

The logo for the 2024 Embedded VISION Summit is centered on the left side of the slide. It features a white octagonal background with a colorful, multi-layered border in shades of purple, blue, green, yellow, and orange. The text "2024" is at the top, "embedded" is below it, "VISION" is in large, bold, dark blue letters with a gradient, and "SUMMIT" is at the bottom.

2024
embedded
VISION
SUMMIT®

Real-Time Retail Product Classification on Android Devices inside the Caper AI Cart

David Scott

Senior Machine Learning Engineer

Instacart (Caper)



Table of Contents



Instacart and Caper



Current challenges for such products



The novel proposed approach



Approach to reduce resource usage



Approach to integrate with Android system



Performance results

Caper @ Instacart



A Seamless Shopping Experience

Caper Cart

1

Sensor Fusion & AI Integrated
Providing customers with a more convenient way to shop

2

Weights & Measures Enabled
Robust and certified

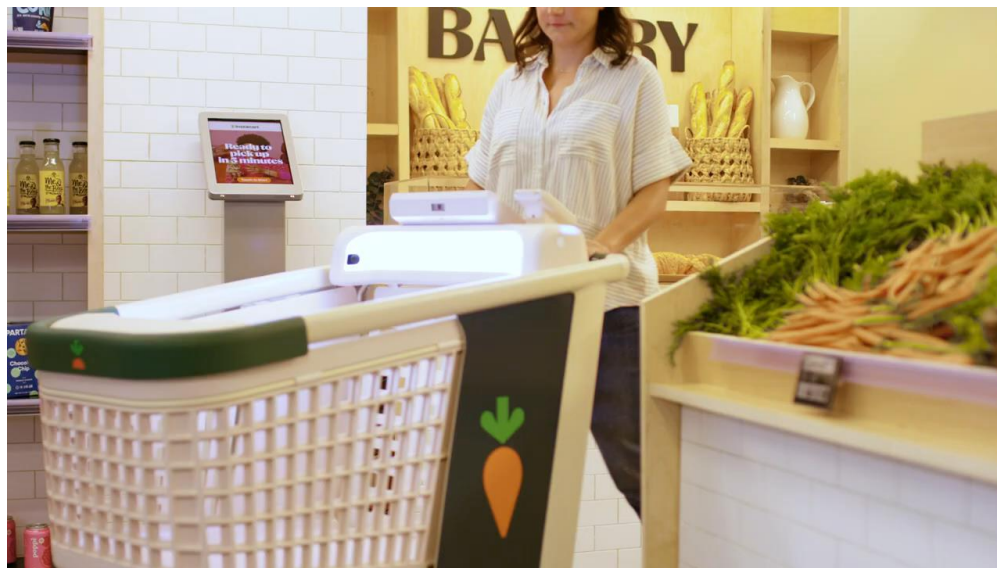
3

Embedded Location Systems
Streamlining the shopper experience



Computer Vision and AI @ Caper

In order to unlock a magical shopping experience, we enable seamless AI that detects and recognizes products in less than 0.5 seconds



Current Challenges for Such Products



1 High-throughput model with high accuracy

Ensure smooth and efficient detection

Avoid delays or disruptions in user experience

2 Limited by system constraints

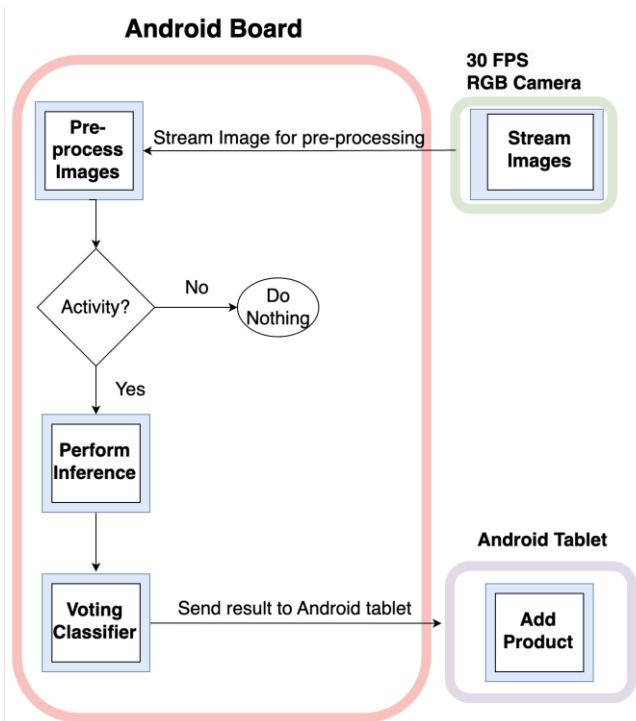
Limited CPU resources shared by Android services

Model speed can impact overall app degradation

3 User experience

How can we utilize the hardware + software resources we have to improve our user experience without requiring re-educating our customer?

