


# **MLOps on the Edge**



# About **difference-engine.ai**

 **difference-engine.ai** is a technology consulting firm specialising in solving business problems through Data Science, Applied Machine Learning and AI Driven Products.

We are a team of 13 ML Engineers.

# About Me



## Prathamesh Sarang

- 6+ years experience in Software Engineering and Data Science
- Ex-Systems Engineer at Infosys
- Ex-Data Scientist at Lemoxo Technologies and Damco (Part of AP Moller Maersk Group)
- Currently working as Machine Learning Engineer (Data Products) at [difference-engine.ai](https://difference-engine.ai)
- Working at the intersection of Data Science and Software Engineering

# What are the interesting things I'm doing?

- Not much traditional Machine Learning these days



# What are the interesting things I'm doing?

- Not much traditional Machine Learning these days
- Working on Web applications



# What are the interesting things I'm doing?

- Not much traditional Machine Learning these days
- Working on Web applications
- Majority my work is into engineering around ML applications 🧐



# Outside of work, what I do!

- Big history, horror and true crime fan, generally listen to podcasts, watch movies/series and read books



# Outside of work, what I do!

- Big history, horror and true crime fan, generally listen to podcasts, watch movies/series and read books
- I teach as well





# **What can you expect from this talk**



**difference-engine.ai**

**I have a few questions**



**difference-engine.ai**

# How many Technical folks?



difference-engine.ai

**What's the exposure of the crowd wrt  
ML?**

# How many have done ML Deployments?



difference-engine.ai

# What is Data Science, Machine Learning and AI?





**Mat Velloso**

@matveloso

Follow



Difference between machine learning  
and AI:

If it is written in Python, it's probably  
machine learning

If it is written in PowerPoint, it's  
probably AI

5:25 PM - 22 Nov 2018

7,622 Retweets 21,608 Likes



187



7.6K



22K



difference-engine.ai

**So many conflicting definitions!**

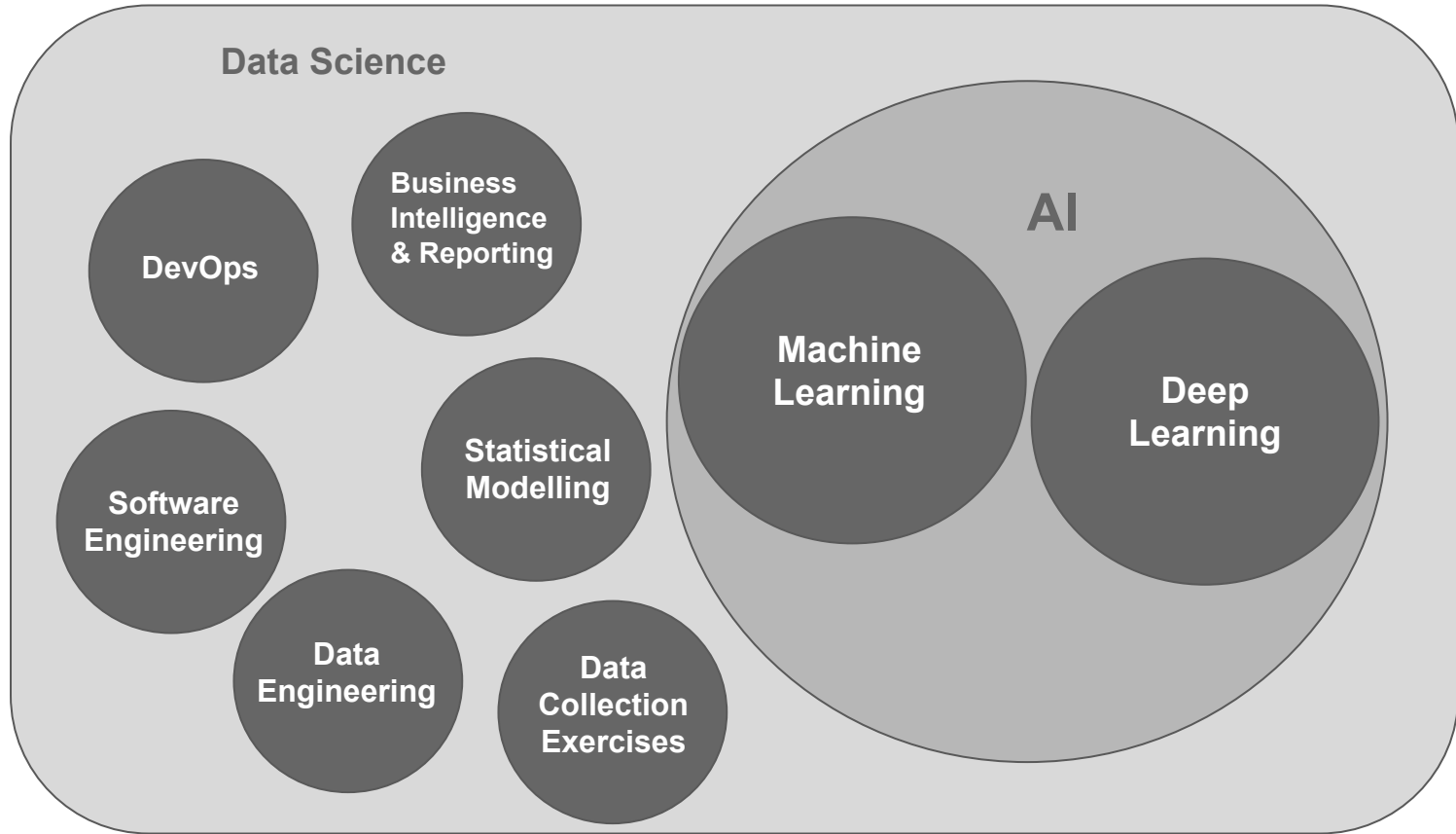


**difference-engine.ai**

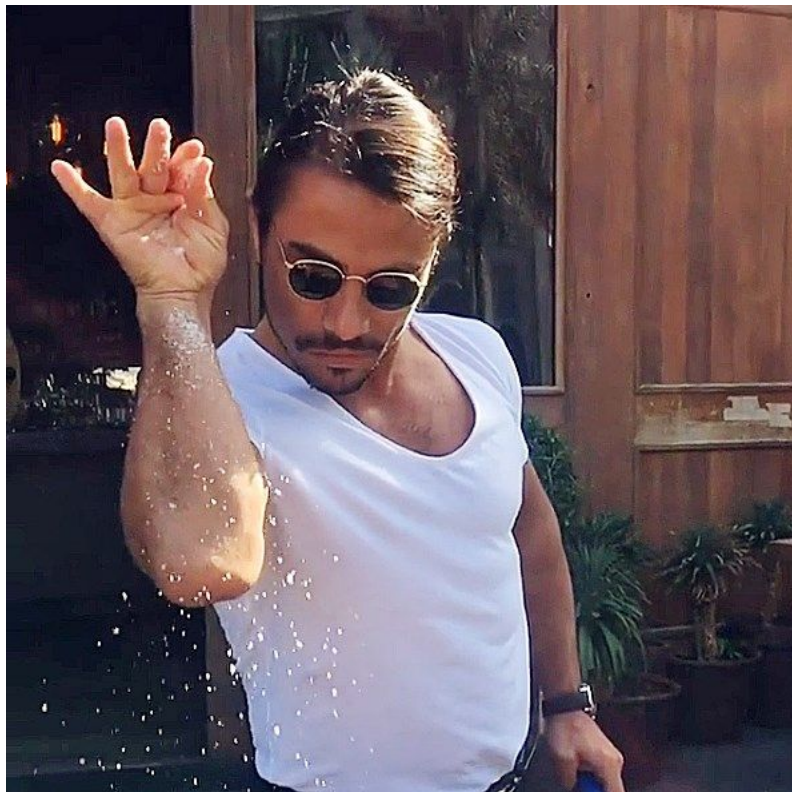


# I'll add one more!





# Take that with a pinch of salt



difference-engine.ai



Google Search

Google Maps

Google Translate

Google Lens



difference-engine.ai



**Content Recommendations for users**

**Artwork Personalization at Netflix**

**Data Science and the Art of Producing Entertainment at Netflix**



**difference-engine.ai**



**Contracting and Procurement**

**Downstream Retail**

**Shell Exploration**

# Closer home!

