
HAIR: A Dataset of Historic Aerial Images of Riverscapes for Semantic Segmentation

Saeid Shamsaliei, Odd Erik Gundersen*, Jo H. Halleraker†, Knut Alfredsen
Norwegian University of Science and Technology, Trondheim, Norway
{saeid.shamsaliei, odderik, knut.alfredsen, jo.halleraker}@ntnu.no

Anders Foldvik
Norwegian Institute for Nature Research, Trondheim, Norway
anders.foldvik@nina.no

Abstract

1 Accurate and reliable semantic segmentation of historical aerial images of land-
2 scapes is crucial for tracking, analyzing, and understanding land use over time.
3 Development of such models is challenging due to the lack of annotated datasets of
4 historical images. We introduce HAIR, the first dataset for semantic segmentation
5 of historical aerial imagery for land cover to address this issue. The dataset contains
6 high resolution, grayscale images of riverscapes with a resolution of 20 cm per
7 pixel, captured from 1947 to 1998. The images of this large-scale dataset are anno-
8 tated into six land types in meticulous detail by domain experts. We benchmark
9 state-of-the-art semantic segmentation models and present both quantitative and
10 qualitative results on in-distribution and out-of-distribution test sets. Our baseline
11 experiments show that pre-training on a recent high-resolution satellite image
12 dataset that is converted to grayscale does not improve performance. They also
13 show that state-of-the-art models do not generalize well on out-of-distribution data.
14 Finally, we characterize four challenges facing the segmentation of historical aerial
15 images, including HAIR, and by this hope to spur interest in developing models
16 that generalize well on historical images to support temporal analysis of land use.

17 1 Introduction

18 Rivers are central for the human condition. Early civilizations were build in river valleys, and today
19 they remain important, i.e. for fresh drinking water and livelihood, such as fishing, agriculture and
20 power production. However, this importance comes at a cost. Riverscapes are under pressure from
21 human development, and this challenges the biodiversity and hydromorphology around the rivers.
22 Given that rivers are home to some of the most diverse and endangered wildlife on Earth, the problem
23 is especially severe. Since 1900 the human population has grown from 1.65 Billion to 7.9 Billion in
24 2022 [51]. The growth has lead to an increasing pressure on all ecosystems on the Earth, including
25 riverscapes. UN has declared the decade from 2021 to 2030 as the UN Decade of Restoration, and it
26 is described as "a rallying call for the protection and revival of ecosystems all around the world".

*Aneo AI Research, Trondheim, Norway

†Norwegian Environment Agency, Trondheim, Norway



Figure 1: Ambiguous areas are marked in red. *Left* and *middle* show cases where gravel and human construction can be confused while *right* is an example where vegetation and water could be confused.

27 To understand the impact of human development in a region, it is crucial to understand the state of the
 28 landscape from a time when the world’s population was significantly smaller. While satellite data has
 29 facilitated the production of land use and land cover maps since the 1990s [40], aerial images have
 30 been systematically captured since the early 1900s in some parts of the world. However, the majority
 31 of these aerial images were captured using an analog film-based camera before 2005. The captured
 32 films were then scanned and converted into a digital format [9]. Additionally, for aerial images
 33 captured before the 2000s, only panchromatic (grayscale) historic photographs are available [27].
 34 Historical aerial images have the potential to be used as a valuable data source for understanding,
 35 monitoring, and analyzing land use over time. However, using aerial images for these purposes
 36 requires automatic and reliable mapping of the aerial grayscale images into desired habitats [44].
 37 Image recognition, specifically semantic segmentation, can be utilized for this mapping but this relies
 38 on the availability of a large dataset of historical images annotated into desired the habitats.

39 Datasets of historical aerial images have four characteristics that impacts semantic segmentation:
 40 1) camera technology has advanced significantly over time, which results in varying image quality
 41 based on when images are captured, 2) lightning conditions that are influenced by factors such as the
 42 time of the day and the airplane’s direction during the capture, 3) class imbalance, as some classes
 43 are underrepresented due to the nature of the aerial images, and 4) grayscale, which means that they
 44 carry less information than satellite images that include RGB channels and sometimes additional
 45 infra-red channels as well.

46 **Contributions:** Our contributions are threefold. First, we release HAIR, the first dataset of high-
 47 resolution, historical aerial images with high-quality annotations of riverscapes made by experts. The
 48 dataset is released under the CC BY-SA 4.0 license³ and contains roughly 8.72 billion annotated
 49 pixels surpassing widely recognized land cover datasets, such as DeepGlobe [22] and Inria [41].
 50 Second, we present a benchmark of state-of-the-art semantic segmentation models to provide as
 51 baselines for future work. Third, experiments show that pre-training on high-resolution satellite
 52 image that are converted to grayscale does not improve performance and that state-of-the-art models
 53 do not generalize well to out-of-distribution data.

54 2 Related work

55 The datasets, LandCover.ai, DeepGlobe and Agriculture-Vision, which are summarized in Table 1 are
 56 most similar to HAIR given that they all contain images of natural landscapes with high resolution.
 57 Long et al.[39] give an overview of many other datasets. LandCover.ai is a semantic segmentation
 58 dataset of aerial images from rural areas across Poland with resolutions between 25 to 50 cm per
 59 pixel and contains 5 classes. Agriculture-Vision [21] is an aerial image dataset for pattern analysis of
 60 agricultural lands in US with nine classes and 10 cm per pixel resolution. DeepGlobe is a Satellite

³<https://creativecommons.org/licenses/>

Table 1: Comparison of HAIR and similar natural landscapes datasets for semantic segmentation. The resolution unit is (meters per pixel).

Dataset	#Classes	#Images	Resolution	#Channels	#Pixel	Size
DeepGlobe Land Cover [22]	7	1147	0.5	RGB	$6.87 * 10^9$	2448x2448
LandCover.ai [5]	3	41	0.25,0.5	RGB	$2.98 * 10^9$	9000x9500;4200x4700
Agriculture-Vision [21]	9	94986	0.1,0.15,0.2	RGB+NIR	$2.49 * 10^{10}$	512x512
HAIR	6	178	0.2	Grayscale	$8.72 * 10^9$	8000x6000;6400x4800;16000x12000

61 Image Understanding Challenge with the three challenges: road extraction, building detection and
 62 land cover classification. DeepGlobe has high resolution images of 50 cm per pixel from India,
 63 Indonesia and Thailand and has 6 classes. Ratajczak et al. [49] introduce a historical aerial image
 64 dataset with grayscale images. They formulate the problem as a classification task where smaller
 65 image patches are given one class labels. However, the low resolution limits the usefulness of the
 66 data for studying land use evolution [13, 25].

67 HAIR has some key differences from other datasets. First, unlike most datasets that consist of recent
 68 aerial and satellite images with at least three channels (RGB), images are exclusively grayscale, as
 69 this is the nature of historical aerial images. Color information could help the segmentation model
 70 to distinguish two otherwise similar classes. In the absence of color information, models must
 71 solely rely on the texture and context of images. However, capturing the global context becomes
 72 challenging due to the high resolution of HAIR images and the limited memory capacity of current
 73 GPUs. Having high resolution images might increase the complexity of semantic segmentation
 74 and hinder the effectiveness of pretraining on common low resolution landscape datasets. However,
 75 for many important applications, such as understanding of river habitat [12] and their evolution
 76 through time [44], ice morphology [3] and fish ecology [59], it is crucial to have high resolution
 77 images. By providing high resolution images, we hope to direct further research into the challenges
 78 related to the context of historical grayscale images. Figure 8 in the appendix provides a visual
 79 comparison of HAIR aerial images and Sentinel-1 SAR images. As it can be seen from the figure,
 80 these SAR images have low resolution which make them undesirable for studying the evolution of
 81 the rivers. Additionally, Sentinel-1 images are only available after 2014 and cannot be used to study
 82 the historical state of the landscape. Second, to the best of our knowledge, HAIR is the only dataset
 83 with *gravel* class. This class is critical in analyzing the evolution of riverscapes and their ecosystems
 84 [4, 30]. Gravel serves as the exclusive habitat for certain insect and plant species, making it crucial to
 85 monitor changes in this class for biodiversity tracking [45]. Finally, even though grayscale images
 86 are more difficult for humans to annotate, as we rely on color information when interpreting images,
 87 annotations in HAIR are finer than in previous works. For example, small vegetation on top of gravel
 88 bars is annotated. These details help to develop models that are useful in ecological applications.
 89 Furthermore, the detailed annotations are made by experts. As illustrated by Table 1, the size of
 90 HAIR exceeds that of many widely recognized datasets in the field.

