

Benchmarking Evaluation Metrics for Tubular Structure Segmentation in Biomedical Images

Meghane Decroocq, Charissa Poon, Matthias Schlachter, Henrik Skibbe

Brain Image Analysis Unit, Center for Brain Science, RIKEN, Wako, Japan



JSPS



RIKEN CBS
脳神経科学研究センター

ShapeMI MICCAI Workshop, 27th Sept. 2025

Introduction

- **Tubular tree-like structures** are common in biology and medicine : vascular networks, airways, neuronal trees
- The **segmentation** of these structures is crucial for many downstream applications : tracing, numerical simulation

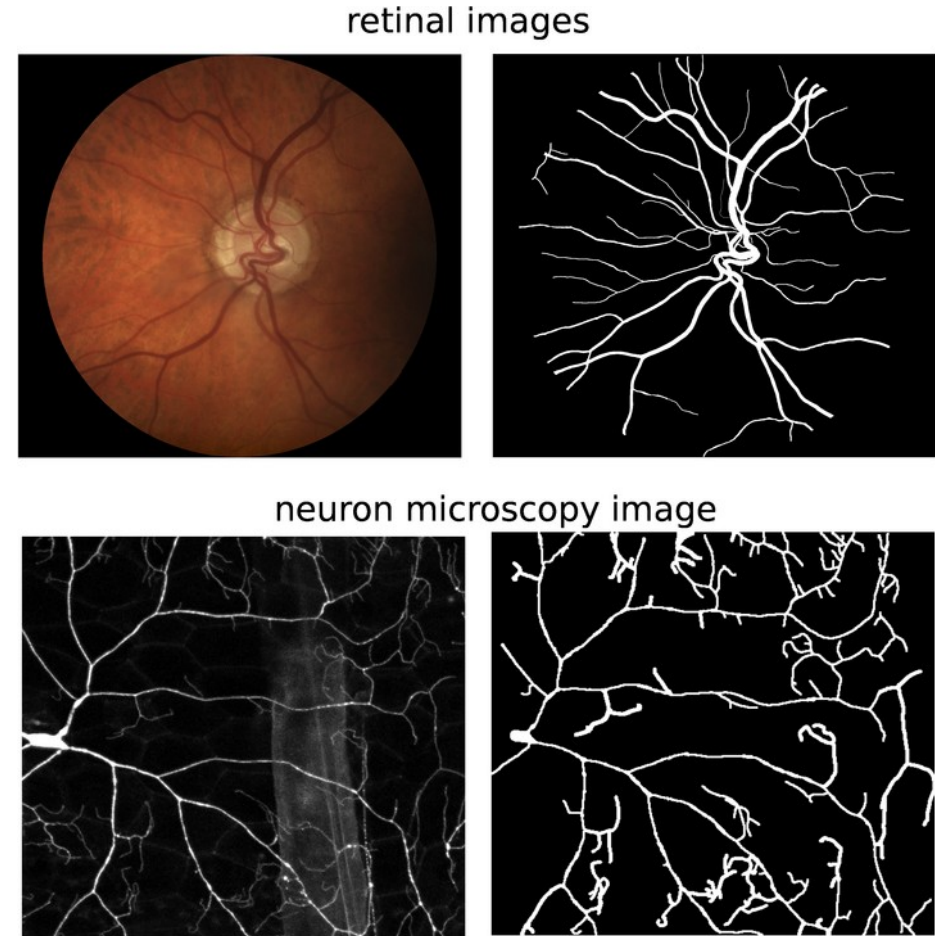
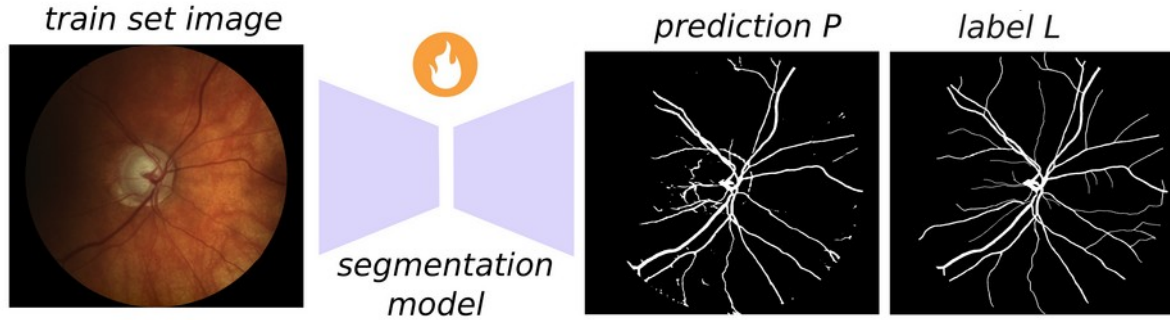


Figure 1. Example of biomedical images showing tubular, tree-like structures and the target segmentation.

Introduction

- **Deep learning** has led to tremendous progress in biomedical image segmentation.

1) model training



2) model evaluation

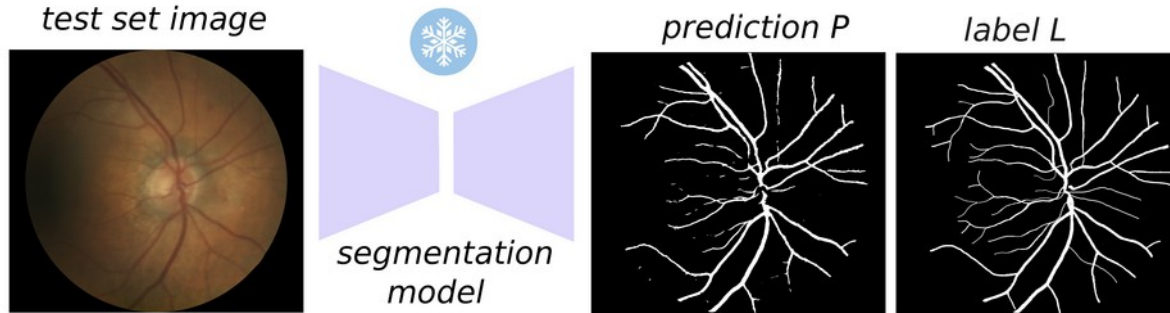


Figure 2. Importance of quality metrics.

Introduction

- To **train and evaluate** segmentation models, we need to design **metrics** measuring the **quality** of a predicted segmentation compared to the ground-truth labels.

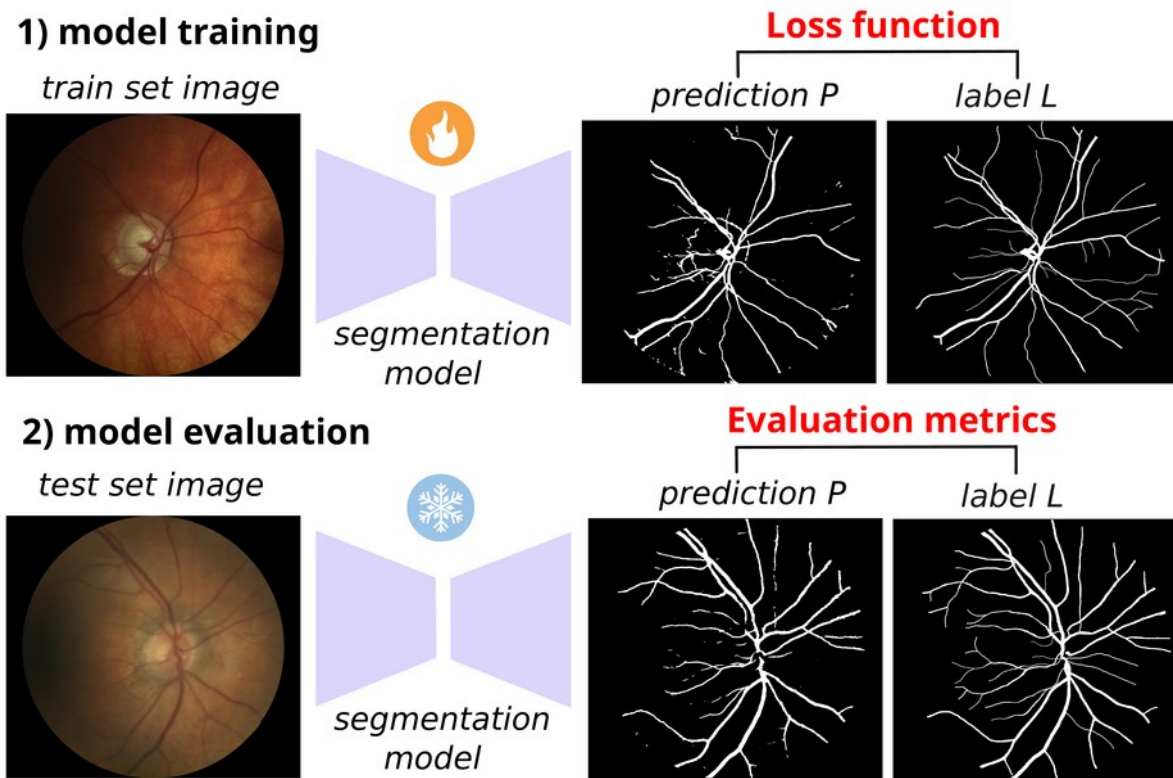


Figure 2. Importance of quality metrics.

