660 A RAGGED BATCHING IN DECODE.

Ragged Batching in Decode. For the purpose of our evaluations, we quantify the heterogeneity of a ragged batch as the ratio of average context length to the maximum context length present in the batch. Figure 11 shows the speedup of LA over FD. We observe that as the heterogeneity of batch increases, LA delivers a higher speedup because of better distribution of work across SMs.



Figure 11. Speedup offered by LA over FD at different batch sizes with heterogeneous context lengths. Batch context ratio(%) shows the
 ratio of average context length over maximum context length

_

_

Alg	orithm 1 LeanTile() for a sequence of lean tile iterations.
1:	function LeanTile(tile_idx, iter_begin, iter_end)
2:	$_shared_{-}O_{acc}[T_m,d]$
3:	$_$ shared $_Q_f[T_m, d]$
4:	$_$ shared $_K_f[T_n, d]$
5:	$_$ shared $_V_f[T_n, d]$
6:	$_$ shared $_ m[T_m, 1]$
7:	$_$ shared $_$ $l[T_m, 1]$
8:	Initialize O_{acc} to $(0)_{T_m \times d} \in \mathbb{R}^{T_m \times d}$ in SMEM.
9:	Initialize m to $(-\infty)_{T_m \times 1}$ and l to $(0)_{T_m \times 1} \in \mathbb{R}^{T_m \times 1}$ in SMEM.
0:	$mm = T_m \times (\text{tile_idx} / 1)$
1:	$nn = d \times (\text{tile_idx \% 1})$
2:	Perform lean tile iterations for this output tile.
3:	for $iter = iter_begin$ to $iter_end$ do
4:	$kk = iter \times T_n$
5:	load fragments from GMEM to SMEM
6:	$Q_f = LoadFragment(Q, mm, nn)$
7:	$K_f = LoadFragment(K, nn, kk)$
8:	$V_f = LoadFragment(V, nn, kk)$
9:	Compute on chip:
20:	$S_f = Q_f K_f$ where $S_f \in R^{T_m \times T_n}$
1:	$m^{new} = max(m, rowmax(S_f))$
22:	$P_f = exp(S_f - m^{new})$ where $P_f \in R^{T_m \times T_n}$
23:	$l^{new} = e^{m - m^{new}} l + rowsum_{new}(P_f)$
24:	$O_{acc} = P_f V_f + diag(e^{m - m^{n \times a}}) O_{acc}$
25:	$l = l^{new}, m = m^{new}$
26:	end for
27:	return O_{acc} , l, m
28:	end function

C PARTITION STRATEGY

Without loss of generality, we describe this process of reduction to obtain one row vector of the attention score matrix **S**, of the form $[\mathbf{S}^{(x)} \ \mathbf{S}^{(y)}]$ consisting of some unequal length vectors $\mathbf{S}^{(x)}, \mathbf{S}^{(y)}$ where $\mathbf{S}^{(x)} \in \mathbb{R}^{1 \times B_c^{(x)}}$ and $\mathbf{S}^{(y)} \in \mathbb{R}^{1 \times B_c^{(y)}}$, where 1 is the query length and $B_c^{(x)}$ and $B_c^{(y)}$ are the unequal context lengths. The vectors $\mathbf{S}^{(x)}$ and $\mathbf{S}^{(y)}$ were computed from $\mathbf{Q} \times (\mathbf{K}^{(x)})^T$ and $\mathbf{Q} \times (\mathbf{K}^{(y)})^T$ (illustrated Figure in Appendix C). Note that, to generalize this procedure for blocks of any size, the context length of $\mathbf{K}^{(x)}$ and $\mathbf{K}^{(y)}$ are $B_c^{(x)}$ and $B_c^{(y)}$ and are not necessarily equal.



Figure 12. Illustrative diagram showing LeanAttention's partitioning strategy with two differently sized work volumes of a head assigned to different CTAs. The un-scaled outputs are independently computed and re-scaled later in a reduction operation. Note that this can be generalized to any arbitrary-sized work volume split.

D PROOF OF ASSOCIATIVITY

827 Proving that f(f(x, y), z) = f(x, y, z):

$$f(x,y) = \tilde{\mathbf{O}}^{(x,y)}$$

$$f(x,y) = \tilde{\mathbf{O}}^{(x,y)}$$

$$f(f(x,y),z) = \text{diag}(e^{m^{(x,y)} - m^{((x,y),z)}})\tilde{\mathbf{O}}^{(x,y)}$$

