

2024.2, Thermo Fisher Scientific, Waltham, USA), utilizing the standard brush and interpolation tools. The cement lines have low contrast and diffuse borders. Importantly, their visibility and area proportion in the 2D view of the 3D image critically depends on the view angle in 3D, which doesn't always align well with the available view axes (xy, xz, zy). When the cement lines slightly change orientation inside the bone they can become much less pronounced in the current 2D view direction and significantly change in the displayed size. To address these issues, we conducted a pre-annotation for all 3D images from all three view angles (axes: xy, xz, zy) to acquire a coarse initial representation of all cement lines. The pre-annotation corresponds to a detailed annotation about every 50 slices for each view axis for every 3D image. The pre-annotation was used to determine the preferred orientation of the majority of the cement lines within the 3D image. The pre-annotation was further used as an aid in the more thorough annotation to not miss any cement lines due to viewing angle effects. The thorough annotation was done along the most suitable view axis (xy, xz , or zy). This most suitable view axis was chosen in a way that the majority of the cement lines had a small area proportion in the 2D view of the 3D image, see Fig. 6. The thorough annotation corresponds to a detailed annotation about every three to five slices along the most suitable view axis for every 3D image. The thorough annotation for each 3D image was then interpolated with the standard Avizo interpolation tool. This interpolation was designed by Avizo to only interpolate between annotation slices without considering any gray values. However, the resulting annotation was a suitable compromise between annotation time and annotation accuracy. Sometimes, the interpolation tool didn't work as expected, and there were holes within the annotation that didn't get interpolated. To close these interpolation holes, we conducted a series of morphological opening and closing operations until the holes were closed. We repeated this step with additional annotated slices when this approach didn't work initially. We note that the interpolation had also severe trouble interpolating more than one cement line at a time. As a solution, we conducted the thorough annotation and interpolation procedure only for selected parts of a single cement line at once, which we later combined into the complete annotation.

A.4. Datasets

Roads contains aerial images (1500^2) of the road system in Massachusetts, similar to settings in [2] and [6] we use 421 images for training and 49 for testing, after excluding images with a white-masked pixel percentage of more than 4%. HRF-Retina consists of high-resolution fundus images (3504×2336) for retinal vessel segmentation, we use 36 images for training and 9 for testing. Vessap contains volumetric scans ($500^2 \times 50$) of brain vessels in two channels, similar to settings in [2] we use 8 3D images for training and 3 for

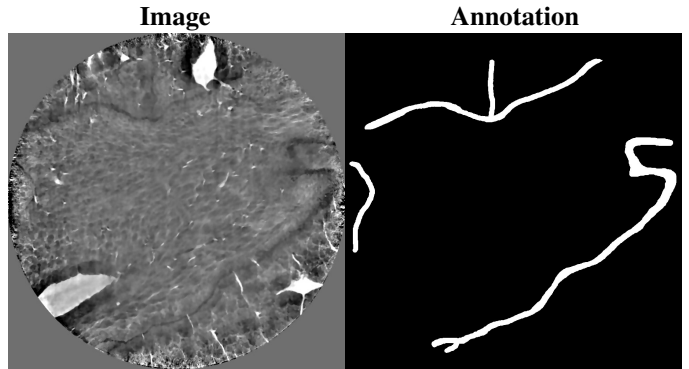


Fig. 6: Cement line Annotation. The cement lines have low contrast and diffuse borders. Also, their appearance is diffuse, note in the left image the changing grayscale characteristics of the cement line in the lower right corner from bottom to right.

testing, using both input channels simultaneously. Our CLD contains volumetric scans ($1024^2 \times 600$) of bone cement lines, we use 13 3D images for training and 4 for testing. We conduct a five-fold cross-validation for all datasets.

We further note that the high-resolution of HRF-Retina is especially suited for examining the influence of the critical pixel mask as the topological errors are represented by more pixels than in lower-resolution images.

A.5. Compared Methods

The nnU-Net [9] is a framework that provides an improved version of the standard U-Net [21] with optimized hyperparameters. The nnU-Net is by default trained with a compound loss, consisting of Dice and Cross-Entropy loss [1].

cIDice [2] is another topology-aware segmentation loss based on critical pixels that focuses on the extracted skeleton of likelihood maps and ground truth masks.

