

| Sample Pruner | Token Pruner | LLaMA2-7B | | | | | | Mistral-7B | | | | | |
|---------------------------|--------------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|-----------------------------|------------------------|------------------------|------------------------|------------------------|
| | | ARC-E | ARC-C | GSM8K | SQuAD | TriviaQA | Avg | ARC-E | ARC-C | GSM8K | SQuAD | TriviaQA | Avg |
| Zero-Shot | | 53.44 | 38.98 | 5.31 | 12.18 | 43.00 | 30.58 | 66.67 | 46.10 | 18.35 | 10.01 | 43.77 | 36.98 |
| 12.5% Samples, 50% Tokens | | | | | | | | | | | | | |
| Random | Random PPL | 59.25 | 41.02 | 8.11 | 12.75 | 48.75 | 33.98 | 70.55 | 48.14 | 22.74 | 19.57 | 52.63 | 42.73 |
| | FastV | 60.49 ^{71.24} | 43.39 ^{72.37} | 7.20 ^{10.91} | 12.20 ^{10.55} | 48.04 ^{40.71} | 34.26 ^{70.28} | 70.72 ^{70.17} | 48.47 ^{40.33} | 25.78 ^{73.04} | 21.36 ^{71.79} | 53.92 ^{71.29} | 44.05 ^{71.32} |
| | FastV | 59.96 ^{70.71} | 42.37 ^{73.15} | 5.76 ^{62.35} | 11.31 ^{11.44} | 46.42 ^{42.33} | 33.17 ^{10.81} | 70.72 ^{70.17} | 46.44 ^{41.70} | 18.80 ^{43.94} | 19.14 ^{40.43} | 51.56 ^{40.17} | 41.33 ^{41.40} |
| | SparseVLM | 54.23 ^{44.93} | 37.97 ^{43.05} | 7.35 ^{10.76} | 12.76 ^{70.01} | 46.65 ^{44.10} | 31.41 ^{42.57} | 67.02 ^{43.53} | 44.75 ^{43.39} | 20.24 ^{42.50} | 10.97 ^{48.60} | 44.61 ^{48.02} | 37.52 ^{45.21} |
| Longest | Random PPL | 59.96 ^{70.71} | 44.41 ^{43.39} | 7.51 ^{10.60} | 15.34 ^{72.59} | 48.91 ^{10.16} | 35.22 ^{71.24} | 74.25 ^{73.70} | 48.81 ^{40.67} | 28.73 ^{75.99} | 17.66 ^{41.91} | 55.73 ^{73.10} | 45.04 ^{72.31} |
| | PPL | 61.19 ^{71.94} | 43.73 ^{72.71} | 6.82 ^{12.29} | 16.33 ^{73.58} | 48.61 ^{40.59} | 35.24 ^{71.26} | 75.49 ^{74.94} | 50.17 ^{72.03} | 27.98 ^{75.24} | 24.49 ^{71.92} | 56.55 ^{73.92} | 46.33 ^{73.60} |
| | FastV | 59.25 ^{70.00} | 43.05 ^{72.03} | 5.69 ^{12.42} | 13.64 ^{70.89} | 46.98 ^{41.77} | 37.72 ^{40.26} | 74.43 ^{73.88} | 49.15 ^{41.01} | 25.70 ^{72.96} | 22.89 ^{73.32} | 54.15 ^{71.52} | 45.26 ^{73.25} |
| | SparseVLM | 54.32 ^{44.93} | 38.31 ^{42.71} | 7.13 ^{10.98} | 10.92 ^{41.83} | 43.77 ^{44.80} | 30.89 ^{43.69} | 69.49 ^{41.06} | 46.10 ^{42.04} | 28.89 ^{75.15} | 8.62 ^{40.05} | 50.30 ^{42.33} | 40.68 ^{42.05} |