$$+ \text{diag}(e^{m^{(z)} - m^{((x,y),z)}})\tilde{\mathbf{O}}^{(z)}$$

$$= \text{diag}(e^{m^{(x)} - m^{(x,y,z)}})\tilde{\mathbf{O}}^{(z)}$$

$$= \text{diag}(e^{m^{(x)} - m^{(x,y)}})\tilde{\mathbf{O}}^{(z)}$$

$$= \text{diag}(e^{m^{(x)} - m^{(x,y)}})\tilde{\mathbf{O}}^{(y)}$$

$$+ \text{diag}(e^{m^{(x)} - m^{(x,y,z)}})\tilde{\mathbf{O}}^{(z)}$$

$$= \text{diag}(e^{m^{(x)} - m^{(x,y,z)}})\tilde{\mathbf{O}}^{(z)}$$

$$+ \text{diag}(e^{m^{(x)} - m^{(x,y,z)}})\tilde{\mathbf{O}}^{(z)}$$

$$= \text{diag}(e^{m^{(x)} - m^{(x,y,z)}})\tilde{\mathbf{O}}^{(z)}$$

$$= \text{diag}(e^{m^{(x)} - m^{(x,y,z)}})\tilde{\mathbf{O}}^{(z)}$$

$$= \text{diag}(e^{m^{(x)} - m^{(x,y,z)}})\tilde{\mathbf{O}}^{(z)}$$

$$= \tilde{\mathbf{O}}^{(x,y,z)} = f(x,y,z)$$

Therefore, f(f(x,y),z) = f(x,y,z) and similarly $\ell^{((x,y),z)} = \ell^{(x,y,z)}$. For brevity, we omit the proof of f(x, f(y,z)) = f(x,y,z), but it can deduced in a similar manner.

This associativity of softmax re-scaling is leveraged in LeanAttention to concurrently calculate the "partial" outputs produced from unequally sized KV blocks and then "reduce" them to obtain exact attention.

E **DETAILED ALGORITHM OF LEANATTENTION** 880 881 882 Algorithm 2 Lean Attention 883 1: $_$ shared $_O[T_m, d]$ 884 2: $_shared_{-}m[T_m, 1]$ 885 3: _shared_ $l[T_m, 1]$ 886 4: Number of output tiles: $C_m = \lceil N_q/T_m \rceil$ 887 5: Number of iterations for each output tile: $C_n = \lfloor N_k/T_n \rfloor$ 888 6: Total number of iterations: $I = C_m C_n$ 889 7: Number of iterations per CTA: $I_G = I/G$ 890 8: fork CTA_g in G do 891 9: cta_start = $g I_G$ and cta_end = cta_start + I_G 892 10: **for** iter = cta_start **to** cta_end **do** 893 Index of current output tile: tile_idx = iter / C_n 11: 894 12: tile_iter = tile_idx $\times C_n$ 895 tile_iter_end = tile_iter + C_n 13: 896 local_iter = iter - tile_iter 14: 897 local_iter_end = min(tile_iter_end, cta_end) - tile_iter 15: 898 16: O, m, l = LeanTile(tile_idx, local_iter, local_iter_end) 899 17: host-block if: iter == tile_iter 900 18: finishing-block if: cta_end >= tile_iter_end 901 if !(host-block) then 19: 902 20: StorePartials(Op[g], O) 903 21: StorePartials(mp[g], m) 904 22: StorePartials(lp[g], l) 905 23: Signal(flags[g]) 906 24: else 907 25: if !(finishing-block) then 908 26: $last_cta = tile_iter_end / C_n$ 909 for cta = (g + 1) to $last_cta$ do 27: 910 28: Wait(flags[cta]) 911 29: $m_{cta} = \text{LoadPartials}(\text{mp[cta]})$ 912 $l_{cta} = \text{LoadPartials}(lp[cta])$ 30: 913 $O_{cta} = \text{LoadPartials}(\text{Op[cta]})$ 31: 914 $m^{new} = max(m_{cta}, m)$ 32: $\begin{array}{l} l^{new} = e^{m_{cta} - m^{new}} l_{cta} + e^{m - m^{new}} l \\ O^{new} = e^{m_{cta} - m^{new}_i} O_{cta} + e^{m - m^{new}_i} O \end{array}$ 915 33: 916 34: 917 Update $m = m_i^{new}, l = l_i^{new}$ 35: 918 end for 36: 919 end if 37: 920 Write $O = diag(l)^{-1}O$ to GMEM. 38: 921 39: Write L = m + log(l) to GMEM. 922 40: end if 923 41: iter = tile_iter_end 924 42: end for 925 43: **join** 926 927 928 929

935 F ENERGY CONSUMPTION EVALUATION

Fixed-split partitioning results in imbalanced workloads across the SMs, leaving many of them idle during the final stages
of computation. This inefficiency makes fixed-split attention mechanisms energy-inefficient. As shown in Figure 13, the
disparity in energy consumption between FlashDecoding, FlashInfer, and LeanAttention increases as context lengths grow
over 128k. LeanAttention, with its well-balanced load partitioning strategy, ensures more consistent utilization of SMs,
making it significantly more energy-efficient.



Figure 13. Ratio of Energy consumed by attention kernel to energy consumed by FlashDecoding kernel of different attention mechanisms
 for batch size = 1, number of heads = 56, head dimension = 64 on a single Nvidia-A100-80GB GPU as measured using NVML APIs.