91 Since the introduction of FCN [38], the first end-to-end deep learning semantic segmentation model,
 92 numerous studies have proposed CNN-based architectures to enhance performance. These archi-
 93 tectures include U-Net [50], ParseNet [36], PSPNet [67], DeepLab [15, 16, 17, 18], and HRNet
 94 [61]. Recently, there has been a growing interest in pure transformer-based architectures inspired
 95 by the success of Visual Transformers [24]. For instance, Segmentor [57], Segformer [62], and
 96 Swin-Unet [11] are pure transformer-based architectures. Additionally, some studies propose models
 97 combining transformers and CNNs, like TransUNet [14]. Minaee et al. [42], provide an overview
 98 of the segmentation models. In remote sensing, many models have been developed for datasets like
 99 DeepGlobe and LandCover.ai. Examples include FPN [54] with ResNet50 [31] as encoder and spatial
 100 dropout, NU-Net [53], and DIResUNet [46]. Some proposed models, such as GLNet [20] and MagNet
 101 [32], leverage both downsampled and patched input images to capture the global context. MagNet
 102 is currently state-of-the-art in DeepGlobe land cover classification. The development of models in the
 103 field has predominantly focused on datasets containing images with three or more channels, so the
 104 research of semantic segmentation for grayscale images have been relatively limited.

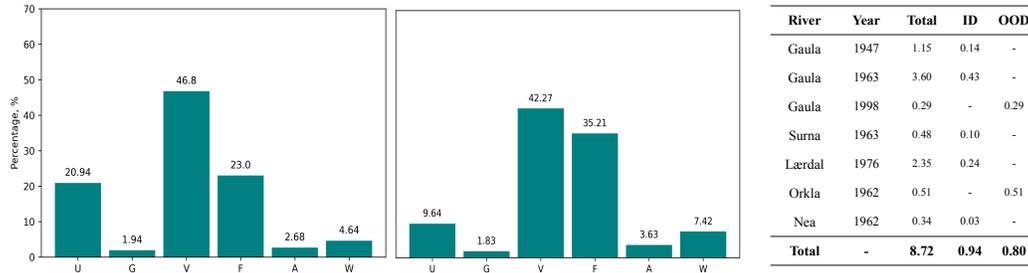


Figure 2: Diagrams: the distribution of classes in the dataset (left) and the OOD test set. Table: overview of the number of pixels taken from the different rivers (Rivers) in the different years (Year) and how many of these are in the ID and OOD test sets (in billion).

105 3 Dataset description

106 HAIR is comprised of 8.72 billion pixels spread over 178 annotated images from the five rivers
 107 Nea, Orkla, Surna, Gaula and Lærdal in Norway. The images come in one of three different sizes
 108 (8000×6000 , 6400×4800 or 16000×12000 pixels), with a resolution of 20 cm per pixel. High
 109 quality annotations are made by experts, which limits the number of images that can be annotated;
 110 quality annotations take long time to make and they are made by experts that is a limited resource.
 111 This effort is inspired by the data-centric movement⁴.

112 Figure 2 provides an overview of the dataset. The dataset contains images from 1947, 1962, 1963,
 113 1978 and 1998 redand therefore, all of them are panchromatic (grayscale) since aerial images before
 114 the 2000 were captured in grayscale [27]. These images are selected so that we have both spatial
 115 and temporal overlap in the dataset. By spatial overlap, we mean when two images capture the same
 116 geographical area in different years, and by temporal overlap, we mean when two images are captured
 117 the same year (but not necessarily over the same geographical area). The images of river Gaula
 118 represent the spatial overlap, and the river Surna represents the temporal overlap with Gaula images
 119 from 1963.

120 The test sets have been designed to test whether both spatial and temporal characteristics are learned
 121 and can be generalized. Approximately 10% of the images from each river in the training set were
 122 chosen randomly for the ID test set, and the rest of images were used for training and validation. The
 123 OOD test set consists of two different rivers Gaula and Orkla. The OOD images of river Gaula are
 124 captured in 1998, which is a different time period than those found in the training set and translates
 125 to 51 years of camera improvement compared to the images from 1947 found in the training set. The
 126 OOD pictures of Orkla are captured in 1962, which are close in time to the Nea 1962, Surna 1963
 127 and Gaula 1963. However, the Orkla images cover a completely different spatial area. Out of the 178
 128 large images, 20 are in the ID test set and nine in the OOD test set. Note that there is a spatial overlap
 129 for areas in different time periods but not for the same period. All images from the same time period
 130 are of different spatial areas.

131 **Data Acquisition:** HAIR consists of images from historical aerial photos used to develop the digital
 132 orthophoto covering the whole of Norway. The images can be accessed through a database of aerial
 133 imagery of the mainland of Norway (www.norgebilder.no) covering both recent and historic photos
 134 that is maintained by the Norwegian mapping authority. The earliest images are from 1937 and the
 135 most recent are from 2022. A large backlog of historic aerial photos of the whole of Norway are
 136 being digitized and will be shared in the database continuously after being completed. All images are
 137 georeferenced in the database. The images in HAIR are projected into EUREF89-UTM33N.

138 The aerial photos are taken at different times in the summer season when there is no snow, but
 139 also during varying periods of the vegetation season. The dataset includes a wide variety of optical

⁴Workshop on Data-centric AI at NeurIPS 2021: <https://datacentricai.org>

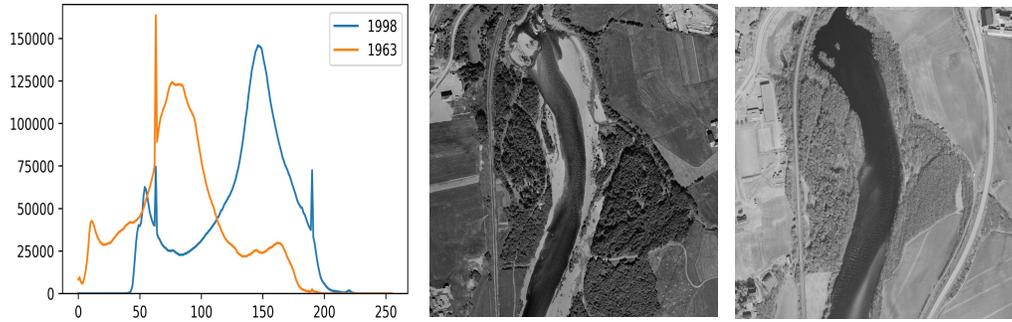


Figure 3: Intensity distribution of images covering the same section of Gaula from 1963 (left) and 1998 (right). It is calculated as the number of pixels with a given color (256 gray levels).

140 conditions including different shadow lengths, angles of sunlight and saturation. As mentioned by
 141 Boguszewski et al. [5], one would expect that this would make the dataset robust. However, the
 142 results on the OOD test data shows that it is not exactly the case.

143 **Annotation:** Traditionally, the most common annotation tools for this purpose is GIS software with
 144 polygon editing, such as QGIS [47]. However, polygon editing lacks the precision required for
 145 the high quality labels. For HAIR, annotations were made manually by the experts using Adobe
 146 Photoshop on an iPad using a pen, as this enables detailed annotations. Each large image of mainly
 147 8000×6000 pixels was loaded as a layer to Adobe Photoshop. Annotations were done on top of
 148 the source image in a layer of its own. A specific color were assigned to each class, and the domain
 149 experts colored the six classes using the corresponding colors. Adobe Photoshop provided many
 150 tools that help facilitate the annotation. Magic Wand was, for example, used as a selection tool for
 151 most of the roads, and the Marching Ants algorithm [60] was used to modify the edges of the objects.

