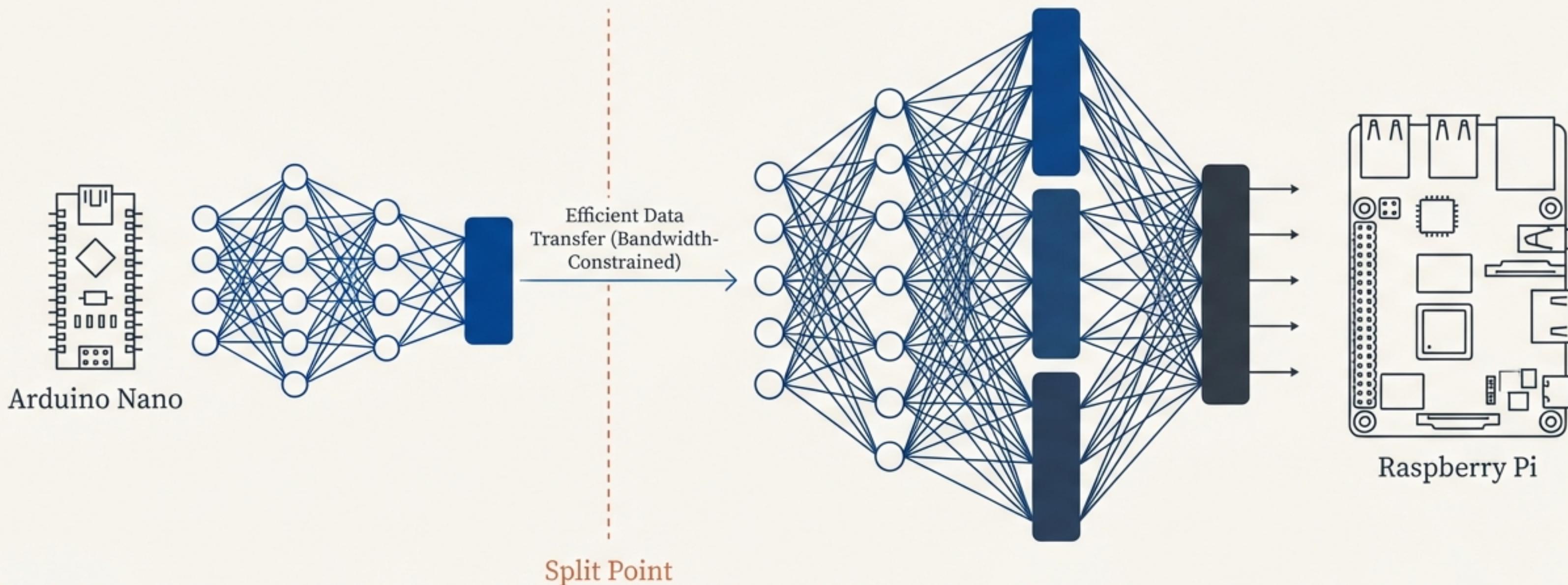


# Hardware-Aware Hybrid Intelligence

A Case Study in Bandwidth-Constrained Split Learning Across an Arduino Nano and Raspberry Pi



# The AI Edge Dilemma: Powerful Models vs. Tiny Devices

## The Goal

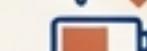
Deploy sophisticated deep learning models on resource-constrained microcontrollers like the Arduino Nano 33 BLE Sense.

