037

038

039

040

041

042

043

044

045

046

047

048

049

050

051

052

053

054

055

# InterConv: Set Interaction for Improved Biomedical Image Segmentation

Anonymous CVPR submission

Paper ID 13795

## Abstract

Image segmentation is an important part of many biomedical research and clinical pipelines. Because images within
a dataset are often similar in appearance and composition,
structures in one image can contain information that is useful for segmenting other images. However, existing image
segmentation models segment each input image independently, limiting their ability to share this information.

We present InterConv, a mechanism that enables seg-008 009 mentation models to interact and share information across a set of structurally related images. InterConv is a layer 010 011 that can be inserted within any network to facilitate set interaction with intermediate sample features without chang-012 ing the fundamental network architecture, and therefore can 013 014 be integrated into most existing segmentation models. We 015 demonstrate the effectiveness of InterConv by applying it to 016 two state-of-the-art image segmentation architectures: UN-017 ets and Vision Transformers, and tackle challenging tasks in both automatic and interactive biomedical image segmenta-018 019 tion. By learning to interact samples through aggregated set 020 features, InterConv consistently improves per-sample seg-021 mentation performance, sometimes by up to 19%.

# **022 1. Introduction**

Image segmentation is a core step in many biomedical im-023 age analysis pipelines. Machine learning-based models 024 025 have been shown to be useful for segmenting images of a wide range of modalities, spanning across diverse biomed-026 ical domains, Typically, these models segment one image 027 at a time. However, in many situations there are multiple 028 029 images acquired using the same modality and of the same 030 anatomy, resulting in similar structural compositions among samples. In some cases, the images have only minor dif-031 ferences, for example, longitudinal scans of the same sub-032 ject. In the same way that pixels can provide contextual in-033 034 formation about neighboring pixels, samples from a set of 035 related images can potentially provide useful information



Figure 1. InterConv helps models achieve substantially higher segmentation accuracy, as seen on the cortex (red). From left to right we show a close-up view of an input T2-FLAIR scan, with corresponding ground truth segmentation map, prediction with our proposed model, and the baseline prediction.

about each other.

We present a novel learning-based segmentation mechanism, InterConv, that enables a model to share information across the set by interacting a set of related input images and making prediction for each of the sample. Our contributions are:

- We introduce InterConv, a novel layer that is easily incorporated in most modern network designs to enable interaction among input images.
- We demonstrate the generalizability of InterConv by applying it to two different state-of-the-art segmentation architecture paradigms: convolutional UNets and Vision Transformers.
- We evaluate the effectiveness of InterConv in two segmentation scenarios: (1) *interactive* segmentation of diverse biomedical image datasets, (2) *automatic* segmentation of low-resolution MRI brain scans. We find that InterConv improves segmentation accuracy in both scenarios.

# 2. Related Work

Medical Image Segmentation.Image segmentation is056widely-studied in many biomedical domains.Most of to-057day's existing methods use a convolutional neural network058

097

098

099

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120



Figure 2. Examples of biomedical image segmentation tasks with spatial consistency. The ground truth segmentation is overlayed in blue [1, 45, 81, 107, 120].

059 to produce label maps for a single image [9, 23, 29, 48, 060 51, 54, 98]. When well-designed, UNet-like convolutional architectures [52, 73, 105] continue to perform similarly 061 to more recent and more complex architectures, based on 062 transformers [16, 37, 114] or state-space [79, 112] models. 063 All of these methods make predictions for a single input 064 065 image, independent of other images in the test set. In contrast, our work provides a mechanism to interact intermedi-066 ate sample features in all these networks, and improve seg-067 mentation accuracy on each of the samples. 068

069 Multimodal Image Inputs. Some methods integrate information from multiple modalities to improve biomedi-070 071 cal segmentation results [40, 57, 86]. These networks use 072 a fixed number of input channels [17, 27, 28], which restricts the amount of input information the network can in-073 tegrate [18, 20, 50, 51, 54, 91, 118]. Importantly, the multi-074 modal inputs are related to a single sample and the networks 075 make independent predictions on each sample. We propose 076 a mechanism that is agnostic to the number of inputs in the 077 set and produces a separate prediction for each set entry. 078

Multiple Predictions. Some segmentation methods gen-079 erate multiple predictions for a given image and aggregate 080 them, to improve accuracy or estimate uncertainty [51, 96, 081 082 103, 108, 110]. Some stochasatic segmentation methods are 083 designed to model the variability in multiple raters [6, 61,95, 96, 115]. While these methods produce, and sometimes 084 interact among, multiple segmentations, they all operate on 085 a single input image. 086

Data-driven (non-parametric) Mechanisms. Data-driven 087 (or non-parametric) methods yield improved performance 088 089 with higher amount of data available at inference. Recent deep-learning methods have combined parametric neural 090 networks with data-driven mechanisms, to draw on data 091 available at inference. For example, in-context learning 092 methods employ a neural network that adapts to new tasks 093 based on a context, often provided as a flexible set of ex-094 095 ample input-output pairs [12, 21, 96, 111]. In InterConv we also propose a non-parametric mechanism for jointly segmenting a set of multiple samples, and to enable interaction across the set. Moreover, InterConv can handle input sets of variables sizes.

**Interactive Segmentation.** Medical researchers and clinicians often need to perform *new* segmentation tasks involving new images, modalities, or labels that require them to perform manual segmentation. Interactive segmentation models reduce this burden by using sparse annotations, such as clicks and scribbles, as an additional model input. Recent foundation models for interactive segmentation yield promising results on many medical image segmentation tasks, but can perform poorly on tasks unseen during training, due to ambiguity in desired segmentation targets because of sparse annotations [60, 78, 113].

We demonstrate that with InterConv such ambiguities can be more easily resolved through set interaction where annotations in one image can be transferred across the set and help to segment other images. We find that this setting improves segmentation accuracy and requires substantially less total user effort than existing methods.

### 3. InterConv

Standard deep learning approaches to segmentation learn a function  $y = g_{\theta}(x)$ , with parameters  $\theta$ , that outputs a prediction y given a single input image x.

We propose InterConv, a mechanism that, when applied 121 to an existing network, enables it to model a function  $\mathcal{Y} =$ 122  $f_{\theta}(\mathcal{X})$  that jointly predicts a set of outputs  $\mathcal{Y} = \{y_i\}$  given 123 a set  $\mathcal{X} = \{x_i\}$  of input samples  $x_i$ , where each element  $y_i$ 124 in the output set  $\mathcal{Y}$  is the prediction corresponding to input 125 image  $x_i$ . The input set  $\mathcal{X}$  can be of flexible size. Input 126 samples interact through the InterConv Layer, so that infor-127 mation from each sample can be shared and contribute to all 128 predictions of the input set. 129

InterConv Layer. The InterConv Layer enables sample in-<br/>teraction across the input set by taking in a set of individual130131



Figure 3. Illustration of InterConv Layer. InterConv Layer takes a set of intermediate features  $\{u_i\}_{i=1}^{K}$  as input, aggregates them to produce the set context s, and fuses the set context representation with original sample features by concatenation and applying a feed-forward layer (FFW).

132sample features  $\mathcal{U}$  from a given network layer, and produc-133ing an updated set of features for each sample  $\mathcal{U}'$  after set in-134teraction. The InterConv Layer first summarizes interacted135information in a set context feature s through the set ag-136gregation step, and combines it with individual information137through the feature fusion step.

138 Set Aggregation. InterConv first takes a set of individ-139 ual image representations  $\mathcal{U} = {\mathbf{u_i}}$ , where each  $\mathbf{u_i}$  is an 140 intermediate feature representation that corresponds to in-141 put image  $x_i$  and has size  $F \times h \times w$ , with height h and 142 width w and F features. The set aggregation step computes 143 set context s as a pixel-wise average across samples for each 144 feature,

145 
$$\mathbf{s} = \frac{1}{K} \sum_{i=1}^{K} \mathbf{u}_i. \tag{1}$$

*Feature Fusion.* Given the set context s and imagespecific feature u<sub>i</sub>, InterConv first applies a channel-wise
concatenation:

$$\mathbf{c_i} = \texttt{Concat}(\mathbf{u_i}; \mathbf{s}), \tag{2}$$

then applies a shared convolution layer to each sample feature  $c_i$  in the set to integrate the interacted information with the individual sample features

$$\mathbf{u}_{\mathbf{i}}' = \operatorname{Conv}(\mathbf{c}_{\mathbf{i}}). \tag{3}$$

The final output features  $u'_i$  are the same size as the output from the original network layer, enabling InterConv to be used after any existing intermediate layer (Fig. 3).

**Training.** With InterConv, the network makes prediction on the entire set of images  $\mathcal{X}$ , and its weights are optimized jointly on the whole set:

$$\mathcal{L}_{\theta}(\mathcal{D}) = \mathbb{E}_{K}\left[\mathbb{E}_{\mathcal{X},\mathcal{Y}\subset\mathcal{D}}\left[\mathcal{L}\left(\hat{\mathcal{Y}},\mathcal{Y}\right)\right]\right]$$
(4) 160

$$= \mathbb{E}_{K} \left[ \mathbb{E}_{\mathcal{X}, \mathcal{Y} \subset \mathcal{D}} \left[ \sum_{i=1}^{K} \mathcal{L}_{seg} \left( \hat{y}_{i}, y_{i} \right) \right] \right]$$
(5) 161

where  $\mathcal{D}$  is a dataset containing sets of image samples  $\mathcal{X}$ and targets  $\mathcal{Y}$  with various set sizes K, and  $\hat{\mathcal{Y}}$  is the collection of network predictions on  $\mathcal{X}$ .  $\mathcal{L}$  is the desired segmentation loss over the input set, which is the sum of per-sample segmentation loss  $\mathcal{L}_{seg}$ .

**Inference.** Once trained, any network with the InterConv layer added can perform joint prediction on input image sets of any size.

