






PseudoClick: Interactive Image Segmentation with Click Imitation

Qin Liu^{1,2}, Meng Zheng²,
Benjamin Planche², Srikrishna Karanam², Terrence Chen², Marc
Niethammer¹, and Ziyang Wu²

¹ University of North Carolina at Chapel Hill, Chapel Hill NC, USA

² United Imaging Intelligence, Cambridge MA, USA
{first.last}@uii-ai.com, qinliu19@cs.unc.edu

Abstract. The goal of click-based interactive image segmentation is to obtain precise object segmentation masks with limited user interaction, *i.e.*, by a minimal number of user clicks. Existing methods require users to provide all the clicks: by first inspecting the segmentation mask and then providing points on mislabeled regions, iteratively. We ask the question: can our model directly predict where to click, so as to further reduce the user interaction cost? To this end, we propose **PseudoClick**, a generic framework that enables existing segmentation networks to propose candidate next clicks. These automatically generated clicks, termed pseudo clicks in this work, serve as an imitation of human clicks to refine the segmentation mask. We build **PseudoClick** on existing segmentation backbones and show how the click prediction mechanism leads to improved performance. We evaluate **PseudoClick** on 10 public datasets from different domains and modalities, showing that our model not only outperforms existing approaches but also demonstrates strong generalization capability in cross-domain evaluation. We obtain new state-of-the-art results on several popular benchmarks, *e.g.*, on the Pascal dataset, our model significantly outperforms existing state-of-the-art by reducing 12.4% number of clicks to achieve 85% IoU.

Keywords: click imitation, interactive image segmentation, pseudo click

1 Introduction

Recent years have seen tremendous progresses in segmentation methods for various applications, *e.g.*, semantic object/instance segmentation [1,2], video understanding [3,4], autonomous driving [5,6,7], and medical image analysis [8,9]. The success of these applications heavily relies on the availability of large-scale pixel-level annotation masks, which are very laborious and costly to obtain. Interactive image segmentation, which aims at extracting the object-of-interest using limited human interactions, is an efficient way to obtain these annotations.

While different interaction types have been investigated, including clicks [10,11], bounding boxes [12,13,14], and scribbles [15,16], we only focus on click-

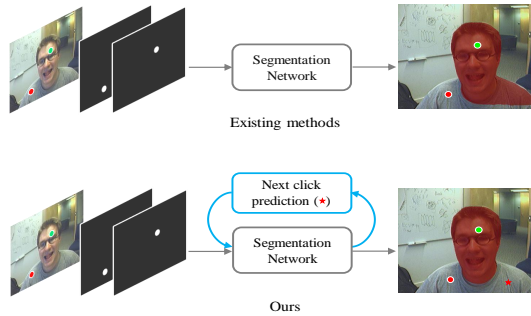


Fig. 1: The key difference between existing methods and our method is the ability for click prediction. (a) Illustration of existing methods. All clicks are provided by users. (b) Illustration of our method, which predicts the next click based on the current segmentation mask. This predicted pseudo click is used to refine the segmentation mask. Note the manual segmentation is not required at test time for the pseudo-click prediction.

based interactive segmentation because it has the simplest interaction, and well-established training and evaluation protocols [10,11]. Compared with scribble-based methods [15,16], click-based interactive segmentation requires no heuristics or complex procedures to simulate user input. Recent work on click-based interactive segmentation has resulted in state-of-the-art segmentation performance using various inference-time optimization schemes [17], which are computationally expensive due to multiple backward passes during inference. More recently, Sofiuk *et al.* proposed RITM [11], a simple feedforward approach requiring no inference-time optimization for click-based interactive segmentation that shows performance superior to all existing models when trained on the combined COCO [18] and LVIS [19] datasets with diverse and high-quality annotations.

While existing methods in this area have shown improved performance, *e.g.*, using inference-time optimization schemes [17] or iterative mask-guided training schemes [11], these methods still require users to provide all the clicks. Hence, users need to interactively inspect the resulting segmentation masks and then provide points for mislabeled regions. We ask the question: can the segmentation model directly predict where to click so as to further reduce the number of human clicks? To this end, we propose **PseudoClick**, a novel framework for interactive segmentation that equips existing segmentation backbones with the ability to predict user clicks automatically. As shown in Fig. 1, the key difference between existing methods and our proposed one is the ability to generate additional, beneficial, and “free” pseudo clicks as the prediction of human clicks for refining the segmentation mask.

PseudoClick is a generic framework that can be built upon existing segmentation backbones, including both CNNs and transformers. We equip an existing segmentation backbone with clicks prediction mechanism in the following pro-

cedures: we first introduce an error decoder, in parallel with the segmentation decoder of the backbone, that produces the false-positive (FP) and false-negative (FN) error maps for the current segmentation mask. We then extract a pseudo click from either the FP or the FN map, depending on which map contains the largest error. Specifically, a positive click should be generated from the FN map, while a negative click should be generated from the FP map. After that, the generated pseudo click (we generate one pseudo click each time) will be updated in the network input to refine the segmentation mask for the next forward pass. Note that the entire process is an imitation of the core human activity in interactive segmentation: visually estimating the segmentation errors, *i.e.*, over-segmentation (FP) or under-segmentation (FN), before determining what and where the next click should be.

We evaluate our method extensively on **10** public datasets (see Sec. 4.1). Evaluation results show that our model not only outperforms existing approaches but also demonstrates strong generalization capability in cross-domain evaluation on medical images. On the Pascal dataset, our model significantly outperforms existing state-of-the-art by reducing 12.4% and 11.4% number of clicks to achieve 85% and 90% IoU, respectively. For the cross-domain evaluation on BraTS [20] and ssTEM [21], our method significantly outperforms existing approaches to a large margin. Our main contributions are:

- 1) We propose **PseudoClick**, a novel interactive segmentation framework that directly imitates human clicks through the segmentation network and refines the segmentation mask with these imitated pseudo clicks. Our proposed framework differs from existing interactive segmentation methods in that it provides additional, beneficial, and “free” clicks during the annotation process for a human-in-the-loop.
- 2) We show that **PseudoClick** is an efficient and generic framework that can be built upon different types of segmentation backbones, including both CNNs and transformers, with little effort in tuning the hyper-parameters and modifying the network architectures.
- 3) We evaluate **PseudoClick** thoroughly on benchmarks from multiple domains and modalities with extensive comparison and cross-domain evaluation experiments. The results show that **PseudoClick** not only outperforms existing state-of-the-art approaches on in-domain benchmarks but also demonstrates strong generalization capability on cross-domain evaluation.

