

tinyML® Talks

Enabling Ultra-low Power Machine Learning at the Edge

“Running and Managing Fleets of Single Board Computers at Scale”

Seth Clark – Co-founder & Head of Product
Modzy

May 2, 2023



www.tinyML.org



Thank you, **tinyML Strategic Partners**,
for committing to take tinyML to the next Level, together



T I N Y



TALKS
webcast

Executive Strategic Partners

T I N Y



TALKS
webcast



EDGE IMPULSE

The Leading Development Platform for Edge ML

edgeimpulse.com

Qualcomm
AI research

Advancing AI research to make efficient AI ubiquitous

Power efficiency

Model design, compression, quantization, algorithms, efficient hardware, software tool

Personalization

Continuous learning, contextual, always-on, privacy-preserved, distributed learning

Efficient learning

Robust learning through minimal data, unsupervised learning, on-device learning

A platform to scale AI across the industry



Perception

Object detection, speech recognition, contextual fusion



Reasoning

Scene understanding, language understanding, behavior prediction



Action

Reinforcement learning for decision making



Edge cloud



Cloud



IoT/IIoT



Automotive



Mobile



Accelerate Your Edge Compute

SYNTIANT

Making Edge AI A Reality

www.syntiant.com

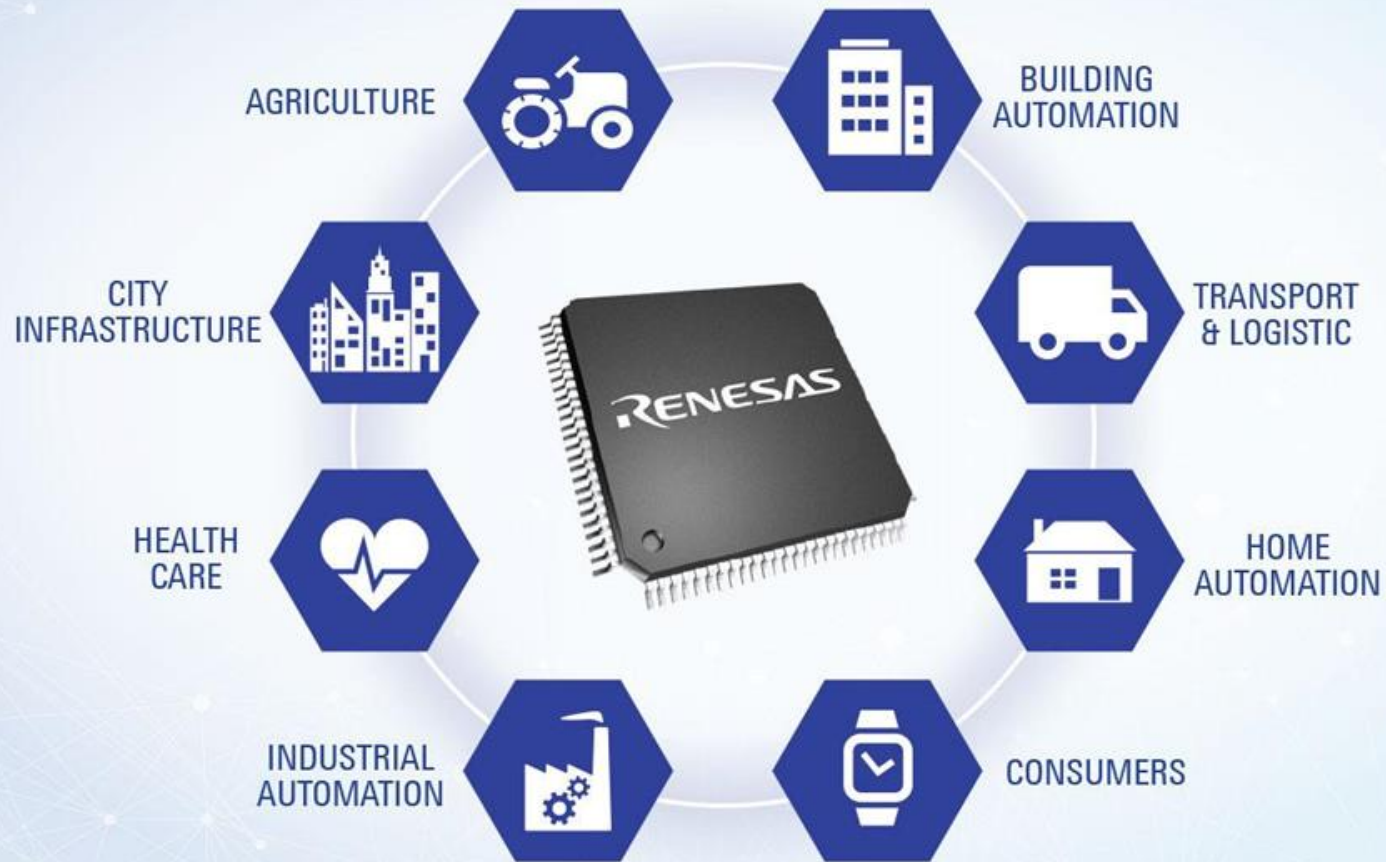
T I N Y



TALKS
webcast

Platinum Strategic Partners

Renesas is enabling the next generation of AI-powered solutions that will revolutionize every industry sector.



[renesas.com](https://www.renesas.com)



**DEPLOY VISION AI
AT THE EDGE AT SCALE**

SONY

Gold Strategic Partners



Where what if
becomes what is.

