

ninon.burgos@cnrs.fr





ENSURING ROBUST AI IN MEDICAL IMAGING:

VALIDATION ON REAL-WORLD DATA, MEANINGFUL METRICS, AND RIGOROUS STATISTICAL ANALYSIS

Ninon Burgos CNRS researcher & PR[AI]RIE chair ARAMIS Lab, Paris Brain Institute, France

Context





Computer-aided diagnosis

- Machine learning / deep learning task: classification
- Disorder targeted: dementia
- Imaging modality: T1-weighted MRI





A (fake & slightly exaggerated) typical MICCAI conclusion

'Our proposed deep learning framework achieves superior accuracy (86.6%) and robust generalizability in distinguishing Alzheimer's disease patients from cognitively normal subjects, outperforming established methods on the ADNI dataset. The model provides anatomically interpretable results by highlighting disease-relevant regions such as the hippocampus, demonstrates high reproducibility and fast inference, and, by eliminating manual feature engineering while aligning with clinical biomarkers, is both practical and explainable-making it ready for immediate integration into routine diagnostic workflows.'



Ensuring Robust AI in Medical Imaging

- Validation on Real-World Data
- Meaningful Metrics
- Rigorous Statistical Analysis



Ensuring Robust AI in Medical Imaging

- Validation on Real-World Data
- Meaningful Metrics
- Rigorous Statistical Analysis



ML/DL for dementia diagnosis & prognosis

• A very active field of research





Motivation

Most machine learning algorithms for computeraided diagnosis of dementia developed and applied on research data



Validation on Real-World Data









Motivation

Most machine learning algorithms for computeraided diagnosis of dementia developed and applied on research data

Objective

- Validate these algorithms on clinical data
- Develop robust algorithms able to handle clinical data





AP-HP: a network of 39 hospitals





AP-HP clinical data warehouse





Objective

 \succ Identify patients with dementia among a clinical data warehouse from anatomical MRI using machine learning

Cohort definition

- **Inclusion criteria**
 - Age \geq 18 years
 - At least one cerebral MRI scan including T1-weighted 3D acquisition Ο
 - Clinical data available (available for ~18% of the subjects with a 3D T1 MRI) Ο
- Definition of the classes of interest using diagnosis (ICD-10) codes

No dementia and no lesions (NDNL) No dementia with lesions (NDL)

Dementia vs



Cohort

Category	N patients	N images	Age (mean ±std [range])	Sex (%F)
Dementia	756	887	71.17 ± 11.58 [18,90]	50.34%
NDNL	756	939	71.17 ± 11.58 [18,90]	50.34%
NDL	756	997	71.17 ± 11.58 [18,90]	50.34%
Total	2268	2823	71.17 ± 11.58 [18,90]	50.34%



Classification algorithms

- Machine learning: support vector machine
 - Input: grey matter density map

- Deep learning: convolutional neural networks
 - Input: minimally processed MRI





Evaluation setting

- Test set of 152 patients/images per class
- 5-fold cross-validation



Results

			Classifica	tion strategy		
	Data set	Task	SVM with grey matter maps	ResNet with minimally processed MRI		
/	Clinical	D vs NDNL	68.75	73.62		
		D vs NDL	73.09	72.24		
ADNI ALZHEIMER'S DISEASE EUROIMAGING INITIATIVE	Research	AD vs CN	86.80	85.30		
			Balanced accuracy (%)			

Decrease of 15 percent points in balanced accuracy: clinical data set more heterogenous



Cohort

Category	N patients	N images	Age (mean ±std [range])	Sex (%F)	% Good/medium quality	% With gadolinium
Dementia	756	887	71.17 ± 11.58 [18,90]	50.34%	57.72%**	24.80%**
NDNL	756	939	71.17 ± 11.58 [18,90]	50.34%	36.42%**	66.13%**
NDL	756	997	71.17 ± 11.58 [18,90]	50.34%	52.25%	63.59%**
Total	2268	2823	71.17 ± 11.58 [18,90]	50.34%	48.71%	52.24%

p-value corrected with Bonferroni < 0.05

- Objectives
 - Reject images that are not proper T1weighted anatomical MR images
 - Recognise images acquired after injection of a contrast agent (gadolinium)
 - Rate the overall image quality (good/medium/low)











