Efficiency Analysis: Quality vs. Speed