

This is your **last** free member-only story this month.
[Sign up for Medium and get an extra one](#)

Prakhar Ganesh · Follow

Jun 24, 2019 · 5 min read ★

...

Deep Learning — Model Optimization and Compression: Simplified

Take a peek into the domain of compression, pruning and quantization of state-of-the-art Machine Learning models

What's this?

The world around us is filled with Neural Networks and Deep Learning models doing *wonders*!! But these models are both computationally expensive and energy intensive. So expensive that people have actually started holding AI/ML accountable for their carbon emission and the numbers are *not pretty*!!

Training a single AI model can emit as much carbon as five cars in their lifetimes

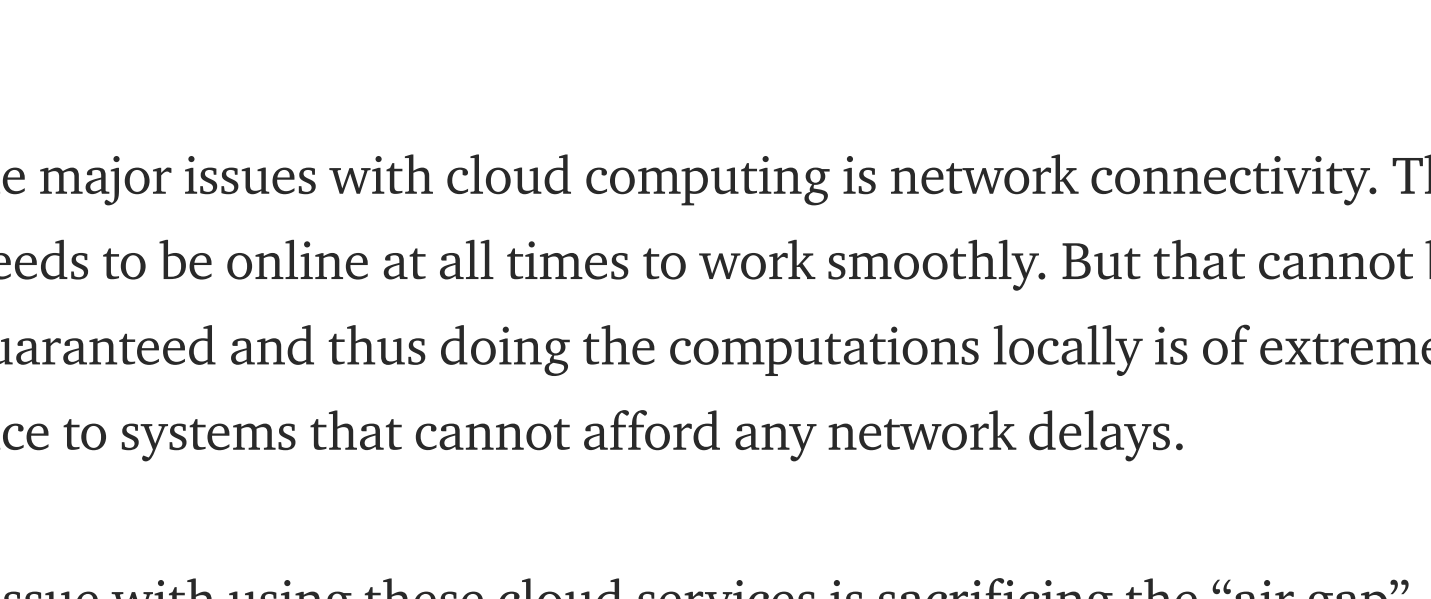
The artificial-intelligence industry is often compared to the oil industry: once mined and refined, data, like oil, can...

www.technologyreview.com

Another major reason why more researchers have turned towards Model compression is the difficulty in deploying these models on systems with limited hardware resources. While these models have been successful in making headlines and achieving extraordinary performances, they require the support of expensive, high speed GPUs to get them working, which limits their applications.

