**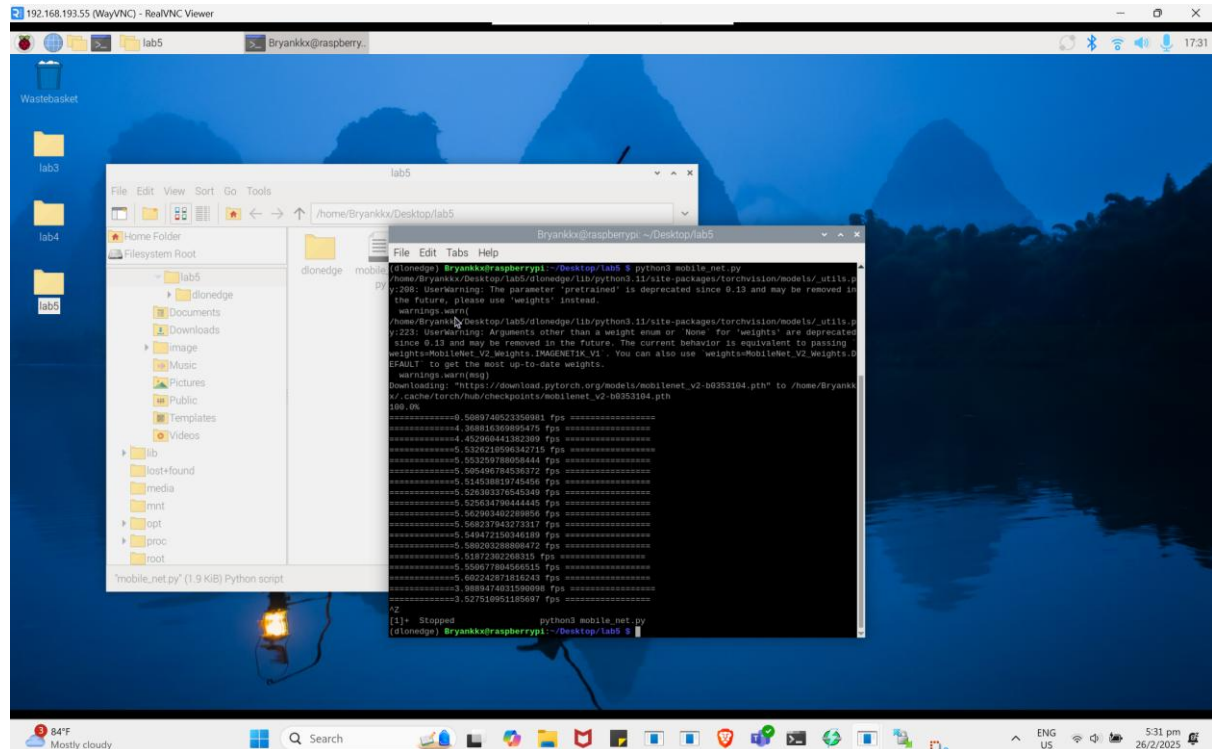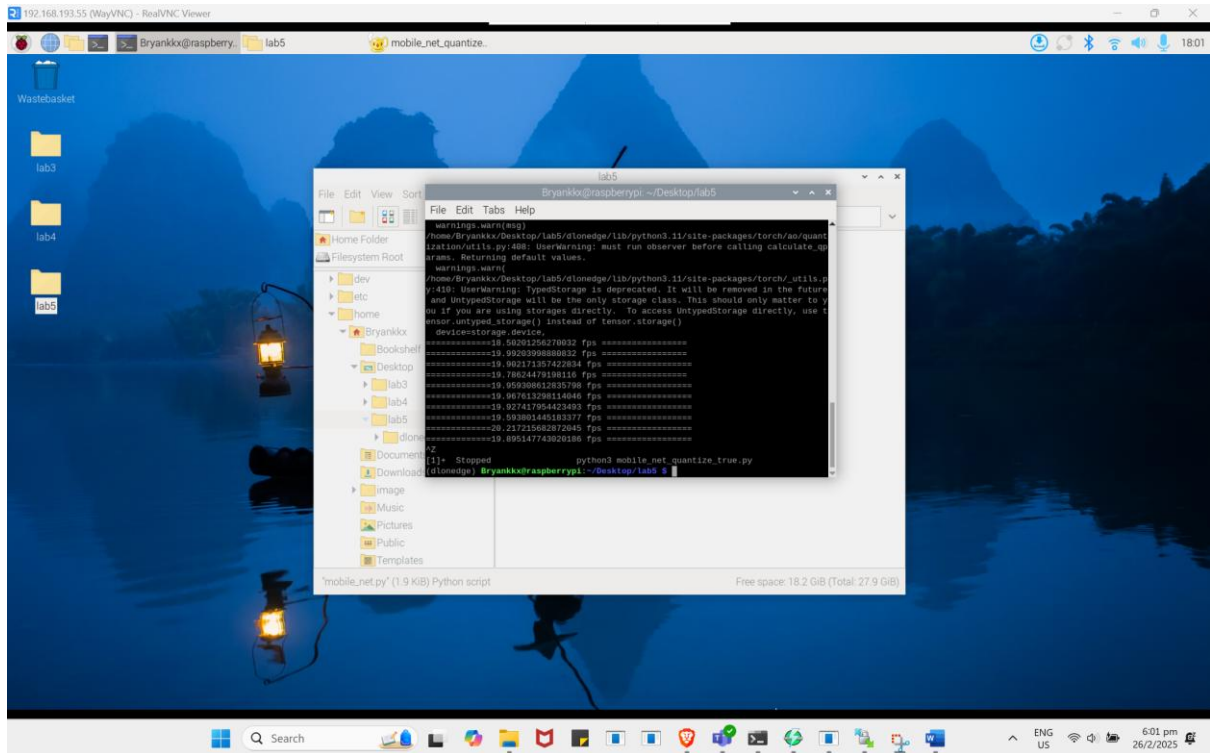Part 1.** sample code is used to directly load pre-trained MobileNetV2 model, doing model inference and finally, Observe the fps as shown in screenshot below when run on RaspberryPi 4B. As shown, with no optimization of model, we could only achieve of 5-6 fps much below our desired target.



**Part 2.** Edit line number 11 as shown below to enable quantization in sample code to use quantized version of MobileNetV2 model.

- **Part 3.** Uncomment lines 57-61 in [sample code](#) to print the top 10 predictions in real-time as shown in below video