### 4. Experiments

We assess the ability of InterConv to productively use information from multiple input images. To do so, we explore two common applications where it is challenging to construct accurate segmentation maps from single input images: 1) interactive segmentation and 2) automatic segmentation of the low-quality clinical-grade MRI scans that are prevalent in clinical practice.

**Baselines.** We demonstrate the effectiveness of InterConv integrating it into two of the most widely used image segmentation frameworks: convolutional UNets and Vision Transformers. Specifically, we apply InterConv to the following architectural designs:

- SimpleUNet, a UNet-like architecture for general image segmentation with a single convolution block at each encoder and decoder step,
- nnUNet [51], a UNet-based biomedical image segmentation network that automatically chooses some network aspects depending on the dataset [51, 52],
- SwinUNet [14], a transformer-based state-of-the-art medical image segmentation model, and
- ScribblePrompt-UNet [113], a UNet-based state-of-theart model for interactive segmentation.

Details about how we integrate InterConv into these architectures appear below.

**Evaluation.** We evaluate model predictions using Dice Score [26] averaged across foreground classes.

# 4.1. Interactive Segmentation

We first evaluate InterConv in the setting of interactive segmentation, focusing on segmentation tasks not seen during training.

Setup. Given an input image  $x_i$  and a collection of interac-201tions  $z_i$ , we predict binary segmentation  $y_i$ . With InterConv202layers, the model *jointly* segments a set of K images of sim-203

159

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

ilar anatomy  $\{x_i\}_{k=1}^K$ , given simulated prompts  $\{z_i\}_{k=1}^K$ , and produces a set of segmentation  $\{\hat{y}_i\}_{k=1}^K$  for each input.

**Interactions.** Following recent methods [113], we consider interactive segmentation using scribbles and clicks. We simulate positive prompts from label area  $y_i$  and negative prompts from  $1 - y_i$ , following the generation method of ScribblePrompt:

- Line scribbles: We randomly sample end points from the label and background areas, create a line segment that joins the end points, and warps the line segment to make sure that all points on the deformed path falls into one of the label or background categories.
- Random clicks: We randomly sample pixels that belong
  to the label or background area.

218 Network Architectures. We follow the architecture of 219 ScribblePrompt-UNet [113], a state-of-the-art interactive segmentation model with an efficient UNet-based architec-220 ture with 192 features at each convolution layer. We con-221 222 struct ScribblePrompt-InterConv by inserting a InterConv Layer with kernel size 3 at the end of each convolution layer 223 224 of the baseline network. We proportionally downscale fea-225 ture sizes in ScribblePrompt-InterConv to 137, so that both networks possess a similar number of parameters ( $\sim$ 3.6M). 226

Data. We use the datasets in MegaMedical [12, 96, 113]: 74
diverse biomedical image datasets, including over 48,000
images, 16 image types, and 602 labels.

We partition the datasets into 65 training datasets and 9 evaluation datasets using the same splits as [113]. Each evaluation dataset is partitioned into validation data, used for model selection, and test data, used for final evaluation. We report results on the test splits of the 9 evaluation datasets, unseen by the models during training [1, 2, 7, 36, 67, 94, 107, 119, 120].

We define a task as a combination of dataset, modality,
axis, and label. For multi-label datasets we consider each label as a separate binary segmentation task. For 3D datasets,
we use the middle slice and slice with maximum label area
to create 2D tasks.

Set Grouping. We sample sets containing different images
from the same segmentation task. This aligns with a realistic scenario, as users often need to segment several images
from the same dataset.

**Training.** *Training Objective.* We minimize Eq. 5 where  $\mathcal{L}_{seg}$  is a combination of Soft Dice Loss [87] and Focal Loss [71] using the Adam optimizer [58].

*Prompt generation.* We train all models with 2 prompt types: scribbles and clicks. For each prompt type we sample a random number of positive prompts ranging from 1 to 6, and a random number of negative prompts ranging from 0 to 6.

We train all models with fixed batch size of 8 sets,

with set sizes  $K \sim Cat[1,8]$  at each iteration. Following the training procedure in ScribblePrompt, we train both ScribblePrompt-UNet and ScribblePrompt-InterConv for 20,000 epochs with 125 batches per each epoch. 258

**Evaluation.** For evaluation, we use the same prompt types used during training. We evaluate model performance using Dice Score [26] on the 9 held-out datasets. We report results averaged across 5 runs with different (random) prompt initializations. At inference, we vary set size  $K \in \{1, 2, 4, 8\}$  and the number of prompt interactions  $n_{pos}, n_{neg} \in \{1, 2, 4, 8\}$ . We filter out tasks with fewer than 8 examples to make sure that scan images are distinct within a set, and that evaluations are performed on the same collection of test scans for each set size.

**Results.** Fig. 4 shows that for both scribbles and clicks, with all prompt sizes, ScibblePrompt- outperforms the baseline for set sizes greater than 1. These improvements are significant with a pairwise t-test (*p*-value ; 7e-4). The improvement in Dice grows larger with larger set sizes. Fig. 5 provides illustrative examples highlighting this improvement. More interactive segmentation visualizations can be found in the Appendix.

**Discussion.** Sample communication is especially useful in interactive segmentation because user-provided annotations on one image from the same task can provide information about the region being segmented on other images. Fig. 5 show that the baseline predictions, which only consider one sample at a time, are localized to the annotated pixels. By considering four annotated samples at once, ScribblePrompt-InterConv is able to gather more information about the target segmentation map from the whole set and propagate this information to each sample, achieving higher dice scores across all four predictions.

That the InterConv improvement over the baseline increases as the set size increases, demonstrates the value of shared information among samples. InterConv also achieves larger improvement with fewer prompts per image, and when there is less information per prompt (clicks as opposed to scribbles). In both situations, the less information there is per image, the more useful the information is from other images.

### 4.2. Automatic Segmentation: Clinical-quality MRI Scans

In this experiment, we evaluate InterConv for automatic segmentation of brain MRIs. MR images acquired in clinical settings often exhibit lower resolution, more noise, and lower tissue contrast than the research scans that are often used to evaluate segmentation methods. We evaluate Inter-Conv on the challenging task of segmenting clinical brain MRI.

**Data.** OASIS-3: OASIS-3 is a collection of 1,378 partic-

337

338

339

340

341

342



(a) Interactive segmentation on Pandental dataset with line scribble prompts. (b) Interactive segmentation on ACDC dataset with line scribble prompts.

Figure 4. ScribblePrompt-InterConv substantially improves interactive segmentation accuracy. We evaluate ScribblePrompt-UNet and comparable-sized ScribblePrompt-InterConv on unseen data. Each plot shows shows mean per-sample Dice score *improvement* of ScribblePrompt-InterConv over baseline on various set sizes and prompt sizes. Shaded regions show 95% CI from bootstrapping with 1,000 runs.

	SimpleUNet	SimpleUNet-InterConv (ours)	SwinUNet	SwinUNet-InterConv (ours)	nnUNet	nnUNet-InterConv (ours)
OASIS-3 T2★	$83.44 \pm 3.26$	$83.65 \pm 3.39$	$81.10 \pm 2.88$	$81.78 \pm 3.06$	$82.62 \pm 2.86$	$83.50 \pm 2.98$
OASIS-3 T2w-tse	$90.28 \pm 1.88$	$90.71 \pm 1.60$	$87.58 \pm 2.40$	$88.02 \pm 2.13$	$89.69 \pm 1.98$	$89.90 \pm 1.67$
ADNI FLAIR	$86.95 \pm 2.33$	$87.95 \pm 2.32$	$85.25 \pm 2.45$	$85.40 \pm 2.38$	$87.23 \pm 2.29$	$88.48 \pm 2.18$

Table 1. **Quantitative results: Dice performance on automatic segmentation.** We show the best performance for each model architecture along with standard deviation across test images.

ipants collected across several ongoing projects [66], con-taining 2,842 MR sessions with multiple modalities.

308 We focus on T2\* and T2-TSE MRI modalities, which have higher noise, lower contrast, and sparser slices than 309 research-quality T1 scans. Many subjects have multiple 310 scan sessions acquired over time. For each session, we ex-311 312 tract T1w scans, as well as T2-TSE and T2\* scans, if they exist. We use segmentations obtained from the paired high-313 314 quality T1w image as ground truth, and use InterConv to predict segmentation maps from the (aligned) low-quality 315 T2 scans. This enables us to assess the value of sharing in-316 formation in low-quality information, while having accurate 317 ground truth targets. 318

319 ADNI: We additionally use a collection of 5117 MR lon-320 gitudinal sessions gathered from 1061 participants in the 321 ADNI dataset. For each session, we extract a high-quality 322 T1w scan and T2-FLAIR scan. The FLAIR images are acquired with large slice separation (5 mm), and therefore ex-323 hibit low-resolution in the slice dimension. As above, we 324 obtain segmentations on the T1w image, and use it as the 325 326 ground truth label map for the aligned FLAIR scans.

**Processing.** We first perform skull-stripping using Synth-327 Strip [46]. Following standard practice in population anal-328 yses, we rigidly align all scans to a common space using 329 SynthMorph [41-44]. Specifically, we first align all modal-330 ities within a session to the T1 scan of that session. For each 331 subject, we then align the T1 scan of each session to the T1 332 scan of the *first* session. Finally, we align the T1 scan of the 333 first session to the Talairach space, and propagate all scans 334 to this space through the predicted transformations. 335

For both OASIS-3 and ADNI data, we use the middle coronal slice from all registered volumes. We filter and keep 20 anatomical regions that are commonly present in all scans. We exclude scans with missing labels. Details about class labels can be found in the Appendix.

We consider an input set to contain all scans from the same subject.