2 Related Work

Click-based interactive image segmentation. Click-based interactive segmentation has the simplest interaction and well-established training and evaluation protocols [10,11]. Xu *et al.* [10] first apply CNNs for interactive segmentation and propose a click simulation strategy for training that has inspired many future works [11,22,23]. Compared with previous click-based approaches [11,24,23], **PseudoClick** is unique because it is the first work that imitates human clicks and refines the segmentation with automatically generated pseudo clicks.

Other types of interactive feedback. Other than clicks, different types of user interaction have been explored in this field, including bounding boxes [13], scribbles [25], and interactions from multiple modalities [26]. The main drawback of bounding box-based approaches is the lack of a specific object reference inside the selected region, as well as a clear approach to correct the predicted mask. Inside outside guidance (IOG) [13] addresses this issue by combining clicks with bounding points of the target object and by allowing corrections of the predicted mask. Our method differs from all these in that we only use clicks as the interaction. Instead of exploring more complicated interactions, we try to imitate user clicks and to decrease their number required to acquire predefined accuracy.

Imitation learning and beyond. Imitation learning aims at mimicking human behavior for a given task [27]. While this field has recently gained attention due to computing and sensing advances as well as the rising demand for intelligent applications [28,29], it has never been explored in the interactive segmentation tasks. We claim that our method is highly related to imitation learning in that it imitates the core user activity in the interactive annotation process: visually estimating the segmentation errors before determining what and where the next point should be. SeedNet [24] first proposes a reinforcement learning approach for automatically generating automatic click. However, their method is not an imitation of human clicks because it automatically generates a sequence of clicks for achieving an implicit long-term reward without human intervention given the initial two clicks. Therefore, it is more like a post-processing approach. In contrast, our method is an imitation of human clicks because it explicitly quantifies the FP&FN errors and then generates the next click based on the estimated errors (just like a human annotator would do). Our method imitates this process by introducing a pseudo click generation mechanism on existing segmentation backbones, as discussed in Sec. 3. Our idea of predicting FP&FN errors for segmentation is also related to the idea of “prediction loss” proposed in [30], but the two ideas are investigated in significantly different problem settings.

3 Method

As shown in Fig. 2, **PseudoClick** is a generic framework that builds upon existing segmentation framework with two additional modules: a segmentation error decoder module and a clicks-encoding module that consists of encoding masks for both human clicks and pseudo clicks. In Sec. 3.1, we first describe the segmentation error decoder that outputs two error maps from which pseudo clicks are generated. In Sec. 3.2, we further describe the pseudo clicks generation process. In Sec. 3.3, we then introduce the encoding mechanism that transforms pseudo clicks into spatial signal for feeding into the segmentation backbone. In Sec. 3.4, we introduce the loss function for training our models. In Sec. 3.5, we conclude the whole section by providing additional implementation details about the proposed method.

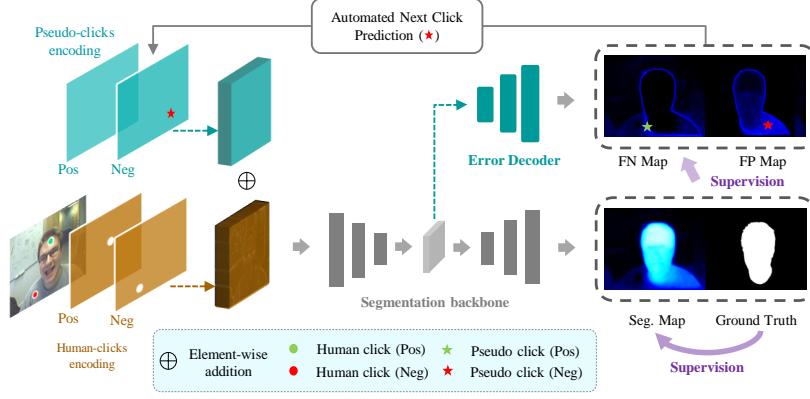


Fig. 2: **PseudoClick** overview. Given the image and clicks, the network outputs the segmentation map, along with two error maps that predict the false-positives and false-negatives of the segmentation mask, respectively. Then, a pseudo click will be generated from either the FP map or the FN map (see Sec. 3.2). After that, the network refines the segmentation mask in the second forward pass by adding the new pseudo click to the click encoding.

3.1 Segmentation Error Decoder

The segmentation error decoder is introduced in parallel with the segmentation decoder of the backbone. It generates two error maps that estimate the false positive (FP) and false negative (FN) errors of the segmentation mask. Both error maps are probability maps³ that can be trained as two binary segmentation tasks. The FP and FN maps generated by the error decoder are supervised by the ground-truth FP map \mathbf{M}_{fp} and ground-truth FN map \mathbf{M}_{fn} that can be obtained from the ground truth mask \mathbf{M} and the segmentation probability map \mathbf{P} in the following way: $\mathbf{M}_{fp} = \neg\mathbf{M} \wedge (\mathbf{P} \geq \tau)$; $\mathbf{M}_{fn} = \mathbf{M} \wedge (\mathbf{P} < \tau)$, where τ is a probability threshold (which is set to 0.5 by default). Since training the segmentation error decoder is formulated as two binary segmentation tasks, the error decoders and the segmentation decoder can be trained end-to-end in a multi-task learning manner (See Sec. 3.4).

