

Finding key biological features for cancer diagnosis from histopathology slides

Cell segmentation within histopathological slides

Naylor Peter

PSL - Mines ParisTech, Institut Curie and Inserm

peter.naylor@mines-paristech.fr



Introduction

Uncovering information from histopathology image data is a difficult task and is mostly unused in cancer research. This data corresponds to thin slices of the tumor and of the surrounding tissue. Histopathology slides can thus be very informative of the cancer subtype and/or of how the patient's immune system is reacting to the cancer. Our ultimate goal would be to quantify, via the extraction of biological features, a patient histopathology data. Biological feature are features that have a true biological meaning, like proportion of cancerous cell, normal cells, etc. See [?] for more insight. I will introduce the data and the methods for segmentation.

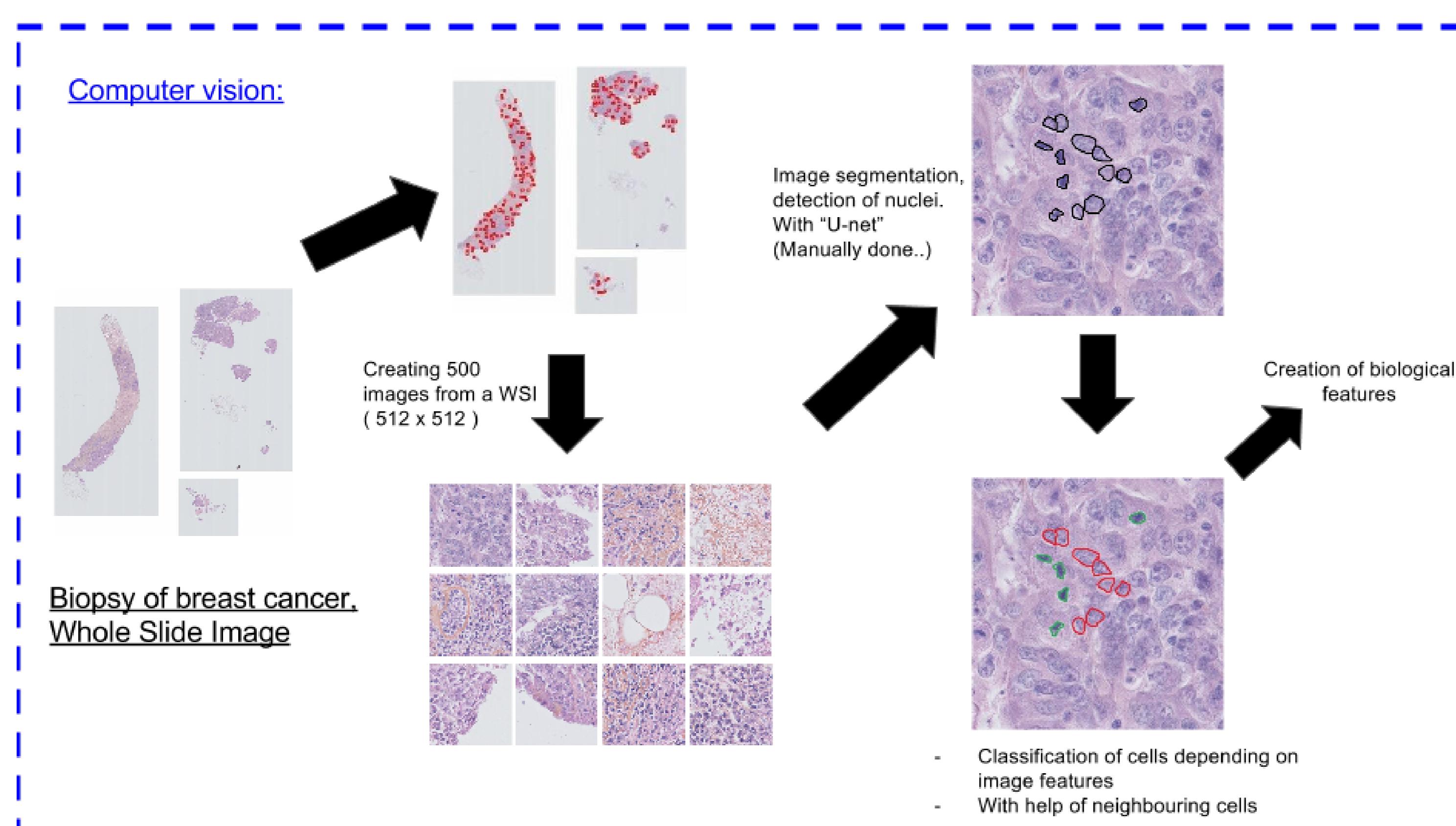
Methods

Our strategy to identify the important features is to first segment the important elements in histopathology slides (such as cells, tumor and stromal tissue, necrotic regions, etc.), second to define physiologically interpretable features for each of these elements and third to build a prediction model in order to assess the importance of each of these features. We propose a method based on fully convolutional network architectures for image segmentation that rely on standard convolutional networks.

Motivation

1. Unused data in cancer research.
2. Residual Cancer Burden Calculator, a tumor grading, is based on the content of histopathology slides, however pathologist only have a limited time per slide.
3. Reproducibility of the RCBC grading, it can vary between hospitals.
4. Defining biology driven features will interpretability for predicting clinical variables, such as outcome, subtype or response to treatment.
5. Allow us to investigate the link between genomic and transcriptomic features.

