# MURA Deep Learning Competition

#### A Case Study



A Platform

For Senior IT Professionals

# Why It Matters

- Is DL competition achievable by us, IT professionals?
- If so, how?

资深人协会

- DL Knowledge
- Data Science Methodologies
- Software: Programing Language / DL Platform / Libraries / Tools
- Hardware: GPU
- Time commitment
- How to go beyond?



MURA: Large Dataset for Abnormality Detection in Musculoskeletal Radiographs

Pranav Rajpurkar<sup>1,\*</sup>, Jeremy Irvin<sup>1,\*</sup>, Aarti Bagul<sup>1</sup>, Daisy Ding<sup>1</sup>, Tony Duan<sup>1</sup>, Hershel Mehta<sup>1</sup>, Brandon Yang<sup>1</sup>, Kaylie Zhu<sup>1</sup>, Dillon Laird<sup>1</sup>, Robyn L. Ball<sup>2</sup>, Curtis Langlotz<sup>3</sup>, Katie Shpanskaya<sup>3</sup>, Matthew P. Lungren<sup>3,†</sup>, Andrew Y. Ng<sup>1,†</sup>

\*, †Equal Contribution

<sup>1</sup> Department of Computer Science Stanford University {pranavsr, jirvin16}@cs.stanford.edu

> <sup>2</sup>Department of Medicine Stanford University rball@stanford.edu

<sup>3</sup>Department of Radiology Stanford University mlungren@stanford.edu

#### Abstract

We introduce MURA, a large dataset of musculoskeletal radiographs containing 40.561 images from 14.863 studies, where each study is manually labeled by radiologists as either normal or abnormal. To evaluate models robustly and to get an estimate of radiologist performance, we collect additional labels from six boardcertified Stanford radiologists on the test set, consisting of 207 musculoskeletal studies. On this test set, the majority vote of a group of three radiologists serves as gold standard. We train a 169-layer DenseNet baseline model to detect and localize abnormalities. Our model achieves an AUROC of 0.929, with an operating point of 0.815 sensitivity and 0.887 specificity. We compare our model and radiologists on the Cohen's kappa statistic, which expresses the agreement of our model and of each radiologist with the gold standard. Model performance is comparable to the best radiologist performance in detecting abnormalities on finger and wrist studies. However, model performance is lower than best radiologist performance in detecting abnormalities on elbow, forearm, hand, humerus, and shoulder studies. We believe that the task is a good challenge for future research. To encourage advances, we have made our dataset freely available at http://stanfordmlgroup.github.io/competitions/mura.







#### The Baseline

ACSIP 近 资深人协会

	Radiologist 1	Radiologist 2	Radiologist 3	Model
Elbow	0.850 (0.830, 0.871)	0.710 (0.674, 0.745)	0.719 (0.685, 0.752)	0.710 (0.674, 0.745)
Finger	0.304 (0.249, 0.358)	0.403 (0.339, 0.467)	0.410 (0.358, 0.463)	0.389 (0.332, 0.446)
Forearm	0.796 (0.772, 0.821)	0.802 (0.779, 0.825)	0.798 (0.774, 0.822)	0.737 (0.707, 0.766)
Hand	0.661 (0.623, 0.698)	0.927 (0.917, 0.937)	0.789 (0.762, 0.815)	0.851 (0.830, 0.871)
Humerus	0.867 (0.850, 0.883)	0.733 (0.703, 0.764)	0.933 (0.925, 0.942)	0.600 (0.558, 0.642)
Shoulder	0.864 (0.847, 0.881)	0.791 (0.765, 0.816)	0.864 (0.847, 0.881)	0.729 (0.697, 0.760)
Wrist	0.791 (0.766, 0.817)	0.931 (0.922, 0.940)	0.931 (0.922, 0.940)	0.931 (0.922, 0.940)
Overall	0.731 (0.726, 0.735)	0.763 (0.759, 0.767)	0.778 (0.774, 0.782)	0.705 (0.700, 0.710)



# **MURA** Competition

- A Good Way To Study
  - Simple Binary Classifications

#### HOW DO WE START FROM HERE?









### MNIST

0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	1	/	l	1	1	1	1	I	- 1	1	1	1	1	l	l	l
1	1	1	7	1	7	7	1	7	7	1	7	1	1	1	2	1
3	3	3	Z	3	3	3	3	3	3	3	3	3	3	3	3	3
R	3	a	2	3	r	Ŷ	2	Ŷ	Ŷ	ત	a	R	a	R	ส	Y
9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9
6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6
8	8	8	8	8	8	8	8	ч	8	8	8	8	8	8	8	8
4	4	4	4	4	ч	4	4	4	4	Ц	4	4	4	4	Ц	Ч
5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5

