

Earth and Space Science



METHOD

10.1029/2022EA002332

Key Points:

- We develop software for segmentation of geoscientific imagery with fully convolutional deep neural network models
- The software presents options for users, but relies on a reusable template that allows for rapid experimentation
- We demonstrate an example workflow with Landsat 8 imagery, and compare loss functions, model size, and model architectures

Correspondence to:

D. Buscombe,
dbuscombe@contractor.usgs.gov

Citation:

Buscombe, D., & Goldstein, E. B. (2022). A reproducible and reusable pipeline for segmentation of geoscientific imagery. *Earth and Space Science*, 9, e2022EA002332. <https://doi.org/10.1029/2022EA002332>

Received 19 MAR 2022

Accepted 23 MAY 2022

Author Contributions:

Conceptualization: D. Buscombe, E. B. Goldstein

Data curation: D. Buscombe

Formal analysis: D. Buscombe

Investigation: D. Buscombe, E. B. Goldstein

Methodology: D. Buscombe, E. B. Goldstein

Project Administration: D. Buscombe, E. B. Goldstein

Resources: D. Buscombe

Software: D. Buscombe, E. B. Goldstein

Validation: D. Buscombe, E. B. Goldstein

Visualization: D. Buscombe

Writing – original draft: D. Buscombe

Writing – review & editing: E. B. Goldstein

© 2022 The Authors. Earth and Space Science published by Wiley Periodicals LLC on behalf of American Geophysical Union.

This is an open access article under the terms of the [Creative Commons Attribution License](#), which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

A Reproducible and Reusable Pipeline for Segmentation of Geoscientific Imagery

D. Buscombe¹  and E. B. Goldstein² 

¹Marda Science, LLC, contracted to USGS Pacific Coastal and Marine Science Center, Santa Cruz, CA, USA, ²Department of Geography, Environment, and Sustainability, University of North Carolina at Greensboro, Greensboro, NC, USA

Abstract Segmentation of Earth science imagery is an increasingly common task. Among modern techniques that use Deep Learning, the UNet architecture has been shown to be a reliable for segmenting a range of imagery. We developed software—Segmentation Gym—to implement a data-model pipeline for segmentation of scientific imagery using a family of UNet models. With an existing set of imagery and labels, the software uses a single configuration file that handles data set creation, as well as model setup and model training. Key benefits of this software are (a) the focus on reproducible data set creation and modeling, and (b) the ability for quick model experimentation through changes to a configuration file. Quick experimentation permits researchers to prototype different model architectures, sizes, and adjust common hyperparameters to find a suitable model. We demonstrate the use of the software using a data set of 419 labeled Landsat-8 scenes of coastal environments and compare results across two model architectures, five model sizes, and three loss functions. This demonstration highlights that our software enables rapid, reproducible experimentation to determine optimal hyperparameters for specific data sets and research questions.

Plain Language Summary A common task for Earth scientists is to divide a satellite or aerial image into specific classes. For example, an image of the coastline might be assigned certain pixels as being water, beach, and land. In the Deep Learning world, this is called segmentation. We wrote a piece of software that helps researchers train Deep Learning models to do segmentation on all types of imagery. A major problem with making Deep Learning models is dealing with all the choices on which model to use and quickly testing many options. We have designed our code in such a way that it can easily be adjusted, and will work in many applications and for many common types of Earth science image data sets.

1. Introduction

Image segmentation has become an increasingly important tool in Earth science research (Pally & Samadi, 2022; Sun et al., 2022; Yuan et al., 2021). In recent years, Deep Learning (LeCun et al., 2015) models based on the UNet (Ronneberger et al., 2015) and the Residual UNet (Liu et al., 2019; Zhang et al., 2018) have become the standard in state-of-the-art Earth science applications involving image segmentation (Chen et al., 2020; Collins et al., 2020; Gupta et al., 2021; Hoese & Kuenzer, 2020; Jin et al., 2022; Kattenborn et al., 2019; Kotaridis & Lazaridou, 2021; Li et al., 2022; Liu et al., 2019; Marangio et al., 2020; Nagi et al., 2021; Nalepa et al., 2019; Rafique et al., 2022; Sáez et al., 2021; Song et al., 2020; van der Meij et al., 2021; Verma et al., 2021; Xiao et al., 2021).

Deep-Learning-based image segmentation (or “semantic segmentation”) starts with a research question and relevant labeled training data, that is, pairs of images and corresponding labels. Training data can come from existing sources (e.g., Wernette et al., 2022) or made from scratch using labeling tools (e.g., Buscombe et al., 2021). With training data in hand, researchers are left to wrangle, preprocess and format data, followed by building, training, and evaluating models using one of several Deep Learning frameworks. This work often requires substantial trial-and-error experimentation; choosing a model and training technique, as well as implementing those techniques, can be challenging (Yuan et al., 2021). Further, guidance in published papers and code repositories often only present the author's best model (in terms of a validation metric), and not the extensive model training trials that might inform other experiments for a researcher to try when developing a suitable model. This points to a gap in the current software landscape, namely an end-to-end pipeline for geoscientific image segmentation that makes quick experimentation relatively easy.