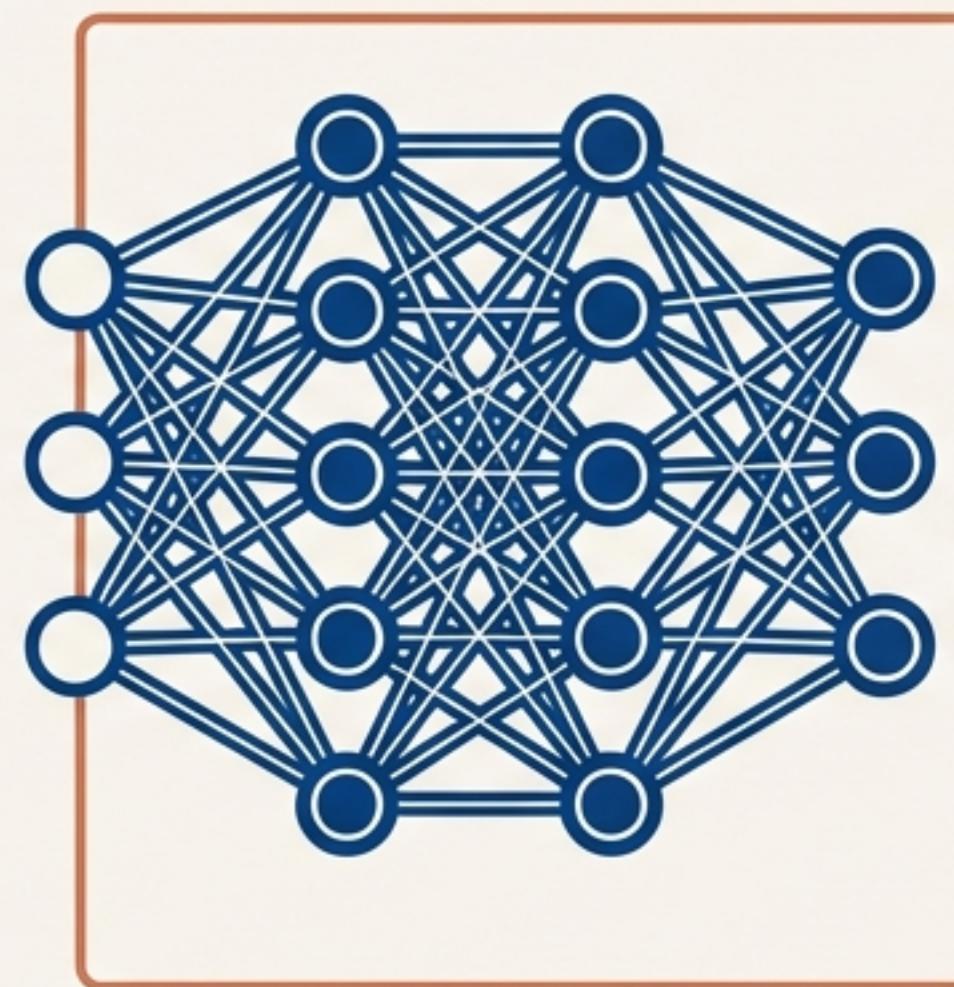
## The Wall

These devices suffer from a “utility gap”—their strict memory (<1MB SRAM/Flash) and processing limits make running complex models locally infeasible.

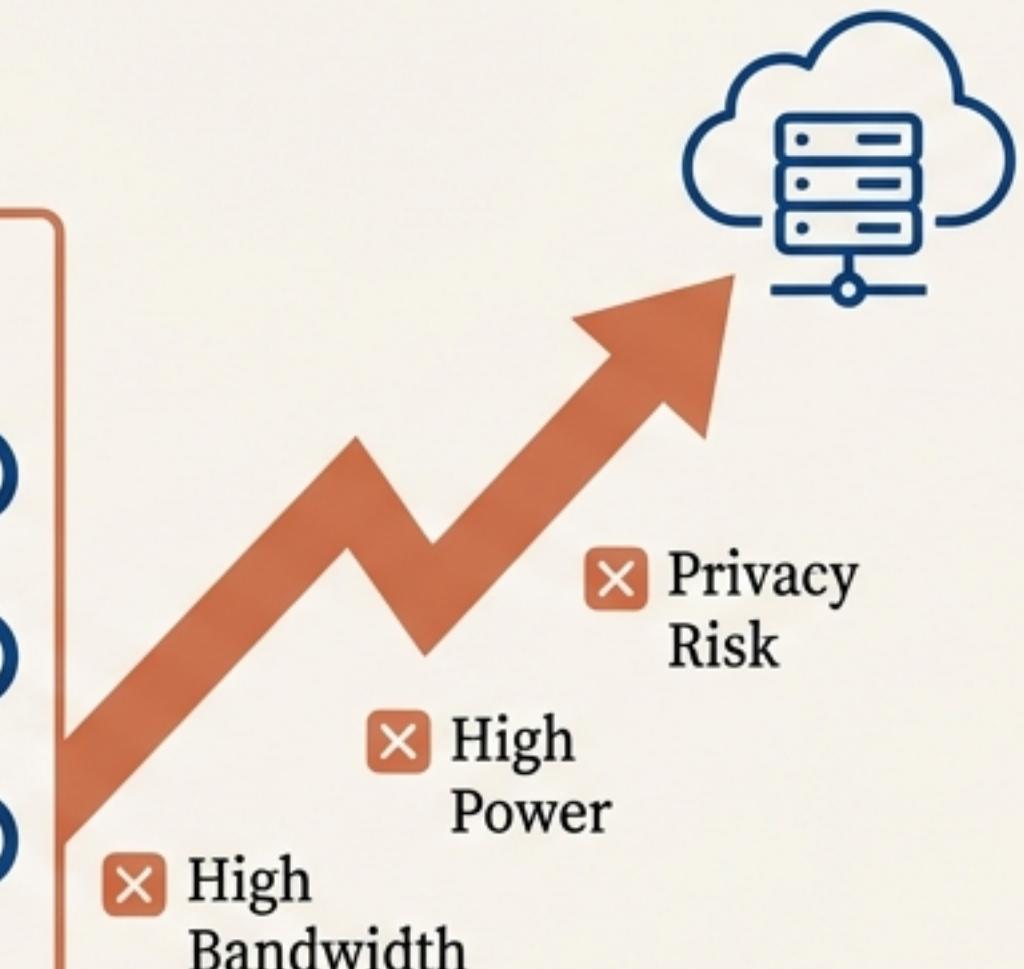
## The Flawed Alternative

Offloading raw sensor data (like images) to a remote server is not a solution. This approach introduces prohibitive costs:

-  Excessive Bandwidth & Power
-  Consumption
-  Significant Data Privacy Risks



**MCU Constraints**  
(<1MB SRAM)



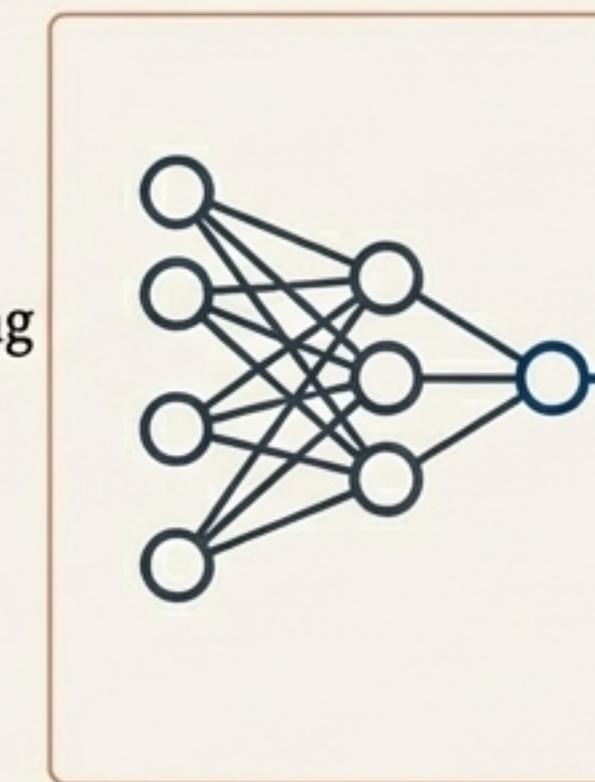
# The Solution: A Hybrid Approach with Split Learning

## Core Concept

Instead of an all-or-nothing choice, we partition the neural network across devices with different capabilities.

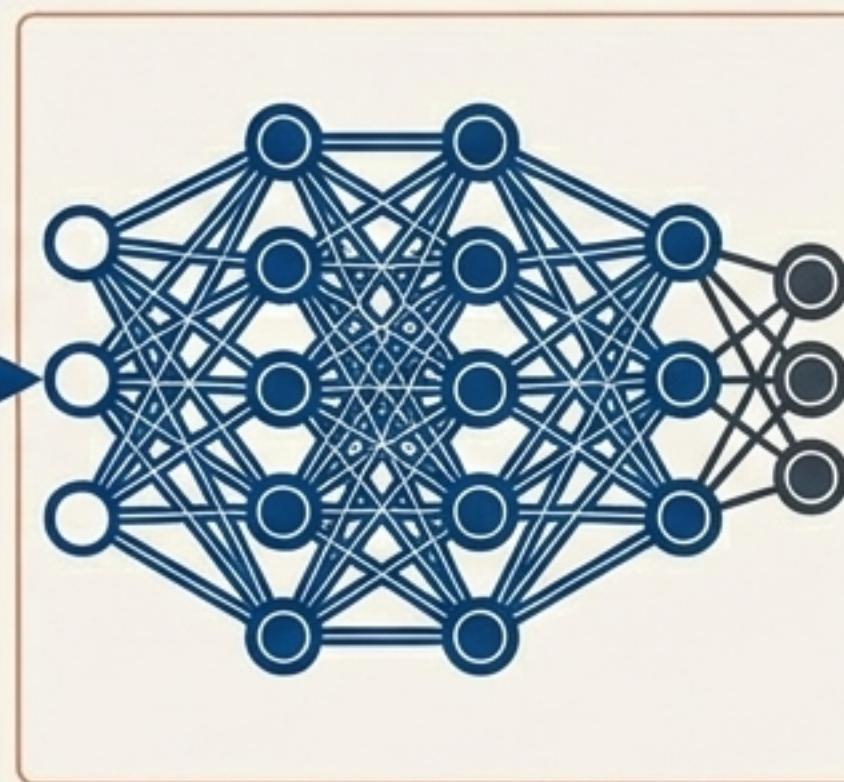
### Arduino (ClientNet)

The Scout



### Raspberry Pi (ServerNet)

The Command Center



## The Critical Advantage

Only a small, compressed feature tensor ('smashed data') is transmitted. Raw data never leaves the edge device, preserving privacy and dramatically reducing bandwidth.