Network Architectures.We evaluate two architectural343designs:UNet-based and Transformer-based.We train344both SimpleUNets and SimpleUNet-InterConvs with vary-345345ing model sizes from  $\sim 60$ K to  $\sim 16$ M parameters, as well346as an nnUNet, which is automatically configured to have347

367

368

369

370

371



Figure 5. ScribblePrompt-InterConv significantly improves interactive segmentation quality with limited prompts. We run ScribblePrompt-InterConv and baseline predictions on 8 input images from Pandental, each provided with 2 positive (green) and 2 negative (red) line scribble prompts, and on 8 images from ACDC, providing each image with 8 positive (green) and 8 negative (red) random click prompts. The input samples come from different subjects within the same task. In each subfigure, **Top row** shows input image with provided prompts, along with ground truth segmentations, **middle row** shows ScribblePrompt-InterConv (ours) predictions, and **bottom row** shows ScribblePrompt (baseline) predictions overlayed in blue.

~20M parameters. We use the original SwinUNet design
with ~26M parameters as a baseline for Transformer-based
models.

We construct SimpleUNet-InterConv and nnUNet-351 InterConv by introducing a InterConv layer with kernel size 352 3 at the end of each convolution layer of SimpleUNets and 353 nnUNet respectively. We construct SwinUNet-InterConv 354 by adding a InterConv layer with kernel size 1 (i.e., point-355 wise feed-forward function) at the end of each multi-head 356 attention layer of SwinUNet. Details about baseline net-357 work sizes and their corresponding InterConv-integrated 358 network sizes can be found in the Appendix. 359

**Training.** *OASIS*: We train separate models for T2 $\star$  and T2-TSE. For both T2 $\star$  and T2wTSE, there are  $\sim$ 1030 subjects with the desired modality, which we first split into  $\sim 830$  362 subjects with  $\sim 1350$  images for model training, and  $\sim 100$  363 held out test subjects with  $\sim 190$  images. Details about train, validation and test data sizes can be found in the Appendix. 365

*ADNI*: With a total of 1061 subjects, we first split with 742 subjects for model training, and 213 subjects for performance evaluation.

We train all models with batches of 8 sets at each iteration, with a combination of Soft Dice Loss and Crossentropy Loss:

$$\mathcal{L}_{seg}(\hat{y}, y) = \mathcal{L}_{dice}(\hat{y}, y) + \mathcal{L}_{CCE}(\hat{y}, y), \tag{6}$$

We train each model configuration of SimpleUNet and SwinUNet with 3 different weight initialization seeds and with 374



Figure 6. **nnUNet-InterConv achieves higher segmentation accuracy on subjects with more scans available.** We show performance improvement on nnUNet-InterConv over baseline nnUNet by averaging per-sample Dice improvement over subjects with a given number of scans. Left figure shows Dice improvement on OASIS T2\* scans. Middle figure shows Dice improvement on OASIS T2\* scans. Shaded regions show 95% CI from bootstrapping with 1,000 runs.



Figure 7. **nnUNet-InterConv generates substantially better prediction in the brain cortex with set interaction mechanism.** We visualize predictions on 4 randomly selected scans from an input subject with 8 scan sessions. We interleave inputs and segmentation maps with a close-up of the temporal lobe. **Top row**: input images, **Second row**: ground truth segmentation maps, **Third row**: nnuNet-InterConv predictions, **Fourth row**: baseline nnUNet predictions.

an early-stopping criterion based on validation loss not decreasing for 40,000 iterations. We initialize all nnUNetstructured models in the same way as the original nnUNet,
and follow the original nnUNet training procedure to train

the models for a fixed 1000 epochs with 250 iterations per gooch. 379

**Evaluation.** We analyze model performance as the average 381

382 Dice score across all test scan samples. For network con383 figurations trained with multiple weight initializations, we
384 average the predicted probabilities from each trained model
385 before computing the Dice Score.

Results. Table 1 shows that all InterConv-integrated archi-386 387 tectures improve segmentation accuracy over corresponding baselines. Fig. 6 shows that subjects with more scans avail-388 able on average achieve better segmentation quality com-389 390 pared to those with fewer scans, demonstrating that larger sets are more informative. Fig. 7 shows that nnUNet-391 InterConv provides substantially better segmentations, e.g., 392 of the cortex, than the baseline nnUNet. We analyze dice 393 394 improvement and visualize results on the nnUNet framework since it consistently produces high segmentation ac-395 curacy and is a widely used baseline in the literature. Ad-396 ditional dice improvement analysis on other network archi-397 tectures can be found in the Appendix. 398

Discussion. We emphasize that if the ground truth label 399 400 maps are easily attainable from an individual image, we 401 do not expect shared information to be helpful. To be able 402 to even assess the value of shared information, we choose data where ground truth segmentation cannot be easily ob-403 tained from the scan itself, like the clinical-quality scans, 404 and use the aligned research-quality scans to extract the 405 ground truth maps. 406

The results in this section emphasize this effect: 407 408 while InterConv outperforms corresponding baselines in all datasets, its improvement in OASIS T2-TSE is minimal, 409 410 where image quality is least affected. In contrast, for the substantially more degraded T2\* and FLAIR modalities, 411 InterConv achieves substantial improvement over the base-412 413 lines, by providing each image with additional shared information through set interaction. 414

# 415 **5.** Conclusion

We introduce InterConv, a mechanism that improves segmentation accuracy by enabling an existing segmentation
model to benefit from jointly segmentating a set of related
images. InterConv leverages intermediate features of each
input sample to share information among the input set.

We demonstrate the value of this interaction by show-421 422 ing that InterConv-integrated networks outperforms baselines with the same architectural designs, in both automatic 423 and interactive segmentation settings. For interactive seg-424 mentation, InterConv can reduce the total amount of labor 425 426 required, since an annotation made on one image provides information that is useful for segmenting other images. For 427 automatic segmentation, since existing segmentation mod-428 els already perform well when the structures in an image 429 are clearly defined, we focus on domains where individual 430 431 images are challenging to segment, such as clinical-quality 432 scans. We find that in these common settings, InterConv

leads to improved segmentation quality.

The proposed InterConv Layer is easy to add to existing434segmentation architectures and promises to help improve435image segmentation in a wide variety of practical settings.436

# References

- Amir Hossein Abdi, Shohreh Kasaei, and Mojdeh Mehdizadeh. Automatic segmentation of mandible in panoramic x-ray. *Journal of Medical Imaging*, 2(4):044003, 2015. 2, 4, 1, 5
- [2] Walid Al-Dhabyani, Mohammed Gomaa, Hussien Khaled, and Aly Fahmy. Dataset of breast ultrasound images. *Data in Brief*, 28:104863, 2020. 4, 1, 5
- [3] Ujjwal Baid, Satyam Ghodasara, Suyash Mohan, Michel Bilello, Evan Calabrese, Errol Colak, Keyvan Farahani, Jayashree Kalpathy-Cramer, Felipe C Kitamura, Sarthak Pati, et al. The rsna-asnr-miccai brats 2021 benchmark on brain tumor segmentation and radiogenomic classification. *arXiv preprint arXiv:2107.02314*, 2021. 6
- [4] Spyridon Bakas, Hamed Akbari, Aristeidis Sotiras, Michel Bilello, Martin Rozycki, Justin S Kirby, John B Freymann, Keyvan Farahani, and Christos Davatzikos. Advancing the cancer genome atlas glioma mri collections with expert segmentation labels and radiomic features. *Scientific data*, 4 (1):1–13, 2017. 6
- [5] Sophia Bano, Francisco Vasconcelos, Luke M Shepherd, 457 Emmanuel Vander Poorten, Tom Vercauteren, Sebastien 458 Ourselin, Anna L David, Jan Deprest, and Danail Stoy-459 anov. Deep placental vessel segmentation for fetoscopic 460 mosaicking. In Medical Image Computing and Computer 461 Assisted Intervention-MICCAI 2020: 23rd International 462 Conference, Lima, Peru, October 4-8, 2020, Proceedings, 463 Part III 23, pages 763-773. Springer, 2020. 6 464
- [6] Christian F Baumgartner, Kerem C Tezcan, Krishna Chai-465 tanya, Andreas M Hötker, Urs J Muehlematter, Khoschy 466 Schawkat, Anton S Becker, Olivio Donati, and Ender 467 Konukoglu. Phiseg: Capturing uncertainty in medical im-468 age segmentation. In Medical Image Computing and Com-469 puter Assisted Intervention-MICCAI 2019: 22nd Interna-470 tional Conference, Shenzhen, China, October 13–17, 2019, 471 Proceedings, Part II 22, pages 119–127. Springer, 2019. 2 472
- [7] Olivier Bernard, Alain Lalande, Clement Zotti, Frederick Cervenansky, Xin Yang, Pheng-Ann Heng, Irem Cetin, Karim Lekadir, Oscar Camara, Miguel Angel Gonzalez Ballester, et al. Deep learning techniques for automatic mri cardiac multi-structures segmentation and diagnosis: is the problem solved? *IEEE transactions on medical imaging*, 37(11):2514–2525, 2018. 4, 1, 5
- [8] Patrick Bilic, Patrick Ferdinand Christ, Eugene Vorontsov, Grzegorz Chlebus, Hao Chen, Qi Dou, Chi-Wing Fu, Xiao Han, Pheng-Ann Heng, Jürgen Hesser, et al. The liver tumor segmentation benchmark (lits). arXiv preprint arXiv:1901.04056, 2019. 6
- [9] Benjamin Billot, Douglas N Greve, Oula Puonti, Axel
  Thielscher, Koen Van Leemput, Bruce Fischl, Adrian V
  Dalca, Juan Eugenio Iglesias, et al. Synthseg: Segmentation of brain mri scans of any contrast and resolution with-

