

Breast Cancer Screening with Deep Neural Networks

Krzysztof J. Geras



NEW YORK UNIVERSITY

Joint work with Kyunghyun Cho, Linda Moy, Gene Kim, Stacey Wolfson and Artie Shen

NEURAL NETWORKS

They are amazing.

NEURAL NETWORKS

They are amazing. Can we save the world?

BREAST CANCER SCREENING

BREAST CANCER SCREENING

- ▶ Data are X-ray images. Basically, an image recognition problem.

BREAST CANCER SCREENING

- ▶ Data are X-ray images. Basically, an image recognition problem.
- ▶ Problem with a high impact on the society.

BREAST CANCER SCREENING

- ▶ Data are X-ray images. Basically, an image recognition problem.
- ▶ Problem with a high impact on the society.
 - In 2015, 232k women in the US were diagnosed with breast cancer, 40k died.

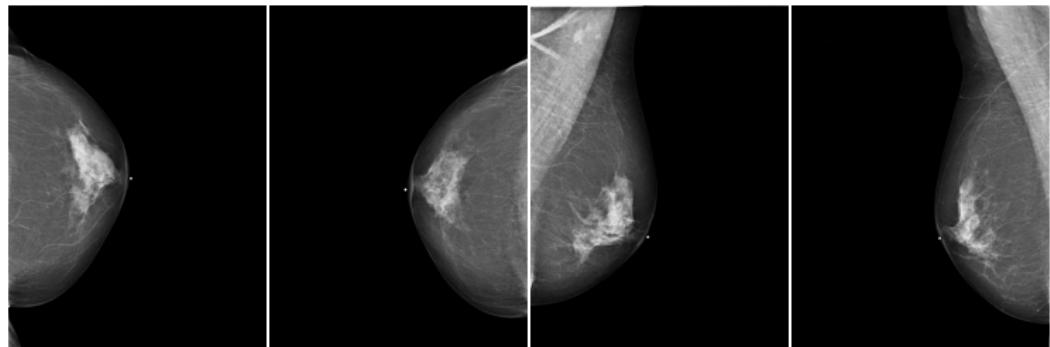
BREAST CANCER SCREENING

- ▶ Data are X-ray images. Basically, an image recognition problem.
- ▶ Problem with a high impact on the society.
 - In 2015, 232k women in the US were diagnosed with breast cancer, 40k died.
- ▶ Results generalizable to other types of diseases.

BREAST CANCER SCREENING

- ▶ Data are X-ray images. Basically, an image recognition problem.
- ▶ Problem with a high impact on the society.
 - In 2015, 232k women in the US were diagnosed with breast cancer, 40k died.
- ▶ Results generalizable to other types of diseases.
- ▶ Very little work done on end-to-end prediction.
Fundamental scientific problems to be solved.

BREAST CANCER SCREENING



L-CC

(left cranial caudal)

R-CC

(right cranial caudal)

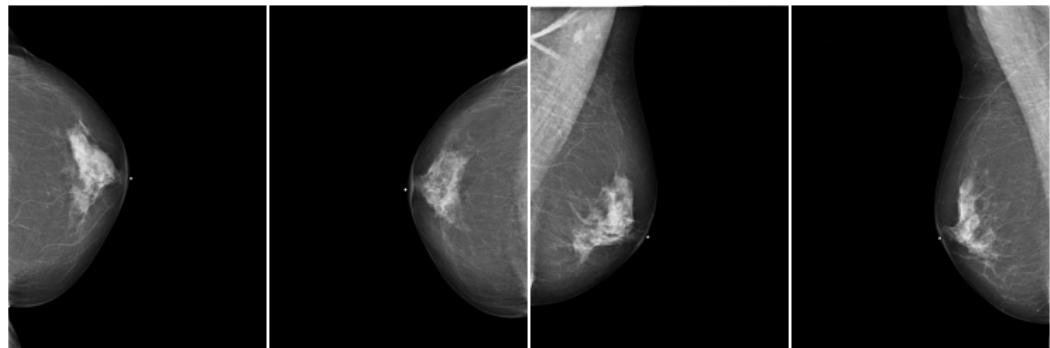
L-MLO

(left mediolateral oblique)

R-MLO

(right mediolateral oblique)

BREAST CANCER SCREENING



L-CC

(left cranial caudal)

R-CC

(right cranial caudal)

L-MLO

(left mediolateral oblique)

R-MLO

(right mediolateral oblique)

+ age, gender, cancer history, ethnicity, religion, ...

