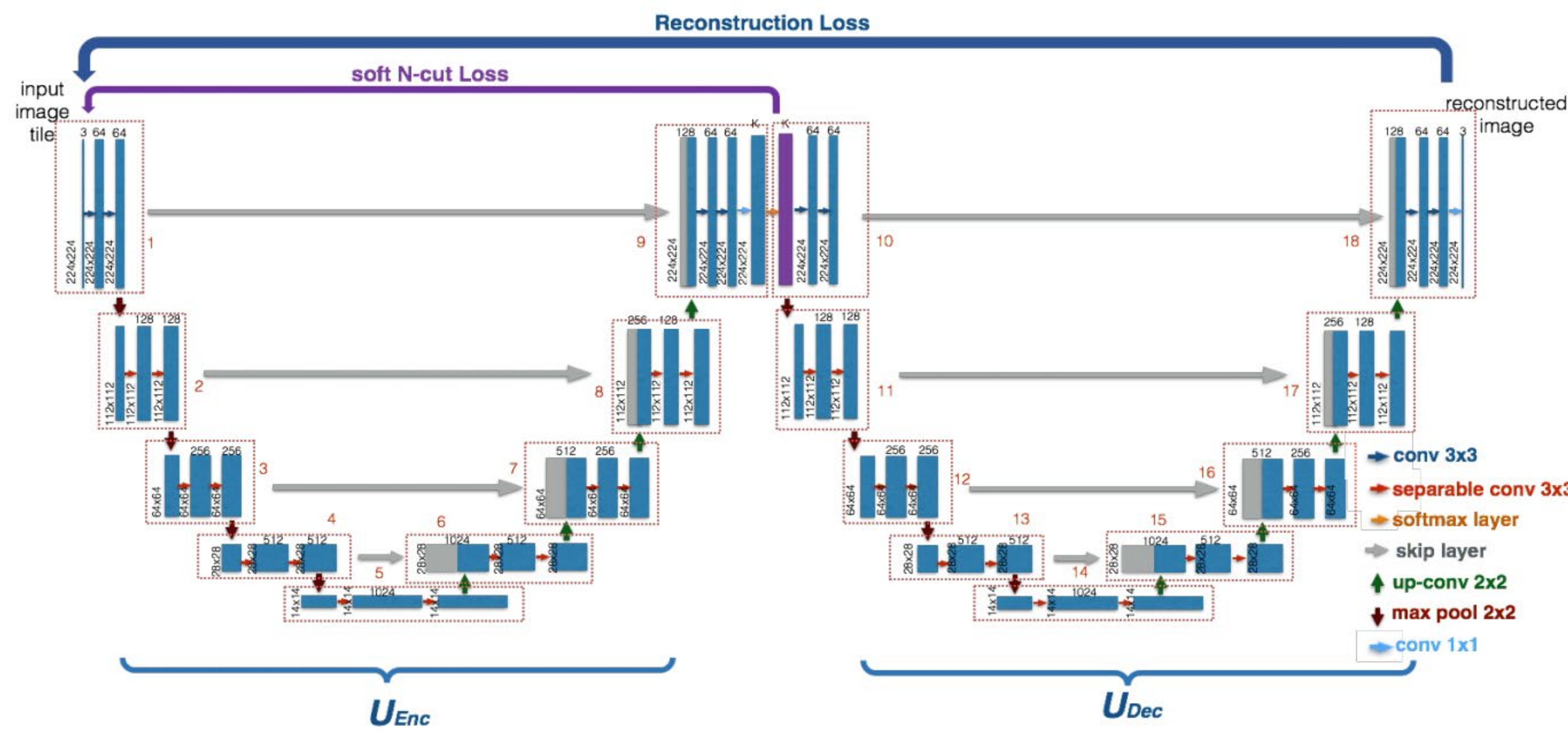


W-Net: A Deep Model for Fully Unsupervised Image Segmentation

Simon Tulling, Natalia Karpova



General W-net architecture



W-Net is an auto-encoder consisting of two auto-encoders. It consists of two connected auto-encoder networks, both having almost similar architecture. The only difference are the amount of channels in the input and output images. The first auto-encoder encodes the image into k segmented image, afterwards the second auto-encoder attempts to reconstruct it back into the original image using the segmented image.

Posprocessing is applied to an image after initial segmentation **for the encoder** is done.

Postprocessing includes 2 main steps, namely:

1. Conditional Random Field
2. Hierarchical Segmentation

CRF is used to increase smoothness within a segmented image such that the final outcome will have sharper boundaries but more even reconstruction inside regions.

