
A DETAILED ALGORITHM OF TurboAttention

 770
 771
 772
 773
 774
 775
 776
 777
 778
 779
 780
 781

Algorithm 1 TurboAttention Prefill

Require: Matrices $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$ in HBM, compressed KV cache $\mathbf{K}^{q2}, \mathbf{V}^{q2}$, integer scale of KV cache $s_{\mathbf{K}^{q2}}^{int}$, floating scale $s_K, s_V, s_{\mathbf{V}^{q2}}$, zero point of KV cache $z_{\mathbf{K}^{q2}}, z_{\mathbf{V}^{q2}}$, block sizes B_c, B_r , Softmax approximate algorithm SAS.
 Divide \mathbf{Q} into $T_r = \lceil \frac{N}{B_r} \rceil$ blocks $\mathbf{Q}_1, \dots, \mathbf{Q}_{T_r}$ of size $B_r \times d$ each, and divide \mathbf{K}, \mathbf{V} in to $T_c = \lceil \frac{N}{B_c} \rceil$ blocks $\mathbf{K}_1, \dots, \mathbf{K}_{T_c}$ and $\mathbf{V}_1, \dots, \mathbf{V}_{T_c}$, of size $B_c \times d$ each.
 Divide the output $\mathbf{O} \in \mathbb{R}^{N \times d}$ into T_r blocks $\mathbf{O}_i, \dots, \mathbf{O}_{T_r}$ of size $B_r \times d$ each, and divide the logsumexp L into T_r blocks L_i, \dots, L_{T_r} of size B_r each.
for $1 \leq i \leq T_r$ **do**
 Load \mathbf{Q}_i from HBM to on-chip SRAM.
 On chip, initialize $\mathbf{O}_i^{(0)} = (0)_{B_r \times d} \in \mathbb{R}^{B_r \times d}, \ell_i^{(0)} = (0)_{B_r} \in \mathbb{R}^{B_r}, m_i^{(0)} = (-\infty)_{B_r} \in \mathbb{R}^{B_r}$.
for $1 \leq j \leq T_c$ **do**
 Load $\mathbf{K}_j, \mathbf{V}_j$ from HBM to on-chip SRAM.
 On chip, compute $s_{\mathbf{Q}_i} = \frac{\max(\text{abs}(\mathbf{Q}_i))}{119}, \mathbf{Q}_i^{q1} = \left\lceil \frac{\mathbf{Q}_i}{s_{\mathbf{Q}_i}} \right\rceil, s_{\mathbf{K}_j} = \frac{\max(\text{abs}(\mathbf{K}_j))}{119}, \mathbf{K}_j^{q1} = \left\lceil \frac{\mathbf{K}_j}{s_{\mathbf{K}_j}} \right\rceil, s_{\mathbf{V}_j} = \frac{\max(\text{abs}(\mathbf{V}_j))}{119}, \mathbf{V}_j^{q1} = \left\lceil \frac{\mathbf{V}_j}{s_{\mathbf{V}_j}} \right\rceil$.
 On chip, compute $\mathbf{S}_i^{(j)} = s_{\mathbf{Q}_i} \cdot s_{\mathbf{K}_j} \cdot \mathbf{Q}_i^{q1} \mathbf{K}_j^{q1} \in \mathbb{R}^{B_r \times B_c}$.
 On chip, compute $m_i^{(j)} = \max(m_i^{(j-1)}, \text{rowmax}(\mathbf{S}_i^{(j)})) \in \mathbb{R}^{B_r}, \tilde{\mathbf{P}}_i^{(j)} = SAS(\mathbf{S}_i^{(j)} - m_i^{(j)}) \in \mathbb{R}^{B_r \times B_c}$ (pointwise), $\ell_i^{(j)} = SAS(m_i^{j-1} - m_i^{(j)}) \ell_i^{(j-1)} + \text{rowsum}(\tilde{\mathbf{P}}_i^{(j)}) \in \mathbb{R}^{B_r}$.