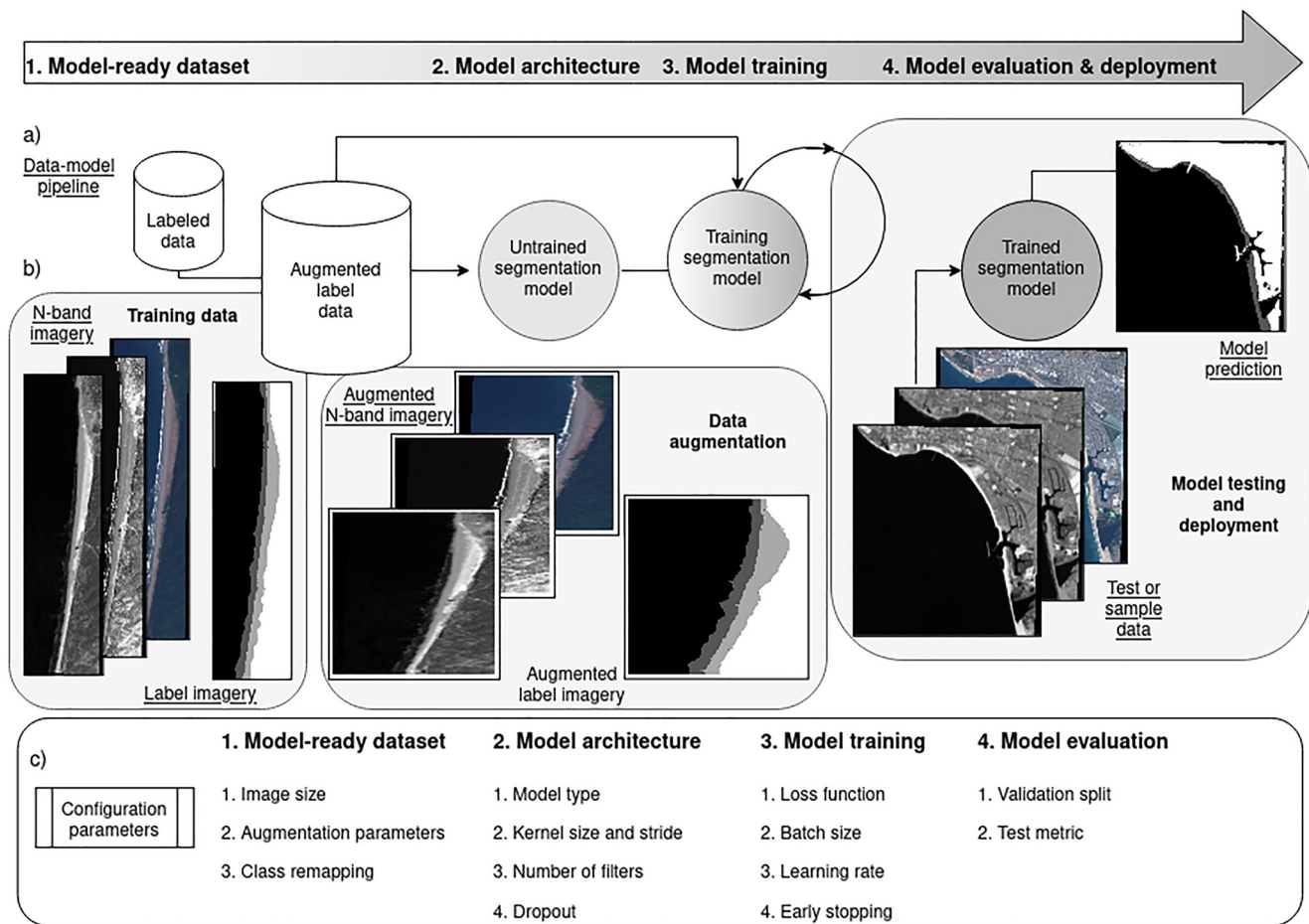


Figure 1. Schematic diagram of the data-model pipeline (a) encoded in the Segmentation Gym software, with examples of model inputs and outputs (b). The four basic pipeline stages are depicted (from left to right): 1. Model-ready data set creation; 2. Model architecture (model building); 3. Model training; and 4. Model evaluation, each with their own set of configuration settings that govern behavior (c).

To fill this gap and aid in the adoption and use of Deep Learning image segmentation, we developed software named “Segmentation Gym” to allow researchers to quickly implement and experiment with segmentation with their own imagery. A fully reproducible workflow enables users to adjust a configuration file and perform their own experimentation with hyperparameters. Segmentation Gym enables its users to fully document the computational provenance of their data models using openly accessible and citable methods (Gil et al., 2016), which moves Earth science segmentation practices closer to realizing a goal of being fully reproducible (Donoho, 2010).

In addition to presenting the design of Segmentation Gym, we demonstrate its use with an example using a data set consisting of 419 image-label pairs. The labeled imagery consist of Landsat-8 scenes of coastal environments from a large collection of labeled images (Buscombe et al., 2022; Wernette et al., 2022). We examine the sensitivity of model outputs to hyperparameter choices that govern model architecture and training strategies. We compare two Deep Learning model architectures (UNets and Residual UNets), five different model sizes, and three different loss functions. The goal of this demonstration is to highlight how researchers can draw insight from the results of this experiment, and adapt a similar modeling campaign with their own data.

2. Implementation

2.1. Overview

We outline the implementation of the software (Figure 1) below in five sections. First, we describe the routines to create a model-ready data set, which ingest data, merge data bands from files, augment the data, remap classes

on-the-fly (if necessary), and format it into batches of tensors of a certain size for model training. Second, we describe the model building process, which involves experimenting with model architecture, such as how large and how numerous feature-extracting kernels are, and experimenting with the use of regularizing layers. Third, model training is described. Fourth, we discuss model evaluation routines, requiring the use of metrics to quantify accuracy on a hold-out data set. Training often involves experimenting with hyperparameters such as the loss function and learning rate. Therefore lastly, we describe how model reproducibility is ensured, for example, by using the same training and validation examples for successive experimentation, such as to test the outcome of retraining with new hyperparameters.

The Python software (Buscombe & Goldstein, 2022) relies on Numpy (Harris et al., 2020), Tensorflow (Abadi et al., 2015), and Keras (Chollet, 2021; Chollet et al., 2015), and is designed to run in an isolated conda (Conda, 2022) environment, a cross-platform and open-source package management system. Models are typically trained using GPUs but a CPU may also be used.

Training a segmentation model requires image-label pairs, and users come to this segmentation workflow with a folder of images and a folder of corresponding labels. Segmentation Gym is part of an ecosystem of tools that includes the labeling program “Doodler,” described by Buscombe et al. (2021). Images labeled using Doodler are readily ingested into the model pipeline (Figure 1), but label images acquired by other means are also supported.

The components outlined in Figure 1 permit rapid exploratory modeling, which is necessary because determining a successful model implementation often can take experimentation. A single configuration file (Figure 1) controls all of the data creation and model training hyperparameters, such as input data size, on-the-fly class remapping, data augmentation behavior, training hyperparameters, regularization, model loss, and model architecture. Configuration files are JSON files containing a data-model pipeline recipe encoded as variables and their values. Therefore each model building trial may be documented by the configuration file that made it.

2.2. Model-Ready Data Set Creation

The routine to convert image-label pairs into formats amenable for training a model is opinionated (Ostblom & Timbers, 2021; Parker, 2017; Peng & Parker, 2021). For example, by design we force users into a certain way of preparing data. In return, users have their data batched and preprocessed in a way that we have found to be successful for training segmentation models using a range of geoscientific imagery.

First, the user is prompted via a Graphical User Interface to enter the path to the images, labels, and a place to store the output files. Imagery can consist of one or more bands in 8-bit unsigned formats. The program allows either (a) 3-band imagery, such as most visible-band photographic data; (b) 1-band imagery, such as other spectral and hyperspectral bands or indices, or bespoke geophysical data bands; or c) merging of 1- or 3-band imagery with any number of additional 1-band images that are coincident in space.

Second, images are standardized by subtracting the mean and dividing by the standard deviation. This is a crucial step toward good model performance on sample imagery whose distributions may differ from the imagery used to train the model (Li et al., 2022; Yuan et al., 2021). Standardizing imagery diminishes the importance of outlier distributions of image values, ensuring better transferability from training to sample imagery.

Third, both images and labels are resized to the target dimensions, and also zero-padded if necessary. Fully convolutional models expect input imagery to all be the same specified size, and it is common practice to resize and reshape imagery to a desired target dimensions. Zero-padding involves placing an image in the center of a large matrix of zeros, such that the boundary pixels of the resulting image are all zero (with a corresponding zero-valued label integer to denote a null class). Padding is necessarily carried out on imagery that has a range of dimensions; if the imagery is smaller than the target dimensions, it is zero-padded. If, however, the imagery is larger than the target dimensions, it is first shrunk (i.e., resampled to a coarser spatial dimension), then zero-padded.