Introduction

- Tubular structure segmentation offers unique **challenges**, leading to the design of **various metrics and loss functions** in recent years.



Figure 3. Challenges of tubular structure segmentation.

Introduction – Overlap-based quality metrics

- The **Dice** score is based on the pixel-wise mismatch between label and prediction.

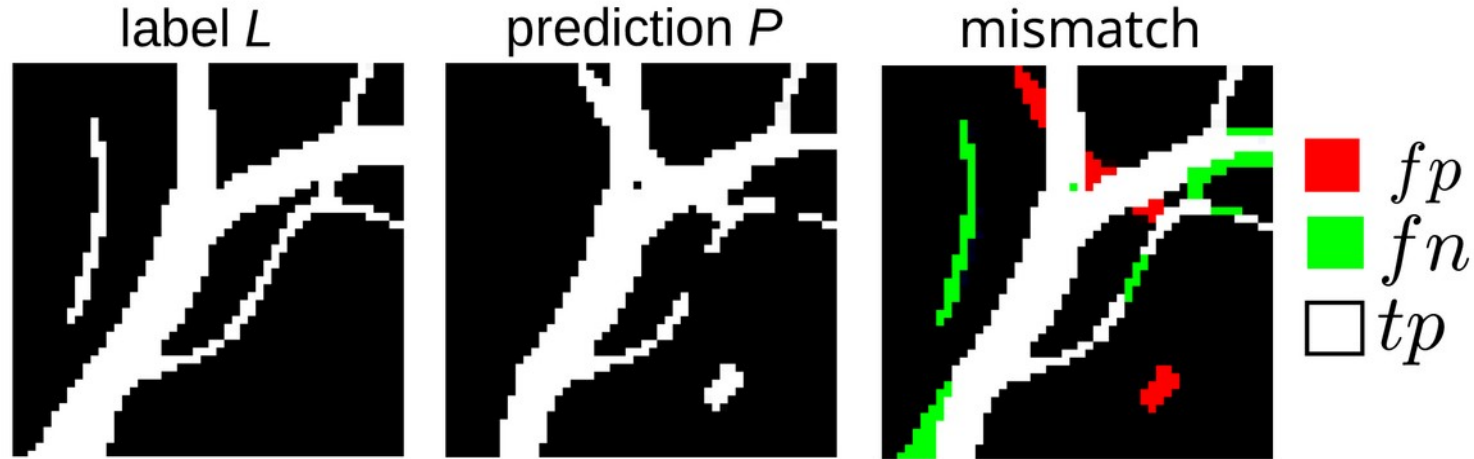


Figure 4. Illustration of the Dice calculation.
(tp = true positives, fn = false, fp = false positives)

Introduction – Skeleton-based quality metrics

- The **clDice**¹ relies on the **skeleton** mismatch.
- The **cbDice**² includes the **distance to boundaries**.

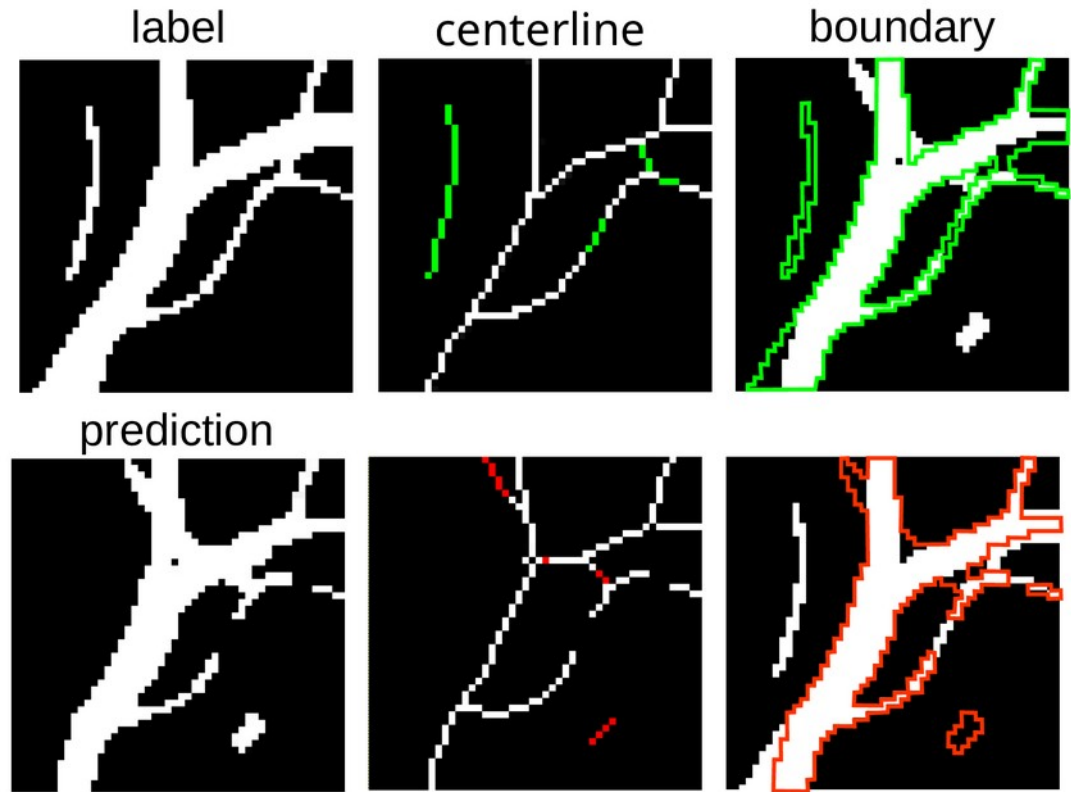


Figure 5. Illustration of the *clDice* and *cbDice*.

Introduction – Algebraic topology quality metrics

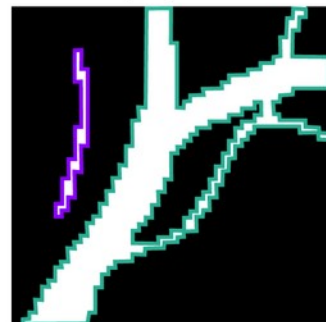
β_0 = number of **connected components CC**

β_1 = number of **holes** in the image

- The **Betti error** β^{err} is the difference between the Betti numbers of the label and prediction.

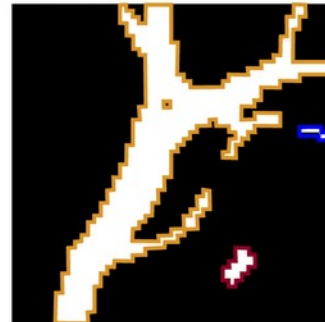
a) connected components

$$\beta_0(L) = 2$$



label L

$$\beta_0(P) = 3$$



prediction P

$$\beta_0^{err} = 1$$

b) holes

$$\beta_1(L) = 1$$



label L

$$\beta_1(P) = 1$$



prediction P

$$\beta_1^{err} = 0$$

Figure 6. Illustration of the Betti error calculation.

Introduction – Algebraic topology quality metrics

- The **Betti matching error**³ μ^{err} and the **ccDice**⁴ are matching CC and holes.

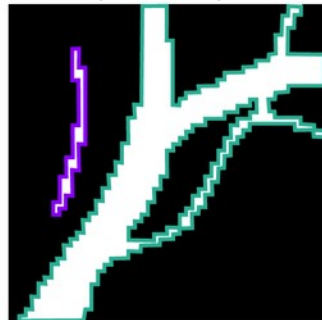
³Stucki, Nico, et al. "Topologically faithful image segmentation via induced matching of persistence barcodes." International Conference on Machine Learning, 2023.

⁴Rougé, Pierre, Odysée Merveille, and Nicolas Passat. "ccDice: A topology-aware Dice score based on connected components." International Workshop on Topology-and Graph-Informed Imaging Informatics, MICCAI 2024.

Figure 7. Illustration of the Betti matching and ccDice calculation.

a) connected components

$$\overline{\mu}_0(L, P) = 1$$



label L

$$\overline{\mu}_0(P, L) = 3$$



prediction P

$$\mu_0^{err} = 3$$

$$ccDice = 0.4$$

b) holes

$$\overline{\mu}_1(L, P) = 1$$



label L

$$\overline{\mu}_1(P, L) = 1$$



prediction P

$$\mu_1^{err} = 2$$

Challenge

- Which metric / loss shall I use? How to be sure than the metric reflects my expectations?
 - Quality is **subjective**, making it difficult to assess the **strengths and weaknesses** of each metric.

