**Architecture Comparison: Training Loss**