Method optimization. To demonstrate the effectiveness and better topological performances of our proposed method, we ensure identical prerequisites of our proposed method and the compared methods so that better performance can be clearly attributed to our proposed critical pixel mask. In that spirit, we conduct hyperparameter optimizations not only for our proposed method but also for the compared methods. For compound cIDice, we conducted a weight hyperparameter search for each dataset ranging from 0.1 to 0.5 in 0.1 steps. For space reasons, we only report the compound cIDice result with the highest cIDice metric in Tab. 1. For cIDice, we use the weight with the highest cIDice metric for each dataset from their paper. For CLD, we choose the same weight for cIDice as for the optimized compound cIDice.

Omitted Methods. We compare against cDice loss [2], a state-of-the-art topology-aware segmentation loss, which defines the skeleton of predictions and ground truth as critical pixel mask. cDice and [7] only differ in L_{pixel} . cDice and [3] only differ in the skeleton algorithm, having the same L_{pixel} and critical pixel mask. [4] is the multi-class adaption of cDice with a different L_{pixel} and a reduced critical pixel mask in comparison to cDice that only includes the slightly dilated ground truth skeleton for computational efficiency. Additionally, [4] removes small structures during training which makes it difficult to point out from which design-choice their reported improvement over cDice (no removed structures) originates. However, as [4] note themselves an inferior skeleton extraction for binary segmentation compared to cDice, we suspect their critical pixel mask selection to be inferior to cDice in binary segmentation as it doesn't consider false positive connections.

Our proposed method can be adapted with any skeleton extraction and L_{pixel} . Hence, we focus on validation against methods with a distinctly different critical pixel mask selection strategy. Therefore, we don't validate against [3, 7] and [4], which use different skeletonization methods or L_{pixel} than cDice but not superior critical pixel masks. CLoss outperforming compound cDice implies CLoss outperforming [3, 4, 7].

We would have liked to compare against [6] but weren't able to reproduce their results. We could only reproduce their proposed critical pixel mask for the 2D datasets but not on our 3D CLD. We suspect that this might be caused by the complex 3D surface structure in our CLD. Initial training attempts with the original code from their repository resulted in numerous errors for the 2D and 3D datasets, which we could only fix for the 3D datasets. However, related work [4, 7, 11, 12] also doesn't report results on [6], so the implementation of [6] might be a common issue.

There also exist numerous other topology-preserving approaches that we don't compare against, since this would go well beyond the scope of this paper. These other approaches are not based on critical pixel masks and include for example graph-based approaches [11] ([11] only works with 2D data) or approaches entirely based on post-processing [12] (omits to improve topology correctness already in the image domain). We note that [12] can be complemented with our proposed method.

A.6. Evaluation Metrics

Evaluation patch sizes. We evaluated the Dice metric respectively over the whole test samples for the used 2D and 3D datasets.

As mentioned in Sec. 2.2, e_0 is not able to capture gaps if it is computed over the whole image for datasets with high overall connectivity. We calculate the metrics on patch sizes which provide meaningful evaluation for the used datasets. We iterate over the whole image shape with the patches as sliding window (similar to [7]) instead of sampling patches randomly for evaluation like [2, 5, 6, 11].

For 2D datasets, we evaluate cDice and AGS on the whole image. We calculate betti-metrics and e_0 -Gt on patches, 375^2 (Roads) and 292^2 (HRF-Retina).

For the 3D dataset Vessap, we calculate cDice, AGS and e_0 -Gt over the whole sample volume. We calculate betti-metrics on full image size along the z-direction.

For our 3D CLD, we calculate cDice, AGS and e_0 -Gt over patches of $1024^2 \times 64$, where we also consider corrections of the cDice metric for its shortcoming with empty patches. We note that this patching in comparison to Vessap is due to memory reasons. Further, we evaluate betti-metrics on full image size along the z-direction.

A.7. Implementation Details

We conduct all our trainings within the nnU-Net framework [9] to ensure a maximum of reproducibility and leading performance [20]. The nnU-Net standard training has a length of 1000 epochs and utilizes Dice and Cross-Entropy as equally weighted compound loss. Training is done with a five-fold cross-validation for all methods. Inference is by standard done with all five folds simultaneously as an ensemble.

For our 2D datasets, we use the nnU-Net configuration 2d. For our 3D datasets, we use the configuration 3d_fullres.

We used the standard nnU-Net configuration for all methods, meaning our results can be reproduced with any GPU that has more than 11 GB VRAM. For increased speed, we conduct all trainings on A100 GPUs. As our used datasets are comparably small, we use the standard nnU-Net configuration without the residual encoder presets. We use PyTorch framework version 2.3.0 to implement our proposed method.