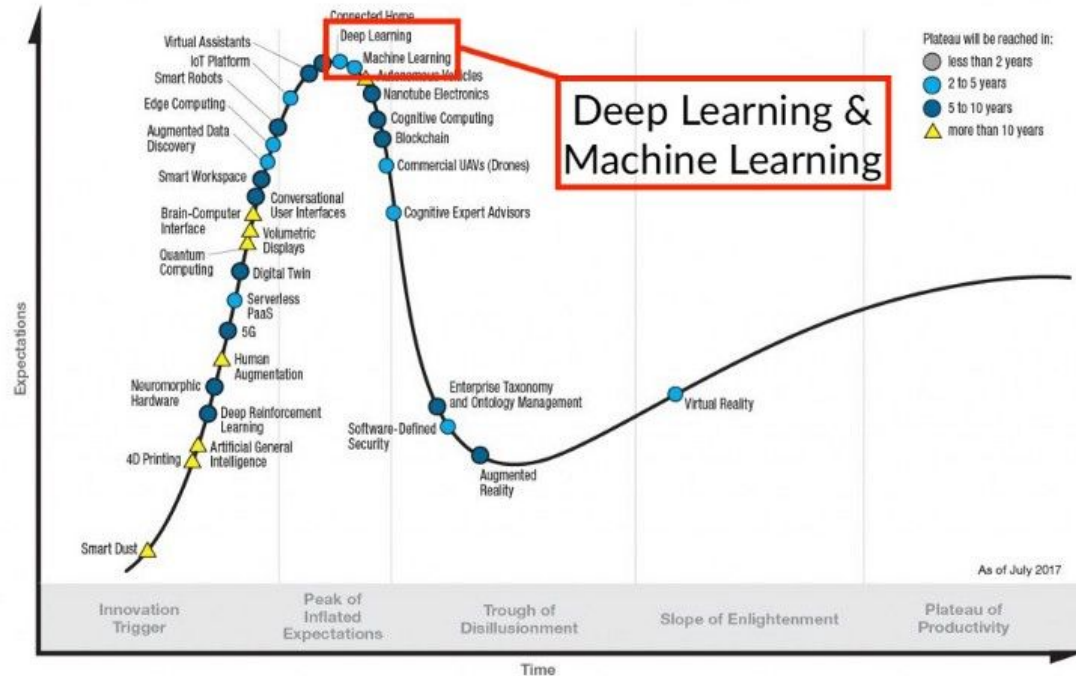


difference-engine.ai

# Machine Learning & Deep Learning hype



## Gartner Hype Cycle for Emerging Technologies, 2017



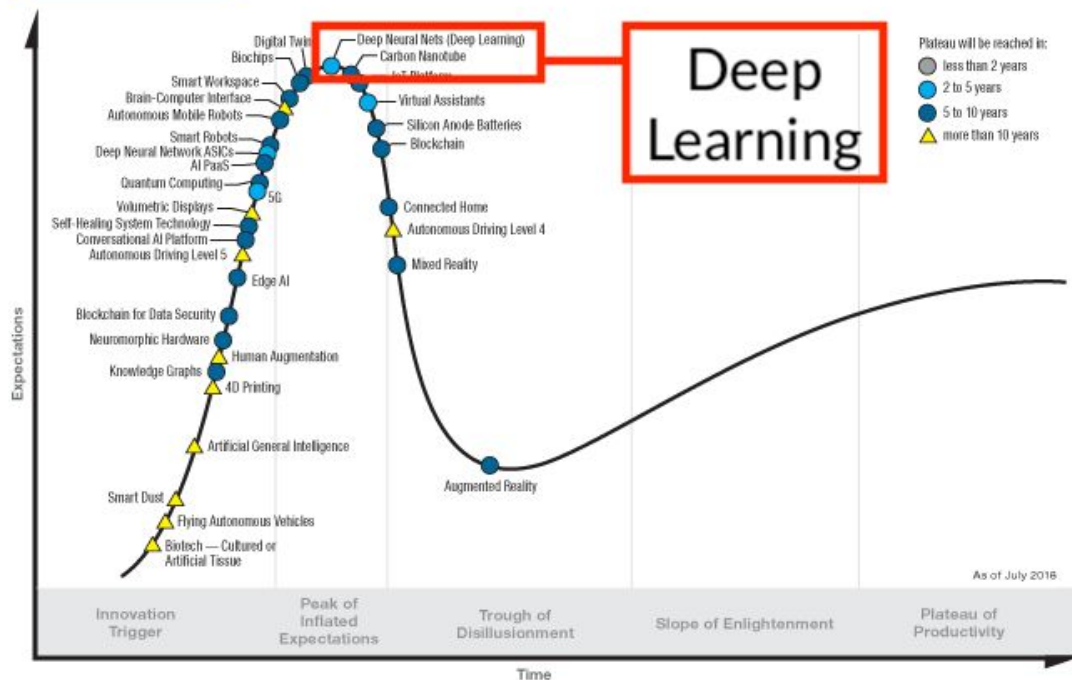
[gartner.com/SmarterWithGartner](https://gartner.com/SmarterWithGartner)

Source: Gartner (July 2017)  
© 2017 Gartner, Inc. and/or its affiliates. All rights reserved.

Gartner

erence-engine.ai

## Hype Cycle for Emerging Technologies, 2018



[gartner.com/SmarterWithGartner](https://gartner.com/SmarterWithGartner)

Source: Gartner (August 2018)  
© 2018 Gartner, Inc. and/or its affiliates. All rights reserved.

