Special for MVA course

# Deep Learning in practice: MammoScreen

Yaroslav Nikulin,
Senior Research Scientist

January, 21 2019

➢Therapixel

➢DL -> radiology

➢Breast cancer

➢DM DREAM Challenge

# Therapixel: Medical Image Understanding

**Therapixel**

**2019** — Clinical Study comparing MammoScreen to radiologists

**2018** — Mamm⚲screen
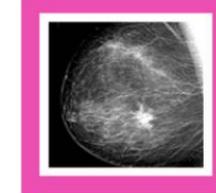
**2017** — 1st place DREAM DM
5th place Kaggle Data Bowl

**Kaggle Bowl**　　**DREAM DM**

**2016** — AI research

**2015** — Visualization SW

**2013** — Founded

inria inside

Olivier Clatz, PhD　　　Pierre Fillard, PhD

# Breast Cancer Screening: some key stats

➢ 33M exams/year = 132M images in US alone
➢ $7.8 billion - cost of mammography screening in US (2010)
➢ 120 sec: average interpretation time.

**1 out of 8**
Woman affected during her lifetime

**10** recall for
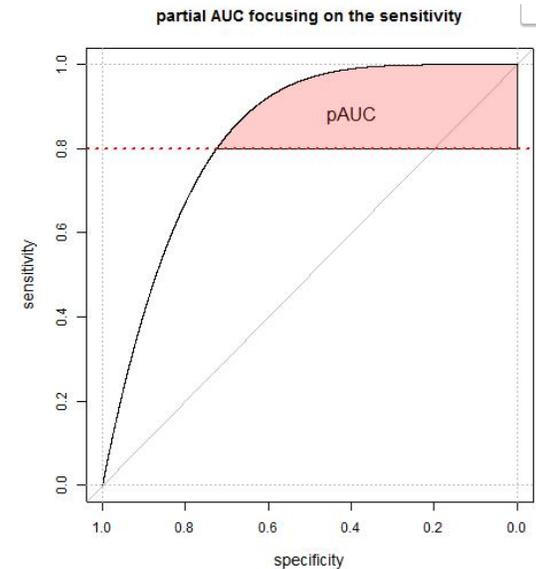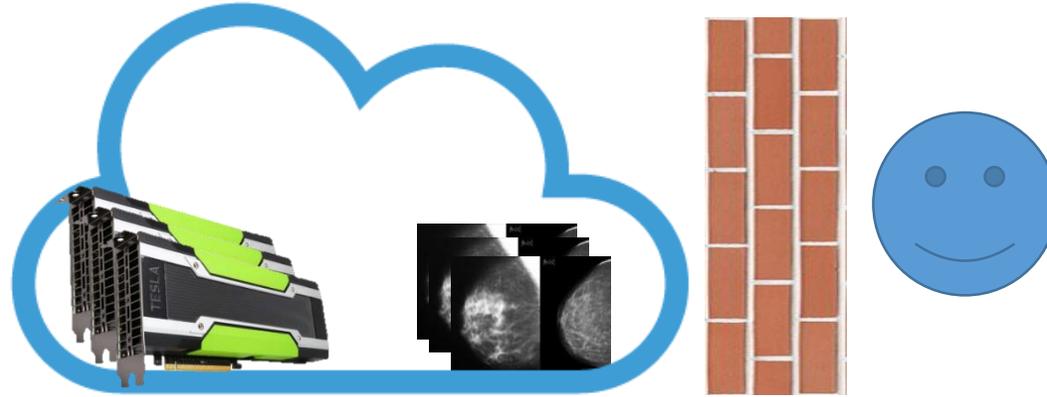**100** screened

**5** cancers for
**1000** screening

"If a typical person can do a mental task with less than one second of thought, we can probably automate it using AI either now or in the near future."

Andrew Ng, 2016
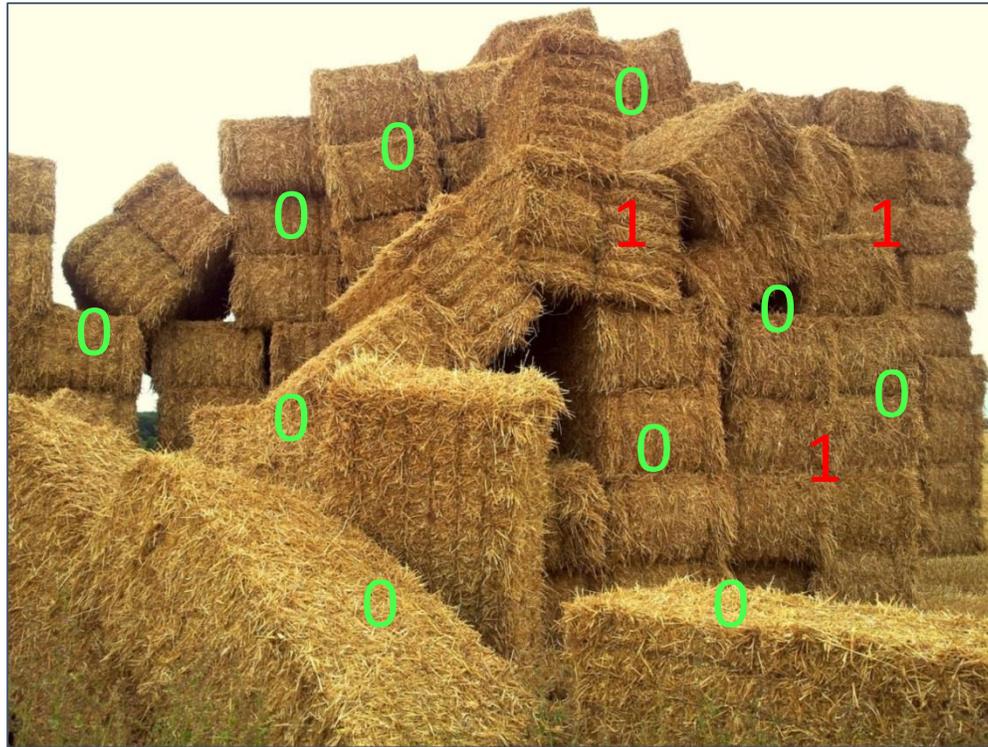
# The Digital Mammography DREAM Challenge

## Challenge setting:

- ➤ Completely in the cloud

- ➤ 22 CPU cores **+ 2 GPUs**

- ➤ 14 days / per team
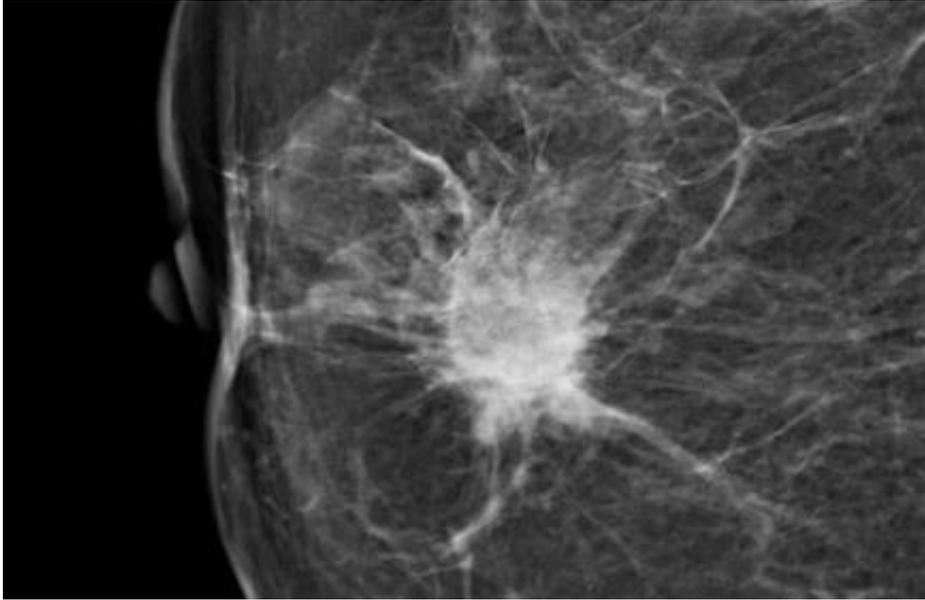
- ➤ Performance measure:
  **AUC** and **partial AUC**



