

The logo for AWS re:Invent features the words "AWS" and "re:Invent" stacked vertically. "AWS" is in a smaller, sans-serif font above "re:Invent", which is in a larger, bold, sans-serif font. The entire logo is white against a red-to-purple gradient background.

AWS
re:Invent

AIM 302

Machine Learning at the Edge

Julien Simon
Principal Tech. Evangelist, AI/ML
Amazon Web Services
[@julsimon](https://twitter.com/julsimon)

Brian Kursar
Chief Data Scientist
Toyota Connected Data Services
Brian.Kursar@toyotaconnected.com
[@briankursar](https://twitter.com/briankursar)

Nimish Amlathe
Data Scientist
Toyota Connected Data Services
Nimish.Amlathe@toyotaconnected.com

Agenda

- Machine Learning at the Edge?
- Leveraging AWS Services
- Case study: Toyota Connected Data Services
- Alternative scenarios
- Optimizing for inference at the Edge
- Getting started

Machine Learning at the Edge?

Most machine data never reaches the cloud



Medical equipment



Industrial machinery

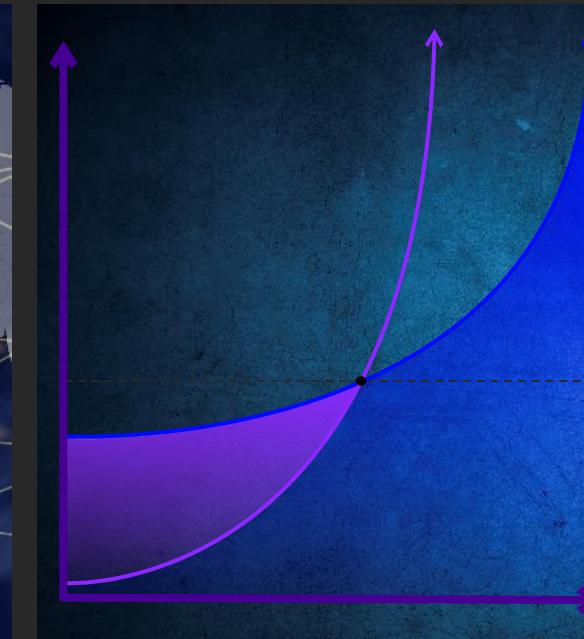


Extreme environments