BREAST CANCER SCREENING

We try to mimic prediction of a radiologist.

- ▶ Class 0: incomplete ($\approx 15\%$).
- ▶ Class 1: negative ($\approx 50\%$).
- ▶ Class 2: benign findings ($\approx 35\%$).

BREAST CANCER SCREENING

We try to mimic prediction of a radiologist.

- ▶ Class 0: incomplete ($\approx 15\%$).
- ▶ Class 1: negative ($\approx 50\%$).
- ▶ Class 2: benign findings ($\approx 35\%$).

Radiologists call these classes BIRADS (short for Breast Imaging-Reporting and Data System).

CHALLENGES (1)

You need a lot of data to do deep learning.

CHALLENGES (1)

You need a lot of data to do deep learning.

Publicly available data sets contain about 1k images.

CHALLENGES (1)

You need a lot of data to do deep learning.

Publicly available data sets contain about 1k images.

We build our own data set:

- ▶ 23k exams,
- ▶ 103k images.

CHALLENGES (1)

You need a lot of data to do deep learning.

Publicly available data sets contain about 1k images.

We build our own data set:

- ▶ 23k exams,
- ▶ 103k images.

(And we are soon going to have $\times 10$ more).

CHALLENGES (2)

CHALLENGES (2)

CHALLENGES (2)



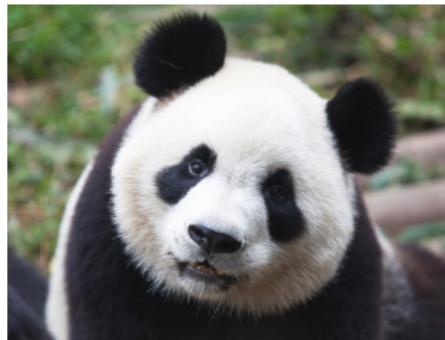
CHALLENGES (2)



CHALLENGES (2)

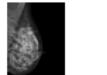
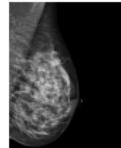
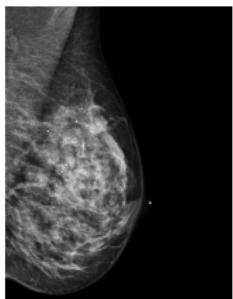
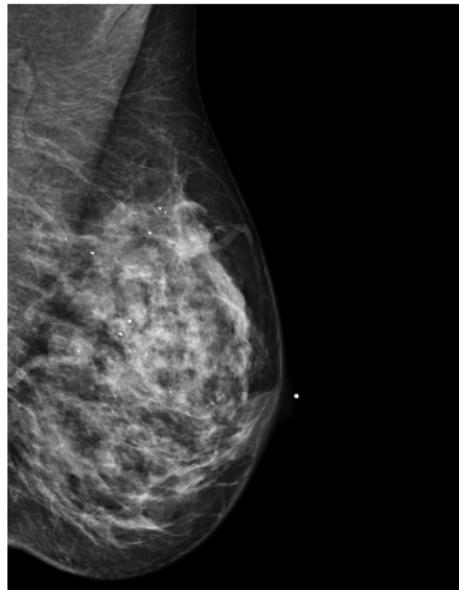


CHALLENGES (2)