| InfoBatch | Random PPL | 60.31 ^{71.06} | 41.36 ^{70.34} | 5.38 ^{12.73} | 15.71 ^{72.96} | 47.74 ^{41.01} | 34.10 ^{70.12} | 69.31 ^{41.24} | 45.76 ^{42.38} | 18.95 ^{43.79} | 21.23 ^{71.66} | 50.39 ^{72.34} | 41.13 ^{40.10} |
| | PPL | 59.43 ^{70.18} | 40.34 ^{70.68} | 5.91 ^{12.20} | 13.18 ^{70.43} | 48.31 ^{40.44} | 34.44 ^{70.54} | 70.72 ^{70.17} | 47.12 ^{41.02} | 18.12 ^{44.62} | 24.10 ^{74.53} | 51.26 ^{71.37} | 42.26 ^{40.47} |
| | FastV | 58.90 ^{40.35} | 43.39 ^{72.37} | 3.34 ^{44.77} | 12.37 ^{10.38} | 46.88 ^{41.87} | 32.98 ^{41.00} | 69.14 ^{41.14} | 45.42 ^{42.72} | 14.86 ^{47.78} | 23.19 ^{73.62} | 50.58 ^{70.25} | 40.64 ^{42.09} |
| | SparseVLM | 54.67 ^{44.58} | 40.00 ^{41.02} | 7.73 ^{10.38} | 12.41 ^{40.34} | 45.07 ^{43.68} | 31.98 ^{42.00} | 68.25 ^{42.30} | 45.08 ^{43.06} | 23.63 ^{70.89} | 10.17 ^{49.40} | 45.34 ^{47.29} | 38.46 ^{44.27} |
| Entropy | Random PPL | 60.31 ^{71.06} | 42.37 ^{73.15} | 6.44 ^{41.67} | 14.10 ^{71.35} | 48.09 ^{40.66} | 34.27 ^{70.29} | 72.13 ^{71.18} | 48.81 ^{40.67} | 20.09 ^{42.65} | 17.55 ^{72.02} | 54.69 ^{72.06} | 42.66 ^{40.07} |
| | PPL | 60.49 ^{71.24} | 43.73 ^{72.37} | 6.90 ^{11.21} | 14.53 ^{71.78} | 48.76 ^{70.01} | 34.88 ^{70.90} | 72.84 ^{72.29} | 47.80 ^{40.34} | 24.18 ^{71.44} | 22.80 ^{73.52} | 54.69 ^{72.06} | 44.64 ^{71.73} |
| | FastV | 58.91 ^{40.34} | 43.05 ^{72.03} | 6.37 ^{41.74} | 13.03 ^{70.28} | 47.05 ^{41.70} | 33.68 ^{40.30} | 73.90 ^{73.05} | 47.12 ^{41.02} | 24.56 ^{71.82} | 23.96 ^{74.39} | 54.67 ^{72.04} | 44.84 ^{72.11} |
| | SparseVLM | 55.20 ^{44.05} | 38.98 ^{42.04} | 7.51 ^{10.60} | 12.65 ^{40.10} | 46.14 ^{42.61} | 32.10 ^{41.88} | 68.08 ^{42.47} | 44.07 ^{44.07} | 24.87 ^{71.23} | 10.72 ^{48.85} | 47.00 ^{45.63} | 38.95 ^{43.78} |
| Q-Tuning (Ours) | | 64.20 ^{74.95} | 42.03 ^{71.01} | 10.54 ^{72.43} | 18.79 ^{76.84} | 53.12 ^{74.37} | 37.74 ^{73.76} | 71.60 ^{71.05} | 48.14 ^{70.00} | 29.34 ^{76.00} | 27.75 ^{78.18} | 57.78 ^{75.15} | 46.92 ^{74.19} |