partial AUC focusing on the sensitivity

pAUC

sensitivity

specificity

➤ 320k images
➤ Only 1548 (**0.47%!**) positive examples
➤ High resolution: from **3328x2560 to 5928x4728**.
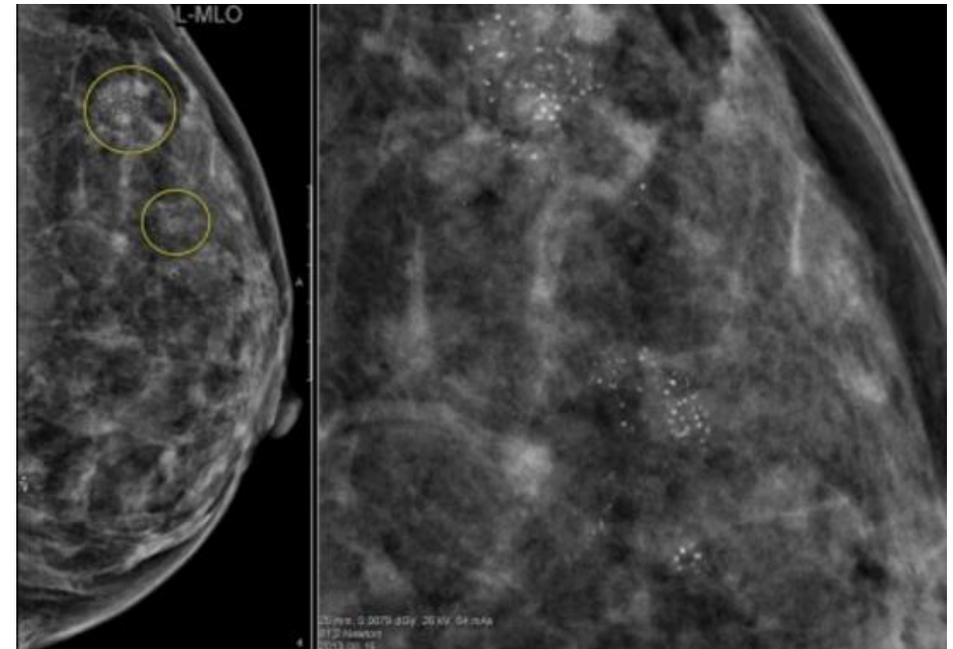➤ One single label per image: 0 or 1



**Now look for a needle in them…**

➢ Different kinds of anomalies: calcifications, masses, distorsions



➢ Different scales of anomalies: from micro-calcifications to big cancerous masses.

**Can be malignant OR benign!**

➢Data specificity

➢Dense annotations

➢Patch model

➢Image model

➢ Visualization

In our approach, limited by several factors.
Actually 3-5 times higher

➢ Resolution: 1200x800 **vs** 224x224

➢ Zone of Interest :  < 1% **vs** > 50%

➢ Number of classes : 2 **vs** 1000

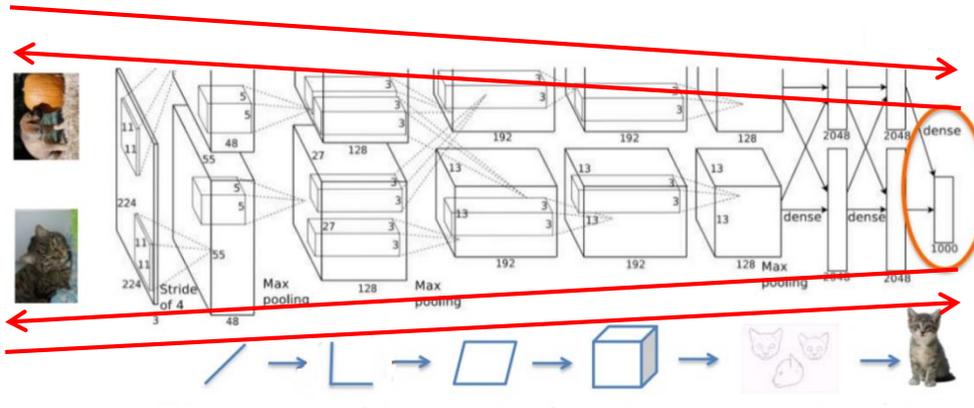➢ Highly imbalanced **vs** roughly balanced
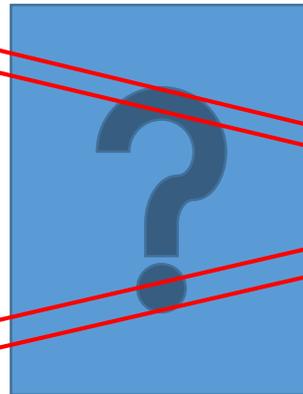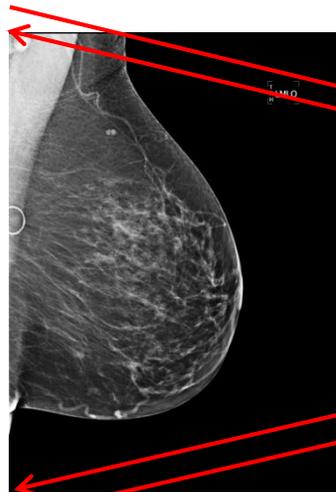
Zone of Interest

AlexNet (Krizhevsky et al. 2012)

Input size:
224x224

Output : 1 out of 1000
~ 10 bits of information

Input size:
~3500x2500

0 or 1

Output : 1 out of 2
= 1 bit of information

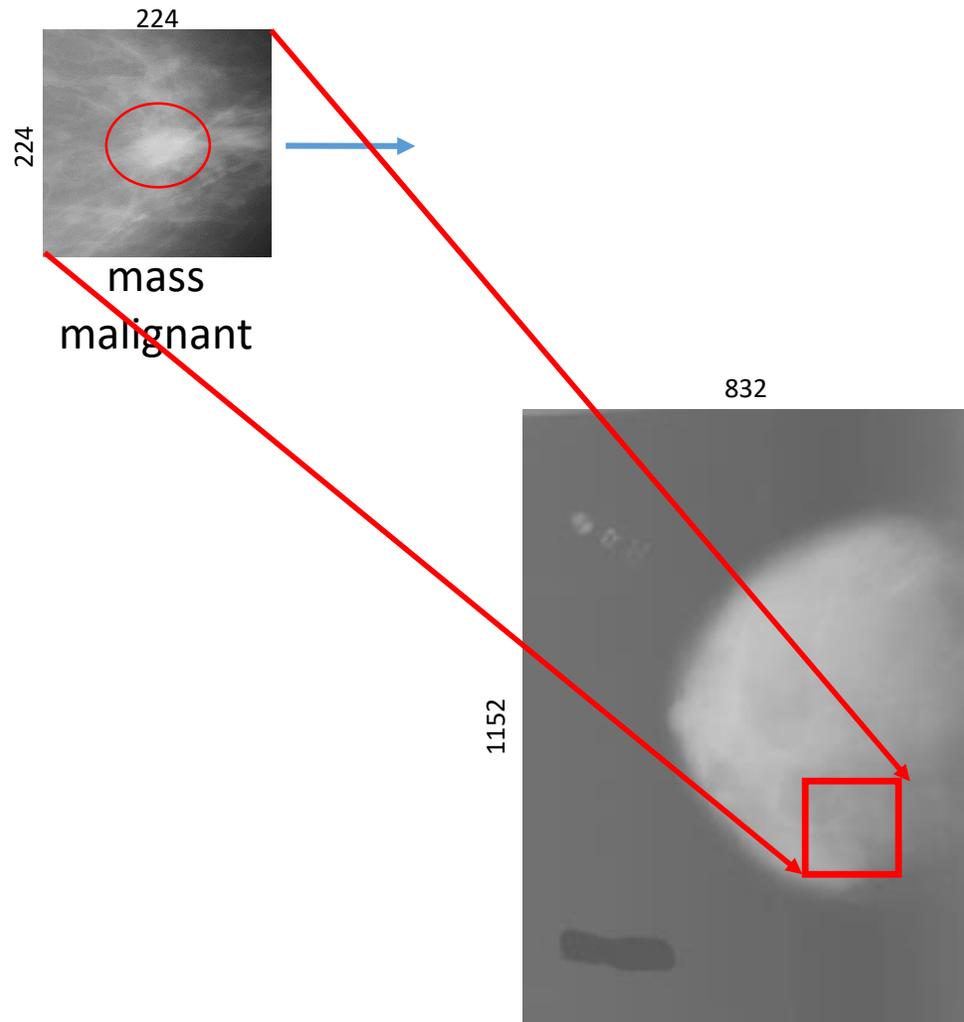| | DDSM | DREAM |
|---|---|---|
| Total im | 10k | 320k |
| Positives | 1807 | 1548 |
| Info | mask&type | 0 or 1 |