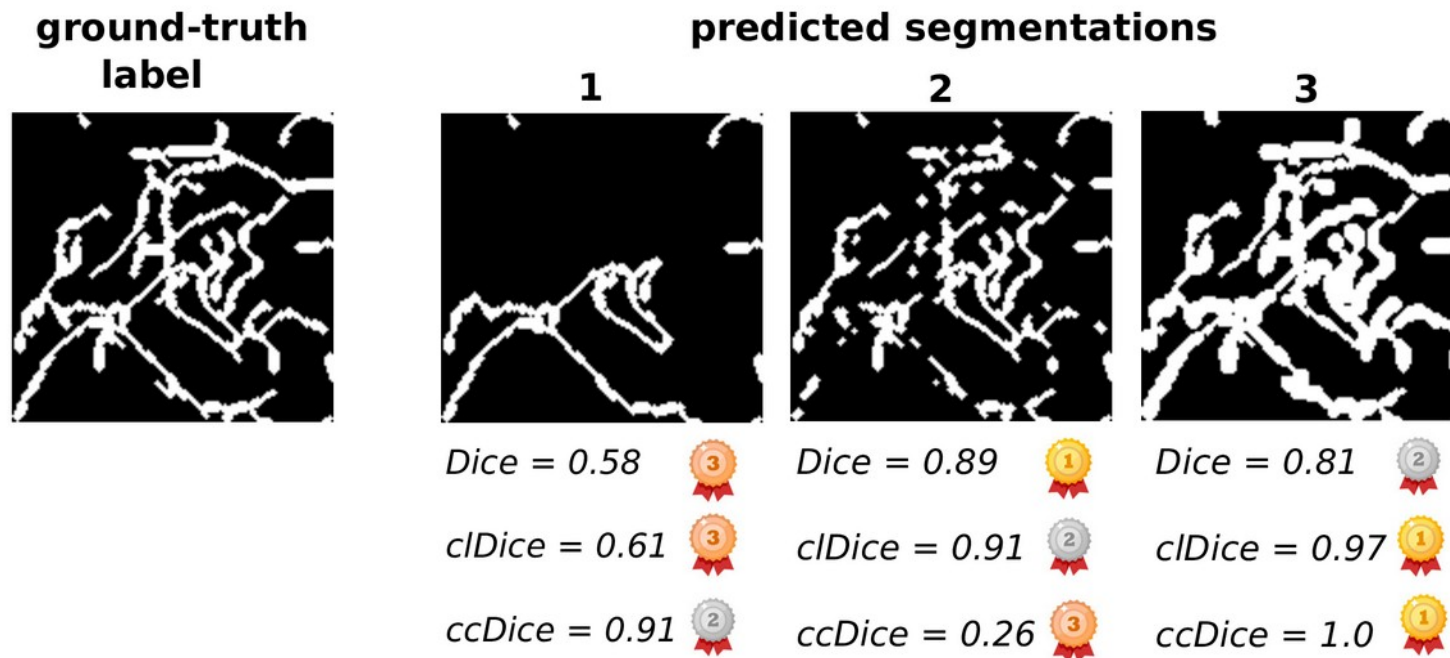


Figure 8. Which is the best segmentation?

How the evaluate quality metrics?

- 1) Use handcrafted examples to show the metric behavior in specific cases^{2,4}
- 2) Ask experts to identify pitfalls and make recommendations (*Metrics Reloaded project*)^{5,6}
- 3) Evaluate metrics based on their correlation to subjective visual scores attributed by experts⁷

⁵Maier-Hein, Lena, et al. "Metrics reloaded: recommendations for image analysis validation." Nature methods, 2024

⁶Reinke, Annika, et al. "Understanding metric-related pitfalls in image analysis validation." Nature methods, 2024

⁷Aydin, Orhun Utku, et al. "An evaluation of performance measures for arterial brain vessel segmentation." BMC medical imaging, 2021

State-of-the-art

How to evaluate quality metrics?

- 1) Use handcrafted examples to show the metric behavior in specific cases^{2,4}
- 2) Ask experts to identify pitfalls and make recommendations (*Metrics Reloaded project*)^{5,6}
- 3) Evaluate metrics based on their correlation to subjective visual scores attributed by experts⁷

Limitations

- No study focusing on **topology-preserving metrics** so far
- Relying on **experts** to grade images or identify pitfalls : **time-consuming**, hard to apply to new metrics
- Limited to the **specific contexts**, such as particular application or dataset.

State-of-the-art

How to evaluate quality metrics?

- 1) Use handcrafted examples to show the metric behavior in specific cases^{2,4}
- 2) Ask experts to identify pitfalls and make recommendations (*Metrics Reloaded project*)^{5,6}
- 3) Evaluate metrics based on their correlation to subjective visual scores attributed by experts⁷

Limitations

- No study focusing on ~~topology-preserving metrics~~ so far
- Relying on ~~experts~~ to grade images or identify pitfalls : ~~time-consuming~~, hard to apply to new metrics
- Limited to the ~~specific contexts~~, such as particular application or dataset.

→ We propose a new approach to benchmark metrics without the need for expert knowledge!

Proposed approach

- Depending on the application, different segmentation errors have different importance or “**weight**”.

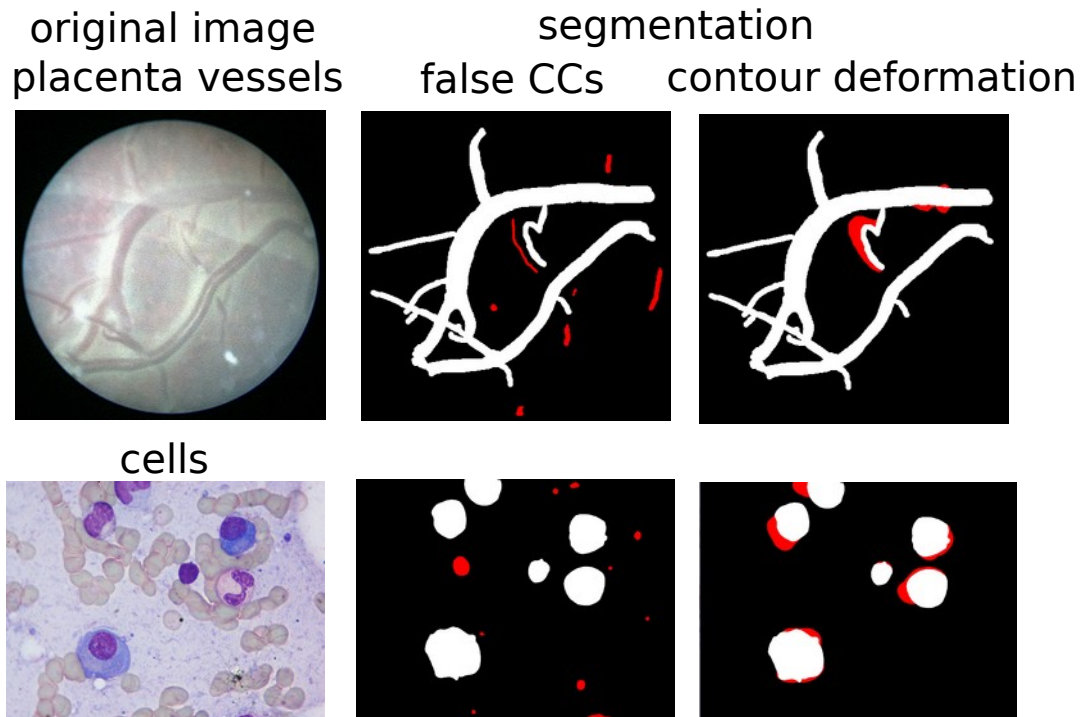


Figure 9. Illustration of different types of segmentation errors. Images from the PSVFM⁸ and Neurips 2022 dataset⁹.

⁸Bano, Sophia, et al. "Deep placental vessel segmentation for fetoscopic mosaicking." MICCAI, 2020.

⁹Ma, Jun, et al. "The multimodality cell segmentation challenge: toward universal solutions." Nature methods, 2024

Proposed approach

- **Break down** error types into **easily interpretable categories** (e.g. false component, deformation...)
- Estimate the “**weight**” that a given metric attributes to each type of errors.