• Examples









- Results
 - Reject images that are not proper T1weighted anatomical MR images

Balanced accuracy > 90%

Recognise images acquired after contrast injection

Balanced accuracy > 95%

- > Rate the overall image quality
 - Good/medium vs low: balanced accuracy > 80%
 - Good vs medium: balanced accuracy > 70%



Good Medium

0W





Cohort

Category	N patients	N images	Age (mean ±std [range])	Sex (%F)	% Good/medium quality	% With gadolinium
Dementia	756	887	71.17 ± 11.58 [18,90]	50.34%	57.72%**	24.80%**
NDNL	756	939	71.17 ± 11.58 [18,90]	50.34%	36.42%**	66.13%**
NDL	756	997	71.17 ± 11.58 [18,90]	50.34%	52.25%	63.59%**
Total	2268	2823	71.17 ± 11.58 [18,90]	50.34%	48.71%	52.24%

p-value corrected with Bonferroni < 0.05



• Results with training subsets

	Training set (n=88 per class)							
Task	<u>Quality</u> : good/medium + low <u>Gadolinium</u> : presence & absence	<u>Quality</u> : good/medium <u>Gadolinium</u> : presence & absence	<u>Quality</u> : good/medium <u>Gadolinium</u> : absence	<u>Quality</u> : good/medium <u>Gadolinium</u> : absence/synthetic				
D vs NDNL	69.47	60.26	51.51	51.71				
D vs NDL	73.03	68.29	50.00	54.08				

Balanced accuracy (%) - Classification strategy: SVM with grey matter maps



Contrast-enhanced to non-contrast-enhanced image translation

• Methods



3D U-Net like generators

- Residual connections
- Attention gates
- Transformer layers

Conditional GANs

- 3D U-Net like generators
- 3D patch discriminator

- Paired T1ce/T1nce MRI
- 230 pairs for training
- 77 pairs for testing



Contrast-enhanced to non-contrast-enhanced image translation



Bottani et al., BMC Medical Imaging, 2024



Contrast-enhanced to non-contrast-enhanced image translation

- Results
 - 1. Image synthesis accuracy
 - Mean absolute error
 - Peak signal to noise ratio
 - Structural similarity
 - 2. Segmentation fidelity
 - Absolute volume difference
 - Volume difference

Reliable feature extraction

Good synthesis accuracy



Bottani et al., BMC Medical Imaging, 2024



• Results with training subsets

	Training set (n=88 per class)							
Task	<u>Quality</u> : good/medium + low <u>Gadolinium</u> : presence & absence	<u>Quality</u> : good/medium <u>Gadolinium</u> : presence & absence	<u>Quality</u> : good/medium <u>Gadolinium</u> : absence	<u>Quality</u> : good/medium <u>Gadolinium</u> : synthetic				
D vs NDNL	69.47	60.26	51.51	51.71				
D vs NDL	73.03	68.29	50.00	54.08				

Balanced accuracy (%) - Classification strategy: SVM with grey matter maps

Classifier exploits biases of the training set

Image translation can be useful to limit bias due to the injection of gadolinium



• Final results with unbiased training sets

		Classification strategy			
Training set	Task	SVM with grey matter maps	ResNet with minimally processed MRI		
Research data set	D vs NDNL	64.08	61.84		
ADNI Alziemen's Digeage Neuroomaging Intranye	D vs NDL	69.47	61.78		
Clinical data set with	D vs NDNL	61.91	63.22		
good/medium quality	D vs NDL	64.61	67.50		
		Balanced	accuracy (%)		

Balanced accuracy still ~20 percent points lower than when testing on research data set



