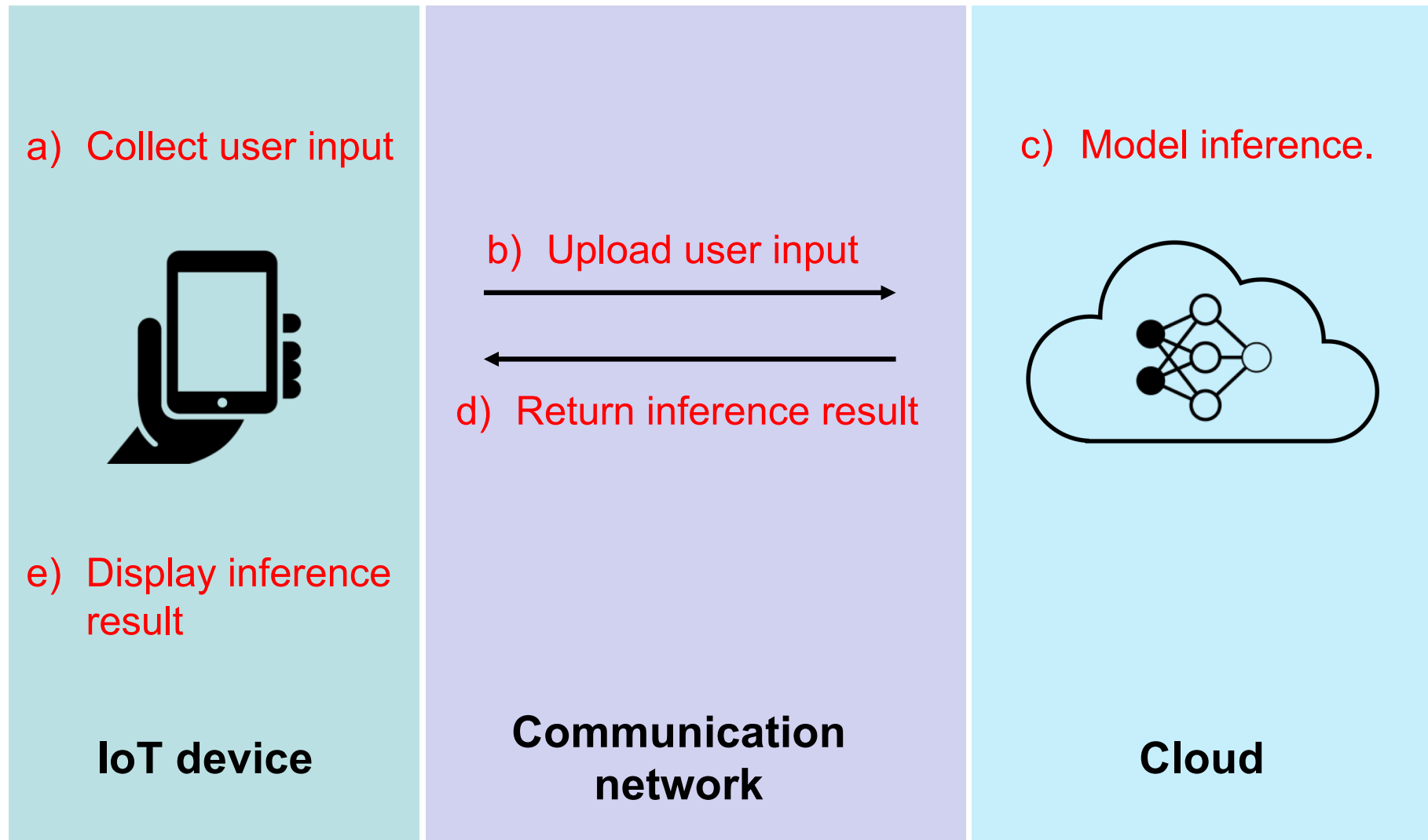# Course Project (35%)

- Offloading AI inference from IoT device to cloud
    - IoT devices are often incapable of AI inference
    - Offload computing from IoT device to server
- Details to be given in Week 3 lecture time
- Project submission
    - Deadline: May 4th 2022
    - Source code package
    - URL to a video demo with max length of 10 minutes (youtube preferred, keep the video at the URL until Oct 31st 2022)