It would be great to:

➢ Make use of local info
➢ Make use of lesion type
➢ Still be able to train on DREAM

mass
malignant

224
224
mass
malignant

832
1152

0 — healthy
1 — calc ben
2 — mass ben
3 — calc mal
4 — mass mal

# Patch model: Fully Convolutional Network



112x112x32  56x56x64  28x28x128  14x14x256  7x7x256  3x3x512

conv1, 32 | conv2, 32 | pool 2x2 | conv3, 64 | conv4, 64 | pool 2x2 | conv5, 128 | conv6, 128 | pool 2x2 | conv7, 256 | conv8, 256 | pool 2x2 | conv9, 256 | conv10, 256 | pool 2x2 | conv11, 512 | conv12, 512 | pool 2x2

FC as FCN

1024  512  5

Who said VGG..?    Total: 11M of parameters

224
224
mass malignant

0  healthy
1  calc ben
2  mass ben
3  calc mal
4  mass mal

832
1152

# Patch model: Fully Convolutional Network



224
224
mass malignant

112x112x32
56x56x64
28x28x128
14x14x256
7x7x256
3x3x512

conv1, 32 | conv2, 32
pool 2x2
conv3, 64 | conv4, 64
pool 2x2
conv5, 128 | conv6, 128
pool 2x2
conv7, 256 | conv8, 256
pool 2x2
conv9, 256 | conv10, 256
pool 2x2
conv11, 512 | conv12, 512
pool 2x2

FC as FCN
1024
512
5

Who said VGG..?

Total: 11M of parameters

0 — healthy
1 — calc ben
2 — mass ben
3 — calc mal
4 — mass mal

832
1152

11
16

5

Long, Shelhamer, Darrell,
Fully Convolutional Networks for Semantic Segmentation, 2014

Information Bottleneck..?

15

All the convolutional kernels have spatial size 3x3

All the pooling layers are max pool

**Intermediate labels:**
0 – healthy
1 – calc benign
2 – mass benign
3 – calc malignant
4 – mass malignant

**Final labels:**
0 – healthy
1 – cancer

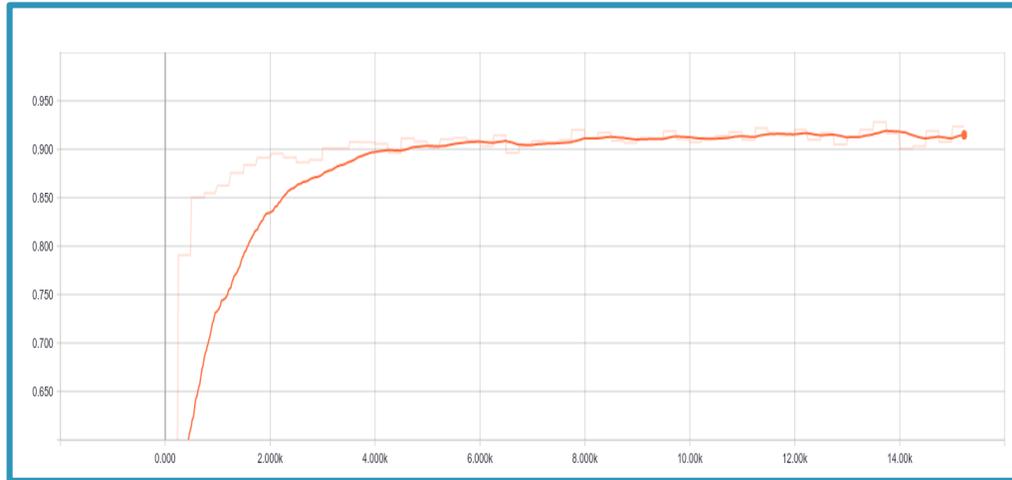- **Detector Net: pretraining by patches, ~2 hours on 4 Titan X**

- **End-to-end finetuning by images, ~20 hours on 4 Titan X**

Important to train on images:
➢ Final pool 5x5
➢ Adjust learning rate
➢ Linear shortcut
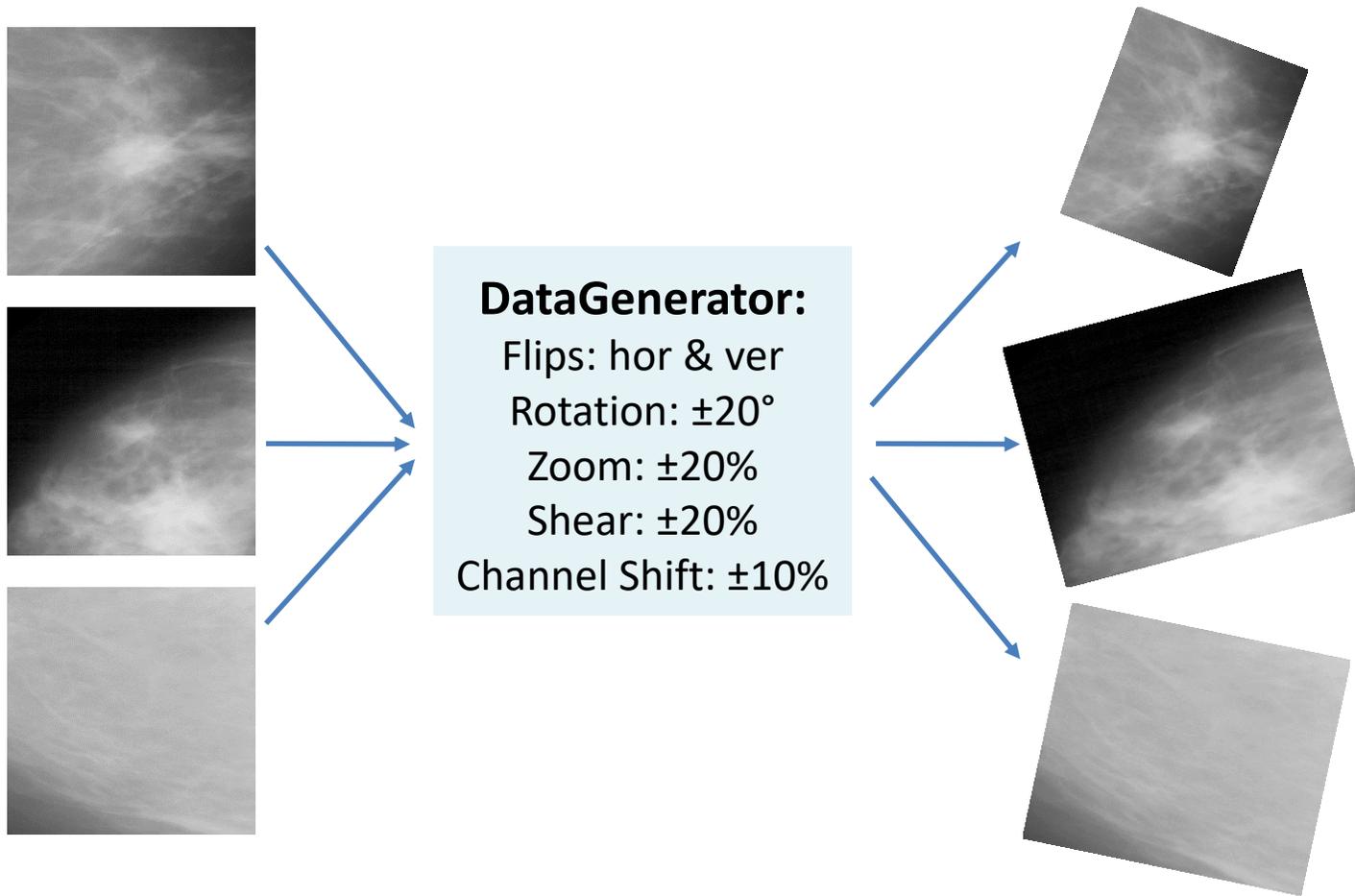
**AUC per breast (DDSM)**