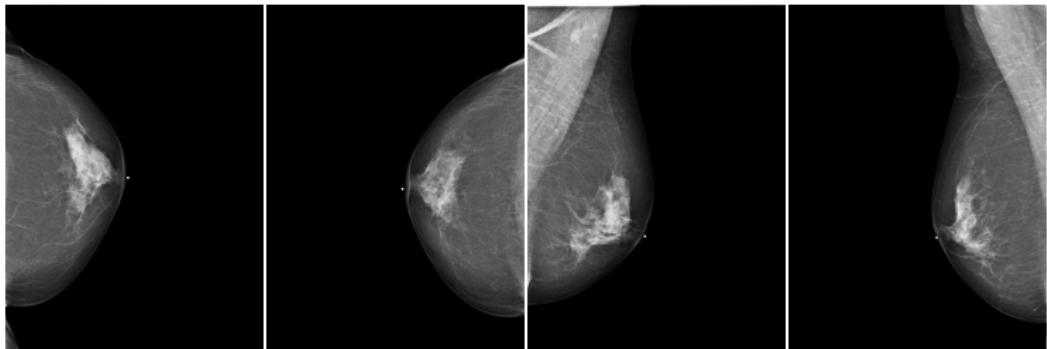
CHALLENGES (2)

High resolution necessary - computational and engineering challenge.



CHALLENGES (3)

Multi-view data. How to integrate information?



OUR MODEL

Classifier $p(y x)$			
Concatenation (256×4 dim)			
DCN	DCN	DCN	DCN
L-CC	R-CC	L-MLO	R-MLO

layer	kernel size	stride	#maps	repetition
global average pooling			256	
convolution	3×3	1×1	256	×3
max pooling	2×2	2×2	128	
convolution	3×3	1×1	128	× 3
max pooling	2×2	2×2	128	
convolution	3×3	1×1	128	× 3
max pooling	2×2	2×2	64	
convolution	3×3	1×1	64	× 2
convolution	3×3	2×2	64	
max pooling	3×3	3×3	32	
convolution	3×3	2×2	32	
input			1	

OUR MODEL

Classifier $p(y x)$			
Concatenation (256×4 dim)			
DCN	DCN	DCN	DCN
L-CC	R-CC	L-MLO	R-MLO

layer	kernel size	stride	#maps	repetition
global average pooling			256	
convolution	3×3	1×1	256	×3
max pooling	2×2	2×2	128	
convolution	3×3	1×1	128	× 3
max pooling	2×2	2×2	128	
convolution	3×3	1×1	128	× 3
max pooling	2×2	2×2	64	
convolution	3×3	1×1	64	× 2
convolution	3×3	2×2	64	
max pooling	3×3	3×3	32	
convolution	3×3	2×2	32	
input			1	

We use dropout, early stopping and data augmentation.

OUR MODEL

Classifier $p(y x)$			
Concatenation (256×4 dim)			
DCN	DCN	DCN	DCN
L-CC	R-CC	L-MLO	R-MLO

layer	kernel size	stride	#maps	repetition
global average pooling			256	
convolution	3×3	1×1	256	×3
max pooling	2×2	2×2	128	
convolution	3×3	1×1	128	× 3
max pooling	2×2	2×2	128	
convolution	3×3	1×1	128	× 3
max pooling	2×2	2×2	64	
convolution	3×3	1×1	64	× 2
convolution	3×3	2×2	64	
max pooling	3×3	3×3	32	
convolution	3×3	2×2	32	
input			1	

We use dropout, early stopping and data augmentation.
Columns share parameters.

OUR MODEL

Classifier $p(y x)$			
Concatenation (256×4 dim)			
DCN	DCN	DCN	DCN
L-CC	R-CC	L-MLO	R-MLO

layer	kernel size	stride	#maps	repetition
global average pooling			256	
convolution	3×3	1×1	256	×3
max pooling	2×2	2×2	128	
convolution	3×3	1×1	128	× 3
max pooling	2×2	2×2	128	
convolution	3×3	1×1	128	× 3
max pooling	2×2	2×2	64	
convolution	3×3	1×1	64	× 2
convolution	3×3	2×2	64	
max pooling	3×3	3×3	32	
convolution	3×3	2×2	32	
input			1	

We use dropout, early stopping and data augmentation.
Columns share parameters. We use Adam for learning.