152 The experts followed a common procedure. Areas that were considered ambiguous by individual
 153 experts were discussed by the group to reduce the noise in the annotation. The annotation has been
 154 taken very seriously and is considered to be of high quality, although ambiguous examples can
 155 without doubt be found in such a large dataset. Figure 1 illustrates three out of the many different
 156 cases that were discussed by the experts.

157 **Classes:** We annotated the images with six different classes that can help understand the human
 158 pressure on river biodiversity and hydromorphology. These six classes are chosen pragmatically
 159 based on ease of manual annotation versus value of analyzing their change over time. More classes
 160 could have been added with high value, but these classes would have been even smaller than the
 161 gravel class with an even higher potential for misclassification.

162 **Water (W):** Water covered areas (not restricted to river).

163 **Gravel (G):** Gravel bars and point bars in the river – vegetation free.

164 **Farmland (F):** Farmland and cultivated land in the river corridor.

165 **Vegetation (V):** Forest and other vegetated areas in the riparian corridor.

166 **Anthropogenic (A):** Anthropogenic structures like houses and roads.

167 **Unknown (U):** Only areas that do not contain any aerial images are labeled as “unknown”.

168 **Statistics for HAIR:** The distribution of different classes in HAIR is shown in Figure 2 and indicates
 169 an imbalanced dataset where the two classes gravel and anthropogenic are underrepresented. The
 170 most common classes, vegetation and forest, cover 69.8% of the images. Gravel, which is the most
 171 important class besides water for the analysis, covers only 1.9% of the images while water covers
 172 4.6%. Figure 3 shows a comparison of the the intensity distribution of images covering the same
 173 section of river Gaula from the years 1963 and 1998. The intensity distribution is calculated as the
 174 number of pixels with a given pixel color, where pixel color is one out of 256 gray levels.

175 4 Two Benchmark Tasks

176 We select U-Net [48], FPN [55], DeepLabV3+ [19], HRNet [61], Magnet [32] and Swin-Unet [11]
177 as baselines to get a broad set of models for the benchmark. The selected models are evaluated on
178 both the ID and OOD test sets.

179 4.1 Method

180 The segmentation of grayscale images is under-studied primarily due to the lack of large-scale
181 datasets. In this benchmark, we try to shed some lights into greyscale segmentation of aerial images
182 by testing a set of approaches to provide insights into these. One promising approach is to employ
183 transfer learning, where a common dataset of colored images is converted to grayscale and used for
184 pre-training selected models. To investigate that, we convert the widely used DeepGlobe dataset to
185 grayscale and employ it to pre-train the baselines. We evaluate two different scenarios for all models:
186 1) all models are trained on HAIR only, and 2) models are first pre-trained on grayscale DeepGlobe
187 dataset and then fine-tuned on HAIR.

188 To convert the images of the DeepGlobe dataset into grayscale, we use the luminance method [33],
189 which is widely used in computer vision [6] and is implemented in several image processing software
190 and libraries such as OpenCV [8]. In luminance method, the function $\mathcal{G}_{luminance}$ is defined as:
191 $\mathcal{G}_{luminance} \leftarrow 0.299 \cdot R + 0.587 \cdot G + 0.114 \cdot B$, where R , G and B represents the red, green
192 and blue channels of the image⁵. Additionally for all the trainings, to leverage the strong feature
193 extraction of encoders pre-trained on RGB Imagenet, we replicate the grayscale channel into the R,G,
194 and B channels, transforming our single-channel input into a standard three-channel image format.
195 ResNet50 pre-trained on ImageNet is used as the encoder for most models, except for HRNet and
196 Swin-Unet that are not pre-trained.

197 The hyperparameters used for training MagNet are the same as in [32]. For other models, we
198 used the Hyperband algorithm [35] with a maximum of 100 epochs and 11 hybrid iteration to
199 find the hyperparameters. Also, validation accuracy was used as the objective. More detail of the
200 hyperparameters for each model is provided in the appendix. The code used to run the experiments
201 are available in <https://github.com/SaeidShamsaliei/HAIR>.

202 4.2 Results

203 We calculate performance using MIOU of the experiments in the same way as Huynh et al. [32].
204 Results are presented in Tables 2 and 3 as mean \pm standard deviation of model performance. Variation
205 is introduced through training five seed replicates for each model. A seed replicate is a model trained
206 using different seeds for the pseudo-random generator [7]. The performance is reported as the mean
207 and standard deviation of the MIOU for the set of five seed replicates. As can be seen, DeepLabV3+
208 achieves the best performance on both datasets. However, it does not perform considerably better
209 than the other models. Gravel is the class that all models struggle the most to identify, and forest
210 is the class for which all models achieve their highest score. For the OOD test set, gravel is still
211 the hardest class to predict for all models, while it is relatively easier for the models to predict the
212 vegetation and farmland.

213 Pre-training on the converted DeepGlobe dataset led to a slight decrease in performance for all models
214 except HRNet and Swin-Unet, which do not perform well in general. Although the limited size of
215 the DeepGlobe dataset could be a factor, the results highlight the need for developing large scale
216 historical grayscale aerial imagery datasets. Another factor that might be a cause of the comparatively
217 poorer performance is the luminance method for converting DeepGlobe to grayscale. Kanan and
218 Cottrell [33] report that the conversion method affects the performance of downstream tasks, which
219 could easily be the case for deep learning as well as pre-processing is known to affect reproducibility
220 of deep learning methods [26, 28].

⁵The method is used in many papers such as [6], however, in these works the name of method is not mentioned.
We got the name from [33]. The coefficients are the central component of the method.

Table 2: The MIoU prediction of different semantic segmentation architectures on the ID test set. MIoU is mean \pm standard deviation over five seed replicates runs

Models	Backbone	Not pre-trained on grayscale DeepGlobe					MIoU	Pre-trained on grayscale DeepGlobe					MIoU
		IoU per class						IoU per class					
		G	V	F	A	W		G	V	F	A	W	
U-Net	ResNet50	60.81	88.02	85.31	65.70	85.54	77.07 \pm 00.48	58.19	87.85	85.01	64.74	84.67	76.10 \pm 00.72
FPN	ResNet50	60.83	87.84	85.11	65.59	82.78	76.43 \pm 01.84	59.40	87.40	84.75	64.70	84.05	76.06 \pm 00.57
DeepLabV3+	ResNet50	61.07	87.92	86.23	65.43	86.19	77.37 \pm 01.02	60.63	87.54	85.14	63.37	84.76	76.28 \pm 01.26
HRNet	HRNetV2-W18	52.56	87.17	82.20	62.97	77.70	72.52 \pm 00.97	56.08	87.04	82.15	64.34	77.87	73.50 \pm 00.79
Swin-Unet	Swin Transformer	43.86	83.50	76.32	47.03	60.35	62.21 \pm 00.40	44.41	83.43	76.31	49.79	61.38	63.06 \pm 01.52
MagNet	FPN	51.75	77.55	76.55	57.75	62.51	65.22 \pm 01.39	44.21	72.53	70.54	44.96	49.36	56.32 \pm 01.12

Table 3: The MIoU prediction of different semantic segmentation architectures on the OOD test set. MIoU is mean \pm standard deviation over five seed replicates.

Models	Backbone	Not pre-trained on grayscale DeepGlobe					MIoU	Pre-trained on grayscale DeepGlobe					MIoU
		IoU per class						IoU per class					
		G	V	F	A	W		G	V	F	A	W	
U-Net	ResNet50	32.09	78.38	75.98	61.20	73.20	64.17 \pm 00.80	30.19	76.49	66.35	51.58	75.86	60.10 \pm 01.03
FPN	ResNet50	29.24	72.43	72.73	56.39	69.78	60.11 \pm 01.51	32.02	76.29	69.42	57.22	63.05	59.60 \pm 01.64
DeepLabV3+	ResNet50	30.55	76.93	73.83	61.57	78.57	64.29 \pm 02.66	32.27	72.93	75.43	57.43	75.34	62.68 \pm 02.81
HRNet	HRNetV2-W18	22.27	67.03	60.27	50.00	68.17	53.44 \pm 02.46	22.09	66.84	55.55	50.95	62.65	51.61 \pm 02.63
Swin-Unet	Swin Transformer	14.29	62.43	50.99	35.12	50.86	42.74 \pm 00.84	14.40	62.71	51.02	35.58	52.09	43.16 \pm 01.82
MagNet	FPN	31.65	64.33	74.00	54.28	61.01	57.06 \pm 03.84	21.01	53.59	54.31	38.88	40.09	41.58 \pm 04.51

Table 4: The MIoU of the DeepLabV3+ on the two different rivers in the OOD test set and Gaula 1963 in ID test set.