**Loss**



default

➤ DetectorNet on patches from scratch: Adam, lr 0.001

➤ Restore DetectorNet weights and Adam variables

➤ On images (partially restored): Adam, lr 0.0001

➤ Send it to the cloud and use as a starting point

➤ Finetuning on DREAM data: Adam, lr 0.0001 and Exponential Moving Averages (0.9)

➤ Restore EMA (0.9), finetune with SGD, lr 0.0001

Why 0.9? Seems to be near optimal for AUC optimization (~+1%) given the number of positives divided by batch size.
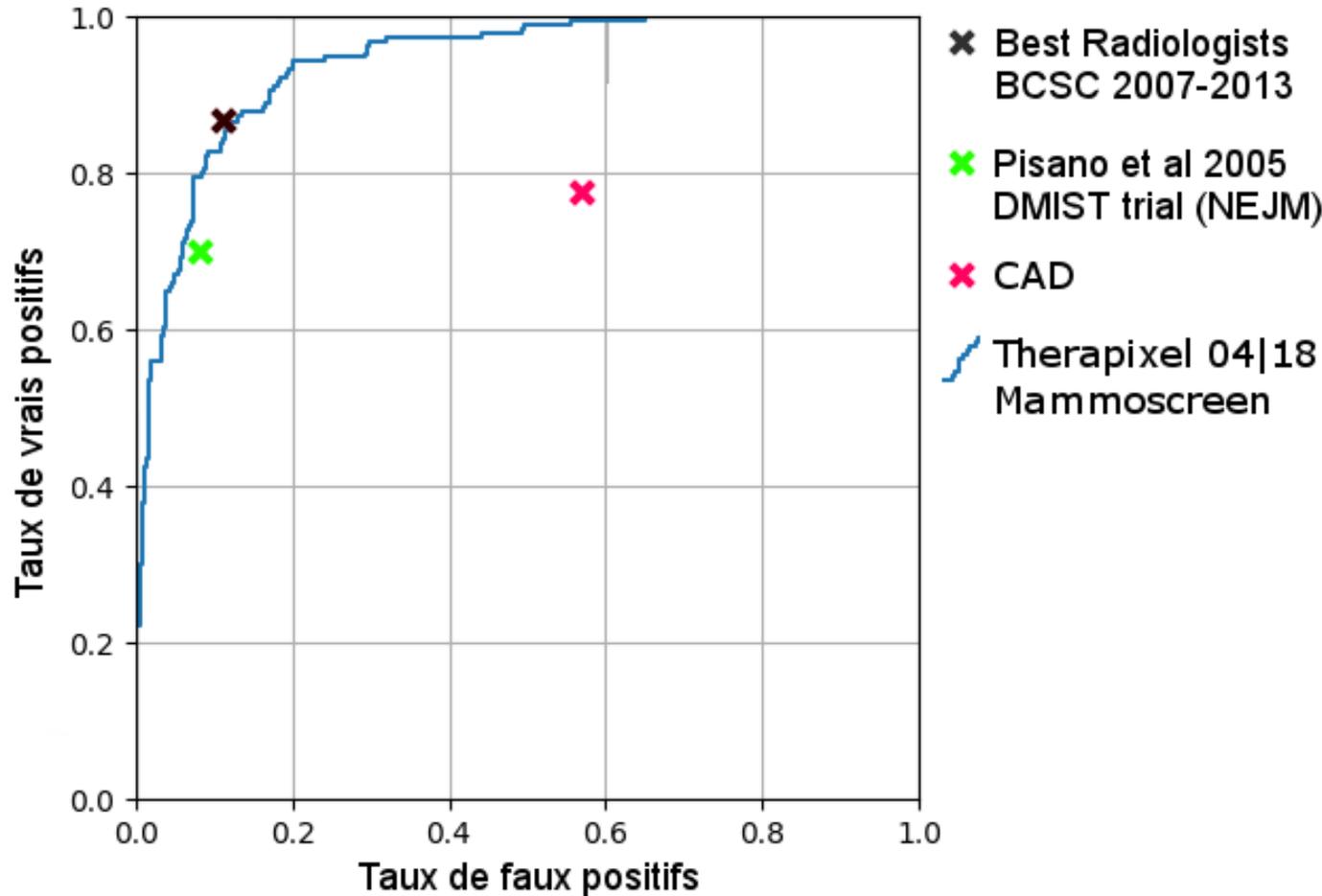
$$0.9^{125}= 2 \cdot 10^{-6}$$
$$0.99^{125} = 0.28$$

**DataGenerator:**
Flips: hor & ver
Rotation: ±20°
Zoom: ±20%
Shear: ±20%
Channel Shift: ±10%

➢ Batches are balanced

➢ Data Augmentation is crucial

➢ It also helps during the inference (4 flips → ~+1%AUC)

➢ Averaging everything works well

**Legend:**
- ✖ Best Radiologists BCSC 2007-2013
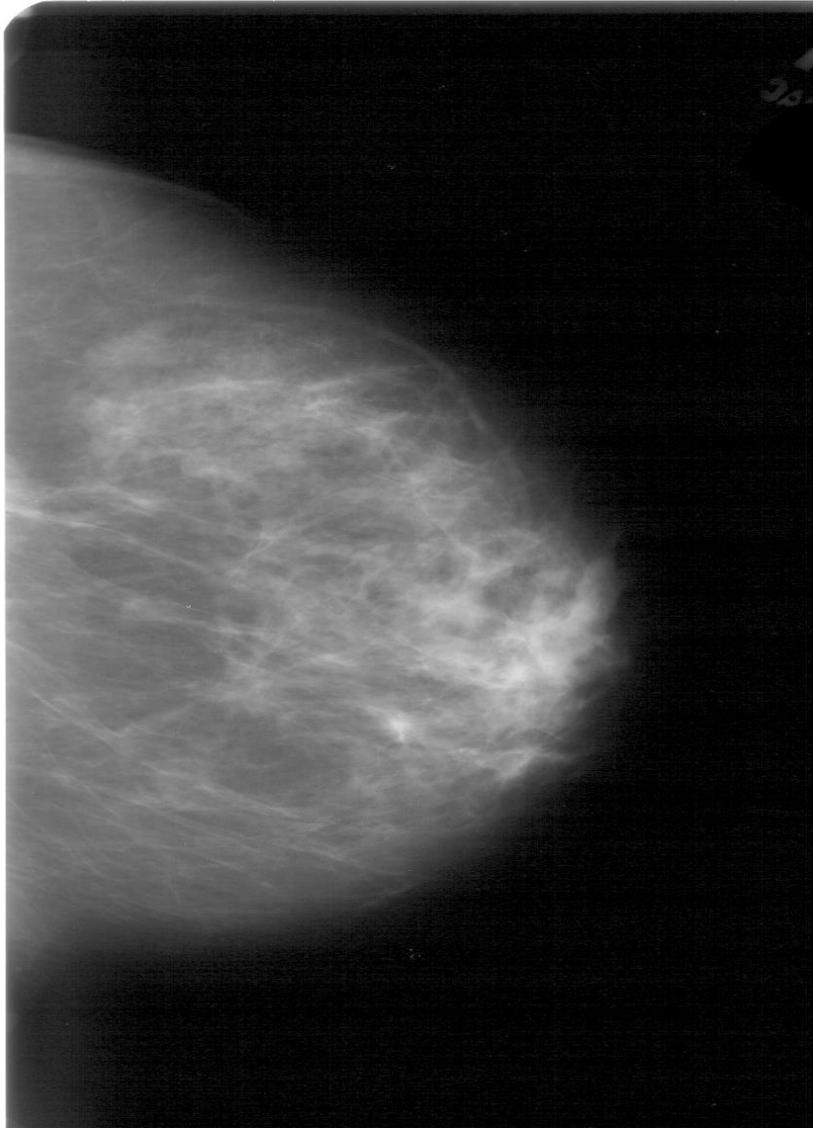- ✖ Pisano et al 2005 DMIST trial (NEJM)
- ✖ CAD
- ╱ Therapixel 04|18 Mammoscreen

