

Inference Custom Float Quantization (4-bit exponent, 4-bit mantissa)