- Summary of the results
 - > Balanced accuracy much lower than in research context
 - Difference in proportions of images with gadolinium and of medium/good quality can bias the results
 - Solution: training the algorithms using only research data or only clinical images of good/medium quality without gadolinium

Conclusion

- Quality control & data set homogenisation essential steps
- Challenge remaining
 - Lack of robustness of the labels



Computer-aided diagnosis of neurodegenerative diseases



Lesion detection





3D FLAIR cohort constitution



Loizillon et al., Medical Image Analysis, 2025



3D FLAIR image quality distribution



Loizillon et al., Medical Image Analysis, 2025



3D FLAIR automatic quality control

- Annotated data available for training:
 - 358 3D FLAIR MRIs
 - 5000 3D T1w MRIs

Task: Low vs Medium/Good

			1
	FLAIR	T1w	\Rightarrow Domain gap
Annotators	86.54	91.56	
Training with FLAIR	69.90 ± 2.38	48.20 ± 0.81	
Training with T1w	50.06 ± 1.71	83.51 ± 0.93	
Training with T1w + FLAIR	58.05 ± 4.45	82.37 ± 0.85	
Palanc	adaccuracy (%)		-

Balanced accuracy (%)

Loizillon et al., DART@MICCAI, 2024



Source domain

5000 3D T1w MRIs manually annotated $D_s = (x_i^s, y_i^s)^{N_s = 5000}$

Target domain

5358 3D FLAIR MRIs, with **358** manually annotated $D_{T_U} = (x_i^{T_U})^{N_{TU}=5000} \quad D_{T_L} = (x_i^{T_L}, y_i^{T_L})^{N_{TL}=358}$



Sundaresan et al., MedIA, 2021

Loizillon et al., Medical Image Analysis, 2025





Loizillon et al., Medical Image Analysis, 2025



Source domain

5000 3D T1w MRIs manually annotated $D_{s} = (x_{i}^{s}, y_{i}^{s})^{N_{s}=5000}$

Target domain

5358 3D FLAIR MRIs, with **358** manually annotated $D_{T_U} = (x_i^{T_U})^{N_{TU}=5000} \quad D_{T_L} = (x_i^{T_L}, y_i^{T_L})^{N_{TL}=358}$





	Straight reject (yes vs no)	Low vs Medium/Good	Medium vs Good
Annotators	94.67	86.31	84.43
	[91.32,97.40]	[84.25,88.25]	[82.54,86.50]
Loizillon et al.	89.27	79.81	73.92
	[79.69,97.08]	[74.65,84.71]	[70.73,77.22]

Balanced accuracy (%)

Loizillon et al., Medical Image Analysis, 2025



Detection of age-related white matter hyperintensities in 3D FLAIR MRI



Loizillon et al., MIDL, 2024



Detection of age-related white matter hyperintensities



Loizillon et al., MIDL, 2024



Detection of age-related white matter hyperintensities in 3D FLAIR MRI





Conclusions

- Quality control & data set homogenisation crucial to enable translation of computer-aided diagnosis tools to clinical practice
- > Not all AI approaches are appropriate for use in clinical practice
- Size is not the only characteristic of data sets that matters



Centralised vs federated learning







Objective:

 $\min_{w} F(w), \quad \text{where } F(w) \coloneqq \sum_{k=1}^{m} p_k F_k(w)$ m: total number of nodes $p_k \ge 0$

 $\sum_k p_k = 1$

 F_k : local objective function for the k^{th} node

Node contribution in federated averaging framework (MacMahan et al., AISTATS, 2017):

 $p_k = \frac{n_k}{n}$ n_k : number of datapoints in the node n: total number of observations studied







Federated learning

Fed-BioMed

'Open, Transparent and Trusted Collaborative Learning for Real-World Healthcare Applications'

fedbiomed.org





Ensuring Robust AI in Medical Imaging

- Validation on Real-World Data
- Meaningful Metrics
- Rigorous Statistical Analysis



A (fake & slightly exaggerated) typical MICCAI conclusion

'Our proposed deep learning framework achieves superior accuracy (86.6%) and robust generalizability in distinguishing Alzheimer's disease patients from cognitively normal subjects, outperforming established methods on the ADNI dataset. The model provides anatomically interpretable results by highlighting disease-relevant regions such as the hippocampus, demonstrates high reproducibility and fast inference, and, by eliminating manual feature engineering while aligning with clinical biomarkers, is both practical and explainable-making it ready for immediate integration into routine diagnostic workflows.'