A note on overfitting and "advertising" stats:

➢ Overfitting happens on several levels:
1. training data
2. validation data
3. test data = overfit dataset
4. overfit a particular problem
5. overfit a particular domain (?)
6. overfit human style of thinking (??)

➢ In particular, performance of DL model on mammographies depends on:
1. Device used for mammography
2. Skills of technician
3. Screening period (1-1.5-2 years)
4. Positive/negative ratio, closely linked to
5. Fraction of truly difficult cases
6. Population (country)
7. ...

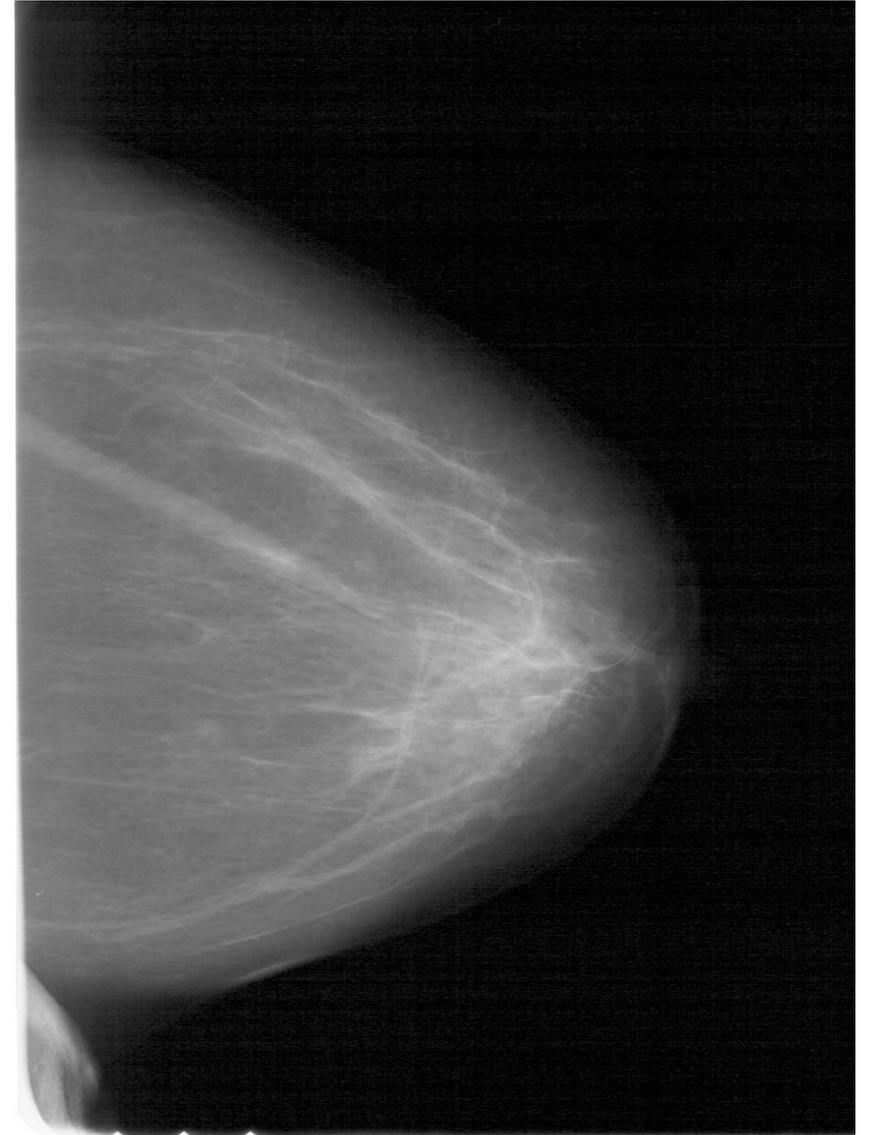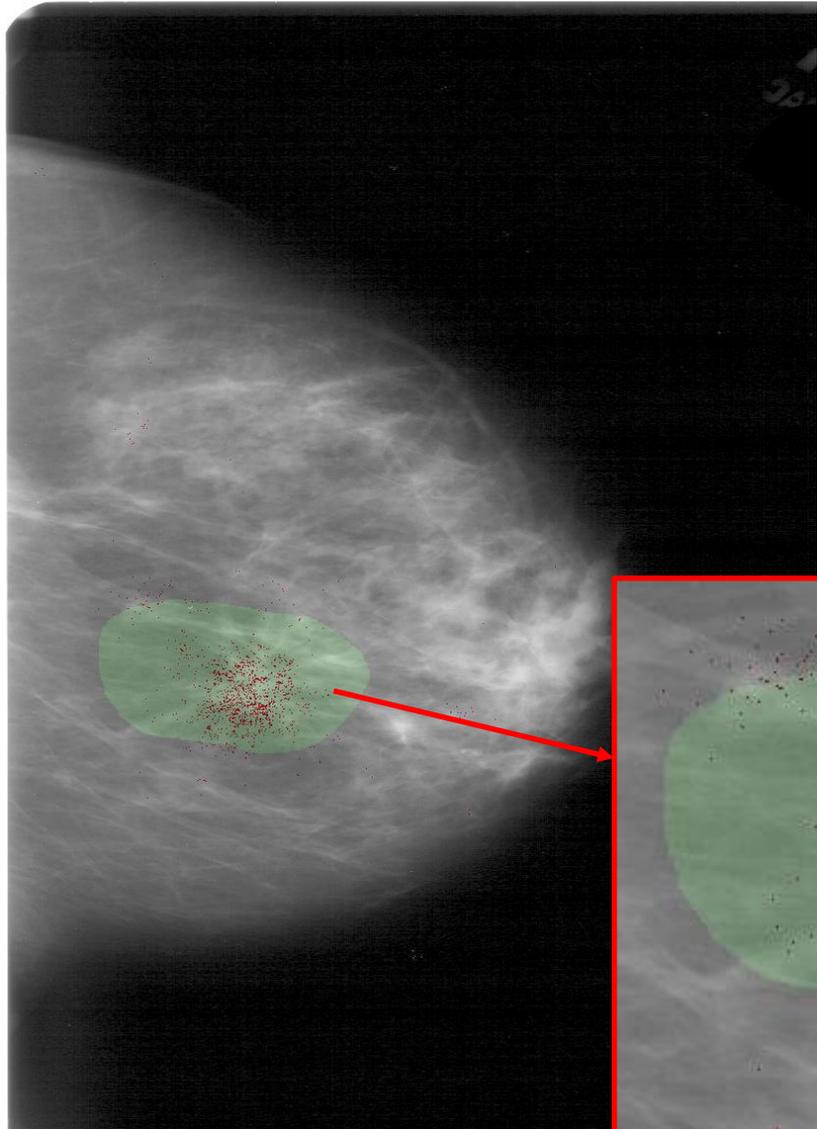# Saliency maps for weak detection



