Feed-Fordward Block Binary Residual Approximation for Salient Weight **Hessian Matrix** Binarized Binarized-FC Projection FC2 0.10 alue 0.000 Residual Activation Binarization Activation 0.006 Binarized 0.05 FC1 2000 MatMul 0.25<u>0</u> Activation Bell-shaped Splitting for Non-salient Weight 2000 MatMul **Binary Weight** Splitting Binarization Binarized Binarized Binarized FC for Q FC for K FC for V **Float Weight** Multi-Head Self-Attrntion Activation **BiLLM Framework**