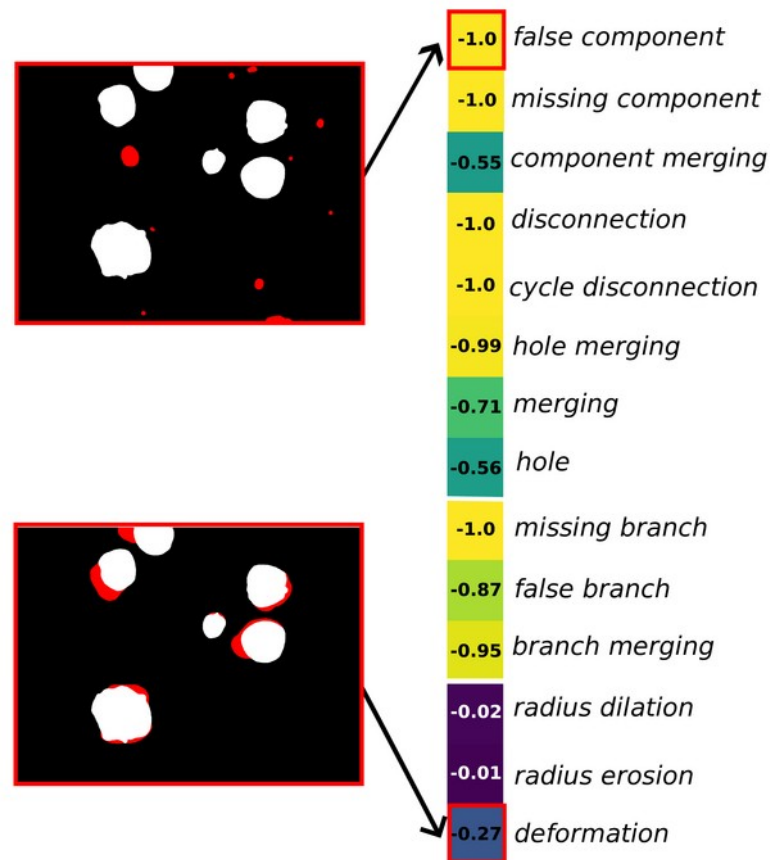
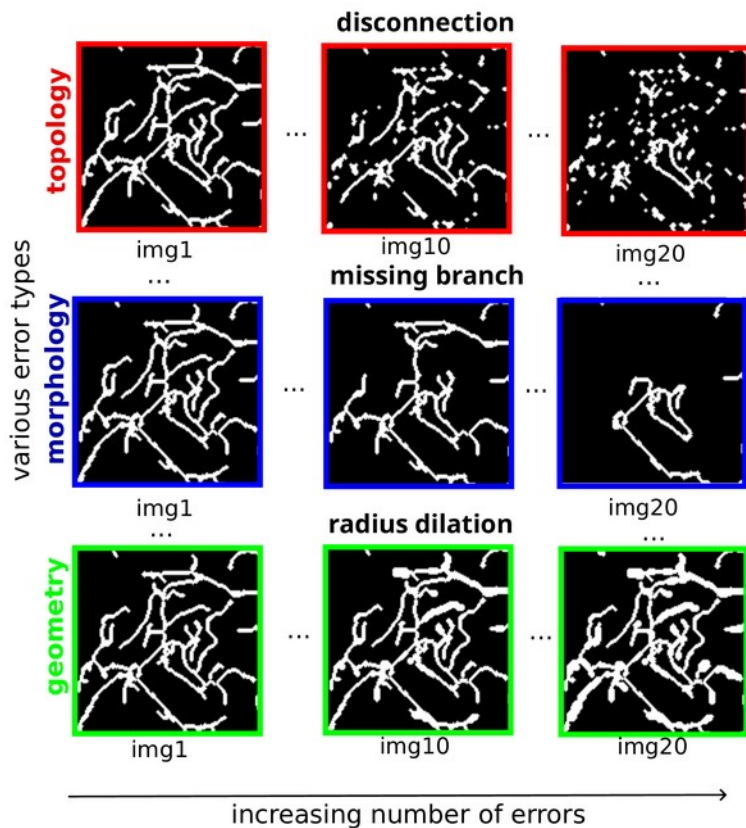


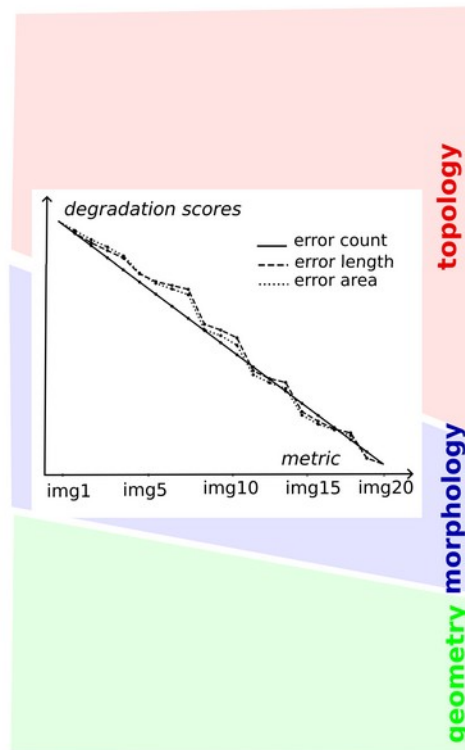
Figure 10. Estimated weights for each type of segmentation errors.

Proposed approach

2 Synthetic benchmark dataset



3 Compare metrics and degradation scores



4 Estimate weights for each error category

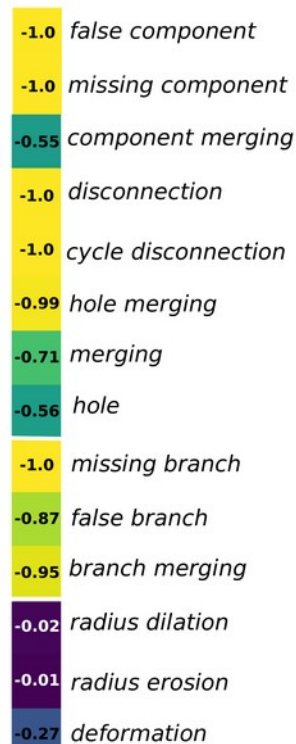


Figure 11. Overview of the proposed method.

Method – Error classification

- We introduce a new distinction between “**topology**”, “**morphology**” and “**geometry**”.

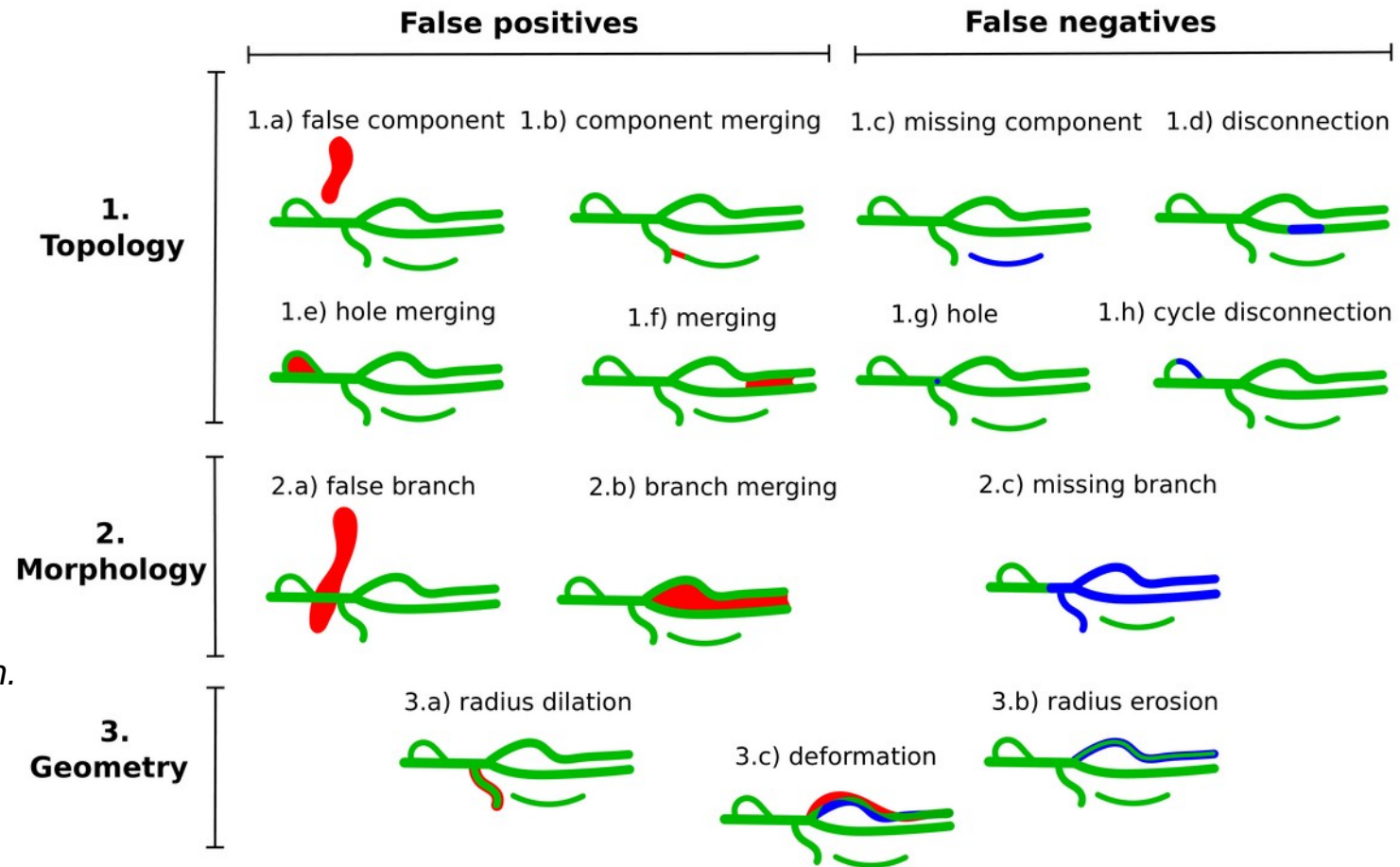
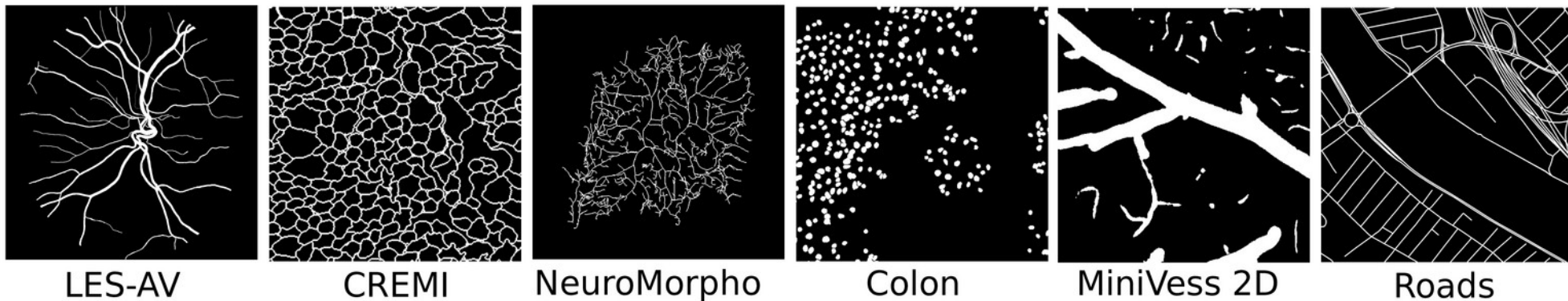


Figure 12. The proposed Segmentation error classification.