Fourth, classes in the labeled imagery can optionally be remapped. Data sets for image segmentation come with labels from pre-determined class sets. Those classes may be merged, split or otherwise remapped from one set of classes to another, depending on the intended application. For example, if the integer 1 is used to encode the class label “ocean,” and the integer 2 is used to denote “river,” those two classes might be merged such that integers 1 and 2 both denote a third common class, “water.” Remapped label imagery is stored directly in the output files.

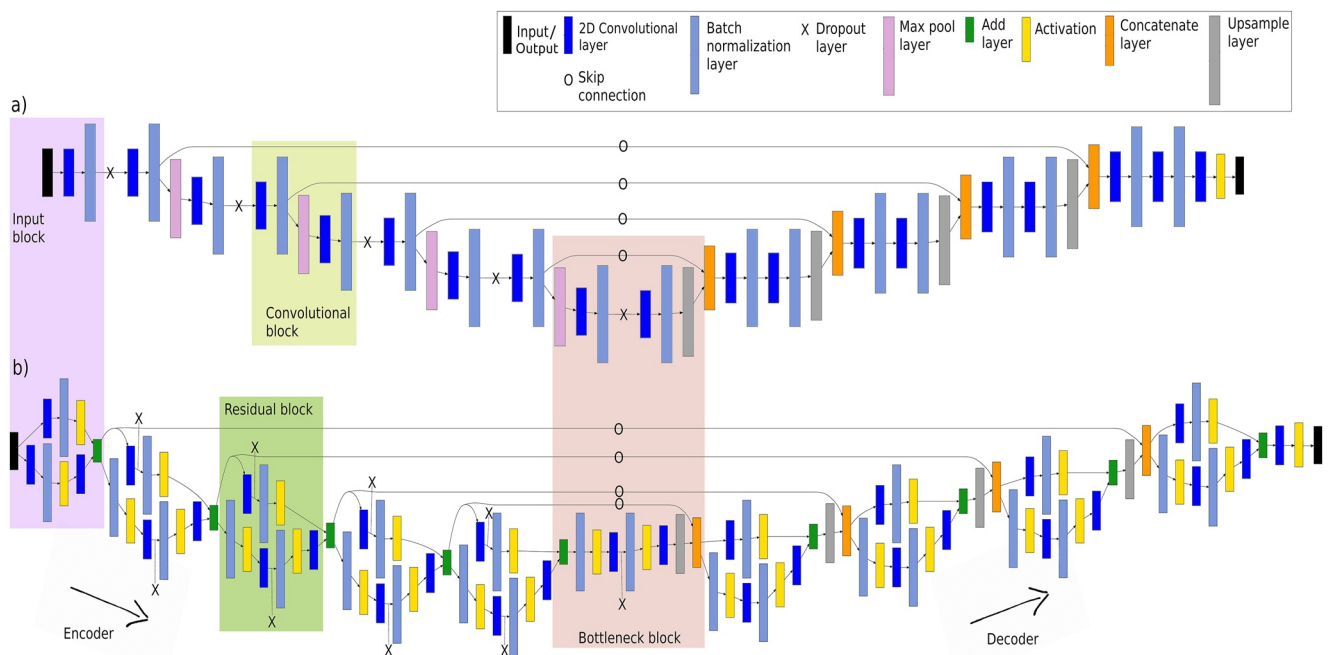


Figure 2. The (a) UNet and (b) Res-UNet fully convolutional model architectures used in the present study. There are several forms of these models available in the software “Segmentation Gym.”

Finally, data undergo augmentation. Augmentation creates transformed versions of the data (Rafique et al., 2022; Stivaktakis et al., 2019; van Lieshout et al., 2020) and is carried out by means of the standard operations available in Keras, such as rotation, width and height shift, zoom, and vertical and horizontal flips. The primary purpose is regularization; by providing alternative versions of the data, the model learns feature representations that are invariant to location, scale, translation, etc. Using augmented imagery to train a model permits oversampling (increasing the size of a data set) without excessive redundancy, because all oversampled augmented imagery will have unique, random augmentations. Optionally, augmentation may be disabled, in which case non-augmented data are used. Examples of augmented and non-augmented images with labels overlayed are printed to file for visual verification.

The model training pipeline stores image/label pairs as a TensorFlow Data set. This format allows for convenient batching, shuffling, and loading of the data set during model training and evaluation. The model training pipeline takes advantage of the TensorFlow Data Application Programming Interface (Tensorflow datasets, 2022), which represents a sequence of single training example, with a pair of tensor components representing the image and its label. Each image and its corresponding label are stored in the compressed numpy binary format, npz (also used by Doodler) as a sequence of binary strings, which allow large data sets to be sequentially loaded to the local (GPU or CPU) memory during model training and evaluation.

2.3. Model Building

Currently, Segmentation Gym implements types of UNet and Residual UNet models. A UNet (Ronneberger et al., 2015) today generally refers to a family of models identified by the following characteristics (Figure 2): (a) fully convolutional (no fully connected layers); (b) four convolutional “blocks” consisting of convolutional layers and Batch Normalization layers connected by ReLu activations, then optionally, Dropout layers; and (c) symmetrical U shape (hence the “U” in the name) with skip connections encoding the Encoder and Decoder branches (Figure 2). Specific UNet implementations often differ in (a) number of filters, (b) stride length (i.e., feature extraction specifics), (c) use of (and type and relative location of) Dropout.

Our UNet (Figure 2a) and Res-UNet (Figure 2b) architecture implementations differ only through the use of residual connections in convolutional blocks in the latter. Residual connections add the outputs of the regular convolutional block with the inputs, so the model learns to map feature representations in context to the inputs

that created those representations (Drozdal et al., 2016). Residual connections (Drozdal et al., 2016) have been shown in numerous contexts to facilitate information flow during model training (Liu et al., 2019; Nagi et al., 2021; Zhang et al., 2018).

The Encoder branch (Figure 2) receives the input image and applies a series of Convolutional and Batch Normalization layers, and optionally Dropout layers, followed by Pooling layers that reduce the spatial size and condense features. Four banks of convolutional filters each use filters that double in size to the previous, thereby progressively downsampling the inputs as features are extracted through pooling. The last set of features (or so-called bottleneck) is a very low-dimensional feature representation of the input imagery. The Decoder upsamples the bottleneck into a label image progressively using convolutional filters, each using filters half in size to the previous, thereby progressively upsampling the inputs as features are extracted through transpose convolutions and concatenation. The sets of features from each of the four levels in the Encoder-Decoder structure are concatenated, which allows learning different features at different levels and leads to spatially well-resolved outputs. The final classification layer maps the output of the previous layer to a single 2D output based on a Sigmoid activation function.

2.4. Model Training

The routine for training a model is opinionated through the specification of the numerical optimizer that guides training, and through the use of a deterministic learning rate scheduler. Neural networks are trained with variations of the stochastic gradient descent (SGD) algorithm, and the specific form of numerical optimizer is a hyperparameter. We use the Adam optimizer (Kingma & Ba, 2014), a variation of SGD. Schmidt et al. (2020) found Adam to be a good choice for almost all Deep Learning models based on Convolutional layers. Other important model training variables specified in the configuration file are (a) batch size, (b) loss function, and (c) learning rate. These tend to have a greater impact on final model accuracy than other tunable hyperparameters such as kernel size, number of convolutional filters, and Dropout, and are therefore described in more detail below.