Two types of losses

The paper makes use of 2 different loss functions. The main reconstruction loss is calculated in the end of W-Net pipeline:

$$J_{reconstr} = \|X - U_{Dec}(U_{Enc}(X; W_{Enc}); W_{Dec})\|_2^2$$

In addition to the reconstruction loss, soft N-Cut Loss is used on the output of encoder:

Final results: white bear image example

Only reconstruction loss

$$J_{soft}(V, K) = \sum_{k=1}^K \frac{cut(A_k, V - A_k)}{assoc(A_k, V)}$$

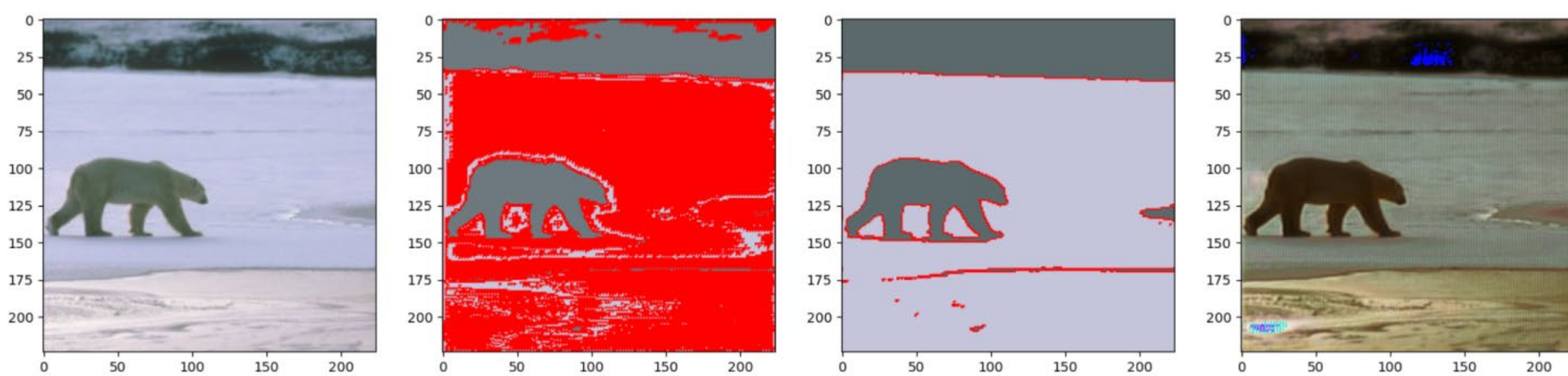
$$= K - \sum_{k=1}^K \frac{assoc(A_k, A_k)}{assoc(A_k, V)}$$

$$= K - \sum_{k=1}^K \frac{\sum_{u \in V, v \in V} w(u, v) p(u = A_k) p(v = A_k)}{\sum_{u \in A_k, t \in V} w(u, t) p(u = A_k)}$$