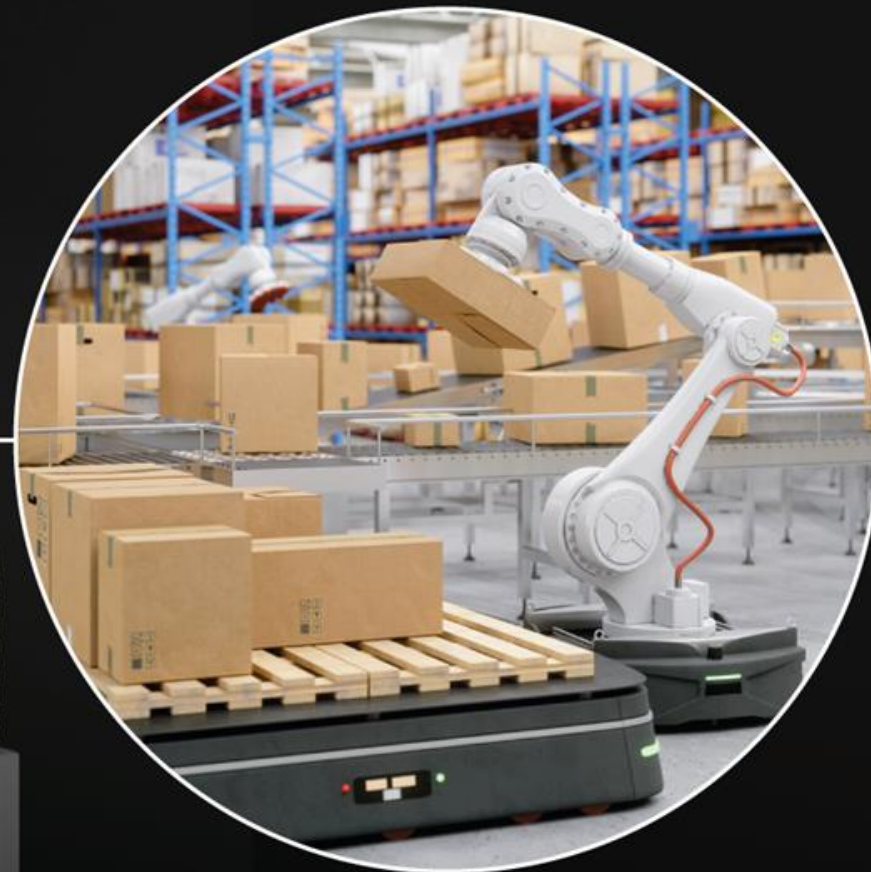
Witness potential made possible at analog.com.



PRO™

Easily deploy your
tinyML solutions with
Arduino Pro

arduino.cc/pro



Made In Italy

arm AI



Powering tinyML Innovation

Arm AI Virtual Tech Talks

The latest in AI trends, technologies & best practices from Arm and our Ecosystem Partners.

Demos, code examples, workshops, panel sessions and much more!

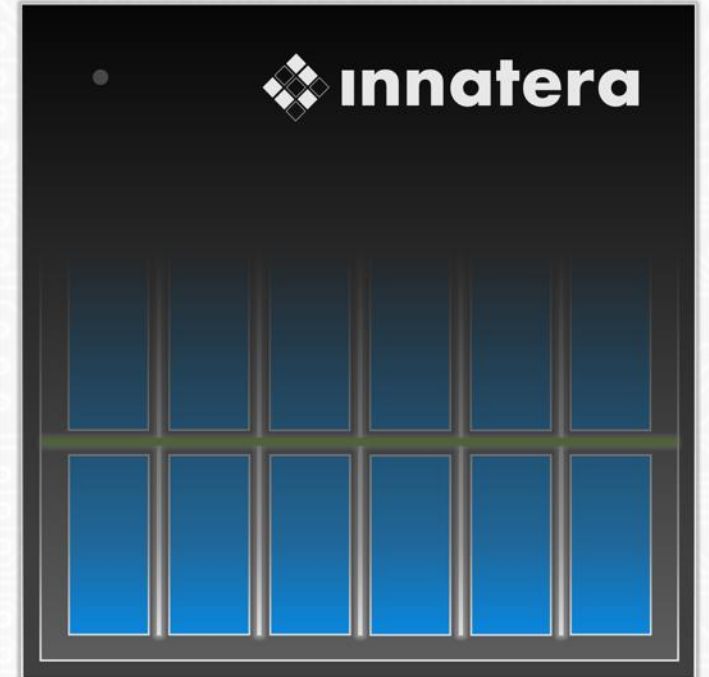
Fortnightly Tuesday @ 4pm GMT/8am PT

Find out more:

www.arm.com/techtalks



NEUROMORPHIC INTELLIGENCE FOR THE SENSOR-EDGE





Microsoft

The Right Edge AI Tools Can Make or Break Your Next Smart IoT Product



Analytics Toolkit Suite





life.augmented

STMicroelectronics provides extensive solutions to make tiny Machine Learning easy



ENGINEERING EXCEPTIONAL EXPERIENCES

We engineer exceptional experiences for consumers in the home, at work, in the car, or on the go.

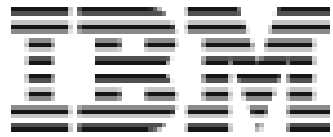
www.synaptics.com



T I N Y



Silver Strategic Partners





Join Growing tinyML Communities:



14.4k members in
47 Groups in 39 Countries

tinyML - Enabling ultra-low Power ML at the Edge

<https://www.meetup.com/tinyML-Enabling-ultra-low-Power-ML-at-the-Edge/>



4k members
&
11.6k followers

The tinyML Community

<https://www.linkedin.com/groups/13694488/>





Subscribe to
tinyML YouTube Channel
for updates and notifications
(including this video)

www.youtube.com/tinyML



tinyML
4.33K subscribers

9.2k subscribers, 555 videos with 322k views

HOME VIDEOS PLAYLISTS COMMUNITY CHANNELS ABOUT

 13:24 On Device Learning Forum - Professors... 106 views · 4 days ago	 33:27 On Device Learning - Manuel Roveri: Is on-... 138 views · 4 days ago	 32:39 On Device Learning Forum - Warren Gros... 54 views · 4 days ago	 36:41 On Device Learning Forum - Yiran Chen: ... 47 views · 4 days ago	 34:03 On Device Learning Forum - Hiroku... 132 views · 4 days ago	 34:58 On Device Learning Forum - Song Han: O... 137 views · 4 days ago
 1:13 tinyML Smart Weather Station Challenge - ... 122 views · 4 days ago	 1:07:43 tinyML Talks Singapore... 262 views · 2 weeks ago	 53:41 tinyML Talks Shenzhen: Data... 511 views · 3 weeks ago	 45:46 tinyML Talks Singapore... 229 views · 3 weeks ago	 51:01 tinyML Smart Weather Station with Syntiant... 265 views · 3 weeks ago	 1:03:24 tinyML Trailblazers August with Vijay... 286 views · 1 month ago
 58:50 tinyML Auto ML Tutorial with SensiML 351 views · 1 month ago	 34:36 tinyML Auto ML Tutorial with Qeexo 462 views · 2 months ago	 55:01 tinyML Talks Germany: Neural network... 374 views · 2 months ago	 59:51 tinyML Trailblazers with Yoram Zylberberg 133 views · 2 months ago	 59:48 tinyML Auto ML Tutorial with Nota AI 287 views · 2 months ago	 58:09 tinyML Auto ML Tutorial with Neuton 336 views · 2 months ago
 1:02:30 tinyML Challenge 2022: Smart weather... 378 views · 2 months ago	 34:31 tinyML Talks South Africa - What is... 214 views · 2 months ago	 1:00:30 tinyML Talks: The new Neuromorphic Anal... 448 views · 2 months ago	 1:06:44 tinyML Talks Shenzhen: 分享主题... 159 views · 2 months ago	 1:53:07 tinyML Auto ML Forum - Paneldiscussion 190 views · 2 months ago	 42:13 tinyML Auto ML Forum - Demos 545 views · 2 months ago



EMEA 2023

<https://www.tinyml.org/event/emea-2023>

More sponsorships are available: sponsorships@tinyML.org

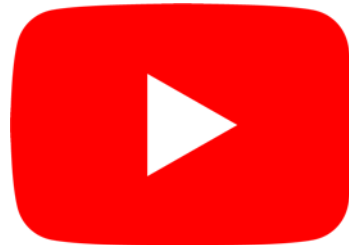


Reminders

Slides & Videos will be posted tomorrow



tinyml.org/forums



youtube.com/tinyml



Please use the Q&A window for your questions





Seth Clark

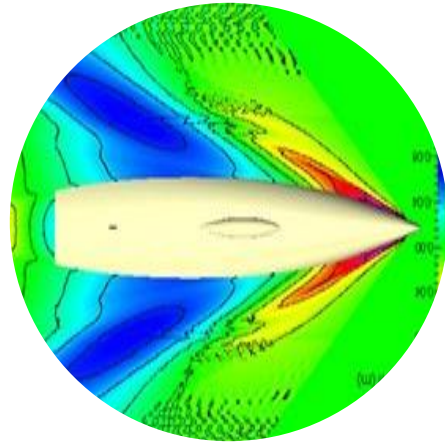
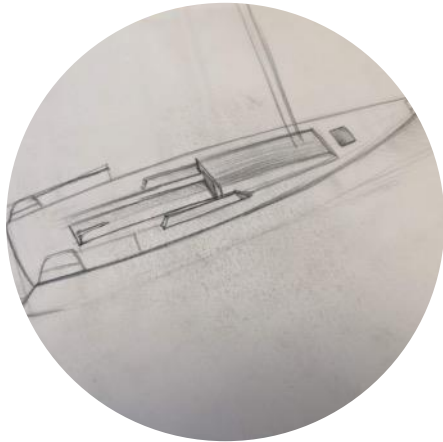


Seth Clark is Head of Product & Co-founder at Modzy. Prior to founding Modzy, Seth served as product manager for a number of successful analytics products. He also served as Principal at Booz Allen Hamilton, where he led complex projects at the intersection of Data Science and Software Development. He has degrees in engineering from the University of Southampton and the Massachusetts Institute of Technology.