Pipeline



Uncovering information is a difficult task

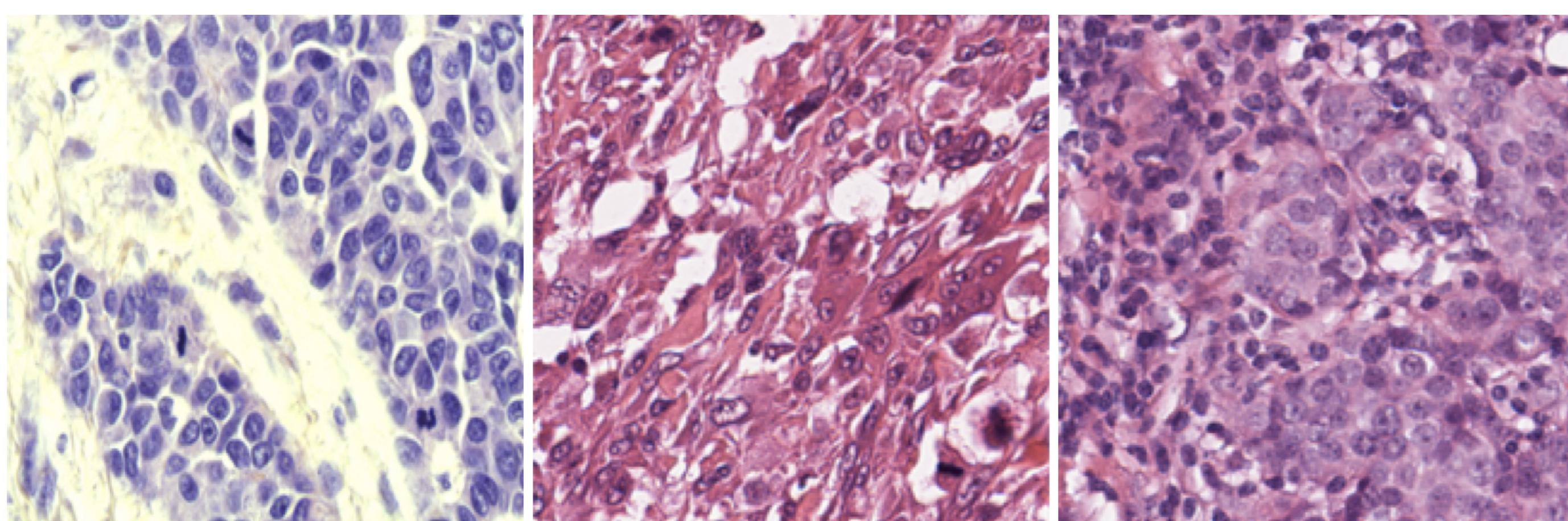


Figure 1: Histopathology data

1. **Data size:** each patient has several slides, one slide is more than 50GB. A typical dataset: hundreds of slides.
2. **Stain variation:** Many reasons for this variability : scanner type, the stain supplier and the stain quality, differences in slide preparation and tissue type.
3. **Variability in the objects:** another variability is biological variability. Many different cells and tissue types.