**Gartner.**

ference-engine.ai

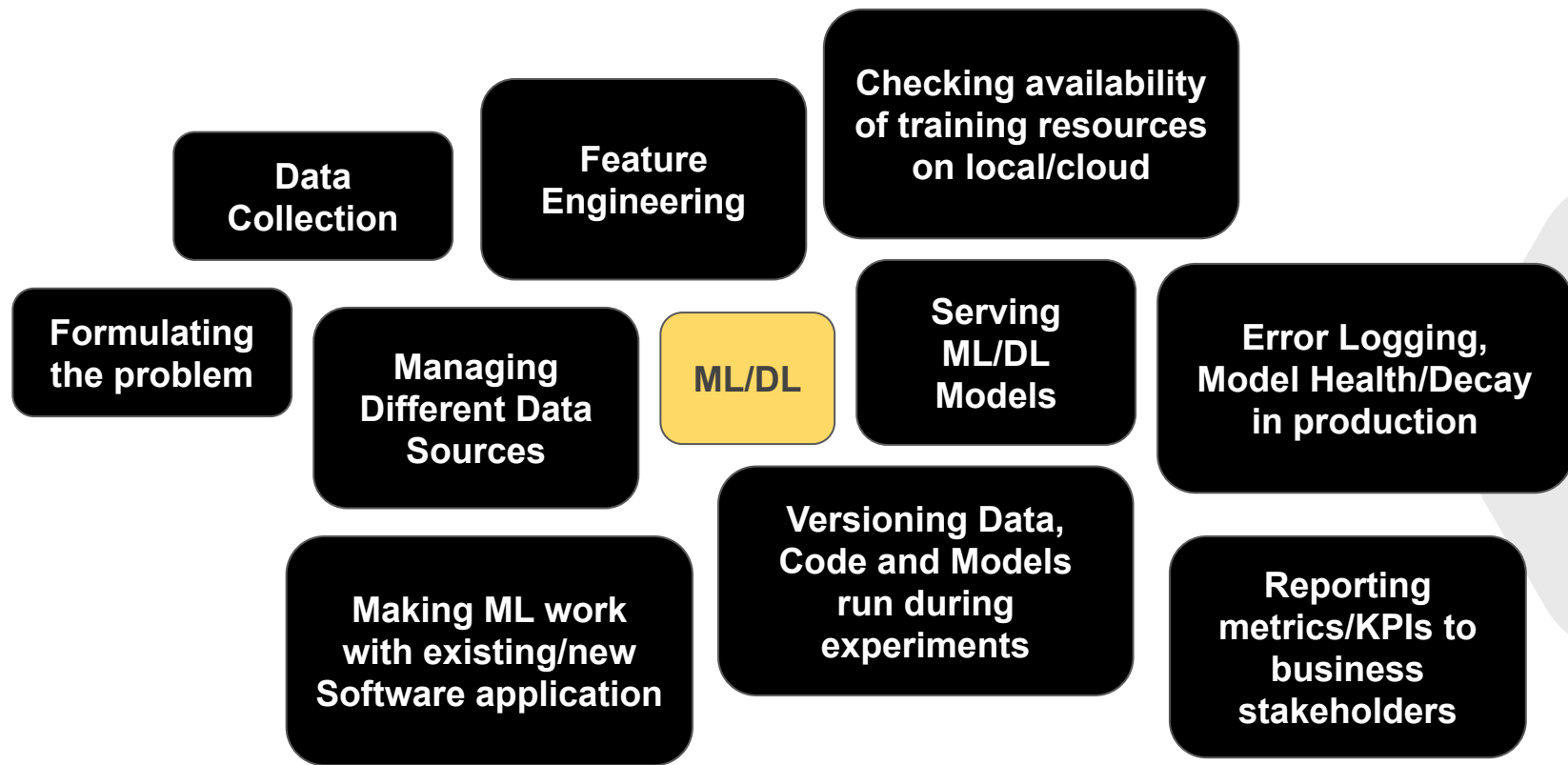
# Why all of this?