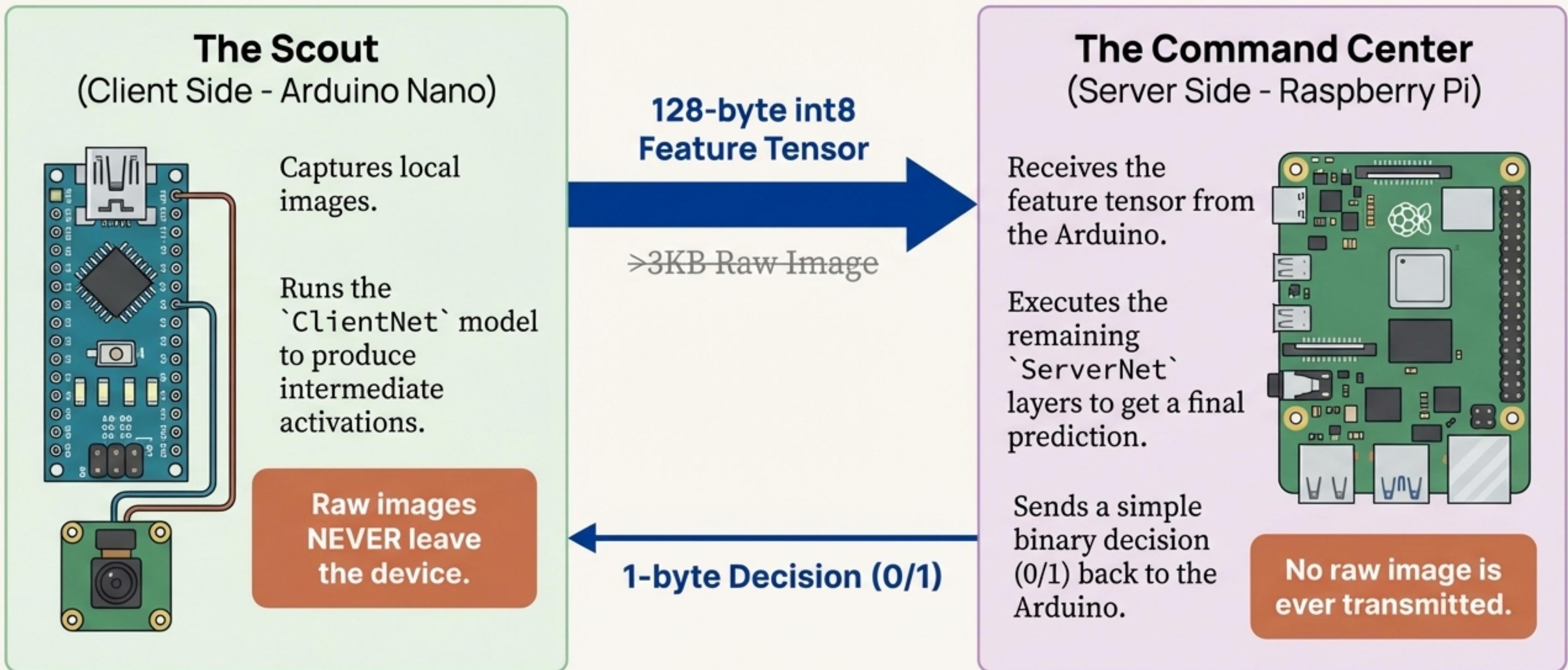
## ClientNet

The first few lightweight layers run on the Arduino. It performs initial feature extraction on-device, identifying basic patterns.