To help readers better understand the function of the proposed error decoder, we show in Fig. 3 the error maps generated by the error decoder. We observe that the error maps provide a meaningful estimation of the true errors, as one can easily estimate by comparing the segmentation mask (red masks in first row) with the ground truth (white masks in the last row). The accuracy of these error maps is essential for reliable pseudo clicks generation, introduced next.

³ Note that we call these *probability maps*, but in general they will likely be miscalibrated. If desired, calibration can be improved, for example, using the approach in [31].

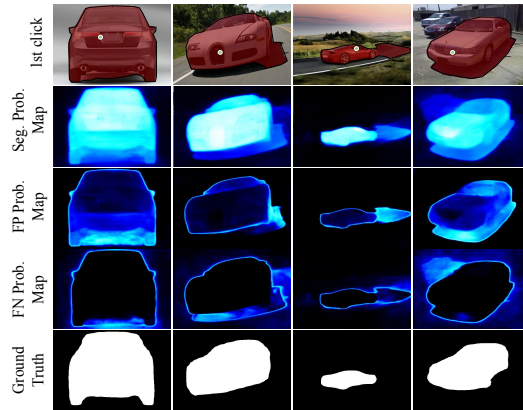


Fig. 3: Illustration of the output error maps of **PseudoClick**. Given the first click (dots in the first row), the network outputs a segmentation map (second row), along with estimated FP and FN error maps (third and fourth rows). Each column represents a test case. All test images are from the Cars dataset [32].

3.2 Pseudo Clicks Generation

Given the predicted FP and FN error maps for the current segmentation mask, our method generates one pseudo click from either the FP or the FN map, depending on which map contains the largest error region. First, we transform the two error maps to two binary masks through a predefined threshold (*i.e.* 0.5), followed by extracting a positive/negative click from one of the two binary masks that contains the largest connected error region—the extracted pseudo click locates at the center of this region. If the largest error region is from the FP mask, then the extracted pseudo click is negative (*i.e.*, indicating that this region should not be segmented); if the largest error region is from the FN mask, then the extracted pseudo click is positive (*i.e.*, indicating that this region should be segmented). Finally, we encode the new pseudo click as a disk on the encoding maps designed exclusively for pseudo clicks, introduced next.

3.3 Pseudo Clicks Encoding

As shown in Fig. 2, human clicks and pseudo clicks are encoded separately as small binary disks in the corresponding encoding maps, resulting in two 2-channel disk maps. Positive clicks are encoded in the positive disk map while the negative clicks are encoded in the negative disk map. Note that our choice of using disk maps instead of Gaussian maps for clicks encoding is inspired by RITM [11], which shows that disk maps are more efficient and effective than Gaussian maps. We perform element-wise addition to merge the feature maps extracted from pseudo clicks and feature maps extracted from the combination of image and human clicks. The merged feature maps will be fed into the segmentation backbone for end-to-end training.

3.4 Loss Function

Although binary cross entropy (BCE) loss is widely used to supervise the interactive segmentation tasks [13,33,34,35,36], we instead use normalized focal loss (NFL) [37], which allows for faster convergence and better accuracy than BCE, as discussed in [11]. The NFL loss L can be written as:

$$L(\mathbf{P}(i,j)) = -\frac{1}{\sum_{i,j} \mathbf{P}(i,j)} (1 - \mathbf{P}(i,j))^\gamma \log \mathbf{P}(i,j) \quad (1)$$

where $\mathbf{P}(i,j)$ denotes the prediction \mathbf{P} at point (i,j) and $\gamma > 0$ is a tunable focusing parameter (as in the focal loss [38]). Since the error branch can be supervised as two binary segmentation tasks (see Sec. 3.1), we also use the NFL loss for them during training. Therefore, the overall loss is a combination of three NFL loss functions:

$$L = \sum_{i,j} (\lambda_1 L_{seg}(i,j) + \lambda_2 L_{fp}(i,j) + \lambda_3 L_{fn}(i,j)) \quad (2)$$

where $\lambda_1, \lambda_2, \lambda_3 > 0$ represent the weights for each component; $L_{seg}(i,j)$, $L_{fp}(i,j)$, and $L_{fn}(i,j)$ denote $L(\mathbf{P}(i,j))$, $L(\mathbf{E}_{fp}(i,j))$, and $L(\mathbf{E}_{fn}(i,j))$, respectively.

3.5 Implementation Details

Click simulation for training and evaluation. We automatically simulate human clicks based on the ground truth and current segmentation for fast training and evaluation. For training, we use a combination of random and iterative click simulation strategies, similar to [11]. The random click simulation strategy generates a set of positive and negative clicks without considering the order between them [10,17,39]. In contrast, the iterative simulation strategy generates clicks sequentially—a new click is generated based on the erroneous region of a prediction produced by a model using the set of previous clicks [36,40,41]. Once the model is trained, there are two modes to evaluate it: automatic evaluation and human evaluation. For automatic evaluation, we adopt the iterative click simulation strategy. Note that the automatically simulated clicks may be different from clicks generated by human evaluation. We present in the supplementary materials some qualitative results obtained via human evaluation.

Previous segmentation as an additional input channel. It is natural to incorporate the output segmentation mask from previous interaction as an input for the next correction, providing additional prior information that can help improve the segmentation quality. The previous segmentation mask is added as an additional channel to the RGB image, resulting in a 4-channel image as the input. This 4-channel image will be concatenated with the human clicks encoding maps, which have two channels (positive and negative clicks are encoded in separate channels). Note that the additional mask input is not shown in Fig. 2 for brevity. For the first interaction, we feed an empty mask to our model.

Post-processing using error maps. Pseudo clicks are extracted from the error maps. Actually, the error maps can be directly used to refine the segmentation mask (*e.g.* via simple post-processing introduced in Sec. 4.4). We argue that post-processing based on FP&FN maps can be regarded as a by-product of our core contribution—**PseudoClick**—which takes a step further towards the general idea of click imitation.

