

Medical Imaging with Deep Learning Overview

Popular image problems:

- Chest X-ray
- Histology

Multi-modality/view

Segmentation

Counting

Incorrect feature attribution

Chapter 1

Radiology and multi-view

Common X-ray projections

Most common



(a) P-A



(b) Lateral



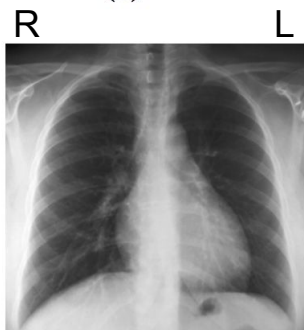
(c) Lordotic



(d) A-P supine



(e) A-P



(f) P-A



(g) Lateral



(h) Lordotic



(i) A-P supine



(j) A-P

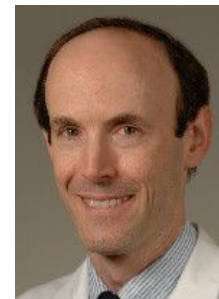
PA = PosteroAnterior

Chest X-ray14 Dataset

Released 2017, first large scale chest X-ray dataset

>100k PA images released without copyright.

Enabled the deep learning radiology revolution



Ronald Summers
NIH Clinical Center

Media Advisory Wednesday, September 27, 2017

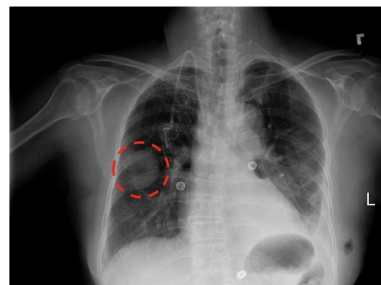
NIH Clinical Center provides one of the largest publicly available chest x-ray datasets to scientific community

The dataset of scans is from more than 30,000 patients, including many with advanced lung disease.

What

The NIH Clinical Center recently released over 100,000 anonymized chest x-ray images and their corresponding data to the scientific community. The release will allow researchers across the country and around the world to freely access the datasets and increase their ability to teach computers how to detect and diagnose disease. Ultimately, this artificial intelligence mechanism can lead to clinicians making better diagnostic decisions for patients.

NIH compiled the dataset of scans from more than 30,000 patients, including



Stanford Pneumonia study

In 2017 Pranav Rajpurkar and Jeremy Irvin trained a DenseNet on NIH data scaled to 224x224 pixels

Set the benchmark performance which has not been significantly improved.

They evaluated pneumonia predictions against 4 radiologists.

"We find that the model exceeds the average radiologist performance on the pneumonia detection task."

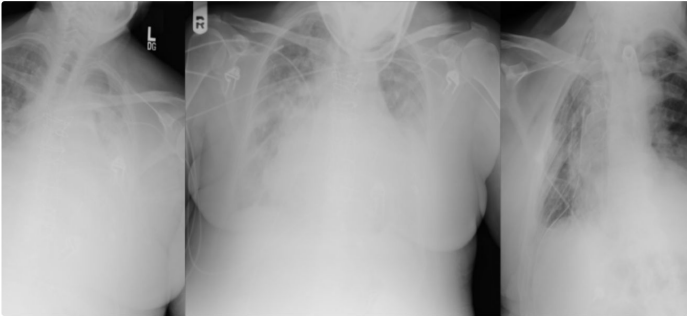
	F1 Score (95% CI)
Radiologist 1	0.383 (0.309, 0.453)
Radiologist 2	0.356 (0.282, 0.428)
Radiologist 3	0.365 (0.291, 0.435)
Radiologist 4	0.442 (0.390, 0.492)
Radiologist Avg.	0.387 (0.330, 0.442)
CheXNet	0.435 (0.387, 0.481)

Criticism of the Chest X-ray14 Dataset

In 2017 Luke Oakden-Rayner published a blog post discussing issues with the labels in the NIH data.

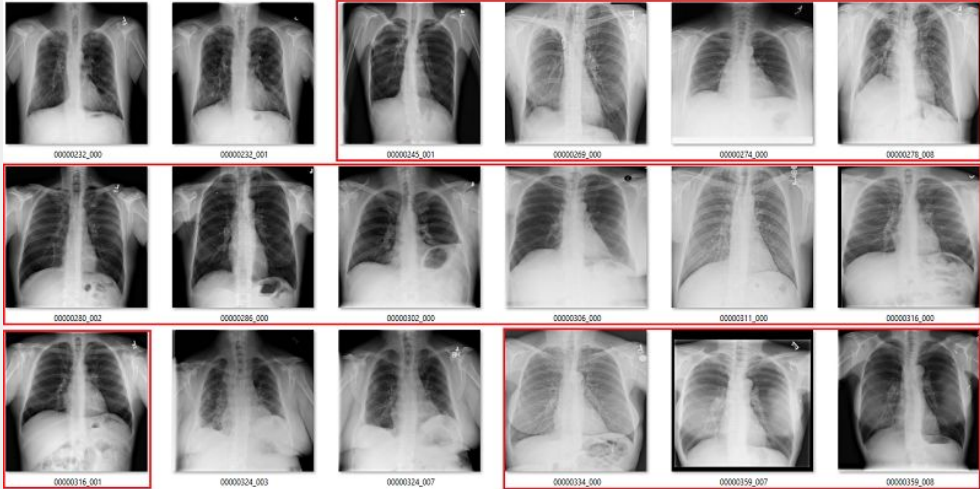
This led to more work on automatic label extraction.

Exploring the ChestXray14 dataset: problems



DECEMBER 18, 2017 - LUKEOAKDENRAYNER

A couple of weeks ago, I mentioned I had some concerns about the ChestXray14 dataset. I said I would come back when I had more info, and since then I have been digging into the data. I've talked with Dr Summers via email a few times as well. Unfortunately, this exploration has only increased my concerns about



In a sample of images red are said to be wrong

2019: the year of X-ray data



PADCHEST
160k images
Multiple views
Almost 200 labels

27% hand labelled, others
using an RNN.

License: Creative Commons
Attribution-ShareAlike



CheXpert
224k images
PA and L views
13 labels.

Automated rule-based
labeler

Non-commercial
research purposes only



MIMIC-CXR
377k images
PA and L views
13 labels.

Automated rule-based labeler.
NIH (NegBio) and CheX
labelers ran.

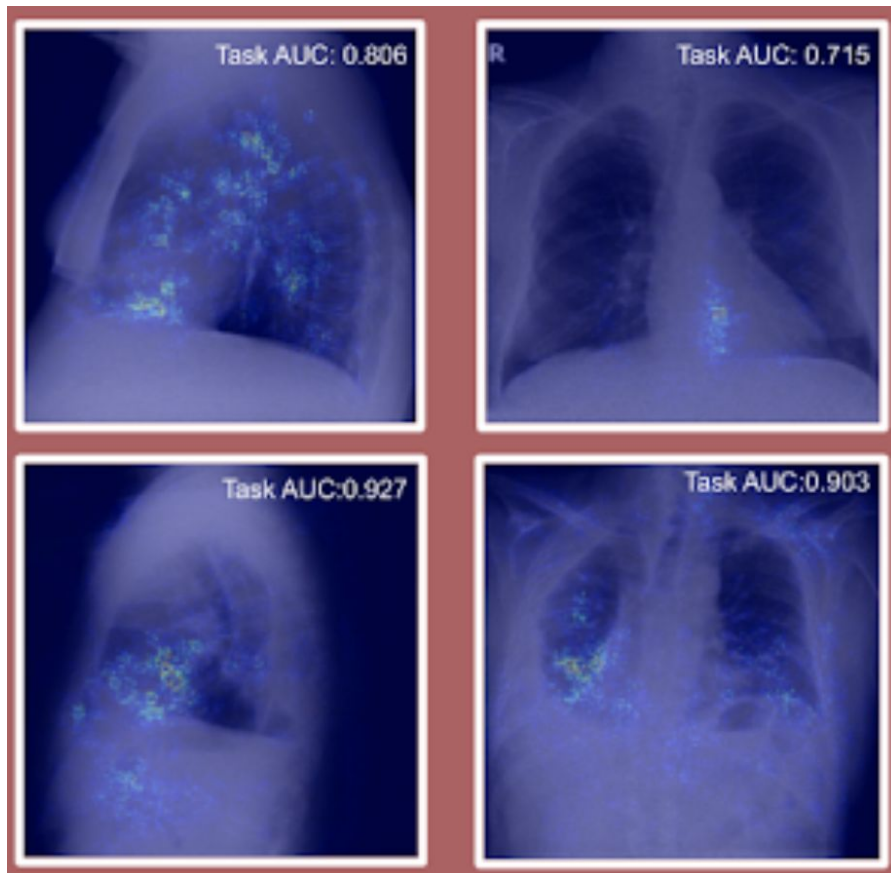
Non-commercial research
purposes only. Confidentially
training required.

Multi-modal/view inference (X-ray use case)

Lateral

PA

Flattened
diaphragm

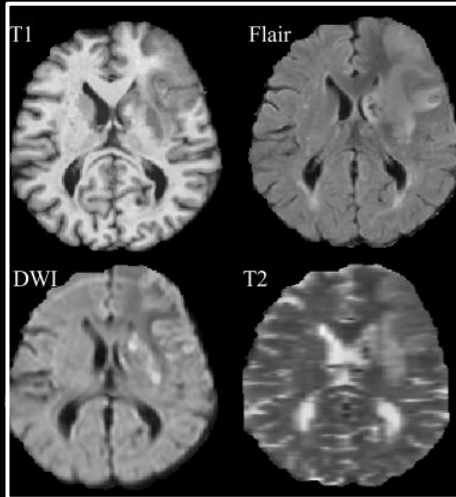


Here saliency maps are from models trained on single views.

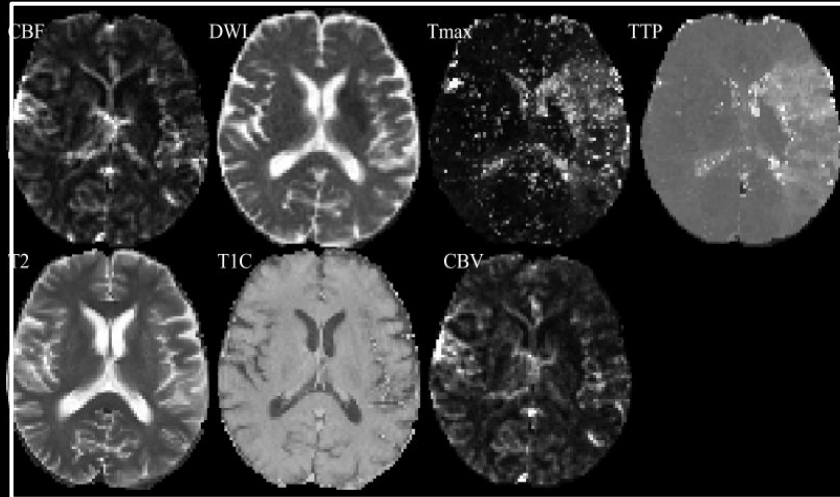
These two tasks perform better when using lateral views.