River	Year	MIoU	
		Not pre-trained	Pre-trained
Orkla	1962	69.64 \pm 03.06	69.17 \pm 01.58
Gaula	1998	61.62 \pm 03.21	59.43 \pm 03.70
Gaula	1963	79.83 \pm 01.19	80.25 \pm 00.26

221 The performance of MagNet experiences a relatively larger decline, compared to other models, when
 222 pre-trained on the grayscale DeepGlobe dataset. The combination of model’s complex architecture
 223 and relatively small scale of DeepGlobe used for pre-training, seems to affect the performance to
 224 a larger extent. The backbones of HRNet and Swin-Unet were not pre-trained on any large-scale
 225 dataset to help with extracting the low level features, like ImageNet. This could be a contributing
 226 factor to their relatively lower performance.

227 All models perform substantially worse on the OOD test set compared to the ID test set. DeepLabV3+
 228 generalizes better than the other models, but only slightly better than U-Net when not pre-trained.
 229 Table 4 shows how DeepLabV3+ performs on the two different types of OOD data and that it
 230 generalizes better to data from the same time-period where the same or a similar camera technology
 231 is used than for the images captured more than 35 years later. The better performance of DeepLabV3+
 232 on OOD might be explained by how Atrous Spatial Pyramid Pooling (ASPP) mechanism fuses the
 233 extracted features from the input image. Since ASPP increases the field of view of the model, it can
 234 efficiently access a larger context of the input image.

235 4.3 Discussion

236 Here we present a short discussion of how the four characteristics of historical aerial images can
 237 affect the results.

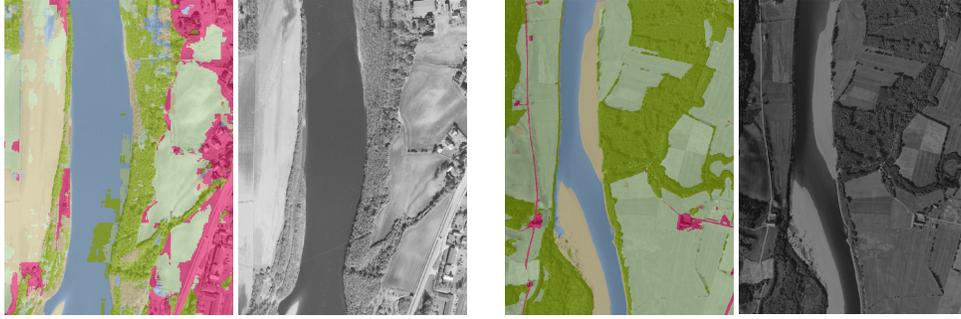


Figure 4: Comparison of predictions for Gaula 1998 (left) and 1963 (right).

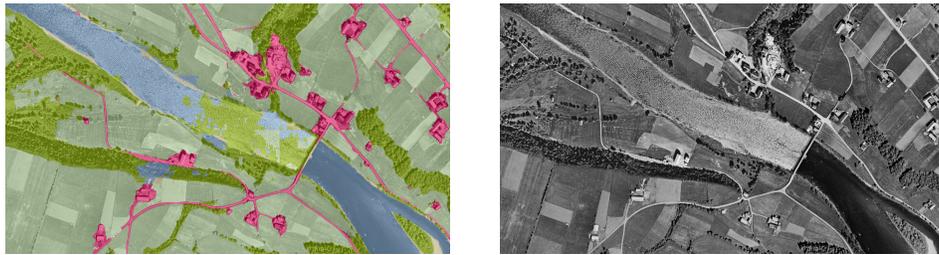


Figure 5: Different light conditions make accurate prediction difficult.

238 **Characteristic #1 - Camera technology:** The results on the OOD test set indicate that camera
 239 technology is an issue. The performance of the models drop considerably for the images of Gaula
 240 1998 compared to Orkla 1962. The Gaula 1998 were captured 35 years after the Orkla 1962 images.
 241 However, DeepLabV3+ performs much better when tested on the ID images of Gaula from 1963, as
 242 can be seen in Table 4. This indicates that it is not river Gaula itself that is challenging, but that it is
 243 the development of camera technology that is the cause of the performance drop. While intensity
 244 differs between the river in the ID test set compared to the older images in the OOD dataset, as
 245 seen in Figures 3 and 4, light intensity alone should not be a problem as contrast and brightness
 246 were changed randomly as part of the training, as described in methodology. More research on data
 247 augmentation methods designed specifically for grayscale aerial images might alleviate this issue.
 248 Examples include modification of brightness and contrast during the training [52, 64], blurring filters,
 249 histogram transformations, special optical and lens distortion [37], and generative models [56].

250 **Characteristic #2 - Lighting conditions:** Light conditions make the prediction more challenging.
 251 In Figure 5, the reflections in the shimmering section represent an unusual visual effect on the water.
 252 Because this effect happens only when the light from the sun hits the water and is reflected in the
 253 camera, there are relatively few examples of such effects in the training data and little data for the
 254 model to learn the pattern. Challenge #1 covers this by suggesting that more research on online
 255 augmentations and generative models, designed specifically for historical grayscale aerial images,
 256 could mitigate this issue.

257 **Characteristic #3 - Class imbalance:** Figure 2 shows the imbalanced class distribution, and the
 258 results presented in Table 2 and Table 3 indicate that the models have the worst performance on
 259 smallest classes. The class distribution clearly poses a problem. More research on effective methods
 260 to mitigate the class imbalance of high resolution aerial images, such as more efficient sampling
 261 methods [58] and more effective loss functions [23] could help overcoming this challenge.

262 **Characteristic #4 - Grayscale:** Texture and context become the main sources of information
 263 for grayscale images. Land covers that are easily distinguishable when having color information
 264 become hard to distinguish in grayscale. Figure 6 shows several cases where texture alone makes

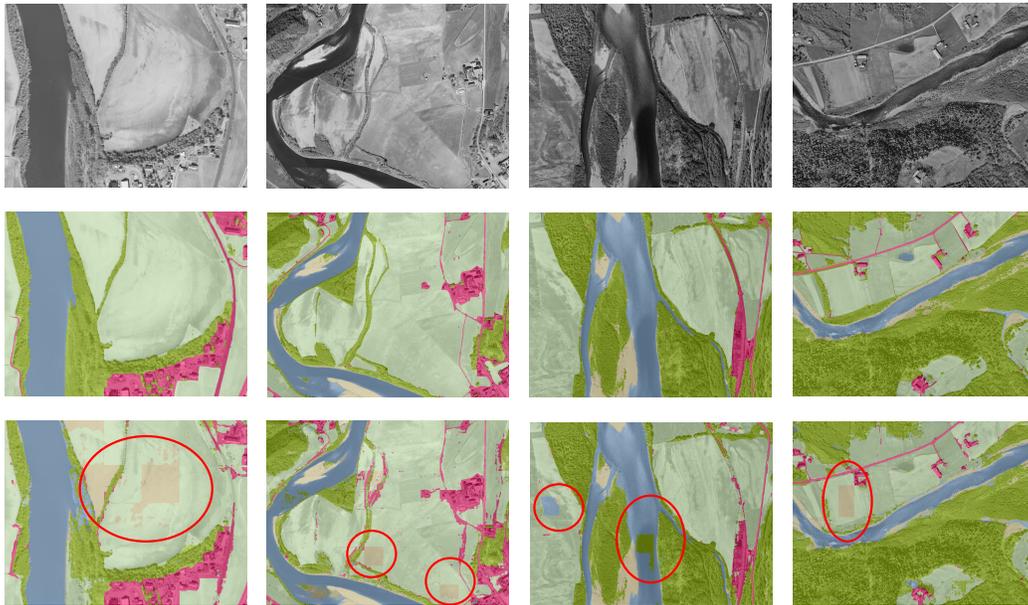


Figure 6: Source images (top row), labels (middle-row) and best predictions (bottom row).