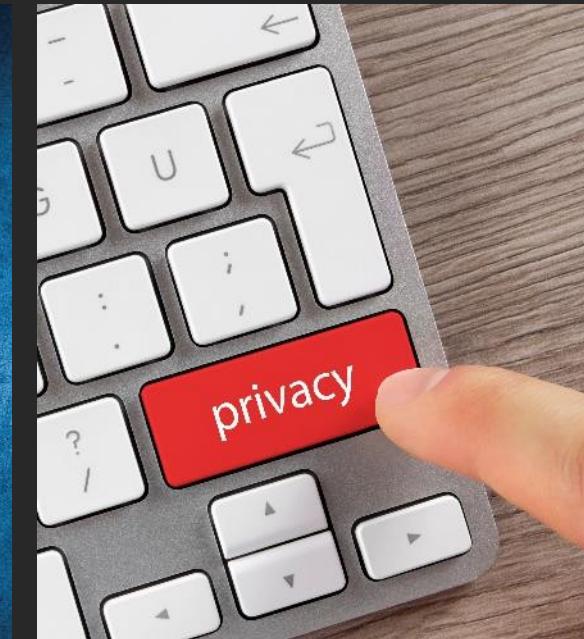
Why this problem isn't going away



Law of physics



Law of economics

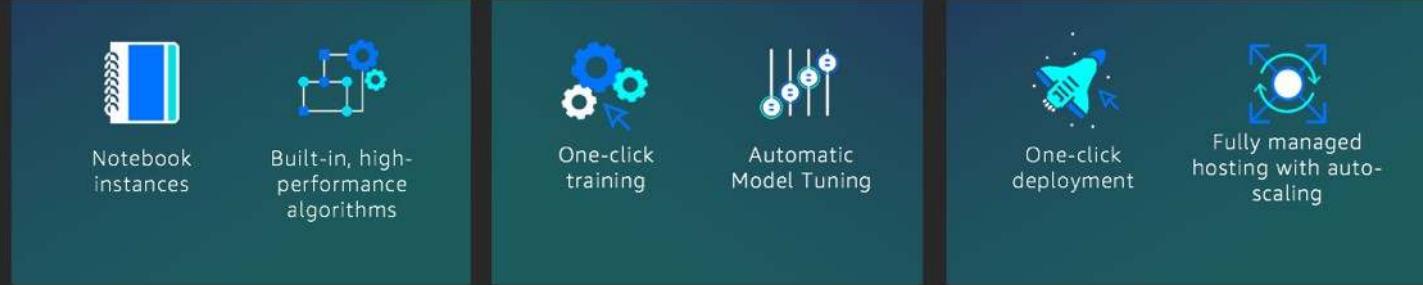


Law of the land

Leveraging AWS Services

Experimenting and training in the Cloud

Amazon SageMaker



Build

Train

Deploy

AWS Deep Learning AMI



Amazon EC2

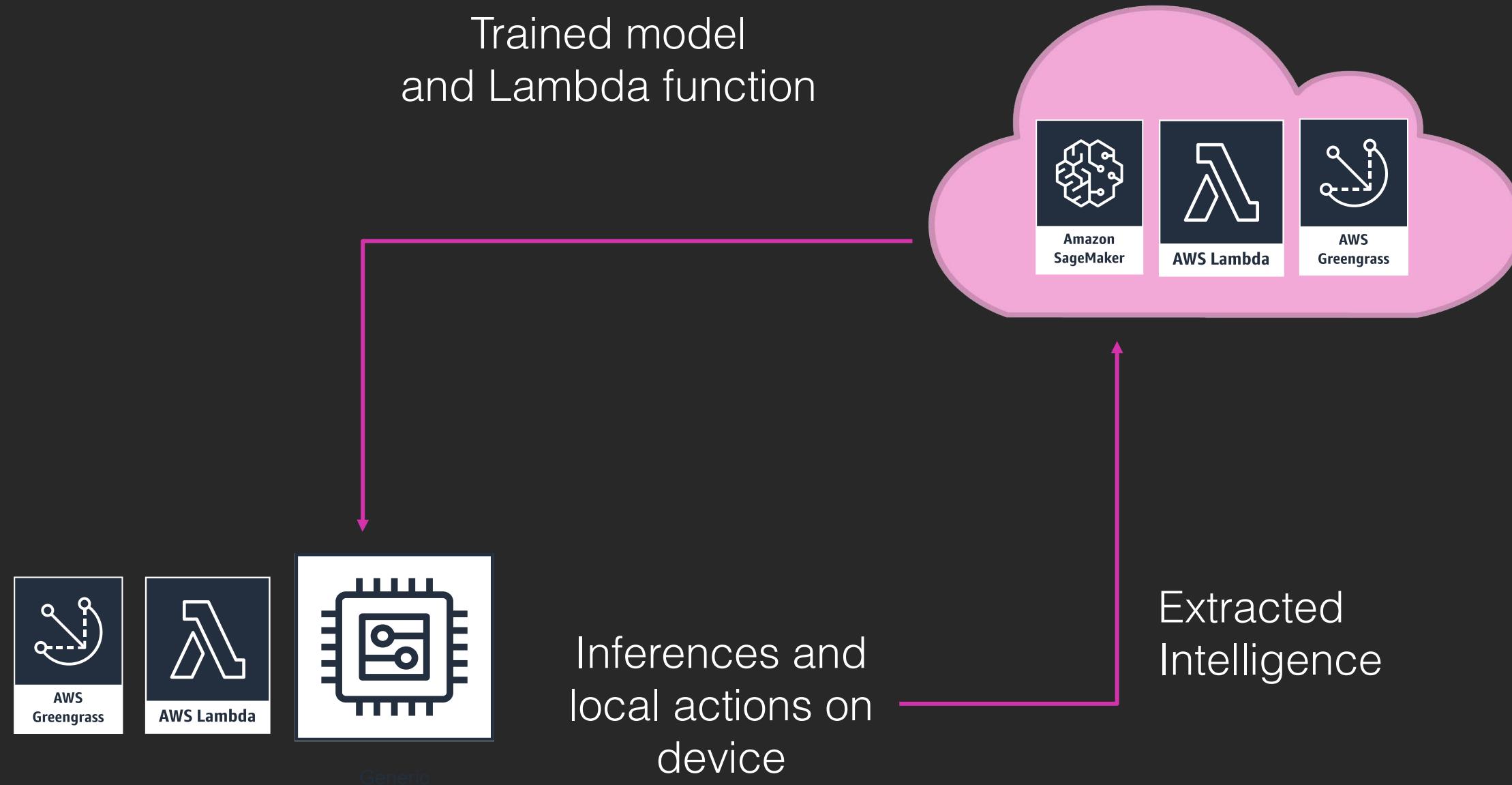


C
5



p3

Deploying and predicting with AWS Greengrass



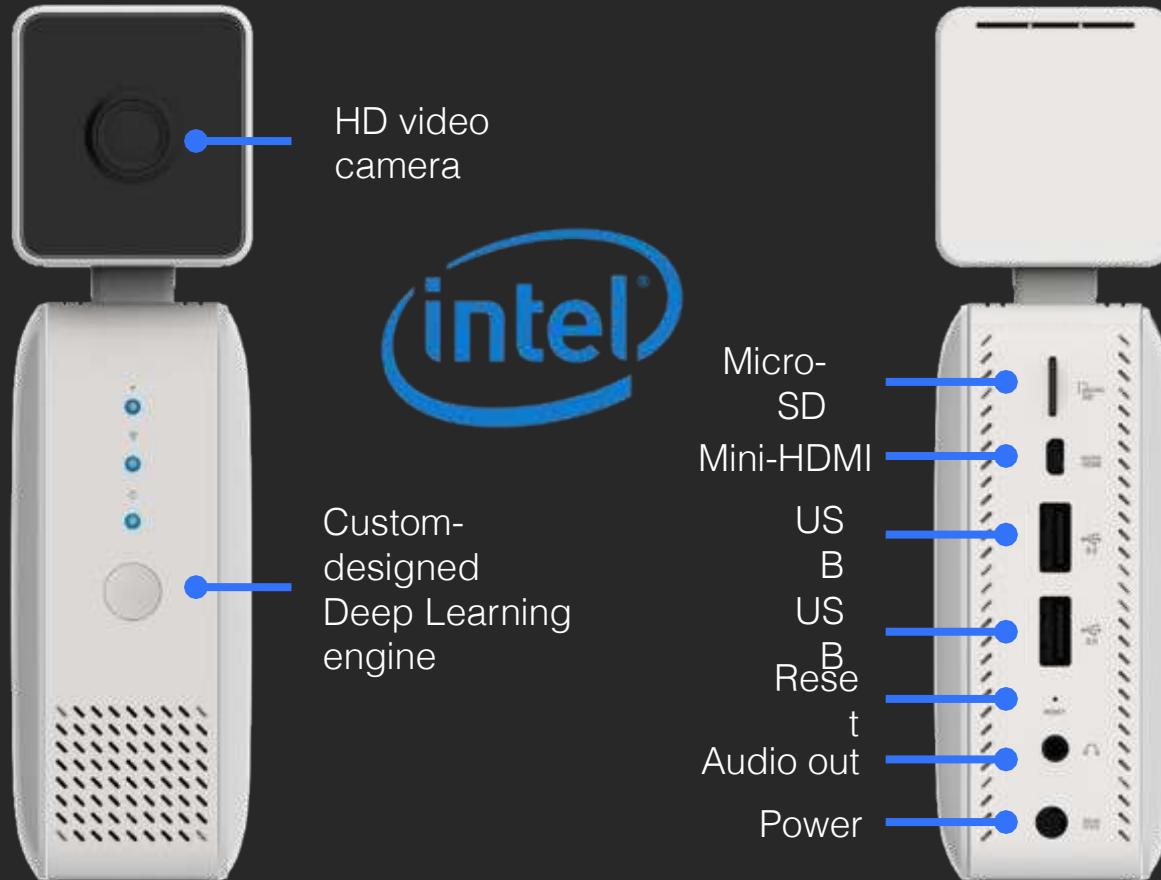
Deploying with AWS Greengrass

- Train a model in Amazon SageMaker (or import your own)
- Write an AWS Lambda function for prediction
- Add both as **resources** in your AWS Greengrass group
- Let AWS Greengrass handle deployment and updates

Best when	Requirements
You want the same programming model in the Cloud and at the Edge.	Devices are powerful enough to run AWS Greengrass
Code and models need to be updated, even if network connectivity is infrequent or unreliable.	Devices are provisioned in AWS IoT Core (certificate, keys).
You need one device in the group to perform prediction on behalf of other devices.	

