

Quantization Backend Configuration

FX Graph Mode Quantization allows the user to configure various quantization behaviors of an op in order to match the expectation of their backend.

In the future, this document will contain a detailed spec of these configurations.

Default values for native configurations

Below is the output of the configuration for quantization of ops in fbgemm and qnnpack (PyTorch's default quantized backends).

Results:

System Message: ERROR/3 (D:\onboarding-resources\sample-onboarding-resources\pytorch-master\docs\source\[pytorch-master] [docs] [source] quantization-backend-configuration.rst, line 20)

Unknown directive type "literalinclude".

```
.. literalinclude:: scripts/quantization_backend_configs/default_backend_config.txt
```