICI score assessment strictly followed  
324 standard recommendations (Zaidat et al., 2013) and was thus not a surrogate for recanal-  
325 ization; (ii) both TICI score and DSC-PWI assess tissue perfusion; similar criteria for both  
326 methods were used to identify reperfusion (TICI  $\geq 2$  and DSC-PWI reperfusion ratio  $\geq 50\%$ );  
327 (iii) in I-KNOW, the follow-up DSC-PWI used to assess reperfusion was performed with a  
328 median delay of 170 min from the baseline MRI, and was thus in a similar ultra-early time  
329 frame as HIBISCUS patients undergoing endovascular treatment. Furthermore, no significant  
330 difference was found between the non-reperfused patients of the two cohorts for the following  
331 baseline variables: gender, age, baseline NIHSS score, time from symptoms onset to MRI,  
332 baseline DWI lesion size. The HIBISCUS cohort had a majority of M1 occlusions (15/18; 3  
333 patients had a M2 occlusion), while most I-KNOW patients had M2 occlusions (12/17; 5 had  
334 a M1 occlusion. This significant difference in occlusion level ( $p=0.002$ , Fisher’s exact test) is  
335 likely related to the distinct inclusion criteria of these two cohorts (HIBISCUS specifically  
336 included patients with proximal intracranial occlusions). Other clinical parameters such as  
337 age and time from symptoms onset to imaging and reperfusion are recognized prognostic  
338 factors. Their integration in predictive CNNs may enhance model performance and warrants  
339 further investigation. Finally, the interval between stroke onset and the follow-up MRI was 6  
340 days. Other studies used different or similar delays: 3 to 7 days (Yu et al., 2020), 1-month  
341 (Nielsen et al., 2018) or 90 days (Winzeck et al., 2018). A previous study has shown that the  
342 24-hour DWI lesion volume was well correlated with day 90 FLAIR lesion volume (Campbell  
343 et al., 2012). Infarct volume at either time points predicted functional outcome. Studies  
344 using different intervals may be compared provided a successful coregistration of baseline  
345 and final images was achieved.

## 346 **5. Conclusion**

347 The performance of deep learning models improved when the reperfusion status was  
348 incorporated in their training. CNN-based models outperformed the clinically-used perfusion-  
349 diffusion mismatch model. Comparing the predicted infarct in case of a successful *vs* failed  
350 reperfusion may help in estimating the treatment effect and guiding therapeutic decisions in  
351 selected patients.

## 352 **Acknowledgments and funding sources**

353 This work was supported by the RHU MARVELOUS (ANR-16-RHUS-0009) of Université  
354 Claude Bernard Lyon-1 (UCBL), and was also performed within the framework of the RHU

355 BOOSTER (ANR-18-RHUS-0001), within the program "Investissements d'Avenir" operated  
356 by the French National Research Agency (ANR).

### 357 **Author Contributions**

358 **Noëlie Debs**: investigation, methodology, writing - original draft. **Tae-Hee Cho**:  
359 conceptualization, data acquisition and annotation, validation, critical revision of the  
360 manuscript. **David Rousseau**: conceptualization, validation, critical revision of the  
361 manuscript. **Yves Berthezène**: data acquisition; critical revision of the manuscript.  
362 **Marielle Buisson**: project administration. **Omer Eker**: data acquisition, critical re-  
363 vision of the manuscript. **Laura Mechtouff**: data acquisition, critical revision of the  
364 manuscript. **Norbert Nighoghossian**: project administration, data acquisition, critical  
365 revision of the manuscript. **Michel Ovize**: project administration, critical revision of  
366 the manuscript. **Carole Frindel**: conceptualization, validation, critical revision of the  
367 manuscript, supervision.

### 368 **Appendix A. Data**

#### 369 *Appendix A.1. Inclusion criteria of HIBISCUS-STROKE and I-KNOW*

370 Inclusion criteria for HIBISCUS-STROKE were: (1) patients with an anterior circulation  
371 stroke related to a proximal intracranial occlusion (internal carotid artery, M1 or M2  
372 occlusion), directly admitted to our comprehensive stroke unit ('mothership' paradigm); (2)  
373 diffusion and perfusion MRI as baseline imaging; (3) patients treated by thrombectomy with  
374 or without intravenous thrombolysis.