Data sets are typically larger than can be held in memory so models are trained in batches. Batch size can be an important hyperparameter; when the model is presented with a batch of, say, six images, and six corresponding labels, the model performance and the magnitude of weight adjustment during backpropagation will be evaluated as the average of the six individual discrepancies between model predictions and ground-truth labels. Therefore the size of the batch has an effect on the model's ability to recognize patterns in the presence of variability, with larger variability given by large batch sizes, and hence its rate of convergence in training. Where possible, we recommend using the largest batch size your available GPU memory will allow. Larger batch sizes tend to promote more stable validation loss curves. This is usually only possible with relatively large hardware, because large batches mean larger amounts of GPU memory required. You may therefore decide to use a smaller model input size to achieve a larger batch size if necessary.

During training, the distribution of accuracy scores over classes are optimized using a loss function. Segmentation Gym provides various options for loss function. In this contribution we compare three loss functions, namely mean Dice, categorical cross-entropy or CCE, and Kullback-Leibler distance or KLD. The mean Dice coefficient is given by $D = 2|Y \cap \hat{Y}| / (|Y| + |\hat{Y}|)$, where Y and \hat{Y} are true and estimated label images, respectively, \cap is intersection. Mean Dice is a spatial metric that is relatively insensitive to class-imbalance, or the tendency for a majority class to dominate over one or more minority classes (Csurka et al., 2004). This is because the numerator is the number of correctly classified pixels, and the denominator is the total number of pixels in a class that is in both estimated and ground truth. We therefore can use $1 - D$ as a loss function during class-imbalanced model training, and D to evaluate model results. Many geoscience data sets are significantly imbalanced, and rare classes may be scientifically important, therefore a loss function such as $1 - D$ that handles this can be important for model accuracy. Categorical cross-entropy, $C = -\sum_c Y \log(\hat{Y})$, is a measure of the difference between two distributions over a class set, c , that is, the target or ground truth and the current model estimate, and is a generalization of log loss to multi-class classification problems. Kullback-Leibler distance measures divergence in class-probability distributions and is given by $KLD = \sum_c Y \log(Y/\hat{Y})$.

We vary the learning rate deterministically using a scheduler function that assigns a specific learning rate value as a function of model training epoch. A model epoch is a full training pass over the entire data set such that

each example has been seen once. Thus, an epoch represents $N/\text{batch size}$ training iterations, where N is the total number of examples in the training set. We make use of a function that starts with a small learning rate, then quickly ramps up to a maximum, then decays exponentially. An advantage of varying the learning rate deterministically using a scheduler is to make training reproducible. Decaying the learning rate as training progresses allows the model to slowly converge on an optimal solution. This procedure prevents the solution from getting stuck in a so-called “saddle point,” which is a local minimum much higher than the global minimum. Without varying the learning rate, large updates to the model can potentially lead to suboptimal convergence or prematurely trigger early stopping criteria. In addition to a scheduler, we also implement an early stopping criterion, where the model ceases training early when no improvement to validation loss is observed over a user-defined number of training epochs. Even when models do end after a differing numbers of epochs, the scheduler ensures the learning rate varied in the same way for each model.

2.5. Model Evaluation

Mean Intersection over Union, given by $IoU = |Y \cap \hat{Y}| / (|Y| + |\hat{Y}| - |Y \cap \hat{Y}|)$, is the canonical metric to evaluate model performance. It is sometimes useful to keep track of multiple metrics (Buscombe et al., 2021). Whereas mean IoU is a spatial measure, KLD measures the difference in observed and estimated class-probability distributions. KLD is used as an alternative metric. The only variable related to model evaluation specified in the configuration file is the validation split, which is the proportion of all data to use for model validation. The remainder will be used to train the model. Starting with a relatively high validation split should be a goal, to avoid overfitting and promote generalization. If the model is under-performing on the training data, the validation split should be reduced by small increments accordingly.

2.6. Reproducibility

The data creation and model training process in Segmentation Gym is reproducible. We take several steps to enable this reproducibility. First, we use a seed value to instantiate any numerical operations that involved random numbers, and operating system environment variables are used to guarantee reproducibility in some of the software routines that accelerate computation on GPUs. These collectively ensure consistency in data set creation and in model training. Second, we have adopted the practice of varying the learning rate deterministically using a scheduler function that assigns a specific learning rate value as a function of model training epoch, rather than an adaptive learning rate.

3. Case Study

We provide a case study to demonstrate the use of Segmentation Gym, the results obtained from the software, and the experimentation that the software permits. In this example, our goal is to develop a segmentation model that is able to operate on 419 coastal images from Landsat-8 and classify pixels into one of four classes: water, whitewater, sand, and other. Our specific target is to develop a segmentation model with acceptable accuracy metrics on a relatively large validation subset, without overfitting to the data, that converges to a solution relatively quickly, and is also parsimonious (with only enough parameters to achieve a desired accuracy threshold).

To develop this model we designed an experimental matrix to independently examine the effects on model performance of the following: (a) five different numbers of model parameters, (b) three alternative loss functions, and (c) the presence/absence of residual connections. We kept track of mean IoU and KLD on training and validation portions during model training, and computed those quantities of the validation subset comprising of a reproducibly random draw of 60% of the data and the remainder for model training. The same validation and training data were used to train each model to ensure repeatability and comparison. We also kept track of the epoch at which training was terminated, to quantify model convergence time.

3.1. Data

The labeled imagery we use in the present contribution are one of the 10 data records that comprise the Coast Train data set (Buscombe et al., 2022; Wernette et al., 2022), specifically 419 Landsat-8 (top-of-atmosphere) images and associated labels consisting of time-series from seven coastal locations around the United States