[Bertrand, 2019]

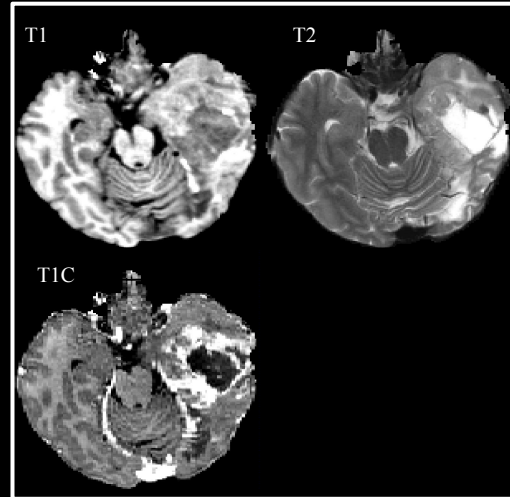
Also: Multi-modal/view inference (MRI use case)



Ischemic stroke
lesion segmentation



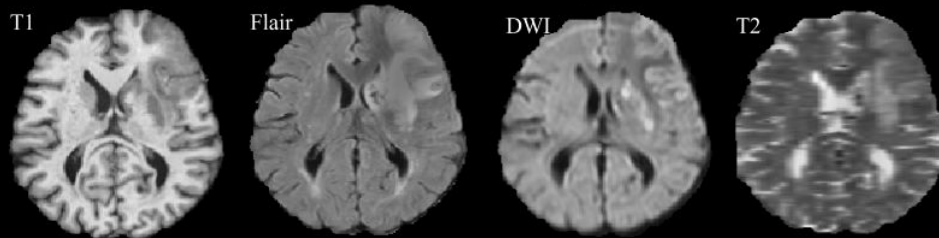
Stroke Perfusion
Estimation



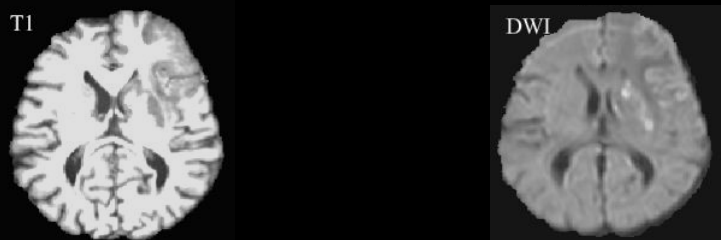
Brain tumor
segmentation

Challenge: missing modalities/views

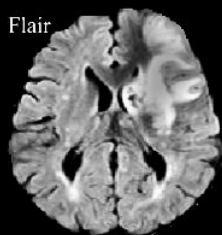
Patient 1



Patient 2



Patient 3



**Incomplete
Input!**

Expected:

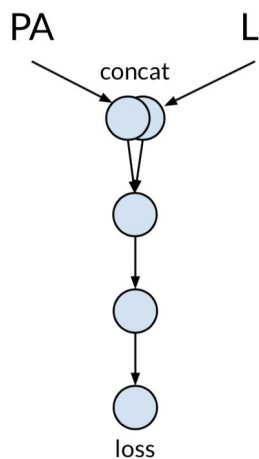
$$f(\text{T1}, \text{Flair}, \text{DWI}, \text{T2})$$

Given:

$$f(\text{T1}, -, \text{DWI}, -)$$

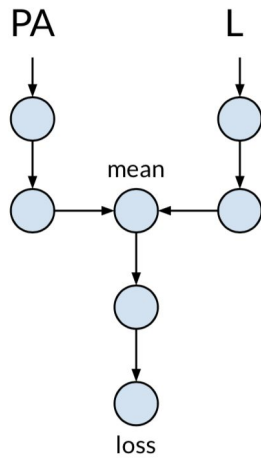
$$f(-, \text{Flair}, -, -)$$

Integrating multiple views



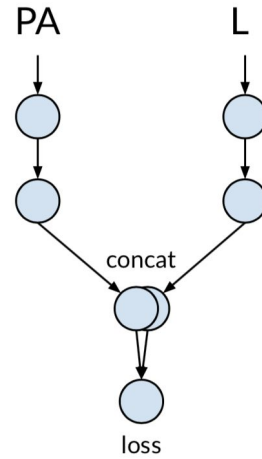
Stacked

Combine images right at the input



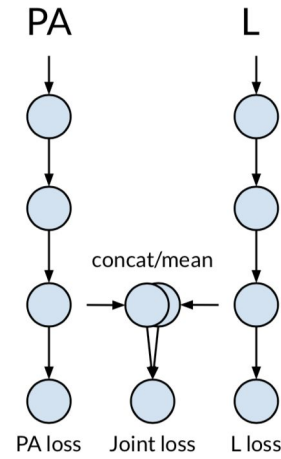
HeMIS

Take mean of activations in the middle of the network



DualNet

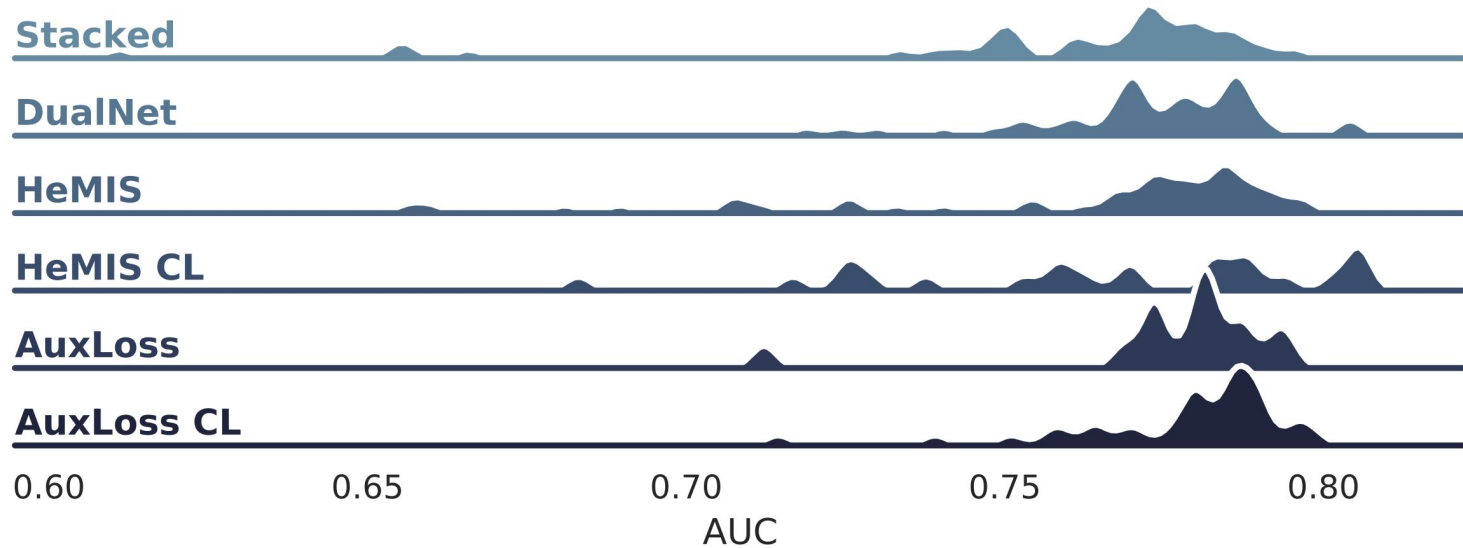
Concat output features of two models with single prediction



AuxLoss

Three losses. A network for each modality with losses that regularize each network.

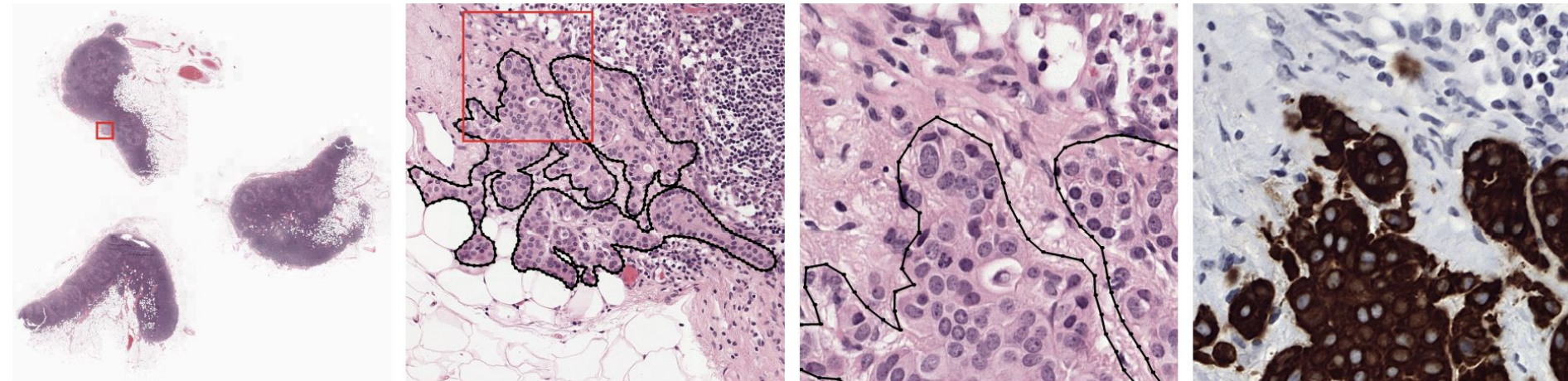
Integrating multiple views (X-ray images)



All models are about equal in performance given the right hyperparameters.
Hyperparameter tuning is easier on some models but not others

Chapter 2

Histology and segmentation



Example of a WSI of a H&E stained section with a delineated micro-metastasis at increasing zoom levels, and the corresponding IHC (cytokeratin 8-18 stained) slide at the same location. The metastasis is outlined with black.

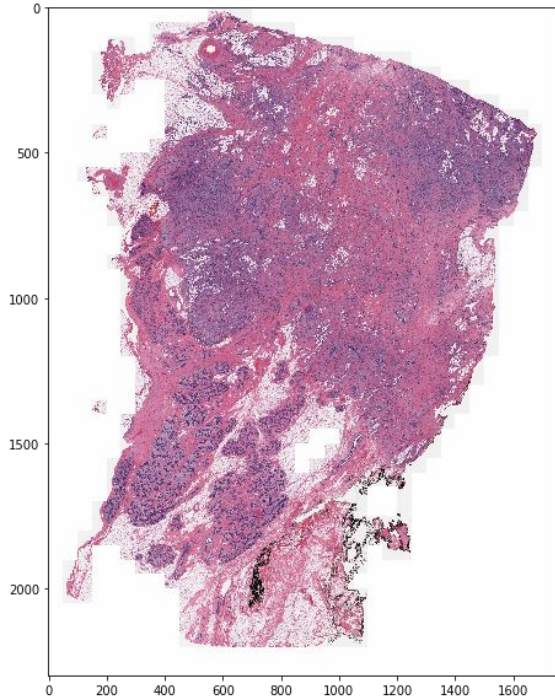
CAMELYON17 Dataset

1000 whole-slide images (WSIs) of sentinel lymph node. (~3GB each!)