375 Inclusion and exclusion criteria for I-KNOW were: (1) NIHSS  $\geq 4$ ; (2) diffusion and  
376 perfusion MRI consistent with an acute anterior circulation ischemic stroke; and (3) admission  
377 MRI completed within 6 hours for patients treated with intravenous thrombolysis, or within 12  
378 hours for those managed without thrombolysis. Patients with lacunar or posterior circulation  
379 stroke, unknown time of onset or intracerebral hemorrhage were excluded. No patient received  
380 intra-arterial therapy. For the present study, additional inclusion criteria were applied, as  
381 follows: (1) both admission and acute follow-up diffusion and perfusion MRI obtained 3  
382 hours after initial imaging (H3) available and assessable; (2) visible occlusion on the baseline  
383 MRA; and (3) H3 perfusion without significant reperfusion.

Table A.5: Baseline characteristics (median with interquartile range, unless otherwise indicated). NIHSS: National Institutes of Health Stroke Scale; DWI: diffusion-weighted imaging; ICA: internal carotid artery.

<b>Clinical variables</b>	
Women, n (percentage)	45 (41.3)
Age	70 (57 - 79)
NIHSS score	15 (10 - 19)
Time from symptoms onset to MRI	105 (78 - 154)
Intravenous tPA, n (percentage)	59 (54.1)
Site of occlusion, n (percentage):	
intracranial ICA+M1	27 (24.8)
M1	54 (49.5)
intracranial ICA+M2	23 (21.1)
M2	5 (4.6)
cervical ICA, n (percentage)	19 (17.4)
DWI lesion size, mL	24.9 (7.4 - 50.9)

384 *Appendix A.2. Patients' baseline characteristics*

385 *Appendix A.3. MRI protocol*

386 All patients underwent DWI (IKNOW : repetition time 6000 ms, field of view 24 cm,  
387 matrix 128×128 (IKNOW) or 192×192 (HIBISCUS-STROKE), slice thickness 5mm), Fluid-  
388 attenuated-inversion-recovery (repetition time 8690 ms, echo time 109 ms, inversion time  
389 2500 ms, field of view 21 cm, matrix 224×256, section thickness 5 mm), T2-weighted gradient  
390 echo (repetition time 800 ms, echo time 28 ms, flip angle 20°, field of view 230 mm, matrix  
391 512×512, section thickness of 5 mm), MRA and DSC-PWI (echo time 40 ms, repetition time  
392 1500 ms, field of view 24 cm, matrix 128×128, slice thickness 5 mm; gadolinium contrast at  
393 0.1 mmol/kg), both for the admission and follow-up MRI.

## 394 **Appendix B. Network training and parameters**

395 Only slices including the final infarct were used to train the U-net and no data augmen-  
396 tation was employed. We used a multi-class Dice function as a loss function (Milletari et al.,  
397 2016), for which the lesion class was assigned a weight 8 times higher than those of healthy  
398 and background classes. We used the Adam optimizer ( $lr = 1 \times 10^{-4}$  and  $decay = 5 \times 10^{-4}$ )  
399 and a batch size of 12. To prevent overfitting, we applied dropout (set to 0.5), used a L2  
400 regularizer  $reg$  at each convolution layer ( $reg = 2 \times 10^{-4}$ ) and the the number of epochs  
401 (set to 500) was regulated by early stopping (*i.e.* the training was stopped once the best  
402 validation multi-class dice did not increase more than 0.005 on 100 epochs). The evaluation of

403 each model was performed using a 5-fold cross-validation. Note that patients from I-KNOW  
404 dataset were added in the training set of the general and the non-reperused models for  
405 data-augmentation purposes, but were not used in the testing set. Specifically, the number  
406 of training patients was, depending on the fold: between 89 and 91 patients for the general  
407 model, between 59 and 60 patients for the reperused model, and between 30 and 31 pa-  
408 tients for the non-reperused model. The number of test patients varied between 17 and 19  
409 (reperused and non-reperused patients combined).

410 The number of parameters is proportional to the number of U-Net path: thus, the number  
411 of trainable parameters is 1997851 for a U-Net architecture with 5 MRI sequence inputs,  
412 1242603 for 3 MRI inputs, and 487355 when using only one input. The higher the number  
413 of paths, the less the information is compressed and the more the architecture offers the  
414 possibility of learning different information on each input data. Thus, we chose not to balance  
415 the number of parameters between each architecture. However, to ensure a fair comparison,  
416 each network’s hyperparameters were independently fine-tuned on a fixed search space. The  
417 best parameters were found to be the same in all tested architectures. We used Keras 2.1.3  
418 library with Python 3.6.3 interface. The training phase took approximately 1 hour on a work  
419 station with an NVIDIA GeForce GTX 1080 GPU with 128 GB memory.