Image level binary classification

Confusion matrix:						Accuracy:	Balanced accuracy:
Predicted diagnosis				osis			
True positive (TP) False positive (FP)		True positive (TP)	False negative (FN)		ative	$ACC = \frac{TP + TN}{TP + TN}$	Bal ACC = $0.5 \times \left(\frac{\text{TP}}{1000} + \frac{\text{TN}}{1000} \right)$
		False positive (FP)	True r (True negative (TN)		TP + FN + FP + TN	$\frac{1}{1000} = \frac{1}{1000} \times \left(\frac{1}{1000} + \frac{1}{1000} + \frac$
	Scenario A		,	AD	CN	45 + 45	$\mathbf{P}_{\mathbf{A}} = \mathbf{A} \mathbf{C} \mathbf{C} + \mathbf{O} \mathbf{E}_{\mathbf{A}} \left(\begin{array}{c} 45 \\ 45 \end{array} \right)$
•	50 AD pation	tc	AD	45	5	$ACC = \frac{1}{45 + 5 + 5 + 45}$	Dat. ACC = $0.5 \times \left(\frac{1}{5 + 45} + \frac{1}{45 + 5}\right)$
•	50 CN subject	cts	CN	5	45	= 0.9	= 0.9
	Scenario B	8		AD	CN	- /	
•	50 AD natien	ts	AD	5	45	$ACC = \frac{5 + 490}{5 + 45 + 10 + 490}$	Bal. ACC = $0.5 \times \left(\frac{5}{5+45} + \frac{490}{490+10}\right)$
•	500 CN subje	ects	CN	10	490	= 0.9	= 0.54



Brain tumour segmentation

Popular voxel-based metrics fail to capture clinical interest Magnetic resonance imaging, same patient, different slices



Sensitivity = 0.94 (voxel-level)

Sensitivity = 0.50 (instance-level)

Missed lesion!



Medical example: brain-tumor segmentation A near-perfect voxel-level sensitivity hides information on missed lesions



A framework for trustworthy image analysis validation

(Quick access to Metrics Reloaded main paper \square)

https://metrics-reloaded.dkfz.de







Maier-Hein et al., Nature Methods, 2024

Meaningful Metrics



Choose your tool

Start your selection

Discover Metrics



(i) Problem Category Selection

Mapping a given research problem to the appropriate image processing task.



Metric Selection

Metric recommendation depending on the task type and specifications of the given problem.



Metric Library

A collection of all metrics in our database including their definitions, references and restrictions.