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

473

474

475

476

477

478

479

480

481

482

483

484

490

529

530

531

532

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

out retraining. *Medical image analysis*, 86:102789, 2023. 2

- [10] Nicholas Bloch, Anant Madabhushi, Henkjan Huisman,
  John Freymann, Justin Kirby, Michael Grauer, Andinet
  Enquobahrie, Carl Jaffe, Larry Clarke, and Keyvan Farahani. Nci-isbi 2013 challenge: automated segmentation of
  prostate structures. *The Cancer Imaging Archive*, 370(6):5,
  2015. 6
- 497 [11] Mateusz Buda, Ashirbani Saha, and Maciej A Mazurowski.
  498 Association of genomic subtypes of lower-grade gliomas
  499 with shape features automatically extracted by a deep learning algorithm. *Computers in biology and medicine*, 109:
  501 218–225, 2019. 6
- 502 [12] Victor Ion Butoi\*, Jose Javier Gonzalez Ortiz\*, Tianyu Ma,
  503 Mert R. Sabuncu, John Guttag, and Adrian V. Dalca. Uni504 verseg: Universal medical image segmentation. *Interna-*505 *tional Conference on Computer Vision*, 2023. 2, 4, 1
- 506 [13] Juan C. Caicedo, Allen Goodman, Kyle W. Karhohs,
  507 Beth A. Cimini, Jeanelle Ackerman, Marzieh Haghighi,
  508 CherKeng Heng, Tim Becker, Minh Doan, Claire McQuin,
  509 Mohammad Rohban, Shantanu Singh, and Anne E. Carpen510 ter. Nucleus segmentation across imaging experiments: the
  511 2018 Data Science Bowl. *Nature Methods*, 16(12):1247–
  512 1253, 2019. 6
- [14] Hu Cao, Yueyue Wang, Joy Chen, Dongsheng Jiang, Xi-aopeng Zhang, Qi Tian, and Manning Wang. Swin-unet:
  Unet-like pure transformer for medical image segmenta-tion. In *Proceedings of the European Conference on Computer Vision Workshops(ECCVW)*, 2022. 3
- [15] Albert Cardona, Stephan Saalfeld, Stephan Preibisch, Benjamin Schmid, Anchi Cheng, Jim Pulokas, Pavel Tomancak, and Volker Hartenstein. An integrated micro-and macroarchitectural analysis of the drosophila brain by computer-assisted serial section electron microscopy. *PLoS biology*, 8(10):e1000502, 2010. 6
- [16] Jieneng Chen, Yongyi Lu, Qihang Yu, Xiangde Luo, Ehsan
  Adeli, Yan Wang, Le Lu, Alan L Yuille, and Yuyin Zhou.
  Transunet: Transformers make strong encoders for medical image segmentation. *arXiv preprint arXiv:2102.04306*,
  2021. 2
  - [17] Lele Chen, Yue Wu, Adora M DSouza, Anas Z Abidin, Axel Wismüller, and Chenliang Xu. Mri tumor segmentation with densely connected 3d cnn. In *Medical imaging* 2018: image processing, pages 357–364. SPIE, 2018. 2
- [18] Albert Clèrigues, Sergi Valverde, Jose Bernal, Jordi Freixenet, Arnau Oliver, and Xavier Lladó. Acute and sub-acute stroke lesion segmentation from multimodal mri. *Computer methods and programs in biomedicine*, 194:105521, 2020.
  2
- 538 [19] Noel C. F. Codella, David A. Gutman, M. Emre Celebi, 539 Brian Helba, Michael A. Marchetti, Stephen W. Dusza, 540 Aadi Kalloo, Konstantinos Liopyris, Nabin K. Mishra, Har-541 ald Kittler, and Allan Halpern. Skin lesion analysis to-542 ward melanoma detection: A challenge at the 2017 international symposium on biomedical imaging (isbi), hosted by 543 544 the international skin imaging collaboration (ISIC). CoRR, 545 abs/1710.05006, 2017. 6

- [20] Shaoguo Cui, Lei Mao, Jingfeng Jiang, Chang Liu, Shuyu Xiong, et al. Automatic semantic segmentation of brain gliomas from mri images using a deep cascaded neural network. *Journal of healthcare engineering*, 2018, 2018. 2
- [21] Steffen Czolbe and Adrian V Dalca. Neuralizer: General neuroimage analysis without re-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6217–6230, 2023. 2
- [22] Adrian V Dalca, John Guttag, and Mert R Sabuncu. Anatomical priors in convolutional networks for unsupervised biomedical segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9290–9299, 2018. 6
- [23] Adrian V Dalca, Evan Yu, Polina Golland, Bruce Fischl, Mert R Sabuncu, and Juan Eugenio Iglesias. Unsupervised deep learning for bayesian brain mri segmentation. In Medical Image Computing and Computer Assisted Intervention– MICCAI 2019: 22nd International Conference, Shenzhen, China, October 13–17, 2019, Proceedings, Part III 22, pages 356–365. Springer, 2019. 2
- [24] Etienne Decenciere, Guy Cazuguel, Xiwei Zhang, Guillaume Thibault, J-C Klein, Fernand Meyer, Beatriz Marcotegui, Gwénolé Quellec, Mathieu Lamard, Ronan Danno, et al. Teleophta: Machine learning and image processing methods for teleophthalmology. *Irbm*, 34(2):196–203, 2013. 6
- [25] Aysen Degerli, Morteza Zabihi, Serkan Kiranyaz, Tahir Hamid, Rashid Mazhar, Ridha Hamila, and Moncef Gabbouj. Early detection of myocardial infarction in lowquality echocardiography. *IEEE Access*, 9:34442–34453, 2021. 6
- [26] Lee R Dice. Measures of the amount of ecologic association between species. *Ecology*, 26(3):297–302, 1945. 3, 4
- [27] Jose Dolz, Christian Desrosiers, and Ismail Ben Ayed. Ivdnet: Intervertebral disc localization and segmentation in mri with a multi-modal unet. In *International workshop and challenge on computational methods and clinical applications for spine imaging*, pages 130–143. Springer, 2018. 2
- [28] Jose Dolz, Karthik Gopinath, Jing Yuan, Herve Lombaert, Christian Desrosiers, and Ismail Ben Ayed. Hyperdensenet: a hyper-densely connected cnn for multi-modal image segmentation. *IEEE transactions on medical imaging*, 38 (5):1116–1126, 2018. 2
- [29] Jose Dolz, Karthik Gopinath, Jing Yuan, Herve Lombaert, Christian Desrosiers, and Ismail Ben Ayed. Hyperdensenet: A hyper-densely connected cnn for multi-modal image segmentation. *IEEE Transactions on Medical Imaging*, 38 (5):1116–1126, 2019. 2
- [30] J Gamper, NA Koohbanani, K Benes, S Graham, M Jahanifar, SA Khurram, A Azam, K Hewitt, and N Rajpoot. Pannuke dataset extension, insights and baselines. arxiv. 2020 doi: 10.48550. ARXIV, 2003. 6
- [31] Stephan Gerhard, Jan Funke, Julien Martel, Albert Cardona, and Richard Fetter. Segmented anisotropic ssTEM dataset of neural tissue. 2013. 6
- [32] Randy L Gollub, Jody M Shoemaker, Margaret D King, Tonya White, Stefan Ehrlich, Scott R Sponheim, Vincent P
   602

625

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

603Clark, Jessica A Turner, Bryon A Mueller, Vince Magnotta,604et al. The mcic collection: a shared repository of multi-605modal, multi-site brain image data from a clinical inves-606tigation of schizophrenia. Neuroinformatics, 11:367–388,6072013. 6

- [33] Ioannis S Gousias, Daniel Rueckert, Rolf A Heckemann,
  Leigh E Dyet, James P Boardman, A David Edwards, and
  Alexander Hammers. Automatic segmentation of brain
  mris of 2-year-olds into 83 regions of interest. *Neuroim- age*, 40(2):672–684, 2008. 6
- [34] Ioannis S Gousias, A David Edwards, Mary A Rutherford, Serena J Counsell, Jo V Hajnal, Daniel Rueckert, and
  Alexander Hammers. Magnetic resonance imaging of the
  newborn brain: manual segmentation of labelled atlases in
  term-born and preterm infants. *Neuroimage*, 62(3):1499–
  1509, 2012. 6
- 619 [35] Endre Grøvik, Darvin Yi, Michael Iv, Elizabeth Tong,
  620 Daniel Rubin, and Greg Zaharchuk. Deep learning enables
  621 automatic detection and segmentation of brain metastases
  622 on multisequence mri. *Journal of Magnetic Resonance*623 *Imaging*, 51(1):175–182, 2020. 6
  - [36] Daniel Gut. X-ray images of the hip joints. 1, 2021. Publisher: Mendeley Data. 4, 1, 5
- [37] Ali Hatamizadeh, Yucheng Tang, Vishwesh Nath, Dong
  Yang, Andriy Myronenko, Bennett Landman, Holger R
  Roth, and Daguang Xu. Unetr: Transformers for 3d medical
  image segmentation. In *Proceedings of the IEEE/CVF win- ter conference on applications of computer vision*, pages
  574–584, 2022. 2
- [38] Nicholas Heller, Fabian Isensee, Klaus H Maier-Hein, Xiaoshuai Hou, Chunmei Xie, Fengyi Li, Yang Nan, Guangrui Mu, Zhiyong Lin, Miofei Han, et al. The state of the art in kidney and kidney tumor segmentation in contrastenhanced ct imaging: Results of the kits19 challenge. *Medical Image Analysis*, page 101821, 2020. 6
- 638 [39] Moritz R Hernandez Petzsche, Ezequiel de la Rosa, Uta
  639 Hanning, Roland Wiest, Waldo Valenzuela, Mauricio
  640 Reyes, Maria Meyer, Sook-Lei Liew, Florian Kofler, Ivan
  641 Ezhov, et al. Isles 2022: A multi-center magnetic resonance
  642 imaging stroke lesion segmentation dataset. *Scientific data*,
  643 9(1):762, 2022. 6
- 644 [40] Moritz R Hernandez Petzsche, Ezequiel de la Rosa, Uta
  645 Hanning, Roland Wiest, Waldo Valenzuela, Mauricio
  646 Reyes, Maria Meyer, Sook-Lei Liew, Florian Kofler, Ivan
  647 Ezhov, et al. Isles 2022: A multi-center magnetic resonance
  648 imaging stroke lesion segmentation dataset. *Scientific data*,
  649 9(1):762, 2022. 2
- [41] Malte Hoffmann, Benjamin Billot, Juan E Iglesias, Bruce
  Fischl, and Adrian V Dalca. Learning mri contrast-agnostic
  registration. In 2021 IEEE 18th International Symposium
  on Biomedical Imaging (ISBI), pages 899–903. IEEE, 2021.
  5
- [42] Malte Hoffmann, Benjamin Billot, Douglas N Greve,
  Juan Eugenio Iglesias, Bruce Fischl, and Adrian V Dalca.
  Synthmorph: learning contrast-invariant registration without acquired images. *IEEE Transactions on Medical Imag- ing*, 41(3):543–558, 2022.