AWS DeepLens

AWS DeepLens



A new way to learn

Custom built for Deep Learning

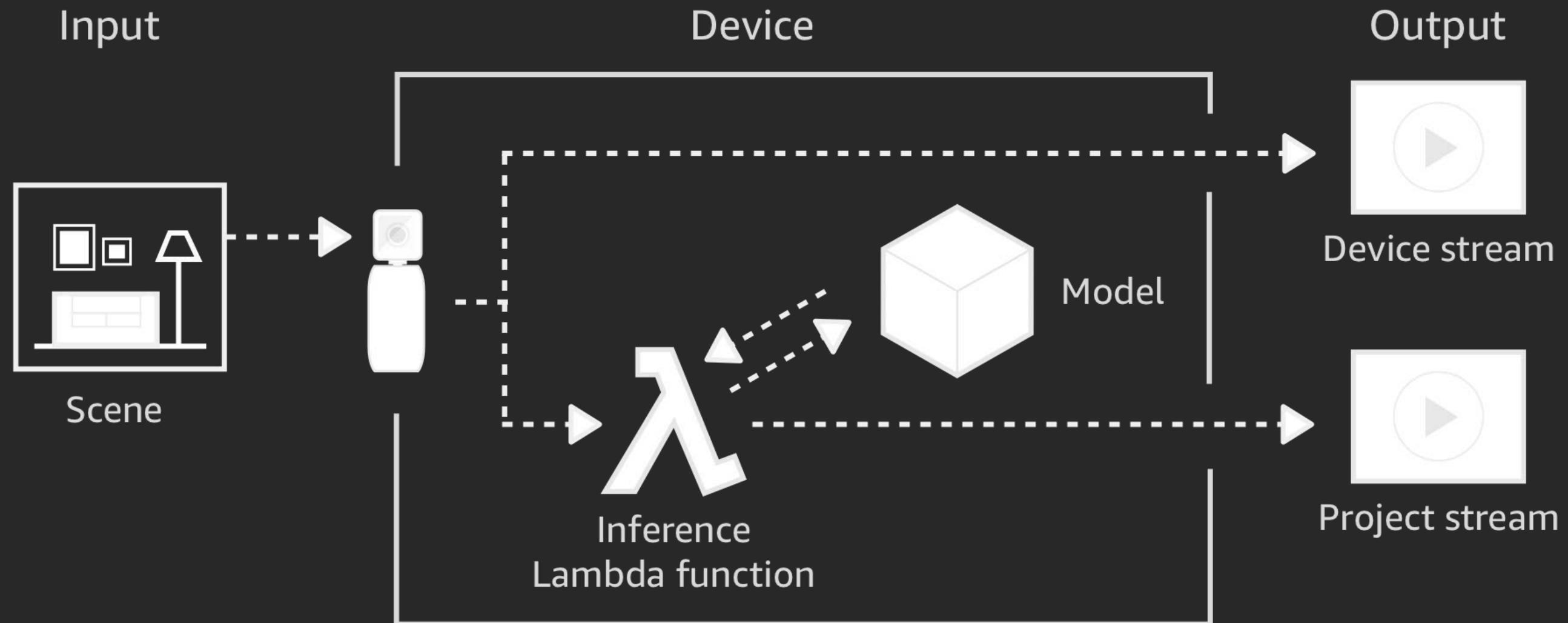
Broad framework support

Deploy models from Amazon SageMaker

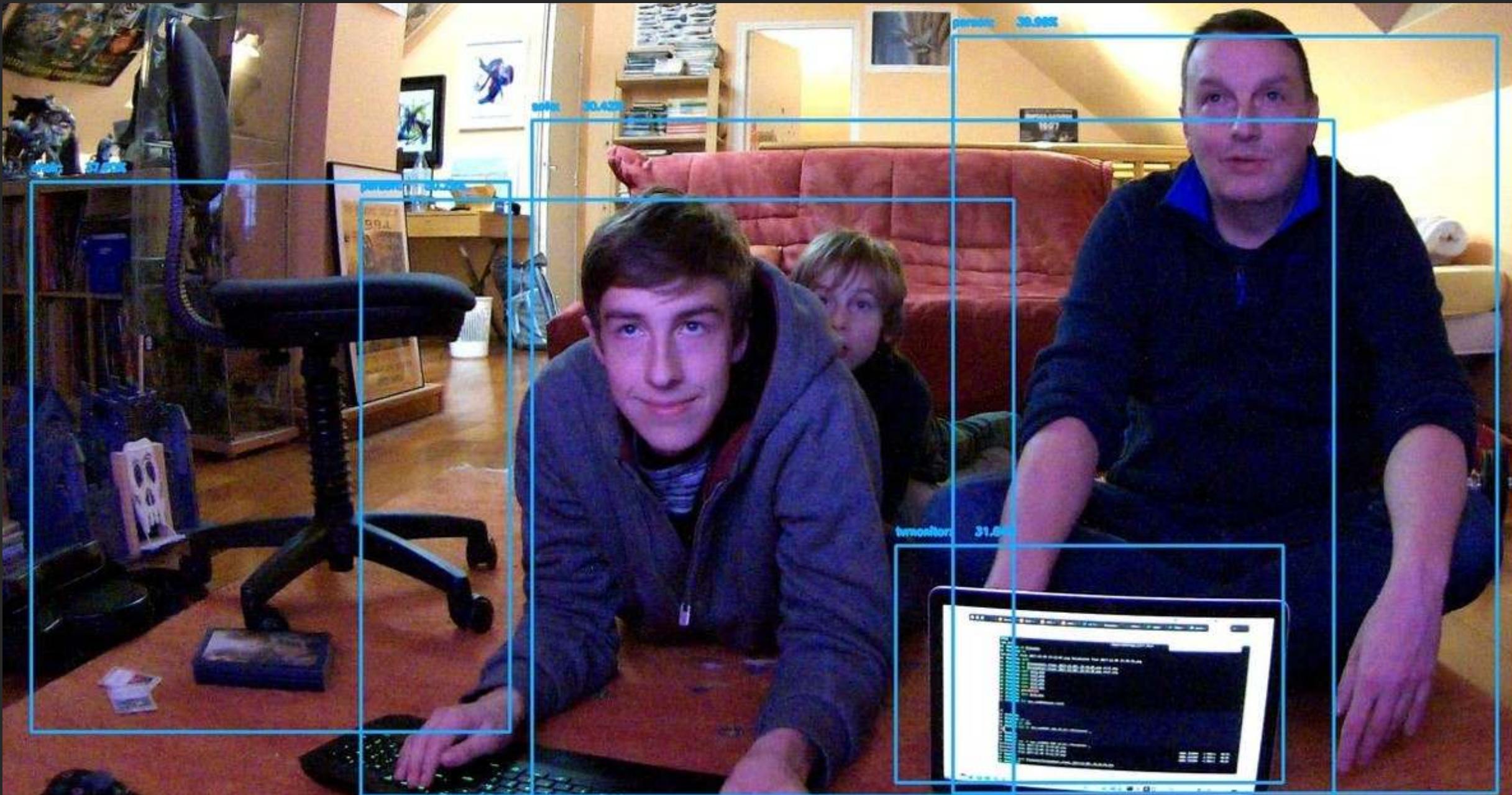
Integrated with AWS

Fully programmable with AWS Lambda

AWS DeepLens



Object detection with AWS DeepLens



Using your own models with AWS DeepLens

- AWS DeepLens can run [TensorFlow](#), [Caffe](#) and [Apache MXNet](#) models.
 - Inception
 - MobileNet
 - NasNet
 - ResNet
 - Etc.
- Train or fine-tune your model on [Amazon SageMaker](#)
- Deploy to DeepLens with [AWS Greengrass](#)