Running and Managing Fleets of Single Board Computers (SBCs) at Scale

tinyML Talks

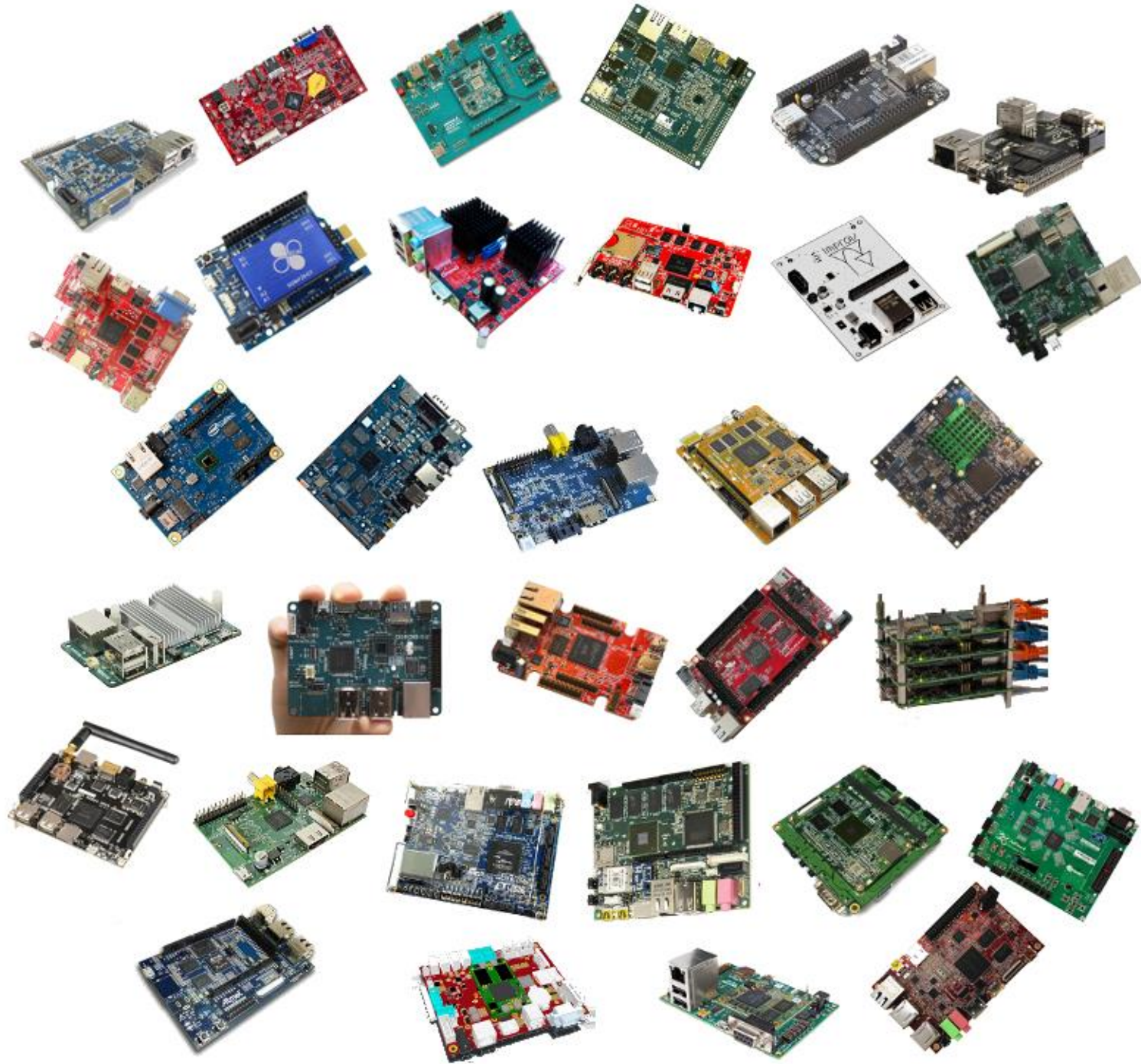
2 May 2023



How I got here



2020s: The Wild West Golden Age of SBCs



SBCs are doing amazing things



ZeroPhone

Pi Zero-based open-source mobile phone (that you can assemble for 50\$ in parts)