- [43] Malte Hoffmann, Andrew Hoopes, Bruce Fischl, and Adrian V Dalca. Anatomy-specific acquisition-agnostic affine registration learned from fictitious images. In *Medical Imaging 2023: Image Processing*, page 1246402. SPIE, 2023.
- [44] Malte Hoffmann, Andrew Hoopes, Douglas N Greve, Bruce Fischl, and Adrian V Dalca. Anatomy-aware and acquisition-agnostic joint registration with synthmorph. *arXiv preprint arXiv:2301.11329*, 2023. 5
- [45] Andrew Hoopes, Malte Hoffmann, Douglas N. Greve, Bruce Fischl, John Guttag, and Adrian V. Dalca. Learning the effect of registration hyperparameters with hypermorph. *Machine Learning for Biomedical Imaging*, 1:1–30, 2022. 2, 6
- [46] Andrew Hoopes, Jocelyn S Mora, Adrian V Dalca, Bruce Fischl, and Malte Hoffmann. Synthstrip: skull-stripping for any brain image. *NeuroImage*, 260:119474, 2022. 5
- [47] AD Hoover, Valentina Kouznetsova, and Michael Goldbaum. Locating blood vessels in retinal images by piecewise threshold probing of a matched filter response. *IEEE Transactions on Medical imaging*, 19(3):203–210, 2000. 6
- [48] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q. Weinberger. Densely connected convolutional networks. In 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 2261–2269, 2017. 2
- [49] Humans in the Loop. Teeth segmentation dataset. 6
- [50] Fabian Isensee, Philipp Kickingereder, Wolfgang Wick, Martin Bendszus, and Klaus H Maier-Hein. Brain tumor segmentation and radiomics survival prediction: Contribution to the brats 2017 challenge. In *Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries: Third International Workshop, BrainLes 2017, Held in Conjunction with MICCAI 2017, Quebec City, QC, Canada, September 14, 2017, Revised Selected Papers 3,* pages 287–297. Springer, 2018. 2
- [51] Fabian Isensee, Paul F Jaeger, Simon AA Kohl, Jens Petersen, and Klaus H Maier-Hein. nnu-net: a selfconfiguring method for deep learning-based biomedical image segmentation. *Nature methods*, 18(2):203–211, 2021. 2, 3
- [52] Fabian Isensee, Tassilo Wald, Constantin Ulrich, Michael Baumgartner, Saikat Roy, Klaus Maier-Hein, and Paul F Jaeger. nnu-net revisited: A call for rigorous validation in 3d medical image segmentation. arXiv preprint arXiv:2404.09556, 2024. 2, 3
- [53] Yuanfeng Ji, Haotian Bai, Jie Yang, Chongjian Ge, Ye Zhu, Ruimao Zhang, Zhen Li, Lingyan Zhang, Wanling Ma, Xiang Wan, et al. Amos: A large-scale abdominal multi-organ benchmark for versatile medical image segmentation. *arXiv* preprint arXiv:2206.08023, 2022. 6
- [54] Konstantinos Kamnitsas, Christian Ledig, Virginia FJ Newcombe, Joanna P Simpson, Andrew D Kane, David K
  Menon, Daniel Rueckert, and Ben Glocker. Efficient multiscale 3d cnn with fully connected crf for accurate brain lesion segmentation. *Medical image analysis*, 36:61–78, 2017. 2
  710
  711
  712
  713
  714
  715
- [55] Rashed Karim, R James Housden, Mayuragoban Balasubramaniam, Zhong Chen, Daniel Perry, Ayesha Uddin, Yosra
   717

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

Al-Beyatti, Ebrahim Palkhi, Prince Acheampong, Samantha Obom, et al. Evaluation of current algorithms for segmentation of scar tissue from late gadolinium enhancement
cardiovascular magnetic resonance of the left atrium: an
open-access grand challenge. *Journal of Cardiovascular Magnetic Resonance*, 15(1):1–17, 2013. 6

- 724 [56] Ali Emre Kavur, M. Alper Selver, Oğuz Dicle, Mustafa
  725 Barış, and N. Sinem Gezer. CHAOS Combined (CT-MR)
  726 Healthy Abdominal Organ Segmentation Challenge Data,
  727 2019. 6
- 728 [57] A. Emre Kavur, N. Sinem Gezer, Mustafa Barış, Sinem Aslan, Pierre-Henri Conze, Vladimir Groza, Duc Duy 729 730 Pham, Soumick Chatterjee, Philipp Ernst, Savaş Özkan, Bora Baydar, Dmitry Lachinov, Shuo Han, Josef Pauli, 731 732 Fabian Isensee, Matthias Perkonigg, Rachana Sathish, Ron-733 nie Rajan, Debdoot Sheet, Gurbandurdy Dovletov, Oliver 734 Speck, Andreas Nürnberger, Klaus H. Maier-Hein, Gözde 735 Bozdağı Akar, Gözde Ünal, Oğuz Dicle, and M. Alper 736 Selver. CHAOS Challenge - combined (CT-MR) healthy 737 abdominal organ segmentation. Medical Image Analysis, 738 69:101950, 2021. 2, 6
- [58] Diederik P Kingma and Jimmy Ba. Adam: A method for
  stochastic optimization. *arXiv preprint arXiv:1412.6980*,
  2014. 4
- [59] Serkan Kiranyaz, Aysen Degerli, Tahir Hamid, Rashid
  Mazhar, Rayyan El Fadil Ahmed, Rayaan Abouhasera,
  Morteza Zabihi, Junaid Malik, Ridha Hamila, and Moncef
  Gabbouj. Left ventricular wall motion estimation by active polynomials for acute myocardial infarction detection. *IEEE Access*, 8:210301–210317, 2020. 6
- 748 [60] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi
  749 Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer
  750 Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollar,
  751 and Ross Girshick. Segment anything. In *Proceedings of*752 *the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4015–4026, 2023. 2
- [61] Simon Kohl, Bernardino Romera-Paredes, Clemens Meyer,
  Jeffrey De Fauw, Joseph R Ledsam, Klaus Maier-Hein, SM
  Eslami, Danilo Jimenez Rezende, and Olaf Ronneberger.
  A probabilistic u-net for segmentation of ambiguous images. Advances in neural information processing systems,
  31, 2018. 2
- [62] Markus Krönke, Christine Eilers, Desislava Dimova,
  Melanie Köhler, Gabriel Buschner, Lilit Schweiger, Lemonia Konstantinidou, Marcus Makowski, James Nagarajah,
  Nassir Navab, et al. Tracked 3d ultrasound and deep neural
  network-based thyroid segmentation reduce interobserver
  variability in thyroid volumetry. *Plos one*, 17(7):e0268550,
  2022. 6
- 767 [63] Hugo J Kuijf, J Matthijs Biesbroek, Jeroen De Bresser,
  768 Rutger Heinen, Simon Andermatt, Mariana Bento, Matt
  769 Berseth, Mikhail Belyaev, M Jorge Cardoso, Adria
  770 Casamitjana, et al. Standardized assessment of automatic
  771 segmentation of white matter hyperintensities and results
  772 of the wmh segmentation challenge. *IEEE transactions on*773 *medical imaging*, 38(11):2556–2568, 2019. 6
- [64] Maria Kuklisova-Murgasova, Paul Aljabar, Latha Srini vasan, Serena J Counsell, Valentina Doria, Ahmed Serag,

Ioannis S Gousias, James P Boardman, Mary A Rutherford,776A David Edwards, et al. A dynamic 4d probabilistic atlas of777the developing brain. NeuroImage, 54(4):2750–2763, 2011.7786779