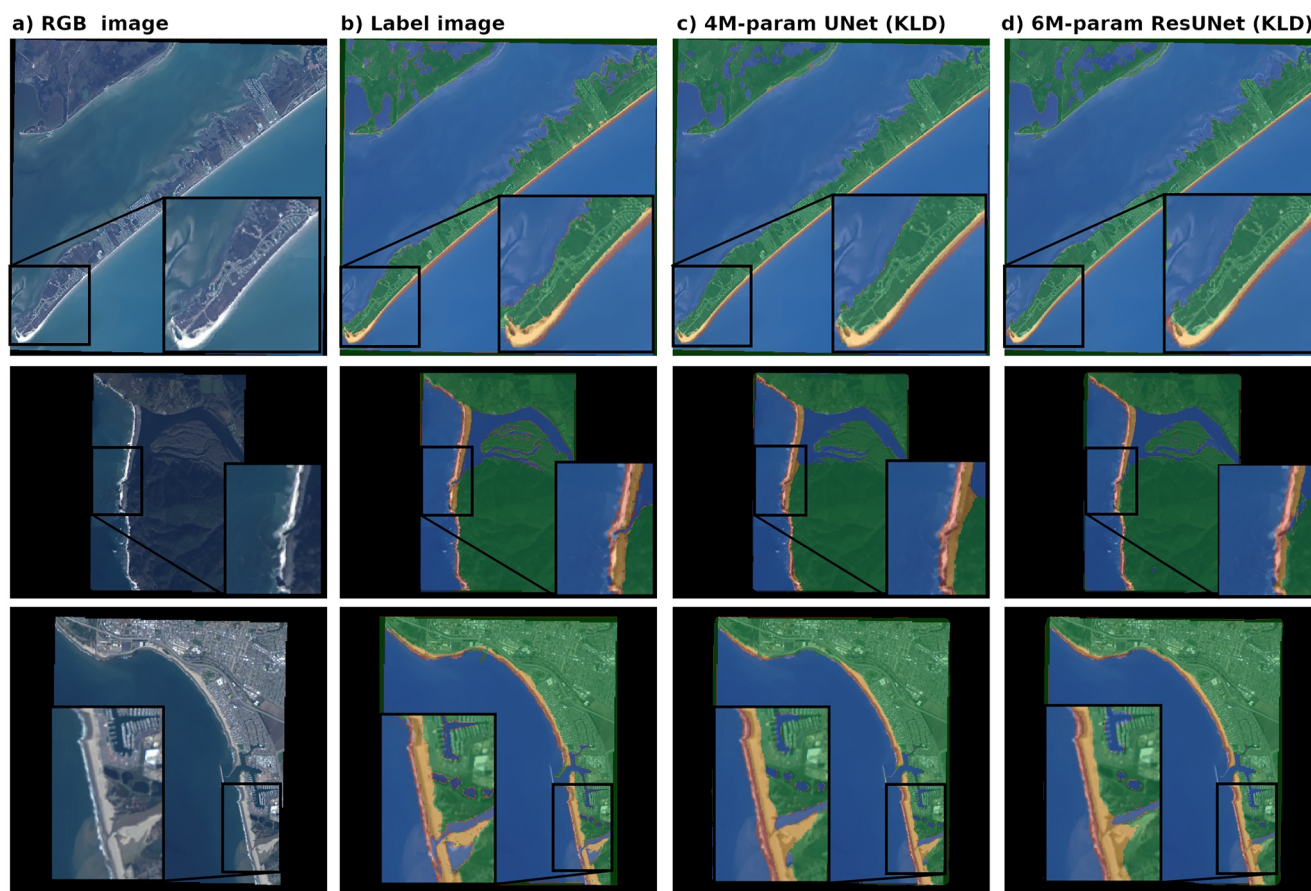


Figure 3. Example model inputs and outputs: (a) RGB imagery, (b) Label imagery (created using Doodler), (c) 4M-parameter UNet model output, (d) 6M-parameter Res-UNet model output. In (b) through (d), labels are shown as colored semi-transparent overlays of the underlying image. Blue indicates water, red is whitewater, yellow is sediment, and green is other. Smaller inset regions are scaled at 200% and better illustrate variability among model outputs, and between model outputs and input label images, as well as error in inputs.

(Figure 3a). The data set consist of visible-band (RGB) imagery, and 2D integer label masks (Figure 3b). The imagery was pre-processed for use in model training. The original 11-classes were remapped into four classes: water, whitewater, sand, and other. Images and corresponding label images were zero-padded to 512×512 pixels if smaller than that dimension, and downsized to 512×512 pixels. We also retrieved the Near Infra-Red (NIR) and Short-wave Infrared (SWIR) bands associated with each RGB image, because spectral indices that contain the NIR and especially the SWIR band have been shown to facilitate more reliable automated classification of water bodies in coastal regions (Luijendijk et al., 2018; Vos et al., 2019). We therefore stacked the three-band visible (RGB) 15-m pan-sharpened imagery with the coincident 15-m pan-sharpened SWIR and NIR bands, and the resulting 5-band raster was used as the model training input.

3.2. Implementation

Images and labels were augmented such that five copies were made each consisting of the original images modified by applying random zoom (up to 5%), rotation (up to 5%), width and height shifts (up to 5%) and horizontal flips. This resulted in 2095 augmented image-label pairs for model training.

We trained 30 different models on the same augmented data set: 15 UNets and 15 ResUNets. Each 15 consist of five model sizes (in terms of parameters), and three loss functions, namely CCE, Dice, and KLD. The number of parameters was varied by adjusting the number of convolutional filters (2, 4, 6, 8, or 12) in the initial convolutional block. The number of filters doubles every subsequent block on the downsampling (Encoder) layer, then halves on every subsequent block on the upsampling (Decoder) layer (Figure 2). Overfitting is countered by three

main model regularization strategies, namely the use of a relatively large validation subset, the use of early stopping, whereby the weights with the smallest validation loss are stored, not from the last training epoch, and the use of Dropout. We used a Dropout rate of 0.1 on each downsample layer but not on upsample layers. We use a 7×7 kernel with a stride of two. Each model was trained with the same learning rate scheduler (varying learning between $1e-7$ and $1e-4$ based on epoch), and batch size (8). Model training stopped early when the validation loss didn't improve upon its previous best value for 10 epochs. As explained above, all of these parameters are specified in a single configuration file. We therefore trained our 30 models using 30 different configuration files.

3.3. Results

To compare the 30 models, we evaluate mean IoU and KLD for the validation subset. Performance statistics (Figure 4) reveal that Res-UNets tend to outperform UNets, as evidenced by a higher average IoU (Figures 4a–4c), lower average KLD (Figures 4d–4f), and also a larger accuracy for smaller number of model parameters. The largest discrepancy between Res-UNet and UNet, and highest model accuracy, is observed when CCE is used for loss (Figure 4b). The performance of Res-UNets demonstrate that residual connections improve statistical measures of success, which is corroborated by visual inspection revealing more realistic model outputs.

In these experiments the best loss is CCE, which promotes the fastest convergence, and highest average and maximum IoU scores. Dice is the worst loss in these trials, as evidenced by very long convergences, and lowest scores. Also, IoU scores do not always appreciably increase with increasing model parameters (Figure 4a). All models improve with more parameters, then plateau or even decline in accuracy (Figures 4b and 4c); an average model size is best overall. Peak Res-UNet model performance tends to occur with fewer parameters compared to an equivalent U-Net.

KLD is a useful model comparative (Figures 4d–4f). It reveals similar (inverse) trends to mean IoU, but larger differences between UNet and Res-UNet trained with Dice loss (Figure 4d), smaller differences between UNet and Res-UNet trained with CCE loss (Figure 4e), and IoU and KLD are most similar when KLD loss is used to train a model (Figure 4f). To provide a representative indication of the variation in accuracies per-site, validation Dice coefficients for the mid-sized model trained using CCE were as follows: Duck: $D = 0.83$ ($N = 76$), Galveston-East: $D = 0.93$ ($N = 40$), Galveston-West: $D = 0.94$ ($N = 40$), Klamath: $D = 0.87$ ($N = 124$), Klamath region: $D = 0.88$ ($N = 69$), Ocean Beach: $D = 0.92$ ($N = 53$), Sunset: $D = 0.90$ ($N = 46$), Ventura: $D = 0.88$ ($N = 40$). These statistics show that the model performs approximately as well across all areas.