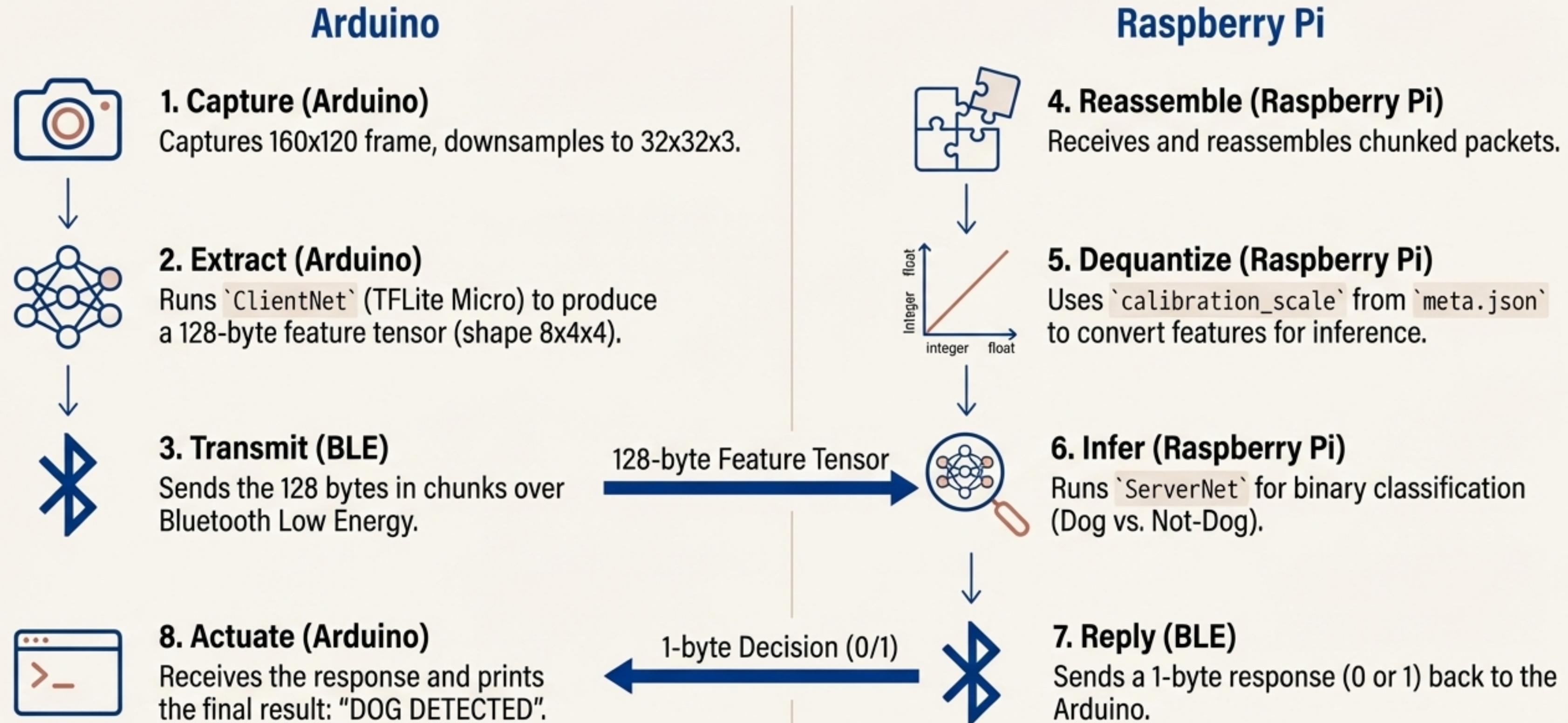
## ServerNet

The remaining, computationally heavy layers run on a more capable local Raspberry Pi, performing the final classification.

# System Architecture: A Tale of Two Chips



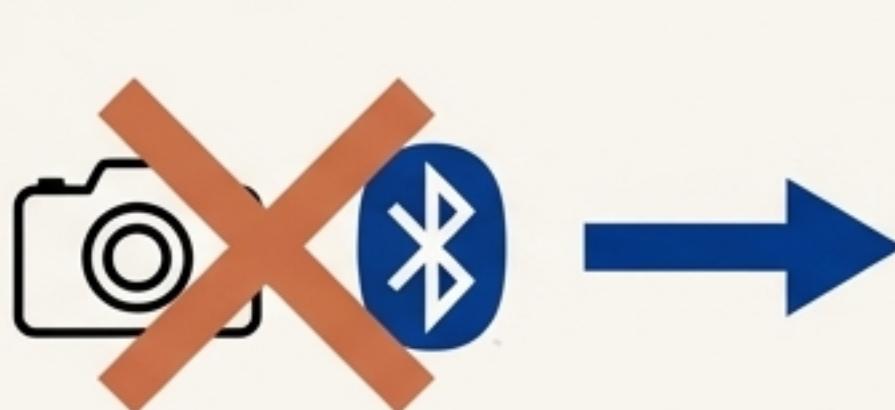
# The End-to-End Data Pipeline



# The Client: Engineering for Extreme Constraints on Arduino

## The Hardware Challenge

- **Memory/Compute:** Limited SRAM prevents running convolutions on data exceeding 20KB.
- **The Hardware Conflict:** The camera and Bluetooth radio cannot operate simultaneously. This discovery forced a sequential, buffered workflow.



## Implementation Details (`ClientNet`)

- **Model:** A `client\_model.h` TFLite Micro model header contains the client-side network.
- **Workflow:**
  1. Capture a frame and populate the input tensor.
  2. Invoke the TFLM interpreter, generating a 128-byte `int8\_t` feature array.
  3. Store these features locally.
  4. *Only after the camera is inactive*, the device begins advertising and sending the stored features over BLE.

## `Workingsketch\_Picture\_BLE\_Reply (1).txt`

```
// Run client model and store intermediate features
bool capture_and_compute_features_when_idle() {
    if (BLE.connected()) return false;

    Serial.println("⚠️ Capturing + client infer (NOT
connected) ...");

    if (!fill_input_from_camera_frame()) return false;

    if (interpreter->Invoke() != kTfLiteOk) {
        Serial.println("✗ Invoke failed");
        return false;
    }
    // ...
    memcpy(feats_int8, output->data.int8, FEAT_LEN);
    feats_ready = true;
    Serial.println("✓ Features ready (stored)");
    return true;
}
```

This line is the critical gate. It explicitly prevents camera operation while a BLE connection is active, enforcing the sequential workflow.

# The Server: Inference and Decision-Making on Raspberry Pi

## Primary Role

To receive fragmented data, reconstruct the feature tensor, and perform the final, heavier part of the inference.

## Implementation Details (`ServerNet`)

- Model: A custom PyTorch `ServerConvBinary` model (`server\_conv\_dogbin.pth`) specifically trained to operate on the 8x4x4 feature tensor from the client.
- Workflow (`recv\_features\_convbin\_and\_reply.py`):
  1. Receive & Reassemble: The `bleak` library listens for chunked BLE notifications and reassembles the 128-byte payload.
  2. Dequantize: Converts the `int8` data to `float32` using the crucial `calibration\_scale` from `meta.json`.
  3. Infer: Feeds the tensor into the `ServerNet` model.
  4. Decide & Reply: Classifies based on the logit output (`logit >= 0 → dog=1`) and writes a single byte response back to the Arduino via the `RESP\_UUID`.