- [65] Zoé Lambert, Caroline Petitjean, Bernard Dubray, and Su Kuan. Segthor: segmentation of thoracic organs at risk in ct images. In 2020 Tenth International Conference on Image Processing Theory, Tools and Applications (IPTA), pages 1–6. IEEE, 2020. 6
  780
  781
  782
  783
  784
- [66] PJ LaMontague, TLS Benzinger, John C Morris, S Keefe, R
  Hornbeck, C Xiong, E Grant, J Hassenstab, K Moulder, AG
  Vlassenko, et al. Oasis-3: Longitudinal neuroimaging, clinical, and cognitive dataset for normal aging and alzheimer
  disease. *medRxiv*, page 2019, 2019. 5
  789
- [67] Bennett Landman, Zhoubing Xu, J Igelsias, Martin Styner, T Langerak, and Arno Klein. Miccai multi-atlas labeling beyond the cranial vault–workshop and challenge. In Proc. MICCAI Multi-Atlas Labeling Beyond Cranial Vault Workshop Challenge, page 12, 2015. 4, 1, 5, 6
- [68] Sarah Leclerc, Erik Smistad, Joao Pedrosa, Andreas Østvik, Frederic Cervenansky, Florian Espinosa, Torvald Espeland, Erik Andreas Rye Berg, Pierre-Marc Jodoin, Thomas Grenier, et al. Deep learning for segmentation using an open large-scale dataset in 2d echocardiography. *IEEE transactions on medical imaging*, 38(9):2198–2210, 2019. 6
- [69] Guillaume Lemaître, Robert Martí, Jordi Freixenet, Joan C Vilanova, Paul M Walker, and Fabrice Meriaudeau. Computer-aided detection and diagnosis for prostate cancer based on mono and multi-parametric mri: a review. *Computers in biology and medicine*, 60:8–31, 2015. 6
- [70] Mingchao Li, Yuhan Zhang, Zexuan Ji, Keren Xie, Songtao Yuan, Qinghuai Liu, and Qiang Chen. Ipn-v2 and octa-500: Methodology and dataset for retinal image segmentation. *arXiv preprint arXiv:2012.07261*, 2020. 6
- [71] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In Proceedings of the IEEE international conference on computer vision, pages 2980–2988, 2017. 4
- [72] Geert Litjens, Robert Toth, Wendy van de Ven, Caroline Hoeks, Sjoerd Kerkstra, Bram van Ginneken, Graham Vincent, Gwenael Guillard, Neil Birbeck, Jindang Zhang, et al. Evaluation of prostate segmentation algorithms for mri: the promise12 challenge. *Medical image analysis*, 18(2):359– 373, 2014. 6
- [73] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *Proceedings of the IEEE/CVF conference* on computer vision and pattern recognition, pages 11976– 11986, 2022. 2
- [74] Vebjorn Ljosa, Katherine L Sokolnicki, and Anne E Carpenter. Annotated high-throughput microscopy image sets for validation. *Nature methods*, 9(7):637–637, 2012. 6
- [75] Maximilian T Löffler, Anjany Sekuboyina, Alina Jacob, Anna-Lena Grau, Andreas Scharr, Malek El Husseini, Mareike Kallweit, Claus Zimmer, Thomas Baum, and Jan S Kirschke. A vertebral segmentation dataset with fracture grading. *Radiology: Artificial Intelligence*, 2(4):e190138, 2020. 6

892

893

894

895

896

897

898

899

900

901

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

- [76] Xiangde Luo, Wenjun Liao, Jianghong Xiao, Tao Song, Xiaofan Zhang, Kang Li, Guotai Wang, and Shaoting Zhang.
  Word: Revisiting organs segmentation in the whole abdominal region. *arXiv preprint arXiv:2111.02403*, 2021. 6
- 838 [77] Jun Ma, Yao Zhang, Song Gu, Xingle An, Zhihe Wang,
  839 Cheng Ge, Congcong Wang, Fan Zhang, Yu Wang, Yinan
  840 Xu, et al. Fast and low-gpu-memory abdomen ct organ seg841 mentation: the flare challenge. *Medical Image Analysis*, 82:
  842 102616, 2022. 6
- [78] Jun Ma, Yuting He, Feifei Li, Lin Han, Chenyu You, and
  Bo Wang. Segment anything in medical images. *Nature Communications*, 15(1):654, 2024. 2
- [79] Jun Ma, Feifei Li, and Bo Wang. U-mamba: Enhancing long-range dependency for biomedical image segmentation. *arXiv preprint arXiv:2401.04722*, 2024. 2
- [80] Yuhui Ma, Huaying Hao, Jianyang Xie, Huazhu Fu, Jiong
  Zhang, Jianlong Yang, Zhen Wang, Jiang Liu, Yalin Zheng,
  and Yitian Zhao. Rose: a retinal oct-angiography vessel
  segmentation dataset and new model. *IEEE Transactions on Medical Imaging*, 40(3):928–939, 2021. 6
- [81] Daniel S Marcus, Tracy H Wang, Jamie Parker, John G
  Csernansky, John C Morris, and Randy L Buckner. Open
  access series of imaging studies (oasis): cross-sectional mri
  data in young, middle aged, nondemented, and demented
  older adults. *Journal of cognitive neuroscience*, 19(9):
  1498–1507, 2007. 2, 6
- [82] Kenneth Marek, Danna Jennings, Shirley Lasch, Andrew
  Siderowf, Caroline Tanner, Tanya Simuni, Chris Coffey,
  Karl Kieburtz, Emily Flagg, Sohini Chowdhury, et al. The
  parkinson progression marker initiative (ppmi). *Progress in neurobiology*, 95(4):629–635, 2011. 6
- [83] Francesco Marzola, Nens Van Alfen, Jonne Doorduin, and
  Kristen M. Meiburger. Deep learning segmentation of
  transverse musculoskeletal ultrasound images for neuromuscular disease assessment. *Computers in Biology and Medicine*, 135:104623, 2021. 6
- [84] Maciej A Mazurowski, Kal Clark, Nicholas M Czarnek,
  Parisa Shamsesfandabadi, Katherine B Peters, and Ashirbani Saha. Radiogenomics of lower-grade glioma:
  algorithmically-assessed tumor shape is associated with tumor genomic subtypes and patient outcomes in a multiinstitutional study with the cancer genome atlas data. *Journal of neuro-oncology*, 133:27–35, 2017. 6
- 877 [85] Bjoern Menze, Leo Joskowicz, Spyridon Bakas, Andras
  878 Jakab, Ender Konukoglu, Anton Becker, Amber Simpson,
  879 and Richard D. Quantification of uncertainties in biomed880 ical image quantification 2021. *4th International Confer-*881 *ence on Medical Image Computing and Computer Assisted*882 *Intervention (MICCAI 2021)*, 2021. 6
- [86] Bjoern H Menze, Andras Jakab, Stefan Bauer, Jayashree
  Kalpathy-Cramer, Keyvan Farahani, Justin Kirby, Yuliya
  Burren, Nicole Porz, Johannes Slotboom, Roland Wiest,
  et al. The multimodal brain tumor image segmentation
  benchmark (brats). *IEEE transactions on medical imaging*,
  34(10):1993–2024, 2014. 2, 6
- [87] Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi.
   V-net: Fully convolutional neural networks for volumetric

medical image segmentation. In 2016 fourth international conference on 3D vision (3DV), pages 565–571. Ieee, 2016.

- [88] Anna Montoya, Hasnin, kaggle446, shirzad, Will Cukierski, and yffud. Ultrasound nerve segmentation, 2016. 6
- [89] Kelly Payette, Priscille de Dumast, Hamza Kebiri, Ivan Ezhov, Johannes C Paetzold, Suprosanna Shit, Asim Iqbal, Romesa Khan, Raimund Kottke, Patrice Grehten, et al. An automatic multi-tissue human fetal brain segmentation benchmark using the fetal tissue annotation dataset. *Scientific Data*, 8(1):1–14, 2021. 6
- [90] Lina Pedraza, Carlos Vargas, Fabián Narváez, Oscar Durán, Emma Muñoz, and Eduardo Romero. An open access thyroid ultrasound image database. In *10th international symposium on medical information processing and analysis*, page 92870W. SPIE / International Society for Optics and Photonics, 2015. 6
  902
  903
  904
  904
  905
  906
  907
- [91] Sérgio Pereira, Adriano Pinto, Victor Alves, and Carlos A Silva. Brain tumor segmentation using convolutional neural networks in mri images. *IEEE transactions on medical imaging*, 35(5):1240–1251, 2016. 2
  911
- [92] Gašper Podobnik, Primož Strojan, Primož Peterlin, Bulat
  Ibragimov, and Tomaž Vrtovec. HaN-Seg: The head
  and neck organ-at-risk CT and MR segmentation dataset. *Medical Physics*, 50(3):1917–1927, 2023. tex.eprint:
  https://aapm.onlinelibrary.wiley.com/doi/pdf/10.1002/mp.16197916
  917
- [93] Prasanna Porwal, Samiksha Pachade, Ravi Kamble, Manesh Kokare, Girish Deshmukh, Vivek Sahasrabuddhe, and Fabrice Meriaudeau. Indian diabetic retinopathy image dataset (idrid), 2018. 6
  918
  919
  920
  921
- [94] Perry Radau, Yingli Lu, Kim Connelly, Gideon Paul, AJWG Dick, and Graham Wright. Evaluation framework for algorithms segmenting short axis cardiac mri. *The MIDAS Journal-Cardiac MR Left Ventricle Segmentation Challenge*, 49, 2009. 4, 1, 5
- [95] Aimon Rahman, Jeya Maria Jose Valanarasu, Ilker Hacihaliloglu, and Vishal M Patel. Ambiguous medical image segmentation using diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11536–11546, 2023. 2
- [96] Marianne Rakic, Hallee E. Wong, Jose Javier Gonzalez Ortiz, Beth Cimini, John V. Guttag, and Adrian V. Dalca. Tyche: Stochastic in-context learning for medical image segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2024. 2, 4, 1
- [97] Blaine Rister, Darvin Yi, Kaushik Shivakumar, Tomomi Nobashi, and Daniel L. Rubin. CT-ORG, a new dataset for multiple organ segmentation in computed tomography. *Scientific Data*, 7(1):381, 2020. 6
- [98] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. Unet: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241.
  944 Springer, 2015. 2
  945
- [99] Adriel Saporta, Xiaotong Gui, Ashwin Agrawal, Anuj Pareek, SQ Truong, CD Nguyen, Van-Doan Ngo, Jayne
   947

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

Seekins, Francis G Blankenberg, AY Ng, et al. Deep learn-948 949 ing saliency maps do not accurately highlight diagnostically 950 relevant regions for medical image interpretation. MedRxiv, 951 2021. 6