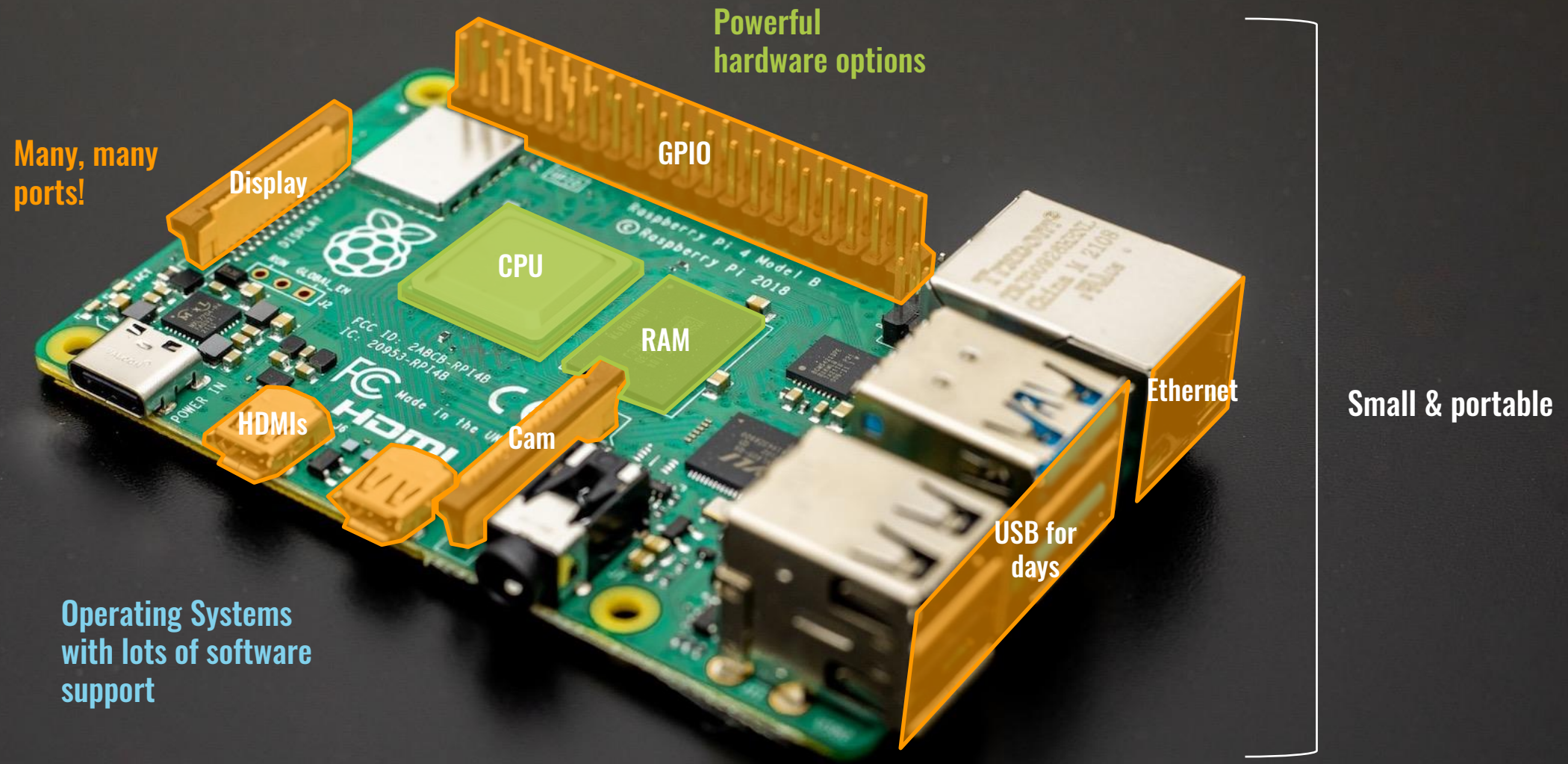
Jetson Clean Water AI

Using AI object detection to detect water contamination



ClippyGPT

Raspberry Pi 3B powered Clippy bot with ChatGPT



Why we love SBCs

SBCs are a convenient alternative to microcontrollers

SBCs	Microcontrollers
Operating System with an optional user interface	No Operating System
Built in networking	No networking
Lots of software options	Limited to software supported by microcontroller's IDE

Businesses are increasingly adopting SBCs for serious applications

- » 3D cameras
- » 5G modules
- » Robotics
- » Smart cities
- » Smart manufacturing
- » IoT
- » Fleet management
- » ...



**\$2.94 billion
market in 2022**

SBCs are also great for machine learning

- » Write, compile, and run custom software in any language you like
- » Download ML frameworks, libraries, and even fully trained models
- » Plug in peripherals like cameras and sensors for computer vision and anomaly detection
- » Use built in hardware accelerators (e.g. GPUs) or plug one in via USB
- » Add an HDMI cable and a monitor and you've got a fully-functional app

SBCs are great, and so is machine learning, so...



?

Setting up an SBC takes time and effort



Running ML apps brings distinctive challenges

- » Minimizing resource needs of models
- » Running on multiple chipsets (arm, x86/AMD, 32-bit & 64-bit)
- » Managing code for pre-processing data & post-processing predictions
- » Working with hardware accelerators (GPUs, TPUs, and FPGAs)
- » Connecting to sensors and local data sources
- » Gracefully handling intermittent network availability
- » Monitoring model performance
- » Improving and updating models overtime

Running models by hand on 1000s of SBCs would be madness, so what if you could...

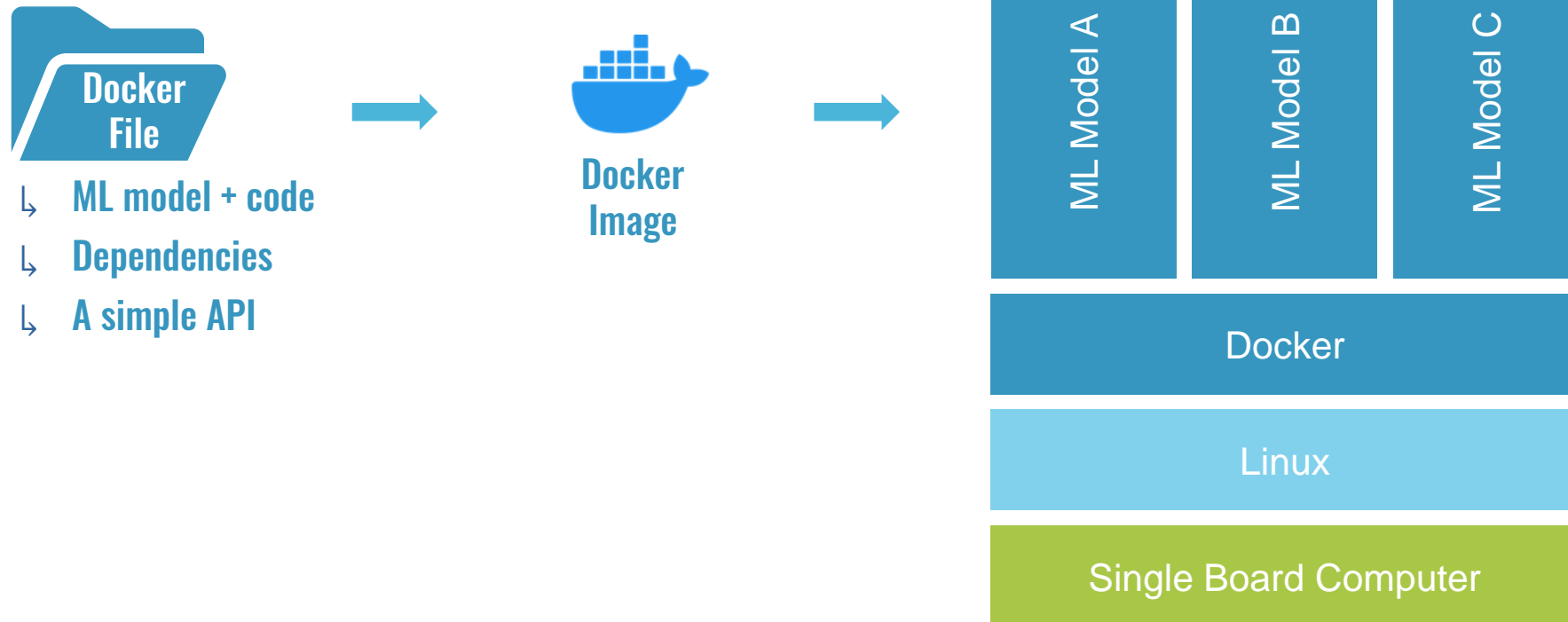
- » Send model code and dependencies to an SBC dependably, repeatably, and scalably
- » Connect to local data sources and other applications
- » Run quickly and use as few resources as possible
- » Monitor models while they're running
- » Will keep working even when the wifi is down