https://metrics-reloaded.dkfz.de



Metrics library



https://metrics-reloaded.dkfz.de

Meaningful Metrics



Metrics Reloaded

Metrics library



https://metrics-reloaded.dkfz.de

MATTHEWS CORRELATION COEFFICIENT (MCC) Synonyms: Phi Coefficient							
$MCC = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}} = \frac{1}{1000000000000000000000000000000000$							
A value of 0 refers to a prediction which is not better than random guessing. DESCRIPTION MCC measures the correlation between the actual and HDLC SemS ObD InS							
DEFINITION [Matthews, 1975]	PREVALENCE DEPENDENCY						
CARDINALITIES TP FP FN TN Counting metric	METRIC FAMILY Multi-threshold Distance- based metric						
IMPORTANT RELATIONS MCC can be rewritten as:							
MCC = $\sqrt{PPV \cdot Sensitivity \cdot Specificity \cdot NPV} - \sqrt{(1 - PPV) \cdot (1 - Sensitivity) \cdot (1 - Specificity) \cdot (1 - NPV)}$ MCC is equivalent to the geometric mean of Markedness and Informedness.							
$\begin{split} & \text{MULTI-CLASS DEFINITION} \\ & \text{For C classes, MCC can be defined as:} \\ & \text{MCC} = \frac{\sum_{i=1}^{C} \sum_{j=1}^{C} \sum_{k=1}^{C} n_{i,1} \cdot n_{j,k} \cdot n_{i,j} \cdot n_{k,i}}{\sqrt{\sum_{i=1}^{C} \left(\sum_{j=1}^{C} n_{i,j}\right) \left(\sum_{i' \mid i' \neq i} \sum_{j'=1}^{C} n_{i'j'}\right)} \sqrt{\sum_{i=1}^{C} \left(\sum_{j=1}^{C} n_{j,j}\right) \left(\sum_{i' \mid i' \neq i} \sum_{j'=1}^{C} n_{j'j'}\right)} \sqrt{\sum_{i=1}^{C} \left(\sum_{j=1}^{C} n_{j,j}\right) \left(\sum_{i' \mid i' \neq i} \sum_{j'=1}^{C} n_{j'j'}\right)}} \\ \end{split}$							
 RELEVANT PITFALLS MCC is prevalence-dependent, thus not comparable across data sets with different prevalences (Fig. SN 2.7 in [Reinke et al., 2023], [Reinke et al., 2021]). MCC does not allow for incorporating unequal severity of confusions across classes. This implies that the metric is not applicable to cost-benefit analyses (Fig. SN 2.9 in [Reinke et al., 2023]), or when target classes are related on an ordinal scale (Fig. 4b in [Reinke et al., 2023]). The theoretical lower bound of MCC (-1) may not always be achievable (Fig. SN 2.35 in [Reinke et al., 2023]). MCC is hard to interpret [Zhu 2020]. Compared to other metrics like EC, MCC lacks a framework to identify and validate the decision rule applied to predicted class scores [Ferrer 2022]. MCC depends on the definition of TN (undefined for ObD and InS). 							
REC depends on the definition of the (undefined for Obb and ins). RECOMMENDATIONS - MCC should not be used/used with care if • class confusions are of unequal severity (example: ordinal target classes). • the provided class prevalences do not reflect the population of interest. • there is a mismatch between class prevalences and class importance. • compensation for class imbalance is not requested. • Otherwise, MCC should be used as a multi-class metric specifically if all basic error rates (Sensitivity, Specificity, PPV, NPV) should be captured in one score. • MCC scores should be carefully interpreted in the presence of class imbalance as the distribution between skewed [7bu 2020].							

Reinke et al., Nature Methods, 2024

Meaningful Metrics

Metrics Reloaded

• Pitfalls related to the inadequate choice of the problem category









- Pitfalls related to poor metric selection
 - Disregard of the domain interest



a Overlap-based metrics disregard structure boundaries





DSC = 0.78



b Common multi-class metrics ignore ordinal grading

		Reference	Prediction 1	Prediction 2
	Patient 1	Class 0	Class 0	Class O
Ordinal classes				
	Patient 2	Class 1	Class 1	Class 1
	Patient 3	Class 2	Class 0 🗙	Class 1 🗸
			Accuracy = 0.67 EC = 0.83	= Accuracy = 0.67



- Pitfalls related to poor metric selection
 - Disregard of the properties of the dataset

Common metrics yield implausible results in the presence of class imbalance



Annotation errors may have huge impact on metric scores





- Pitfalls related to poor metric application
 - **b** Simple averaging disregards non-independence of test data