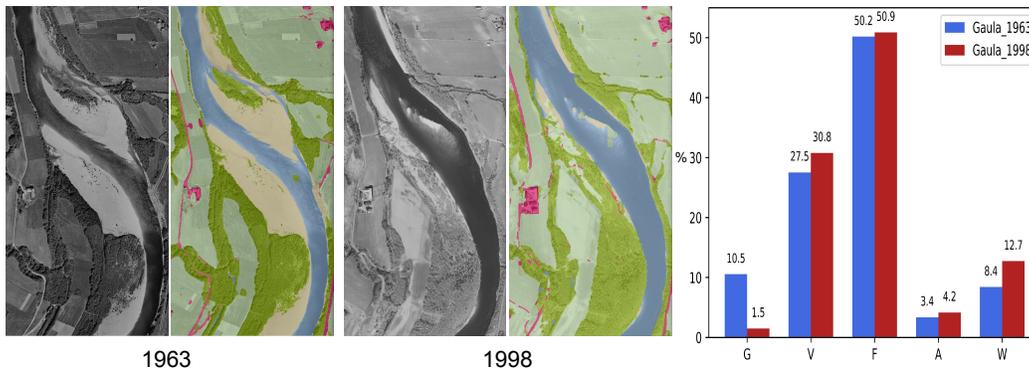


Figure 7: A section of River Gaula from 1963 and 1998 with distribution of classes in percentage.

265 the prediction hard. The deep learning segmentation methods typically divide the larger images into
 266 smaller patches and run these through the deep nets. Context is the underlying issue when the method
 267 miss-classifies two patches that clearly are of the same class. Examples of this can easily be spotted
 268 in Figure 6 where patches of the wrong class can be seen very clearly in larger regions that are mostly
 269 predicted correctly. These errors happen when there are examples of patches in the test set that have
 270 similar texture but different class than patches in the training data. It is easy to imagine that color
 271 information would help solve these cases. We believe that continued development of architectures
 272 similar to MagNet and GLNet that use different scales of images to provide the context to smaller
 273 patches could mitigate this issue.

274 5 Future Work and Conclusions

275 Manual analyses of aerial images of riverscapes is the standard method for understanding the
276 development of individual rivers [29, 66], and they include spatio-temporal dynamics of parameters
277 such as river width, bank cover, sinuosity and other geomorphological features. Figure 7 shows raw
278 images and segmentations of the years 1963 and 1998 as well as a bar chart that enables us to perform
279 a basic version of such an analysis on the current dataset. The images show that the reduction in
280 gravel is substantial and so is the increase in water. The bar chart confirms not only this but that all
281 classes except for gravel increase. The analysis could indicate that the conditions for the plants and
282 insects living in gravel has worsened and that their existence is threatened. However, as these images
283 could have been taken at different seasons we cannot draw this conclusion so easily. Hence, we need
284 to semantically segment images of River Gaula for more years and the findings for one river must be
285 compared to findings in other rivers to draw more certain conclusions. This is an example for a large
286 scale analysis. However, such large scale analyses covering multiple rivers have not been done before.
287 To achieve this, it is essential to automate the land cover classification of historical aerial images [2],
288 and the segmentation model must generalize well to OOD data. Given that such land area cover maps
289 made by segmentation models exist for multiple rivers over a long period of time (1940’s through
290 2020’s), a completely new type of analysis can be done.

291 In this paper, we present a free and open dataset for advancing semantic segmentation of historic, high
292 resolution, grayscale aerial images of land cover. We provide baselines for selected state-of-the-art
293 deep learning models on the two benchmark tasks and show several issues with applying state-of-
294 the-art models developed for color images on grayscale images. It is clear that the training is not
295 robust and that the models do not generalize well beyond the training data. Pre-training on a recent
296 dataset of converted satellite images is shown to not improve the performance of the models, which
297 highlights the importance of generating large scale datasets of historical aerial images. Additionally,
298 more research can be done to improve the effectiveness of pre-training on low resolution RGB
299 datasets to improve semantic segmentation of high resolution grayscale images. For example, one
300 can explore the impact of converting the aerial resolution of HAIR images (e.g., from 20 cm per pixel
301 to 50 cm per pixel as in DeepGlobe images) to determine whether the conversion would improve
302 when using models pre-trained on lower resolution datasets. Finally, we identify and show how four
303 characteristics of historical aerial images negatively affect the performance of the tested models.

304 Grayscale aerial images are a potentially important data source that has not received much attention.
305 Our hope is that bringing attention to this data source will support the development of generalizable
306 models that can overcome the challenges of the dataset and be used to perform large scale land-
307 mapping of historical grayscale aerial images accurately. This will enable longitudinal studies
308 and quantified analyses of land use, and support large scale investigations into which practices are
309 sustainable and which are not. Such analyses require far too much manual work of experts to be
310 feasible without automatic solutions.

311 Acknowledgments and Disclosure of Funding

312 Thanks to Statkart-Geovekst for supporting the sharing of the enhanced dataset publicly under the
313 Creative Commons 4.0 SA-BY license. J. H. H. was funded by the Norwegian Research Council
314 through the OFFPHD program, Grant No: 289725 – Ecosystem-based management.

315 References

- 316 [1] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado,
317 Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey
318 Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg,
319 Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens,
320 Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda
321 Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng.

- 322 TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from
323 tensorflow.org.
- 324 [2] Knut Alfredsen, Arild Dalsgård, Saeid Shamsaliei, Jo Halvard Halleraker, and Odd Erik Gundersen.
325 Towards an automatic characterization of riverscape development by deep learning. *River Research and*
326 *Applications*, 38(4):810–816, December 2021.
- 327 [3] Knut Alfredsen, Christian Haas, Jeffrey A Tuhtan, and Peggy Zinke. Brief communication: Mapping river
328 ice using drones and structure from motion. *The Cryosphere*, 12(2):627–633, 2018.
- 329 [4] Bjørn T Barlaup, Sven Erik Gabrielsen, Helge Skoglund, and Tore Wiers. Addition of spawning gravel—a
330 means to restore spawning habitat of atlantic salmon (*salmo salar* l.), and anadromous and resident brown
331 trout (*salmo trutta* l.) in regulated rivers. *River Research and Applications*, 24(5):543–550, 2008.
- 332 [5] Adrian Boguszewski, Dominik Batorski, Natalia Ziemba-Jankowska, Tomasz Dziedzic, and Anna Zam-
333 brzycka. Landcover.ai: Dataset for automatic mapping of buildings, woodlands, water and roads from
334 aerial imagery. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*
335 *(CVPR) Workshops*, pages 1102–1110, June 2021.
- 336 [6] Anna Bosch, Andrew Zisserman, and Xavier Munoz. Image classification using random forests and ferns.
337 In *2007 IEEE 11th international conference on computer vision*, pages 1–8. Ieee, 2007.
- 338 [7] Xavier Bouthillier, César Laurent, and Pascal Vincent. Unreproducible research is reproducible. In
339 *International Conference on Machine Learning*, pages 725–734. PMLR, 2019.
- 340 [8] Gary Bradski. The opencv library. *Dr. Dobb’s Journal: Software Tools for the Professional Programmer*,
341 25(11):120–123, 2000.
- 342 [9] Hardy Buller. personal communication with Norwegian Mapping Authority, Aug. 14 2023.
- 343 [10] Alexander Buslaev, Vladimir I. Iglovikov, Eugene Khvedchenya, Alex Parinov, Mikhail Druzhinin, and
344 Alexandr A. Kalinin. Albumentations: Fast and flexible image augmentations. *Information*, 11(2), 2020.
- 345 [11] Hu Cao, Yueyue Wang, Joy Chen, Dongsheng Jiang, Xiaopeng Zhang, Qi Tian, and Manning Wang.
346 Swin-unet: Unet-like pure transformer for medical image segmentation. *arXiv preprint arXiv:2105.05537*,
347 2021.
- 348 [12] Patrice E Carbonneau, Stephen J Dugdale, Toby P Breckon, James T Dietrich, Mark A Fonstad, Hi-
349 toshi Miyamoto, and Amy S Woodget. Adopting deep learning methods for airborne rgb fluvial scene
350 classification. *Remote Sensing of Environment*, 251:112107, 2020.
- 351 [13] Patrice E. Carbonneau and Hervé Piégay. Introduction: The growing use of imagery in fundamental and
352 applied river sciences, August 2012.
- 353 [14] Jieneng Chen, Yongyi Lu, Qihang Yu, Xiangde Luo, Ehsan Adeli, Yan Wang, Le Lu, Alan L. Yuille, and
354 Yuyin Zhou. Transunet: Transformers make strong encoders for medical image segmentation. *CoRR*,
355 abs/2102.04306, 2021.
- 356 [15] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Semantic
357 image segmentation with deep convolutional nets and fully connected crfs. *arXiv preprint arXiv:1412.7062*,
358 2014.
- 359 [16] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L. Yuille. Deeplab:
360 Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs,
361 2017.
- 362 [17] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution
363 for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017.
- 364 [18] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder
365 with atrous separable convolution for semantic image segmentation. In *Proceedings of the European*
366 *conference on computer vision (ECCV)*, pages 801–818, 2018.
- 367 [19] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder
368 with atrous separable convolution for semantic image segmentation. In *Proceedings of the European*
369 *conference on computer vision (ECCV)*, pages 801–818, 2018.
- 370 [20] Wuyang Chen, Ziyu Jiang, Zhangyang Wang, Kexin Cui, and Xiaoning Qian. Collaborative global-
371 local networks for memory-efficient segmentation of ultra-high resolution images. In *2019 IEEE/CVF*
372 *Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, June 2019.