4. **Projection artefacts:** a slice is actually a 3D slice, we can have overlapping cells / nuclei and other artefacts.
5. **No explicit formalization of information:** No clear descriptions of the objects we are trying to identify.

Manual annotation of histopathology slides

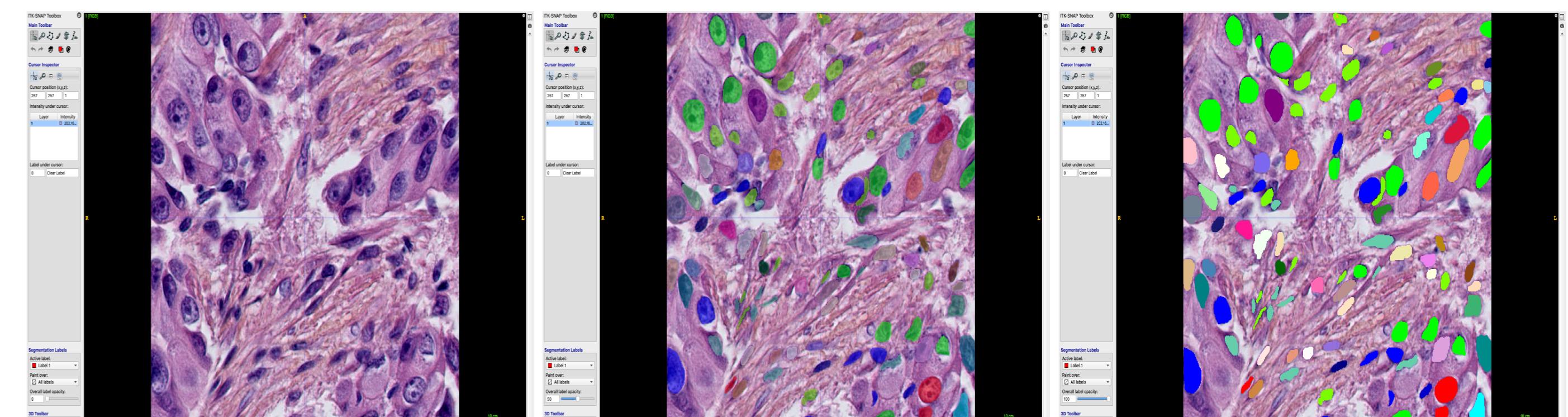


Figure 2: Histopathology data

We use ITK-snap to manually annotate our histopathology slide. We wish to detect cells vs background. We have 33 images of size 512×512 over 7 different patients. Cells can be very tricky to annotate:

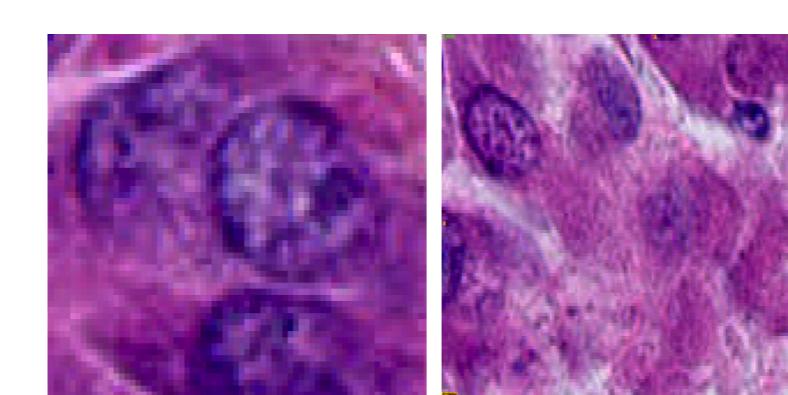


Figure 3: 3D Slice

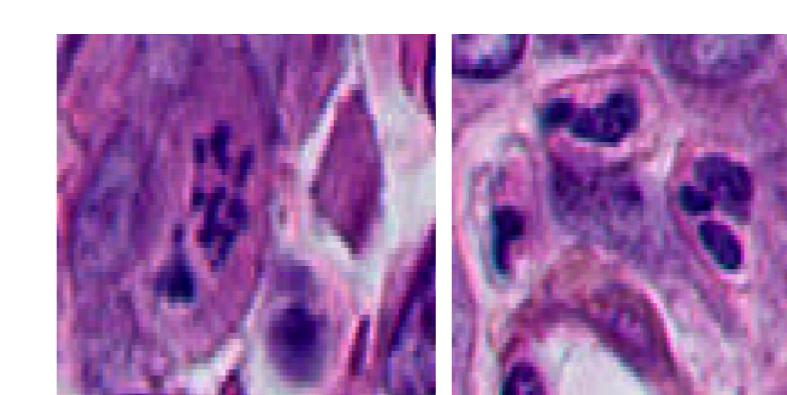


Figure 4: Weird looking nuclei

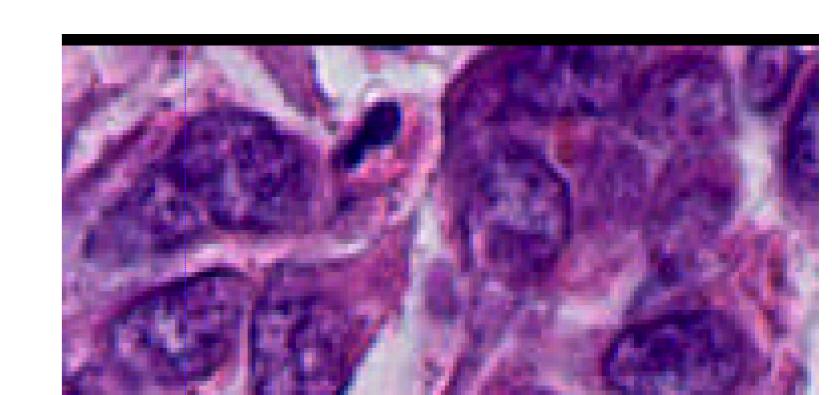


Figure 5: Dense region

Fully Convolutional Networks (FCN)