- [100] Constantin Seibold, Simon Reiß, Saguib Sarfraz, 952 953 Matthias A. Fink, Victoria Mayer, Jan Sellner, Moon Sung 954 Kim, Klaus H. Maier-Hein, Jens Kleesiek, and Rainer 955 Stiefelhagen. Detailed annotations of chest x-rays via ct 956 projection for report understanding. In Proceedings of the 957 33th British Machine Vision Conference (BMVC), 2022. 6
- 958 [101] Ahmed Serag, Paul Aljabar, Gareth Ball, Serena J Counsell, James P Boardman, Mary A Rutherford, A David Edwards, 959 960 Joseph V Hajnal, and Daniel Rueckert. Construction of a 961 consistent high-definition spatio-temporal atlas of the de-962 veloping brain using adaptive kernel regression. Neuroim-963 age, 59(3):2255-2265, 2012. 6
- 964 [102] Arnaud Arindra Adiyoso Setio, Alberto Traverso, Thomas 965 De Bel, Moira SN Berens, Cas Van Den Bogaard, Piergior-966 gio Cerello, Hao Chen, Qi Dou, Maria Evelina Fantacci, Bram Geurts, et al. Validation, comparison, and combi-967 968 nation of algorithms for automatic detection of pulmonary 969 nodules in computed tomography images: the luna16 chal-970 lenge. Medical image analysis, 42:1-13, 2017. 6
- 971 [103] Divya Shanmugam, Davis Blalock, Guha Balakrishnan, 972 and John Guttag. Better aggregation in test-time augmen-973 tation. In Proceedings of the IEEE/CVF international con-974 ference on computer vision, pages 1214-1223, 2021. 2
- 975 [104] Amber L Simpson, Michela Antonelli, Spyridon Bakas, 976 Michel Bilello, Keyvan Farahani, Bram Van Ginneken, An-977 nette Kopp-Schneider, Bennett A Landman, Geert Litjens, 978 Bjoern Menze, et al. A large annotated medical image 979 dataset for the development and evaluation of segmentation algorithms. arXiv preprint arXiv:1902.09063, 2019. 6 980
- 981 [105] Samuel L Smith, Andrew Brock, Leonard Berrada, and So-982 ham De. Convnets match vision transformers at scale. arXiv 983 preprint arXiv:2310.16764, 2023. 2
- 984 [106] Yuxin Song, Jing Zheng, Long Lei, Zhipeng Ni, Baoliang 985 Zhao, and Ying Hu. CT2US: Cross-modal transfer learning 986 for kidney segmentation in ultrasound images with synthe-987 sized data. Ultrasonics, 122:106706, 2022. 6
- 988 [107] Joes Staal, Michael D Abràmoff, Meindert Niemeijer, 989 Max A Viergever, and Bram Van Ginneken. Ridge-based 990 vessel segmentation in color images of the retina. IEEE 991 transactions on medical imaging, 23(4):501-509, 2004. 2, 992 4, 1, 5
- [108] Vajira Thambawita, Steven A Hicks, Pål Halvorsen, 993 994 and Michael A Riegler. Divergentnets: Medical im-995 age segmentation by network ensemble. arXiv preprint 996 arXiv:2107.00283, 2021. 2
- [109] Santiago Vitale, José Ignacio Orlando, Emmanuel Iarussi, 997 998 and Ignacio Larrabide. Improving realism in patient-999 specific abdominal ultrasound simulation using cyclegans. 1000 International journal of computer assisted radiology and surgery, 15(2):183-192, 2020. 6 1001
- 1002 [110] Guotai Wang, Wenqi Li, Michael Aertsen, Jan Deprest, 1003 Sébastien Ourselin, and Tom Vercauteren. Aleatoric uncer-1004 tainty estimation with test-time augmentation for medical

image segmentation with convolutional neural networks. 1005 Neurocomputing, 338:34-45, 2019. 2 1006

- [111] Xinlong Wang, Xiaosong Zhang, Yue Cao, Wen Wang, 1007 Chunhua Shen, and Tiejun Huang. Seggpt: Segmenting 1008 everything in context. arXiv preprint arXiv:2304.03284, 1009 2023. 2 1010
- [112] Ziyang Wang, Jian-Qing Zheng, Yichi Zhang, Ge Cui, and 1011 Lei Li. Mamba-unet: Unet-like pure visual mamba for medical image segmentation. arXiv preprint arXiv:2402.05079, 2024. 2
- [113] Hallee E. Wong, Marianne Rakic, John Guttag, and Adrian V. Dalca. Scribbleprompt: Fast and flexible interactive segmentation for any biomedical image. arXiv:2312.07381, 2024. 2, 3, 4, 1, 6
- [114] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. Advances in neural information processing systems, 34:12077-12090, 2021. 2
- [115] Lukas Zbinden, Lars Doorenbos, Theodoros Pissas, Adrian Thomas Huber, Raphael Sznitman, and Pablo Márquez-Neila. Stochastic segmentation with conditional categorical diffusion models. In Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), pages 1119-1129, 2023. 2
- [116] Yingtao Zhang, Min Xian, Heng-Da Cheng, Bryar Shareef, Jianrui Ding, Fei Xu, Kuan Huang, Boyu Zhang, Chunping Ning, and Ying Wang. Busis: A benchmark for breast ultrasound image segmentation. In Healthcare, page 729. MDPI, 2022. 6
- [117] Oi Zhao, Shuchang Lvu, Wenpei Bai, Linghan Cai, Binghao Liu, Meijing Wu, Xiubo Sang, Min Yang, and Lijiang Chen. A multi-modality ovarian tumor ultrasound image dataset for unsupervised cross-domain semantic segmentation. CoRR, abs/2207.06799, 2022. 6
- [118] Xiaomei Zhao, Yihong Wu, Guidong Song, Zhenye Li, Yazhuo Zhang, and Yong Fan. A deep learning model integrating fcnns and crfs for brain tumor segmentation. Medical image analysis, 43:98–111, 2018. 2
- [119] Guoyan Zheng, Chengwen Chu, Daniel L Belavy, Bulat Ibragimov, Robert Korez, Tomaž Vrtovec, Hugo Hutt, Richard Everson, Judith Meakin, Isabel Lŏpez Andrade, et al. Evaluation and comparison of 3d intervertebral disc localization and segmentation methods for 3d t2 mr data: A grand challenge. Medical image analysis, 35:327-344, 2017. 4. 1. 5
- [120] Xin Zheng, Yong Wang, Guoyou Wang, and Jianguo Liu. 1051 Fast and robust segmentation of white blood cell images by 1052 self-supervised learning. Micron, 107:55-71, 2018. 2, 4, 1, 1053 5 1054

# InterConv: Set Interaction for Improved Biomedical Image Segmentation

# Supplementary Material

Label	Class
0	Background
1	Left cerebral white matter
2	Left cerebral cortex
3	Left lateral ventricle
4	Left inferior lateral ventricle
5	Left thalamus
6	Left caudate
7	Left putamen
8	3rd Ventricle
9	Brain stem
10	Left hippocampus
11	Left ventral DC
12	Right cerebral white matter
13	Right cerebral cortex
14	Right lateral ventricle
15	Right inferior lateral ventricle
16	Right thalamus
17	Right caudate
18	Right putamen
19	Right hippocampus
20	Right ventral DC

Table 2. OASIS and ADNI class labels.

#### **1055 6.** Automatic Segmentation

### 1056 6.1. Data

1059

1060

1063

1064

1065

1066

1067

1057 We provide training, validation and test sample sizes in Ta-1058 ble 3. Each dataset is split into four subsets:

- Train: during training stage used by the model for gradient descent.
- Stopping-Val: during training stage used for early stop-ping.
  - Performance-Val: during performance evaluation stage used for choosing the best model configuration.
    - Test: used for reporting model performance.

Set size distribution in each of the data split is shown in Fig. 8.

1068We perform the experiment in 2D since we thoroughly1069evaluate many variants of our method and the baseline,1070which would be prohibitive in 3D. After pre-processing,1071each scan volume is of size: (160, 192, 224). We take1072the 109th coronal slice. The resulting slices are of size1073 $160 \times 192$ . We list the class labels in Table 2.

# 6.2. Model Configuration

1	0	7	4
1	0	7	5

1076

1077

1078

1079

1080

1081

1089

1096

1097

We apply InterConv to three architectural designs: Simple-UNet, nnUNet, and SwinUNet.

SimpleUNets have 4 encoder layers with the same feature size at each layer. We train baselines and SimpleUNet-InterConv with comparable number of parameters. We detail various network sizes and corresponding output feature sizes per layer in Table 4.

We designed nnUNet-InterConv and SwinUNet-<br/>InterConv using the same network backbone, with both the<br/>same output feature size per layer as the original model, and<br/>with adjusted feature sizes to align InterConv-integrated<br/>models with baselines. We include detailed layer-wise<br/>feature configurations and their corresponding network<br/>sizes in Table 5 and 6.1082<br/>1083

**6.3.** Additional Experiments

Fig. 9 and 10 show that larger set size enables sharing of<br/>more informative accross the set for the SimpleUNet and<br/>SwinUNet architectures as well (along with the nnUNet<br/>shown in the main paper). Subjects with more scans avail-<br/>able on average achieve better segmentation quality com-<br/>pared to those with fewer scans.1090<br/>1091

### 7. Interactive Segmentation

7.1. Data

Datasets. Building upon large data collection efforts like 1098 MegaMedical [12, 96, 113], we use a collection of 74 1099 biomedical image segmentation datasets. We divide the col-1100 lection into 65 training datasets (Table 8) and 9 evaluation 1101 datasets (Table 7) following the same partitions as [113]. 1102 The evaluation datasets include a diverse array of biomed-1103 ical domains including eyes [107], abdominal [67], car-1104 diac [7, 94], bones [36], teeth [1], spine [119], cells [120], 1105 and lesions [2]. 1106

Each dataset was split into 60% training, 20% validation1107and 20% test as in [12, 96, 113], although not all splits were1108used. Each model was trained on the training splits of the1109training datasets. We report final results on the test split of1110the evaluation datasets.1111

Tasks. As in [12, 96, 113], we define a 2D segmentation1112task as a combination of (sub)dataset, axis (for 3D modali-<br/>ties), and label. For datasets with sub-datasets (e.g., malig-<br/>nant vs. benign lesions), each cohort is considered a sepa-<br/>rate task. Each segmentation label in multi-label datasets111211131114

### CVPR 2025 Submission #13795. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

С	V	Ρ	F	2
#1	3	7	9	5

	Subjects	Images
Train	777	1332
Stopping-Val	50	82
Performance-Val	47	82
Test	102	187
(a) OA	SIS T2*	

(b) OASIS T2-TSE

Subjects

780

49

48 99 Images

1381

81

87

188

	Subjects	Images
Train	636	3075
Stopping-Val	106	493
Performance-Val	106	504
Test	213	1045

(c) ADNI T2-FLAIR

Table 3	Sam	nle Sizes	s for	automatic	segmentatio	n experiments.
rubic 5	. Dum	pic Dille	, 101	uutomutic	beginemuno	n experimento.