## 420 **Appendix C. Impact of the multiple MRI fusion configuration**

421 We compared our proposed late fusion deep learning architecture to an early fusion  
422 one, where all patient input images are combined at the beginning of the CNN. This fusion  
423 strategy reduces both the computational complexity and training parameters (Chen et al.,  
424 2019). Each patient being represented by DWI, ADC,  $T_{max}$ , CBV, CBF, the early fusion  
425 architecture stacks channel-wise these 5 MRI inputs and does not process them independently.  
426 Results are shown in Table C.6.

427 It appears that best metric values are obtained when performing a late fusion strategy  
428 rather than an early fusion: average values of DSC, VS, precision and recall are higher  
429 whatever the training set (all, reperused, non reperused). However, lowest values for HD  
430 metric are obtained when performing early fusion. Early fusion seems to offer a better spatial  
431 delineation of the final lesion: fewer outliers seem to be predicted, which drastically decreases  
432 HD values.

433 With early fusion configuration, differences observed between the global model and the  
434 reperused and non-reperused submodels are smaller and not significant. This type of  
435 architecture seems less adapted to take into account the status of reperfusion.

Table C.6: Evaluation metrics after training models on different training set (all, reperfused, and non-reperfused) with different fusion strategies (early and late) and evaluating them on reperfused testing patients (a) and non-reperfused testing patients (b) (average values  $\pm$  standard deviation). Bold values correspond to the best value of the respective evaluation metric (column-wise). A two-sided wilcoxon signed-rank test was performed between global model and the two other models (reperfused and non-reperfused) for a given fusion strategy, with (.) indicating  $P < 0.10$ , (\*) indicating  $P < 0.05$ , (\*\*) indicating  $P < 0.01$  and (\*\*\*) indicating  $P < 0.001$ .

(a) Evaluation on reperfused testing patients

Fusion	Training	DSC	VS	Precision	Recall	HD
early	all	0.39 $\pm$ 0.25	0.59 $\pm$ 0.30	0.56 $\pm$ 0.31	0.40 $\pm$ 0.26	29.51 $\pm$ 16.26
early	reperfused	0.41 $\pm$ 0.25	0.64 $\pm$ 0.30	0.46 $\pm$ 0.29 (***)	0.49 $\pm$ 0.30 (***)	31.24 $\pm$ 15.61
early	non-reperfused	0.36 $\pm$ 0.22 (*)	0.63 $\pm$ 0.27	0.54 $\pm$ 0.26	0.33 $\pm$ 0.24 (***)	<b>26.64 <math>\pm</math> 11.16</b>
late	all	0.43 $\pm$ 0.24	0.69 $\pm$ 0.27	0.55 $\pm$ 0.28	0.43 $\pm$ 0.25	<b>33.23 <math>\pm</math> 15.64</b>
late	reperfused	<b>0.44 <math>\pm</math> 0.25</b>	<b>0.70 <math>\pm</math> 0.27</b>	0.50 $\pm$ 0.27	<b>0.50 <math>\pm</math> 0.26 (***)</b>	38.58 $\pm$ 18.15
late	non-reperfused	0.35 $\pm$ 0.21 (***)	0.57 $\pm$ 0.28 (***)	<b>0.60 <math>\pm</math> 0.25 (***)</b>	0.31 $\pm$ 0.24 (***)	40.05 $\pm$ 15.66 (**)