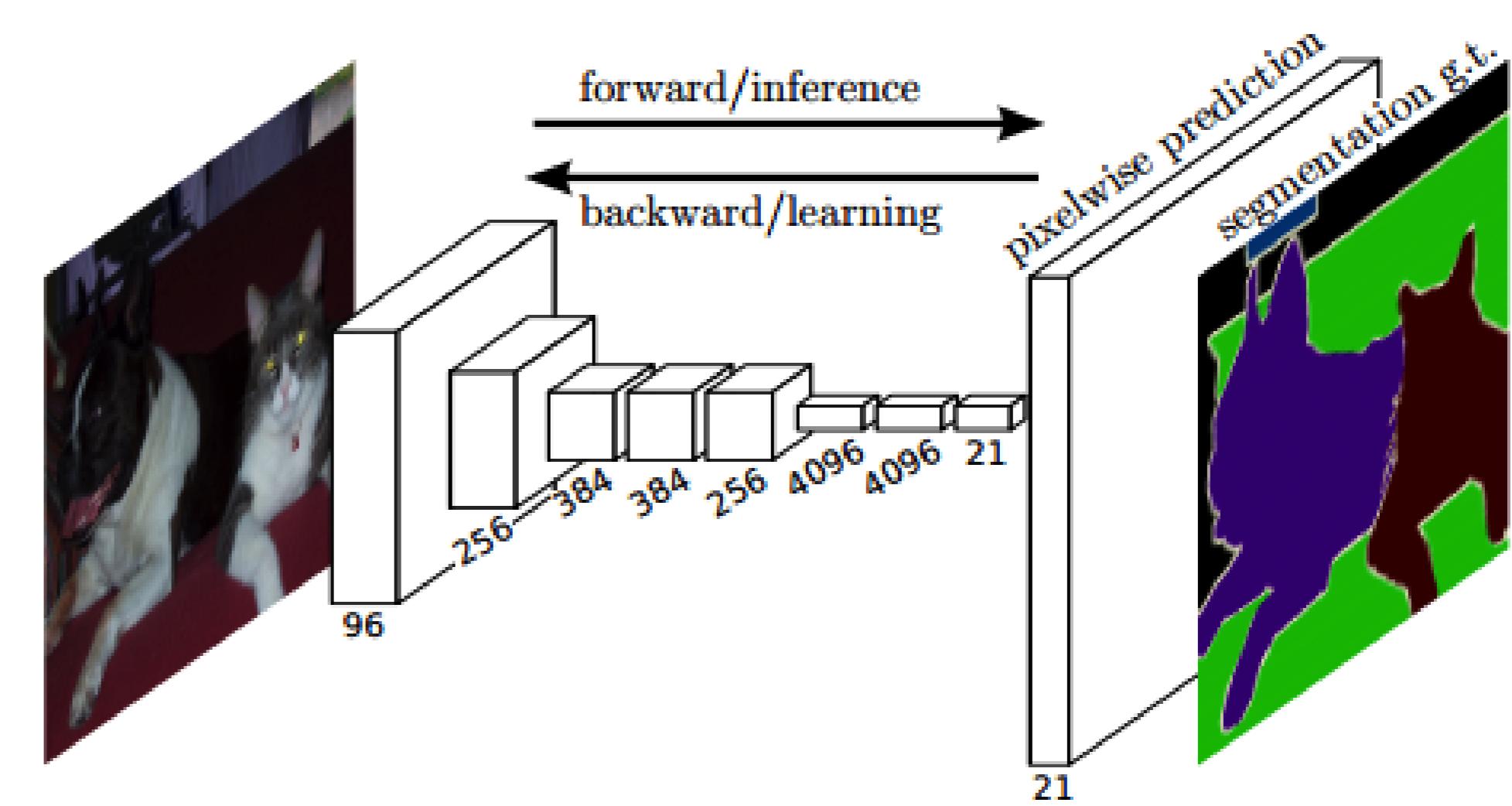


Figure 6: FCN for semantic segmentation (image taken from [?])

FCN for segmentation are an extension of standard Convolutional network for image recognition. These models take an end to end tuned CNN, and cast it into a FCN. Upsampling and deconvolutional layers are added, the standard part of the CNN provides the "what" while the added layers try to provide the "where". Finally, skip path are added between layers to provide final layers with information from the first layers.

Results

The fully convolutional network is fine tuned with a lot of use of data augmentation: rotation, flips, blurring and elastic deformation. The size of the input images can be of size 512×512 or of size 256×256 . Several metrics were kept, especially: mean accuracy, intersection over union, the Jaccard Index, recall and precision. Finally training was performed on 21 images and test on 5 images across 6 patients. 7 validation images were used for reporting the validation scores, these validation images were provided by one patient.

Crop size	Network Name	MA	IU	Recall	Precision
512	FCN8	0.54	0.52	0.09	0.12
256	FCN8	0.63	0.53	0.28	0.09
256	FCN8_200	0.71	0.53	0.45	0.10
256	FCN8_2000	0.65	0.53	0.31	0.11

Table 1: First results

Conclusions

Mediocre results with only 33 annotated patches of whole slide tumors. The segmentation method is not precise enough. Dense cell region will lead to connexe regions.

Forthcoming Research

- Setting up slightly different architectures: U-net [?].
- Changing the loss to make it more adapted to segmentation, [?].
- Changing the channels on the input, HE standardization [?].
- Incorporating different prior knowledge about cell segmentation directly in the architecture. [?].