## To solve business problems!



**But ML isn't the most important part!**





# It's already hard!



# THE DATA SCIENCE HIERARCHY OF NEEDS

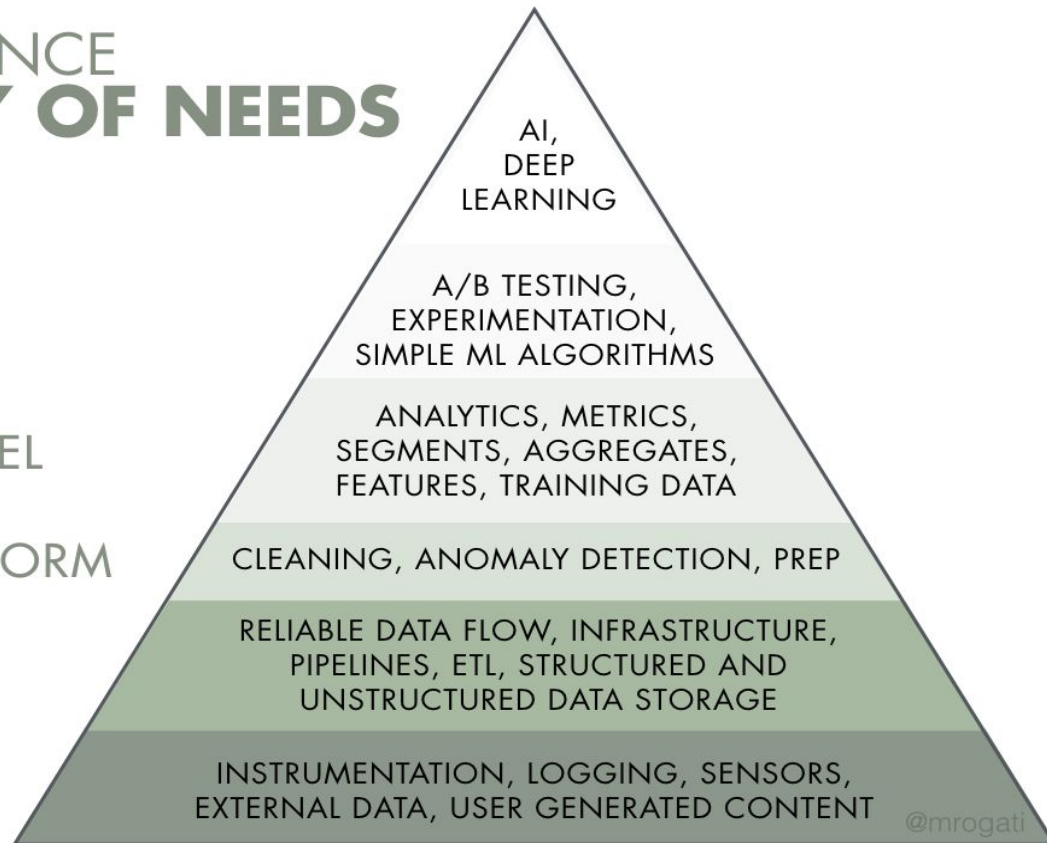
LEARN/OPTIMIZE

AGGREGATE/LABEL

EXPLORE/TRANSFORM

MOVE/STORE

COLLECT



Source: Monica Rogati's AI Needs of Hierarchy



difference-engine.ai

# **Mantra to do Data Science in an organization**

**(Shamelessly copied from Dr. D.J. Patil)**



**difference-engine.ai**



# 1x

**Prototyping Phase: Picking  
up problems and solving  
them**

**Just letting the business  
know that we are capable**



# 10x

**Real Deployment Strategy  
comes into play**



**difference-engine.ai**

# 100x

**In simple words, this is pure  
scale! Organization wide ML  
Adoption**



**difference-engine.ai**

**0 to 1x is cool!**



**Let's assume you were able to get a  
pretty compelling AUC 0.834**

**Or the Deep Learning architecture you  
stumbled on is great!**



# Your Data Scientist be like!



Kabhi Kabhi lagta hai ki  
apun hi Bhagwan hai.