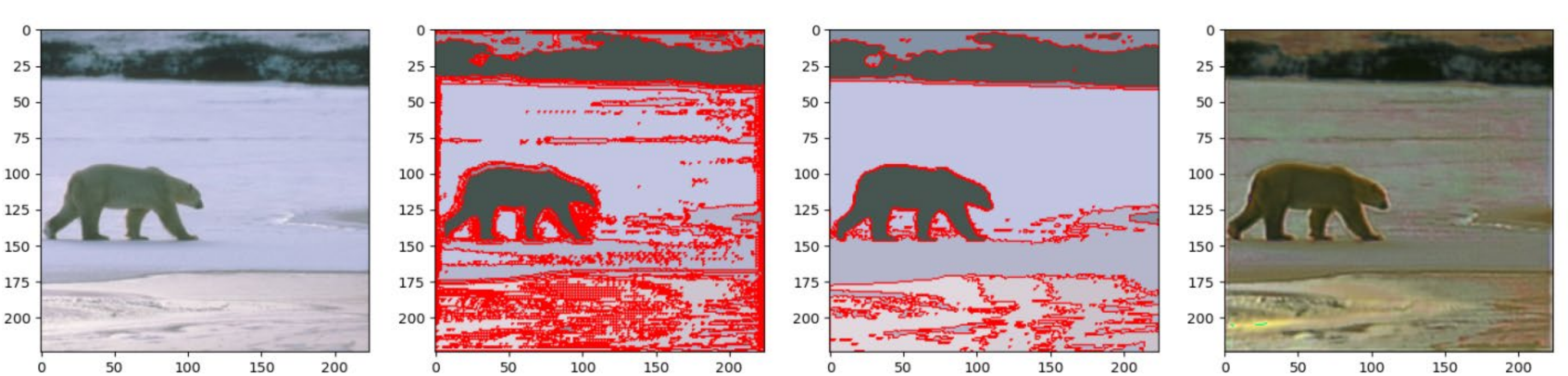
Reconstruction and Soft N-Cut loss

Input After U encoder After CRF Output

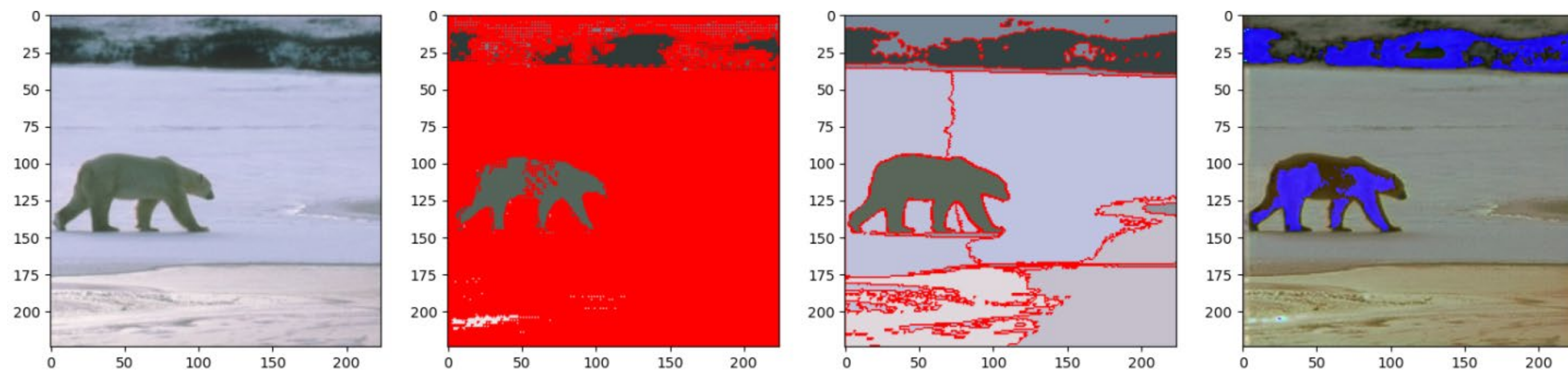
k=5



k=20

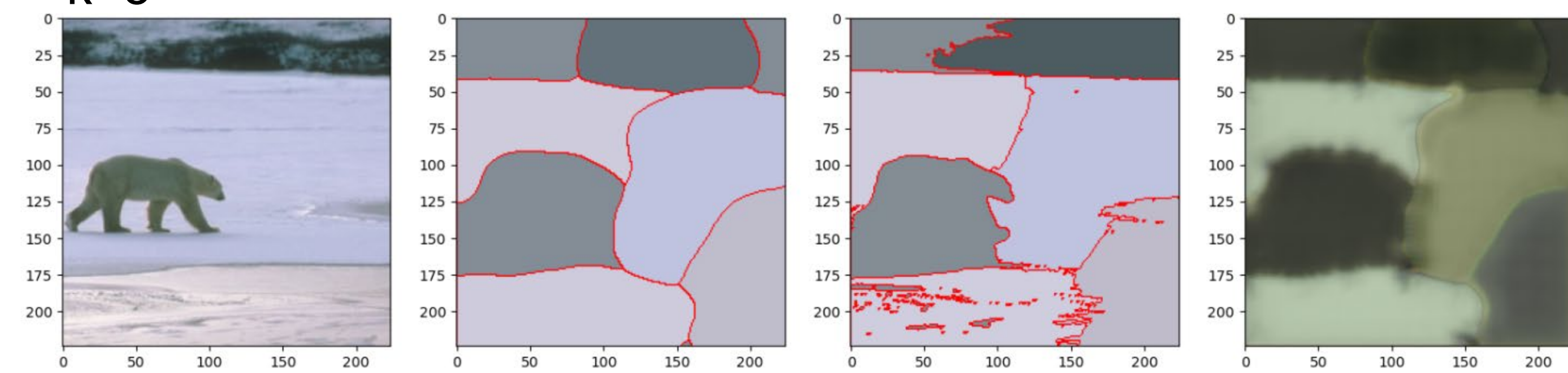


k=64

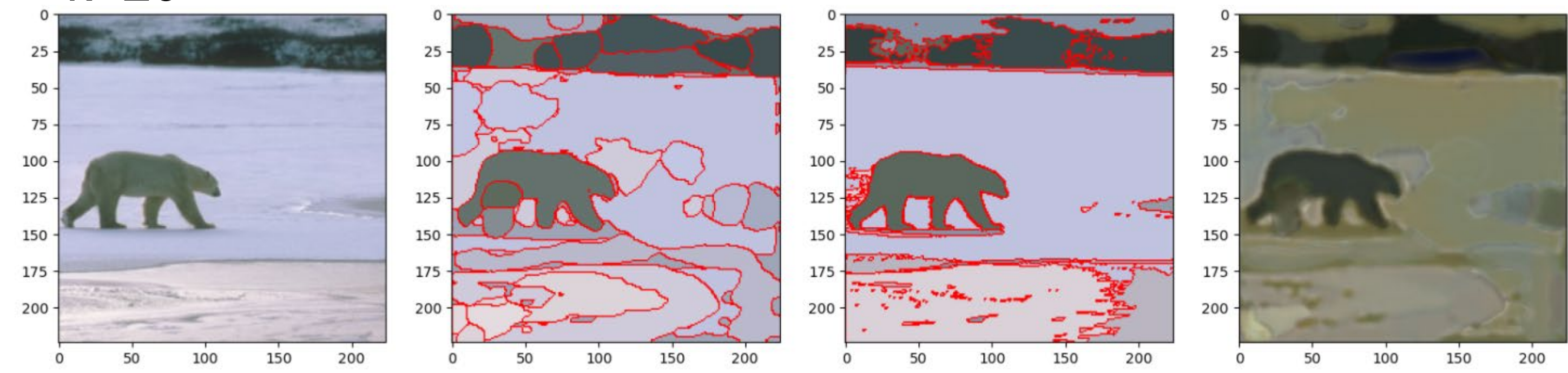


Input After U encoder After CRF Output

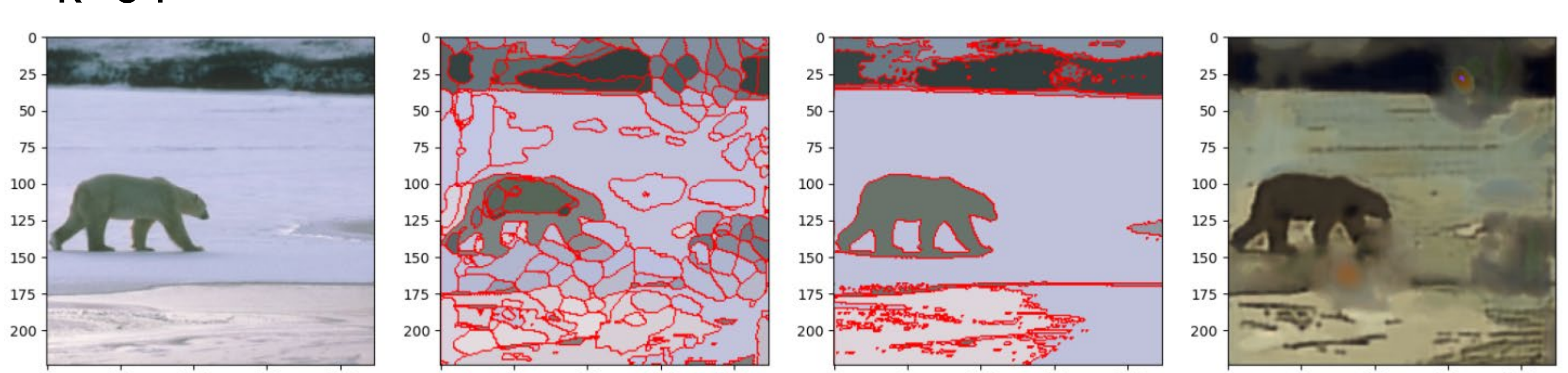
k=5



k=20



k=64



Conclusions:

1. When using reconstruction loss only, image reconstruction is better, yet segmentation is not
2. Hierarchical segmentation is missing and might be crucial for more accurate reconstruction
3. **Amount of classes around 20 seems to be the best k values**
4. CRF does heavy lifting

We would like to thank our supervisor Attila Lengyel for his help during this project

Email: n.karpova@tudelft.nl