5 medical centers. 40 patients from each center. 5 whole-slide images per patient.

Patch wise segmentation

Use case: Invasive Ductal Carcinoma (most common subtype of all breast cancers)



Starting with a full slide image of breast tissue.

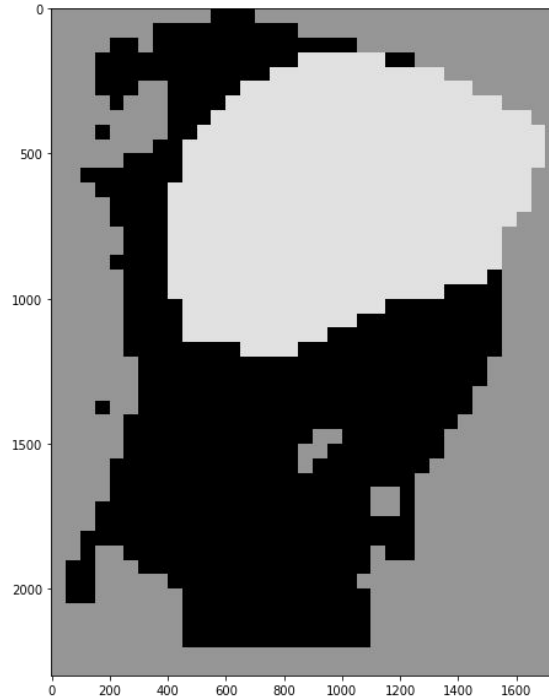


Image is labelled as IDC or not

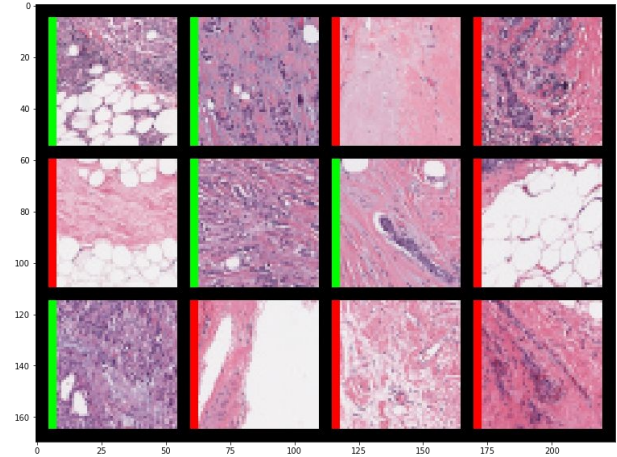
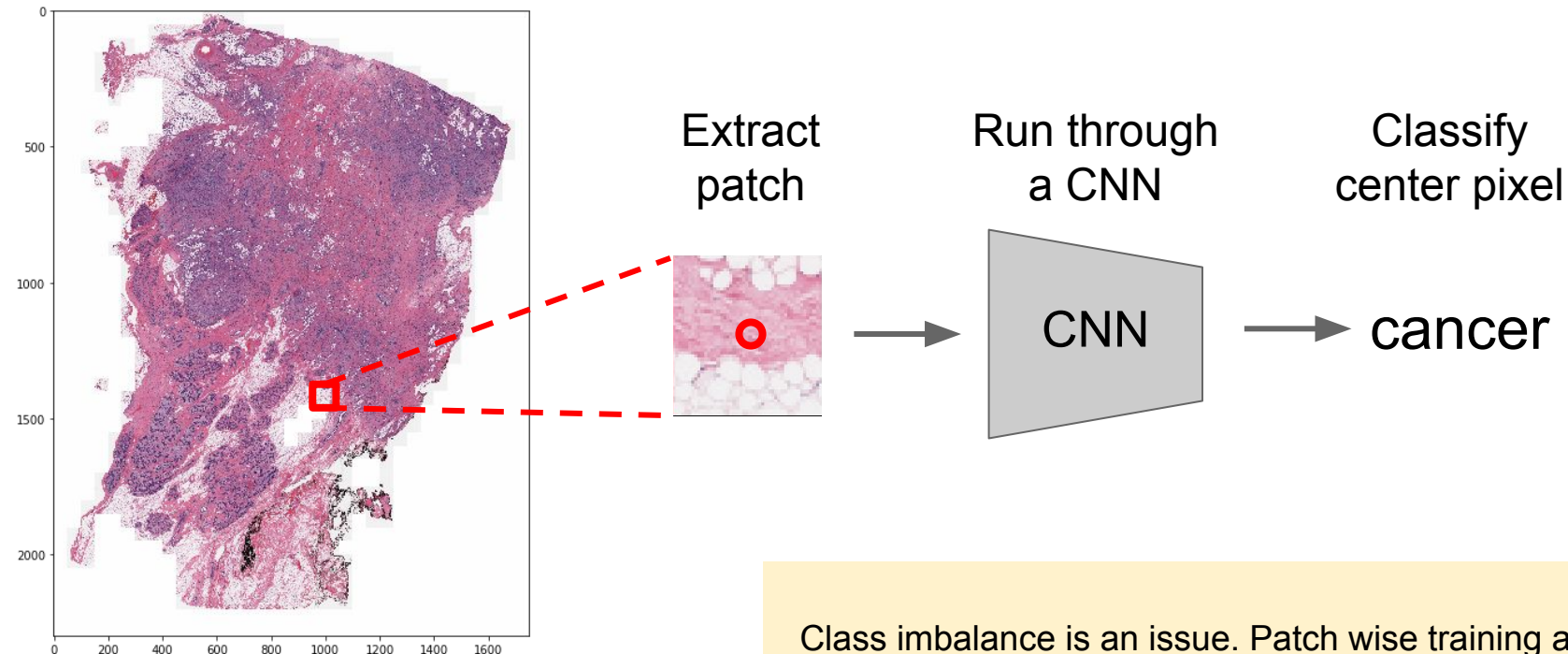


Image is chopped into patches and labelled as IDC or not

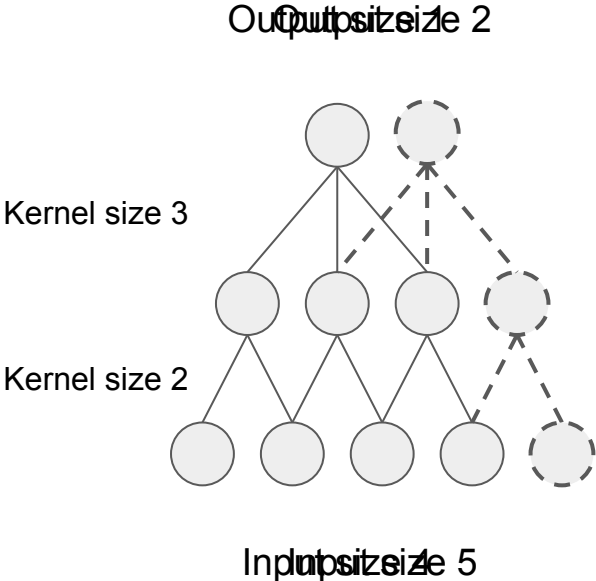
Patch wise segmentation

Use case: Invasive Ductal Carcinoma (most common subtype of all breast cancers)



Class imbalance is an issue. Patch wise training allows easy balancing of classes using standard methods.

Fully convolutional processing (FCN)



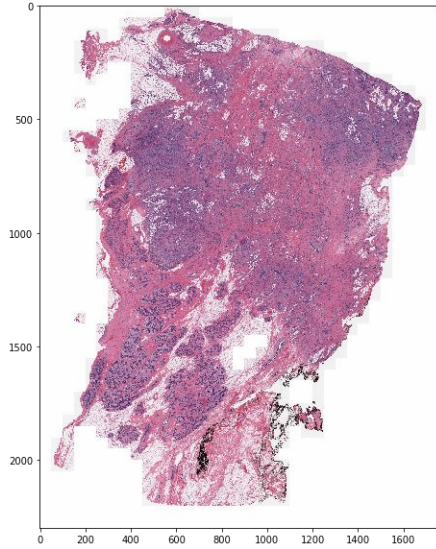
- What is this model's receptive field? 4 nodes
- How many multiplications were saved? 4
- How many saved for an input size of 6? 8

Allows for very fast inference.

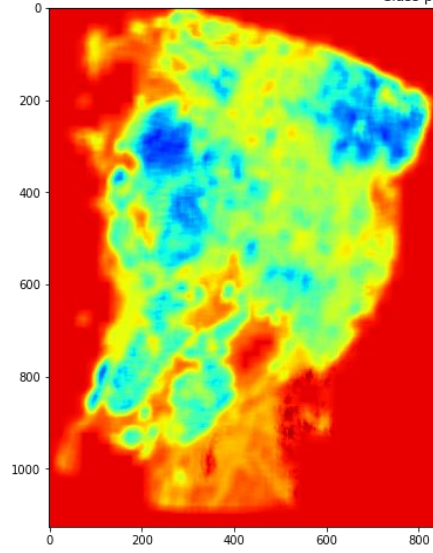
However, training this way requires a lot of memory. Need to save past outputs.

Patch wise training together with FCN inference is a good balance.