4. Discussion

Deep Learning is a powerful set of tools to develop models to segment imagery because of the lack of restrictions imposed on the input variables (such as their distributions or covariance structure). Despite the rapid pace of Deep Learning research, UNets are likely to continue to have considerable application for a number of reasons. First, their success has been demonstrated across many domains and tasks and are already widely known and effectively used in a wide range of scientific applications. Second, they converge well in training even with relatively small amounts of data. Third, they are easily and predictably scalable; catering to larger data sets can be accommodated by increasing the number of filters or batch size and/or lowering the learning rates. Finally, they are easily modified for regression tasks e.g. predicting subgrid wave or current fields from gridded weather variables (Sha et al., 2020), or estimating surf zone bathymetry (Collins et al., 2020).

A key aspect of Segmentation Gym is its link to an existing data labeling software, “Doodler” (Buscombe et al., 2021). Labeled Earth science imagery can be rare and data set creation can be costly and it is not yet possible to know a priori how much data will be needed to train a segmentation model for a given task. Segmentation Gym is specifically designed to quickly and easily build models from Doodler output. This interoperability facilitates interactive data labeling and model building cycle, and permitting researchers to quickly evaluate whether they have enough labeled data to develop a model able to reach the desired level of a test metric.

The Segmentation Gym software can be used across a wide range of remotely sensed imagery, where there is a growing availability, size and relevancy of labeled data sets. Across Earth science fields, segmentation of aerial and satellite imagery is especially a common task (e.g., Bishop-Taylor et al., 2021; Vos et al., 2019). The software system we describe here allows for experimentation with many hyperparameters, and enables researchers

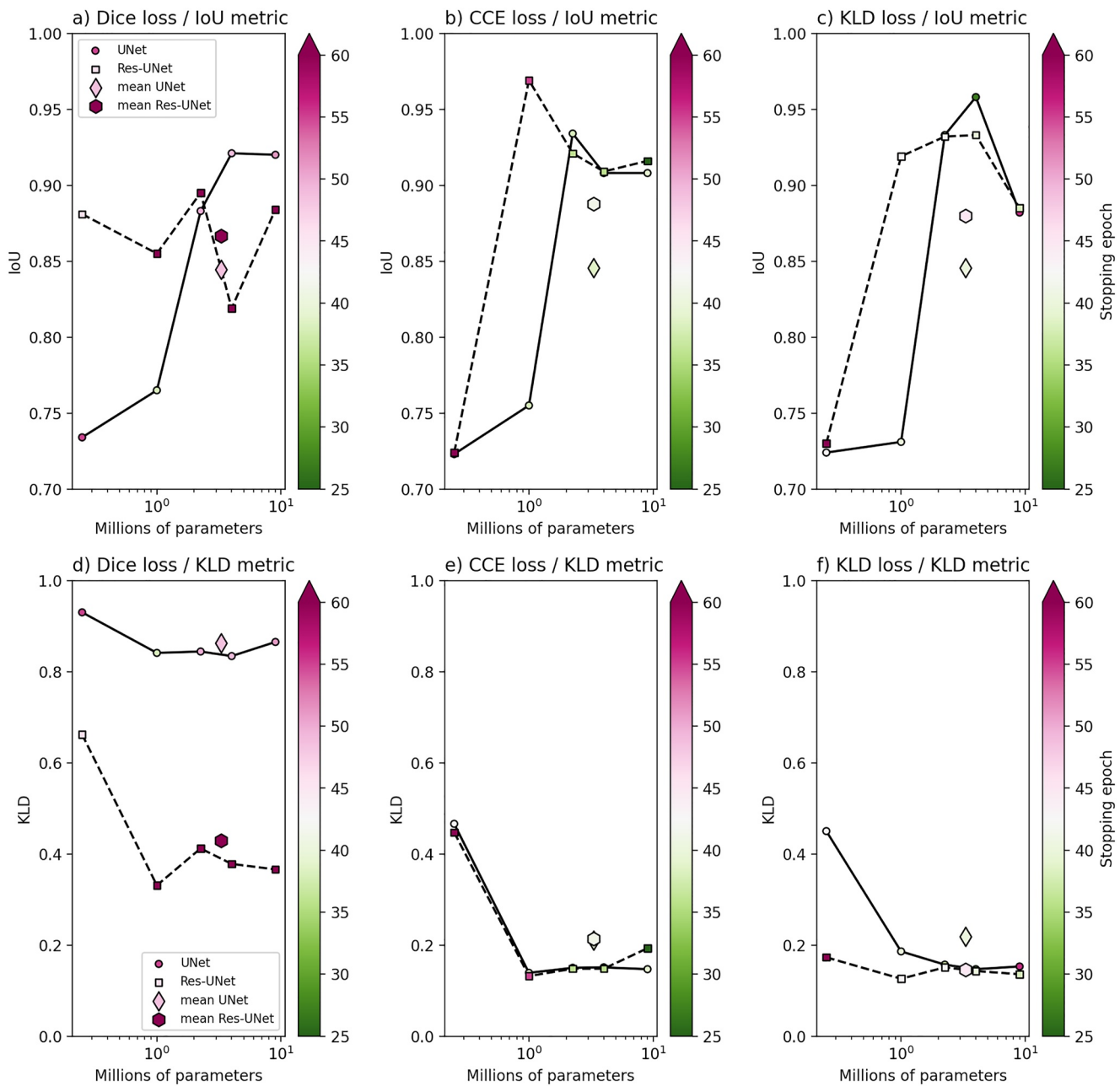


Figure 4. A summary of validation metrics for all 30 models. Each plot shows a model evaluation metric (mean IoU on the top row and Kullback-Leibler distance (KLD) on the bottom) as a function of the number of model parameters. Dashed lines connecting square markers represent Res-UNets, and solid lines connecting circular markers represent UNets. Columns from left to right represent models trained using respectively Dice, categorical cross-entropy, and KLD for a loss function.

to build and use performant models on any suitable data set. We developed Segmentation Gym as an end-to-end reproducible workflow, whereby both labels and model results may be perfectly reproduced in another computing environment by an automated process. Donoho (2010) mentions several important advantages of computational reproducibility, including improved transparency, improved continuity since others can build on the work, greater reusability that leads to greater impact, and obliging scientific funders that the work is preserved.

We envision future work can build from Segmentation Gym. For example, developing a place to share models trained using Segmentation Gym, as well as relevant metadata that is, the configuration file, an example data set, and a “model card” for basic model inventory, reporting, and dissemination (Mitchell et al., 2019). Additionally, we provide a script to segment all images in a directory with a Gym model, but future work could expand this

functionality so that models can be deployed in a range of settings (e.g., as a web or edge computing application, as an internal-facing research tool, etc.).