# Ultimate Goal

a) Collect user input



e) Display inference result

**IoT device**

b) Upload user input

d) Return inference result

**Communication network**

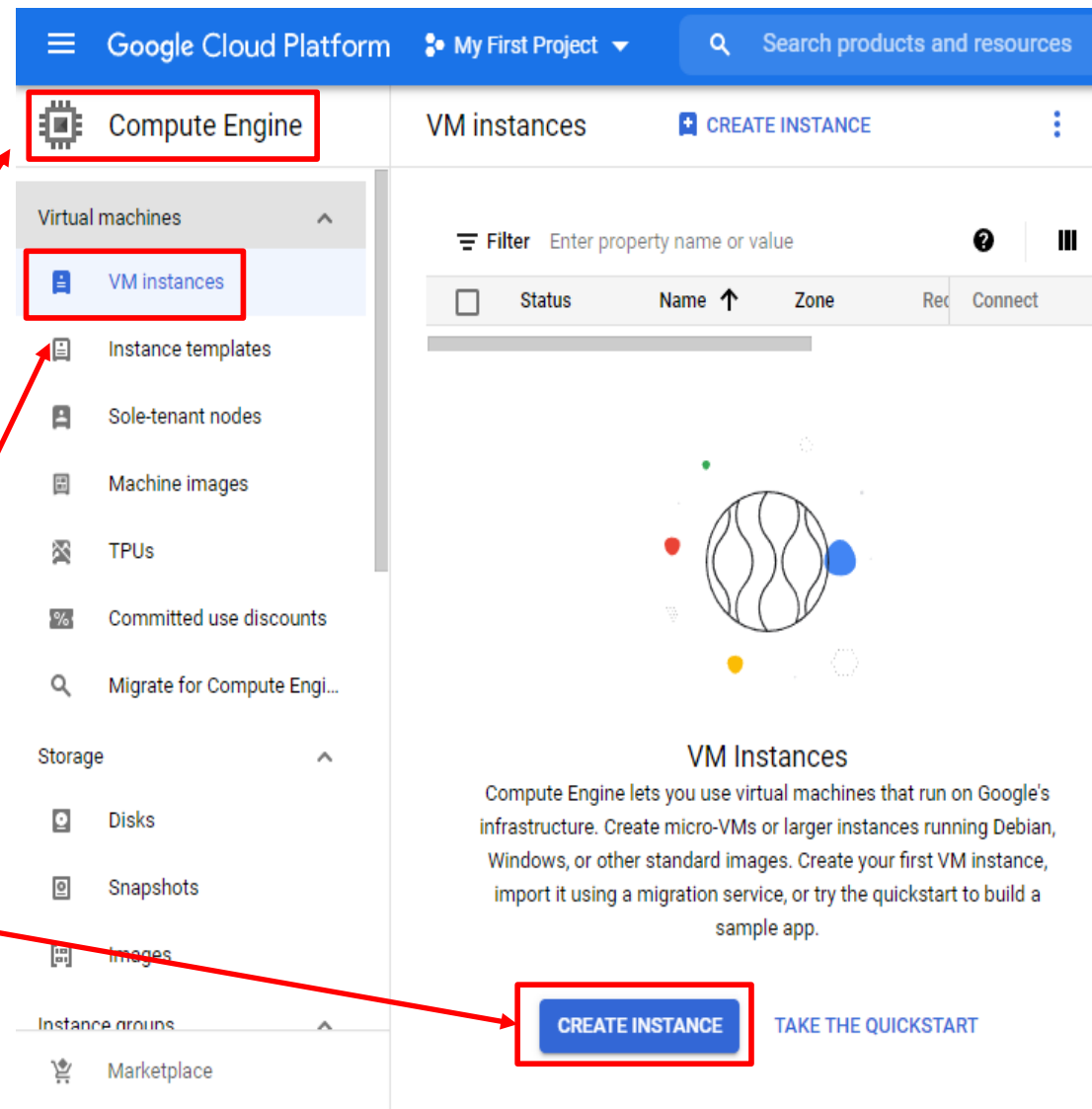c) Model inference.



**Cloud**

# Tips

- You may choose any applications.
  - E.g., speech recognition, image classification, etc.
- Any mobile platform for the IoT app, e.g., Android, iOS.
  - If you don't have a physical device, you may use the phone emulator.
- You may use the VM on Google Could Platform or your own PC to host the cloud service.
  - Google Cloud: 90-day, $300 Free Trial
  - **You're responsible for managing your Google Cloud account to avoid any monetary charge**
  - **In particular, close/remove the billing account after the project**
  - **If you're unsure, use your own PC to host the service**
- Any programming language for the server app
  - Python + Flask is recommended.

# Creating a VM on Google Cloud Platform

- Visit the following URL and register an account
  - https://cloud.google.com/compute
- Follow the guides to create the VM of Linux / Windows
  - https://cloud.google.com/compute/docs/quickstarts

# Creating a VM on Google Cloud Platform

i.   Log in to home page

ii.  Find "Computing Engine"

iii. Search for "VM instances"

iv.  Choose "Create Instance".

