**Task Performance by Output Length and Dispersion**
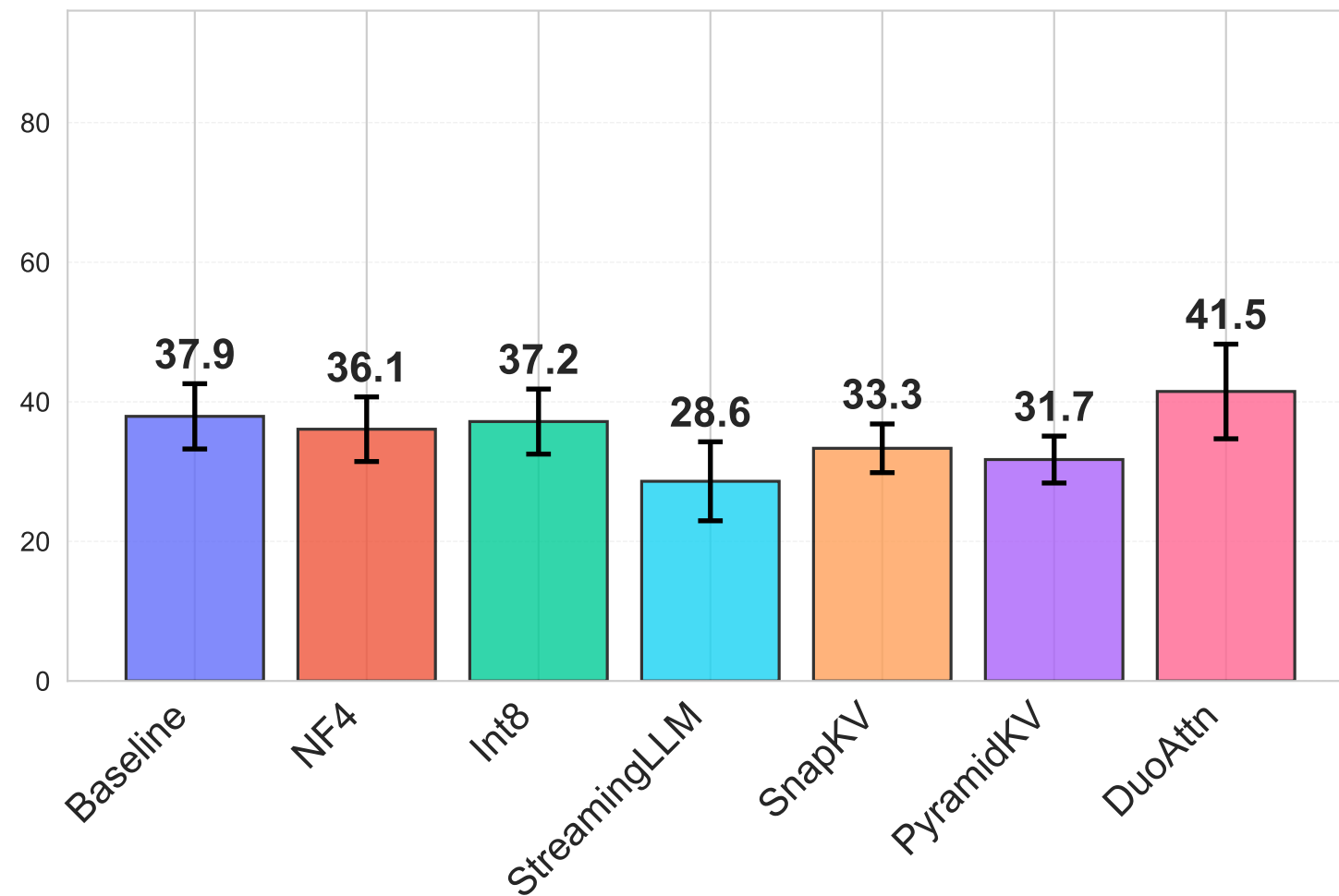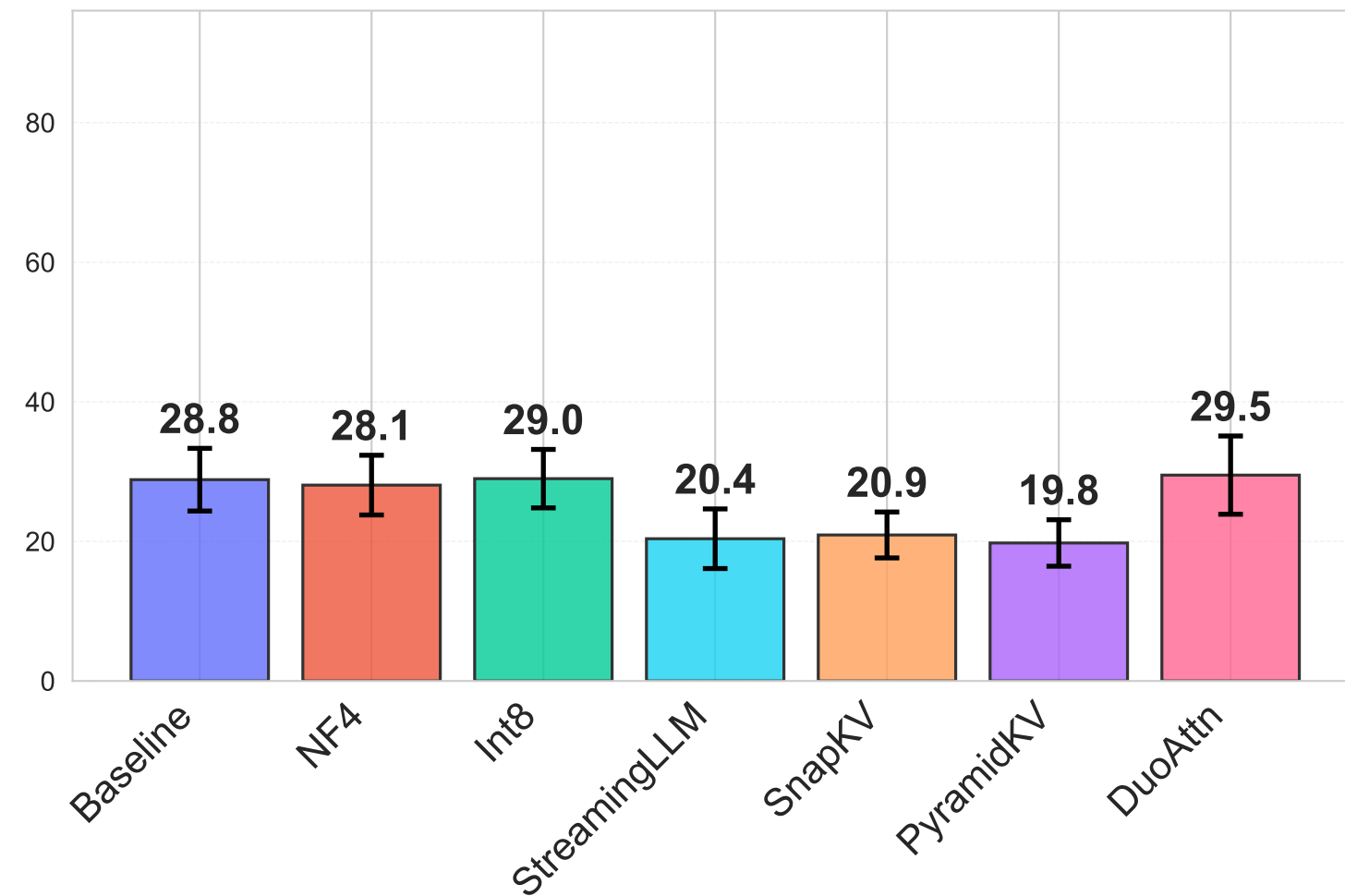**Averaged Across All Models and Context Lengths**
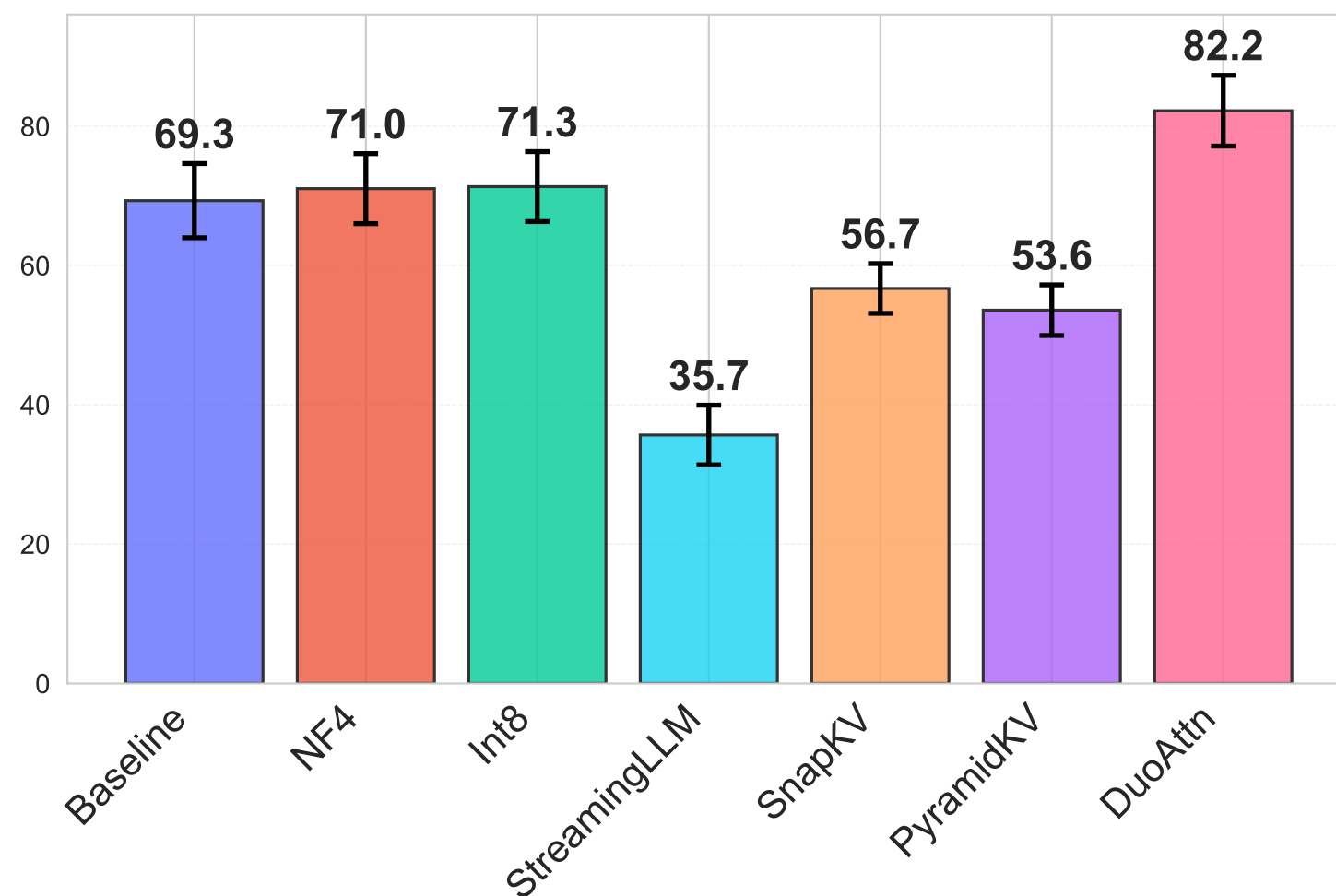
### Short Output, High Dispersion

| Model | Value |
|---|---|
| Baseline | 37.9 |
| NF4 | 36.1 |
| Int8 | 37.2 |
| StreamingLLM | 28.6 |
| SnapKV | 33.3 |
| PyramidKV | 31.7 |
| DuoAttn | 41.5 |

### Long Output, High Dispersion

| Model | Value |
|---|---|
| Baseline | 28.8 |
| NF4 | 28.1 |
| Int8 | 29.0 |
| StreamingLLM | 20.4 |
| SnapKV | 20.9 |
| PyramidKV | 19.8 |
| DuoAttn | 29.5 |

### Short Output, Low Dispersion

| Model | Value |
|---|---|
| Baseline | 69.3 |
| NF4 | 71.0 |
| Int8 | 71.3 |
| StreamingLLM | 35.7 |
| SnapKV | 56.7 |
| PyramidKV | 53.6 |
| DuoAttn | 82.2 |

*No tasks in this category*