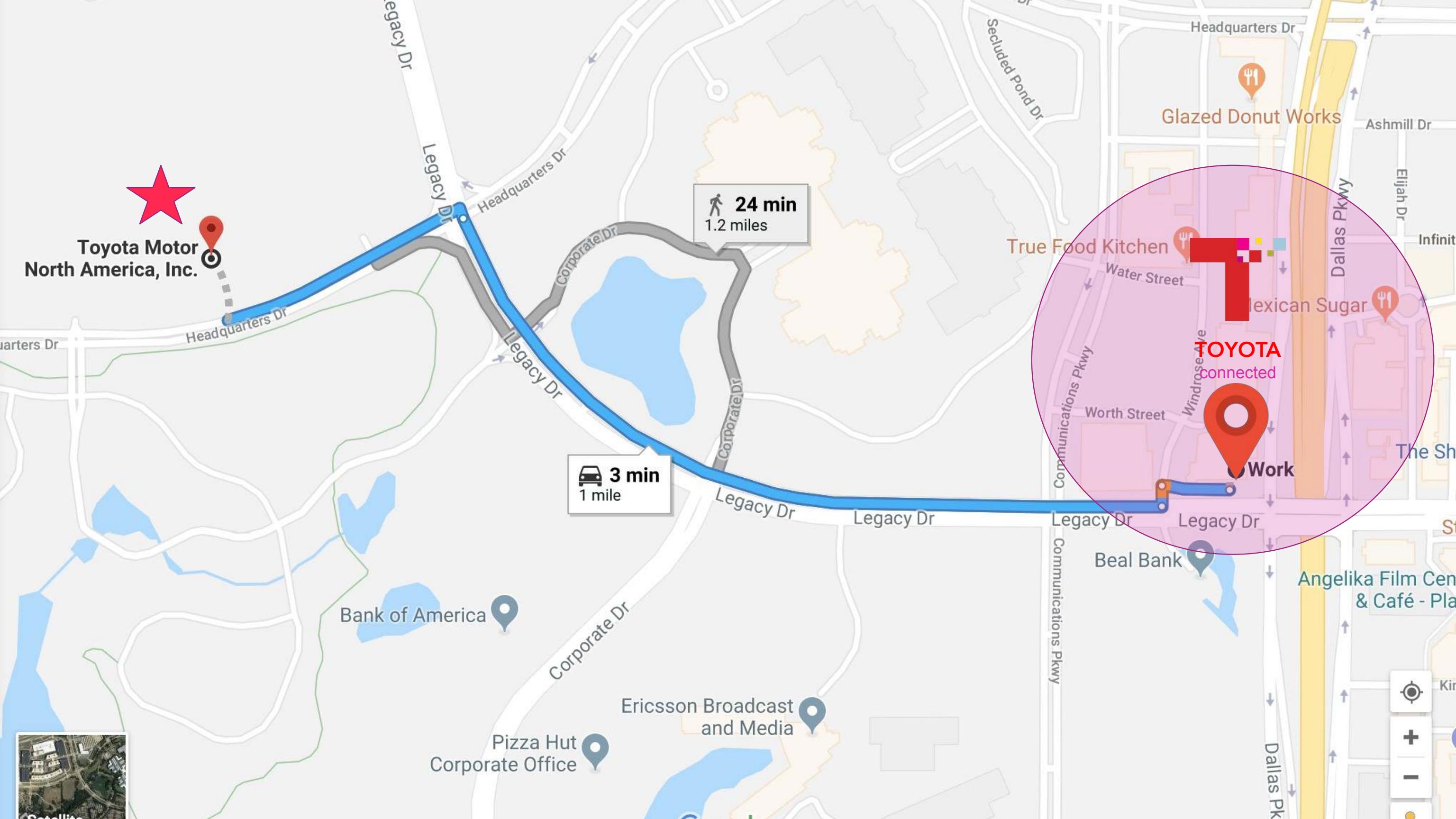
TOYOTA
connected

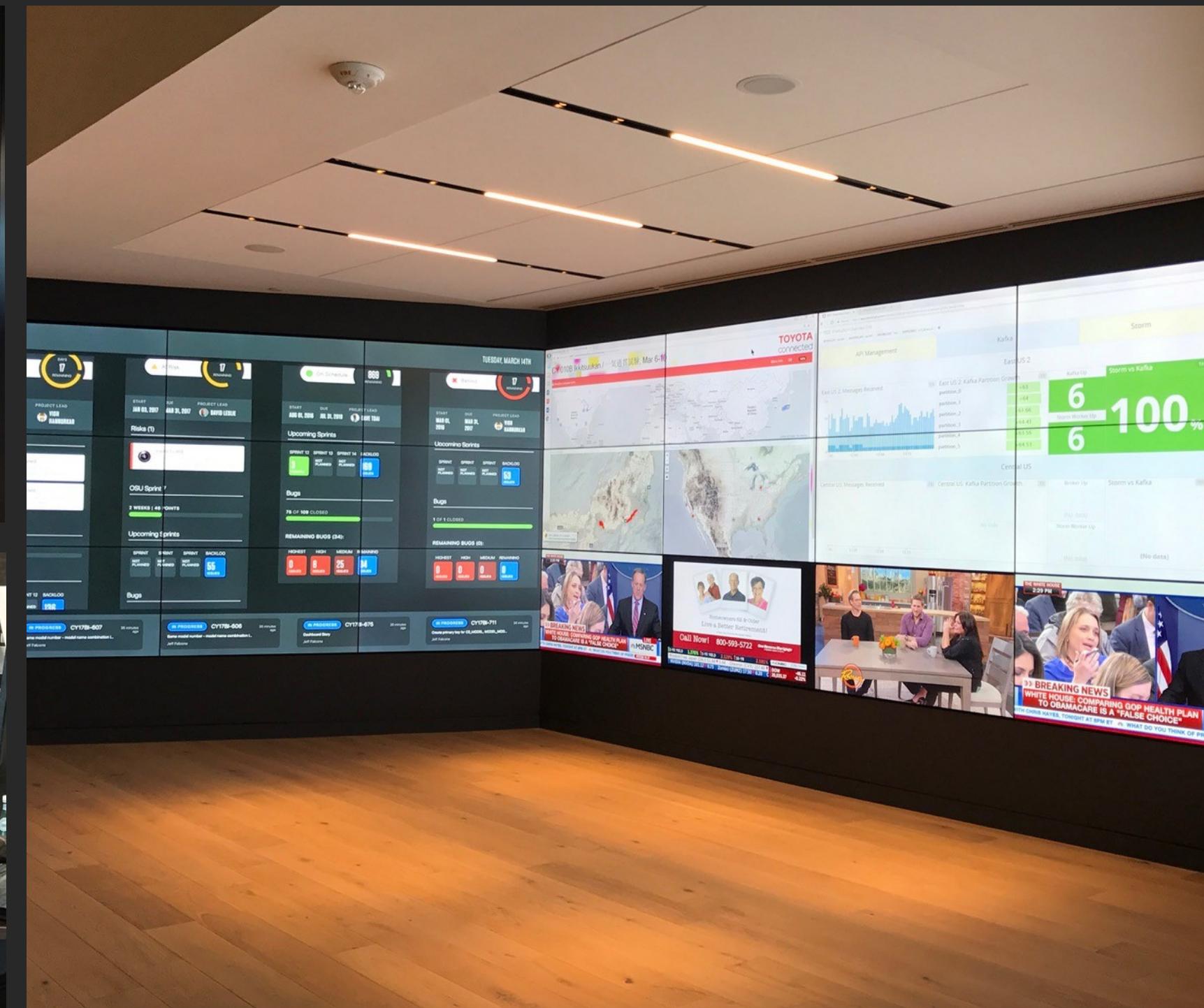
DEEP LEARNING AT THE EDGE

BRIAN KURSAR
VICE PRESIDENT AND CHIEF DATA SCIENTIST

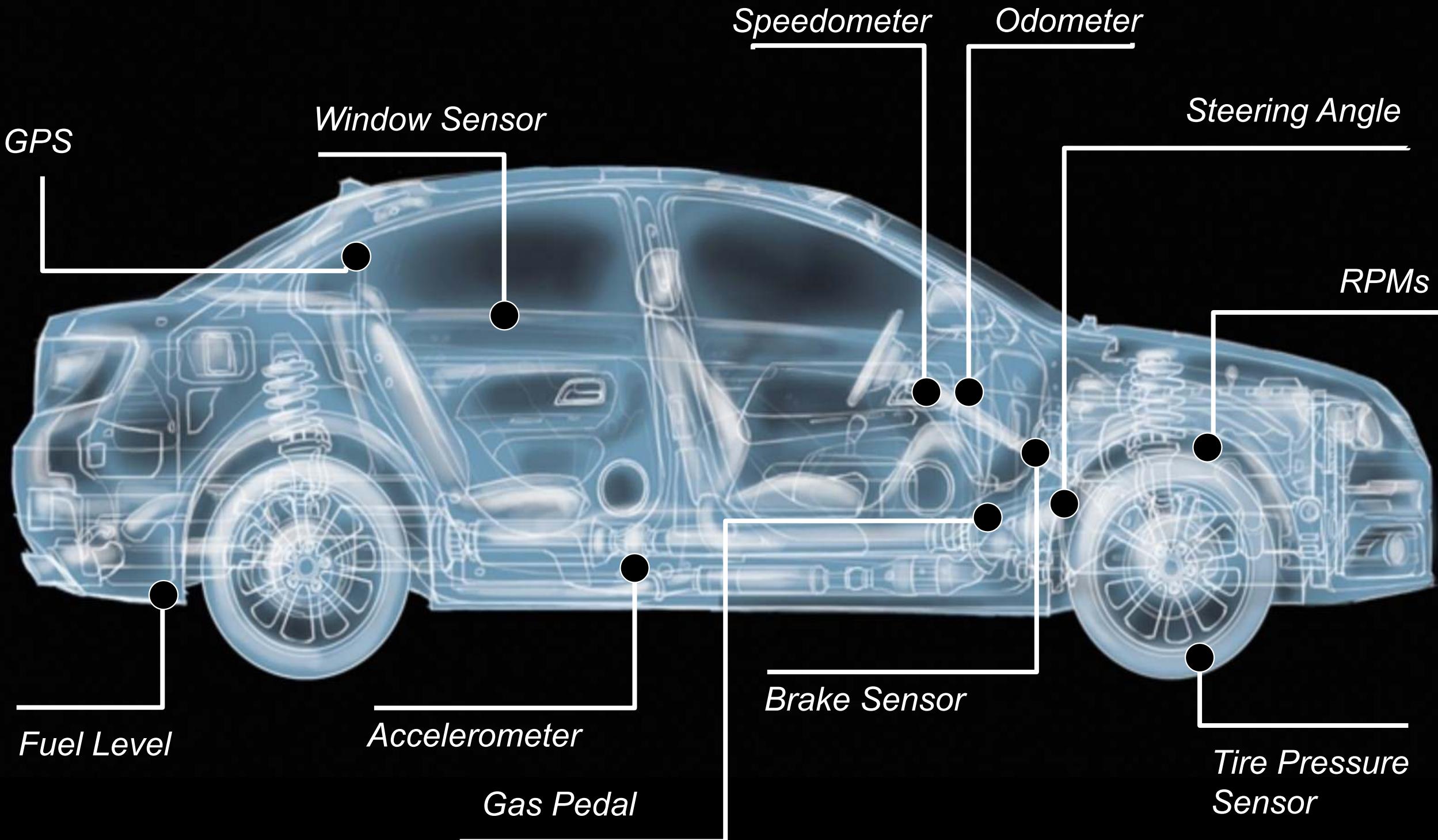
NIMISH AMLATHE
DATA SCIENTIST

TOYOTA CONNECTED DATA SERVICES









48 average driving minutes per day

500+ unique data points
generated every 200 milliseconds
captured in **real-time**

7.2 million data points
per connected vehicle
per day





Data services that drive
customer satisfaction



Data insights that
improve our products



2021: vehicles will be able to upload video in addition to telemetry

CARS FOR GOOD

Can we use video from vehicle cameras
to do **positive things** for the World?



CARS FOR GOOD



Detect and alert other drivers on dangerous weather conditions

CARS FOR GOOD



Keep our environment clean

CARS FOR GOOD



Reduce graffiti

CARS FOR GOOD



Automatically improve **maps** in real-time

CARS FOR GOOD



4 People

Optimize public transportation routes

CARS FOR GOOD



4 Occupants

Optimize Emergency services

CARS FOR GOOD



Improve **safety outside** the car

CARS FOR GOOD



Improve safety inside the car

We can,
but...