4 Experiments

4.1 Evaluation Details

Datasets. We evaluate **PseudoClick** on **10** public datasets: GrabCut [16], Berkeley [42], DAVIS [43], Pascal [44], Semantic Boundaries Dataset (SBD) [45], Brain Tumor Segmentation challenge (BraTS) [20], ssTEM [21], Cars [32], COCO [18], and LVIS [19]. Since COCO and LVIS datasets share the same set of images and are complementary to each other in terms of annotation quality and object categories, they can be combined as an ideal training set for the interactive image segmentation task [11]. Therefore, we use the combined COCO+LVIS for training and the remaining 8 datasets for evaluation. Specifically, we use the training set of the COCO+LVIS dataset for training and its validation set for model selection. The Cars [32] dataset is only used for qualitative evaluation. We test the trained **PseudoClick** models on the test set of the remaining 7 datasets; no finetuning is conducted on these datasets. For the DAVIS and BraTS datasets, we do not use the original videos or volumes. Instead, we extract 345 and 369 2D slices from the two 3D datasets, respectively. We extracted from each volume in the BraTS the slice that contains the largest tumor area. The two medical image datasets, BraTS [20] and ssTEM [21], are used for cross-domain evaluation (see Sec. 4.3) because our models are trained with natural images, which are significantly different from images from medical domain. We refer the readers to the supplementary material for more details.

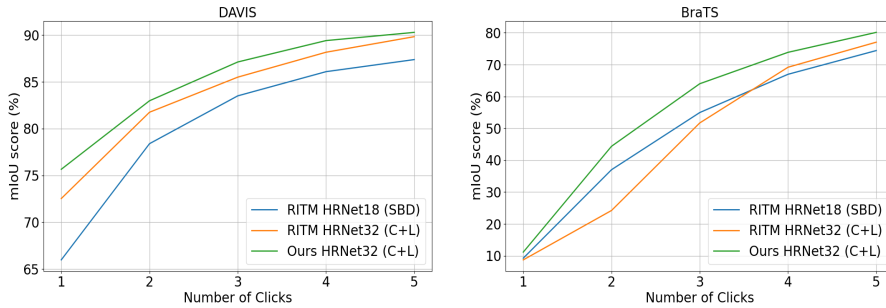


Fig. 4: Comparison of mIoU for the first five clicks on DAVIS and BraTS datasets.

Segmentation backbone. We choose HRNet-18 and HRNet-32 [46] as backbones for our PseudoClick model. To show that PseudoClick is a generic framework that can be built upon most of existing segmentation backbones, we also implement PseudoClick on two recently proposed segmentation transformers: SegFormer-B5 [47] and HRFormer-base [48]. While the two transformers have shown promising preliminary results in our experiments, HRNet still outperforms them to a large margin under the same experimental settings.

Evaluation modes and metrics. We evaluate the trained models with two modes: automatic evaluation and human evaluation. For automatic evaluation, we simulate human clicks based on the ground truth and current segmentation mask (see Sec. 3.5); for human evaluation, a human-in-the-loop will provide clicks based on his/her subjective evaluation. We use the standard Number of Clicks (NoC) metric that measures the number of clicks required to achieve predefined Intersection over Union (IoU). Specifically, we use NoC@85 and NoC@90 as two main metrics to measure the number of clicks required to obtain 85% and 90% IoU, respectively.

More implementation details. We adopt the same data augmentation techniques as in RITM [11] for fair comparison. Pseudo-clicks were not counted in the NoC@85% or NoC@90% metrics as they are “free” in terms of human labor. We implement our models in Python, using PyTorch [49]. We train the models for 200 epochs on COCO+LVIS dataset with an initial learning rate 5×10^{-4} , which will decrease to 5×10^{-5} after the epoch 50. During training, we crop images with a fixed size of 320×480 and set the batch size to 32. We optimize using Adam with $\beta_1 = 0.9$, $\beta_2 = 0.999$. All models are trained and tested on a single NVIDIA RTX A6000 GPU.

4.2 Comparison with State-of-the-Art

We compare our results with existing state-of-the-art methods, including f-BRS [17], IA+SA [40], FCA [54], and RITM [11]. Tab. 1 shows the quantitative results. Our proposed PseudoClick approach outperforms all the methods across the five benchmark datasets. For example, compared with RITM on the Pascal dataset, our model uses 12.4% and 11.4% fewer number of clicks for achieving 85% and 90% IoU, respectively. Fig. 4 shows plots for mean IoU with respect to the first several clicks for the DAVIS and BraTS datasets. Our approach shows continuous improvement in terms of accuracy and stability. We also visualize the evaluation process of the proposed method on some images in Fig. 6. From the qualitative results, we can see that the pseudo clicks automatically generated from our model can accurately focus on false positive and false negative regions. Hence, they are able to refine the predicted segmentation masks and thereby alleviate human annotation effort. Computational analysis is shown in Fig. 7.

Method		GrabCut Berkeley		SBD	DAVIS		Pascal	
		NoC@90	NoC@90	NoC@90	NoC@85	NoC@90	NoC@85	NoC@90
GC [15]	ICCV01	10.00	14.22	15.96	15.13	17.41	-	-
RW [50]	PAMI06	13.77	14.02	15.04	16.71	18.31	-	-
GM [51]	IJCV09	14.57	15.96	17.60	18.59	19.50	-	-
GSC [52]	CVPR10	9.12	12.57	15.31	15.35	17.52	-	-
ESC [52]	CVPR10	9.20	12.11	14.86	15.41	17.70	-	-
DIOS [10]	CVPR16	6.04	8.65	-	-	12.58	6.88	-
LD [53]	CVPR18	4.79	-	10.78	5.05	9.57	-	-
BRS [33]	CVPR19	3.60	5.08	9.78	5.58	8.24	-	-
f-BRS [17]	CVPR20	2.72	4.57	7.73	5.04	7.41	-	-
IA+SA [40]	ECCV20	3.07	4.94	-	5.16	-	3.18	-
FCA [54]	CVPR20	2.08	3.92	-	-	7.57	2.69	-
SBD RITM-H18		2.04	3.22	<u>5.43</u>	4.94	6.71	2.51	3.03
C+L RITM-H32		1.56	<u>2.10</u>	5.71	4.11	5.34	2.19	2.57
SBD H18	w/ PC	2.04	3.23	5.40	4.81	6.57	2.34	2.74
SBD H32	w/ PC	1.84	2.98	5.61	4.74	6.16	2.37	2.78
C+L H32	w/o PC	<u>1.55</u>	2.11	5.68	<u>4.09</u>	<u>5.27</u>	<u>2.14</u>	<u>2.52</u>
C+L H32	w/ PC	1.50	2.08	5.54	3.79	5.11	1.94	2.25