The nnU-Net configured a batch size of two for all datasets. The patch sizes for the 2D datasets were configured by nnU-Net to 1280×1024 (Roads), 1536×1024 (HRF-Retina) and for the 3D datasets to $256^2 \times 32$ (Vessap) and $160^2 \times 90$ (CLD).

A.8. Additional Discussion of Quantitative Results

HRF-Retina. CLoss has the smallest error on the number of connected components (e_0) in all datasets except for HRF-Retina. We attribute the value for HRF-Retina to the mentioned susceptibility of e_0 to artifacts (Sec. 2.2) from the patch-based evaluation. The artifact distortion of e_0 is indicated by e_0 -Gt, AGS, and the cIDice metric. Our proposed metrics e_0 -Gt and AGS clearly show a better gap closing of CLoss. Additionally, also the cIDice metric of CLoss for HRF-Retina is significantly improved (about 2%) over its closest competing method. The cIDice metric is especially reliable for HRF-Retina due to the smooth surface structure of the vessels. This indicates a better topology performance of CLoss for HRF-Retina despite the seemingly unfavorable e_0 value.

Vessap. For the Vessap dataset, we observe that the best cIDice score is achieved by the topology-insensitive nnU-Net (pretraining). This seems to contradict with the other topology-sensitive metrics. Therefore, we suspect that the thin network 3D structures of Vessap could lead to many seemingly false positive predictions of the prediction skeleton, which distorts the topological sensitivity of the cIDice metric on this dataset, see Fig. 4. This hypothesis is supported by an increasing AGS score for the topological fine-tuned methods. We also note that the seemingly high e_0 values for Vessap come from our patch-based evaluation of (full 2D image sizes stepping along the z-direction) in combination with the 3D network structure of Vessap. We note that e_1 in this context can especially contain artifacts from the slicing. Therefore, a minimized e_0 indicates topological performance on Vessap more reliable than minimized e_1 .

CLD. The unmodified cIDice loss only has a better gap closing (better e_0 -Gt and AGS) than our standard CLoss on CLD. Importantly, this is only due to a different pixel-wise loss and not due to the different critical pixel mask. We verify this by adding results of CLoss (Dice) to CLD, which has the same pixel-wise loss function as cIDice and only differs in the critical pixel mask to the cIDice loss. CLoss (Dice) visibly outperforms cIDice on the topology metrics. Interestingly, CLoss (Dice) is better than our proposed standard CLoss in e_1 , e_0 -Gt, and AGS but not in cIDice, e , and e_0 . This indicates that the target structure of our CLD favors different L_{pixel} for different metrics, which is not the case for the other datasets (compound cIDice consistently better or comparable to original cIDice). The favoring of different L_{pixel} for different metrics further illustrates the increased complexity of CLD over the other datasets.

A.9. Ablations

Additional results of the compound cIDice optimization are displayed in Tab. 2. We show the results on Roads as an

example, but conducted the optimization for all datasets.

Additional results of the post-processing are displayed in Tab. 3. All metric values in this example improve except for e_1 , which remains the same, and AGS, which slightly decreases. All methods benefit from our post-processing in a similar quantity. It can be observed for all results, that the change in the Dice value is comparably small to the change in e . This indicates, that our post-processing primarily removes noise in the form of small separated structures.

Ablations for the fine-tuning length are displayed in Tab. 4. A shorter fine-tuning length seems preferable, although we note that there is no consistent trend between the different lengths. For more detailed ablations on the fine-tuning between 50 epochs and 100 epochs, there was no clear winner, as might be falsely suggested by Tab. 4. Hence, we chose 50 epochs for computational efficiency.

Additional ablations for the critical pixel mask are included in Tab. 5. We include Thin CLoss, which has a critical pixel extraction analog to CLoss except for the context extraction. Therefore, Thin CLoss only considers the skeleton at the topological errors. Our proposed critical pixel mask with context extraction yields significantly better topological correctness. CLoss seems only second best in e_1 to compound cIDice ($\gamma = 0.1$), but the other bad metric values of compound cIDice ($\gamma = 0.1$) indicate that this is likely not due to overall topological correctness but rather related to artifacts.