label = 1

↓

Cancer. But where?
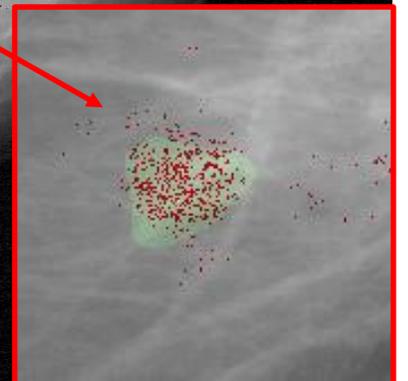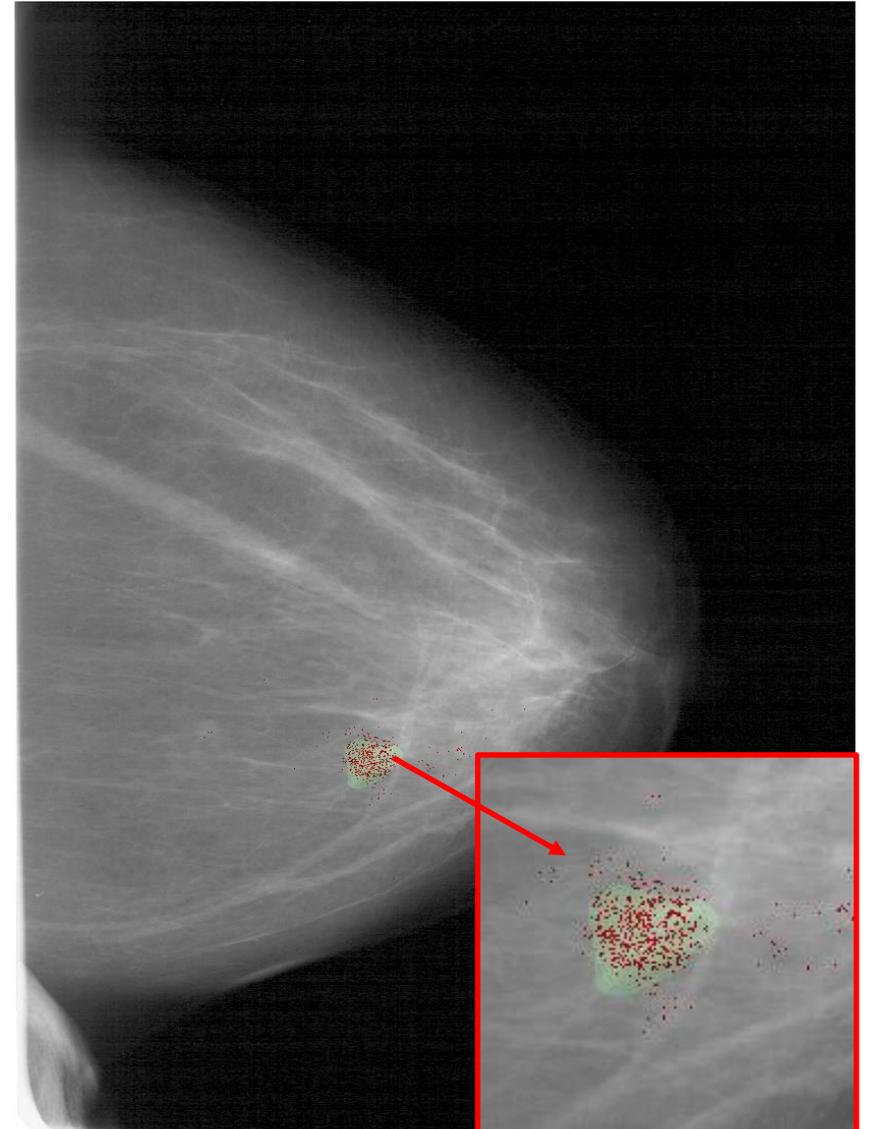
↓

$$\frac{\partial O_1}{\partial Im}$$
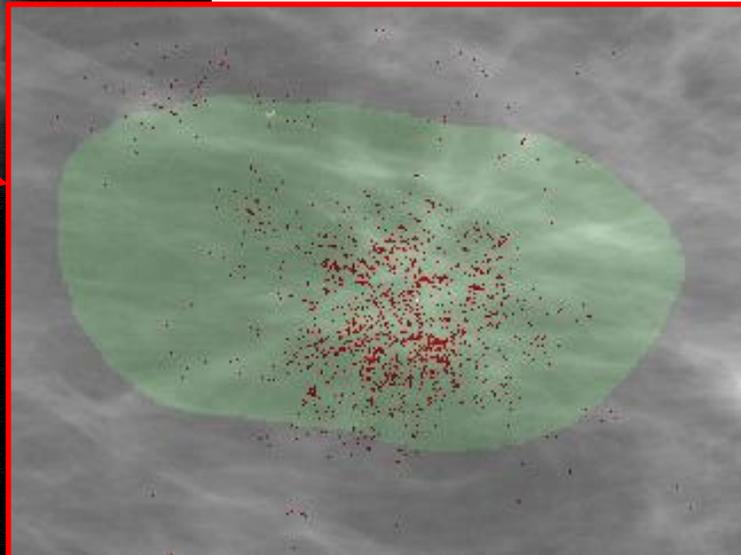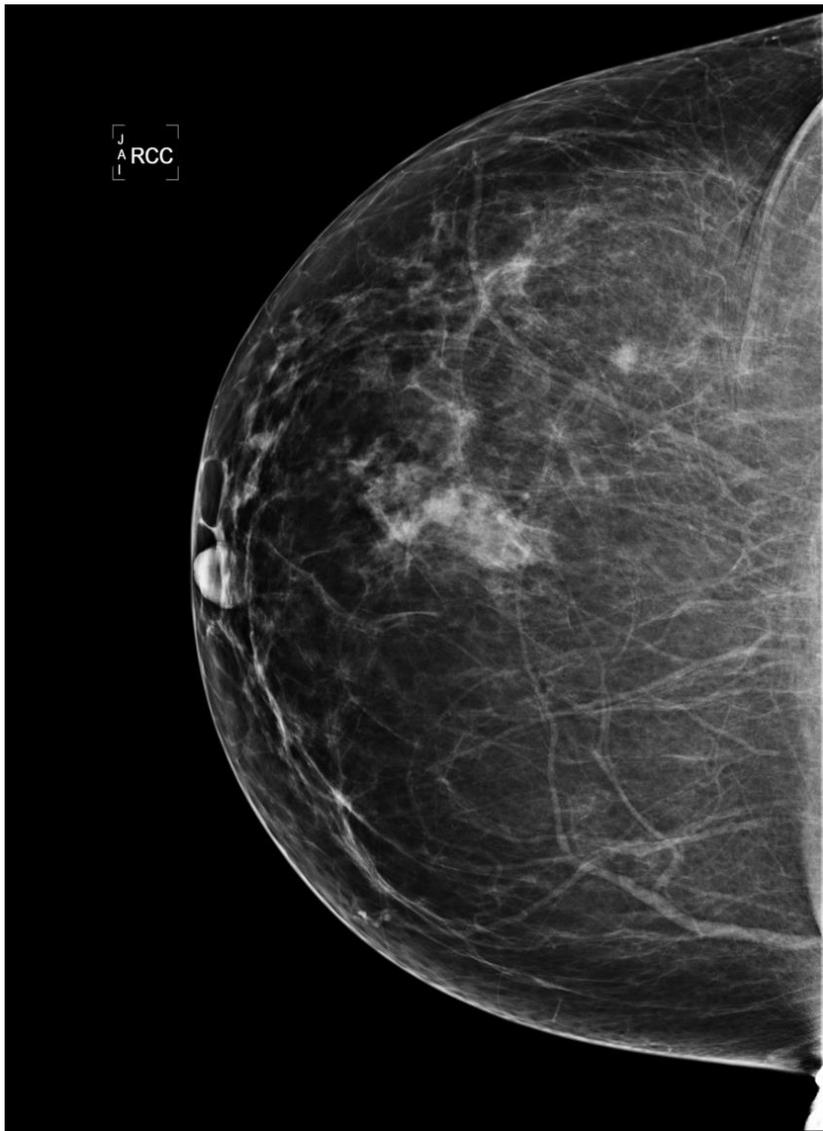
gradient of the output "1" with respect to the input