OUR MODEL

Classifier $p(y x)$			
Concatenation (256×4 dim)			
DCN	DCN	DCN	DCN
L-CC	R-CC	L-MLO	R-MLO

layer	kernel size	stride	#maps	repetition
global average pooling			256	
convolution	3×3	1×1	256	×3
max pooling	2×2	2×2	128	
convolution	3×3	1×1	128	× 3
max pooling	2×2	2×2	128	
convolution	3×3	1×1	128	× 3
max pooling	2×2	2×2	64	
convolution	3×3	1×1	64	× 2
convolution	3×3	2×2	64	
max pooling	3×3	3×3	32	
convolution	3×3	2×2	32	
input			1	

We use dropout, early stopping and data augmentation.
Columns share parameters. We use Adam for learning.
A highly optimised implementation needs a week to train.

EVALUATION METRICS

- ▶ Per class AUC.

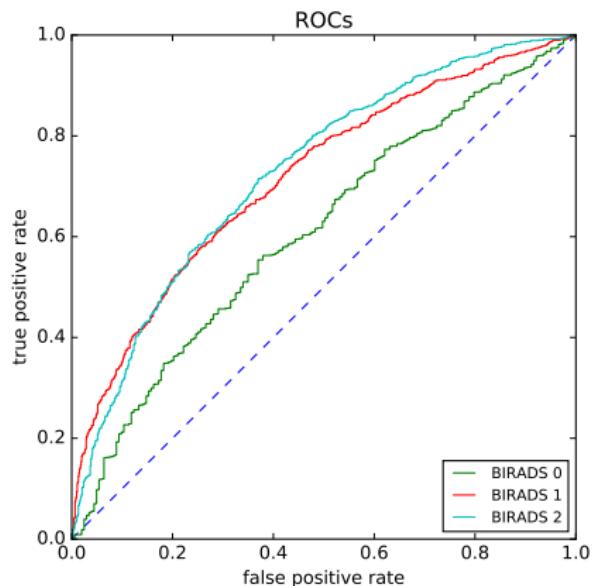
EVALUATION METRICS

- ▶ Per class AUC.
- ▶ Average AUC.

EVALUATION METRICS

- ▶ Per class AUC.
- ▶ Average AUC.
- ▶ Average AUC for 30% most confident examples.

RESULTS



0 vs. others: 0.609

1 vs. others: 0.717

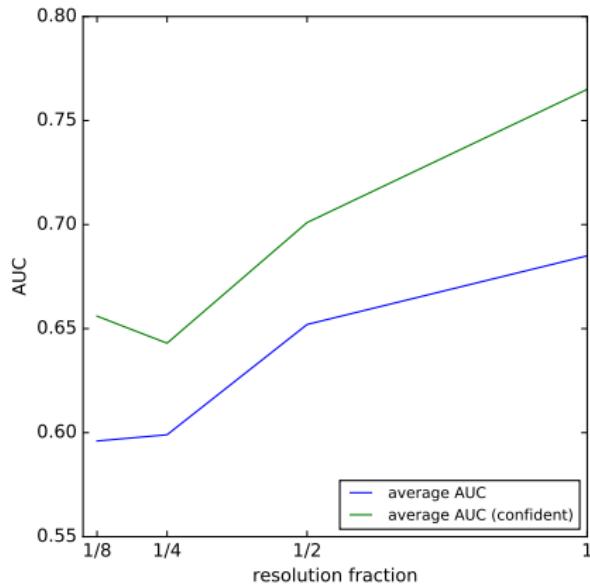
2 vs. others: 0.728

average: **0.685**

average (confident): **0.765**

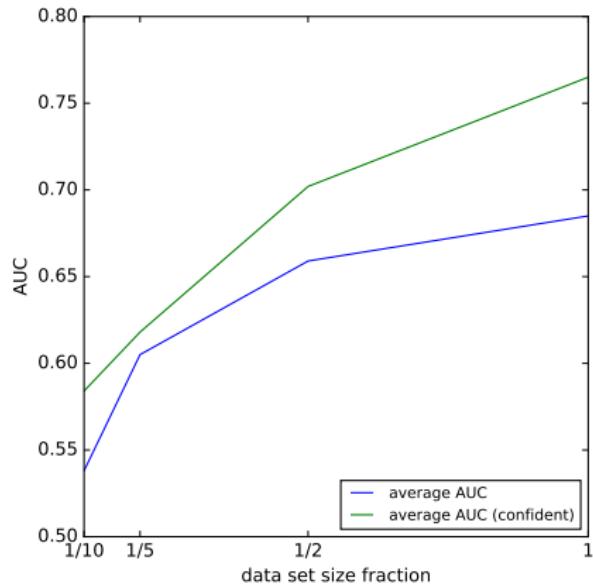
IMPACT OF DOWNSCALING

Original image size - 2600×2000 .