CARS FOR GOOD

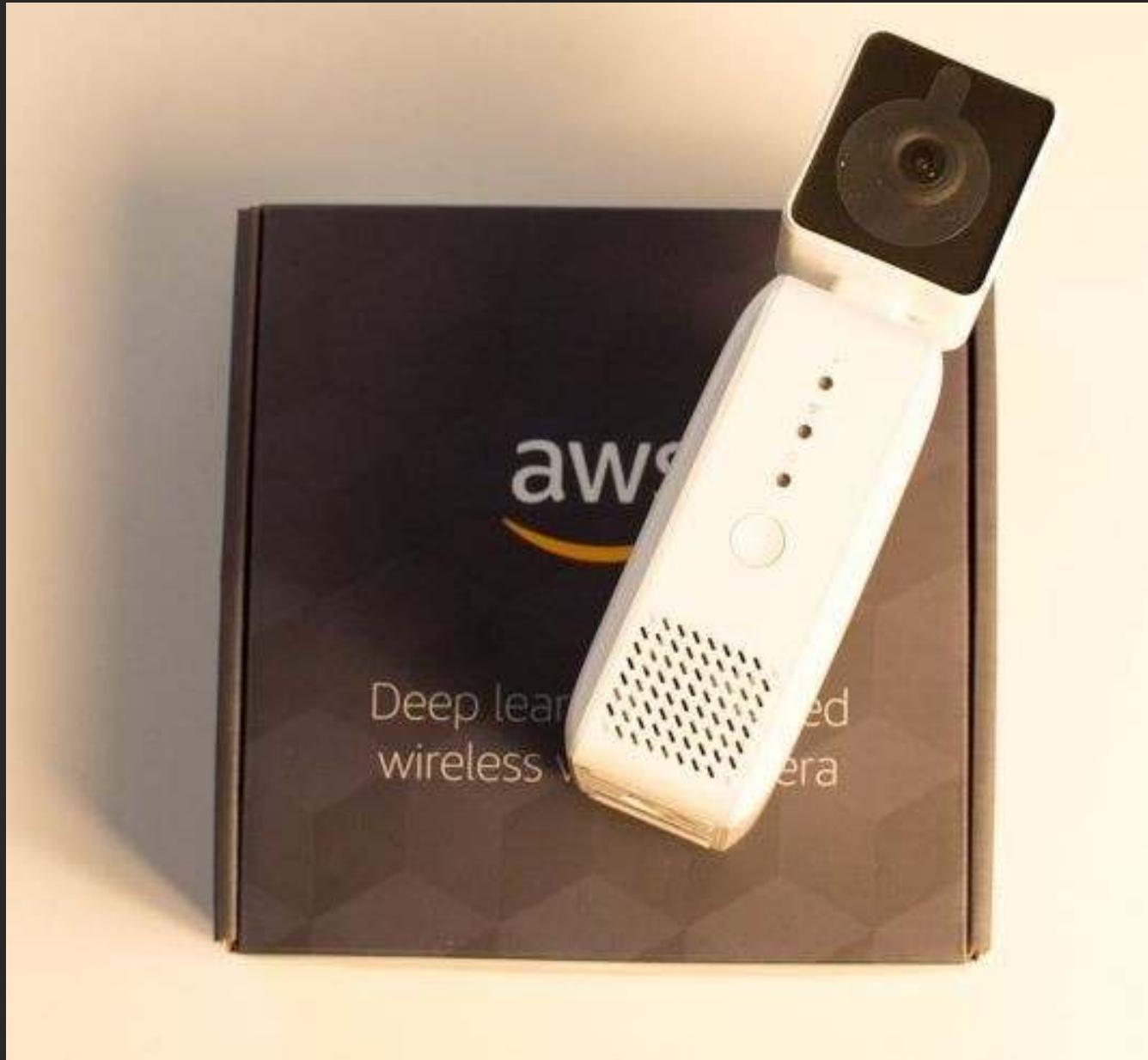
Transferring to the cloud HD video from multiple cameras
is too expensive to support real-time Computer Vision services



7.5 million vehicles,
1 camera in each vehicle,
1 minute of HD video from each vehicle

The cost of uploading is...