# Creating a VM on Google Cloud Platform

- Follow the steps to create the VM

# Creating a VM on Google Cloud Platform

- Check status of VM
  - Find VM's IP for your setup
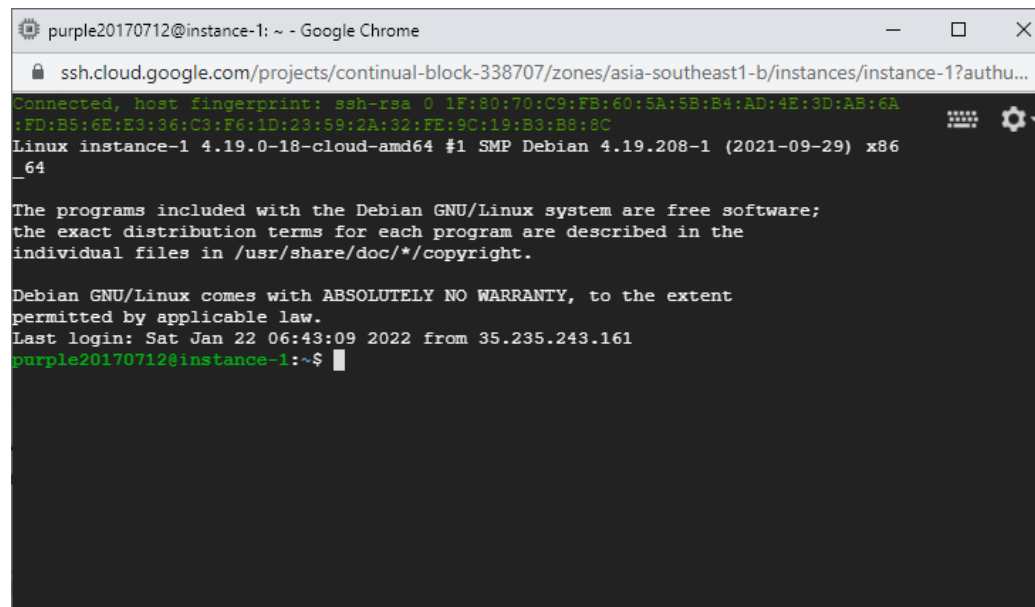  - Click "SSH" to connect the VM

# Tasks (100 marks)

- Local inference                                               50 marks
  - Collect user input                                        (15)
  - Infer locally and display result                          (20)
  - Run on emulated/physical IoT device                       (15)
- Cloud inference                                               30 marks
  - Run inference in cloud virtual machine                     (10)
  - Communicate btw IoT device & cloud                         (20)
- Advanced tasks                            capped at   20 marks
  - **If** train your own model                               (10)
  - **If** Support multiple concurrent users                  (10)
  - **If** support online model updating                      (10)

# Marking Criteria (1)

- Collect user input (15)
  - **If** load input from storage (8)
  - **If** collect real-time input by touch screen, microphone, or camera (15)

# Marking Criteria (2)

- Infer locally and display result      (20)
  - **If** run model inference on the mobile app    (15)
    - If your app offloads the inference to cloud, you have the marks for this part automatically
  - **If** run heuristic algorithm (not neural network) on the mobile app      (9)
    - If your app offloads the execution to cloud, you have the marks for this part automatically
  - Display the inference result in real time by screen or synthetic voice      (5)

# Marking Criteria (3)

- Run on emulated/physical IoT device      (15)
    - **If** the program runs natively on a desktop or laptop computer      (5)
    - **If** the program runs on an emulated or a physical IoT device      (15)
        - iOS emulator provided by Xcode
        - Android emulator
        - Real smartphone
        - Raspberry Pi + add-on sensors (camera, microphone, etc)

# Marking Criteria (4)

- Run inference in cloud virtual machine (10)
  - **If** deploy server program on a cloud virtual machine, e.g., Azure. (10)
  - **If** deploy server program on your own computer (10)

# Marking Criteria (5)

- Communication btw IoT device and cloud    (20)

    - Send the user input from the mobile app to the cloud for inference                                           (10)

    - Send the inference result from the cloud to the mobile app                                                          (10)

# Marking Criteria (6)

- ## Train your own model  (10)
  - **If** train the used ML model by yourself  (10)
    - Training program should be in the code package
    - Training results (e.g., accuracy) in a readme file
    - The training data can be any publicly available dataset
  - **If** use downloaded pre-trained model  (8)
  - **If** use heuristic algorithm (not neural network)  (6)

# Marking Criteria (7)

- Support multiple concurrent users       (10)
  - Demonstrate multiple IoT devices can use the cloud service simultaneously

# Marking Criteria (8)

- Online model updating <span style="color:red">(10)</span>
  - Demonstrate that the model in the cloud can be updated at run time using newly collected user data
    - User input and the corresponding label are transmitted to cloud
    - The cloud retrains the model

# Grading

- Project should be independently accomplished by each student.
- Q&A and investigation may be conducted on similar topics and implementations
- 3-day grace period
  - No penalty if a valid excuse provided
  - Otherwise, a penalty of 20% reduction will be applied to the mark of the late submissions
- Submissions after the grace period
  - Zero mark