Enter the container

Matson LEASING COMPANY
MLCU 249978 3 2210
IPYU 214850 6 2261
MSC MEDU 163664 3 2261
FCIU 604407 2 2261
GOLD
GLDU 526446 8 2261
MSC MEDU 292017 1 2261
MSC GLDU 395 226
M.G.W. TARE NET CU.CAP

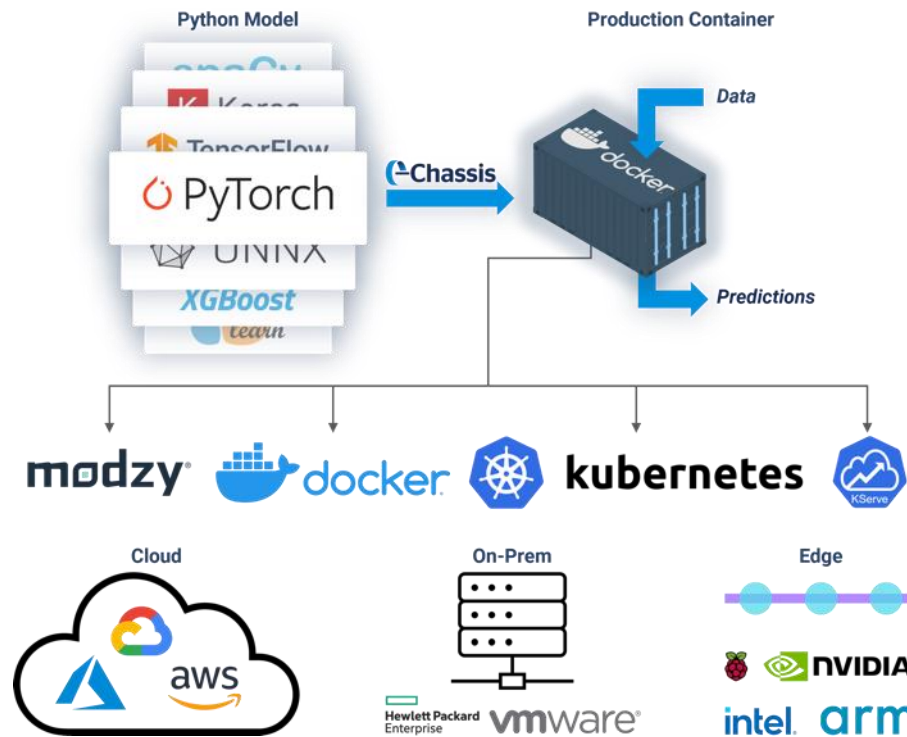
The basic idea



A hypothetical ML container image

API	gRPC APIs: Status, Run, Shutdown			
Inference script	inference.py			
Trained model	model-weights.pkl			
Dependencies	pickle	cv2	torch	numpy
Base Image	armv6 Debian OS			