Figure 8. Set Size distribution within splits for OASIS and ADNI data. We show the various number of scans for train, validation and test subjects, as well as the distribution in the test, stopping-val, performance-val and test splits.

is treated as a separate binary task. For multi-annotator
datasets, each annotator is considered a separate label. In
instance segmentation datasets, we train on one instance at
a time. For 3D modalities, we use the slice with the maximum label area ("maxslice") for each subject.

**Image Processing.** We resize images to  $128^2$  and rescale **1122** image intensities to [0,1]. **1123** 

### CVPR 2025 Submission #13795. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

	# ouput features per layer		
Network Size (# of parameters)	SimpleUNet-InterConv	Baseline	
60K	16	24	
120K	23	34	
240K	33	48	
500K	49	70	
1M	70	100	
2M	100	142	
4M	141	200	
8M	200	284	
16M	282	400	

Table 4. Network sizes for SimpleUNet architectures.

	<b>Encoder Feature Dimensions</b>	Network Size
nnUNet	[32, 64, 128, 256, 512, 512]	20.6M
nnUNet-InterConv	$\left[23, 46, 92, 184, 368, 368 ight]$	20.4M
	$\left[ 32, 64, 128, 256, 512, 512 \right]$	34.9M

Table 5. Network sizes for nnUNet architectures.



Figure 9. SimpleUNet-InterConv achieves higher segmentation accuracy on subjects with more scans available. We show performance improvement on SimpleUNet-InterConv over baseline SimpleUNet by averaging per-sample Dice improvement over subjects with a given number of scans. We evaluate with models with the best average test performance. Left figure shows Dice improvement on OASIS T2\* scans. Middle figure shows Dice improvement on OASIS T2-TSE scans. Right figure shows Dice improvement on ADNI T2-FLAIR scans. Shaded regions show 95% CI from bootstrapping with 1,000 runs.

### **1124 7.2. Additional Experiments**

We provide additional interactive segmentation visualizations with random click prompts in Fig. 11.
ScribblePrompt-InterConv is able to leverage the information across the input set to substantially improve the seg-

mentations compared to the baseline models, especially

1130 when given few user interactions.

	Embedding Dimension	Network Size
SwinUNet	96	27.1M
SwinUNet-InterConv	84	26.8M
	96	34.9M

Table 6. Network sizes for SwinUNet architectures.



Figure 10. SwinUNet-InterConv generally achieves higher segmentation accuracy on subjects with more scans available. We show performance improvement on SwinUNet-InterConv over baseline SwinUNet by averaging per-sample Dice improvement over subjects with a given number of scans. We evaluate with models with the best average test performance. Left figure shows Dice improvement on OASIS T2\* scans. Middle figure shows Dice improvement on ADNI T2-FLAIR scans. Shaded regions show 95% CI from bootstrapping with 1,000 runs.



(a) Segmentation visualization with 1 positive and 1 negative prompt

(b) Segmentation visualization with 8 positive and 8 negative prompts

Figure 11. ScribblePrompt-InterConv significantly improves interactive segmentation quality with random click prompts. We run ScribblePrompt-InterConv and baseline predictions on 4 input images from DRIVE, providing each image with 1 positive (green) and 1 negative (red) random click prompt, and on the same set of images with 8 positive (green) and 8 negative (red) random click prompts. The input samples come from different subjects within the same task. In each subfigure, **Top row** shows input image with provided prompts, along with ground truth segmentations, **middle row** shows ScribblePrompt-InterConv (ours) predictions, and **bottom row** shows ScribblePrompt (baseline) predictions overlayed in blue.

Dataset Name	Description	Scans	Labels	Modalities
ACDC [7]	Left and right ventricular endocardium	99	3	cine-MRI
BTCV	Bladder, uterus, rectum, small bowel	30	4	CT
Cervix [67]				
BUID [2]	Breast tumors	647	2	Ultrasound
DRIVE [107]	Blood vessels in retinal images	20	1	Optical camera
HipXRay [36]	Ilium and femur	140	2	X-Ray
PanDental [1]	Mandible and teeth	215	2	X-Ray
SCD [94]	Sunnybrook Cardiac Multi-Dataset Collec-	100	1	cine-MRI
	tion			
SpineWeb [119]	Vertebrae	15	1	T2-weighted
				MRI
WBC [120]	White blood cell cytoplasm and nucleus	400	2	Microscopy

Table 7. **Evaluation datasets**. For the relative size of datasets, we include the number of unique scans (subject and modality pairs) that each dataset has. These datasets were useen by the models during training. The validation splits of the datasets were used for model selection. We report final results on the test splits of these datasets.

Dataset Name	Description	Scans	Modalities
AbdominalUS [109]	Abdominal organ segmentation	1,543	Ultrasound
AMOS [53]	Abdominal organ segmentation	240	CT, MRI
BBBC003 [74]	Mouse embryos	15	Microscopy
BBBC038 [13]	Nuclei instance segmentation	670	Microscopy
BrainDev [33, 34, 64, 101]	Adult and neonatal brain atlases	53	Multimodal MRI
BrainMetShare[35]	Brain tumors	420	Multimodal MRI
BRATS [3, 4, 86]	Brain tumors	6,096	Multimodal MRI
BTCV Abdomi-	13 abdominal organs	30	СТ
BUSIS [116]	Breast tumors	163	Ultrasound
CAMUŠ [68]	Four-chamber and Apical two-chamber heart	500	Ultrasound
CDemris [55]	Human left atrial wall	60	CMR
CHAOS [56, 57]	Abdominal organs (liver, kidneys, spleen)	40	CT, T2-weighted MRI
CheXplanation [99]	Chest X-Ray observations	170	X-Ray
CT2US [106]	Liver segmentation in synthetic ultrasound	4,586	Ultrasound
CI-ORG[97]	Abdominal organ segmentation (overlap with LiTS)	140	CT
DDTI [90]	Thyroid segmentation	472	Ultrasound
EOphtha [24]	Eye microaneurysms and diabetic retinopathy	102	Optical camera
FeTA [89]	Fetal brain structures	80	Fetal MRI
FetoPlac [5]	Placenta vessel	6	Fetoscopic optical camera
FLARE [77]	Abdominal organs (liver, kidney, spleen, pan- creas)	361	CT
HaN-Seg [92]	Head and neck organs at risk	84	CT, T1-weighted MRI
HMC-QU [25, 59]	4-chamber (A4C) and apical 2-chamber (A2C)	292	Ultrasound
I2CVB [69]	Prostate (peripheral zone, central gland)	19	T2-weighted MRI
IDRID [93]	Diabetic retinopathy	54	Optical camera
ISBI-EM [15]	Neuronal structures in electron microscopy	30	Microscopy
ISIC [19]	Demoscopic lesions	2,000	Dermatology
ISLES [39]	Ischemic stroke lesion	180	Multimodal MRI
KiTS [38]	Kidney and kidney tumor	210	СТ
LGGFlair [11, 84]	TCIA lower-grade glioma brain tumor	110	MRI
LiTS [8]	Liver tumor	131	CT
LUNA [102]	Lungs	888	CT
MCIC [32]	Multi-site brain regions of schizophrenic patients	390	I I-weighted MRI
	Ovarian tumors	1,140	Oltrasound CT Multimedal MDI
MSD [104]	datasets	3,225	CI, Multimodal MRI
MuscleUS [83]	Muscle segmentation (biceps and lower leg)	8,169	Ultrasound
NCI-ISBI [10]	Prostate	30	T2-weighted MRI
NerveUS [88]	Nerve segmentation	5,635	Ultrasound
OASIS [45, 81]	Brain anatomy	414	TI-weighted MRI
OC1A500 [70]	Retinal vascular	500	OCT/OCTA
PanNuke [30]	Nuclei instance segmentation	/,901	Microscopy
PAXRay [100]	sub-diaphram in Chest X-Ray	852	х-кау
PROMISE12 [72]	Prostate	37	T2-weighted MRI
PPMI [22, 82]	Brain regions of Parkinson patients	1,130	T1-weighted MRI
QUBIQ [85]	Collection of 4 multi-annotator datasets (brain, kidney pancreas and prostate)	209	T1-weighted MRI, Multi- modal MRI CT
ROSE [80]	Retinal vessel	117	OCT/OCTA
SegTHOR [65]	Thoracic organs (heart, trachea, esophagus)	40	СТ
SegThy [62]	Thyroid and neck segmentation	532	MRI, Ultrasound
ssTEM [31]	Neuron membranes, mitochondria, synapses and	20	Microscopy
STARE [47]	Blood vessels in retinal images (multi annotator)	20	Optical camera
ToothSeg $[40]$	Individual teeth	508	X-Ray
VerSe [75]	Individual vertebrae	55	CT
WMH [63]	White matter hyper-intensities	60	Multimodal MRI
WORD [76]	Abdominal organ segmentation	120	CT

Table 8. **Training datasets**. We use the same set of 65 training datasets as [113]. For the relative size of datasets, we show the number of unique scans (subject and modality pairs) that each dataset has. 6