Table 1: Evaluation results on GrabCut, Berkeley, DAVIS, SBD, and Pascal datasets. Our models are trained on either SBD or COCO+LVIS (denoted as C+L above) datasets. The best results are set in bold; the second best results are shown underlined. “H18” and “H32” represent “HRNet-18” and “HRNet-32”, respectively. “PC” means the model is implemented with pseudo clicks. The metrics are NoC@85% and NoC@90%, representing the number of clicks required to achieve 85% and 90% IoU, respectively.

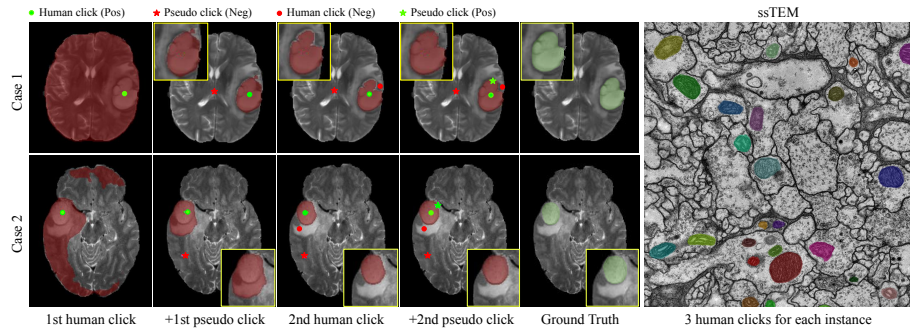


Fig. 5: Cross-domain evaluation on two medical image datasets: BraTS [20] (left) and ssTEM [21] (right). The evaluation is performed by a human annotator through our internally developed interactive segmentation GUI. For the BraTS dataset, we use two clicks for each image. For the ssTEM dataset, we strictly use three clicks for each instance. Note that our model is trained on natural images, but shows very robust results on the two medical datasets.

4.3 Cross-Domain Evaluation

To evaluate the generalization capability of the proposed method, we conduct cross-domain evaluation on two medical image datasets: BraTS [20] and ssTEM [21]. Specifically, we directly apply our **PseudoClick** models trained on SBD or COCO+LVIS datasets to the medical images without finetuning (medical images in grayscale are replicated 3 times channel-wise to be of the same channel dimension with RGB images). We report cross-domain evaluation results in Tab. 2 and Tab. 3. We observe that our models generalize very well to medical images without fine-tuning. Note that these results are evaluated by a human annotator. Some qualitative results on the two medical datasets are shown in Fig. 5.

Method	Train	Finetune	Backbone	mIoU@3	mIoU@5
RITM [11]	SBD	N/A	HRNet18	54.9	74.4
RITM [11]	C+L	N/A	HRNet32	51.7	77.1
Ours	SBD	N/A	HRNet18	54.5	74.6
Ours	C+L	N/A	HRNet32	64.0	80.1

Table 2: Cross-domain evaluation on the BraTS dataset. The evaluation measure is mean IoU (%) given 3 or 5 human clicks.

Method	Train	Finetune	Backbone	#Clicks	mIoU
Curve-GCN [55]	CityScapes	N/A	ResNet-50	2	60.9
I0G [13]	Pascal	N/A	ResNet-101	3	83.7
RITM [11]	SBD	N/A	HRNet18	3	77.3
RITM [11]	C+L	N/A	HRNet32	3	86.4
Ours	SBD	N/A	HRNet18	3	80.9
Ours	C+L	N/A	HRNet32	3	87.2

Table 3: Cross-domain evaluation on the ssTEM [21] dataset. The evaluation measure is mean IoU (%) given 2 or 3 human clicks. The results for IOG and Curve-GCN methods are copied from the corresponding papers.

4.4 Comparison Study

Segmentation backbone comparison. We have demonstrated in Tab. 1 that our method outperforms existing state-of-the-art when using HRNet-32 as its backbone. In this study, we implement other backbones including two recently proposed vision transformers, SegFormer [47] and HRFormer [48], that show encouraging results when compared with CNNs. Tab. 4 shows the comparison results. All models are pre-trained on ImageNet [56] and finetuned on COCO+LVIS dataset with NFL loss function. We use SegFormer-B5 [47] and HRFormer-Base [48] for the transformers. Models are trained and evaluated