 On chip, compute $s_{\tilde{\mathbf{P}}_i^{(j)}} = \frac{\max(\text{abs}(\tilde{\mathbf{P}}_i^{(j)}))}{119}, Q(\tilde{\mathbf{P}}_i^{(j)}) = \left\lceil \frac{\tilde{\mathbf{P}}_i^{(j)}}{s_{\tilde{\mathbf{P}}_i^{(j)}}} \right\rceil$,
 On chip, compute $\mathbf{O}_i^{(j)} = \text{diag}(SAS(m_i^{(j-1)} - m_i^{(j)}))^{-1} \mathbf{O}_i^{(j-1)} + s_{\tilde{\mathbf{P}}_i^{(j)}} \cdot s_{\mathbf{V}_j} \cdot Q(\tilde{\mathbf{P}}_i^{(j)}) \mathbf{V}_j^{q1}$.
 On chip, compute $s_{\mathbf{K}^{q2}}^{int} = \left\lceil \frac{\max(\mathbf{K}_j^{q2}) - \min(\mathbf{K}_j^{q2})}{2^{bit-1}} \right\rceil, z_{\mathbf{K}^{q2}}^{int} = \left\lceil \frac{\min(\mathbf{K}_j^{q2})}{s_{\mathbf{K}^{q2}}^{int}} \right\rceil, \mathbf{K}_j^{q2} = \left[\frac{\mathbf{K}_j^{q1}}{s_{\mathbf{K}^{q2}}^{int}} - z_{\mathbf{K}^{q2}}^{int} \right]$ (channelwise)
 On chip, compute $s_{\mathbf{V}^{q2}}^{int} = \left\lceil \frac{\max(\mathbf{V}_j^{q2}) - \min(\mathbf{V}_j^{q2})}{2^{bit-1}} \right\rceil, z_{\mathbf{V}^{q2}}^{int} = \left\lceil \frac{\min(\mathbf{V}_j^{q2})}{s_{\mathbf{V}^{q2}}^{int}} \right\rceil, \mathbf{V}_j^{q2} = \left[\frac{\mathbf{V}_j^{q1}}{s_{\mathbf{V}^{q2}}^{int}} - z_{\mathbf{V}^{q2}}^{int} \right]$ (channelwise)
 Write \mathbf{K}_j^{q2} to HBM as the j -th block of \mathbf{K}^{q2} .
 Write \mathbf{V}_j^{q2} to HBM as the j -th block of \mathbf{V}^{q2} .
 Write $s_{\mathbf{K}^{q2}}, s_{\mathbf{V}^{q2}}$ to HBM as the j -th block of $s_{\mathbf{K}^{q2}}, s_{\mathbf{V}^{q2}}$.
 Write $s_{\mathbf{K}_j}, s_{\mathbf{V}_j}$ to HBM as the j -th block of $s_{\mathbf{K}}, s_{\mathbf{V}}$.
 Write $z_{\mathbf{K}^{q2}}, z_{\mathbf{V}^{q2}}$ to HBM as the j -th block of $z_{\mathbf{K}^{q2}}, z_{\mathbf{V}^{q2}}$.
end for
 On chip, compute $\mathbf{O}_i = \text{diag}(\ell_i^{(T_c)})^{-1} \mathbf{O}_i^{(T_c)}$.
 On chip, compute $L_i = m_i^{(T_c)} + \log(\ell_i^{(T_c)})$.
 Write \mathbf{O}_i to HBM as the i -th block of \mathbf{O} .
 Write L_i to HBM as the i -th block of L .
end for
 Return the output \mathbf{O} and the logsumexp L .

