# Homework 2

## Zhengxiang Huang

### 520021910014

# Contents

# 1 Introduction

In this Homework, you shall yourself explore more on how to implement quantization, and deploy LLM on Nano with the multiple online resources.

We sincerely hope you enjoy your journey through this homework as we do. Have fun!

# 2 Q1: Quantization in PyTorch (85')

## 2.1 Problem Formulation

We've discussed quantization in class.

Now, try to quantize your own network in Homework 1 Q5 on your own laptop computer.

You shall then be able to answer the following questions:

1. (5') Q: Screenshot the results of "print(qconfig)". Explain the parameters of it. e.g., explain what is "per_channel_symmetric".

2. (10') Q: Implement "fuse_model" function for your model, explain how "Conv2d", "Batch-Norm2d", and "ReLU" are fused. Screenshot the results of "print(model)" after you fuse it.

3. (10') Q: Calibrate and quantize your model in PyTorch. Screenshot the results of "print(model)" after you quantize it. Explain why we need calibration phase, and what we are observing during that phase.

4. (20') Q: Measure the size, inference accuracy, and inference time of your model after quantization. Compare the size, time, and accuracy before and after quantization.

5. (40') Q: Submit your quantization code, whether it's runnable or not. (Even it can't work, we attribute some points to you.)

6. (Bonus +15') Q: Quantized model can't be directly deployed on nano, we need onnx workaround. We have a line of "torch.onnx.export" in our "example.py". Search online how to install and run the "quant.onnx" you export, and deploy that stuff onto your Nano board with onnx. If onnx is impossible on nano, write in the report to inform me, and I will attributes bouns to you too. (No hints nor guides, search yourself! hahahaha!)

(**Hint**: It's suggested that you install a CPU version of Pytorch on your own computer if you don't have a powerful Nvidia GPU installed.

"pip install torch==2.0.1 torchvision==0.15.2 –index-url https://download.pytorch.org/whl/cpu" to install a cpu version pytorch on your computer, if you previously have no pytorch installed.)

(**Hint**: Refer to https://pytorch.org/docs/stable/quantization.html#post-training-static-quantization and https://pytorch.org/tutorials/advanced/static_quantization_tutorial.html for more information of how to quantize a model in PyTorch. Also, see our "example.py" and "model.py" provided.)
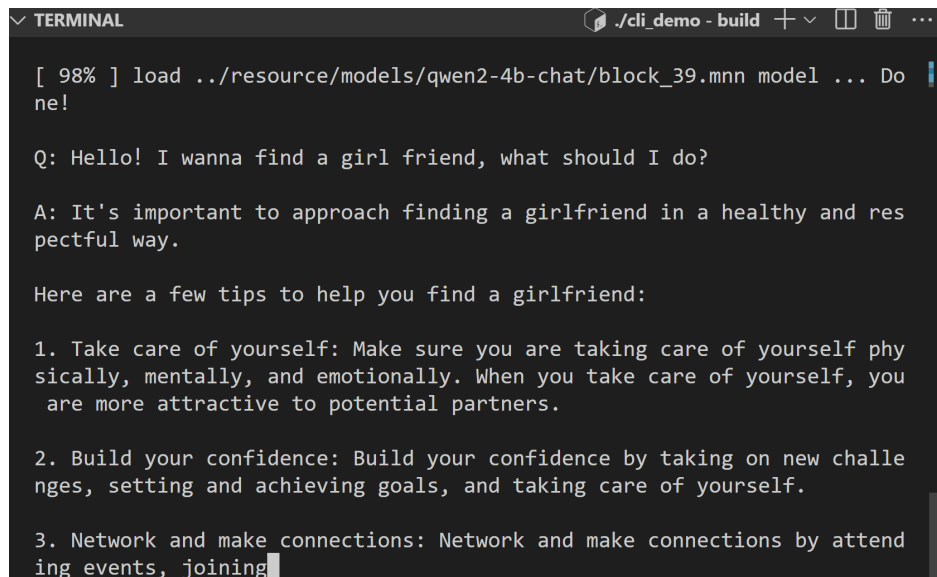
## 2.2 Solution

(Your answers go here!)

# 3 Q2: Advanced Topic, LLM on Nano (15')

## 3.1 Problem Formulation

1. (12') Follow our instructions in LLM-README.md to deploy a LLM (Qwen1.5-4B) on Nano. Measure time and space it consumes. Have fun chatting with it! Screenshot your conversions.



Figure 1: Nano Chat Bot

The resource LLM takes up.



Figure 2: LLM resource

(Hint: If you can't deploy it, screenshot the step you get stuck and we will attribute points to you!)

2. (3') Tell me how many bit are this model's weights quantized into? 4 bits? 8 bits? No quantization? Where did you find that info?

## 3.2 Solution

(Your answer here!)

# 4 What have you learned?

## 4.1 Problem Formulation

(Your answer here!)

## 4.2 Tell me what you have learned

(Your answer here!)

# 5 Acknowledgement

(Your answer here!)

# References