IMPACT OF THE DATA SET SIZE

Original data set size - 23k exams.

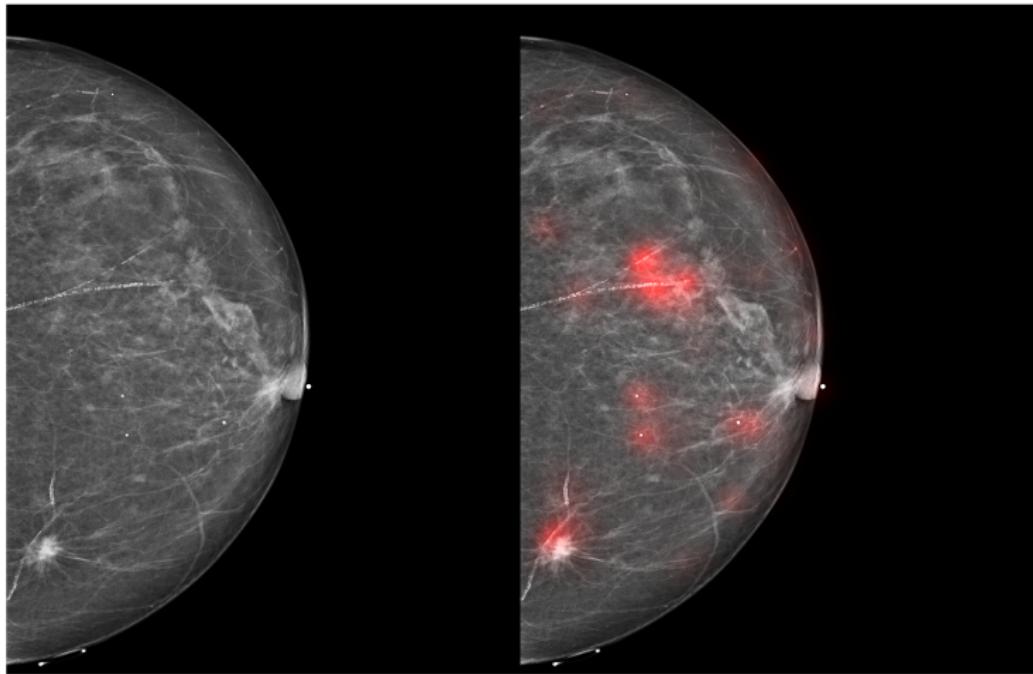


VISUALISATION

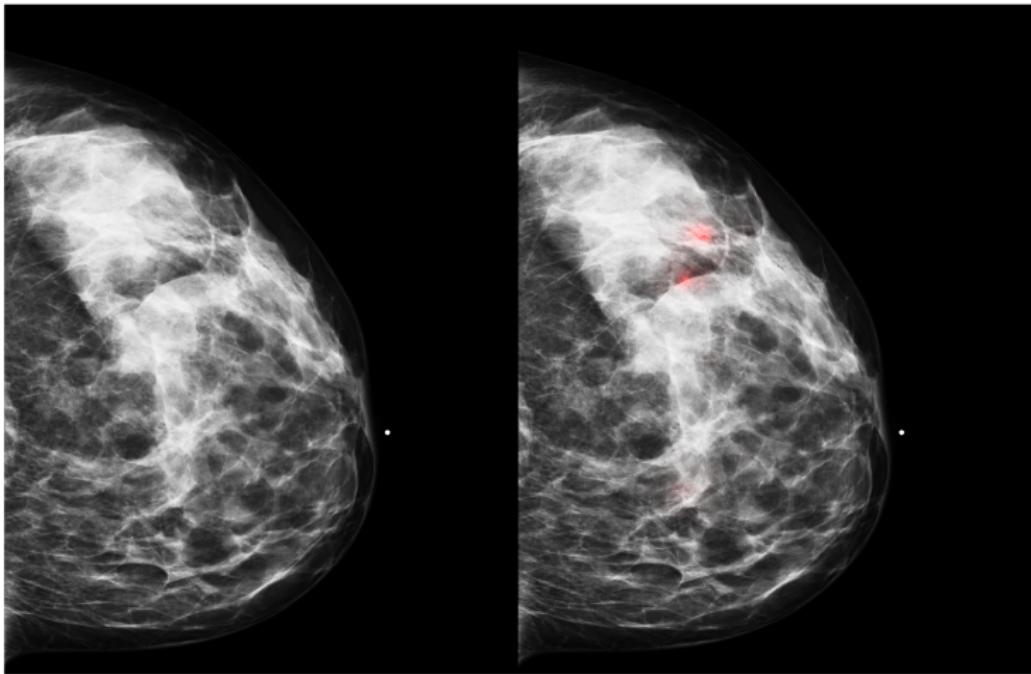
We visualise $\left| \frac{\partial H(y|\mathbf{x})}{\partial \mathbf{x}_{(i,j)}^v} \right|$,