- 373 [21] Mang Tik Chiu, Xingqian Xu, Yunchao Wei, Zilong Huang, Alexander G. Schwing, Robert Brunner,
374 Hrant Khachatryan, Hovnatan Karapetyan, Ivan Dozier, Greg Rose, David Wilson, Adrian Tudor, Naira
375 Hovakimyan, Thomas S. Huang, and Honghui Shi. Agriculture-vision: A[jenssen] large aerial image
376 database for agricultural pattern analysis. In *2020 IEEE/CVF Conference on Computer Vision and Pattern
377 Recognition (CVPR)*, pages 2825–2835, 2020.
- 378 [22] Ilke Demir, Krzysztof Koperski, David Lindenbaum, Guan Pang, Jing Huang, Saikat Basu, Forest Hughes,
379 Devis Tuia, and Ramesh Raskar. DeepGlobe 2018: A challenge to parse the earth through satellite images.
380 In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*. IEEE,
381 June 2018.
- 382 [23] Rongsheng Dong, Xiaoquan Pan, and Fengying Li. Denseu-net-based semantic segmentation of small
383 objects in urban remote sensing images. *IEEE Access*, 7:65347–65356, 2019.
- 384 [24] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
385 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth
386 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- 387 [25] Kurt D Fausch, Christian E Torgersen, Colden V Baxter, and Hiram W Li. Landscapes to riverscapes:
388 bridging the gap between research and conservation of stream fishes: a continuous view of the river is
389 needed to understand how processes interacting among scales set the context for stream fishes and their
390 habitat. *BioScience*, 52(6):483–498, 2002.
- 391 [26] Maurizio Ferrari Dacrema, Simone Boglio, Paolo Cremonesi, and Dietmar Jannach. A troubling analysis
392 of reproducibility and progress in recommender systems research. *ACM Transactions on Information
393 Systems (TOIS)*, 39(2):1–49, 2021.
- 394 [27] Sarah E Gergel and Monica G Turner. *Learning landscape ecology: a practical guide to concepts and
395 techniques*. Springer, 2017.
- 396 [28] Odd Erik Gundersen, Kevin Coakley, Christine Kirkpatrick, and Yolanda Gil. Sources of irreproducibility
397 in machine learning: A review, 2023.
- 398 [29] Angela M Gurnell. Channel change on the river dee meanders, 1946–1992, from the analysis of air
399 photographs. *Regulated Rivers: Research & Management: An International Journal Devoted to River
400 Research and Management*, 13(1):13–26, 1997.
- 401 [30] F Richard Hauer, Harvey Locke, Victoria J Dreitz, Mark Hebblewhite, Winsor H Lowe, Clint C Muhlfeld,
402 Cara R Nelson, Michael F Proctor, and Stewart B Rood. Gravel-bed river floodplains are the ecological
403 nexus of glaciated mountain landscapes. *Science Advances*, 2(6):e1600026, 2016.
- 404 [31] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.
405 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- 406 [32] Chuong Huynh, Anh Tuan Tran, Khoa Luu, and Minh Hoai. Progressive semantic segmentation. In *IEEE
407 Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*, pages
408 16755–16764. Computer Vision Foundation / IEEE, 2021.
- 409 [33] Christopher Kanan and Garrison W Cottrell. Color-to-grayscale: does the method matter in image
410 recognition? *PloS one*, 7(1):e29740, 2012.
- 411 [34] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and
412 Yann LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego,
413 CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- 414 [35] Lisha Li, Kevin Jamieson, Giulia DeSalvo, Afshin Rostamizadeh, and Ameet Talwalkar. Hyperband: A
415 novel bandit-based approach to hyperparameter optimization. *Journal of Machine Learning Research*,
416 18(185):1–52, 2018.
- 417 [36] Wei Liu, Andrew Rabinovich, and Alexander C Berg. Parsenet: Looking wider to see better. *arXiv preprint
418 arXiv:1506.04579*, 2015.
- 419 [37] Yan Liu, Qirui Ren, Jiahui Geng, Meng Ding, and Jiangyun Li. Efficient patch-wise semantic segmentation
420 for large-scale remote sensing images. *Sensors*, 18(10):3232, 2018.
- 421 [38] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic seg-
422 mentation. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, June
423 2015.