Í	A( 丁 资	CSIP 深人协会	MANUST Load Data	
			IVITAIST - LOUG DUTG	
		Load MNIS	ST data	Python
		<pre># Load pr (X_train,</pre>	re-shuffled MNIST data into train and test sets , y_train), (X_test, y_test) = mnist.load_data()	
		print X_1 # (60000	train.shape , 28, 28)	Python



#### MNIST – Data Pre-Processing

Reshape input data

X\_train = X\_train.reshape(X\_train.shape[0], 1, 28, 28)

 $X_{test} = X_{test.reshape}(X_{test.shape}[0], 1, 28, 28)$ 

Convert data type and normalize values
X\_train = X\_train.astype('float32')
X\_test = X\_test.astype('float32')
X\_train /= 255
X\_test /= 255

```
Preprocess class labels

# Convert 1-dimensional class arrays to 10-dimensional class matrices

Y_train = np_utils.to_categorical(y_train, 10)

Y_test = np_utils.to_categorical(y_test, 10)
```

Python

Python







```
CSIP
17 资深人协会
                    MNIST - Model
       Model architecture
                                                                                           Python
       model = Sequential()
       model.add(Convolution2D(32, 3, 3, activation='relu', input_shape=(1,28,28)))
       model.add(Convolution2D(32, 3, 3, activation='relu'))
       model.add(MaxPooling2D(pool_size=(2,2)))
       model.add(Dropout(0.25))
       model.add(Flatten())
       model.add(Dense(128, activation='relu'))
       model.add(Dropout(0.5))
       model.add(Dense(10, activation='softmax'))
```







# **GPU** Battle

- Build Your Own vs Cloud

Build Your Own:

- Full control
- Fixed cost (as opposed to Cloud)
- OS set up is preserved over time
- Cheaper than cloud over the long term
- Looks good

#### Cloud:

- Much more scalable
- Don't need to pay for electricity consumed by GPUs
- Don't need to worry about hardware issues
- Easy and secure remote access

#### Standard practice:

Set up model training pipeline and do experiments on local GPU, then spin up hundreds of cloud instances to perform hyperparameter tuning in parallel



## GPU - Cloud

		AWS	Paperspace	FloydHub	Suggested platform
	Ease of setup	Complex. Manual setup required	Pre-installed libraries and frameworks	Provides almost all major DL environments	FloydHub/Paperspace
	User Experience	Smooth. Data upload/download tricky	Quite okay apart from the server lag	Not very intuitive. Takes time getting used to	AWS
	Hardware/Software	Tesla K80s	Latest Pascal and Volta series GPUs	Tesla K80s	Paperspace
	Performance	Decent	Maxwell series – almost like AWS; Pascal series – 3x of AWS	0.75x of AWS	AWS/Paperspace
$\mathbb{N}/$	Additional features	Offers multi-GPU systems	Desktop environment; Shareable drives	Public datasets; concurrent job runs; support for MOOCs	FloydHub
	Pricing	Starting at \$0.9/hr + \$13/mo for storage + IP	Starting at \$0.4/hr + \$7/mo for storage + IP	\$14/mo + Powerups on base plan	Paperspace

# 资深人协会 The Winner – Google Cloud US\$300 Credit = 500 + hours Free GPU time if pick Tesla K80 (0.45 per hour) Excellent customer service Easy to setup and scalable Simple Pricing Model Using it 25% of a month, standard discount kicks in If running for full month, 30% discount



# Software Installation In the Cloud

- Tensorflow 1.8.0 or newer
- Numpy, Scipy, Scikit Learn, Scikit Image, Pandas, Matplotlib, Pillow
- If you prefer to use docker, there are plenty of docker image files to choose
- Keras docker https://github.com/keras-team/keras/tree/master/docker
- All in one DL docker https://github.com/floydhub/dl-docker



# Things to consider if you want to build your own workstation

Nvidia vs AMD

NVIDIA:

- Much larger DL community
- Easy to find resources for
- Supports DL with 9 series + cards

AMD:

- More cost efficient
- Started DL cards rather recently so smaller community and less support
- Supports DL with Vega and later cards

Why not Intel?

Intel compute cards (Xeon Phi) are just way too expensive...





