

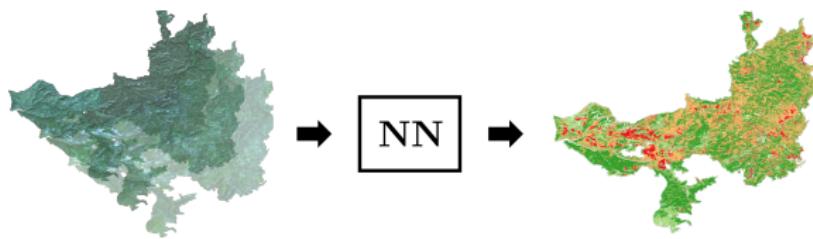
# Semi-Supervised Learning for Spatio-Temporal Segmentation of Satellite Images

Antonín Hruška

June 5, 2023

# Motivation

ESA project<sup>1</sup> in collaboration with Mapradix s.r.o.  
Enhanced Spatiotemporal Land Change Monitoring Based on  
Sentinel-2 Time Series and VHR Images



<sup>1</sup>At <https://eo4society.esa.int/projects/spatiotemporal-sen2vhr/>

# Remote Sensing and its issues

## Sentinel-2 imagery

- ▶ 13 spectral bands
- ▶ period of 10 days
- ▶ different spatial resolution ( $10/20/60\ m^2$ )

## Issues

- ▶ Partial Occlusions: Cloud cover or snow cover (seasonal dependence)
- ▶ Different viewing angle
- ▶ Atmospheric calibration etc.
- ▶ Annotation process is time demanging and requires high expertise.

Semi Supervised Learning can reduce the need for large number of annotations. Generative models can impaint the occluded regions.

# Semi Supervised Learning

$$\mathcal{D} = \mathcal{X} \cup \mathcal{U} : \mathcal{X} = \{(x_1, y_1), \dots, (x_l, y_l)\}, \quad \mathcal{U} = \{(x_{l+1}), \dots, (x_u)\},$$

where  $x_i$  are features and  $y_i$  are labels.

Assumptions from unsupervised learning:

- ▶ Smoothness:  $x_1$  close to  $x_2$  in high density region imples  $y_1$  close to  $y_2$
- ▶ Cluster: Points in one cluster are likely to be of the same class
- ▶ Manifold: The data lie along low-dimensional latent manifolds inside that high-dimensional space.

# MixMatch

**Discriminative model** based on empirical risk minimization

- ▶ Data Augmentation
- ▶ Pseudo labeling
- ▶ Sharpening
- ▶ MixUp

$$\mathcal{X}', \mathcal{U}' = \text{MixMatch}(\mathcal{X}, \mathcal{U}, T, K, \alpha)$$

$$\mathcal{L}_{\mathcal{X}} = \frac{1}{|\mathcal{X}'|} \sum_{x', p' \in \mathcal{X}'} H(p', f_{\theta}(x')), \quad \mathcal{L}_{\mathcal{U}} = \frac{1}{L|\mathcal{U}'|} \sum_{u', q' \in \mathcal{U}'} \|q' - f_{\theta}(u')\|_2^2$$

$$\mathcal{L} = \mathcal{L}_{\mathcal{X}} + \lambda_{\mathcal{U}} \mathcal{L}_{\mathcal{U}}$$

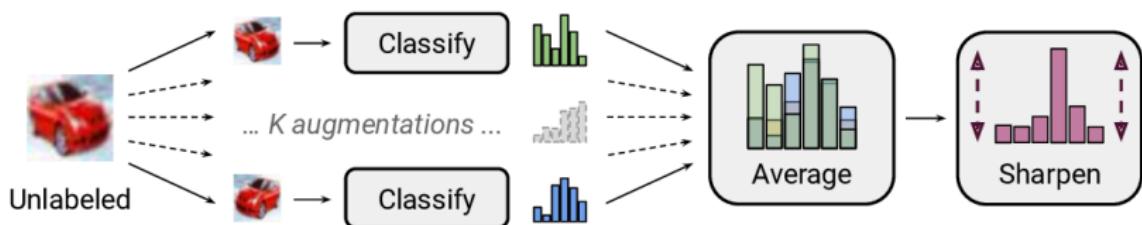


Figure: Mixmatch up to the point of sharpening. Source: MixMatch

# MixMatch

After sharpening :