Proposed Approach



- 1 Generate image stream and process directly on Android board**
Utilize 30 FPS red-green-blue (RGB) camera to stream images for processing
- 2 Skip processing images on Jetson GPU**
Utilize the processing power of our Android board to stream and process images directly on Android
- 3 Deploy using TensorFlow Lite (TFLite)**
Train and convert our PyTorch models to TFLite for model inference

Proposed Approach



- 1 **User adds product under our camera**
Get a stream of images that are background removed
- 2 **Perform inference on each image**
Android board takes in images as an input and returns embeddings
- 3 **Utilize a voting classifier with embeddings**
 - Generate one set of embeddings per frame from a single network
 - Utilize ground truth embeddings to calculate highest similarity in class

Approach to Reduce Resource Usage

Approach to Reduce Resource Usage: Camera

Frame skipping

Do we have an accuracy reduction or lose important information by skipping frames?

Re-using circular buffers

Other services utilize the same camera; we can re-use buffers



Reducing the resolution of incoming image stream

Since our model utilizes 224 [w, h] images, is it possible to modify the camera firmware to give us smaller images?

Approach to Reduce Resource Usage: System

Threading

How many threads is necessary to maximize the performance of the system?

Perfetto traces

Utilized to help pinpoint heavy system resource usages on CPU and RAM



Reducing embeddings dimensions

By decreasing the dimensions and # of our output embeddings, we reduced system usage

Approach to Integrate with Android System

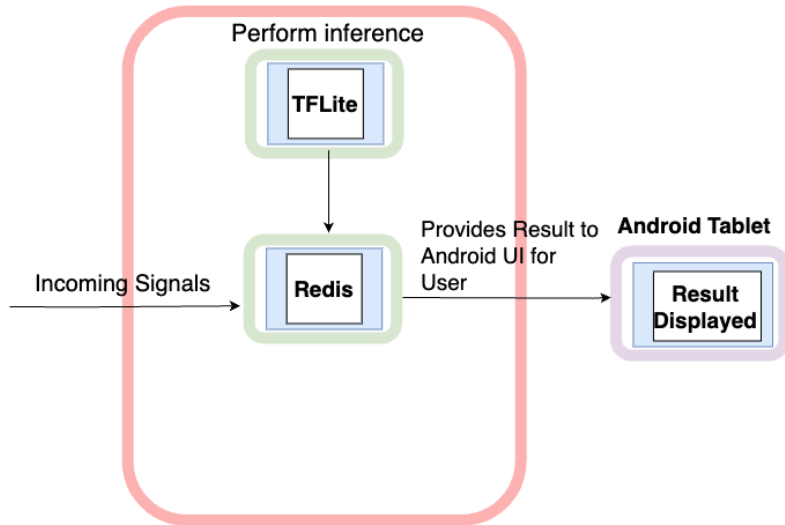


Android



Approach to Integrate with Android System

Android Board



- 1 **Utilizing co-routines and state machine learning model**
Break thread if issue arises, avoiding impact on other apps
- 2 **Incoming signals via Redis**
Combine information across services to give best result
- 3 **TFLite integration**
 - Feed model and perform embedding similarity
 - Utilize TFLite built in optimizations for on-device inference

Results and Next Steps

CPU usage

Get the best performance while minimizing CPU usage for other Apps



Model throughput

How can we optimize model throughput while maintaining precision and recall?

Camera speed

Higher FPS gets us better recognition
More FPS = More CPU;
how to balance?

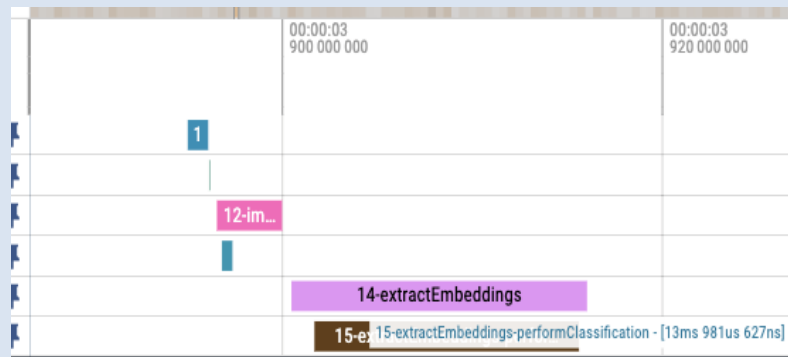


Model throughput

We optimized our throughput by the following:

- Using INT8 instead of FLOAT32
- Utilizing 112x112 to train model instead of 224x224
- Decreasing output embedding size

Increase in speed of ~96% from 400ms!