| Full Dataset | | 61.55 | 42.37 | 8.64 | 13.80 | 50.45 | 35.36 | 71.25 | 45.76 | 26.68 | 31.81 | 53.67 | 45.84 |
| 12.5% Samples, 70% Tokens | | | | | | | | | | | | | |
| Random | Random PPL | 59.43 | 41.02 | 6.97 | 13.64 | 47.97 | 33.81 | 71.08 | 47.46 | 24.34 | 21.64 | 53.15 | 43.53 |
| | PPL | 60.14 ^{70.71} | 43.39 ^{72.37} | 6.22 ^{40.75} | 12.18 ^{41.46} | 48.18 ^{70.21} | 34.02 ^{70.21} | 70.72 ^{70.36} | 47.80 ^{70.34} | 25.09 ^{70.75} | 21.28 ^{40.36} | 53.83 ^{70.68} | 43.74 ^{70.21} |
| | FastV | 58.20 ^{41.23} | 41.02 ^{70.00} | 6.29 ^{40.68} | 13.42 ^{40.22} | 45.32 ^{42.65} | 32.85 ^{40.96} | 70.72 ^{70.36} | 46.44 ^{41.02} | 19.56 ^{47.48} | 21.38 ^{40.26} | 53.34 ^{70.19} | 42.29 ^{41.24} |
| | SparseVLM | 54.67 ^{44.76} | 37.97 ^{43.05} | 8.04 ^{10.17} | 13.06 ^{40.58} | 43.77 ^{42.10} | 31.72 ^{42.69} | 67.72 ^{43.36} | 44.75 ^{42.71} | 23.65 ^{40.69} | 11.76 ^{49.88} | 44.90 ^{48.25} | 38.58 ^{49.45} |
| Longest | Random PPL | 59.44 ^{70.01} | 43.39 ^{72.37} | 7.35 ^{70.38} | 15.59 ^{71.95} | 50.02 ^{72.05} | 35.15 ^{71.34} | 73.37 ^{72.29} | 48.81 ^{71.35} | 27.82 ^{73.48} | 21.31 ^{40.33} | 55.77 ^{72.62} | 45.42 ^{71.89} |
| | PPL | 60.85 ^{71.42} | 43.39 ^{72.37} | 7.73 ^{70.76} | 16.21 ^{72.57} | 48.76 ^{70.10} | 35.35 ^{71.54} | 74.96 ^{73.88} | 49.83 ^{72.37} | 28.73 ^{74.39} | 21.62 ^{40.02} | 56.59 ^{73.44} | 46.35 ^{73.82} |
| | FastV | 59.44 ^{70.01} | 42.71 ^{71.69} | 6.29 ^{40.68} | 14.53 ^{70.89} | 47.46 ^{40.51} | 34.09 ^{70.28} | 74.07 ^{72.99} | 49.83 ^{72.37} | 24.18 ^{40.66} | 25.74 ^{74.10} | 55.86 ^{72.71} | 45.94 ^{72.41} |
| | SparseVLM | 54.85 ^{44.58} | 37.97 ^{43.05} | 7.05 ^{70.08} | 11.20 ^{42.44} | 44.16 ^{43.81} | 31.04 ^{42.77} | 69.14 ^{41.94} | 44.75 ^{42.71} | 31.01 ^{76.67} | 6.25 ^{43.63} | 52.94 ^{40.21} | 40.82 ^{41.51} |
| InfoBatch | Random PPL | 59.26 ^{40.17} | 42.37 ^{73.15} | 6.22 ^{40.75} | 16.10 ^{72.46} | 47.72 ^{40.25} | 34.33 ^{70.52} | 70.19 ^{40.69} | 47.80 ^{70.34} | 20.77 ^{43.57} | 19.03 ^{42.61} | 52.13 ^{41.02} | 41.98 ^{41.55} |
| | PPL | 60.49 ^{71.06} | 39.32 ^{41.70} | 7.56 ^{41.21} | 14.47 ^{70.83} | 48.06 ^{70.09} | 33.62 ^{70.19} | 70.72 ^{70.36} | 46.44 ^{41.02} | 19.03 ^{45.31} | 23.20 ^{71.56} | 51.75 ^{74.10} | 42.23 ^{41.30} |