Method – Synthetic benchmark dataset

- We selected 3 labels from 6 public datasets.



¹⁰CreMi. miccai challenge on circuit reconstruction from electron microscopy images <https://cremi.org/> (2016), accessed: 2024-02-12

¹¹Ljosa, V., Sokolnicki, K.L., Carpenter, A.E.: Annotated high-throughput microscopy image sets for validation. *Nature methods*, 2012

¹²Tecuatl, C., Ljungquist, B., Ascoli, G.A.: Accelerating the continuous community sharing of digital neuromorphology data. *FASEB BioAdvances*, 2024

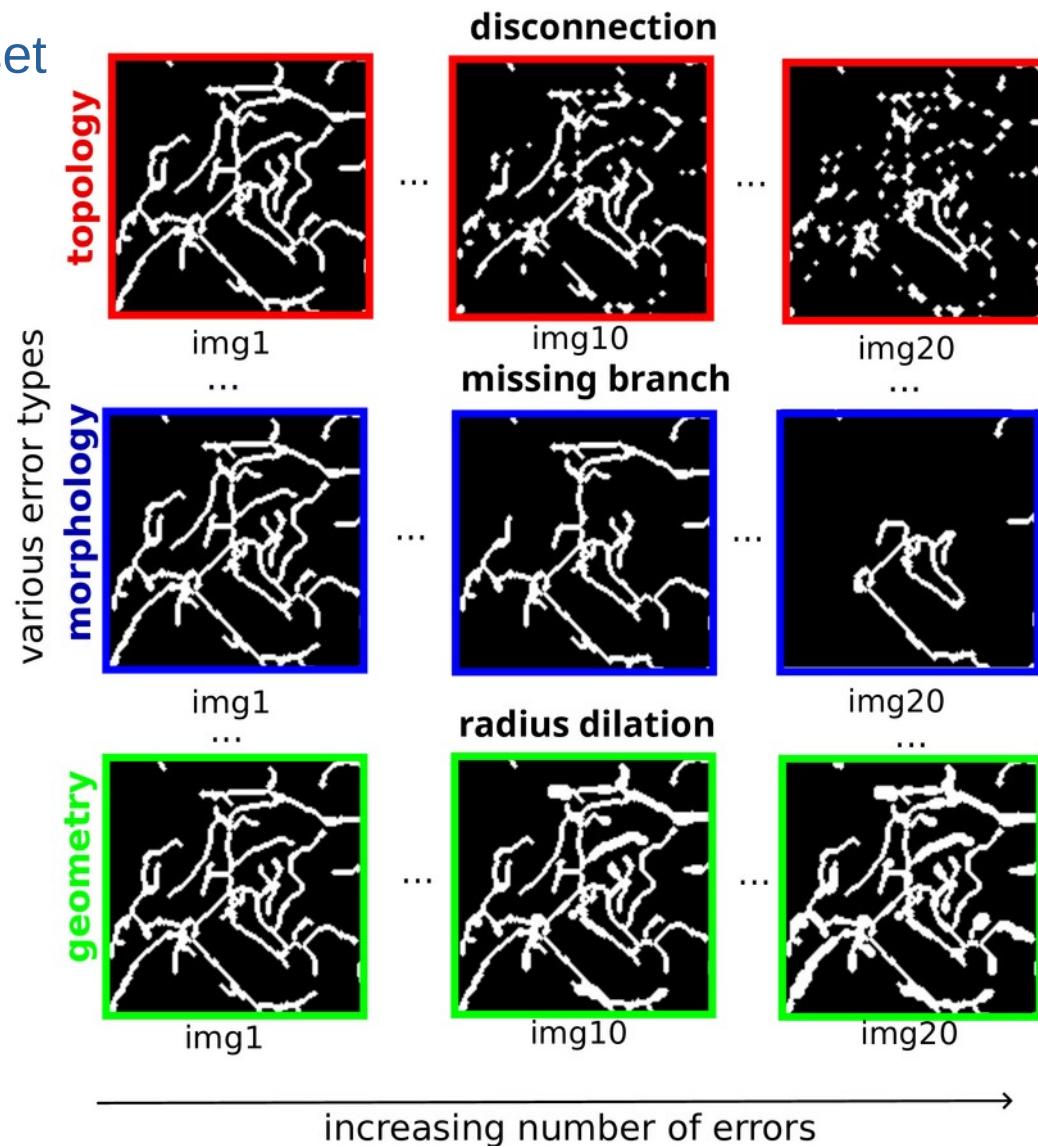
¹³Mnih, V.: Machine Learning for Aerial Image Labeling. Ph.D. thesis, University of Toronto, 2013

¹⁴Odstřilík, J., Kolar, R., Budai, A., Hornegger, J., Jan, J., Gazarek, J., Kubena, T., Cernosek, P., Svoboda, O., Angelopoulou, E.: Retinal vessel segmentation by improved matched filtering: evaluation on a new high-resolution fundus image database. *IET Image Processing*, 2013

¹⁵Poon, C., Teikari, P., Rachmadi, M.F., Skibbe, H., Hynynen, K.: A dataset of rodent cerebrovasculature from in vivo multiphoton fluorescence microscopy imaging. *Scientific Data*, 2023

Method – Synthetic benchmark dataset

- We **cumulatively add errors** from each category to mimic predicted segmentations.



Method – Degradation scores

- We design **degradation scores** to quantify the degradation of synthetic images. We consider three **error properties** : the error “*area*”, the error “*length*” and the error “*count*”.

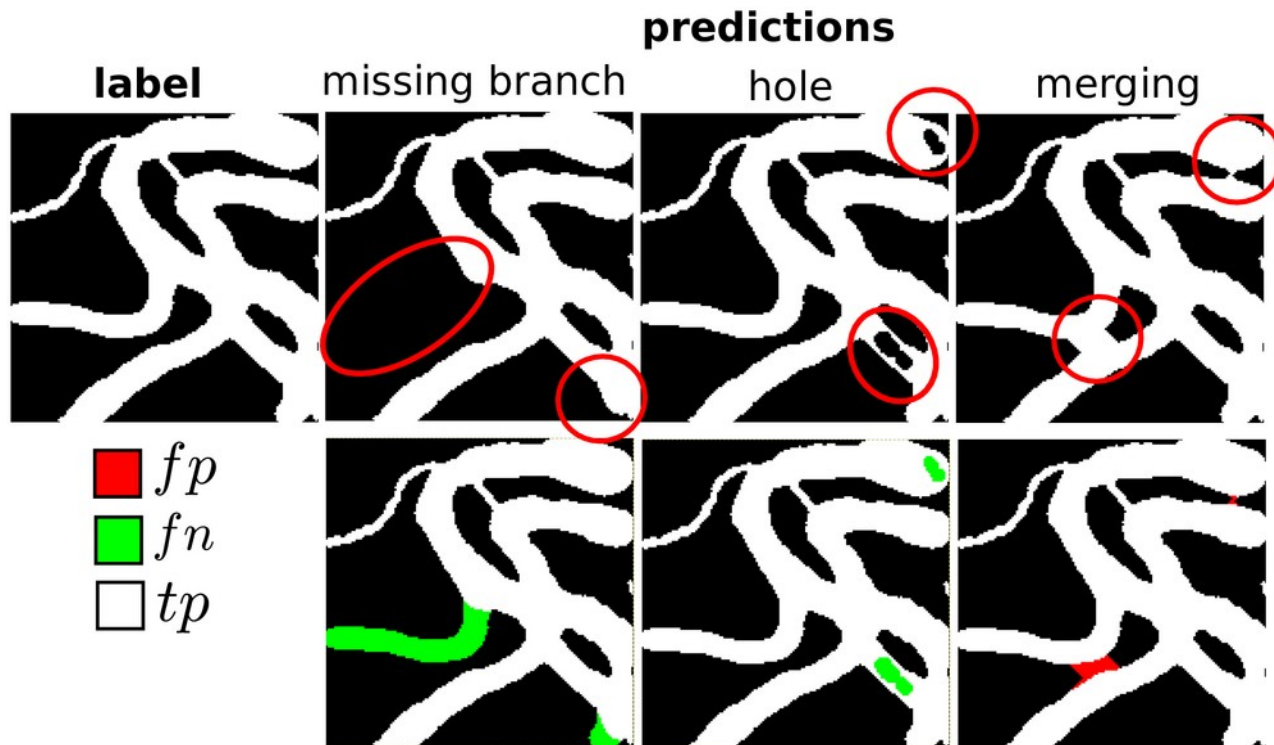


Figure 13. Calculation of the degradation scores for the property “area”.

Method – Degradation scores

- We design **degradation scores** to quantify the degradation of synthetic images. We consider three **error properties** : the error “*area*”, the error “*length*” and the error “*count*”.