Perspective Published: 12 February 2024

Understanding metric-related pitfalls in image analysis validation

Annika Reinke [⊠], <u>Minu D. Tizabi</u> [⊠], <u>Michael Baumgartner</u>, <u>Matthias Eisenmann</u>, <u>Doreen Heckmann-</u> <u>Nötzel</u>, <u>A. Emre Kavur</u>, <u>Tim Rädsch</u>, <u>Carole H. Sudre</u>, <u>Laura Acion</u>, <u>Michela Antonelli</u>, <u>Tal Arbel</u>, <u>Spyridon</u> <u>Bakas</u>, <u>Arriel Benis</u>, <u>Florian Buettner</u>, <u>M. Jorge Cardoso</u>, <u>Veronika Cheplygina</u>, <u>Jianxu Chen</u>, <u>Evangelia</u> <u>Christodoulou</u>, <u>Beth A. Cimini</u>, <u>Keyvan Farahani</u>, <u>Luciana Ferrer</u>, <u>Adrian Galdran</u>, <u>Bram van Ginneken</u>, <u>Ben Glocker</u>, ... <u>Lena Maier-Hein</u> [⊠] + Show authors

Nature Methods 21, 182–194 (2024) Cite this article

Perspective Published: 12 February 2024

Metrics reloaded: recommendations for image analysis validation

Lena Maier-Hein [⊠], Annika Reinke [⊠], Patrick Godau, Minu D. Tizabi, Florian Buettner, Evangelia Christodoulou, Ben Glocker, Fabian Isensee, Jens Kleesiek, Michal Kozubek, Mauricio Reyes, Michael A. Riegler, Manuel Wiesenfarth, A. Emre Kavur, Carole H. Sudre, Michael Baumgartner, Matthias Eisenmann, Doreen Heckmann-Nötzel, Tim Rädsch, Laura Acion, Michela Antonelli, Tal Arbel, Spyridon Bakas, Arriel Benis, … Paul F. Jäger [⊠] + Show authors

Nature Methods 21, 195–212 (2024) | Cite this article





- Metric implementation
 - Use reference implementations
 - Read metric-specific recommendations in the cheat sheets
- Aggregation
 - Address the potential correlation between classes when aggregating
 - Complement validation with multi-class metrics
 - Respect the hierarchical data structure when aggregating metrics
 - Leverage metadata to reveal potential algorithmic bias
 - Follow category-specific aggregation strategy
- Interpretation
 - Read metric-related recommendations to obtain awareness of the pitfall
 - Report on the quality of the reference (e.g. intra-rater and inter-rater variability).

Maier-Hein et al., Nature Methods, 2024



Ensuring Robust AI in Medical Imaging

- Validation on Real-World Data
- Meaningful Metrics
- Rigorous Statistical Analysis



A (fake & slightly exaggerated) typical MICCAI conclusion

'Our proposed deep learning framework achieves superior accuracy (86.6%) and robust generalizability in distinguishing Alzheimer's disease patients from cognitively normal subjects, outperforming established methods on the ADNI dataset. The model provides anatomically interpretable results by highlighting disease-relevant regions such as the hippocampus, demonstrates high reproducibility and fast inference, and, by eliminating manual feature engineering while aligning with clinical biomarkers, is both practical and explainable-making it ready for immediate integration into routine diagnostic workflows.'



Common practice in medical imaging algorithm performance analysis

Commonly encountered results tables





Common practice in medical imaging algorithm performance analysis





Standard deviation & confidence interval

Standard deviation (SD)

Measure of the dispersion or spread of data points from the mean value

• Confidence interval (CI)

Range within which a population parameter is expected to lie with a certain level of confidence.





Common practice with respect to variability reporting





Approximation of the standard deviation





Performance differences versus CI widths





Conclusion

 Current publications typically do not provide sufficient evidence to support which models could potentially be translated into clinical practice

Recommendations

- Incorporate robust statistical analyses
- Report performance variability
- Use sufficiently large test sets to substantiate claims of outperformance

Ensuring Robust AI in Medical Imaging: Validation on Real-World Data,

Meaningful Metrics, and Rigorous Statistical Analysis



Research vs clinical data