difference-engine.ai

# While deployment!



**1x to 10x is hard for a lot of people**  
**Friction between Software/Ops and  
Data Science**



**Code or algorithms don't make ML hard,  
people do!**



**difference-engine.ai**

# 1x solutions



difference-engine.ai

# **1x Solution**

## **Dashboards**

## **Jupyter Notebooks**

## **Powerpoint Presentation**



# Problems with 1x solution



## Ease of Use

**Scaling it to actual users**

**Are you really serious about doing ML?**



# 10x solution



difference-engine.ai

# **10x Solutions**

**Dashboards**

**Data Engineering**

**APIs**

**On Device ML**



**difference-engine.ai**

# Where 10x solutions might fail



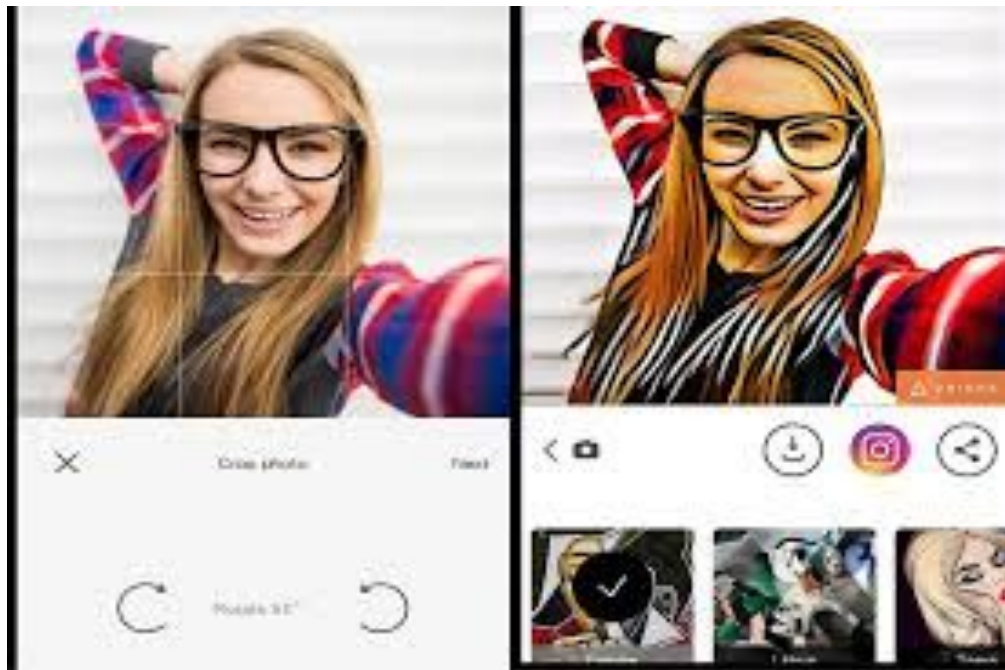


**Latency**

**Cost per inference**

**Security & Privacy**

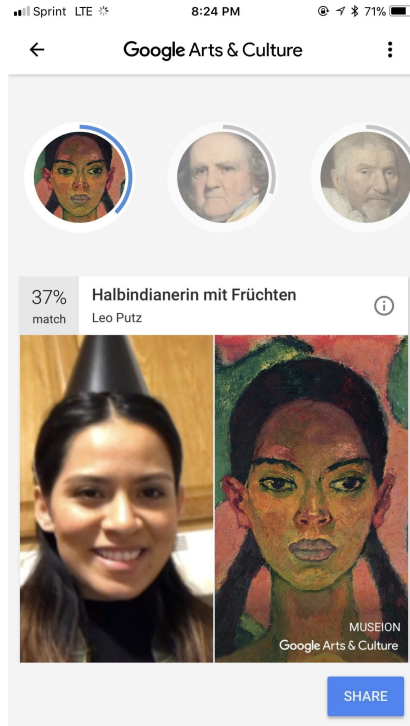
# Prisma



**Prisma suffered from Latency when it started  
off**



# Google Arts and Culture App



# **Big backlash on Social Media due to the intrusive nature of the application**

# IoT/Edge devices

Why am I awakened by a freezing house in 14 degree weather? Furnace is working. Vents are well maintained.

Oh, I see why.

The goddamn @ecobee server for the networked thermostat is down. And what could go wrong, in depending on the internet to keep the child warm?



This site can't be reached

www.ecobee.com refused to connect.

Try:

- Checking the connection
- [Checking the proxy and the firewall](#)