Idea credit: Simonyan, Vedaldi, Zisserman, Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps, 2013

> Red dots – slightly post-processed saliency maps

> Green area – mask suggested by radiologist
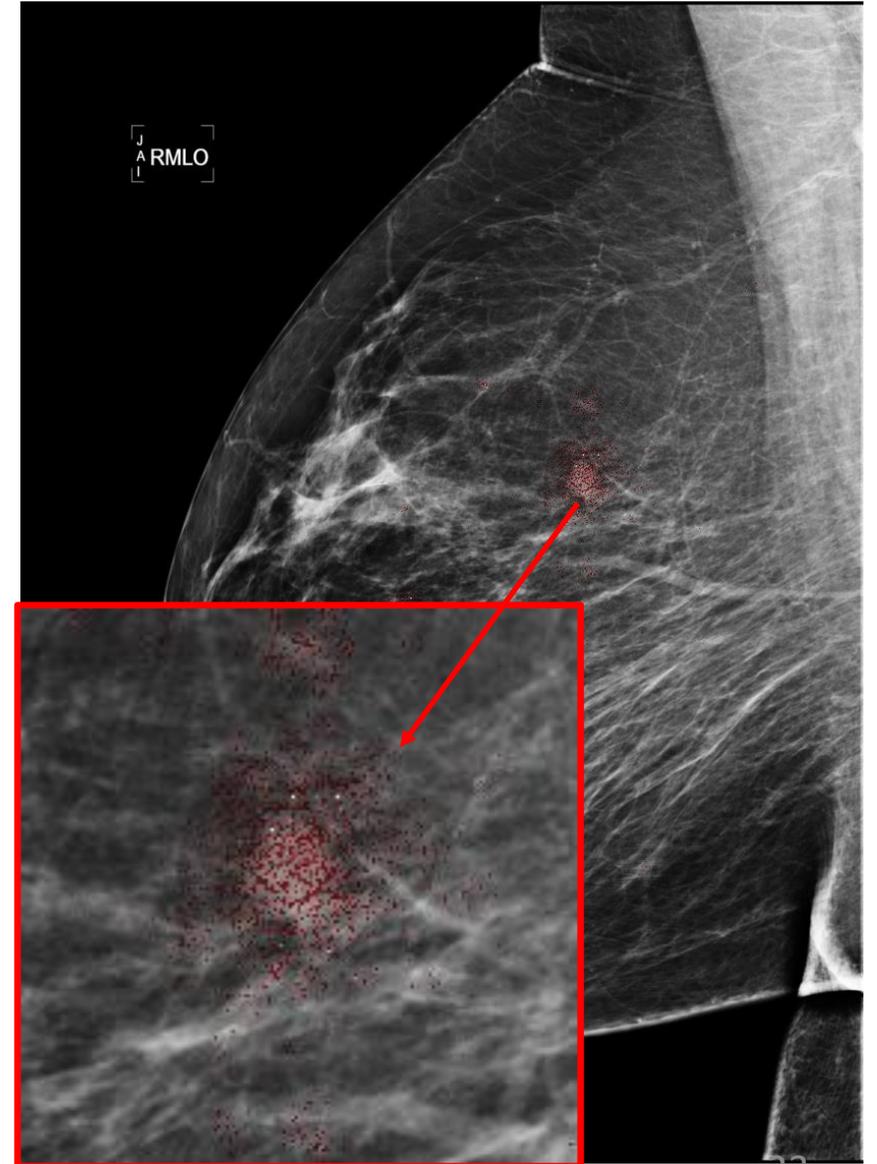
➢ DREAM data is much sharper

➢ Red dots – slightly post-processed saliency maps
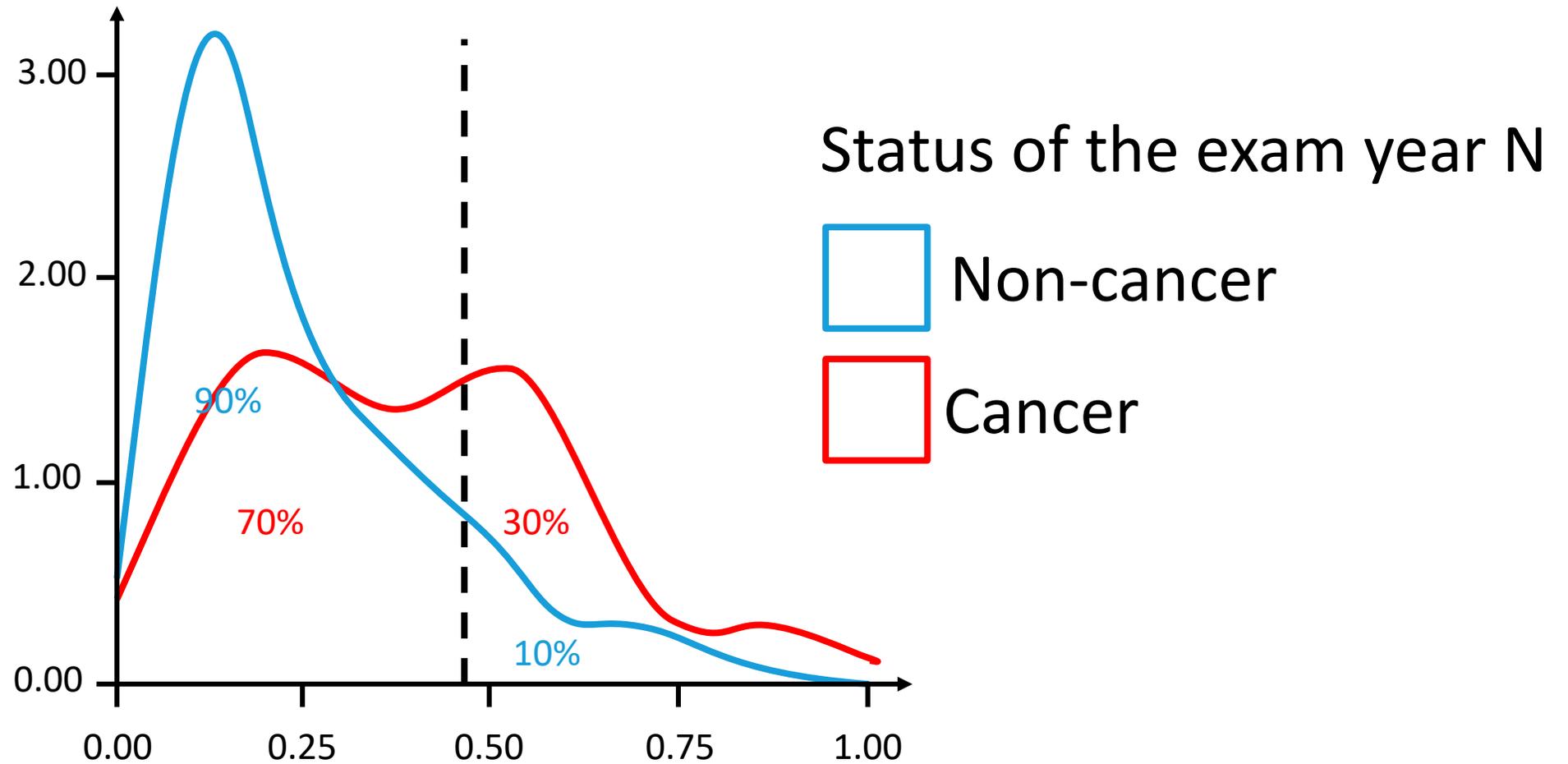
➢ The same lesion is highlighted on both views

➢ML projects need new paradigm

➢How we work at Therapixel

➢Specific advices

# Data Science 2019 = Software Engineering 1999

➢Visual Studio 1st release: 1997

➢Development process and paradigm evolving

➢Data becomes 2nd part of your code

➢Software 2.0 stack (©Andrej Karpathy)

➢IDEs for ML models are yet to come?