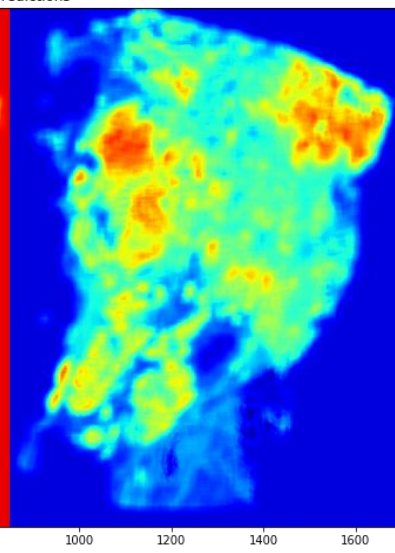
Input image



Output class 0

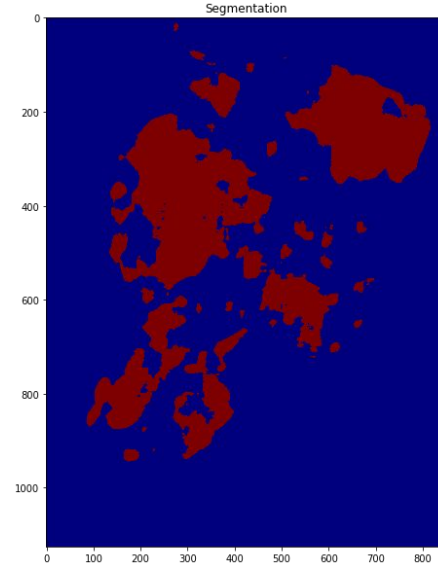


Output class 1



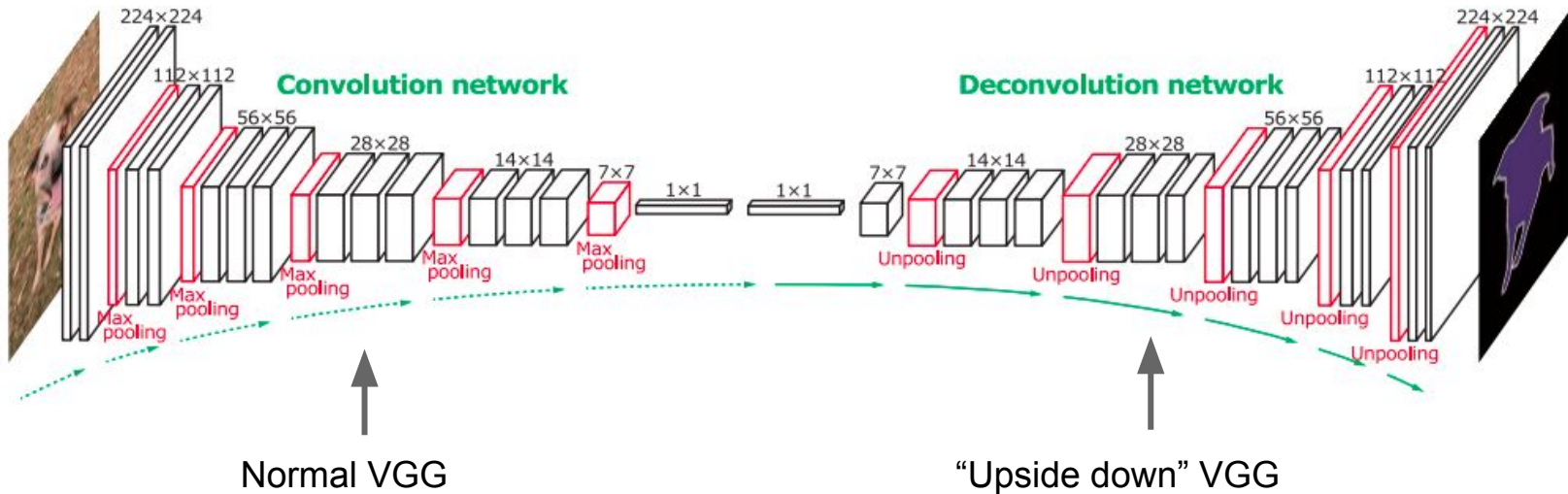
Class predictions

class 1 > class 0



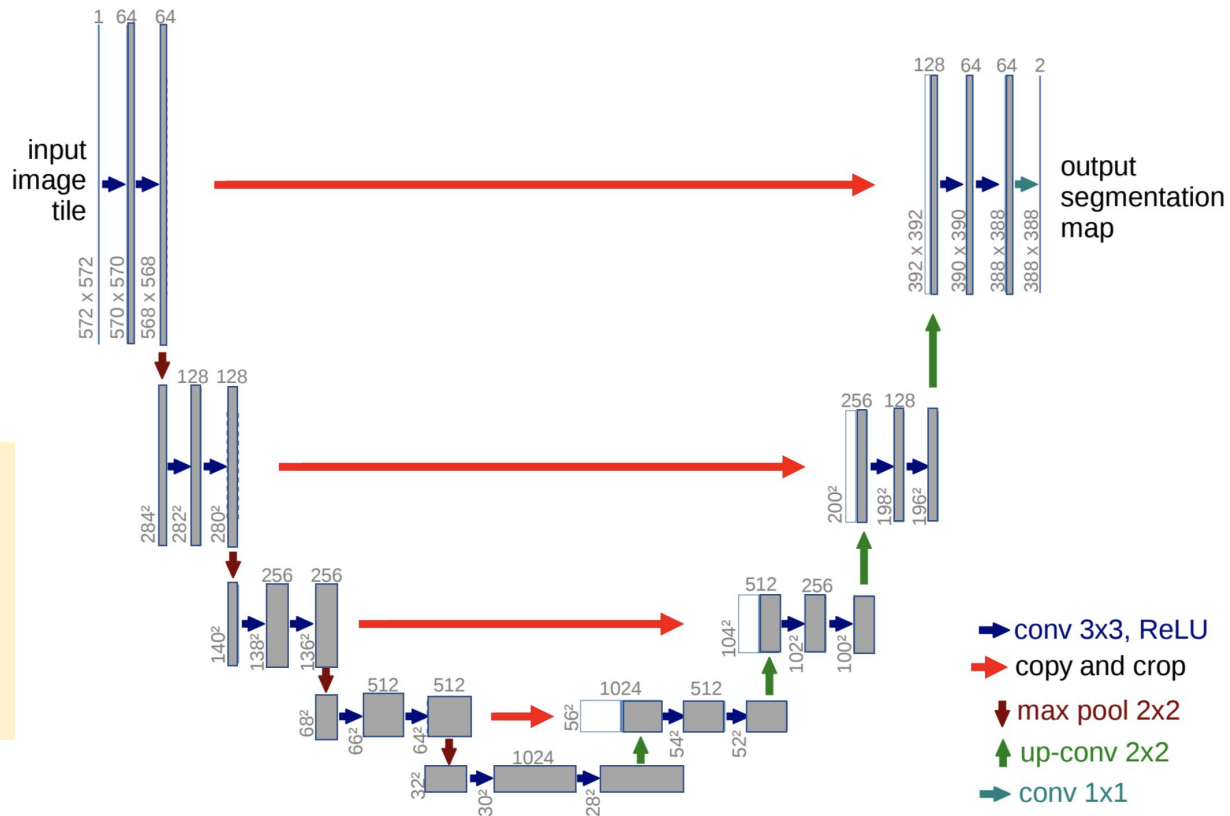
Segmentation

Recap: Segmentation using a bottleneck



- Upsampling possible with
- Unpooling
 - Transposed convolutions

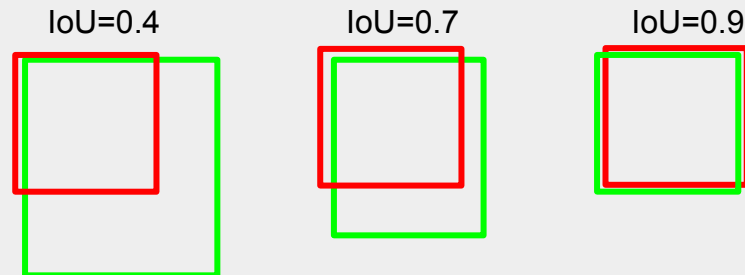
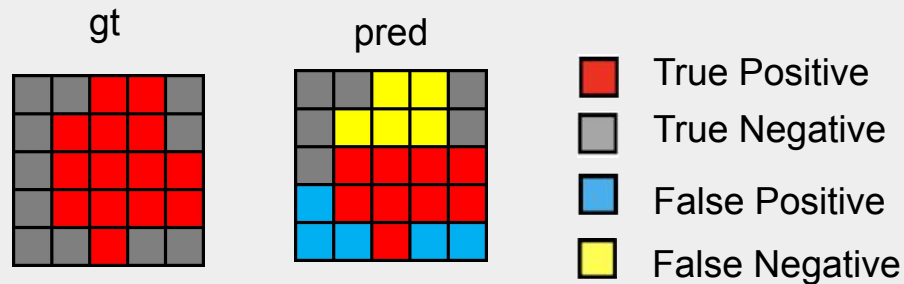
Recap: U-NET



Difference:
Skip connections (like resnet)

Dogma: skips carry spatial information, bottleneck carries high level structure.

Segmentation metrics



$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{IoU} = \text{Jaccard Index} = \frac{TP}{TP + FN + FP} = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|}$$

$$\text{Dice Coefficient} = \frac{2 \cdot TP}{(2 \cdot TP) + FN + FP} = \frac{2|X \cap Y|}{|X| + |Y|}$$

Training with dice

What maximizes the numerator?

Using the dot product to compute the intersection allows for a differentiable loss.

$$DL(p, \hat{p}) = \frac{2 \sum_i p_i \hat{p}_i}{\sum_i (p_i + \hat{p}_i)}$$

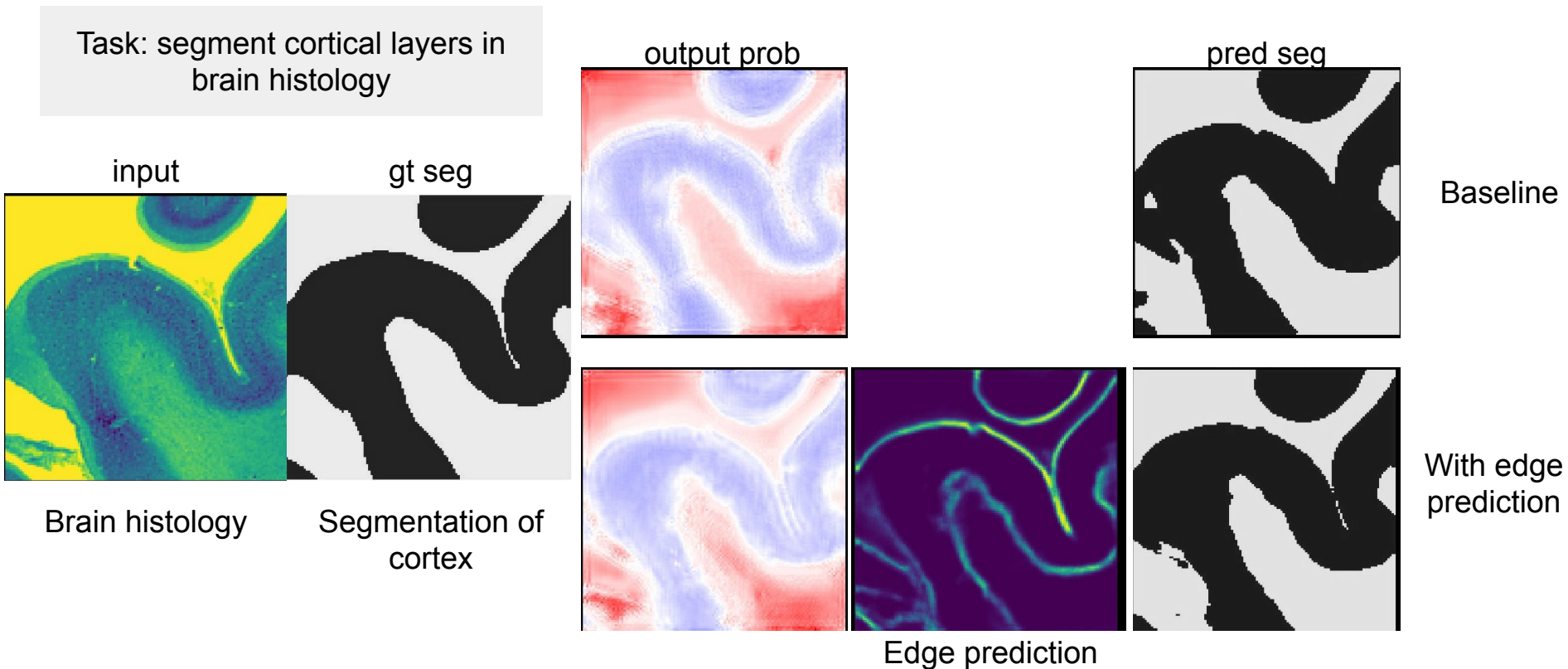
For multiple classes a basic approach is to average over all classes

$$DL_{mean}(p, \hat{p}) = \frac{1}{|C|} \sum_{c \in C} \frac{2 \sum_i p_i^c \hat{p}_i^c}{\sum_i (p_i^c + \hat{p}_i^c)}$$

Use a sigmoid or a softmax to restrict output.

