# The Distributions of LLM's
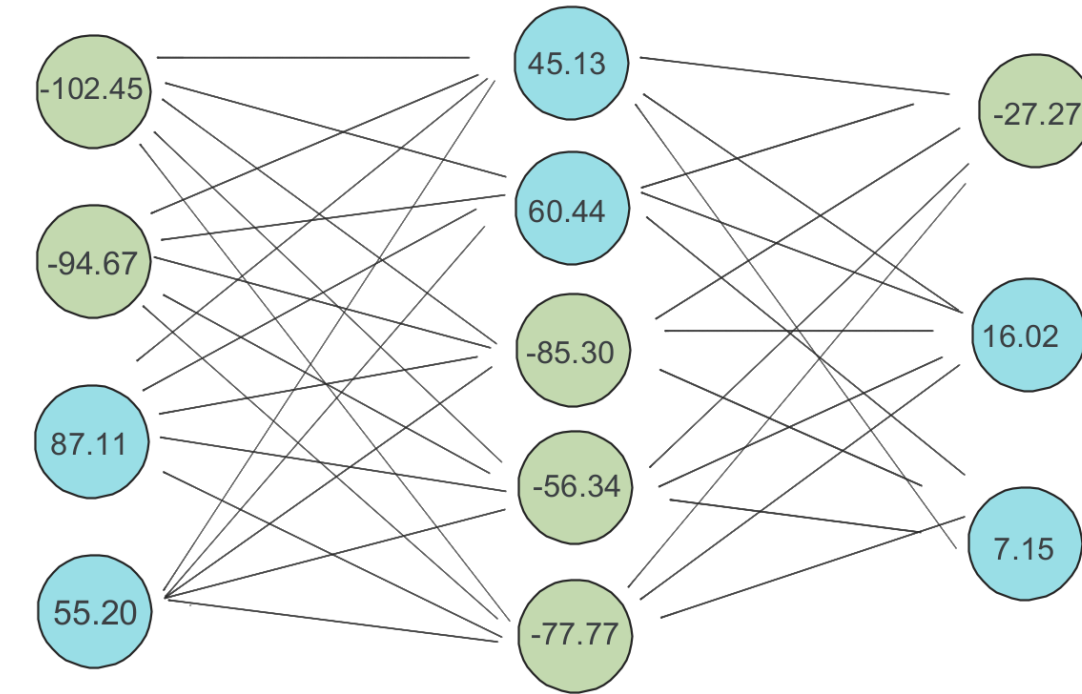
Yosvel Baez, UNC Charlotte
Tyler Allen, College of Computing and Informatics

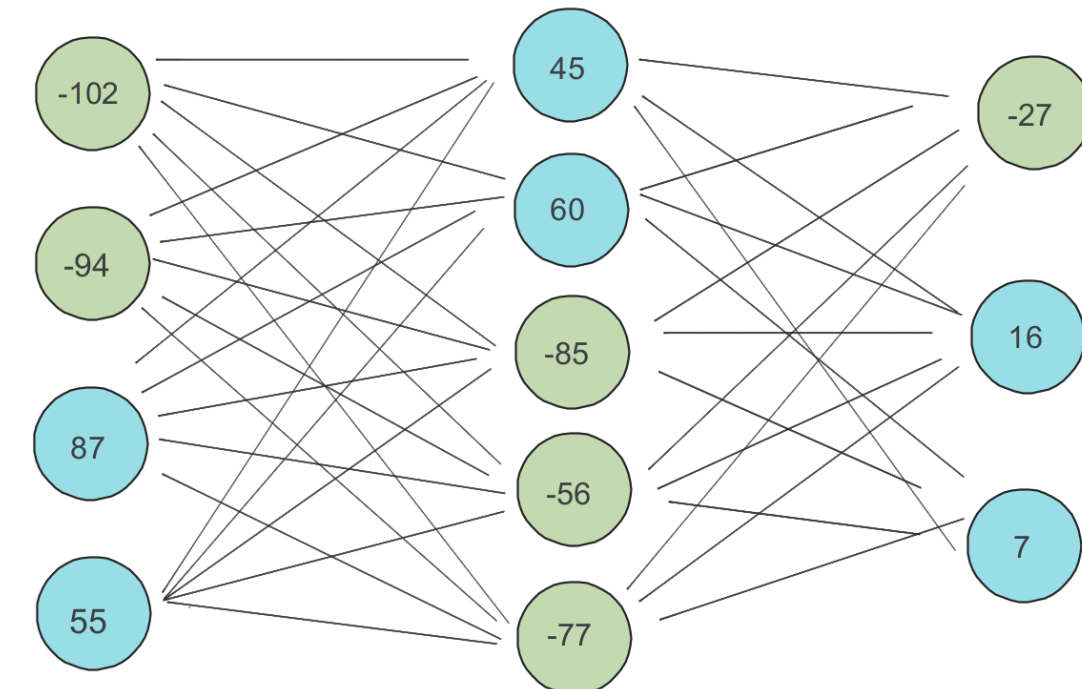UNIVERSITY OF NORTH CAROLINA
CHARLOTTE

## Introduction

**Activation-aware Weight Quantization (AWQ),** is a method for reducing the size and increasing the speed of LLMS by quantizing their weights which considers the activation patterns of the neural network. This approach enables the deployment of powerful AI models on devices with limited computational resources.

LLM's utilizing AWQ

- LLaMA
- OPT
- Vicuna

Distribution of AWQ ensures that LLM's are both memory and computation-efficient, enabling their deployment on devices with limited resources without sacrificing performance.

## Methodology

**Beginning Stage**

- Familiarizing myself with command line language.
- Exploring GitHub and the repository chosen to conduct research on.

**Mid-Stage**

- Failed setting up my environment in VS Code through the instructions provided within the repository
- Moved onto the HPC server provided by UNCC due to complications from working on my own computer

**Current Stage**

- Ran all command line prompts presented in repository
- Waiting to be approved for an OPT data set from hugging face to complete my results.

## Current Progress



*Screenshot of the work conducted through UNCC HPC server*

## Why utilize AWQ?

- To evaluate the efficacy of Activation-aware Weight Quantization (AWQ) in enhancing the execution efficiency of Large Language Models (LLMs) across diverse computational environments.

- To benchmark the execution performance of LLMs pre-AWQ and post-AWQ application, measuring improvements in speed, memory usage, and energy consumption on both standard and resource-constrained hardware platforms.

## Motivation

- Advancing the field of machine learning by improving the efficiency of model deployment, especially for LLM's.

- Contributing to the development of more sustainable AI by reducing the computational and energy requirements for LLM's.

- Enabling the use of advanced AI models in real-world applications where resources are limited, such as mobile devices or other edge computing scenarios.

## References

- MIT HAN LAB. (n.d.). tinychat. LLM-AWQ. Retrieved November 29, 2023, from https://github.com/mit-han-lab/llm-awq/tree/main/tinychat