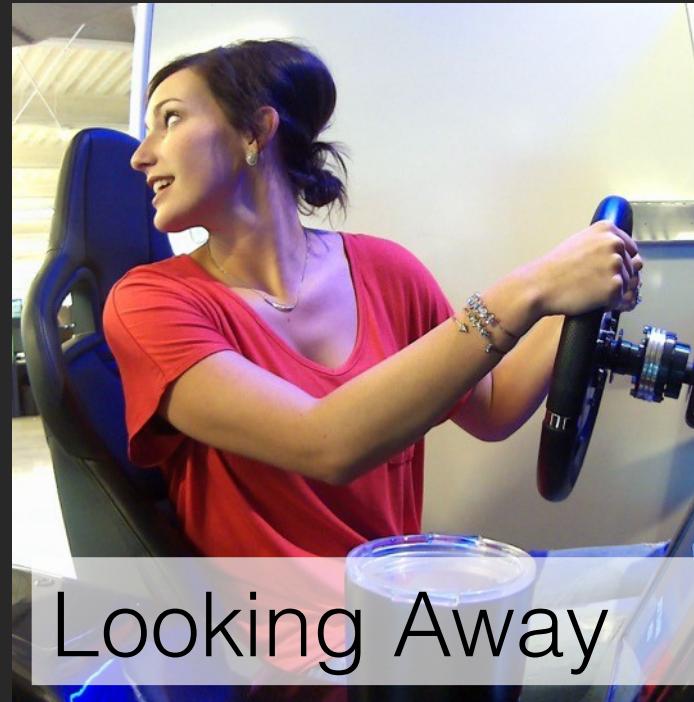
Billions of dollars per year



Innovation



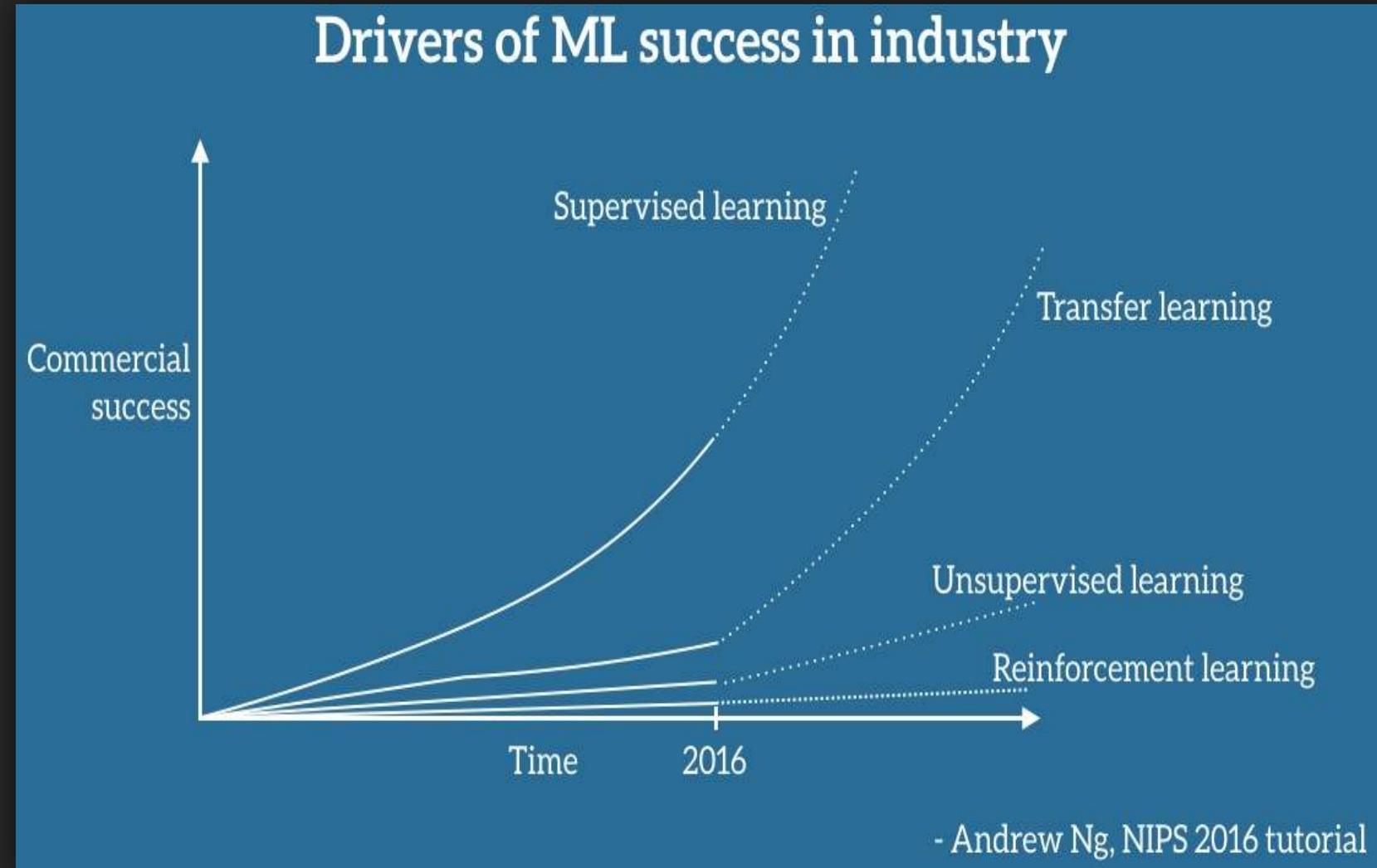
Can we predict driver distraction?



Training dataset

- Multiple drivers performing 4 complex actions
- Labelled as: Normal Driving, Texting, Looking Away, On the Phone

Transfer Learning



“Don’t be a Hero”

Instead of rolling your own network architecture, look at whatever currently works best on ImageNet, download a **pretrained model** and **finetune** it on your data.

You should **rarely ever** have to design a model from scratch or even train from scratch or.

ResNet-152

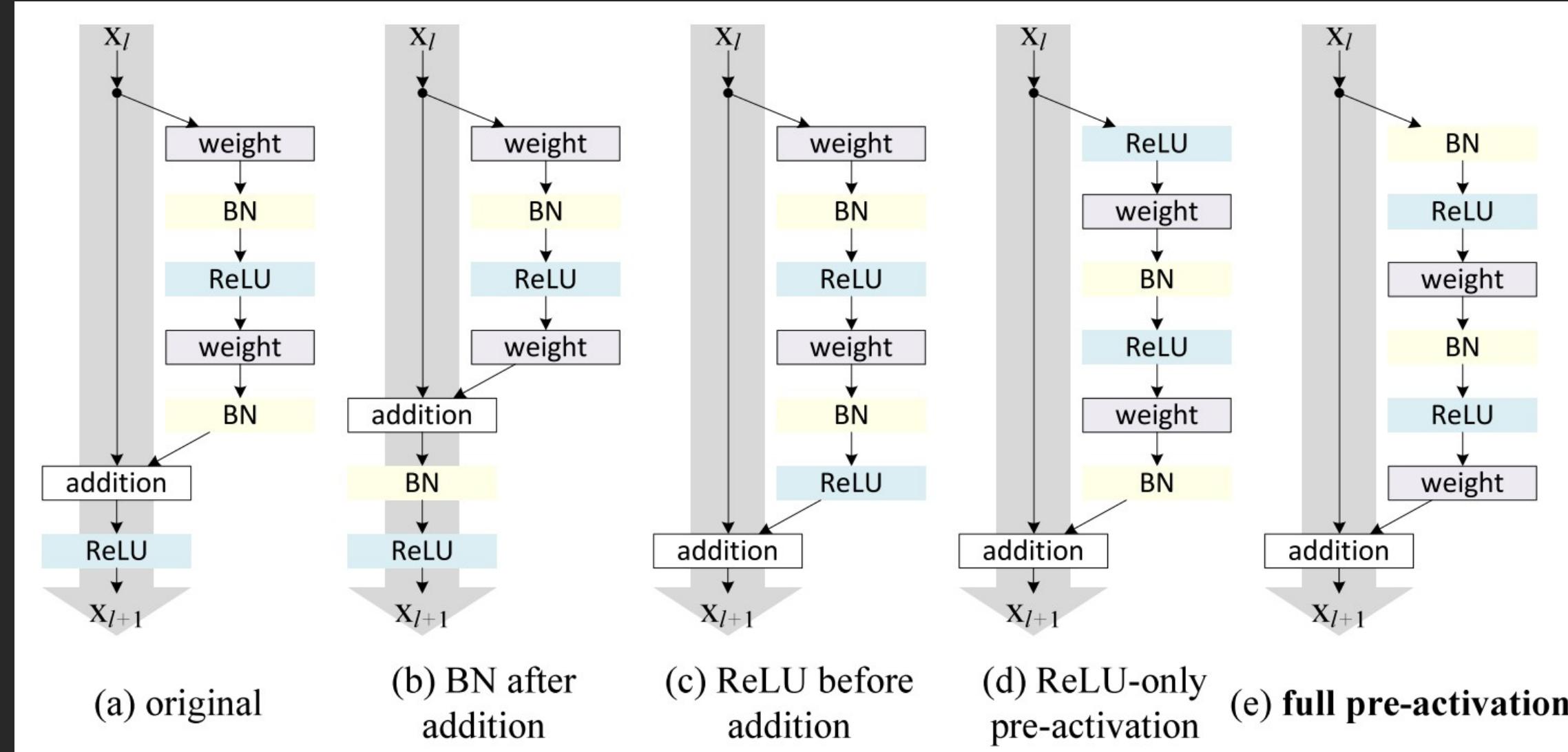
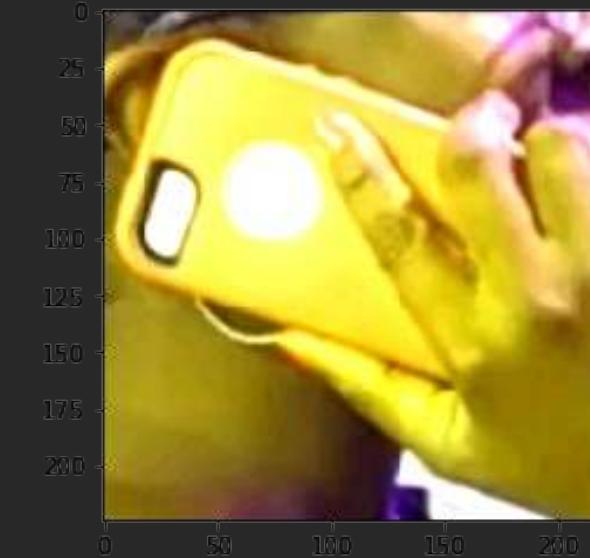


Image Augmentation

We performed random rotations, random crop, RGB mean centering and translations on each image sample.