$$\mathcal{X}^* = \{(\tilde{x}_i, p_i) \mid i \in \{1, \dots, n\}\}, |\mathcal{X}^*| = n$$

$$\mathcal{U}^* = \{(\tilde{u}_{j,k}, q_j) \mid j \in \{1, \dots, n\}, k \in \{1, \dots, K\}\}, |\mathcal{U}^*| = Kn$$

$$\mathcal{W} = \text{Shuffle}(\text{Concat}(\mathcal{X}^*, \mathcal{U}^*))$$

**MixUp:**

$$\begin{bmatrix} \mathcal{X}' \\ \mathcal{U}' \end{bmatrix} = \text{MixUp}\left(\begin{bmatrix} \mathcal{X}^* \\ \mathcal{U}^* \end{bmatrix}, \mathcal{W}\right)$$

where the  $\text{MixUp}(\cdot, \cdot)$  corresponds to following operation applied elementwise:

$$\lambda' = \max(\lambda, 1 - \lambda), \text{ where } \lambda \sim \text{Beta}(\alpha, \alpha)$$

$$\begin{bmatrix} x' \\ p' \end{bmatrix} = \lambda' \begin{bmatrix} x_1 \\ p_1 \end{bmatrix} + (1 - \lambda') \begin{bmatrix} x_2 \\ p_2 \end{bmatrix}$$

# VAE

**Generative model** based on ML maximization

Simple bayesian network with hidden variables  $z$ :

$$p_{\theta}(x, z) = p_{\theta}(z)p_{\theta}(x | z)$$

where  $p_{\theta}(z)$  or  $p_{\theta}(x | z)$  is specified and parametrized by NN.

The intractability of evidence:

$$p_{\theta}(x) = \int_z p_{\theta}(x, z) dz = \frac{p_{\theta}(x, z)}{p_{\theta}(z|x)}$$

do not allow us to maximize the log-likelihood directly, instead we introduce the *inference model*  $q_{\phi}(z|x)$  and optimize the ELBO proxy:

$$\begin{aligned} \log p_{\theta}(x) &= \mathbb{E}_{q_{\phi}(z|x)} \left[ \log \frac{p_{\theta}(x, z)q_{\phi}(z|x)}{q_{\phi}(z|x)p_{\theta}(z|x)} \right] \\ &= \underbrace{\mathbb{E}_{q_{\phi}(z|x)} \left[ \log \frac{p_{\theta}(x, z)}{q_{\phi}(z|x)} \right]}_{\text{ELBO}} + \underbrace{\mathbb{E}_{q_{\phi}(z|x)} \left[ \log \frac{q_{\phi}(z|x)}{p_{\theta}(z|x)} \right]}_{\text{KL}(q_{\phi}(z|x)||p_{\theta}(z|x))} \end{aligned}$$

# VAE

**Learning :** Stochastic approach - unbaised estimator and reparametrization trick allow us to use SGD

$$\mathcal{L} = \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x, z)] - \mathbb{E}_{q_\phi(z|x)} [\log q_\phi(z|x)]$$

## Issues and limitations:

- ▶ Parametrization trick requires continuous latent variable
- ▶ Posterior collapse: Optimization can end in an undesirable stable equilibrium:  $q_\phi(z|x) \approx p_\theta(z)$
- ▶ Bluriness of images: The ELBO enforces the consistent encoder and decoder pairs. This disallow the model to learn the complex data distribution. Solved by more flexible model (such as hierarchical VAE).