Data Availability Statement

The Coast Train data used in this study are available from Wernette et al. (2022). The cross-platform open-source application “Segmentation Gym” is available at https://github.com/Doodleverse/segmentation_gym (Buscombe & Goldstein, 2022). Case study model weights and configuration files are available from Buscombe (2022a, 2022b).

Acknowledgments

This work has been supported by the U.S. Geological Survey Coastal/Marine Hazards and Resources Program and by Congressional appropriations through the Additional Supplemental Appropriations for Disaster Relief Act of 2019 (H.R. 2157). EBG acknowledges support from USGS (G20AC00403). The work described here began in preparation for “ML Mondays,” an online course for application of deep-learning-based image analyses, supported by the U.S. Geological Survey Community of Data Integration. See <https://mlmondays.github.io/MLMONDAYS/> for more information. Thanks to Leslie Hsu, Jon Warrick, Sharon Fitzpatrick, Chris Magirl, Kristen Splinter, Simon Topp, and an anonymous reviewer for valuable inputs. Any use of trade, firm, or product names is for descriptive purposes only and does not imply endorsement by the U.S. Government.

References

- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., et al. (2015). TensorFlow: Large-scale machine learning on heterogeneous systems. Retrieved from <https://www.tensorflow.org/>
- Bishop-Taylor, R., Nanson, R., Sagar, S., & Lymburner, L. (2021). Mapping Australia's dynamic coastline at mean sea level using three decades of Landsat imagery. *Remote Sensing of Environment*, 267, 112734. <https://doi.org/10.1016/j.rse.2021.112734>
- Buscombe, D. (2022a). Segmentation Zoo Res-UNet models for Landsat-8 satellite imagery, Coast train v1 Landsat-8 4-class subset [Software]. Zenodo. <https://doi.org/10.5281/zenodo.6229071>
- Buscombe, D. (2022b). Segmentation Zoo UNet models for Landsat-8 satellite imagery, Coast train v1 Landsat-8 4-class subset [Software]. Zenodo. <https://doi.org/10.5281/zenodo.6230083>
- Buscombe, D., Goldstein, E., Sherwood, C., Bodine, C., Brown, J., Favela, J., et al. (2021). Human-in-the-Loop segmentation of Earth surface imagery. *Earth and Space Science*, 9(3), e2021EA002085. <https://doi.org/10.1029/2021EA002085>
- Buscombe, D., & Goldstein, E. B. (2022). Segmentation Gym: Zenodo release for journal article submission March 2022 [Software]. Zenodo. <https://doi.org/10.5281/zenodo.6349591>
- Buscombe, D., Wernette, P., Fitzpatrick, S., Favela, J., Goldstein, E. B., & Enwright, N. (2022). A 1.2 Billion pixel human-labeled dataset for data-driven classification of coastal environments. *EarthArXiv*. <https://doi.org/10.31223/X5Z06C>
- Chen, Z., Liu, X., Yang, J., Little, E., & Zhou, Y. (2020). Deep learning-based method for SEM image segmentation in mineral characterization, an example from Duvernay Shale samples in Western Canada Sedimentary Basin. *Computers & Geosciences*, 138, 104450. <https://doi.org/10.1016/j.cageo.2020.104450>
- Chollet, F. (2021). *Deep learning with python*. Simon and Schuster.
- Chollet, F., et al. (2015). Keras. Retrieved from <https://keras.io>
- Collins, A. M., Brodie, K. L., Spicer, B. A., Hesser, T. J., Farthing, M. W., Lee, J., & Long, J. W. (2020). Bathymetric inversion and uncertainty estimation from synthetic surf-zone imagery with machine learning. *Remote Sensing*, 12(20), 3364. <https://doi.org/10.3390/rs12203364>
- Conda (2022). Continuum analytics, Inc. Retrieved from <https://docs.conda.io/projects/conda/en/latest/>
- Csurka, G., Larlus, D., Perronnin, F., & Meylan, F. (2004). What is a good evaluation measure for semantic segmentation. *IEEE PAMI*, 26(1), 1–11. <https://doi.org/10.5244/C.27.32>
- Donoho, D. L. (2010). An invitation to reproducible computational research. *Biostatistics*, 11(3), 385–388. <https://doi.org/10.1093/biostatistics/kxq028>
- Drozdzal, M., Vorontsov, E., Chartrand, G., Kadoury, S., & Pal, C. (2016). The importance of skip connections in biomedical image segmentation. In *Deep learning and data labeling for medical applications* (pp. 179–187). Springer.
- Gil, Y., David, C. H., Demir, I., Essawy, B. T., Fulweiler, R. W., Goodall, J. L., et al. (2016). Toward the geoscience paper of the future: Best practices for documenting and sharing research from data to software to provenance. *Earth and Space Science*, 3(10), 388–415. <https://doi.org/10.1002/2015ea000136>
- Gupta, A., Watson, S., & Yin, H. (2021). Deep learning-based aerial image segmentation with open data for disaster impact assessment. *Neuro-computing*, 439, 22–33. <https://doi.org/10.1016/j.neucom.2020.02.139>
- Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., et al. (2020). Array programming with NumPy. *Nature*, 585(7825), 357–362. <https://doi.org/10.1038/s41586-020-2649-2>
- Hoeser, T., & Kuenzer, C. (2020). Object detection and image segmentation with deep learning on Earth observation data: A review. Part I: Evolution and recent trends. *Remote Sensing*, 12(10), 1667. <https://doi.org/10.3390/rs12101667>
- Jin, C., Wang, K., Han, T., Lu, Y., Liu, A., & Liu, D. (2022). Segmentation of ore and waste rocks in borehole images using the multi-module densely connected U-net. *Computers & Geosciences*, 159, 105018. <https://doi.org/10.1016/j.cageo.2021.105018>
- Kattenborn, T., Eichel, J., & Fassnacht, F. E. (2019). Convolutional Neural Networks enable efficient, accurate and fine-grained segmentation of plant species and communities from high-resolution UAV imagery. *Scientific Reports*, 9(1), 1–9. <https://doi.org/10.1038/s41598-019-53797-9>
- Kingma, D. P., & Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*. <https://doi.org/10.48550/arXiv.1412.6980>
- Kotaridis, I., & Lazaridou, M. (2021). Remote sensing image segmentation advances: A meta-analysis. *ISPRS Journal of Photogrammetry and Remote Sensing*, 173, 309–322. <https://doi.org/10.1016/j.isprsjprs.2021.01.020>
- LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, 521(7553), 436–444. <https://doi.org/10.1038/nature14539>
- Li, S., Liu, N., Li, F., Gao, J., & Ding, J. (2022). Automatic fault delineation in 3D seismic images with deep learning: Data augmentation or ensemble learning? *IEEE Transactions on Geoscience and Remote Sensing*, 60, 1–14. <https://doi.org/10.1109/tgrs.2022.3150353>
- Liu, C.-C., Zhang, Y.-C., Chen, P.-Y., Lai, C.-C., Chen, Y.-H., Cheng, J.-H., & Ko, M.-H. (2019). Clouds classification from Sentinel-2 imagery with deep residual learning and semantic image segmentation. *Remote Sensing*, 11(2), 119. <https://doi.org/10.3390/rs11020119>
- Luijendijk, A., Hagenaars, G., Ranasinghe, R., Baart, F., Donchyts, G., & Aarninkhof, S. (2018). The state of the world's beaches. *Scientific Reports*, 8(1), 1–11. <https://doi.org/10.1038/s41598-018-24630-6>
- Marangio, P., Christodoulou, V., Filgueira, R., Rogers, H. F., & Beggan, C. D. (2020). Automatic detection of Ionospheric Alfvén Resonances in magnetic spectrograms using U-net. *Computers & Geosciences*, 145, 104598. <https://doi.org/10.1016/j.cageo.2020.104598>
- Mitchell, M., Wu, S., Zaldivar, A., Barnes, P., Vasserman, L., Hutchinson, B., et al. (2019). Model cards for model reporting. In *Proceedings of the conference on fairness, accountability, and transparency* (pp. 220–229). <https://doi.org/10.1145/3287560.3287596>
- Nagi, A. S., Kumar, D., Sola, D., & Scott, K. A. (2021). RUF: Effective sea Ice Floe segmentation using end-to-end RES-UNET-CRF with dual loss. *Remote Sensing*, 13(13), 2460. <https://doi.org/10.3390/rs13132460>