Tricks: Improving edges in segmentations by predicting edges

Task: segment cortical layers in brain histology



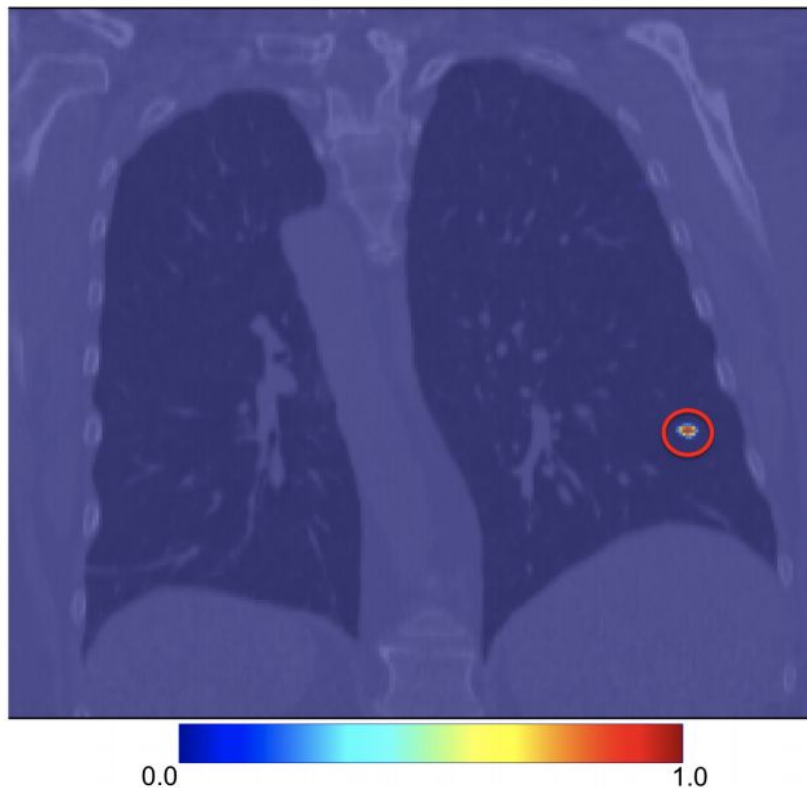
Images provided by Konrad Wagstyl (University College London) 2020

More reading about idea: [Polzounov, WordFence: Text Detection in Natural Images with Border Awareness, 2017]

Challenge: extreme class imbalance (e.g. lung nodule)

Background classes can dominate the loss and cause learning instability due to large gradients.

Balanced sampling may not work as well because patches which could yield false positives are rarely seen to train on.



CASED importance sampling for large images

General Idea:

Store a probability for each patch.

Generate patches based on this probability.

Probability is inverse of how well your model performs on that patch.

Samples are stratified by class.

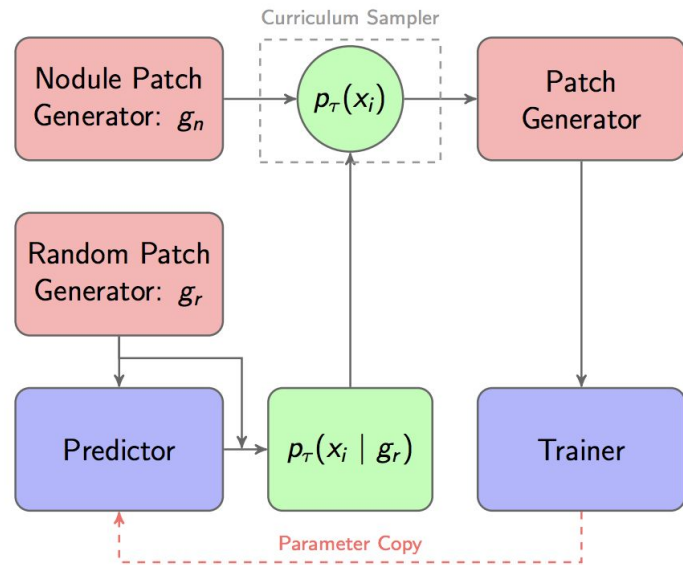


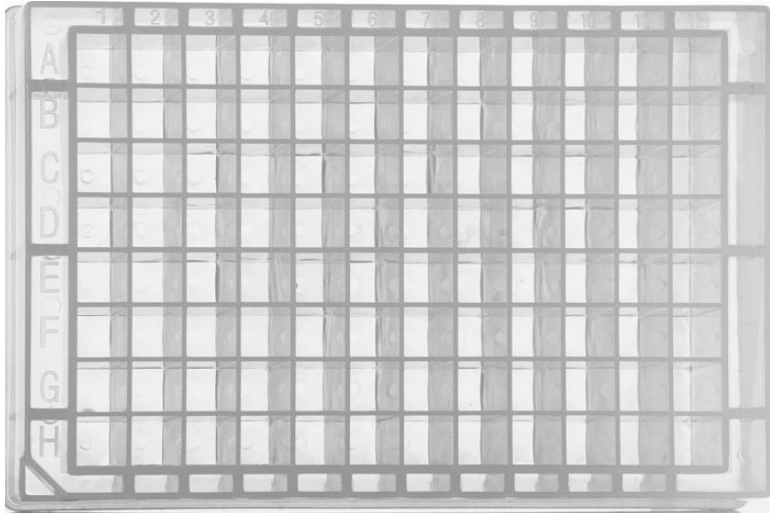
Fig. 1. Schematic diagram of CASED framework

Chapter 3

Counting

Use case: Proliferation/Cell growth studies

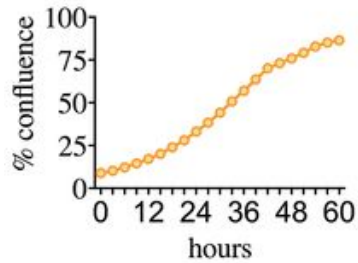
Treat cells with different compounds and observe proliferation over time



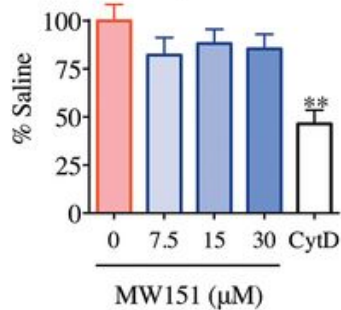
Standard 96-well plate



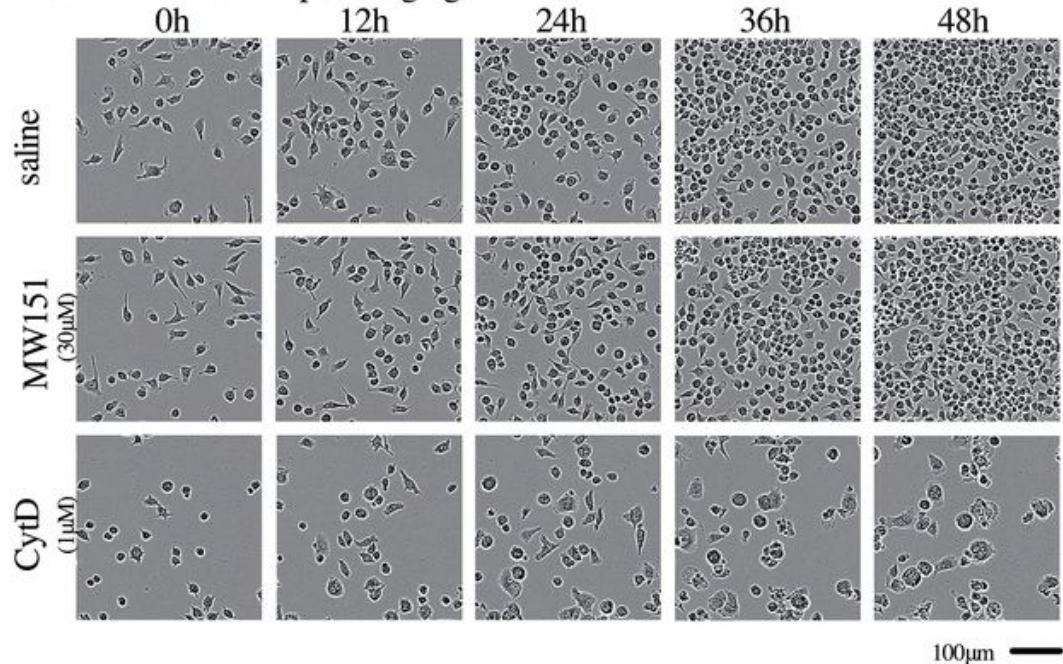
A: BV2 cell growth curve



B: cell density at 30h

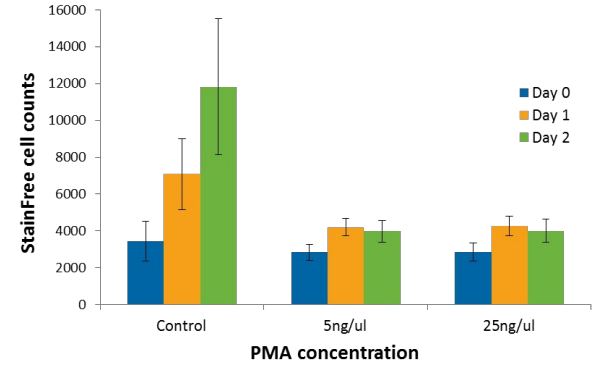
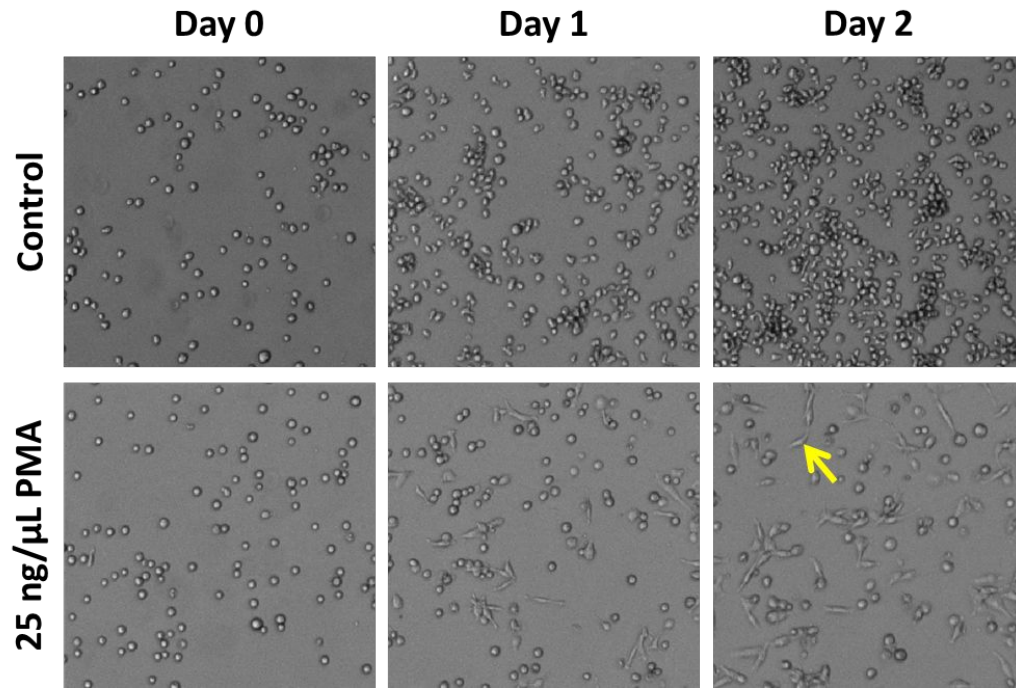