(b) Evaluation on non-reperfused testing patients

Fusion	Training	DSC	VS	Precision	Recall	HD
early	all	0.42 $\pm$ 0.24	0.62 $\pm$ 0.27	0.42 $\pm$ 0.28	0.55 $\pm$ 0.29	30.98 $\pm$ 18.23
early	reperfused	0.41 $\pm$ 0.26	0.51 $\pm$ 0.31	0.36 $\pm$ 0.29	0.69 $\pm$ 0.24	30.94 $\pm$ 16.30
early	non-reperfused	0.42 $\pm$ 0.18	0.66 $\pm$ 0.17	0.42 $\pm$ 0.24	0.55 $\pm$ 0.22	<b>28.48 <math>\pm</math> 13.63</b>
late	all	0.44 $\pm$ 0.21	0.66 $\pm$ 0.26	0.39 $\pm$ 0.25	0.63 $\pm$ 0.21	<b>30.61 <math>\pm</math> 16.15</b>
late	reperfused	0.44 $\pm$ 0.22	0.63 $\pm$ 0.25	0.36 $\pm$ 0.23	<b>0.69 <math>\pm</math> 0.22 (*)</b>	44.53 $\pm$ 16.79 (**)
late	non-reperfused	<b>0.47 <math>\pm</math> 0.17</b>	<b>0.74 <math>\pm</math> 0.13</b>	<b>0.49 <math>\pm</math> 0.22 (**)</b>	0.52 $\pm$ 0.21 (***)	37.70 $\pm$ 17.74

## References

- Albers, G., et al., 2018. Thrombectomy for stroke at 6 to 16 hours with selection by perfusion imaging. *N Engl J Med* 378, 708–718.
- Arulkumaran, K., Deisenroth, M.P., Brundage, M., Bharath, A.A., 2017. Deep reinforcement learning: a brief survey. *IEEE Signal Processing Magazine* 34, 26–38.
- Avants, B.B., Tustison, N.J., Song, G., Cook, P.A., Klein, A., Gee, J.C., 2011. A reproducible evaluation of ANTs similarity metric performance in brain image registration. *NeuroImage* 54, 2033–2044.
- Aygün, M., Şahin, Y.H., Ünal, G., 2018. Multi modal convolutional neural networks for brain tumor segmentation. *arXiv preprint arXiv:1809.06191* .
- Barber, P., Darby, D., Desmond, P., Yang, Q., Gerraty, R., Jolley, D., Donnan, G., Tress, B., Davis, S.M., 1998. Prediction of stroke outcome with echoplanar perfusion- and diffusion-weighted MRI. *Neurology* 51, 418–426.
- Campbell, B.C.V., Tu, H.T.H., Christensen, S., Desmond, P.M., Levi, C.R., Bladin, C.F., Hjort, N., Ashkanian, M., Sølling, C., Donnan, G.A., Davis, S.M., Ostergaard, L., Parsons, M.W., 2012. Assessing response to stroke thrombolysis: validation of 24-hour multimodal magnetic resonance imaging. *Arch Neurol* 69, 46–50.
- Chen, Y., Wang, K., Liao, X., Qian, Y., Wang, Q., Yuan, Z., Heng, P.A., 2019. Channel-UNet: a spatial

453 channel-wise convolutional neural network for liver and tumors segmentation. *Frontiers in Genetics* 10.

454 Cho, T.H., Nighoghossian, N., Mikkelsen, I.K., Derex, L., Hermier, M., Pedraza, S., Fiehler, J., Østergaard,  
455 L., Berthezène, Y., Baron, J.C., 2015. Reperfusion within 6 hours outperforms recanalization in predicting  
456 penumbra salvage, lesion growth, final infarct, and clinical outcome. *Stroke* 46, 1582–1589.

457 Christensen, S., Mouridsen, K., Wu, O., Hjort, N., Karstoft, H., Thomalla, G., Röther, J., Fiehler, J.,  
458 Kucinski, T., Østergaard, L., 2009. Comparison of 10 perfusion MRI parameters in 97 sub-6-hour stroke  
459 patients using voxel-based receiver operating characteristics analysis. *Stroke* 40, 2055–2061.

460 Dolz, J., Ayed, I.B., Desrosiers, C., 2018a. Dense multi-path U-Net for ischemic stroke lesion segmentation  
461 in multiple image modalities, in: *International MICCAI Brainlesion Workshop*, Springer. pp. 271–282.

462 Dolz, J., Desrosiers, C., Ayed, I.B., 2018b. IVD-Net: Intervertebral disc localization and segmentation in  
463 MRI with a multi-modal UNet, in: *International Workshop and Challenge on Computational Methods  
464 and Clinical Applications for Spine Imaging*, Springer. pp. 130–143.

465 Goyal, M., Menon, B.K., van Zwam, W.H., Dippel, D.W., Mitchell, P.J., Demchuk, A.M., Dávalos, A., Majoie,  
466 C.B., van der Lugt, A., De Miquel, M.A., et al., 2016. Endovascular thrombectomy after large-vessel  
467 ischaemic stroke: a meta-analysis of individual patient data from five randomised trials. *The Lancet* 387,  
468 1723–1731.