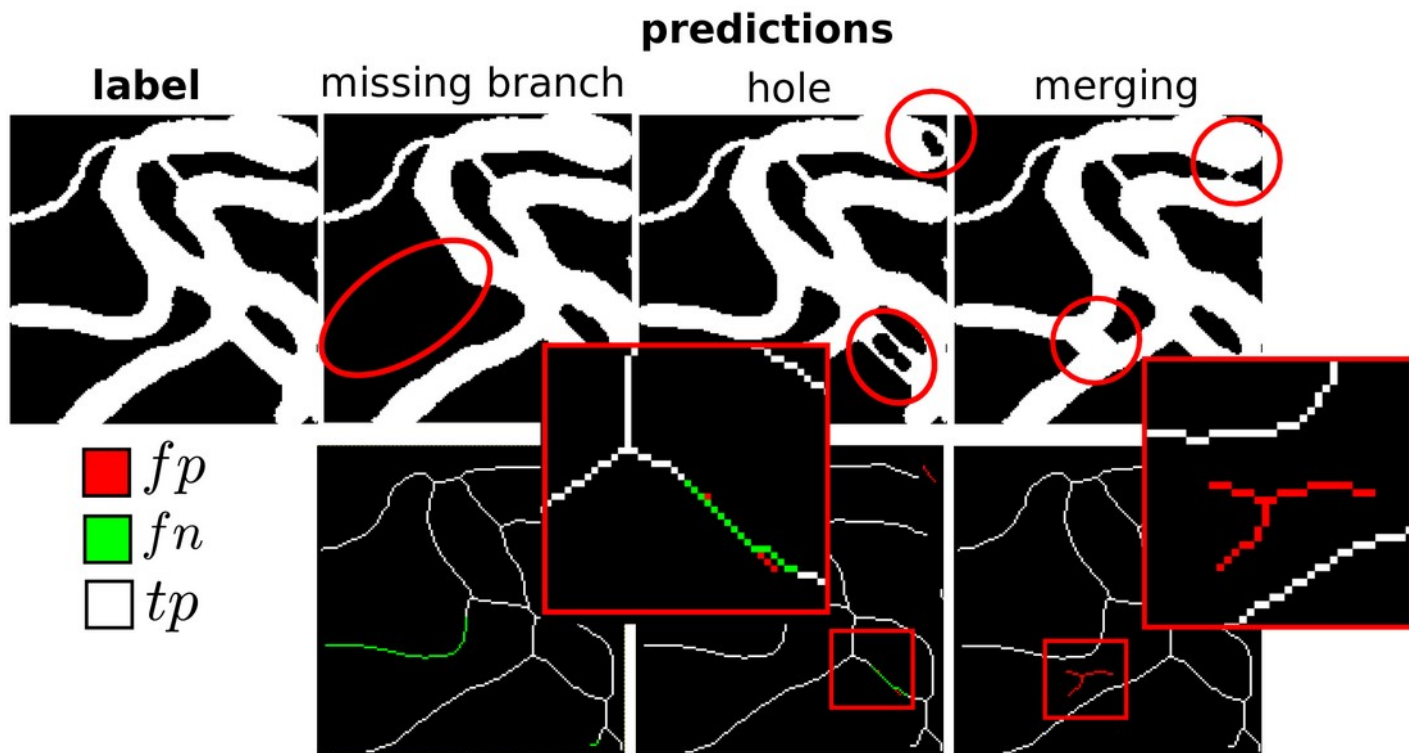
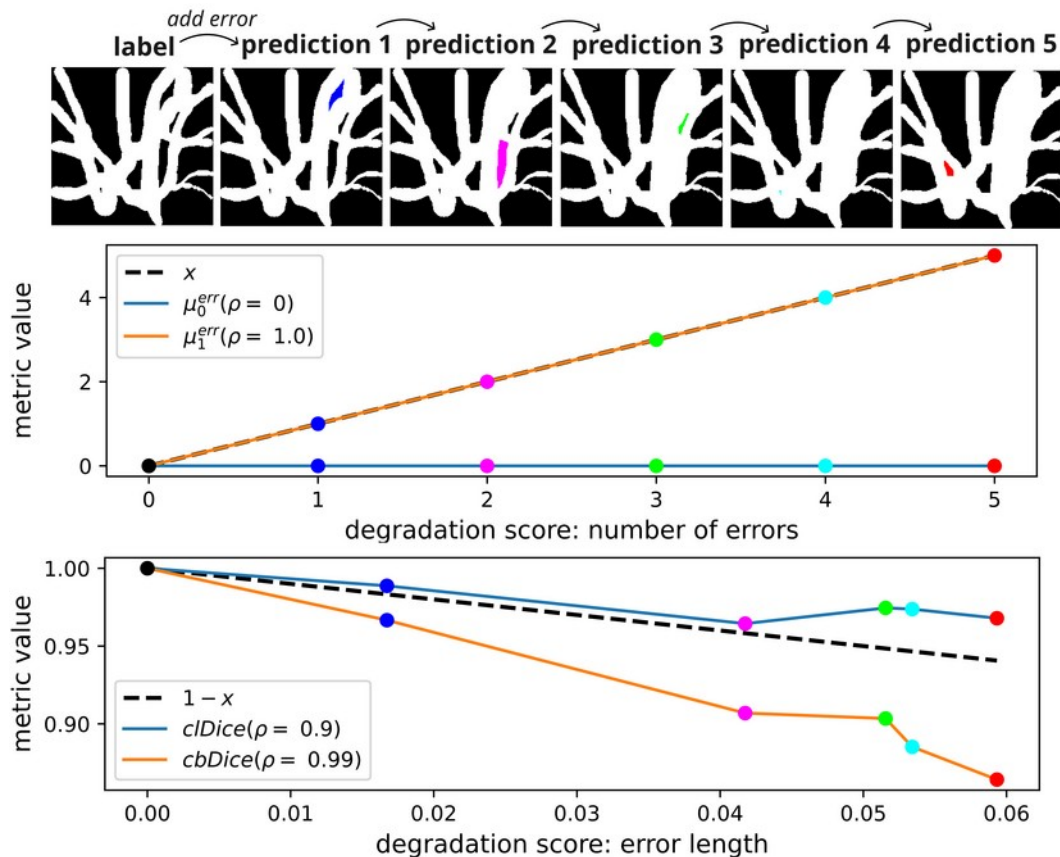


Figure 14. Calculation of the degradation scores for the property “length”.

Method – Correlation and weights estimation



- **Correlation** (absolute Pearson coefficient)

= Is the metric sensitive to one type of error ($\rho = 1$) or not ($\rho = 0$)?

- **Weight** (slope of fitted linear function)

= How strongly does the metric penalize a given type of error?

Results – Correlation analysis

($\rho = 0$) : Does not account for this type of error
 ($\rho > 0.9$) : Good estimator of the type of error
 ($\rho < 0.9$) : Influenced by the error (limitation?)

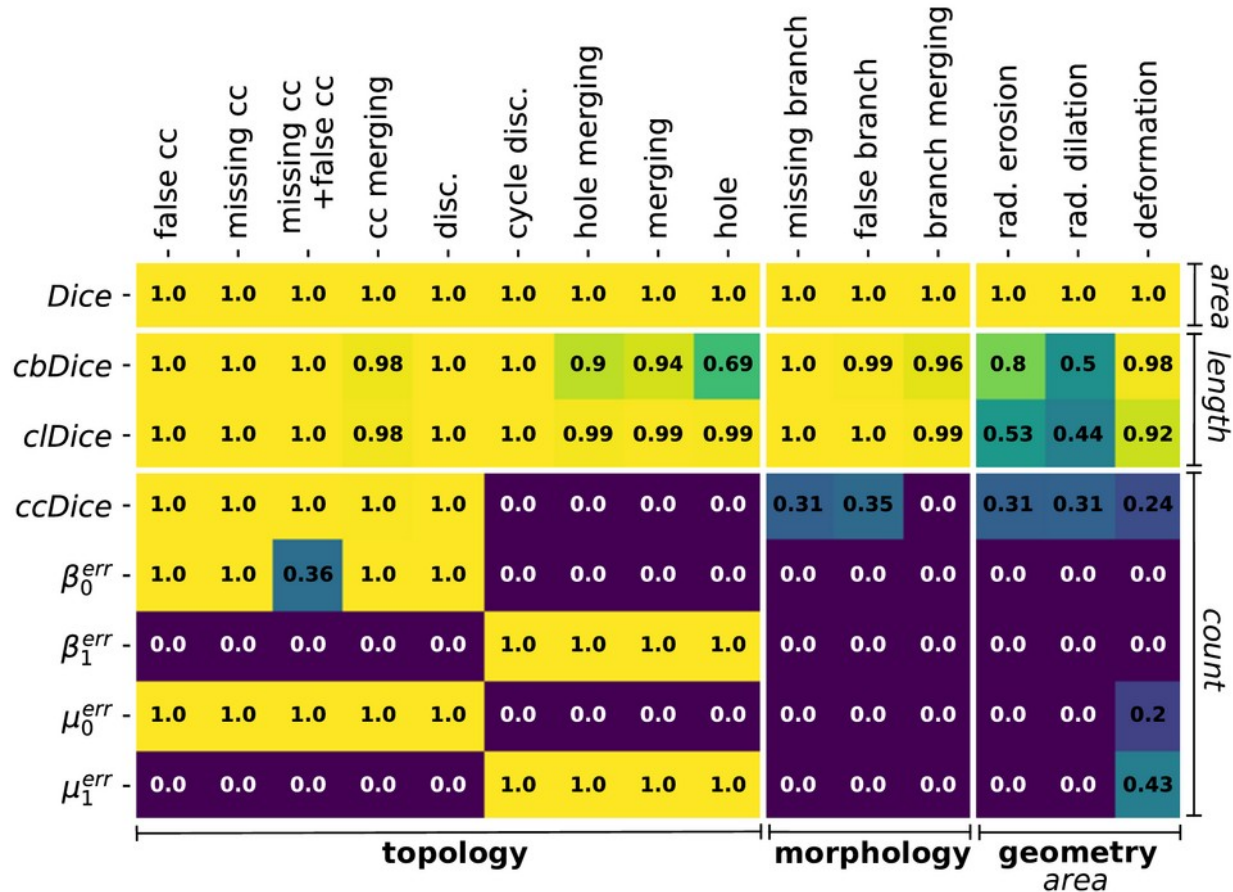


Figure 15. Correlation array.