C: BV2 cell time-lapse imaging



Bachstetter, MW151 Inhibited IL-1 β Levels after Traumatic Brain Injury with No Effect on Microglia Physiological Responses, PLOS ONE, 2017

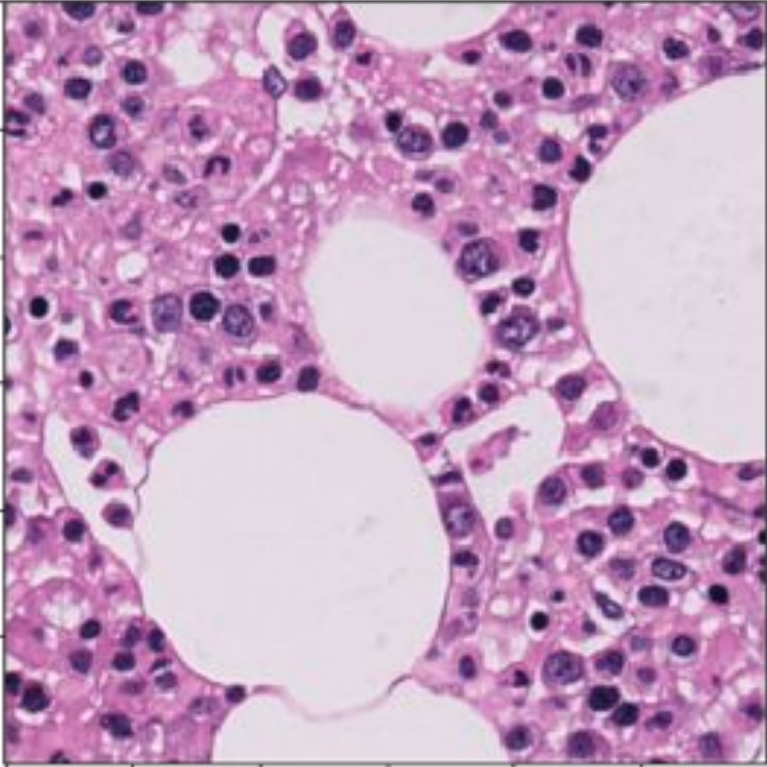
Use case: Proliferation/Cell growth studies



At the Cell Counter: THP-1 Cells, Molecular Devices
<https://www.moleculardevices.com/cell-counter-thp-1-cells>

Use case: Proliferation/Cell growth studies

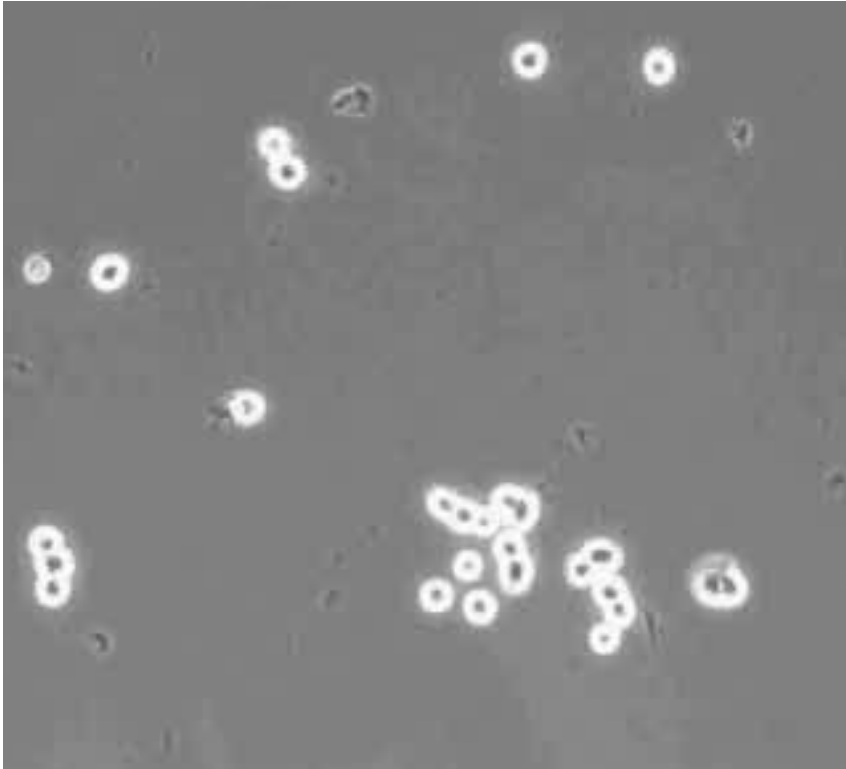
Use case: Counting in histology slides



Complicated cell structure

A 2x3 grid of six small histology images, each showing a different cell or a specific region of a cell. The cells exhibit complex internal structures, including prominent nuclei, nucleoli, and various organelles, illustrating the difficulty of counting individual cells in such a dense and complex tissue.

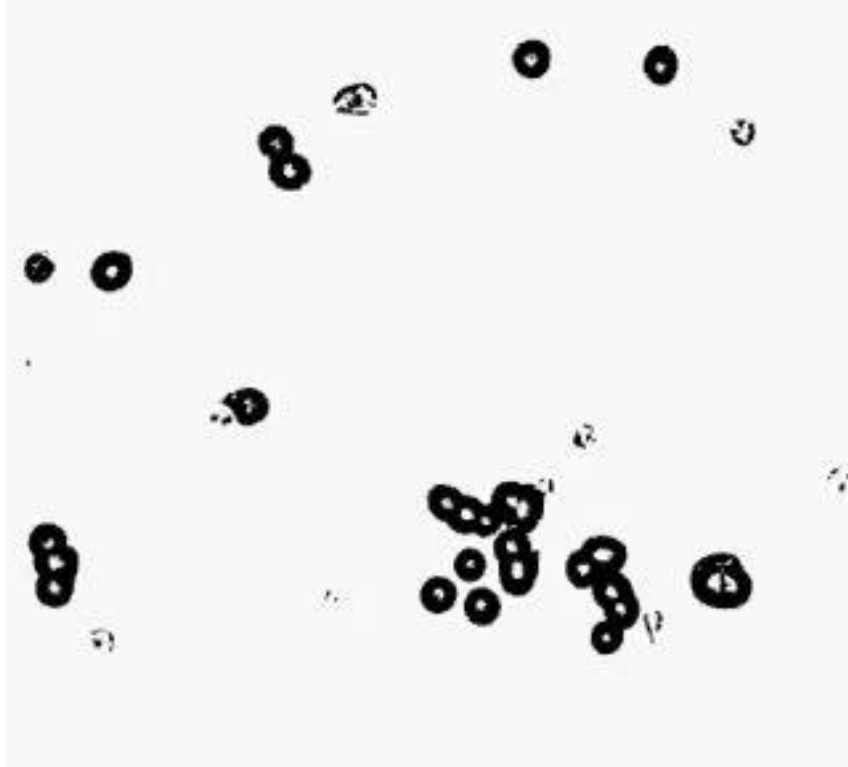
Cell counting (classic CV)



1. Create binary segmentation image
2. Watershed segmentation
3. Isolate and count

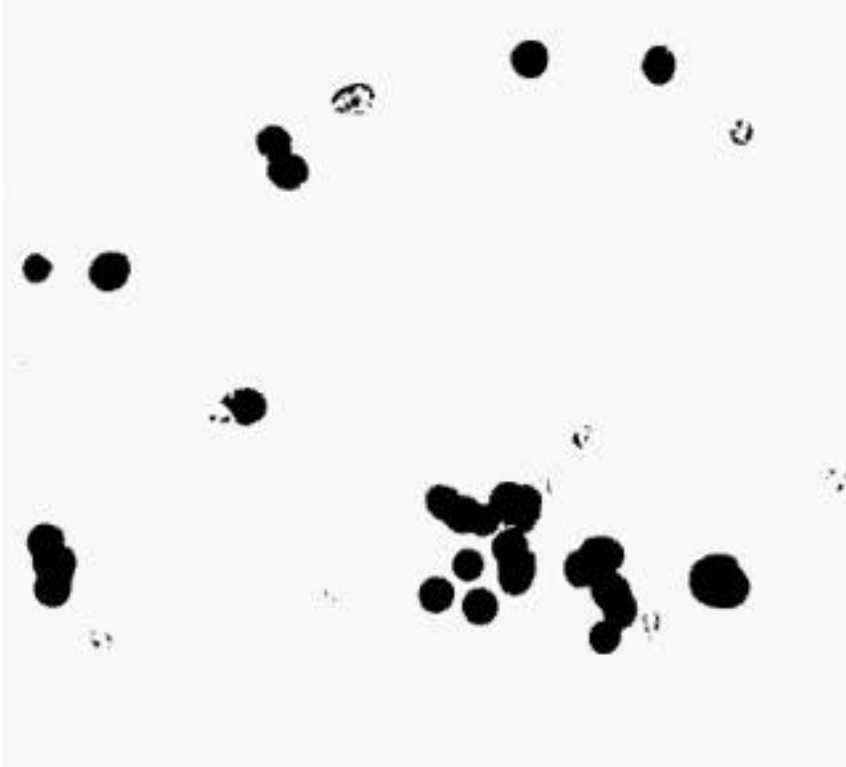


Cell counting (classic CV)