# Symmetric equilibrium leaning in VAE

New objectives are

$$\mathcal{L}_p(\theta, \phi) = \mathbb{E}_{\pi(x,z)} [\log p_\theta(x, z)] + \mathbb{E}_{\pi(z)} [\log p_\theta(z)] + \\ \mathbb{E}_{\pi(x)} \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x, z)]$$

$$\mathcal{L}_q(\theta, \phi) = \mathbb{E}_{\pi(x,z)} [\log q_\phi(z|x)] + \mathbb{E}_{\pi(z)} \mathbb{E}_{p_\theta(x|z)} [\log q_\phi(z|x)]$$

for semi-supervised training instead of ELBO:

$$\mathbb{E}_{\pi(x)} [\log p_\theta(x) - \text{KL}(q_\phi(z|x) || p_\theta(z|x))]$$

The optimization does not require the reparametrization trick and allows for broader families of distributions, such as exponential family. *Symmetric Equilibrium Learning of VAE* [submitted to NeurIPS 23]

## Methods

- ▶ CityScape dataset (simple 2D segmentation)
- ▶ U-net backbone
- ▶ plain accuracy and IoU metric
- ▶ hierarchical VAE for symmetric equilibrium learning

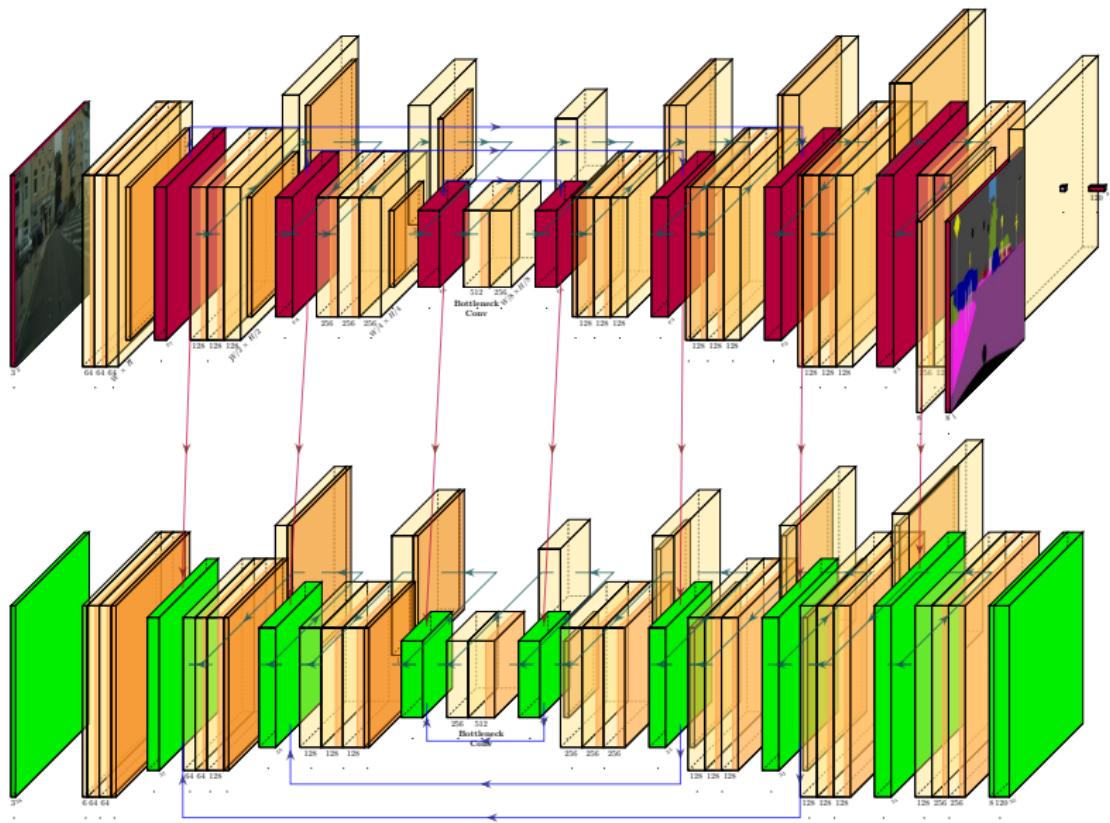
$$z = (z_0, z_1, \dots, z_m), \quad z_0 = (s, l), \quad z_m = x$$