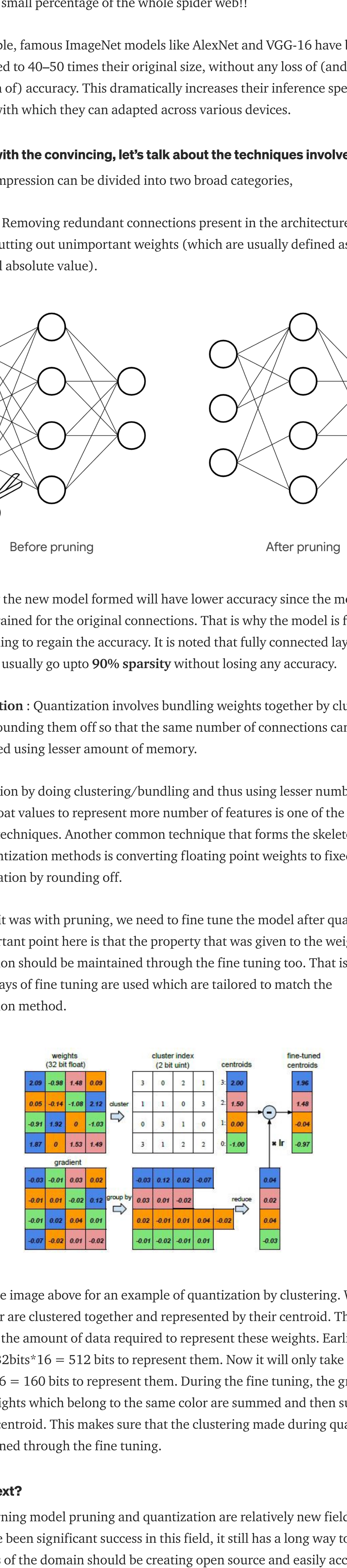
Why not just use GPU servers?

Yes *offcourse!!* With the internet giants like Google and Amazon offering computational services online, one does wonder if doing remote computations is the way to go. While people have started using these cloud services like a crutch for heavy computations, it does come with it's own set of problems.



One of the major issues with cloud computing is network connectivity. The system needs to be online at all times to work smoothly. But that cannot be always guaranteed and thus doing the computations locally is of extreme importance to systems that cannot afford any network delays.

Another issue with using these cloud services is sacrificing the “air gap”. The air gap is a technical term used to represent systems that are not connected to the internet and thus cannot be breached remotely. Getting access to the data present in these configurations needs to be done physically, *Mission Impossible style*!! :P



For systems which are extremely protective regarding their privacy and security, giving up this “air gap” is not ideal and so they prefer local computations over cloud services.

But none of that actually affects me!!

That's what the majority of ML community believes but is not true!! If you are a beginner in ML looking to develop state-of-the-art models and are not bound by processing capacities, you might think that highly complex and Deep models are always the way to go.