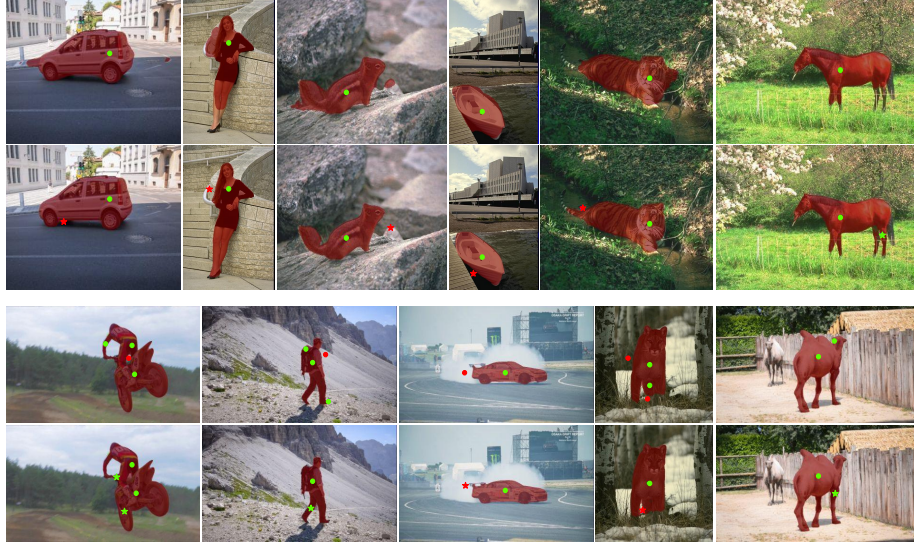


Fig. 6: Qualitative evaluation of PseudoClick model on natural images. Top: segmentation results with one pseudo click. The first row shows the results of the first human click; the second row shows the improvement of the segmentation mask with the help of the first pseudo click. Bottom: segmentation results with IoU greater than 90% given several clicks. The first row shows the segmentation using only human clicks. The second row shows the segmentation with both human and pseudo clicks. The color and shape of a click follow the same rule shown in Fig. 5.

using pseudo clicks. The evaluation measure is NoC@85%. While the transformers achieve decent results, we actually spend little time in tuning the hyperparameters and modifying the architecture when transferring from CNNs to transformers. This demonstrates the flexibility and generalization capability of our framework.

Loss functions comparison. In this study, we train PseudoClick models using four different loss functions: binary cross entropy (BCE) loss, focal loss (FL) [38], Soft IoU loss [57], and normalized focal loss (NFL) [37]. Each experiment uses the HRNet32 model. All the four models are trained on the COCO+LVIS dataset. Results in Tab. 5 show that training with NFL leads to the best accuracy.

Training datasets comparison. In this study, we train PseudoClick models on four different training datasets: Pascal, SBD, LVIS, and COCO+LVIS. We test on four datasets: Pascal, SBD, Berkeley, and DAVIS. For each experiment, our model is based on HRNet32 and is trained with NFL loss function. We

Backbone	Berkeley	SBD	DAVIS	Pascal
ResNet51	1.94	3.49	4.87	2.18
ResNet101	1.85	3.44	4.14	2.02
SegFormer	2.54	4.10	4.11	2.26
HRFormer	1.84	4.53	4.80	2.62

Table 4: Comparison study for different segmentation backbones.

Loss	Berkeley	SBD	DAVIS	Pascal
BCE	1.44	3.53	3.97	1.98
FL	1.43	3.54	3.79	2.01
Soft IoU	1.44	3.63	3.96	2.10
NFL	1.40	3.46	3.79	1.94

Table 5: Comparison study for different loss functions.

Train	Berkeley	SBD	DAVIS	Pascal
Pascal	2.33	5.87	5.67	2.66
SBD	1.67	3.51	4.74	2.37
LVIS	2.63	5.40	6.97	3.14
C+L	1.40	3.46	3.79	1.94

Table 6: Comparison study for different training datasets.

Model	Param/M	FLOPs/G	Speed/s
RITM-H32	30.95	16.57	0.137
Ours-H32-PC	36.79	18.43	0.185

Table 7: Computational analysis. Speed is measured as second per click (including a pseudo click for ours) with a NVIDIA A6000 GPU.

report results in Tab. 6. We observe that the model trained on COCO+LVIS shows the best performance, highlighting the benefit of combining COCO and LVIS for training interactive segmentation models. We also notice that on the Pascal dataset, model trained on COCO+LVIS dataset is even better than the model trained on Pascal dataset. This, again, highlights the strengths of COCO+LVIS dataset: 1) large dataset size. The number of annotated instances in COCO+LVIS dataset is $50\times$ and $170\times$ times larger than SBD and Pascal, respectively; 2) diverse and high annotation quality.

Post-processing vs. pseudo clicks. In this study, we directly use the two error maps for refining the segmentation mask. The two error maps serve as a regularization for the segmentation branch during training. As shown in Fig. 3, they quantize the segmentation mask reasonably well, and thus can be used for refining the segmentation mask. To achieve this goal, we simply subtract the two error maps from the segmentation map (all three maps are probability maps). We compare the post-processing with adding one pseudo click. The comparison results are shown in Tab. 8. We emphasize that the post-processing based on FP&FN maps can also be regarded as a contribution of our work as it is a by-product of our core contribution. Given two human clicks on the BraTS dataset, the relative mIoU obtained by adding one pseudo-click is 5.2% higher than the mIoU by post-processing, which is substantial considering the strong performance of post-processing.

mIoU@Human-clks	BraTS			DAVIS		
	2	3	5	2	3	5
Baseline (BL)	23.2	51.2	77.3	80.2	85.6	89.3
BL+post-processing	42.6	63.5	79.7	81.3	86.2	90.1
BL+1 pseudo-click	44.8	64.0	80.1	83.7	87.4	90.8

Table 8: Adding one pseudo-click vs. post-processing. Comparison of post-processing and pseudo clicks for mask refinement given a fixed number of human clicks. The Baseline model above is our best **PseudoClick** model (C+L HRNet32) on the BraTS dataset.

5 Limitations

The major limitation of the proposed method is that pseudo clicks may not be as accurate as human clicks, and therefore may cause the segmentation accuracy to drop. This may lead to extra work for users to withdraw the poorly placed pseudo clicks or to correct the error by putting more points. Fortunately, this issue has been greatly alleviated by separating the encoding maps for the two types of clicks. By separating the two encoding maps, the inaccurate pseudo clicks are tolerated during the training and less likely to cause accuracy to drop during evaluation. In the early stage of this project, we discovered this issue. After separating the two types of encoding maps, as implemented in our current architecture, this issue has been significantly eliminated.