469 Hajian-Tilaki, K., 2013. Receiver operating characteristic (ROC) curve analysis for medical diagnostic test  
470 evaluation. *Caspian journal of internal medicine* 4, 627.

471 Hougaard, K.D., Hjort, N., Zeidler, D., Sørensen, L., Nørgaard, A., Hansen, T.M., von Weitzel-Mudersbach,  
472 P., Simonsen, C.Z., Damgaard, D., Gottrup, H., Svendsen, K., Rasmussen, P.V., Ribe, L.R., Mikkelsen,  
473 I.K., Nagenthiraja, K., Cho, T.H., Redington, A.N., Bøtker, H.E., Østergaard, L., Mouridsen, K., Andersen,  
474 G., 2013. Remote ischemic preconditioning as an adjunct therapy to thrombolysis in patients with acute  
475 ischemic stroke: a randomized trial. *Stroke* 45, 159–167.

476 Jonsdottir, K.Y., Østergaard, L., Mouridsen, K., 2009. Predicting tissue outcome from acute stroke magnetic  
477 resonance imaging: improving model performance by optimal sampling of training data. *Stroke* 40,  
478 3006–3011.

479 Kidwell, C.S., Wintermark, M., De Silva, D.A., Schaewe, T.J., Jahan, R., Starkman, S., Jovin, T., Hom, J.,  
480 Jumaa, M., Schreier, J., et al., 2013. Multiparametric MRI and CT models of infarct core and favorable  
481 penumbral imaging patterns in acute ischemic stroke. *Stroke* 44, 73–79.

482 Livne, M., Boldsen, J.K., Mikkelsen, I.K., Fiebach, J.B., Sobesky, J., Mouridsen, K., 2018. Boosted tree  
483 model reforms multimodal magnetic resonance imaging infarct prediction in acute stroke. *Stroke* 49,  
484 912–918.

485 Marks, M.P., Lansberg, M.G., Mlynash, M., Kemp, S., McTaggart, R., Zaharchuk, G., Bammer, R., Albers,  
486 G.W., 2014. Correlation of AOL recanalization, TIMI reperfusion and TICI reperfusion with infarct  
487 growth and clinical outcome. *Journal of neurointerventional surgery* 6, 724–728.

488 McKinley, R., Häni, L., Gralla, J., El-Koussy, M., Bauer, S., Arnold, M., Fischer, U., Jung, S., Mattmann, K.,  
489 Reyes, M., et al., 2017. Fully automated stroke tissue estimation using random forest classifiers (FASTER).  
490 *Journal of Cerebral Blood Flow & Metabolism* 37, 2728–2741.

491 Milletari, F., Navab, N., Ahmadi, S.A., 2016. V-Net: Fully convolutional neural networks for volumetric  
492 medical image segmentation, in: *2016 Fourth International Conference on 3D Vision (3DV)*, IEEE. pp.  
493 565–571.



494 Nie, D., Wang, L., Gao, Y., Shen, D., 2016. Fully convolutional networks for multi-modality iso-intense infant  
495 brain image segmentation, in: 2016 IEEE 13th international symposium on biomedical imaging (ISBI),  
496 IEEE. pp. 1342–1345.

497 Nielsen, A., Hansen, M.B., Tietze, A., Mouridsen, K., 2018. Prediction of tissue outcome and assessment of  
498 treatment effect in acute ischemic stroke using deep learning. *Stroke* 49, 1394–1401.

499 Nogueira, R.G., et al., 2018. Thrombectomy 6 to 24 hours after stroke with a mismatch between deficit and  
500 infarct. *N Engl J Med* 378, 11–21.

501 Olivot, J.M., Mlynash, M., Thijs, V.N., Kemp, S., Lansberg, M.G., Wechsler, L., Bammer, R., Marks, M.P.,  
502 Albers, G.W., 2009. Optimal Tmax threshold for predicting penumbral tissue in acute stroke. *Stroke* 40,  
503 469–475.

504 Pinto, A., McKinley, R., Alves, V., Wiest, R., Silva, C.A., Reyes, M., 2018. Stroke lesion outcome prediction  
505 based on MRI imaging combined with clinical information. *Frontiers in neurology* 9, 1060.