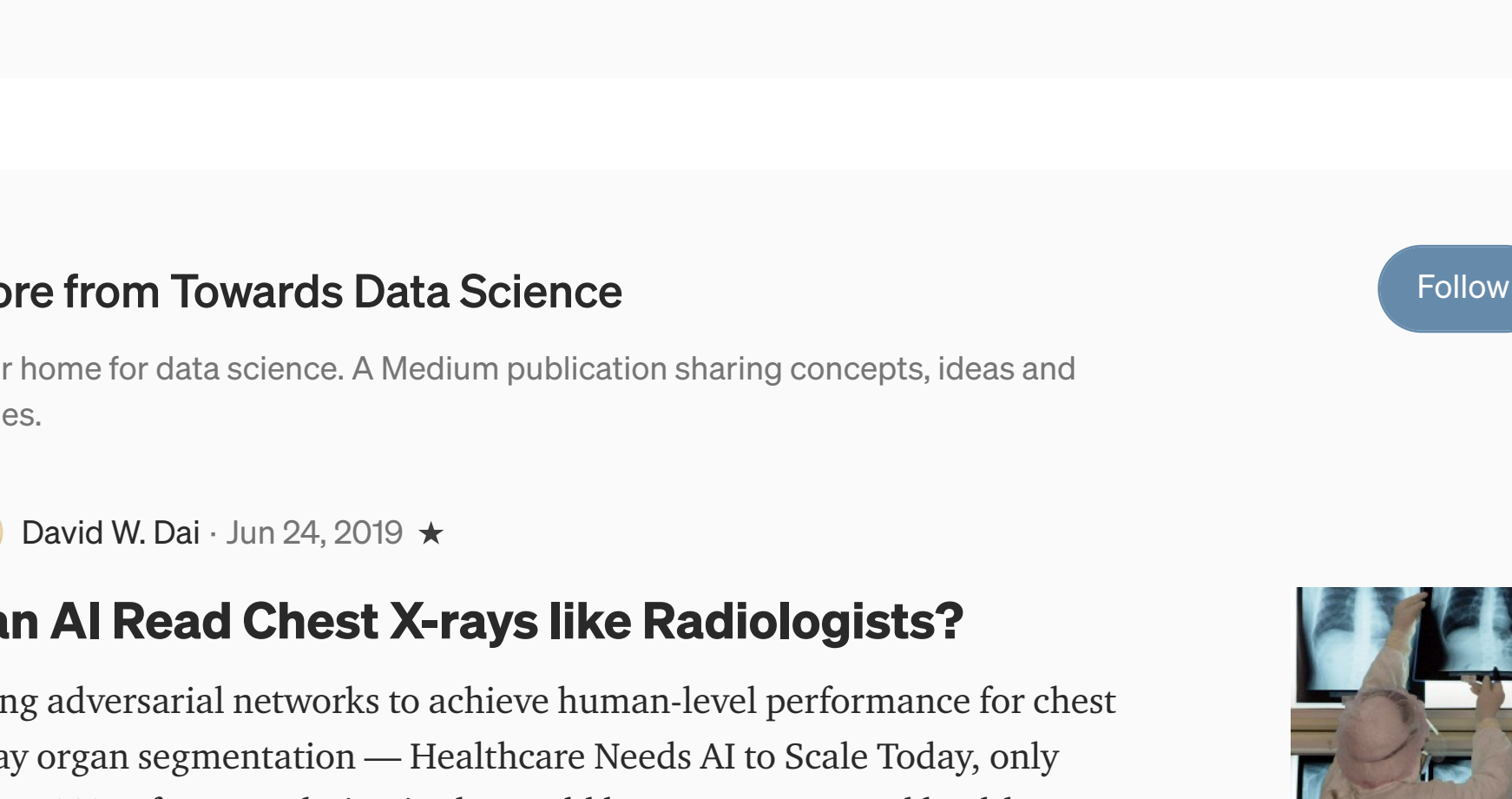
But that's a huge misconception. Highly complex and Deep models does not guarantee performance. Not to mention, these models can take hours or sometimes even days to train (even on GPUs). Research into pruning and quantization has shown that the connections that actually matter in the model are only a small percentage of the whole spider web!!

For example, famous ImageNet models like AlexNet and VGG-16 have been compressed to 40–50 times their original size, without any loss of (and actually a slight gain of) accuracy. This dramatically increases their inference speed and the ease with which they can be adapted across various devices.

Enough with the convincing, let's talk about the techniques involved!!

Model compression can be divided into two broad categories,

Pruning : Removing redundant connections present in the architecture. Pruning involves cutting out unimportant weights (which are usually defined as weights with small absolute value).