with blockwise optimization:

$$\begin{aligned}\mathcal{L}_p(\theta, \phi) &= \mathbb{E}_{\pi(x,s)} \mathbb{E}_{q_{\theta,\phi}(z_{>0},l|x,s)} [\log p_{\theta}(x,z)] \\ &\quad + \mathbb{E}_{\pi(x)} \mathbb{E}_{q_{\theta,\phi}(z|x)} [\log p_{\theta}(x,z)] \\ \mathcal{L}_q(\theta, \phi) &= \mathbb{E}_{\pi(x,s)} [\log q_{\phi}(s | x)] \\ &\quad + \mathbb{E}_{\pi(s)} \mathbb{E}_{p(l)} \mathbb{E}_{p_{\theta}(x,z_{>0}|z_0)} [\log q_{\theta,\phi}(z|x)]\end{aligned}$$

$\pi(x, s)$  represents the underlying distribution with marginals  $\pi(x)$  and  $\pi(s)$ .

# Architecture of the HVAE



## Mixmatch results (tables)

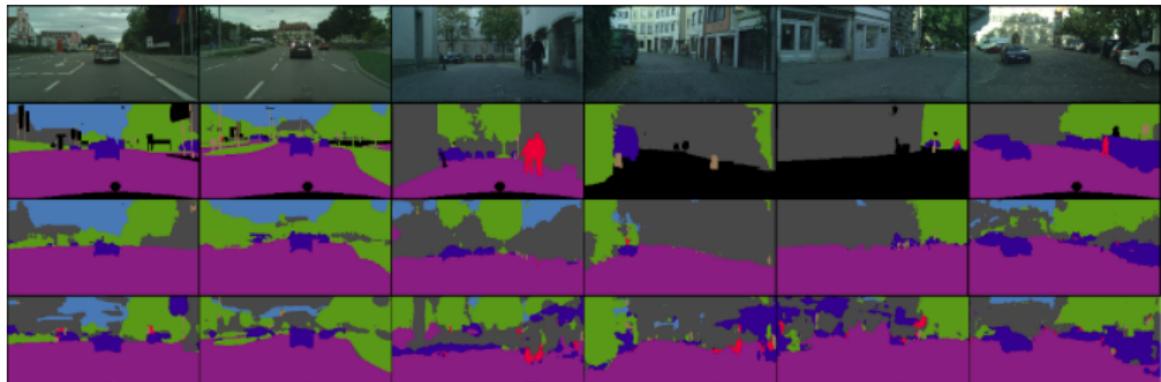
Labels [#]	10	100	500	1000	All
Mixmatch	84.53	90.76	93.42	94.32	94.84
Supervised	76.85	87.59	93.50	94.71	95.58

Table: Mixmatch accuracy rate (%) on CityScape dataset

Labels [#]	10	100	500	1000	All
Mixmatch	48.10	61.98	69.24	71.97	73.08
Supervised	41.78	54.38	68.25	71.93	73.84