506 Powers, W.J., et al., 2019. Guidelines for the early management of patients with acute ischemic stroke: 2019  
507 update to the 2018 guidelines for the early management of acute ischemic stroke: A guideline for healthcare  
508 professionals from the American Heart Association/American Stroke Association. *Stroke* 50, e344–e418.

509 Qiu, W., Kuang, H., Teleg, E., Ospel, J.M., Sohn, S.I., Almekhlafi, M., Goyal, M., Hill, M.D., Demchuk,  
510 A.M., Menon, B.K., 2020. Machine learning for detecting early infarction in acute stroke with non-contrast-  
511 enhanced CT. *Radiology*, 191193.

512 Reikik, I., Allassonnière, S., Carpenter, T.K., Wardlaw, J.M., 2012. Medical image analysis methods in  
513 MR/CT-imaged acute-subacute ischemic stroke lesion: segmentation, prediction and insights into dynamic  
514 evolution simulation models. A critical appraisal. *NeuroImage: Clinical* 1, 164–178.

515 Robben, D., Boers, A.M., Marquering, H.A., Langezaal, L.L., Roos, Y.B., van Oostenbrugge, R.J., van Zwam,  
516 W.H., Dippel, D.W., Majoie, C.B., van der Lugt, A., et al., 2020. Prediction of final infarct volume from  
517 native CT perfusion and treatment parameters using deep learning. *Medical image analysis* 59, 101589.

518 Smith, S., Bannister, P.R., Beckmann, C., Brady, M., Clare, S., Flitney, D., Hansen, P., Jenkinson, M.,  
519 Leibovici, D., Ripley, B., et al., 2001. FSL: New tools for functional and structural brain image analysis.  
520 *NeuroImage* 13, 249.

521 Taha, A.A., Hanbury, A., 2015. Metrics for evaluating 3D medical image segmentation: analysis, selection,  
522 and tool. *BMC medical imaging* 15, 29.

523 Tsai, J.P., Albers, G.W., 2015. Reperfusion versus recanalization: the winner is...

524 Wheeler, H.M., Mlynash, M., Inoue, M., Tipirneni, A., Liggins, J., Zaharchuk, G., Straka, M., Kemp, S.,  
525 Bammer, R., Lansberg, M.G., et al., 2013. Early diffusion-weighted imaging and perfusion-weighted  
526 imaging lesion volumes forecast final infarct size in DEFUSE 2. *Stroke* 44, 681–685.

527 Winder, A.J., Siemonsen, S., Flottmann, F., Thomalla, G., Fiehler, J., Forkert, N.D., 2019. Technical  
528 considerations of multi-parametric tissue outcome prediction methods in acute ischemic stroke patients.  
529 *Scientific reports* 9, 1–12.

530 Winzeck, S., Hakim, A., McKinley, R., Pinto, J.A., Alves, V., Silva, C., Pisov, M., Krivov, E., Belyaev, M.,  
531 Monteiro, M., et al., 2018. ISLES 2016 and 2017-benchmarking ischemic stroke lesion outcome prediction  
532 based on multispectral MRI. *Frontiers in neurology* 9, 679.

533 Yoo, J., Choi, J.W., Lee, S.J., Hong, J.M., Hong, J.H., Kim, C.H., Kim, Y.W., Kang, D.H., Kim, Y.S.,  
534 Hwang, Y.H., Ovbiagele, B., Demchuk, A.M., Lee, J.S., Sohn, S.I., 2019. Ischemic diffusion lesion reversal

535 after endovascular treatment. *Stroke* 50, 1504–1509.

536 Yu, Y., Xie, Y., Thamm, T., Gong, E., Ouyang, J., Huang, C., Christensen, S., Marks, M.P., Lansberg,  
537 M.G., Albers, G.W., et al., 2020. Use of deep learning to predict final ischemic stroke lesions from initial  
538 magnetic resonance imaging. *JAMA Network Open* 3, e200772–e200772.

539 Zaidat, O.O., Yoo, A.J., Khatri, P., Tomsick, T.A., Von Kummer, R., Saver, J.L., Marks, M.P., Prabhakaran,  
540 S., Kallmes, D.F., Fitzsimmons, B.F.M., et al., 2013. Recommendations on angiographic revascularization  
grading standards for acute ischemic stroke: a consensus statement. *Stroke* 44, 2650–2663.