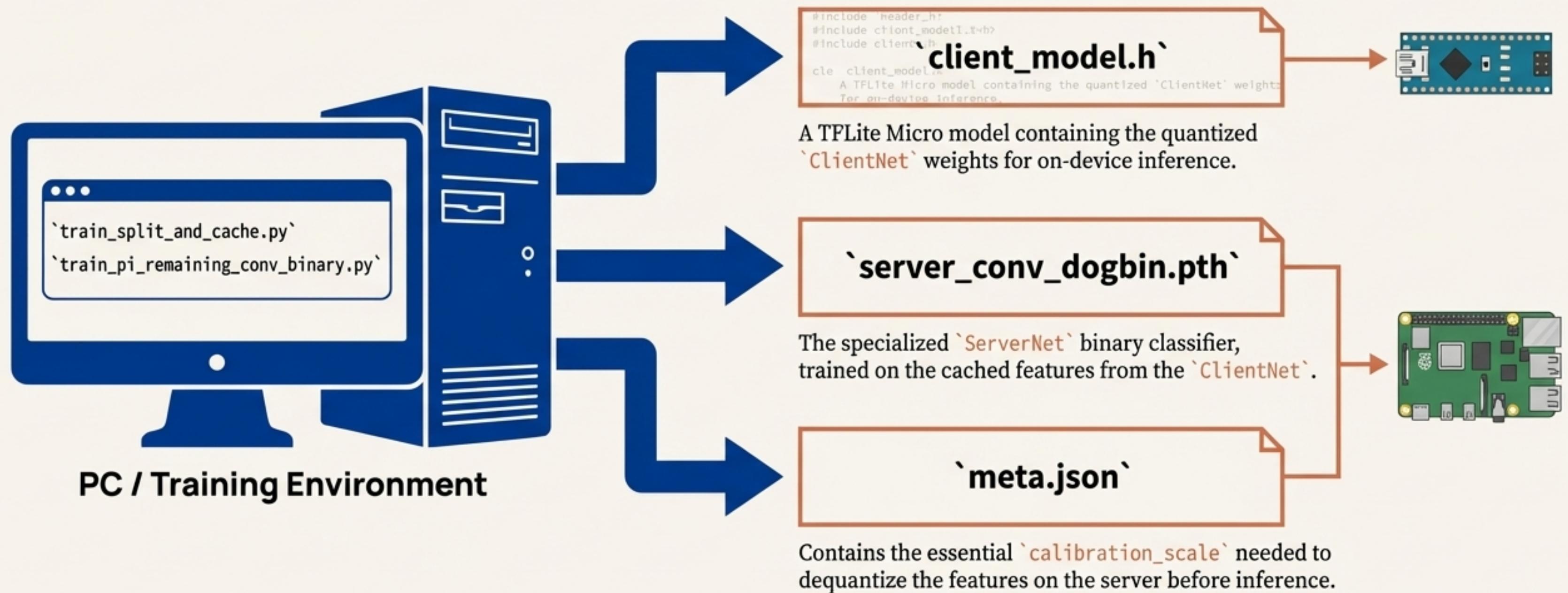
recv\_features\_convbin\_reply\_pi (1).txt

```
class ServerConvBinary(nn.Module):  
    def __init__(self):  
        super().__init__()  
        self.net = nn.Sequential(  
            nn.Conv2d(8, 16, 3, padding=1), ←  
            nn.ReLU(inplace=True),  
            nn.Conv2d(16, 16, 3, padding=1),  
            nn.ReLU(inplace=True),  
            nn.AdaptiveAvgPool2d((1, 1)),  
            nn.Flatten(),  
            nn.Linear(16, 1),  
        )  
  
    def forward(self, x):  
        return self.net(x).squeeze(1)
```

The first convolutional layer is explicitly designed with `in\_channels=8`, proving this model is purpose-built to accept the 8x4x4 feature tensor from the Arduino, not a raw image.

# The Training Strategy: Engineering a Cohesive Split Model

The entire split network was trained *end-to-end* on a PC using PyTorch and the CIFAR-10 dataset. This ensures the client's feature representation is optimally designed for the server's classification task.



# Performance Validated: Key System Metrics



## Bandwidth Reduction **>99%**

A **128-byte** quantized feature tensor was sent instead of the **38,400-byte** raw camera frame, resulting in a dramatic reduction in data transmission per inference.



## Data Privacy **100%**

Raw visual data never leaves the Arduino. Only an abstract, non-human-readable feature representation is transmitted, ensuring on-device privacy.



## Compute Offloading **Successful**

The most computationally expensive CNN layers were successfully offloaded from the memory-constrained Arduino to the more capable Raspberry Pi.



## System Accuracy **83.40%**

The final Pi-side binary classifier achieved a test accuracy of 83.40% on the dog vs. not-dog classification task.

# Evidence of a Successful Round Trip

## 1. Arduino: Capture & Actuate

```
14:32:57.788 -> RAC: 1c:c8:ef:c5:c0:15
14:32:57.788 -> 🚫 Capturing + client infer (NOT connected) ...
14:32:59.699 -> ✅ Features ready (stored for BLE send)
16:49:58.490 -> Initializing camera...
16:49:59.621 -> ✅ Camera initialised
16:49:59.621 -> Client output bytes: 128
16:49:59.789 -> ✅ BLE Advertising: Nano33BLE-ClientNet
16:49:59.789 -> RAC: 1c:c8:ef:c5:c0:15
16:49:59.625 -> 🚫 Capturing + client infer (NOT connected) ...
16:50:01.921 -> ✅ Features ready (stored for BLE send)
21:13:41.443 -> ✅ BLE Advertising
21:13:41.443 -> 🚫 Capturing + client infer (NOT connected) ...
21:13:43.548 -> ✅ Features ready (stored)
21:15:16.294 -> 🚫 P1 response: 0
21:15:16.294 -> 🚫 NOT A DOG
```



Arduino captures the image, computes 128 bytes of features, and prints the final result received from the Pi.