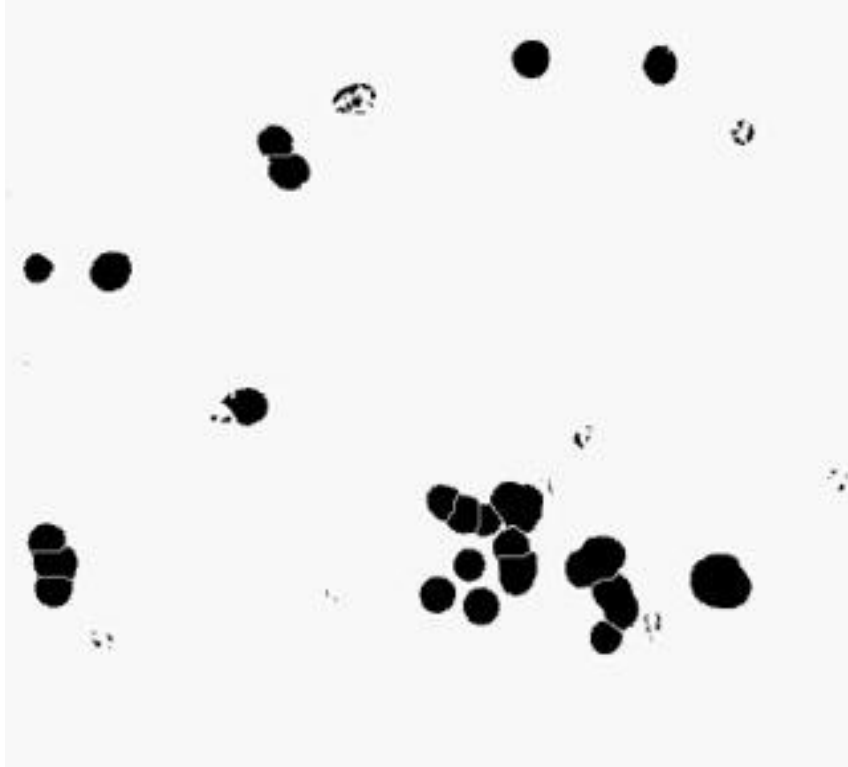
1. **Create binary segmentation image**
2. Watershed segmentation
3. Isolate and count

Cell counting (classic CV)



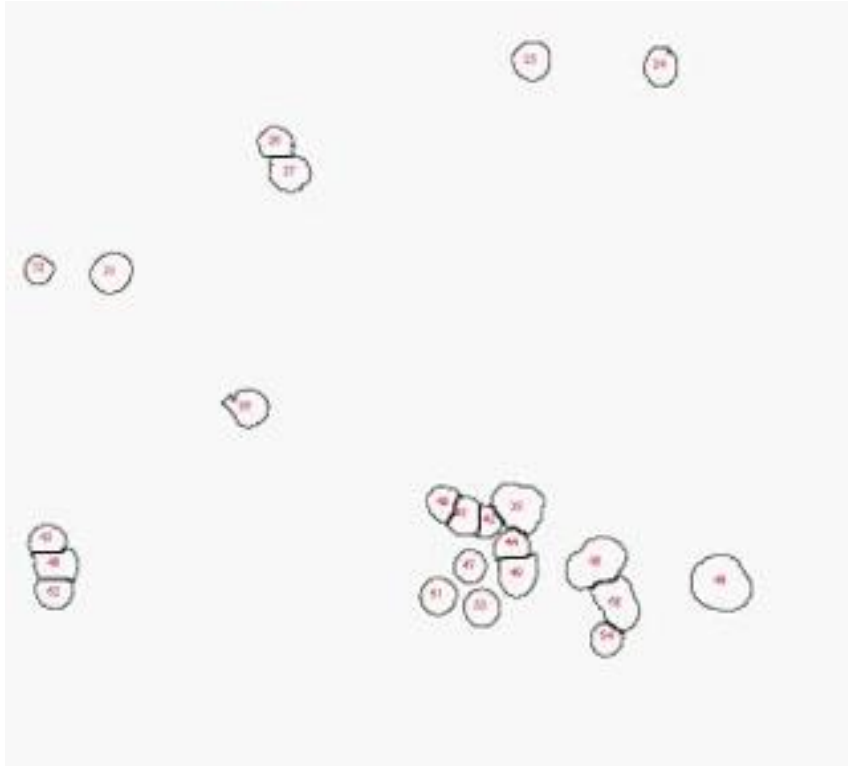
1. **Create binary segmentation image**
2. Watershed segmentation
3. Isolate and count

Cell counting (classic CV)



1. Create binary segmentation image
2. **Watershed segmentation**
3. Isolate and count

Cell counting (classic CV)

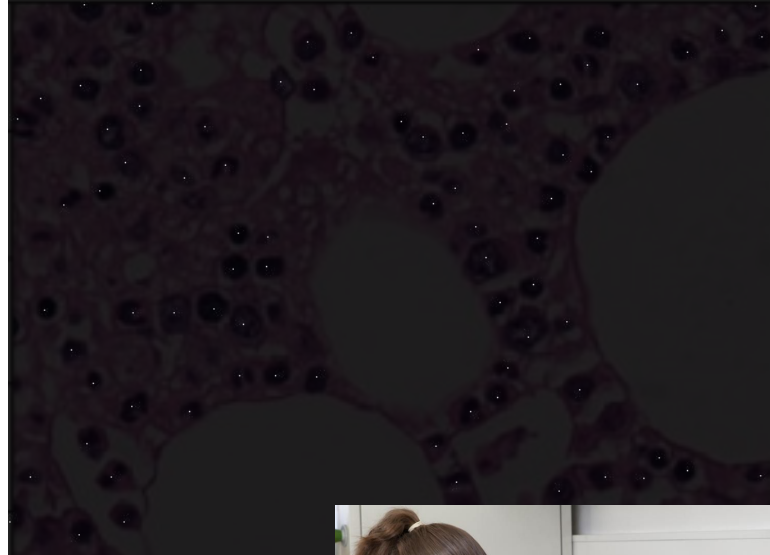


1. Create binary segmentation image
2. Watershed segmentation
3. **Isolate and count**

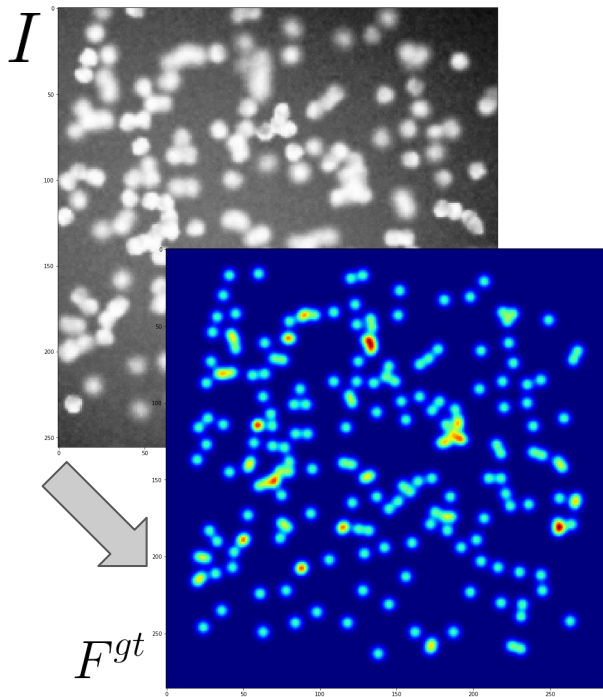
This works well on easy tasks but doesn't scale.

"Pipelines" end up breaking on new images with different lighting or stain.

How to get labels?



Counting via Segmentation



Targets for regression

$$F^{gt}(x, y) = \sum_i^{\# \text{ Cells}} \mathcal{N}([x, y]; [x_i, y_i], \sigma^2)$$

Sigma is typically small
like a few pixels

Train model to regress

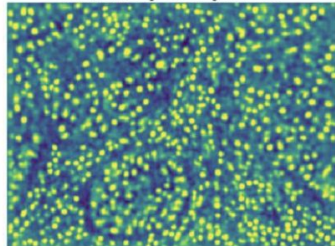
$$L = \sum_{x,y} |F^{gt}(x, y) - F(I)(x, y)|$$

To recover count:

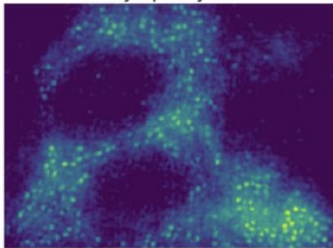
$$\text{count} = \sum_{x,y} F(I)(x, y)$$

Multiple output classes

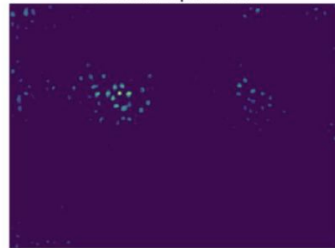
Probability of any nuclei



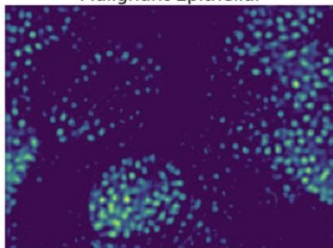
Lymphocyte



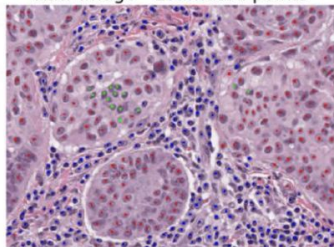
Normal Epithelial



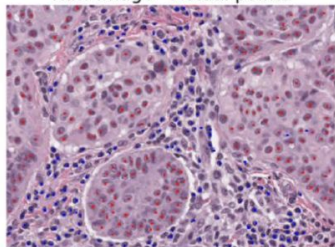
Malignant Epithelial



Raw Image - Estimated points



Raw Image - Actual points



Counting and classifying also possible using multiple output channels.