ERR\_CONNECTION\_REFUSED

Details

Reload



difference-engine.ai

**Right now everyone is lingering around  
this area**



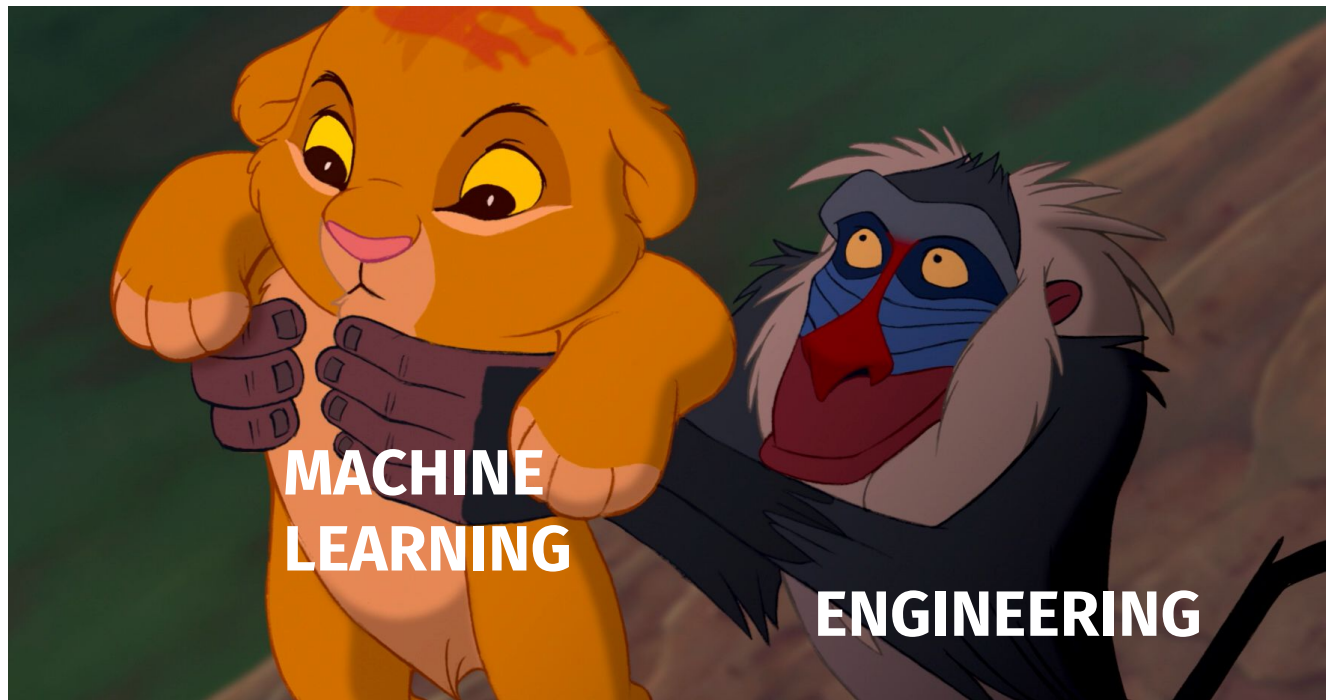
**100x**



**difference-engine.ai**



# 100x Solution



# Why are we all here?



**ML on IoT/Edge devices is a great use case**



**difference-engine.ai**

**'Edge' refers to the computing infrastructure that exists close to the sources of data**

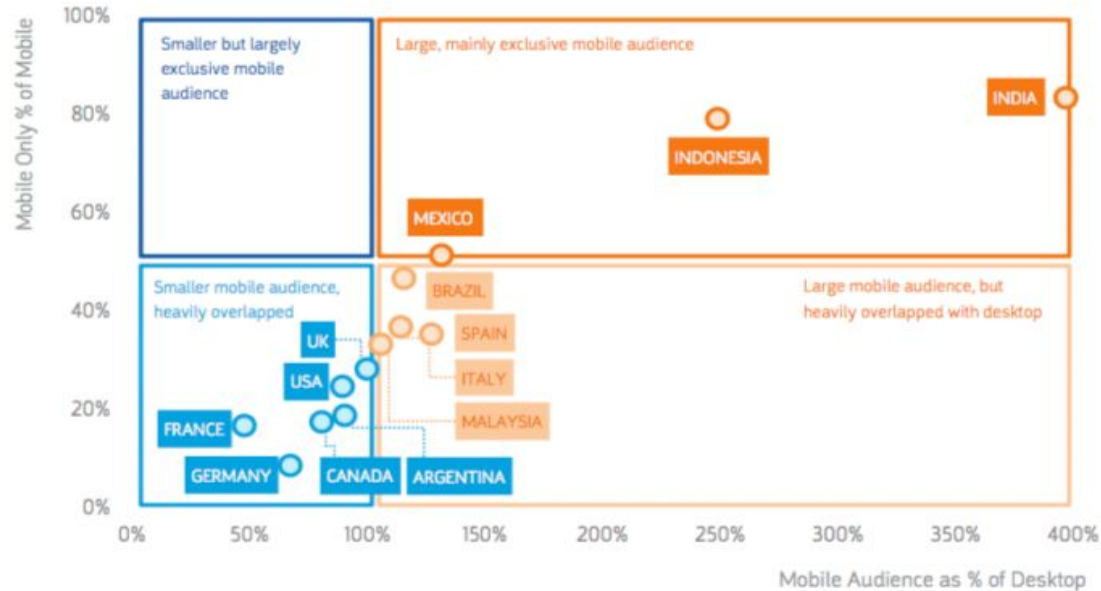


# Can you name some edge devices?



# Impact of mobile devices isn't to be ignored

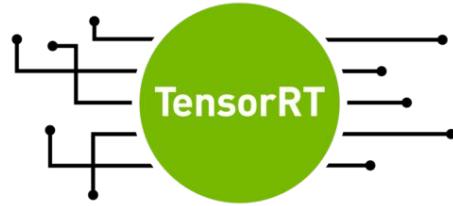
'Mobile-Firstness' of Markets' Total Digital Populations