## 2. Raspberry Pi: Receive & Infer

```
(base) eypi13@eypi13:~/artifacts $ client_loop: send disconnect: Connect
C:\Users\Hovssh eypi13@eypi13.local
eypi13@eypi13.local's password:
Linux eypi13 6.12.47+rp1-rpi-v8 81 SMP PREEMPT Bebian 1:6.12.47-1+rp1 (
The programs included with the Debian GNU/Linux system are free software;
the exact distribution terms for each program are described in the
individual files in /usr/share/doc/4/copyright.

Debian GNU/Linux comes with ABSOLUTELY NO WARRANTY, to the extent
permitted by applicable law.
Last login: Sat Dec 20 13:40:08 2025 from 199.168.8.100
(base) eypi13@eypi13:~ $ ls -lh /hexe/oypi13/art2facts | grep server_
-rw-r--r-- 1 eypi13 eypi13 17M Dec 20 21:06 server_cenv_degbtn.pth
(base) eypi13@eypi13:~ $ ls -lh /hexe/oypi13/art2facts/eeta.json
(base) eypi13@eypi13:~ $ ls -lh /hexe/oypi13/art2facts/meta
(base) eypi13@eypi13:~ $ nuno cecv_features_ronebih_and_reply.py
(base) eypi13@eypi13:~ $ sudo -/miniforge3/envs/splitenv/bin/python recv
✅ Redef loaded
✅ calibration_scale = 0.B4724494488188976
✖ Scanning for Arduino...
✅ Found device: 1C:C8:EF:C5:C0:15
✅ Connected to Arduino
⚠ Feature S measured (128 bytes)
✖ logit=-12.9388 → dsgr8
✖ Sent response to Arduino: 0
```

Pi connects, receives the 128-byte frame, runs inference to get a logit value, decides the class, and sends the 1-byte response.

## 3. PC: Training & Calibration

```
== Phase A: Float training ==
Epoch 01 | loss=1.7461 | acc=25.21%
Epoch 02 | loss=1.6712 | acc=48.82%
Epoch 03 | loss=1.5553 | acc=74.82%
Epoch 04 | loss=1.4702 | acc=97.60%
Epoch 05 | loss=1.4089 | acc=99.86%
Epoch 06 | loss=1.3098 | acc=82.32%
Epoch 07 | loss=1.3154 | acc=83.36%
Epoch 08 | loss=1.2794 | acc=85.81%
Epoch 09 | loss=1.2632 | acc=88.21%
Epoch 10 | loss=1.2182 | acc=97.21%
Epoch 11 | loss=1.1891 | acc=98.12%
Epoch 12 | loss=1.1760 | acc=98.76%
Epoch 13 | loss=1.1272 | acc=99.60%
Epoch 14 | loss=1.1115 | acc=99.12%
Epoch 15 | loss=1.0813 | acc=99.68%
Epoch 16 | loss=1.0672 | acc=99.38%
Epoch 17 | loss=1.0645 | acc=99.38%
Epoch 18 | loss=1.0645 | acc=99.38%
Epoch 19 | loss=1.0482 | acc=99.21%
Epoch 20 | loss=1.0195 | acc=99.26%
Epoch 21 | loss=1.0151 | acc=99.86%
Epoch 22 | loss=1.0088 | acc=99.82%
Epoch 23 | loss=0.9988 | acc=99.54%
Epoch 24 | loss=0.9780 | acc=99.72%
```

Computing calibration scale...
calibration\_scale = 0.B4724494488188976

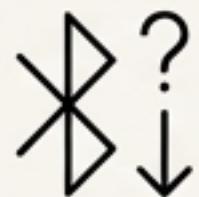
PC-side training generates not only the models but also the critical `calibration\_scale` value, which is essential for the Pi to correctly interpret the Arduino's data.

# Navigating Real-World Implementation Challenges

## Challenge 1: Unstable Bluetooth Low Energy

### Problem

Frequent, unexplained disconnections during continuous data transmission.



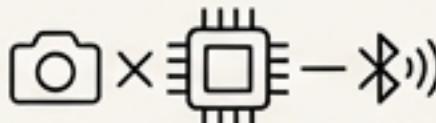
### Solution

Required extensive debugging of UUIDs, tuning transmission delays (`CHUNK_DELAY_MS`), and implementing robust packet reassembly logic on the Pi.

## Challenge 2: Camera vs. BLE Hardware Conflict

### Problem

The Arduino Nano 33 BLE Sense hardware cannot operate the camera and radio simultaneously.