Algorithm 2 TurboAttention Decode

825
 826 **Require:** Matrices $\mathbf{Q} \in \mathbb{R}^{1 \times d}$ in HBM, compressed KV cache $\mathbf{K}^{q2}, \mathbf{V}^{q2} \in \mathbb{R}^{N \times d}$, integer scale of KV cache $s_{\mathbf{K}^{q2}}^{int}, s_{\mathbf{V}^{q2}}^{int}$,
 827 floating scale s_K, s_V , zero point of KV cache $z_{\mathbf{K}^{q2}}^{int}, z_{\mathbf{V}^{q2}}^{int}$, block size B_c , Softmax approximate algorithm SAS.
 828 Divide $\mathbf{K}^{q2}, \mathbf{V}^{q2}$ in to $T_c = \left\lceil \frac{N}{B_c} \right\rceil$ blocks $\mathbf{K}_1^{q2}, \dots, \mathbf{K}_{T_c}^{q2}$ and $\mathbf{V}_1^{q2}, \dots, \mathbf{V}_{T_c}^{q2}$, of size $B_c \times d$ each.
 829 Divide $s_{\mathbf{K}^{q2}}^{int}, s_{\mathbf{V}^{q2}}^{int}, z_{\mathbf{V}^{q2}}^{int}$, block size B_c into $T_c = \left\lceil \frac{N}{B_c} \right\rceil$ blocks accordingly .
 830 Divide the output $\mathbf{O} \in \mathbb{R}^{1 \times d}$ into T_c blocks $\mathbf{O}_1, \dots, \mathbf{O}_{T_c}$ of size $1 \times \frac{d}{T_c}$ each, and divide the logsumexp L into T_r blocks
 831 L_1, \dots, L_{T_r} of size 1 each.
 832 Load \mathbf{Q} from HBM to on-chip SRAM.
 833 On chip, initialize $\mathbf{O}_j = (0)_{1 \times \frac{d}{T_c}} \in \mathbb{R}^{1 \times d}, \ell_i^{(0)} = (0)_{B_c} \in \mathbb{R}^{B_c}, m_i^{(0)} = (-\infty)_{B_c} \in \mathbb{R}^{B_c}$.
 834 **for** $1 \leq j \leq T_c$ **do**
 835 Load $\mathbf{K}_j^{q2}, \mathbf{V}_j^{q2}$ from HBM to on-chip SRAM.
 836 On chip, compute $s_{\mathbf{Q}} = \frac{\max(\text{abs}(\mathbf{Q}))}{119}, \mathbf{Q}^{q1} = \left\lceil \frac{\mathbf{Q}}{s_{\mathbf{Q}}} \right\rceil, \mathbf{K}_j^{q1} = \mathbf{K}_j^{q2} \cdot s_{\mathbf{K}^{q2}}^{int} + z_{\mathbf{K}^{q2}}, \mathbf{V}_j^{q1} = \mathbf{V}_j^{q2} \cdot s_{\mathbf{V}^{q2}}^{int} + z_{\mathbf{V}^{q2}}$.
 837 On chip, compute $\mathbf{S}_j = s_{\mathbf{Q}} \cdot s_{\mathbf{K}_j} \cdot \mathbf{Q}^{q1} \mathbf{K}_j^{q1} \in \mathbb{R}^{1 \times B_c}$.
 838 On chip, compute $m_j = \max(m_{(j-1)}, \text{rowmax}(\mathbf{S}_j)) \in \mathbb{R}^1, \tilde{\mathbf{P}}^{(j)} = SAS(\mathbf{S}_j - m_j) \in \mathbb{R}^{1 \times B_c}$ (pointwise), $\ell^{(j)} = SAS(m^{(j-1)} - m^{(j)}) \ell^{(j-1)} + \text{rowsum}(\tilde{\mathbf{P}}^{(j)}) \in \mathbb{R}^1$.
 839 On chip, compute $s_{\tilde{\mathbf{P}}^{(j)}} = \frac{\max(\text{abs}(\tilde{\mathbf{P}}^{(j)}))}{119}, Q(\tilde{\mathbf{P}}^{(j)}) = \left\lceil \frac{\tilde{\mathbf{P}}^{(j)}}{s_{\tilde{\mathbf{P}}^{(j)}}} \right\rceil$,
 840 On chip, compute $\mathbf{O}^{(j)} = \text{diag}(SAS(m^{(j-1)} - m^{(j)}))^{-1} \mathbf{O}^{(j-1)} + s_{\tilde{\mathbf{P}}^{(j)}} \cdot s_{\mathbf{V}_j} \cdot Q(\tilde{\mathbf{P}}^{(j)}) \mathbf{V}_j^{q1}$.
 841 **end for**
 842 On chip, compute $\mathbf{O} = \text{diag}(\ell^{(T_c)})^{-1} \mathbf{O}^{(T_c)}$.
 843 On chip, compute $L = m^{(T_c)} + \log(\ell^{(T_c)})$.
 844 Write \mathbf{O} to HBM.
 845 Return the output \mathbf{O} and the logsumexp L .

B DETAILED ALGORITHM OF SAS

Algorithm 3 Sparsity-based Softmax Approximation

855
 856 **Require:** Input matrix $X \in \mathbb{R}^{N \times d}$, integer threshold n_r , look-up table T with $T[n_r + 1]$ set to 0, polynomial approximator
 857 function $POLY$.
 858
 859 **Step 1: Normalize input values**
 860 $x_{\max} = \text{rowmax}(X)$ // Max value per row
 861 $X = X - x_{\max}$
 862
 863 **Step 2: Apply sparsity threshold**
 864 $X[X < n_r] = n_r + 1$
 865
 866 **Step 3: Separate integer and decimal parts**
 867 $X = X_{\text{int}} + X_{\text{decimal}}$
 868
 869 **Step 4: Approximate exponential values**
 870 $X_{\text{lut}} = T[X_{\text{int}}]$ // Lookup for integer part
 871 $X_{\text{approx}} = POLY(X_{\text{decimal}})$ // Polynomial approx. for decimal part
 872 $X = X_{\text{lut}} \times X_{\text{approx}}$
 873
 874 **Step 5: Normalize with row-wise sum**
 875 $X_{\text{sum}} = \text{rowsum}(X)$
 876 $X = \frac{X}{X_{\text{sum}}}$
 877 Return X

C DETAILED ILLUSTRATION OF FlashQ AND TurboAttention

878 In this section, we separately evaluate the accuracy degradation caused by *FlashQ* and *SAS*. Results are shown in [Table 4](#). It
 879 demonstrated that both techniques are near loss-less.

880
 881
 882
 883
 884
 885
 886
 887
 888
 889
 890
 891
 892
 893
 894
 895
 896
 897
 898
 899
 900
 901
 902
 903
 904
 905
 906
 907
 908
 909
 910
 911
 912
 913
 914
 915
 916
 917
 918
 919
 920
 921
 922
 923
 924
 925
 926
 927
 928
 929
 930
 931
 932
 933
 934

Table 4. Separate discussion about accuracy degradation of FlashQ and SAS

Model	Dataset	Method	Acc
LLaMA3-8B-inst	AQuA	FP16	50.79
LLaMA3-8B-inst	AQuA	FlashQ-4bit	49.60
LLaMA3-8B-inst	AQuA	SAS	50.12
LLaMA3-8B-inst	AQuA	FlashQ-4bit + SAS	48.03