Chassis.ml: An easy way to build model containers

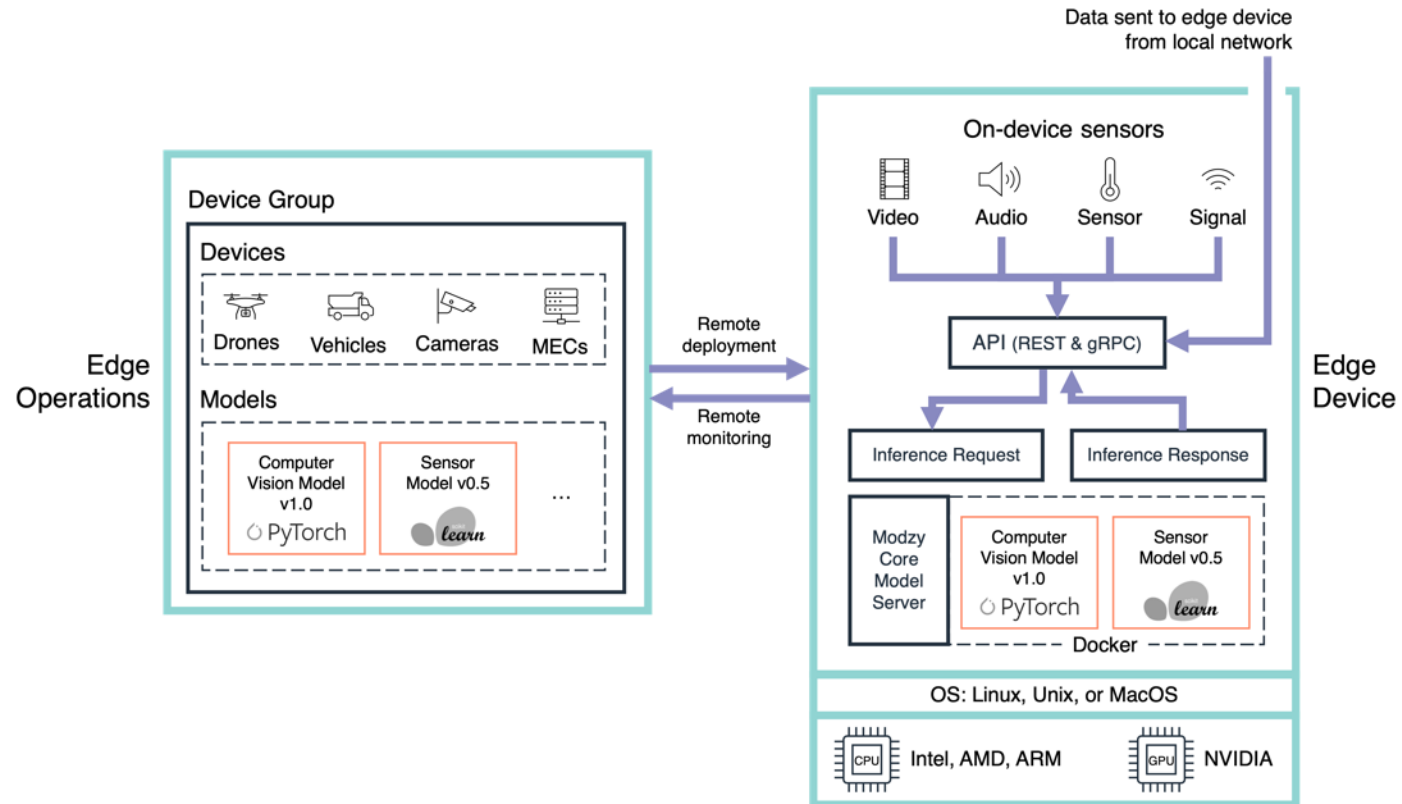


Website: <https://chassis.ml>
 Github: [modzy/chassis](https://github.com/modzy/chassis)

A few more ingredients for managing containers on lots of SBCs

- » Chassis.ml containers
- » Dockerhub
- » Docker
- » Raspberry Pi, Jetson Nano, Intel Upboard
- » Modzy Core
- » Nats.io + Jetstream

How it works



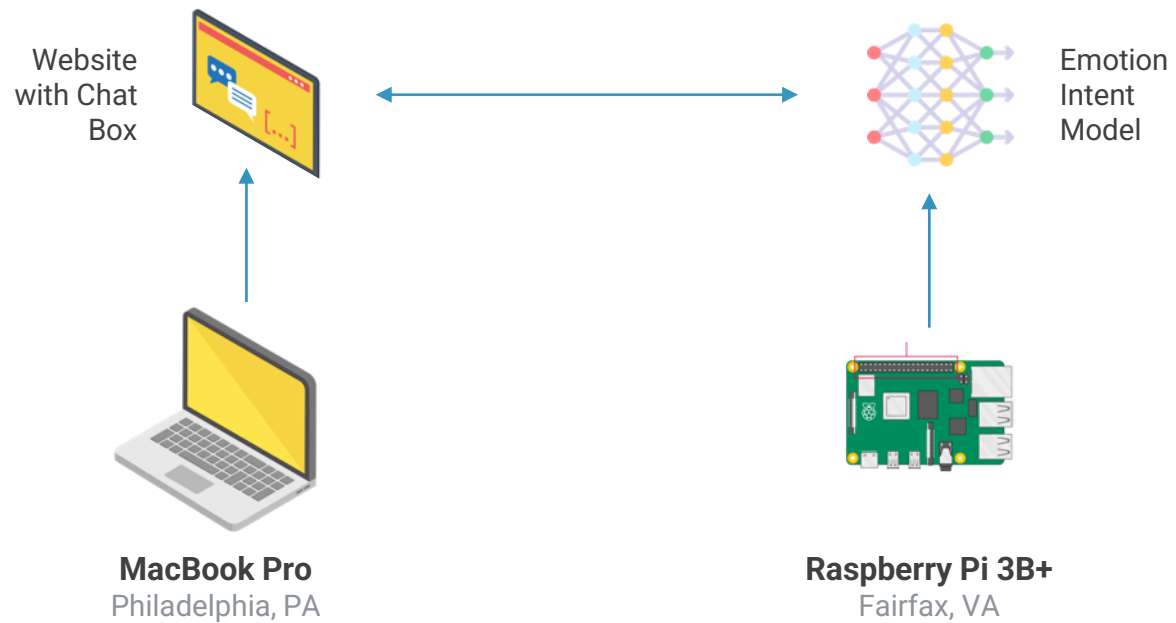
Benefits of container-based ML on SBCs

- » Supports both ARM and x86/AMD chipsets, 32 or 64-bit
- » Built-in support for GPUs
- » Models never break due to missing dependencies
- » Network connectivity is only needed initially to deploy models
- » Models can be monitored at any time
- » New models and versions can be deployed centrally to many, many devices at once