- 424 [39] Yang Long, Gui-Song Xia, Shengyang Li, Wen Yang, Michael Ying Yang, Xiao Xiang Zhu, Liangpei
425 Zhang, and Deren Li. Dirs: On creating benchmark datasets for remote sensing image interpretation.
426 *CoRR*, abs/2006.12485, 2020.
- 427 [40] Thomas R Loveland, Bradley C Reed, Jesslyn F Brown, Donald O Ohlen, Zhiliang Zhu, LWMJ Yang, and
428 James W Merchant. Development of a global land cover characteristics database and igbp discover from 1
429 km avhrr data. *International journal of remote sensing*, 21(6-7):1303–1330, 2000.
- 430 [41] Emmanuel Maggiori, Yuliya Tarabalka, Guillaume Charpiat, and Pierre Alliez. Can semantic labeling
431 methods generalize to any city? the inria aerial image labeling benchmark. In *IEEE International
432 Geoscience and Remote Sensing Symposium (IGARSS)*. IEEE, 2017.
- 433 [42] Shervin Minaee, Yuri Boykov, Fatih Porikli, Antonio Plaza, Nasser Kehtarnavaz, and Demetri Terzopoulos.
434 Image segmentation using deep learning: A survey. *IEEE transactions on pattern analysis and machine
435 intelligence*, 44(7):3523–3542, 2021.
- 436 [43] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
437 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang,
438 Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie
439 Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In
440 *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.
- 441 [44] Hervé Piégay, Fanny Arnaud, Barbara Belletti, Mélanie Bertrand, Simone Bizzi, Patrice Carbonneau,
442 Simon Dufour, Frédéric Liébault, Virginia Ruiz-Villanueva, and Louise Slater. Remotely sensed rivers in
443 the anthropocene: State of the art and prospects. *Earth Surface Processes and Landforms*, 45(1):157–188,
444 2020.
- 445 [45] Francesca Pilotto, Christer Nilsson, Lina E Polvi, and Brendan G McKie. First signs of macroinvertebrate
446 recovery following enhanced restoration of boreal streams used for timber floating. *Ecological Applications*,
447 28(2):587–597, 2018.
- 448 [46] Priyanka, Sravya N, Shyam Lal, J Nalini, Chintala Sudhakar Reddy, and Fabio Dell’Acqua. DIResUNet:
449 Architecture for multiclass semantic segmentation of high resolution remote sensing imagery data. *Applied
450 Intelligence*, March 2022.
- 451 [47] QGIS Development Team. *QGIS Geographic Information System*. Open Source Geospatial Foundation,
452 2009.
- 453 [48] Alexander Rakhlin, Alex Davydow, and Sergey Nikolenko. Land cover classification from satellite imagery
454 with u-net and lovász-softmax loss. In *Proceedings of the IEEE Conference on Computer Vision and
455 Pattern Recognition Workshops*, pages 262–266, 2018.
- 456 [49] Rémi Ratajczak, Carlos Fernando Crispim-Junior, Elodie Faure, Béatrice Fervers, and Laure Tougne.
457 Automatic land cover reconstruction from historical aerial images: An evaluation of features extraction
458 and classification algorithms. *IEEE Transactions on Image Processing*, 28(7):3357–3371, 2019.
- 459 [50] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical
460 image segmentation. In *International Conference on Medical image computing and computer-assisted
461 intervention*, pages 234–241. Springer, 2015.
- 462 [51] Max Roser, Hannah Ritchie, and Esteban Ortiz-Ospina. World population growth, 2013.
463 <https://ourworldindata.org/world-population-growth>.
- 464 [52] Dimitrios Sakkos, Hubert PH Shum, and Edmond SL Ho. Illumination-based data augmentation for robust
465 background subtraction. In *2019 13th International Conference on Software, Knowledge, Information
466 Management and Applications (SKIMA)*, pages 1–8. IEEE, 2019.
- 467 [53] Mohamed Samy, Karim Amer, Kareem Eissa, Mahmoud Shaker, and Mohamed ElHelw. NU-net: Deep
468 residual wide field of view convolutional neural network for semantic segmentation. In *2018 IEEE/CVF
469 Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*. IEEE, June 2018.
- 470 [54] Selim Seferbekov, Vladimir Iglovikov, Alexander Buslaev, and Alexey Shvets. Feature pyramid network
471 for multi-class land segmentation. In *2018 IEEE/CVF Conference on Computer Vision and Pattern
472 Recognition Workshops (CVPRW)*, pages 272–2723, 2018.
- 473 [55] Selim Seferbekov, Vladimir Iglovikov, Alexander Buslaev, and Alexey Shvets. Feature pyramid network
474 for multi-class land segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern
475 Recognition Workshops*, pages 272–275, 2018.

- 476 [56] Connor Shorten and Taghi M Khoshgoftaar. A survey on image data augmentation for deep learning.
477 *Journal of big data*, 6(1):1–48, 2019.
- 478 [57] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segmenter: Transformer for semantic
479 segmentation. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*. IEEE, October
480 2021.
- 481 [58] Antonio Tavera, Edoardo Arnaudo, Carlo Masone, and Barbara Caputo. Augmentation invariance and
482 adaptive sampling in semantic segmentation of agricultural aerial images. In *Proceedings of the IEEE/CVF*
483 *Conference on Computer Vision and Pattern Recognition*, pages 1656–1665, 2022.
- 484 [59] Robin L. Vannote, G. Wayne Minshall, Kenneth W. Cummins, James R. Sedell, and Colbert E. Cushing.
485 The river continuum concept. *Canadian Journal of Fisheries and Aquatic Sciences*, 37(1):130–137, January
486 1980.
- 487 [60] Alberto Viseras, Rafael Ortiz Losada, and Luis Merino. Planning with ants. *International Journal of*
488 *Advanced Robotic Systems*, 13(5):172988141666407, September 2016.
- 489 [61] Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong Mu,
490 Mingkui Tan, Xinggang Wang, et al. Deep high-resolution representation learning for visual recognition.
491 *IEEE transactions on pattern analysis and machine intelligence*, 43(10):3349–3364, 2020.
- 492 [62] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M. Alvarez, and Ping Luo. Segformer:
493 Simple and efficient design for semantic segmentation with transformers. In Marc Aurelio Ranzato, Alina
494 Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan, editors, *Advances in Neural*
495 *Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021,*
496 *NeurIPS 2021, December 6-14, 2021, virtual*, pages 12077–12090, 2021.
- 497 [63] Pavel Yakubovskiy. Segmentation models. https://github.com/qubvel/segmentation_models,
498 2019.
- 499 [64] Bo Yang, Kaiyong Xu, Hengjun Wang, and Hengwei Zhang. Random transformation of image brightness
500 for adversarial attack. *Journal of Intelligent & Fuzzy Systems*, (Preprint):1–12, 2022.
- 501 [65] Yuhui Yuan, Xiaokang Chen, Xilin Chen, and Jingdong Wang. Segmentation transformer: Object-
502 contextual representations for semantic segmentation. *arXiv preprint arXiv:1909.11065*, 2019.
- 503 [66] Luca Zanoni, Angela Gurnell, Nick Drake, and Nicola Surian. Island dynamics in a braided river from
504 analysis of historical maps and air photographs. *River Research and Applications*, 24(8):1141–1159, 2008.
- 505 [67] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing
506 network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages
507 2881–2890, 2017.

508 Checklist

- 509 1. For all authors...
- 510 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contribu-
511 tions and scope? [Yes]
- 512 (b) Did you describe the limitations of your work? [Yes] Limitations to size of dataset and bench-
513 marking methodology were discussed in the text.
- 514 (c) Did you discuss any potential negative societal impacts of your work? [Yes] We mentioned that
515 we do not see any negative social impacts to our work. On the contrary we see only potential in
516 positive change of ecosystem restoration policies.
- 517 (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
- 518 2. If you are including theoretical results...
- 519 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 520 (b) Did you include complete proofs of all theoretical results? [N/A]
- 521 3. If you ran experiments (e.g. for benchmarks)...
- 522 (a) Did you include the code, data, and instructions needed to reproduce the main experimental
523 results (either in the supplemental material or as a URL)? [Yes] Code, data and datasheet can be
524 found here: <https://folk.idi.ntnu.no/odderik/HAIR/> All above will be shared on a website and
525 through Zenodo if accepted with a DOI.

- 526 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)?
527 [Yes]
- 528 (c) Did you report error bars (e.g., with respect to the random seed after running experiments
529 multiple times)? [Yes] We did only run experiments with five random seeds. Results only for
530 comparison purposes, not main deliverable.
- 531 (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs,
532 internal cluster, or cloud provider)? [Yes] This is shared in the appendix.
- 533 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 534 (a) If your work uses existing assets, did you cite the creators? [Yes] All models used for bench-
535 marking have been cited.
- 536 (b) Did you mention the license of the assets? [Yes] We share them using CC 4.0 BY-SA
- 537 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 538 (d) Did you discuss whether and how consent was obtained from people whose data you're us-
539 ing/curating? [Yes]
- 540 (e) Did you discuss whether the data you are using/curating contains personally identifiable in-
541 formation or offensive content? [No] These are land area images, which do not contain such
542 information.
- 543 5. If you used crowdsourcing or conducted research with human subjects...
- 544 (a) Did you include the full text of instructions given to participants and screenshots, if applicable?
545 [N/A]
- 546 (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB)
547 approvals, if applicable? [N/A]
- 548 (c) Did you include the estimated hourly wage paid to participants and the total amount spent on
549 participant compensation? [N/A]

550 A Appendix

551 Include extra information in the appendix. This section will often be part of the supplemental material. Please
552 see the call on the NeurIPS website for links to additional guides on dataset publication.