6 Conclusion

We proposed **PseudoClick**, a novel interactive segmentation framework that automatically imitates human clicks and efficiently segments objects with the imitated pseudo clicks. **PseudoClick** is a general framework that can be built upon different types of segmentation backbones, including both CNNs and transformers, with little effort in tuning the hyper-parameters and modifying the network architectures. We evaluated **PseudoClick** thoroughly on benchmarks from multiple domains and modalities with extensive comparison and cross-domain evaluation experiments that demonstrated the effectiveness as well as the generalization capability of the proposed method.

Acknowledgements

Research reported in this publication was supported by the National Institutes of Health (NIH) under award number NIH 1R01AR072013. The content is solely the responsibility of the authors and does not necessarily represent the official views of the NIH.

References

1. A. Garcia-Garcia, S. Orts-Escolano, S. Oprea, V. Villena-Martinez, and J. Garcia-Rodriguez, “A review on deep learning techniques applied to semantic segmentation,” *arXiv preprint arXiv:1704.06857*, 2017.
2. S. Minaee, Y. Y. Boykov, F. Porikli, A. J. Plaza, N. Kehtarnavaz, and D. Terzopoulos, “Image segmentation using deep learning: A survey,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
3. L. Yang, Y. Fan, and N. Xu, “Video instance segmentation,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 5188–5197, 2019.
4. N. Xu, L. Yang, Y. Fan, D. Yue, Y. Liang, J. Yang, and T. Huang, “Youtube-vos: A large-scale video object segmentation benchmark,” *arXiv preprint arXiv:1809.03327*, 2018.
5. M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, “The cityscapes dataset for semantic urban scene understanding,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3213–3223, 2016.
6. A. Geiger, P. Lenz, and R. Urtasun, “Are we ready for autonomous driving? the KITTI vision benchmark suite,” in *2012 IEEE conference on computer vision and pattern recognition*, pp. 3354–3361, IEEE, 2012.
7. G. Neuhold, T. Ollmann, S. Rota Bulò, and P. Kotschieder, “The mapillary vistas dataset for semantic understanding of street scenes,” in *Proceedings of the IEEE international conference on computer vision*, pp. 4990–4999, 2017.
8. D. Shen, G. Wu, and H.-I. Suk, “Deep learning in medical image analysis,” *Annual review of biomedical engineering*, vol. 19, pp. 221–248, 2017.
9. G. Litjens, T. Kooi, B. E. Bejnordi, A. A. A. Setio, F. Ciompi, M. Ghafoorian, J. A. Van Der Laak, B. Van Ginneken, and C. I. Sánchez, “A survey on deep learning in medical image analysis,” *Medical image analysis*, vol. 42, pp. 60–88, 2017.
10. N. Xu, B. Price, S. Cohen, J. Yang, and T. S. Huang, “Deep interactive object selection,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 373–381, 2016.
11. K. Sofiiuk, I. A. Petrov, and A. Konushin, “Reviving iterative training with mask guidance for interactive segmentation,” *arXiv preprint arXiv:2102.06583*, 2021.
12. N. Xu, B. Price, S. Cohen, J. Yang, and T. Huang, “Deep grabcut for object selection,” *arXiv preprint arXiv:1707.00243*, 2017.
13. S. Zhang, J. H. Liew, Y. Wei, S. Wei, and Y. Zhao, “Interactive object segmentation with inside-outside guidance,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12234–12244, 2020.
14. J. Wu, Y. Zhao, J.-Y. Zhu, S. Luo, and Z. Tu, “Milcut: A sweeping line multiple instance learning paradigm for interactive image segmentation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 256–263, 2014.
15. Y. Y. Boykov and M.-P. Jolly, “Interactive graph cuts for optimal boundary & region segmentation of objects in ND images,” in *Proceedings eighth IEEE international conference on computer vision. ICCV 2001*, vol. 1, pp. 105–112, IEEE, 2001.
16. C. Rother, V. Kolmogorov, and A. Blake, “Grabcut: Interactive foreground extraction using iterated graph cuts,” *ACM transactions on graphics (TOG)*, vol. 23, no. 3, pp. 309–314, 2004.

17. K. Sofiiuk, I. Petrov, O. Barinova, and A. Konushin, “f-brs: Rethinking backpropagating refinement for interactive segmentation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8623–8632, 2020.
18. T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft COCO: Common objects in context,” in *European conference on computer vision*, pp. 740–755, Springer, 2014.
19. A. Gupta, P. Dollar, and R. Girshick, “Lvis: A dataset for large vocabulary instance segmentation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5356–5364, 2019.
20. U. Baid, S. Ghodasara, M. Bilello, S. Mohan, E. Calabrese, E. Colak, K. Farahani, J. Kalpathy-Cramer, F. C. Kitamura, S. Pati, *et al.*, “The RSNA-ASNR-MICCAI BRATS 2021 benchmark on brain tumor segmentation and radiogenomic classification,” *arXiv preprint arXiv:2107.02314*, 2021.
21. S. Gerhard, J. Funke, J. Martel, A. Cardona, and R. Fetter, “Segmented anisotropic ssTEM dataset of neural tissue,” *figshare*, pp. 0–0, 2013.
22. Q. Liu, Z. Xu, Y. Jiao, and M. Niethammer, “isegformer: Interactive image segmentation with transformers,” *arXiv preprint arXiv:2112.11325*, 2021.
23. X. Chen, Z. Zhao, Y. Zhang, M. Duan, D. Qi, and H. Zhao, “Focalclick: Towards practical interactive image segmentation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1300–1309, 2022.
24. G. Song, “Seednet,” in *CVPR*, 2018.
25. H. K. Cheng, Y.-W. Tai, and C.-K. Tang, “Modular interactive video object segmentation: Interaction-to-mask, propagation and difference-aware fusion,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5559–5568, 2021.
26. H. Ding, S. Cohen, B. Price, and X. Jiang, “Phraseclick: toward achieving flexible interactive segmentation by phrase and click,” in *European Conference on Computer Vision*, pp. 417–435, Springer, 2020.
27. A. Hussein, M. M. Gaber, E. Elyan, and C. Jayne, “Imitation learning: A survey of learning methods,” *ACM Computing Surveys (CSUR)*, vol. 50, no. 2, pp. 1–35, 2017.
28. T. Osa, J. Pajarinen, G. Neumann, J. A. Bagnell, P. Abbeel, and J. Peters, “An algorithmic perspective on imitation learning,” *arXiv preprint arXiv:1811.06711*, 2018.
29. T. Zhang, Z. McCarthy, O. Jow, D. Lee, X. Chen, K. Goldberg, and P. Abbeel, “Deep imitation learning for complex manipulation tasks from virtual reality teleoperation,” in *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 5628–5635, IEEE, 2018.
30. D. Yoo and I. S. Kweon, “Learning loss for active learning,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 93–102, 2019.
31. Z. Ding, X. Han, P. Liu, and M. Niethammer, “Local temperature scaling for probability calibration,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 6889–6899, 2021.
32. J. Krause, M. Stark, J. Deng, and L. Fei-Fei, “3d object representations for fine-grained categorization,” in *4th International IEEE Workshop on 3D Representation and Recognition (3dRR-13)*, (Sydney, Australia), 2013.
33. W.-D. Jang and C.-S. Kim, “Interactive image segmentation via backpropagating refinement scheme,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5297–5306, 2019.