#### An example of an overkill build (And this is why you should define your budget early)



- Motherboard:
  - ASUS ROG Rampage V Edition 10
- CPU:
  - Intel i7 6850K
  - RAM:
    - Corsair Dominator Platinum 4 x 16GB 2400GHz
- GPU:
  - 2 x EVGA Titan X Hybrid (Maxwell)
  - 2 x EVGA 1080 Ti FTW3 Hybrid
- Hard Drive:
  - Samsung 850 Evo M.2. SSD
  - 3 x Western Digital 4TB hard drive (RAID 5)



#### **GPU - BRANDS**

- Brand
  - You may see these GPUs having a lot of different brands, such as
    - EVGA
    - ASUS
    - MSI
    - GIGABYTE
    - Etc.
  - Essentially NVIDIA/AMD only provides the chip, and each brand deals with
    - Overclocking
    - Cooling
    - Card Structure
  - Hence, check the price and review across different brands to decide which one to buy



However, if you would like to explore speeding up your training by using half-precision, check out the GPUs with TPUs (Tensor Processing Units)





- The Lowest GPU is your bottleneck
- Cooling System is Important!
  - 4 x GTX 1080 Ti = 1120W, a heater 24/7
- Motherboard is important
- CPU (40 vs 28 PCIe Lanes)
- PSU (1300W for 4 GPUs)



### **GPU** Power Consumptions

ri Ju	1.0s: nvidia	-smi 9 2018			Fri Jul
NVID	IA-SMI 390.6	7	Driver Version: 390	.67	
GPU	Name	Persistence-M	Bus-Id Disp.A	Volatile	Uncorr. ECC
Fan	Temp Perf	Pwr:Usage/Cap	Memory-Usage	GPU-Util	Compute M.
0	GeForce GTX	TIT On	00000000:02:00.0 On		N/A
29%	69C P2	191W / 250W	11734MiB / 12211MiB	59%	Default
1	GeForce GTX	108 On	00000000:04:00.0 Off	50%	N/A
9%	55C P2	208W / 280W	10795MiB / 11178MiB		Default
2	GeForce GTX	108 On	00000000:05:00.0 Off	45%	N/A
0%	49C P2	209W / 280W	10795MiB / 11178MiB		Default
3	GeForce GTX 56C P2	TIT On   186W / 250W	00000000:06:00.0 Off   11730MiB / 12212MiB	72%	N/A Default

- My electricity bill did go up by ~\$100 for the month that I trained models non-stop...
- And my room temperature went up by 2 degrees
  - The corner where the computer sits was always warm



#### Nvidia GTX 1080 Ti vs AWS Tesla K80: ~4X

To gain performance of one GTX 1080TI equivalent in the cloud running 24x7 = \$2500/ mth



## Is accuracy a good measurement

- A Disease: 1 out of 1 million probability
- A test: 99.99% accuracy

资深人协会

- If you are tested positive, should you be worried?
- Out of 1 million prediction, 100 are predicted wrong
- Out of 1 million population, there is only 1 person who has it







# Replacing the TOP

资深人协会

- vgg\_conv = VGG16(weights='imagenet', include\_top=False)
- model = models.Sequential().add(vgg\_conv)
- model.add(layers.Flatten())
- model.add(layers.Dense(2, activation='softmax'))



# Data Augmentation



## **Tuning Hyper Parameters**

MobileNet, VGG, Densenet, Resnet ......

- Learning Rate Decay, Cos Annealing, and Early Stopping
- Backprop Algorithm (SGD, Adam, Adagrad ... )



## Visualization – Correct Reasoning

CSIP 资深人协会

Prediction: [[12.194778]], Label: 1



# CSIP <sub>资深人协会</sub> Visualization – Wrong Reasoning



Image: D:\Code\mura\dataset\MURA-v1.1/valid/XR\_ELBOW/patient11342/study1\_positive/image1.png





- Submit your code
- Run through some sample data
- Tag it

CSIP

资深人协会

They will run your model with real validation data





## Next Steps

- Taking Online Courses
  - Andrew Ng's Deep Learning specialization on Coursera
  - CS229, CS231n and CS224d from Stanford
  - Deep Reinforcement Learning from UC Berkeley
  - Fast.AI Deep Learning Course
- Attending Conferences
  - NIPS, ICML, ICLR, ACL, ReWork Deep Learning Submit
- Reading Papers
- Working On Projects (Competitions, Side Projects etc)

#### MOST IMPORTANT: TIME COMMITMENT



### **Future References**

# http://forum.aisquaredforum.ca



- https://github.com/DeepMachineLearning/mura-team1
- <u>https://github.com/DeepMachineLearning/mura-team2</u>
- <u>https://github.com/DeepMachineLearning/mura-team3</u>