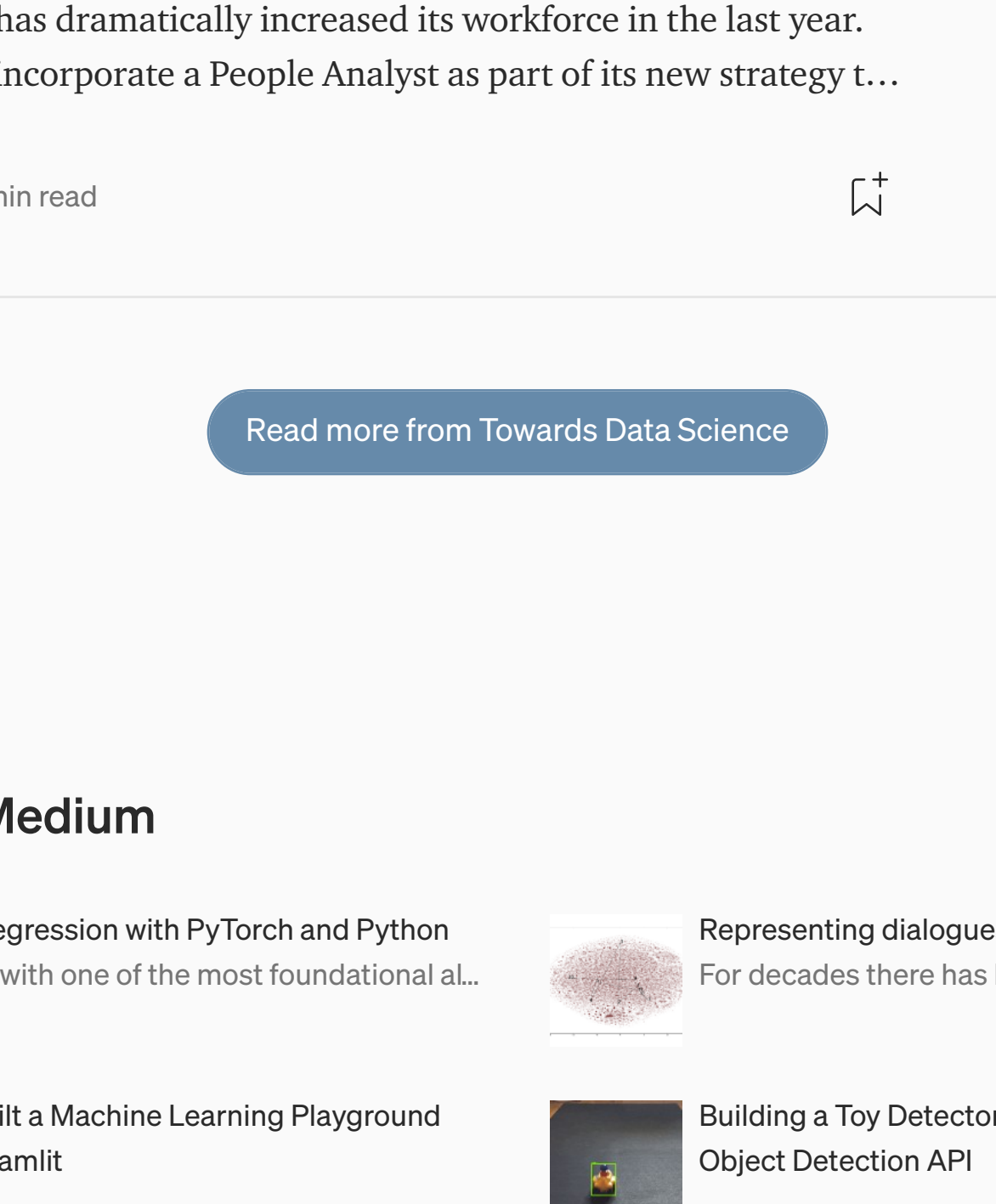


Obviously the new model formed will have lower accuracy since the model was actually trained for the original connections. That is why the model is fine tuned after pruning to regain the accuracy. It is noted that fully connected layers and CNNs can usually go upto **90% sparsity** without losing any accuracy.

Quantization : Quantization involves bundling weights together by clustering them or rounding them off so that the same number of connections can be represented using lesser amount of memory.

Quantization by doing clustering/bundling and thus using lesser number of distinct float values to represent more number of features is one of the most common techniques. Another common technique that forms the skeleton for a lot of quantization methods is converting floating point weights to fixed point representation by rounding off.

Again, as it was with pruning, we need to fine tune the model after quantization. The important point here is that the property that was given to the weights while quantization should be maintained through the fine tuning too. That is why specific ways of fine tuning are used which are tailored to match the quantization method.