Combine losses together

$$\lambda_1 L_{lymph} + \lambda_2 L_{norm} + \lambda_3 L_{mal}$$

Max prediction over output channels for each cell identified

Count and classify different cell types [Bidart 2018]

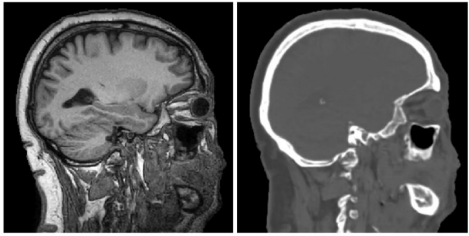
Chapter 4

GANs

Medical image-to-image translation considered harmful

Many papers have proposed methods that can "translate between modalities"

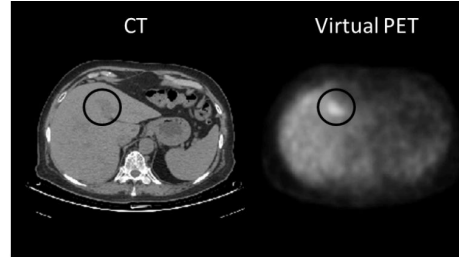
MR -> CT



I_{MR}

$Syn_{CT}(I_{MR})$

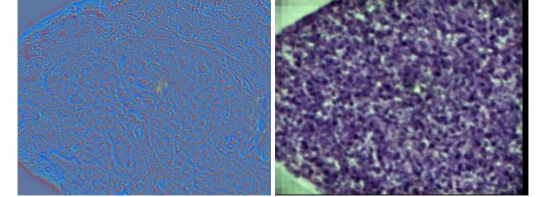
CT -> PET



CT

Virtual PET

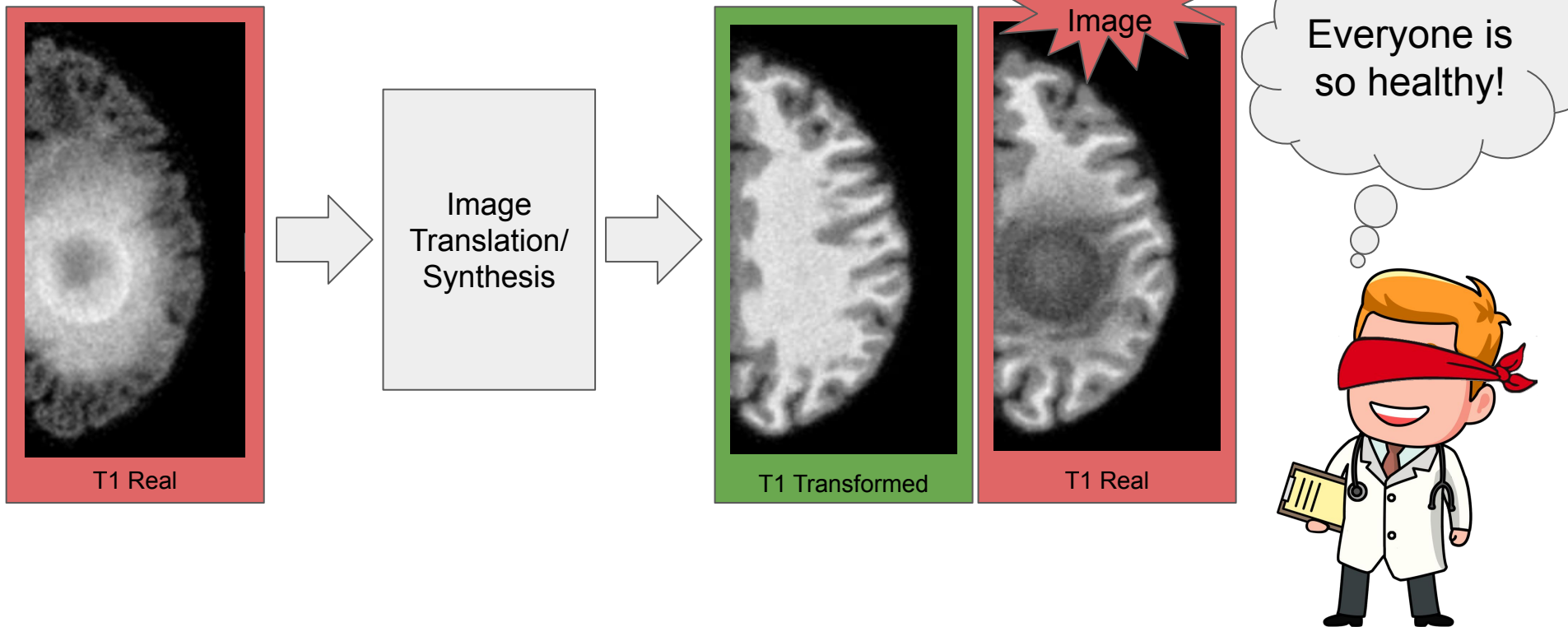
Synthesized H&E staining



Adversarial losses are very good at distribution matching (e.g. CycleGAN).
But artifacts could be introduced and then used in diagnosis which can be dangerous.

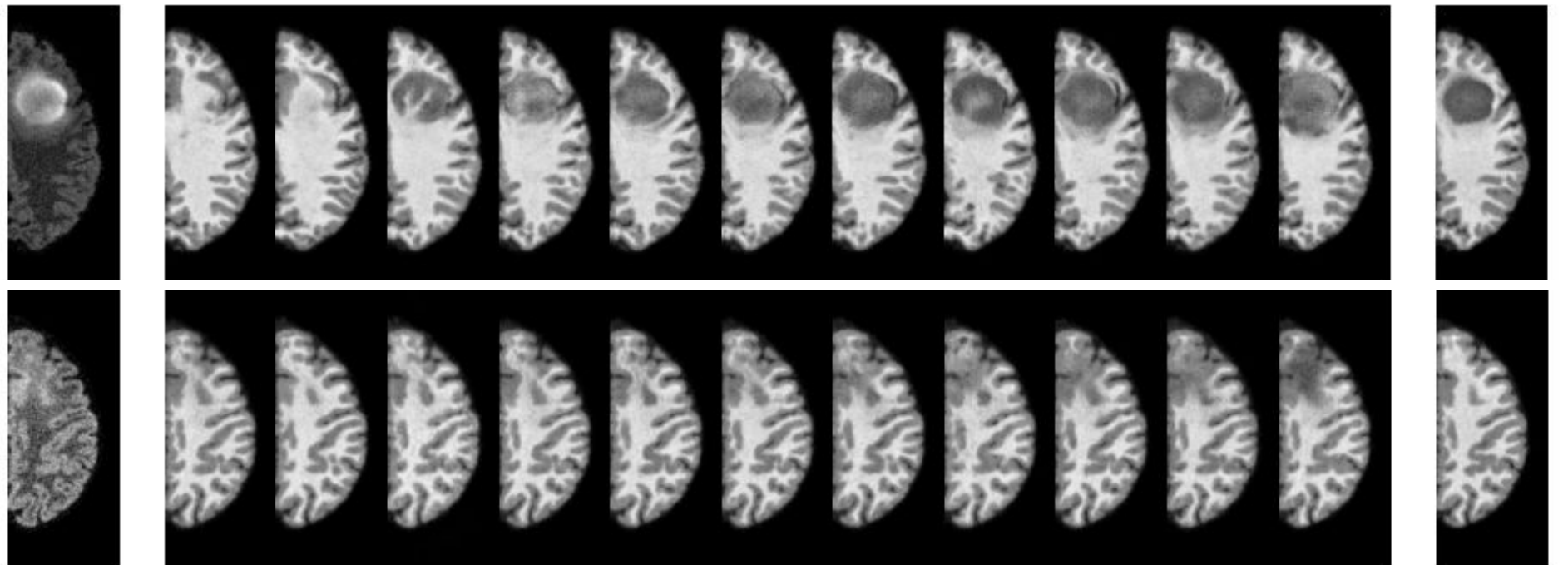


Use case: MRI modality transformation



But a bias in training data can lead to incorrect translation

CycleGAN results



Flair
(source)

0% 10% 20% 30% 40% 50% 60% 70% 80% 90% 100%

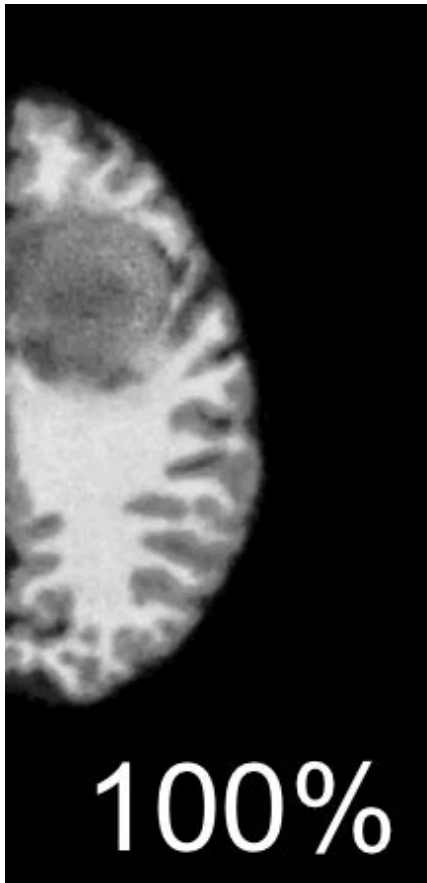
T1
(target)

% training data with tumor

Real Flair



Biased Transformations



Real T1



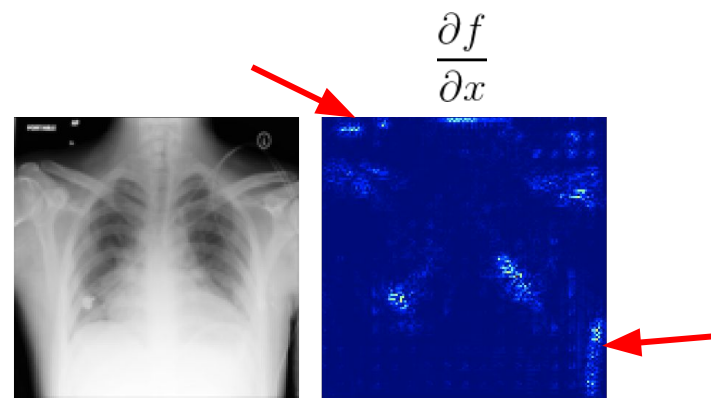
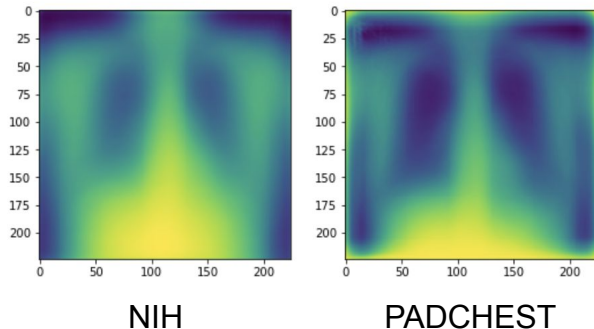
Chapter 5

Right for the right reasons

Incorrect feature attribution

Models can overfit to confounding variables in the data.

Example: Systematic discrepancy between average image in datasets



Overfitting while predicting Emphysema [Vivano 2019]

- Merging datasets with different class imbalance (confounding artifacts from each hospital)
- Labels confounding with each other
- Demographics confounding with labels

[Zeck, Confounding variables can degrade generalization performance of radiological deep learning models, 2018]

[Vivano, Underwhelming Generalization Improvements From Controlling Feature Attribution, 2019]

[Simpson, GradMask: Reduce Overfitting by Regularizing Saliency, 2019]

[Ross, Right for the Right Reasons, 2017]

Mitigation approaches

Feature engineering

- **Range normalization** ($/\max$)
- **Subspace alignment** (align data using their eigenbases)

During training

- **Reverse gradient** [Ganin & Lempitsky, Unsupervised Domain Adaptation by Backpropagation, 2014]
- **Right for the Right Reasons regularization** [Ross, Hughes, & Finale Doshi-Velez, 2017]
- **GradMask contrast loss** [Simpson, 2019]
- **ActivDiff** [Viviano, 2019]

What if feature artifact is correlated with target label?
Is the reason that should be used for prediction known?
What if it is not known?

Not discussed

Image Registration

Cell morphology representation (e.g. BBBC021)