- Nalepa, J., Myller, M., & Kawulok, M. (2019). Validating hyperspectral image segmentation. *IEEE Geoscience and Remote Sensing Letters*, 16(8), 1264–1268. <https://doi.org/10.1109/lgrs.2019.2895697>
- Ostblom, J., & Timbers, T. (2021). Opinionated practices for teaching reproducibility: Motivation, guided instruction and practice. *arXiv preprint arXiv:2109.13656*. <https://doi.org/10.1080/26939169.2022.2074922>
- Pally, R., & Samadi, S. (2022). Application of image processing and convolutional neural networks for flood image classification and semantic segmentation. *Environmental Modelling & Software*, 148, 105285. <https://doi.org/10.1016/j.envsoft.2021.105285>
- Parker, H. (2017). Opinionated analysis development. *PeerJ Preprints*, 5, e3210v1. <https://doi.org/10.7287/peerj.preprints.3210v1>
- Peng, R. D., & Parker, H. S. (2021). Perspective on data science. *Annual Review of Statistics and Its Application*, 9(1), 1–20. <https://doi.org/10.1146/annurev-statistics-040220-013917>
- Rafique, M. U., Zhu, J., & Jacobs, N. (2022). Automatic segmentation of sinkholes using a convolutional neural network. *Earth and Space Science*, 9(2), e2021EA002195. <https://doi.org/10.1029/2021ea002195>
- Ronneberger, O., Fischer, P., & Brox, T. (2015). U-net: Convolutional networks for biomedical image segmentation. In *International conference on medical image computing and computer-assisted intervention* (pp. 234–241). <https://doi.org/10.48550/arXiv.1505.04597>
- Sáez, F. J., Catalan, P. A., & Valle, C. (2021). Wave-by-wave nearshore wave breaking identification using U-Net. *Coastal Engineering*, 170, 104021. <https://doi.org/10.1016/j.coastaleng.2021.104021>
- Schmidt, R. M., Schneider, F., & Hennig, P. (2020). Descending through a crowded valley: Benchmarking deep learning optimizers. *arXiv preprint arXiv:2007.01547*. <https://doi.org/10.48550/arXiv.2007.01547>
- Sha, Y., Gagne, D. J., II, West, G., & Stull, R. (2020). Deep-learning-based gridded downscaling of surface meteorological variables in complex terrain. Part I: Daily maximum and minimum 2-m temperature. *Journal of Applied Meteorology and Climatology*, 59(12), 2057–2073. <https://doi.org/10.1175/jamc-d-20-0057.1>
- Song, Q., Cui, Z., & Liu, P. (2020). An efficient solution for semantic segmentation of three ground-based cloud datasets. *Earth and Space Science*, 7(4), e2019EA001040. <https://doi.org/10.1029/2019ea001040>
- Stivaktakis, R., Tsagkatakis, G., & Tsakalides, P. (2019). Deep learning for multilabel land cover scene categorization using data augmentation. *IEEE Geoscience and Remote Sensing Letters*, 16(7), 1031–1035. <https://doi.org/10.1109/lgrs.2019.2893306>
- Sun, Z., Sandoval, L., Crystal-Ornelas, R., Mousavi, S. M., Wang, J., Lin, C., et al. (2022). A review of Earth Artificial Intelligence. *Computers & Geosciences*, 159, 105034. <https://doi.org/10.1016/j.cageo.2022.105034>
- Tensorflow datasets. (2022). Tensorflow datasets. Retrieved from <https://www.tensorflow.org/guide/data>
- van der Meij, W. M., Meijles, E. W., Marcos, D., Harkema, T. T., Candel, J. H., & Maas, G. J. (2021). Comparing geomorphological maps made manually and by deep learning. *Earth Surface Processes and Landforms*, 47(4), 1089–1107. <https://doi.org/10.1002/esp.5305>
- van Lieshout, C., van Oeveren, K., van Emmerik, T., & Postma, E. (2020). Automated river plastic monitoring using deep learning and cameras. *Earth and Space Science*, 7(8), e2019EA000960. <https://doi.org/10.1029/2019ea000960>
- Verma, U., Chauhan, A., Pai Manohara, M. M., & Pai, R. (2021). DeepRivWidth: Deep learning based semantic segmentation approach for river identification and width measurement in SAR images of Coastal Karnataka. *Computers & Geosciences*, 154, 104805. <https://doi.org/10.1016/j.cageo.2021.104805>
- Vos, K., Splinter, K. D., Harley, M. D., Simmons, J. A., & Turner, I. L. (2019). CoastSat: A Google Earth Engine-enabled Python toolkit to extract shorelines from publicly available satellite imagery. *Environmental Modelling & Software*, 122, 104528. <https://doi.org/10.1016/j.envsoft.2019.104528>
- Wernette, P., Buscombe, D., Favela, J., Fitzpatrick, S., Goldstein, E., Enwright, N. M., & Dunand, E. (2022). Coast Train—Labeled imagery for training and evaluation of data-driven models for image segmentation [Dataset]. U.S. Geological Survey Data Release. <https://doi.org/10.5066/P91NP871>
- Xiao, H., Zhang, F., Shen, Z., Wu, K., & Zhang, J. (2021). Classification of weather phenomenon from images by using deep convolutional neural network. *Earth and Space Science*, 8(5), e2020EA001604. <https://doi.org/10.1029/2020ea001604>
- Yuan, X., Shi, J., & Gu, L. (2021). A review of deep learning methods for semantic segmentation of remote sensing imagery. *Expert Systems with Applications*, 169, 114417. <https://doi.org/10.1016/j.eswa.2020.114417>
- Zhang, Z., Liu, Q., & Wang, Y. (2018). Road extraction by deep residual U-Net. *IEEE Geoscience and Remote Sensing Letters*, 15(5), 749–753. <https://doi.org/10.1109/lgrs.2018.2802944>