| | FastV | 58.55 ^{40.88} | 43.39 ^{72.37} | 5.53 ^{41.44} | 13.13 ^{40.51} | 47.64 ^{40.63} | 33.65 ^{40.16} | 69.49 ^{41.59} | 43.39 ^{40.47} | 16.68 ^{71.67} | 25.27 ^{73.63} | 51.47 ^{41.68} | 41.16 ^{42.27} |
| | SparseVLM | 56.61 ^{42.82} | 38.31 ^{42.71} | 5.76 ^{41.21} | 12.47 ^{41.17} | 44.49 ^{43.48} | 31.53 ^{42.28} | 68.25 ^{42.83} | 44.41 ^{43.05} | 23.73 ^{40.61} | 9.07 ^{41.57} | 45.73 ^{47.42} | 38.24 ^{45.29} |
| Entropy | Random PPL | 61.02 ^{71.59} | 43.05 ^{72.03} | 7.66 ^{70.69} | 14.11 ^{70.47} | 48.44 ^{40.47} | 34.86 ^{71.65} | 73.37 ^{72.29} | 49.83 ^{72.37} | 23.05 ^{41.29} | 16.52 ^{45.12} | 55.18 ^{72.83} | 43.59 ^{70.08} |
| | PPL | 61.02 ^{71.59} | 43.39 ^{72.37} | 6.97 ^{70.00} | 14.94 ^{71.30} | 48.94 ^{70.97} | 35.05 ^{71.24} | 73.02 ^{71.94} | 47.46 ^{70.37} | 24.03 ^{70.68} | 22.82 ^{74.21} | 54.89 ^{71.74} | 44.45 ^{70.92} |
| | FastV | 58.73 ^{40.70} | 43.39 ^{72.37} | 6.14 ^{40.83} | 14.23 ^{70.59} | 47.03 ^{40.84} | 33.90 ^{70.69} | 74.07 ^{72.99} | 50.85 ^{73.39} | 24.94 ^{70.60} | 23.79 ^{72.15} | 55.94 ^{72.72} | 45.92 ^{72.39} |
| | SparseVLM | 54.85 ^{44.58} | 37.29 ^{43.73} | 6.52 ^{40.45} | 12.73 ^{40.91} | 46.24 ^{41.73} | 31.53 ^{42.28} | 68.08 ^{43.00} | 44.41 ^{43.05} | 26.38 ^{72.04} | 11.06 ^{49.58} | 46.68 ^{46.47} | 39.32 ^{42.41} |
| Q-Tuning (Ours) | | 64.37 ^{74.94} | 42.37 ^{71.35} | 10.84 ^{73.87} | 17.63 ^{73.99} | 52.17 ^{74.20} | 37.48 ^{73.67} | 71.78 ^{70.70} | 48.14 ^{70.68} | 30.33 ^{76.00} | 28.59 ^{76.95} | 57.93 ^{74.78} | 47.35 ^{73.82} |
| Full Dataset | | 61.55 | 42.37 | 8.64 | 13.80 | 50.45 | 35.36 | 71.25 | 45.76 | 26.68 | 31.81 | 53.67 | 45.84 |
| 25% Samples, 50% Tokens | | | | | | | | | | | | | |
| Random | Random PPL | 60.32 | 41.69 | 5.76 | 13.43 | 48.41 | 33.92 | 70.19 | 46.10 | 20.62 | 24.07 | 53.74 | 42.95 |
| | PPL | 60.32 ^{70.00} | 42.03 ^{70.34} | 5.71 ^{71.75} | 15.94 ^{72.51} | 48.58 ^{40.76} | 34.87 ^{70.95} | 69.66 ^{40.53} | 47.46 ^{71.36} | 19.86 ^{40.76} | 19.51 ^{44.56} | 53.74 ^{70.00} | 42.05 ^{40.90} |