34. K.-K. Maninis, S. Caelles, J. Pont-Tuset, and L. Van Gool, “Deep extreme cut: From extreme points to object segmentation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 616–625, 2018.
35. J. H. Liew, S. Cohen, B. Price, L. Mai, S.-H. Ong, and J. Feng, “Multiseg: Semantically meaningful, scale-diverse segmentations from minimal user input,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 662–670, 2019.
36. S. Majumder and A. Yao, “Content-aware multi-level guidance for interactive instance segmentation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11602–11611, 2019.
37. K. Sofiiuk, O. Barinova, and A. Konushin, “Adaptis: Adaptive instance selection network,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 7355–7363, 2019.
38. T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, “Focal loss for dense object detection,” in *Proceedings of the IEEE international conference on computer vision*, pp. 2980–2988, 2017.
39. R. Benenson, S. Popov, and V. Ferrari, “Large-scale interactive object segmentation with human annotators,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11700–11709, 2019.
40. T. Kontogianni, M. Gygli, J. Uijlings, and V. Ferrari, “Continuous adaptation for interactive object segmentation by learning from corrections,” in *European Conference on Computer Vision*, pp. 579–596, Springer, 2020.
41. S. Mahadevan, P. Voigtlaender, and B. Leibe, “Iteratively trained interactive segmentation,” *arXiv preprint arXiv:1805.04398*, 2018.
42. D. Martin, C. Fowlkes, D. Tal, and J. Malik, “A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics,” in *Proceedings Eighth IEEE International Conference on Computer Vision. ICCV 2001*, vol. 2, pp. 416–423, IEEE, 2001.
43. F. Perazzi, J. Pont-Tuset, B. McWilliams, L. Van Gool, M. Gross, and A. Sorkine-Hornung, “A benchmark dataset and evaluation methodology for video object segmentation,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 724–732, 2016.
44. M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, “The PASCAL visual object classes (VOC) challenge,” *International journal of computer vision*, vol. 88, no. 2, pp. 303–338, 2010.
45. B. Hariharan, P. Arbeláez, L. Bourdev, S. Maji, and J. Malik, “Semantic contours from inverse detectors,” in *2011 International Conference on Computer Vision*, pp. 991–998, IEEE, 2011.
46. J. Wang, K. Sun, T. Cheng, B. Jiang, C. Deng, Y. Zhao, D. Liu, Y. Mu, M. Tan, X. Wang, *et al.*, “Deep high-resolution representation learning for visual recognition,” *IEEE transactions on pattern analysis and machine intelligence*, 2020.
47. E. Xie, W. Wang, Z. Yu, A. Anandkumar, J. M. Alvarez, and P. Luo, “Segformer: Simple and efficient design for semantic segmentation with transformers,” *arXiv preprint arXiv:2105.15203*, 2021.
48. Y. Yuan, R. Fu, L. Huang, W. Lin, C. Zhang, X. Chen, and J. Wang, “Hrformer: High-resolution transformer for dense prediction,” *arXiv preprint arXiv:2110.09408*, 2021.
49. A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimeshein, L. Antiga, *et al.*, “Pytorch: An imperative style, high-performance deep learning library,” *Advances in neural information processing systems*, vol. 32, pp. 8026–8037, 2019.

50. L. Grady, “Random walks for image segmentation,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 28, no. 11, pp. 1768–1783, 2006.
51. X. Bai and G. Sapiro, “A geodesic framework for fast interactive image and video segmentation and matting,” in *2007 IEEE 11th International Conference on Computer Vision*, pp. 1–8, IEEE, 2007.
52. V. Gulshan, C. Rother, A. Criminisi, A. Blake, and A. Zisserman, “Geodesic star convexity for interactive image segmentation,” in *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pp. 3129–3136, IEEE, 2010.
53. Z. Li, Q. Chen, and V. Koltun, “Interactive image segmentation with latent diversity,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 577–585, 2018.
54. Z. Lin, Z. Zhang, L.-Z. Chen, M.-M. Cheng, and S.-P. Lu, “Interactive image segmentation with first click attention,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13339–13348, 2020.
55. H. Ling, J. Gao, A. Kar, W. Chen, and S. Fidler, “Fast interactive object annotation with curve-GCN,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5257–5266, 2019.
56. J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255, Ieee, 2009.
57. M. A. Rahman and Y. Wang, “Optimizing intersection-over-union in deep neural networks for image segmentation,” in *International symposium on visual computing*, pp. 234–244, Springer, 2016.