Look at the image above for an example of quantization by clustering. Weights of same color are clustered together and represented by their centroid. This decreases the amount of data required to represent these weights. Earlier it required 32bits*16 = 512 bits to represent them. Now it will only take 32bits*4 + 2bits*16 = 160 bits to represent them. During the fine tuning, the gradient for all the weights which belong to the same color are summed and then subtracted from the centroid. This makes sure that the clustering made during quantization is maintained through the fine tuning.

What's next?

Deep Learning model pruning and quantization are relatively new fields. While there have been significant success in this field, it still has a long way to go. The next focus of the domain should be creating open source and easily accessible pipelines for transferring common Deep Learning models to embedded systems like FPGAs.

This blog is a part of an effort to create simplified introductions to the field of Machine Learning. Follow the complete series here

Machine Learning : Simplified

Know it before you dive in

towardsdatascience.com

Or simply read the next blog in the series

High Frequency Trading (HFT) with AI : Simplified

Take a peek into the ever-competing world of HFTs and how is AI becoming a part of it.

towardsdatascience.com

References

- [1] Han, Song, Huizi Mao, and William J. Dally. “Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding.” *arXiv preprint arXiv:1510.00149* (2015).
- [2] Jia, Haipeng, et al. “DropPruning for Model Compression.” *arXiv preprint arXiv:1812.02035* (2018).
- [3] Wang, Shuo, et al. “C-lstm: Enabling efficient lstm using structured compression techniques on fpgas.” *Proceedings of the 2018 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*. ACM, 2018.

Sign up for The Variable

By Towards Data Science

Every Thursday, The Variable delivers the very best of Towards Data Science: from hands-on tutorials and cutting-edge research to original features you don't want to miss. [Take a look.](#)

✉

Get this newsletter

More from Towards Data Science

Your home for data science. A Medium publication sharing concepts, ideas and codes.

Follow

David W. Dai · Jun 24, 2019 ★

Can AI Read Chest X-rays like Radiologists?

Using adversarial networks to achieve human-level performance for chest x-ray organ segmentation — Healthcare Needs AI to Scale Today, only about 10% of 7B population in the world have access to good healthcare...

Machine Learning · 8 min read

🔖

+

Post a quick thought or a long story. It's easy and free.

Write with the app

Tarun Acharya · Jun 24, 2019

Prime numbers and Goldbach's conjecture visualization.

Mathematics is everywhere, from the pattern in the leaves of a sunflower to the reflecting angles of the sun rays. It is merely the language of nature...

Mathematics · 3 min read

🔖

+

Mathanraj Sharma · Jun 24, 2019 ★

Hypothesis Testing with Numpy

Howdy fellow Data Science enthusiasts, we all do hypothesis testing to infer the results from a sample data from a larger population. The test helps us to find whether our Null Hypothesis is True or false. First, let us understand...

Data Science · 3 min read

🔖

+

Raj · Jun 24, 2019 ★

Forming a Machine Learning/AI team?

Don't forget other key skills you'll need — Enterprises are overwhelmed by the hype created around Machine Learning & Artificial Intelligence, yet I'm sure no one wants to miss this bus and are striving to get their act together...

Data Science · 5 min read

🔖

+

Pablo Alvarez · Jun 24, 2019 ★

People Analytics

Human Bias in recruitment selection — Case Study I — INTRODUCTION
The startup XYZ has dramatically increased its workforce in the last year. Now, it wants to incorporate a People Analyst as part of its new strategy t...

Data Science · 8 min read

🔖

+

Read more from Towards Data Science

More from Medium

Linear Regression with PyTorch and Python

Working with one of the most foundational al...

Representing dialogues

For decades there has been a lot of interest ...

How I Built a Machine Learning Playground with Streamlit

Building a Toy Detector with Tensorflow Object Detection API

Sentiment analysis with BERT in PyTorch

That day in autumn of 2018 behind the walls ...

Theseus, the magnetic mouse

A life-sized magnetic mouse learns its way around a ma...

ML Kit: Boosting capabilities of an Android or iOS app

RescueForest: Predicting Emergency Response with Random Forests

🏠

🔍

✍