Results – Correlation analysis

($\rho = 0$) : Does not account for this type of error
 ($\rho > 0.9$) : Good estimator of the type of error
 ($\rho < 0.9$) : Influenced by the error (limitation?)

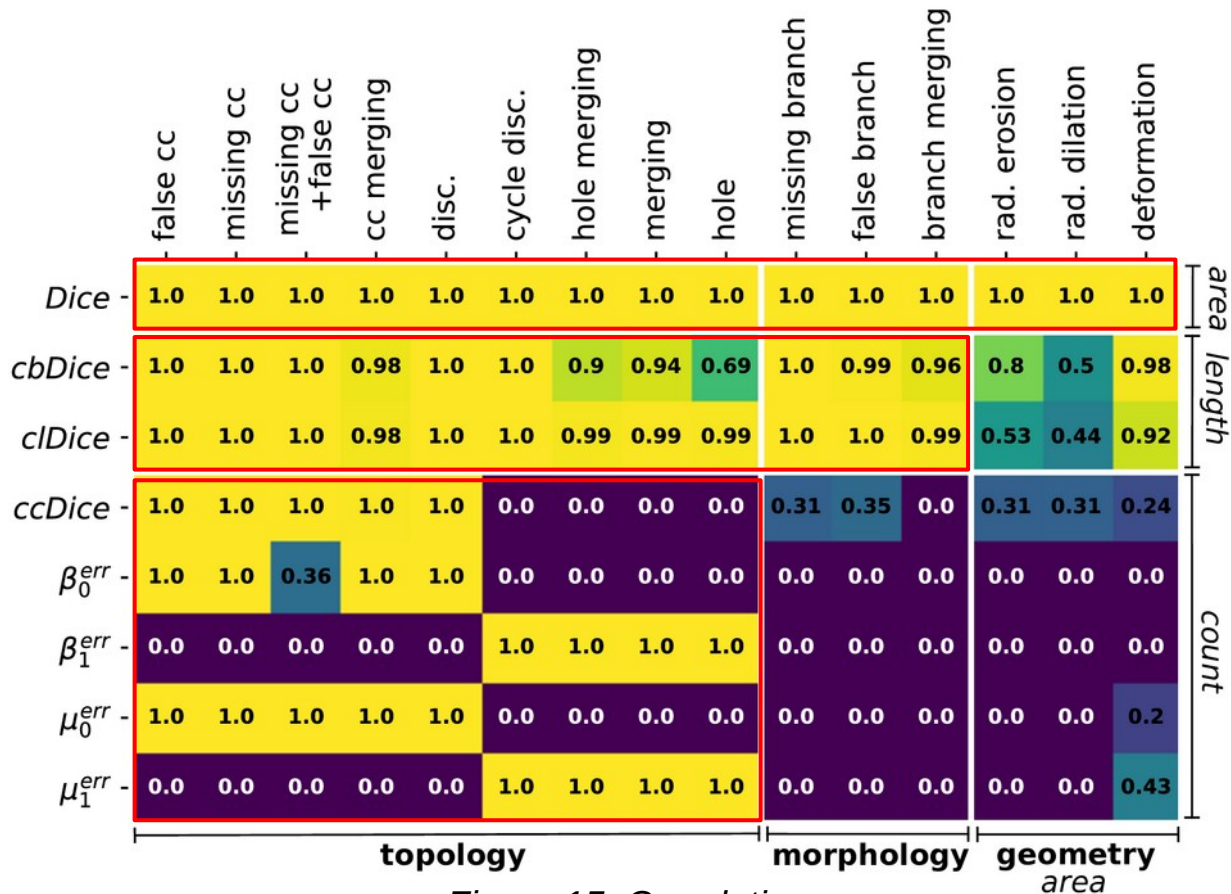


Figure 15. Correlation array.

Results – Correlation analysis

- The cIDice is affected by **radius change and deformations**.

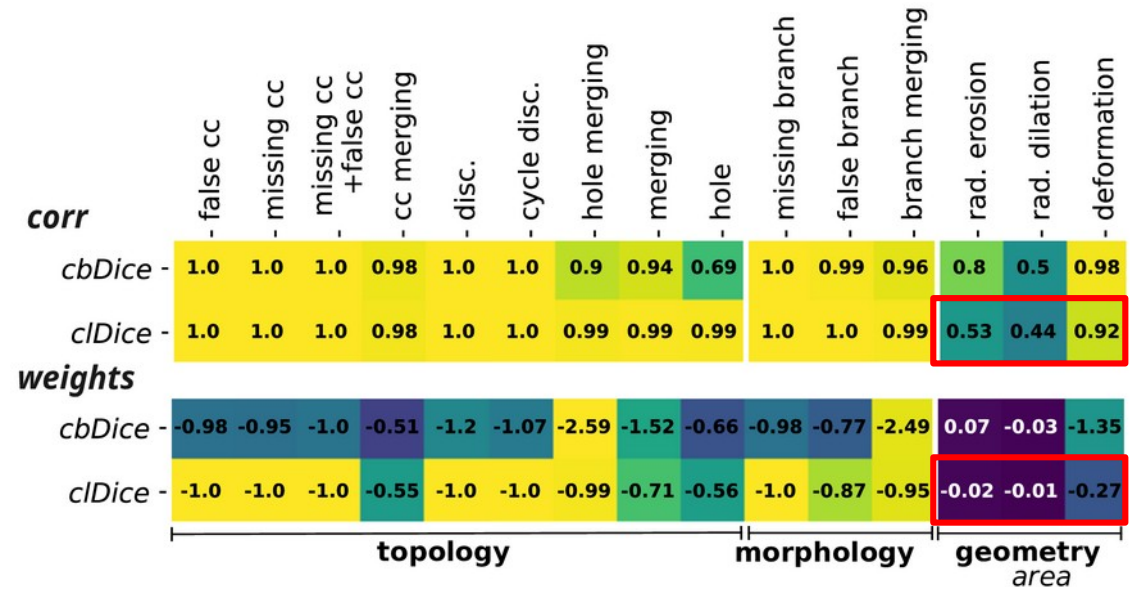


Figure 16. Correlations and weights for the cIDice and cbDice.

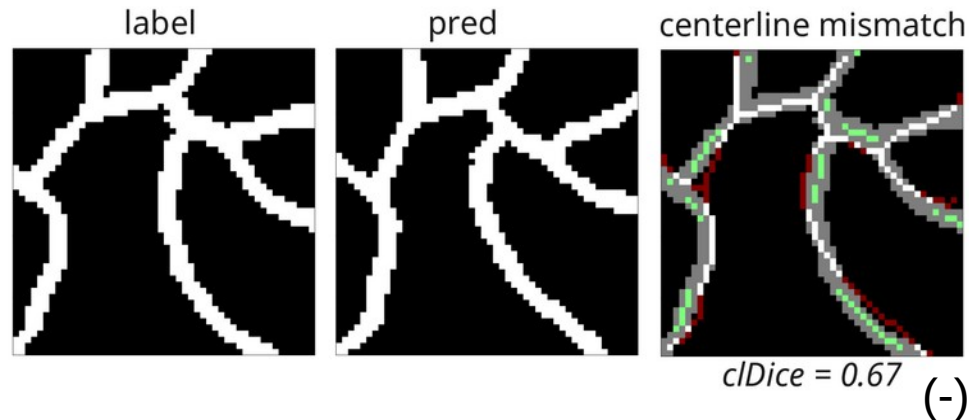


Figure 17. Illustration of the cIDice behavior.

Results – Correlation analysis

- The cIDice may overlook **holes** and **merging**.