where $H(y|\mathbf{x}) = - \sum_{y' \in \mathcal{C}} p(y'|\mathbf{x}) \log p(y'|\mathbf{x})$.

VISUALISATION



VISUALISATION



VISUALISATION

- ▶ Activations of the penultimate layer.

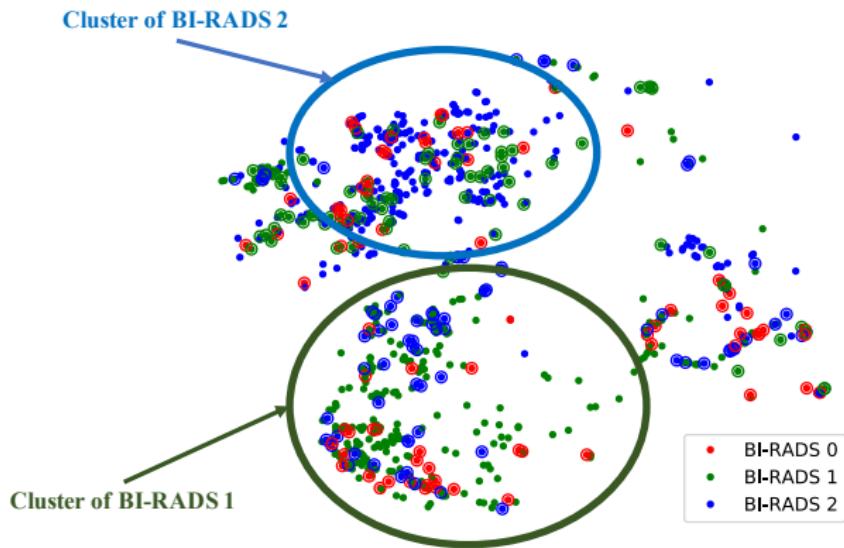
VISUALISATION

- ▶ Activations of the penultimate layer.
- ▶ t-SNE.

VISUALISATION

- ▶ Activations of the penultimate layer.
- ▶ t-SNE.
- ▶ Only confident examples shown.

VISUALISATION



CONCLUSIONS

CONCLUSIONS

- ▶ It is much harder to learn the “incomplete” (0) class than other classes.

CONCLUSIONS

- ▶ It is much harder to learn the “incomplete” (0) class than other classes.
- ▶ We need to use the full resolution.

CONCLUSIONS

- ▶ It is much harder to learn the “incomplete” (0) class than other classes.
- ▶ We need to use the full resolution.
- ▶ We need more data.

WAYS TO GO FORWARD

WAYS TO GO FORWARD

- ▶ Incorporating side information.

WAYS TO GO FORWARD

- ▶ Incorporating side information.
- ▶ Taking advantage of previous exams.

WAYS TO GO FORWARD

- ▶ Incorporating side information.
- ▶ Taking advantage of previous exams.
- ▶ Learning where to look.