- 553 1. Submission introducing new datasets must include the following in the supplementary materials:
 - 554 (a) Dataset documentation and intended uses. Recommended documentation frameworks include
555 datasheets for datasets, dataset nutrition labels, data statements for NLP, and accountability
556 frameworks.
 - 557 (b) URL to website/platform where the dataset/benchmark can be viewed and downloaded by the
558 reviewers. *URL: <https://folk.idi.ntnu.no/odderik/HAIR/hair.zip>*
 - 559 (c) Author statement that they bear all responsibility in case of violation of rights, etc., and confir-
560 mation of the data license. *We confirm that we bear all responsibility in case of violations of*
561 *rights. We confirm that the data will be shared under CC 4.0 BY-SA.*
 - 562 (d) Hosting, licensing, and maintenance plan. The choice of hosting platform is yours, as long as you
563 ensure access to the data (possibly through a curated interface) and will provide the necessary
564 maintenance. *Code is shared on GitHub and data will be shared through Zenodo. a website will*
565 *be made to host all information about the data. The website will also link to new versions of the*
566 *data set if new ones are made.*
- 567 2. To ensure accessibility, the supplementary materials for datasets must include the following:
 - 568 (a) Links to access the dataset and its metadata. This can be hidden upon submission if the
569 dataset is not yet publicly available but must be added in the camera-ready version. In
570 select cases, e.g when the data can only be released at a later date, this can be added af-
571 terward. Simulation environments should link to (open source) code repositories. *Data:*
572 *<https://folk.idi.ntnu.no/odderik/HAIR/hair.zip>, code: <https://github.com/SaeidShamsaliei/HAIR>*
 - 573 (b) The dataset itself should ideally use an open and widely used data format. Provide a detailed ex-
574 planation on how the dataset can be read. For simulation environments, use existing frameworks
575 or explain how they can be used. *The image data and annotations are stored as both tiff and png*
576 *files, which are standard formats that are supported by basically all machine learning libraries.*
 - 577 (c) Long-term preservation: It must be clear that the dataset will be available for a long time, either
578 by uploading to a data repository or by explaining how the authors themselves will ensure this.
579 *We share the data through Zenodo.org, which will store the data safely for the future in CERN's*
580 *Data Centre for as long as CERN exists.*
 - 581 (d) Explicit license: Authors must choose a license, ideally a CC license for datasets, or an open
582 source license for code (e.g. RL environments). *We have chosen Creative Commons 4.0 BY-SA.*
 - 583 (e) Add structured metadata to a dataset's meta-data page using Web standards (like schema.org and
584 DCAT): This allows it to be discovered and organized by anyone. If you use an existing data
585 repository, this is often done automatically. *This will be done.*
 - 586 (f) Highly recommended: a persistent dereferenceable identifier (e.g. a DOI minted by a data
587 repository or a prefix on identifiers.org) for datasets, or a code repository (e.g. GitHub, GitLab,...)
588 for code. If this is not possible or useful, please explain why. *DOIs are supported by zenodo.org*
589 *and we will get one when the dataset is stored there.*
- 590 3. For benchmarks, the supplementary materials must ensure that all results are easily reproducible.
591 Where possible, use a reproducibility framework such as the ML reproducibility checklist, or otherwise
592 guarantee that all results can be easily reproduced, i.e. all necessary datasets, code, and evaluation
593 procedures must be accessible and documented.
- 594 4. For papers introducing best practices in creating or curating datasets and benchmarks, the above
595 supplementary materials are not required.

596 B Runtime Environment

597 All experiments were run on one node the Yoda cluster of Computer Science department of Norwegian University
598 of Science and Technology. The information of this node is described below:

599 **GPU:** NVIDIA Tesla V100

600 **CPU:** Intel(R) Xeon-Gold 6240

601 **Number of Cores:** 18 cores @ 2.6 Ghz

602 **RAM:** 32 GiB

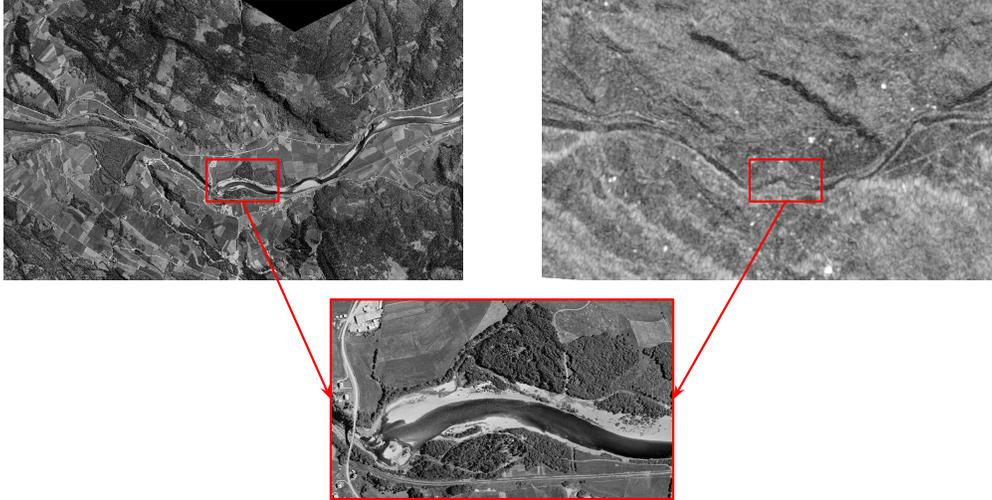


Figure 8: A visual comparison of HAIR images of Gaula captured in 1963 (top left) and SAT image of Gaula captured by Sentinel-1 in 2015 (top right). A subsection of HAIR images of Gaula 1963 is magnified to illustrate the high resolution of the HAIR images.

603 Models were implemented using Tensorflow [1] and Pytorch [43] packages. In addition, Albumentations
 604 library [10] was used for Online Augmentation, and SegmentModel [63] library was used in some of the
 605 implementations. In order to train the MagNet [32], the script provided by the paper’s github page was used.

606 C Hyperparameters of benchmark models

607 The batch size for the models is 16, except for HRNet, Swin-Unet and backbone of MagNet, which is 12, and
 608 MagNet refinement module, which is 8. Models were trained using a weighted categorical cross entropy loss
 609 function to mitigate the class imbalance of the dataset. For training the MagNet, stochastic gradient decent
 610 (SGD) is used for with a weight decay of 0.9, For the other architectures the Adam optimizer [34] is used.

611 The ReduceLROnPlateau algorithm was applied in Adam to reduce the learning rate by a factor of 0.5 if value
 612 loss did not decrease for more than 5 epochs. Except for MagNet, L2 regularization is used for convolutional
 613 layers. The FPN backbone of the MagNet is pretrained on DeepGlobe dataset and the output layer is changed
 614 to have 6 outputs instead of 7. MagNet is trained for 484 epochs while the other architectures are trained until
 615 convergence .

616 Gradient clipping with clipping value of 0.5 is applied for training the Swin-Unet. During training, images were
 617 randomly flipped and transposed, and their contrast and brightness were randomly transformed with the changing
 618 factor set to 0.1, and the probability of applying the changes was set to 20%. For training the MagNet, images
 619 were sampled as 2448×2448 px patches, and the rest are sampled as 512×512 px patches. To sample the
 620 patches, each large image is first divided into smaller patches with no overlap. Afterwards, images were flipped
 621 and patches with center pixel of gravel and water were added to the training set. The input size of MagNet used
 622 is the same as in the original paper for DeepGlobe. For the rest of the models, 512×512 is used as the input size.
 623 For inference, large images were divided into patches with the same size used to train the corresponding model,
 624 and each of these patches was used for inference. Due to low GPU memory, we did not used object-contextual
 625 representation for HRNet [65].

626 D Visualization of the Synthetic Aperture Radar (SAR) images

627 As illustrated in the Figure 8, the resolution of the SAR images from Sentinel-1 are comparably lower than the
 628 aerial images in the HAIR dataset. Additionally, Sentinel-1 started the mission in 2014, so evolution of the
 629 landscape during previous years is not captured.