Table: Mixmatch average IoU (%) on CityScape dataset

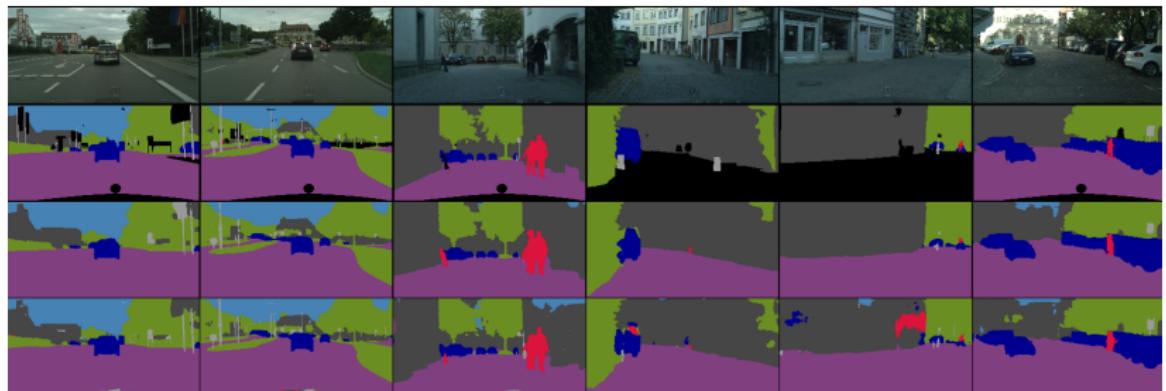
## Mixmatch results (images)



**Figure:** Models were trained on **ten** images. The first two rows contain the image and its ground truth segmentation. The MixMatch predictions are shown in the third row, while the supervised baseline predictions are in the fourth.

The colors used in the segmentation are as follows: flat (purple), human (red), vehicle (dark blue), construction (dark grey), object (light gray), nature (green), sky (light blue), and void (black). The predictions on the "void" class are not penalized nor evaluated.

## Mixmatch results (images)



**Figure:** Models were trained on **all** images. The first two rows contain the image and its ground truth segmentation. The MixMatch predictions are shown in the third row, while the supervised baseline predictions are in the fourth.

The colors used in the segmentation are as follows: flat (purple), human (red), vehicle (dark blue), construction (dark grey), object (light gray), nature (green), sky (light blue), and void (black). The predictions on the "void" class are not penalized nor evaluated.

# Symmetric equilibrium training results

Labels [#]	10	100	500	1000	All
HVAE	74.01	86.18	90.23	92.63	93.51
Supervised	74.82	87.73	90.81	44.05*	43.37*

Table: HVAE accuracy rate (%), \* (learning has failed)

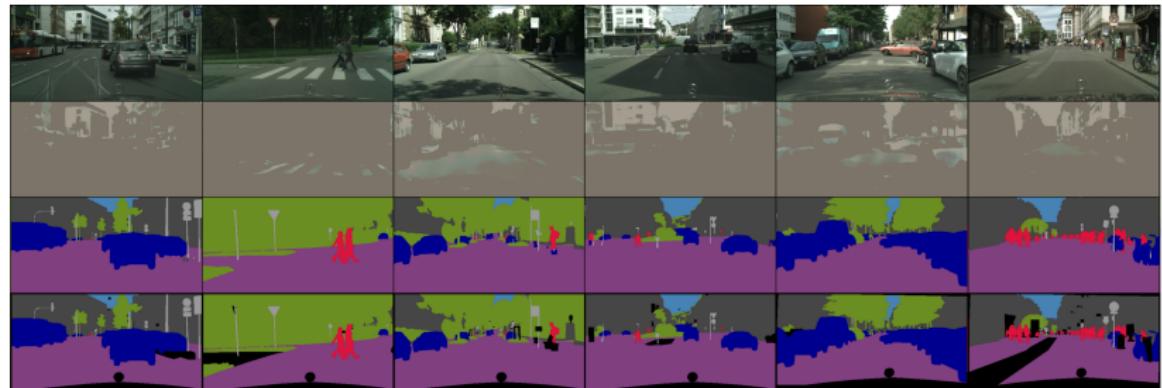


Figure: The first row contains the original image. The second row is filled with the reconstructed images from encoding the original image into  $z_0$  and decoding. The third and fourth rows contain model and ground truth segmentation, respectively.