| | FastV | 59.08 ^{41.24} | 41.69 ^{70.00} | 3.56 ^{42.20} | 12.78 ^{40.65} | 46.50 ^{42.81} | 32.54 ^{41.38} | 71.78 ^{71.59} | 47.12 ^{41.02} | 15.77 ^{44.85} | 26.97 ^{72.90} | 50.84 ^{42.90} | 42.50 ^{40.45} |
| | SparseVLM | 54.50 ^{45.82} | 38.64 ^{43.05} | 6.44 ^{70.68} | 12.04 ^{41.39} | 44.79 ^{43.62} | 31.28 ^{42.64} | 67.55 ^{42.64} | 46.44 ^{70.34} | 24.41 ^{73.79} | 11.80 ^{41.27} | 48.14 ^{45.60} | 39.67 ^{43.38} |
| Longest | Random PPL | 61.20 ^{70.88} | 42.03 ^{70.34} | 7.88 ^{72.12} | 15.40 ^{71.97} | 48.29 ^{40.12} | 34.96 ^{71.04} | 73.54 ^{73.38} | 48.14 ^{72.04} | 23.73 ^{73.11} | 26.34 ^{72.27} | 54.06 ^{70.32} | 45.16 ^{72.21} |
| | PPL | 60.85 ^{70.53} | 43.39 ^{71.70} | 7.20 ^{71.44} | 13.88 ^{40.45} | 48.48 ^{70.07} | 34.76 ^{71.04} | 72.31 ^{72.12} | 48.14 ^{72.04} | 24.34 ^{73.72} | 24.84 ^{72.03} | 55.22 ^{71.48} | 44.77 ^{71.82} |
| | FastV | 59.08 ^{41.24} | 42.71 ^{71.02} | 5.16 ^{40.60} | 14.00 ^{70.57} | 47.47 ^{40.84} | 33.68 ^{40.24} | 72.66 ^{72.47} | 46.10 ^{70.00} | 18.88 ^{41.74} | 31.52 ^{72.75} | 52.13 ^{41.61} | 44.26 ^{71.31} |
| | SparseVLM | 56.61 ^{43.71} | 37.29 ^{41.40} | 7.58 ^{12.82} | 12.09 ^{41.34} | 44.76 ^{43.65} | 31.66 ^{42.26} | 66.84 ^{43.55} | 44.41 ^{43.69} | 24.69 ^{72.80} | 11.22 ^{42.85} | 48.47 ^{45.27} | 40.07 ^{42.88} |
| InfoBatch | Random PPL | 58.73 ^{41.59} | 40.68 ^{41.01} | 6.67 ^{70.91} | 9.95 ^{43.48} | 48.98 ^{70.57} | 33.00 ^{70.92} | 70.55 ^{70.36} | 46.44 ^{70.34} | 21.53 ^{70.91} | 23.93 ^{40.14} | 52.14 ^{41.80} | 42.92 ^{40.03} |
| | PPL | 59.96 ^{40.87} | 42.71 ^{71.02} | 6.52 ^{70.76} | 14.58 ^{71.15} | 48.47 ^{70.16} | 34.47 ^{70.55} | 71.08 ^{70.97} | 47.80 ^{71.70} | 20.62 ^{70.00} | 24.88 ^{70.81} | 51.61 ^{42.13} | 43.20 ^{40.75} |
| | FastV | 59.08 ^{41.24} | 42.37 ^{70.68} | 3.02 ^{43.73} | 11.13 ^{40.20} | 47.50 ^{40.91} | 32.63 ^{41.29} | 69.31 ^{40.88} | 44.41 ^{41.69} | 14.48 ^{41.84} | 23.63 ^{40.44} | 49.16 ^{45.88} | 40.20 ^{42.75} |
| | SparseVLM | 55.73 ^{44.59} | 39.66 ^{42.03} | 5.31 ^{40.45} | 11.66 ^{41.37} | 43.25 ^{45.16} | 31.12 ^{42.80} | 67.20 ^{41.00} | 45.76 ^{40.34</} | | | | |