### Solution

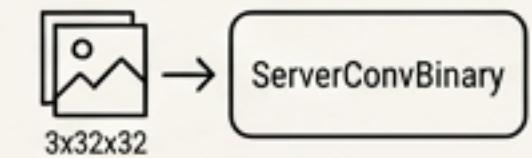
Architected a sequential “capture-then-send” workflow, buffering features locally before initiating any BLE communication.

capture → buffer → send

## Challenge 3: Model Input Incompatibility

### Problem

Early Pi-side models failed because they expected a raw image (e.g., 3x32x32) instead of the 8x4x4 feature tensor from the Arduino.



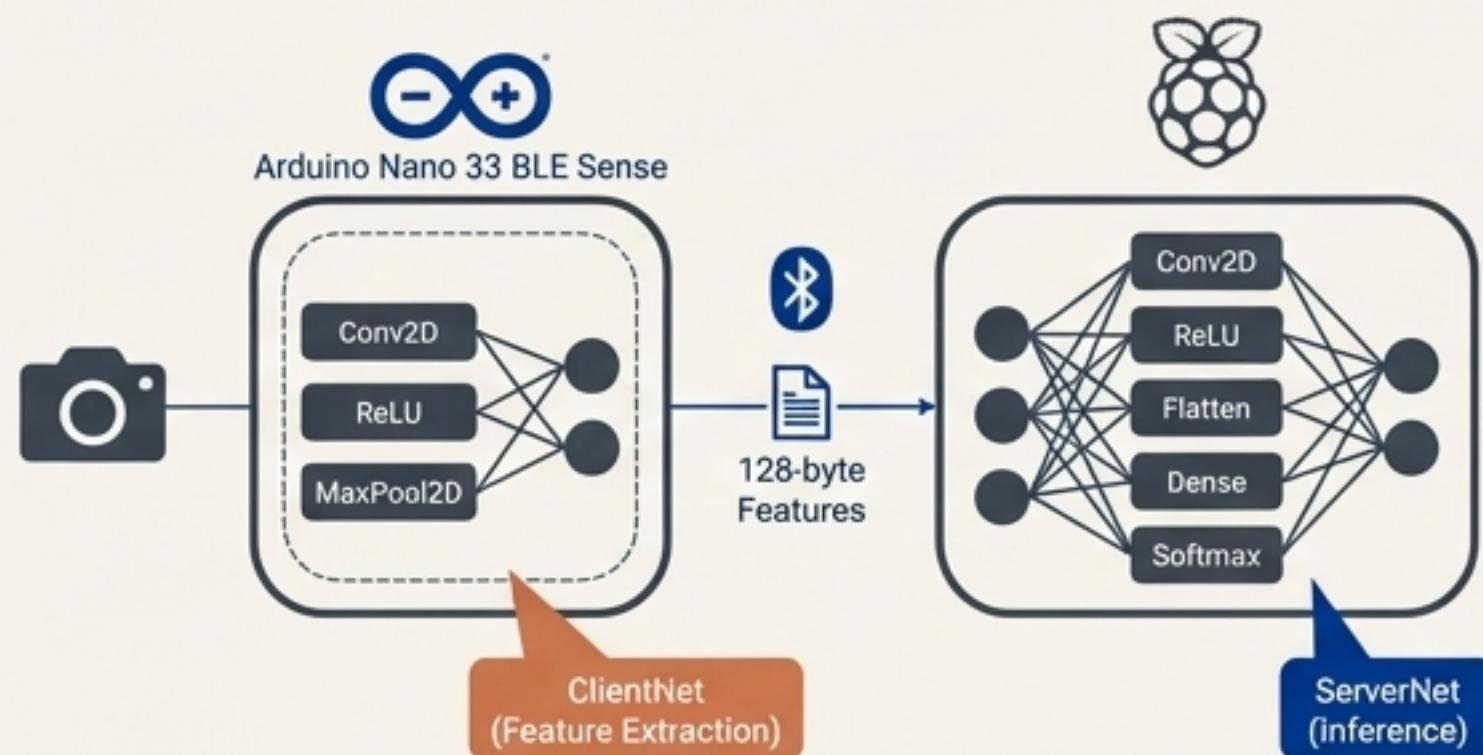
### Solution

Trained a new, custom “conv head” (`ServerConvBinary`) for the Pi specifically designed to match the client’s output contract.

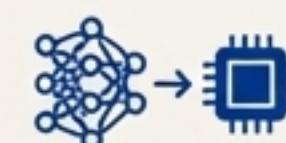
# A Successful Scaffold for Hybrid Edge Intelligence

## Key Accomplishment

This project successfully demonstrates a practical, low-bandwidth, and privacy-preserving solution for edge inference by intelligently splitting a neural network across constrained, heterogeneous hardware.



## Future Directions & Improvements



- Improve Client Accuracy:** Port the fully trained `ClientNet` weights (not just an approximation) to the Arduino's C++ environment. This is a complex but valuable task to maximize system accuracy.



- Latency Optimization:** Profile and reduce the end-to-end time delay introduced by BLE transmission and Pi-side inference to improve performance for real-time applications.



- Expand Functionality:** Move beyond a binary classifier to the full 10-class CIFAR-10 model, requiring modifications to the `ServerNet` and the Arduino's response handling.