The augmented dataset makes the resulting model more robust and less prone to overfitting.



Pre-processing and hyper parameter tuning

Optimized Hyper Parameters

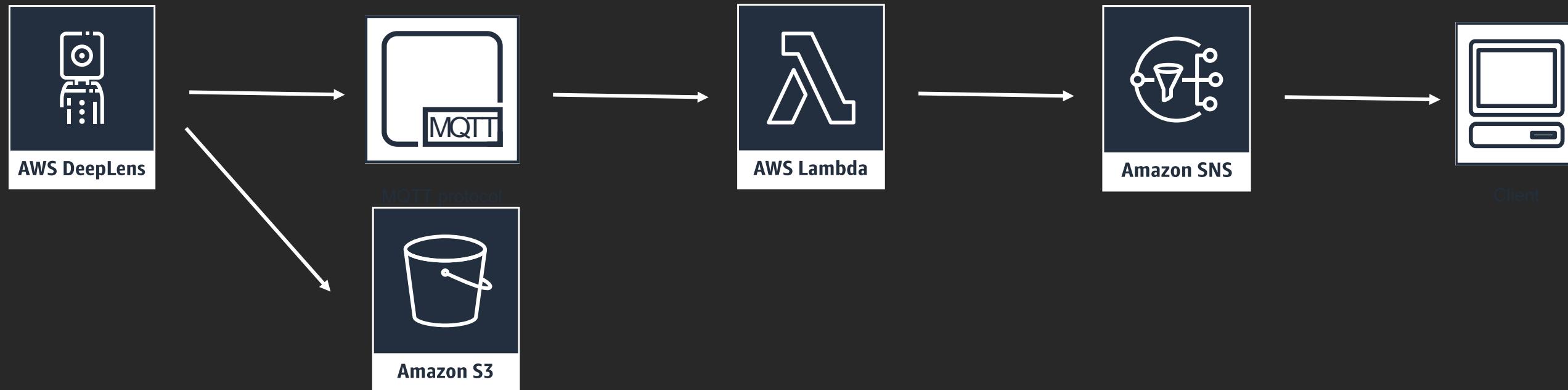
- Epochs : 50
- Image Shape : 3,224,224
- Learning_rate : 0.001
- Lr_scheduler_factor : 0.0001
- Lr_scheduler_step : 5
- Mini_batch_size : 64
- Num_Layers : 152
- Num_training_samples : 16180
- Optimizer : sgd

Data preprocessing : Down sample the images from 640x480 to 224x224 so that their dimension matches the one used by the pre-trained model.

Fine-tune pretrained network : after multiple iterations on training sample data to find hyper parameters, train using the pre-learned weights and our image data.

Deep Learning at the Edge

- Train on [Amazon Sagemaker](#): Amazon EC2 P2 instances, 8 NVIDIA Tesla K80 GPUs
- Deploy model artifacts on [AWS DeepLens](#) with [AWS Greengrass](#)
- Use the [Intel Model Optimizer](#) for faster GPU inferences
- Upload predicted images to [Amazon S3](#) in real-time
- Send inference results to [AWS IoT](#)



Video

Our volunteers



Dany Benjamin



Tori Salido



Steve King



Blake Wong



Trish Batson



Will Dombrowski



Toshi Tanaka



Michelle Haught

Lessons Learnt

Use of other state-of-the-art models at the Edge?

Ensemble is a very powerful tool but performance on constrained devices is a challenge.

Image Preprocessing is the key.

Focus on the part of the image that contains the driver and cut out the useless part of the image.

Clipping does not work.

Clipping the loss with various threshold does not help. The reason is that the model has different level of confidence on different classes .

How often you run the inference and how you interpret it is more important than model metrics like accuracy.

Optimizing performance for inference at the edge

- Constrained devices require **smaller** and **faster** models
 - Reduce **memory** usage
 - Reduce **computational complexity**
 - Reduce **power consumption**
- **Code optimization:** Intel Inference Engine, Intel MKL, NNPACK
- **Quantization / mixed mode training:** smaller weights and activations
- **Pruning:** remove useless connections
- **Compression:** encode weights
- Some of these techniques are available in Apache MXNet or TensorFlow

What if local inference is not possible?

Device can support HTTPS

Invoke an Amazon SageMaker endpoint

- Train a model in SageMaker (or import your own).
- Deploy it to a prediction endpoint.
- Device invokes the HTTPS endpoint and receives prediction results.

Best when	Requirements
Devices are not powerful enough for local inference.	Network is available.
Models can't be easily deployed to devices.	Devices support HTTP.
Additional cloud-based data is required for prediction.	
Prediction activity must be centralized.	

Device cannot support HTTPS

Invoke an AWS Lambda function through AWS IoT Core

- Train a model in Amazon SageMaker (or import your own).
- Deploy it to a prediction endpoint.
- Write an AWS Lambda function performing prediction.
- Device sends an MQTT message to AWS IoT Core, triggers the function and receives prediction results.

Best when	Requirements
Devices can support neither HTTP nor local inference (e.g. Arduino).	Network is available (MQTT is less demanding than HTTP).
Costs must be kept as low as possible.	Devices are provisioned in AWS IoT Core (certificate, keys).

Getting started

Resources

<https://ml.aws/>

<https://aws.amazon.com/sagemaker>

<https://aws.amazon.com/machine-learning/amis/>

<https://aws.amazon.com/greengrass>

<https://aws.amazon.com/greengrass/ml>

<https://aws.amazon.com/deeplens>

<https://medium.com/@julsimon>



Please complete the
session survey in the mobile
app.

Thank you!

Julien Simon
Principal Tech. Evangelist, AI/ML
Amazon Web Services
[@julsimon](https://twitter.com/julsimon)

Brian Kursar
Chief Data Scientist
Toyota Connected Data Services
Brian.Kursar@toyotaconnected.com
[@briankursar](https://twitter.com/briankursar)

Nimish Amlathe
Data Scientist
Toyota Connected Data Services
Nimish.Amlathe@toyotaconnected.com