An ablation on L_{pixel} is included in Tab. 6. cIDice and compound cIDice differ only in L_{pixel} . cIDice has L_{Dice} , and compound cIDice the equally weighted combination of L_{Dice} and L_{CE} , analogue to our CLoss implementation. Compound cIDice performs better on average for Roads, HRF-Retina, and Vessap. For CLD we note an advantage for cIDice.

Table 2: Compound cIDice Fine-Tuning. Example for Roads dataset. All results are post-processed.

| Method | Weight γ | Dice \uparrow | cIDice \uparrow [2] | e \downarrow | e $_1$ \downarrow | e $_0$ \downarrow | e $_0$ -Gt \downarrow | AGS \uparrow |
|--------------------------|-----------------|-----------------|-----------------------|----------------|---------------------|---------------------|-------------------------|----------------|
| nnU-Net [9] Dice & CE | | 79.69 | 89.34 | 1.181 | 0.895 | 0.286 | 0.702 | 86.46 |
| | | 79.77 | 89.37 | 1.156 | 0.949 | 0.207 | 0.699 | 86.41 |
| Compound cIDice | 0.1 | 79.69 | 89.26 | 1.089 | 0.866 | 0.223 | 0.691 | 86.78 |
| | 0.2 | 79.59 | 89.12 | 1.082 | 0.880 | 0.202 | 0.708 | 86.87 |
| | 0.3 | 79.73 | 89.34 | 1.108 | 0.909 | 0.199 | 0.673 | 87.05 |
| | 0.4 | 79.68 | 89.26 | 1.125 | 0.915 | 0.210 | 0.677 | 87.05 |
| | 0.5 | 79.59 | 89.21 | 1.065 | 0.870 | 0.195 | 0.673 | 87.25 |
| <i>C</i> Loss | 0.08 | 79.82 | 89.47 | 1.065 | 0.880 | 0.185 | 0.656 | 87.70 |
| | 0.1 | 79.57 | 89.21 | 0.990 | 0.810 | 0.180 | 0.617 | 87.88 |
| | 0.2 | 79.12 | 89.13 | 0.994 | 0.788 | 0.205 | 0.494 | 89.48 |

Table 3: Post-Processing. Results for CLD.

| Post-Processing | Method | Weight γ | Dice \uparrow | cIDice \uparrow [2] | e \downarrow | e $_1$ \downarrow | e $_0$ \downarrow | e $_0$ -Gt \downarrow | AGS \uparrow |
|-----------------|-----------------|-----------------|-----------------|-----------------------|----------------|---------------------|---------------------|-------------------------|----------------|
| w/o | Dice & CE | | 70.23 | 85.17 | 3.615 | 1.075 | 2.540 | 2.479 | 82.03 |
| with | | | 70.27 | 85.23 | 3.422 | 1.075 | 2.347 | 2.447 | 82.02 |
| w/o | Compound cIDice | 0.5 | 70.58 | 85.39 | 3.325 | 1.065 | 2.259 | 2.176 | 82.67 |
| with | | | | 70.61 | 85.45 | 3.165 | 1.065 | 2.099 | 2.150 |
| w/o | <i>C</i> Loss | 0.2 | 70.41 | 86.75 | 3.224 | 1.047 | 2.177 | 1.916 | 85.99 |
| with | | | | 70.44 | 86.83 | 3.020 | 1.047 | 1.973 | 1.899 |

Table 4: Fine-tuning length. Results for CLD. All results are post-processed.

| Epochs | Method | Weight γ | Dice \uparrow | cIDice \uparrow [2] | e \downarrow | e $_1$ \downarrow | e $_0$ \downarrow | e $_0$ -Gt \downarrow | AGS \uparrow | | |
|--------|-----------------|-----------------|-----------------|-----------------------|----------------|---------------------|---------------------|-------------------------|----------------|--------------|--------------|
| 50 | Compound cIDice | 0.5 | 70.61 | 85.45 | 3.165 | 1.065 | 2.099 | 2.150 | 82.66 | | |
| 100 | | | | | 70.67 | 85.37 | 3.108 | 1.062 | 2.046 | 2.086 | 82.79 |
| 150 | | | | | 70.58 | 85.22 | 3.097 | 1.064 | 2.034 | 2.180 | 82.19 |
| 300 | | | | | 70.65 | 85.44 | 3.115 | 1.053 | 2.062 | 2.097 | 82.71 |
| 50 | <i>C</i> Loss | 0.1 | 70.75 | 86.26 | 3.109 | 1.069 | 2.040 | 2.040 | 84.44 | | |
| 100 | | | | | 70.77 | 86.25 | 3.046 | 1.050 | 1.996 | 2.010 | 84.78 |
| 150 | | | | | 70.46 | 85.64 | 3.173 | 1.062 | 2.110 | 2.218 | 83.02 |
| 300 | | | | | 70.51 | 85.34 | 3.075 | 1.060 | 2.015 | 2.079 | 82.39 |