➢ Therapixel:

1. Development team
   - Cloud infrastructure
   - Integration with PACS in hospitals
   - Visualization & User Interface
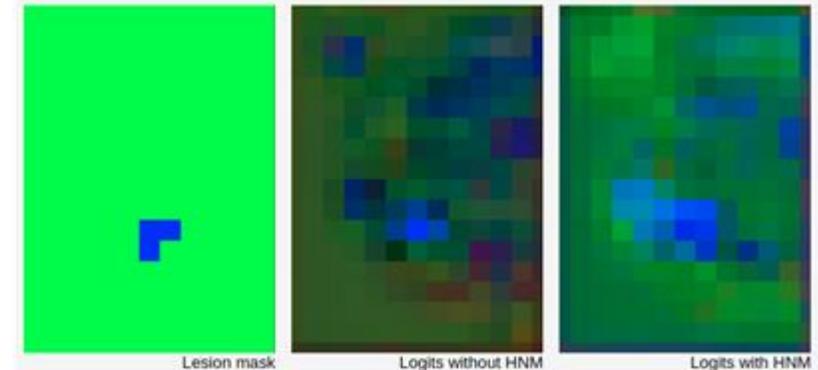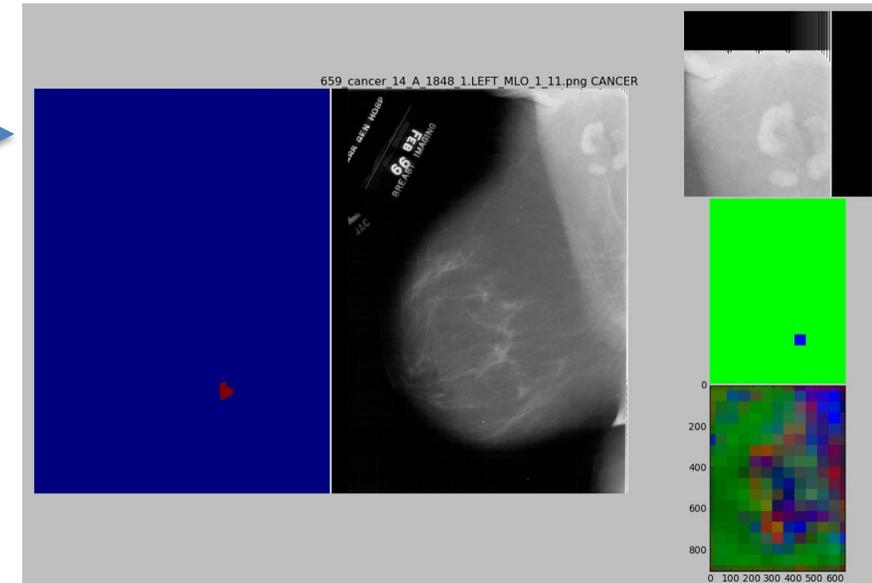
2. Data team
   - Partnerships with hospitals
   - Raw data extraction
   - Data clearing and structuring

3. Research team
   - Interfacing of structured data
   - Running experiments, reporting errors
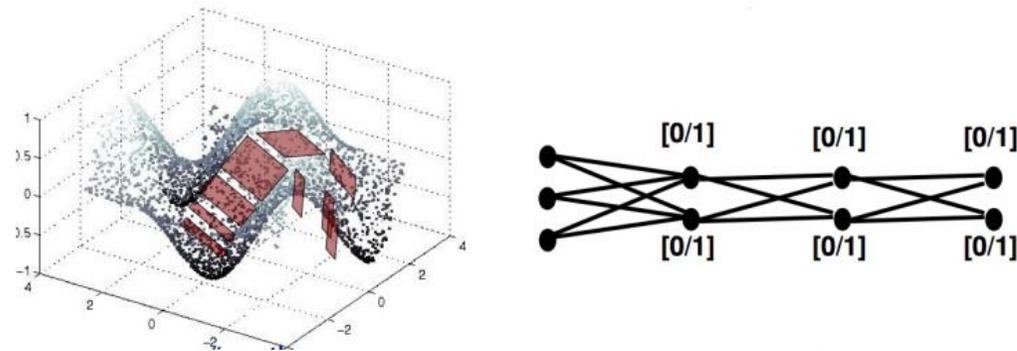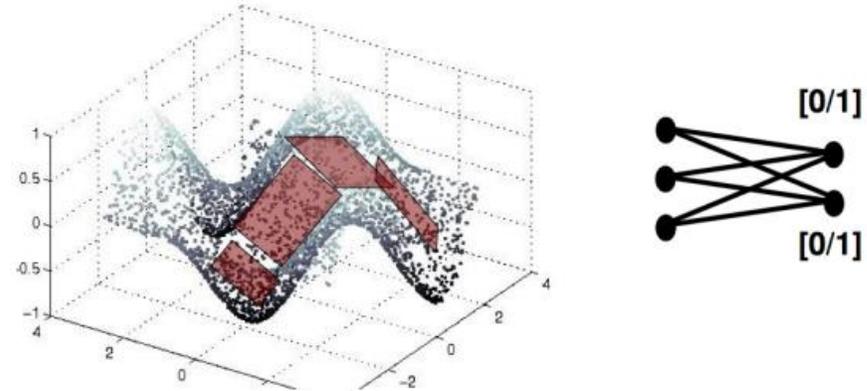   - Testing new ideas and extensions

# Some specific advices and practical moments

➤ Know your data
- If you don't understand your data – DL won't either
- Total nb of images, nb of images per class, typical resolution, RoI…
- Regularly examine worst offenders, manually guide your model
- Metadata is also under git (and dumped at each experiment)

➤ Enforce reproducibility
- No more binary reproducibility – GPUs
- Each experiment has an output folder
- For each experiment git hash and git diff are dumped
- Unit tests where applicable (example: complex stats calculations)

➤ Work in team
- Development cycles: 1-2 week
- Regular meetings with discussions
- Issue tracking tool
- Code review



659_cancer_14_A_1848_1.LEFT_MLO_1_11.png CANCER

Lesion mask          Logits without HNM          Logits with HNM

➢ Adapt model to your problem

➢ good data and gradient flow: "well-wired net"

➢ Adjust architecture !

➢ Deep = complex, but cheap

Slide credit: 1)G. Montúfar et al, On the Number of Linear Regions of Deep Neural Networks 2) Marc'Aurelio Ranzato slides 3) Introduction to Deep Learning by Iasonas Kokkinos

➢ Pretrain on balanced batches – make your network distinguish.

➢ Fight the overfitting: early stopping is simple but undesirable.
- Smaller model
- More data/data augmentation
- Regularization



Find the next number of the sequence

1, 3, 5, 7, ?

Correct solution
217341
because when

$$f(x) = \frac{18111}{2} x^4 - 90555 x^3 + \frac{633885}{2} x^2 - 452773 x + 217331$$

f(1)=1
f(2)=3          much solution
f(3)=5              very logic
        wow
f(4)=7
f(5)=217341
such function
many maths
        wow

# Thank you for your attention!

## Q&A session