# Current state for OnDevice ML



ML Kit



ONNX



TensorFlow

Lite



difference-engine.ai

**Memory constraints on device**

**Models are huge!**

**Trade Offs: Metrics vs Usability**

**No platform agnostic solutions**





# Better workflows for mobile deployments



OR

+  TensorFlowLite +



# Demo



difference-engine.ai

# Better workflows for client-side deployments



ONNX

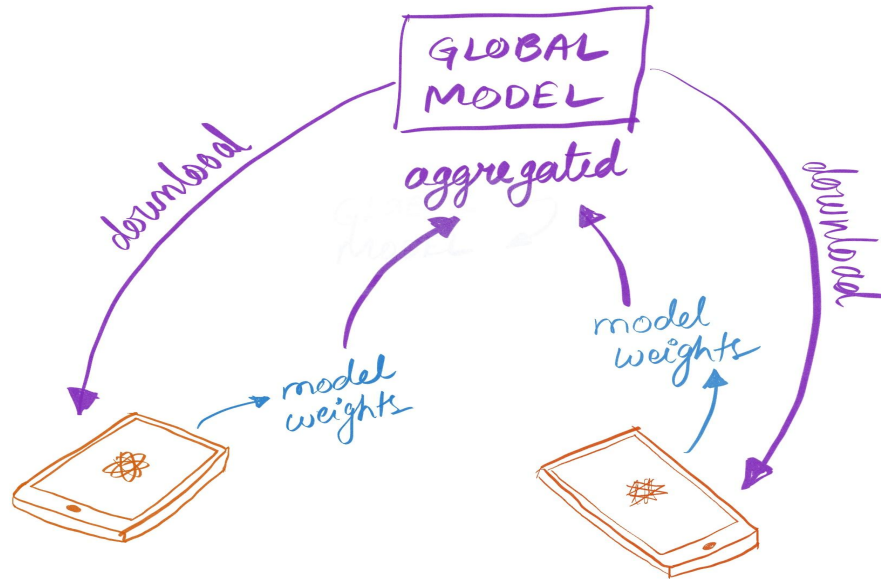


difference-engine.ai

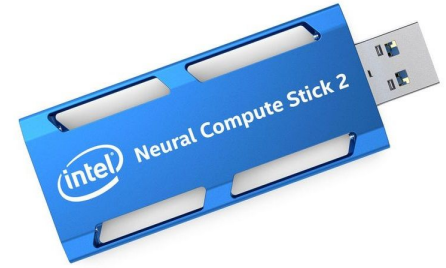
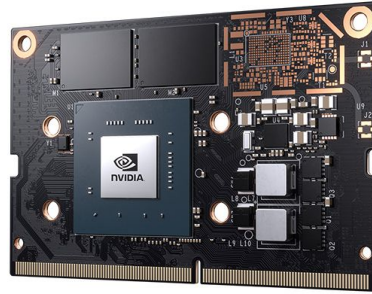
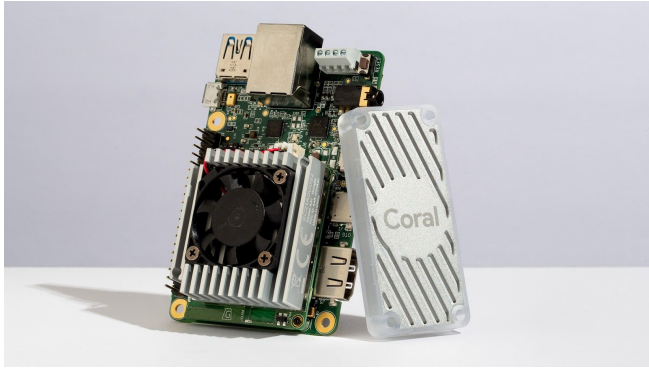
# Recent developments



# Federated Learning



# Custom Hardware to do on-device ML



# Improved Model Compression techniques

**Quantization**

**Pruning**

**Knowledge Distillation**

**Low Rank approximation**



# Improved Model Architectures





# Where can I read?

Recent Advances in Efficient Computation of Deep Convolutional Neural Networks:

<https://arxiv.org/pdf/1802.00939.pdf>

Awesome Model compression and abstraction:

<https://github.com/memoiry/Awesome-model-compression-and-acceleration/blob/master/README.md>

Model Compression Papers: <https://paperswithcode.com/task/model-compression>



**Thank You**



# Catching hold of me

Email: [prathamesh.b.sarang@gmail.com](mailto:prathamesh.b.sarang@gmail.com)

## Shameless promotion

LinkedIn: [Prathamesh Sarang](#)

Personal blog: <https://pratos.github.io>

Twitter: [prat0s](#)



difference-engine.ai