Table 5: Critical pixel mask. Results for CLD. All results are post-processed. All methods differ only in the critical pixel mask.

| Critical pixel mask | Method | Weight γ | Dice \uparrow | cIDice \uparrow [2] | e \downarrow | e $_1$ \downarrow | e $_0$ \downarrow | e $_0$ -Gt \downarrow | AGS \uparrow |
|--------------------------------|--------------------|-----------------|-----------------|-----------------------|----------------|---------------------|---------------------|-------------------------|----------------|
| w/o | Dice & CE | | 70.27 | 85.23 | 3.422 | 1.075 | 2.347 | 2.447 | 82.02 |
| Full skeleton | Compound cIDice | 0.1 | 68.93 | 83.11 | 3.522 | 1.029 | 2.493 | 2.679 | 79.06 |
| | | 0.2 | 70.47 | 84.96 | 3.361 | 1.077 | 2.283 | 2.352 | 81.81 |
| | | 0.3 | 70.61 | 85.41 | 3.230 | 1.077 | 2.153 | 2.249 | 82.80 |
| | | 0.4 | 70.51 | 85.16 | 3.287 | 1.063 | 2.223 | 2.283 | 82.26 |
| | | 0.5 | 70.61 | 85.45 | 3.165 | 1.065 | 2.099 | 2.150 | 82.66 |
| Skeleton at topological errors | <i>Thin C</i> Loss | 0.08 | 70.33 | 85.06 | 3.353 | 1.087 | 2.266 | 2.438 | 81.53 |
| | | 0.2 | 69.76 | 84.39 | 3.356 | 1.095 | 2.261 | 2.492 | 79.66 |
| | | 0.5 | 69.62 | 84.05 | 3.475 | 1.091 | 2.383 | 2.643 | 79.29 |
| Context at topological errors | <i>C</i> Loss | 0.08 | 70.98 | 86.22 | 3.158 | 1.063 | 2.095 | 2.061 | 84.47 |
| | | 0.1 | 70.75 | 86.26 | 3.109 | 1.069 | 2.040 | 2.040 | 84.44 |
| | | 0.2 | 70.44 | 86.83 | 3.020 | 1.047 | 1.973 | 1.899 | 85.98 |

Table 6: Pixel-wise loss. We focus on ablation between L_{Dice} and our L_{pixel} used for CLoss. All results are post-processed.

| Dataset | Method | Weight γ | Dice \uparrow | clDice \uparrow [2] | $e\downarrow$ | $e_1\downarrow$ | $e_0\downarrow$ | $e_0-Gt\downarrow$ | AGS \uparrow |
|------------|-----------------|-----------------|-----------------|-----------------------|---------------|-----------------|-----------------|--------------------|----------------|
| Roads | clDice [2] | 0.5 | 79.50 | 89.11 | 1.126 | 0.897 | 0.230 | 0.710 | 87.07 |
| | Compound clDice | 0.5 | 79.59 | 89.21 | 1.065 | 0.870 | 0.195 | 0.673 | 87.25 |
| HRF-Retina | clDice [2] | 0.5 | 82.15 | 83.09 | 0.426 | 0.256 | 0.170 | 2.475 | 82.23 |
| | Compound clDice | 0.5 | 82.33 | 82.98 | 0.405 | 0.250 | 0.155 | 2.429 | 81.25 |
| Vessap | clDice [2] | 0.4 | 92.70 | 94.86 | 29.00 | 1.240 | 27.76 | 9.44 | 97.74 |
| | Compound clDice | 0.4 | 92.98 | 95.42 | 26.880 | 1.220 | 25.660 | 11.040 | 97.28 |
| <i>CLD</i> | clDice [2] | 0.5 | 70.88 | 86.22 | 3.374 | 1.042 | 2.333 | 1.825 | 87.11 |
| | Compound clDice | 0.5 | 70.61 | 85.45 | 3.165 | 1.065 | 2.099 | 2.150 | 82.66 |