Android Performance Metrics — Camera FPS

Camera FPS

We optimized our image resolution + camera FPS by

- Using 640P images @ 30 FPS
- Re-using existing circular buffers from another service
- Worked with camera vendor to get custom firmware

Throughput + resolution investigation

```
sorted supported format(MJPEG) w*h(320x240)@30.000000 fps
sorted supported format(MJPEG) w*h(352x288)@30.000000 fps
sorted supported format(YUYV) w*h(640x480)@30.000000 fps
sorted supported format(MJPEG) w*h(800x600)@30.000000 fps
sorted supported format(MJPEG) w*h(1280x720)@30.000000 fps
sorted supported format(MJPEG) w*h(1280x960)@30.000000 fps
sorted supported format(MJPEG) w*h(1920x1080)@30.000000 fps
sorted supported format(MJPEG) w*h(1600x1200)@15.000000 fps
sorted supported format(MJPEG) w*h(2048x1536)@15.000000 fps
sorted supported format(MJPEG) w*h(2592x1944)@15.000000 fps
sorted supported format(MJPEG) w*h(2592x1944)@15.000000 fps
```

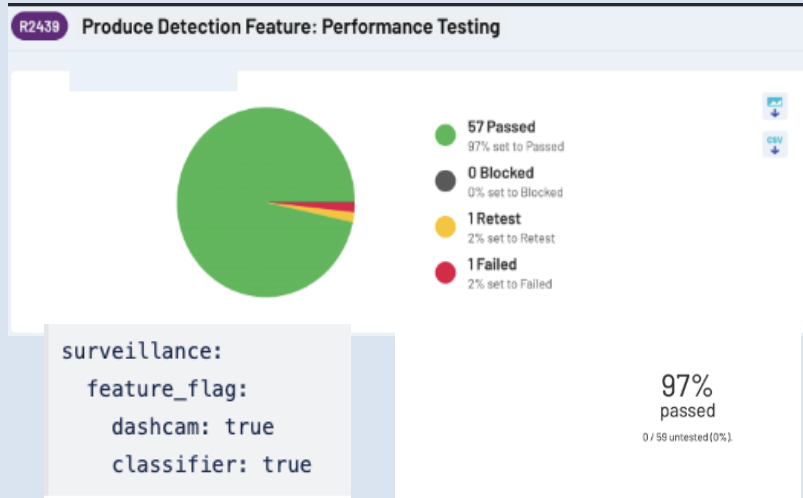
Android Performance Metrics — CPU Usage

CPU usage

We optimized CPU usage by optimizing our model & camera

- Skip intermediate steps, like image resizing
- Feed system with a queue of images
- Stress test across our system to ensure performance

Feature Testing from QA





1 Further optimizations

Continuous improvement of model accuracy and speed

2 Improving TFLite ecosystem and documentation

Contribute our learnings and findings to the sparse documentation for TFLite.

3 Utilizing research to enhance user experience

Provide a seamless and delightful experience for users



- 1 Solve your problems first with hardware**
Optimizing hardware selection is a great place to start before software
- 2 Optimizing throughput leads to better overall performance**
Gains in throughput speed helped decrease CPU usage (through frame skipping)
- 3 Deeply investigating pre-processing and model steps to maximize throughput**
Spending time investigating pre-processing and model optimization was fruitful for the team

TensorFlow Lite Documentation:

- [TensorFlow Lite Model Optimization Documentation](#)
- [TensorFlow Lite Quantization Documentation](#)
- [TensorFlow Lite Model Analyzer Documentation](#)

Android Development for TensorFlow Lite:

- [Quickstart for Android Development](#)

The logo for the 2024 Embedded VISION Summit is centered on the left side of the slide. It features a white octagonal background with the text "2024 embedded VISION SUMMIT" inside. The word "VISION" is in a large, bold, blue font with a gradient effect, while "2024", "embedded", and "SUMMIT" are in a smaller, black, sans-serif font. The octagon is surrounded by a colorful, geometric border of overlapping triangles in shades of purple, blue, green, yellow, and orange.

2024
embedded
VISION
SUMMIT®

Thank you for your time!
Look forward to answering your questions

David Scott

Senior Machine Learning Engineer

Instacart (Caper)

<https://www.linkedin.com/in/davejscott/>