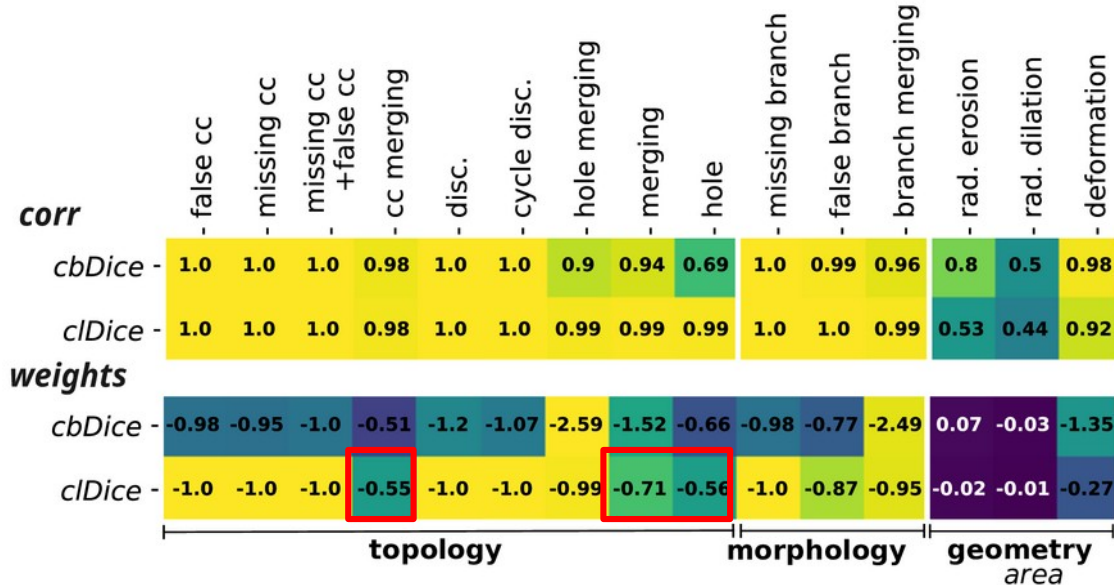


Figure 18. Correlations and weights for the cIDice and cbDice.

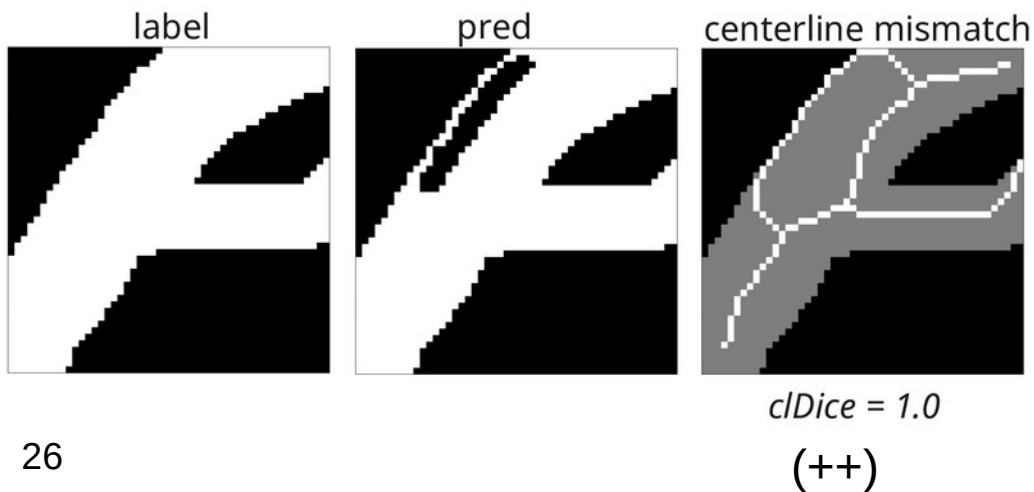


Figure 19. Illustration of the cIDice behavior.

Results – Correlation analysis

- The ccDice and Betti metrics **ignores morphological errors** (e.g. missing branches)

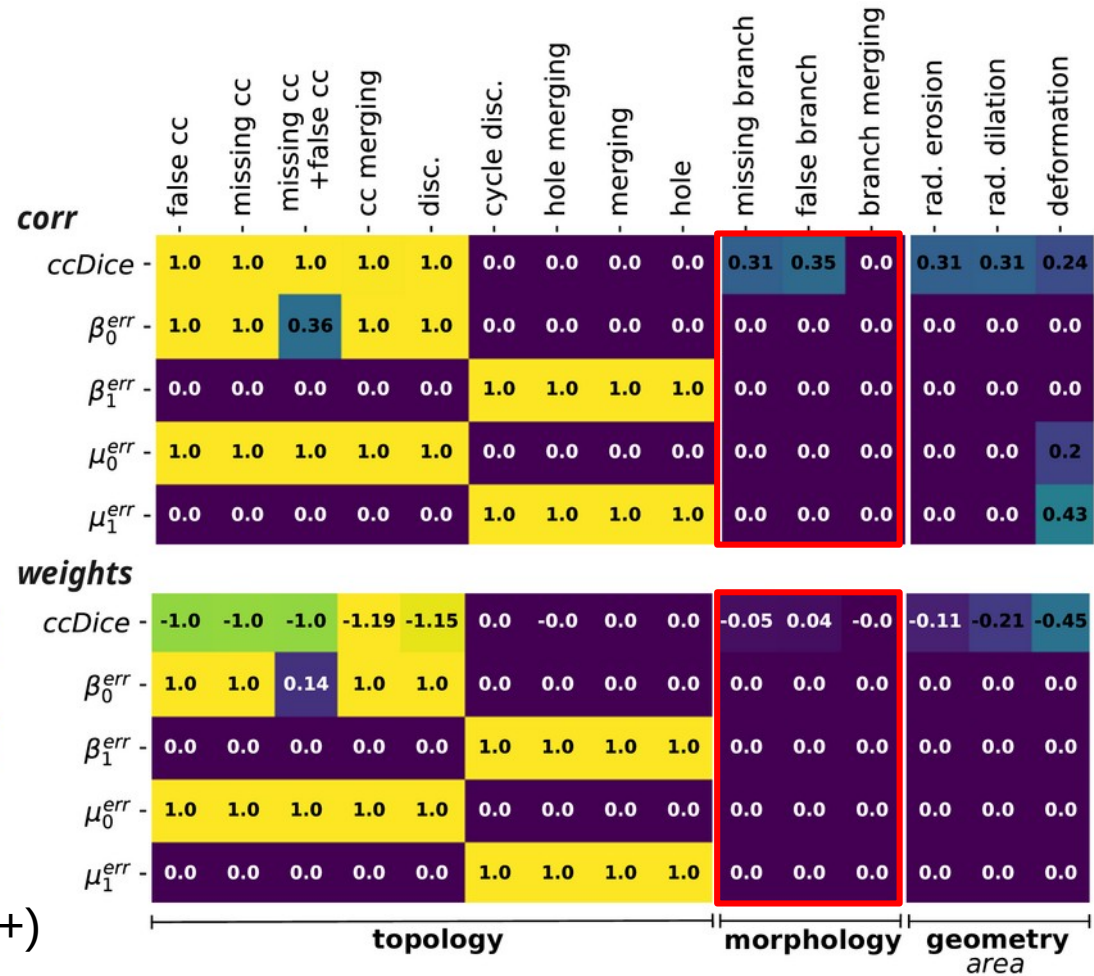
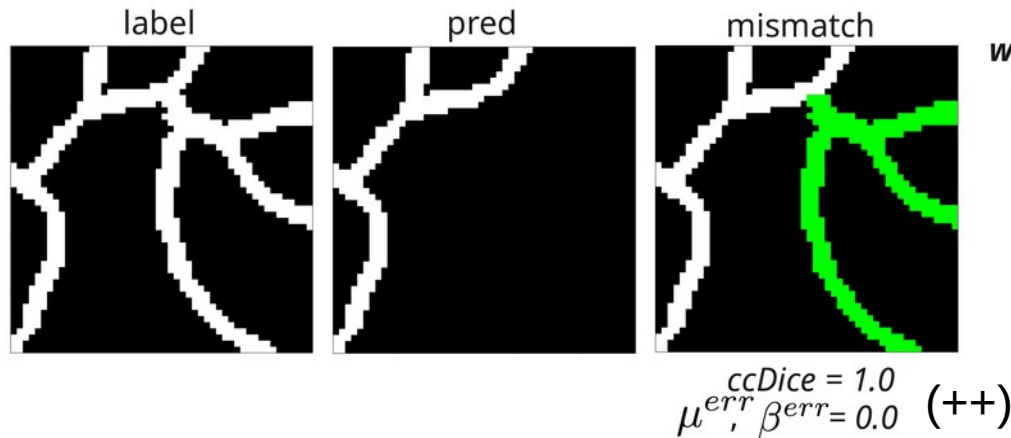


Figure 20. Illustration of the metric's behavior.

Figure 21. Correlations and weights for the ccDice and Betti metrics.

Conclusion

Contributions

- First classification of segmentation errors for tubular structures
- Method to generate synthetic segmentation with a given type of error
- New approach to visualize and interpret the metrics behavior

Conclusion

Contributions

- First classification of segmentation errors for tubular structures
- Method to generate synthetic segmentation with a given type of error
- New approach to visualize and interpret the metrics behavior

Advantages +

- No need for expert knowledge on the metrics
- Automated (easy to apply to new metrics)
- Covers a large range of applications and contexts (can find unexpected pitfalls!)

Conclusion

Contributions

- First classification of segmentation errors for tubular structures
- Method to generate synthetic segmentation with a given type of error
- New approach to visualize and interpret the metrics behavior

Advantages +

- No need for expert knowledge on the metrics
- Automated (easy to apply to new metrics)
- Covers a large range of applications and contexts (can find unexpected pitfalls!)

Limitations -

- Degradation scores may not reflect the desired metric behavior.
- Necessary to consider more error properties (boundary, center-of-mass)