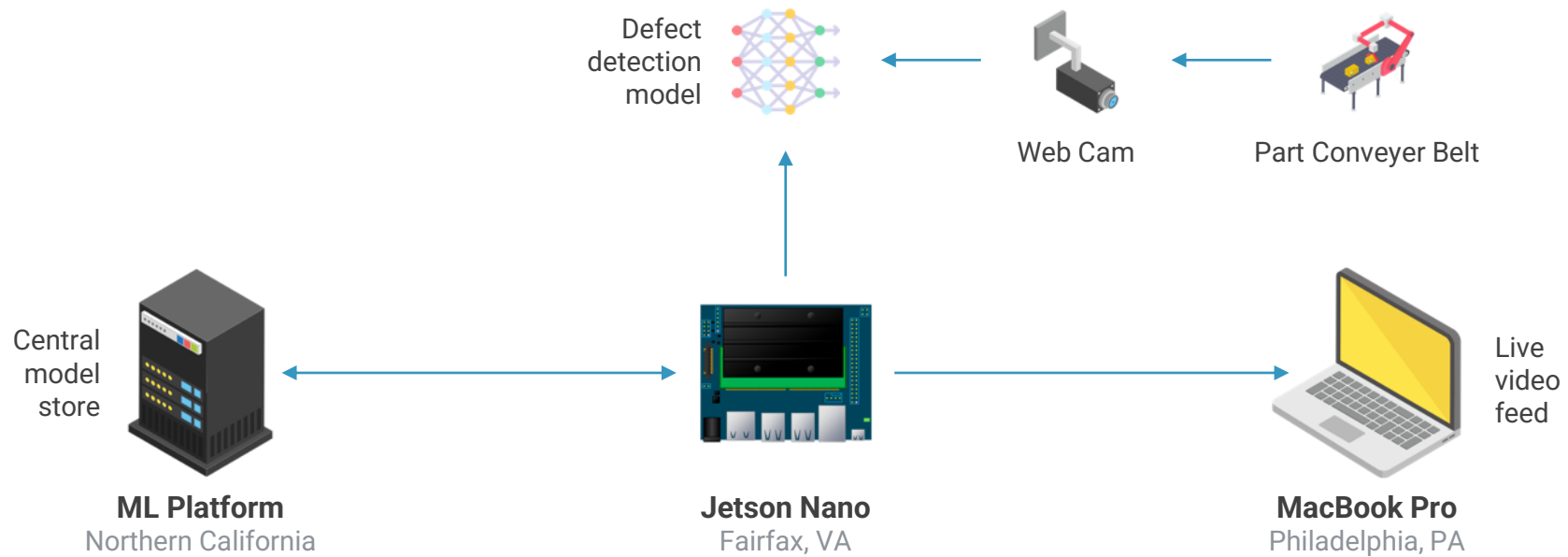
Examples



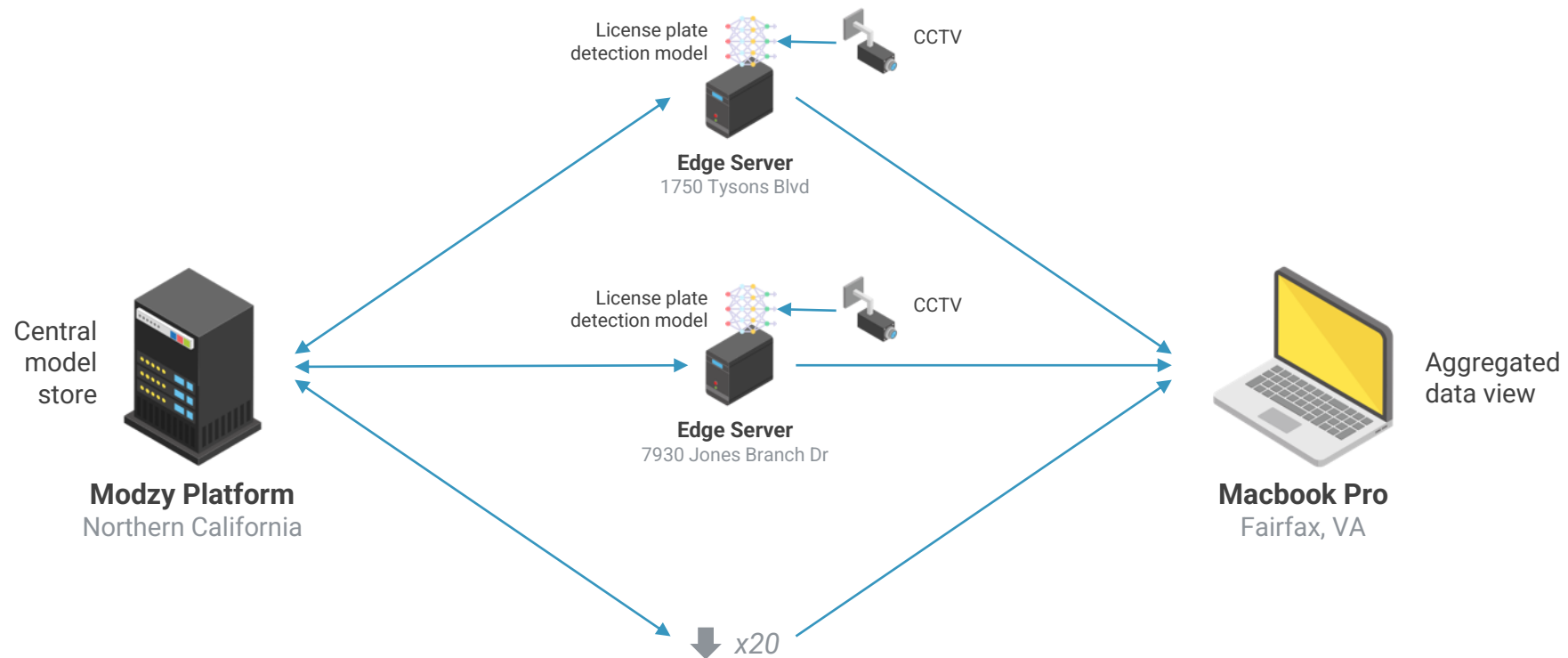
Example: Edge-centric NLP App



Example: Industrial Defect Detection



Example: Computer Vision Parking Lots



Final thoughts

- » Use an edge-centric architecture anytime you need low-latency, non-stop performance on or offline, or models that need to run in more than one place
- » Start small, but plan for massive scale

modzy[®]

info@modzy.com

 [@getModzy](https://twitter.com/getModzy)

 [getModzy](https://www.linkedin.com/company/getModzy)

 [discord.gg](https://discord.gg/modzy)



Seth Clark

Thanks for joining our webinar!
Scan the QR code to join our Discord server for
more resources.





Copyright Notice

This multimedia file is copyright © 2023 by tinyML Foundation. All rights reserved. It may not be duplicated or distributed in any form without prior written approval.

tinyML[®] is a registered trademark of the tinyML Foundation.

www.tinyml.org



Copyright Notice

This presentation in this publication was presented as a tinyML® Talks webcast. The content reflects the opinion of the author(s) and their respective companies. The inclusion of presentations in this publication does not constitute an endorsement by tinyML Foundation or the sponsors.

There is no copyright protection claimed by this publication. However, each presentation is the work of the authors and their respective companies and may contain copyrighted material. As such, it is strongly encouraged that any use reflect proper acknowledgement to the appropriate source. Any questions regarding the use of any materials presented should be directed to the author(s) or their companies.

tinyML is a registered trademark of the tinyML Foundation.

www.tinyml.org