

	MM: Mat Mul MFS: mask, fill, softmax ZSC: zeros, slice, cat TD: to(device)		Mem_Selfattn		Naive_Selfattn	
k, B, Tc, Tg, H	0th forward pass: q.kt_cached, y_cached = None					
x: B * Tc * H	forward_uncached: $q.kt = q \cdot k^{T_{c+1}}$ $(\frac{q.kt}{\sqrt{d_k}}) \in O(1)$		$BK \frac{T}{K} T^2 = BT^3$ —	MM	$BK \frac{T}{K} T^2 = BT^3$ —	MM
	self.mask[:, :, T, T] + att.masked_fill + softmax		$BT^2$	MFS	$BT^2$	MFS
	y = att.v $(B \times T), (H \times \frac{T}{K})$ to_device		$BK T^2 \frac{T}{K} = BT^3$ $BKT^2$	MM TD	$BK T^2 \frac{T}{K} = BT^3$ $BKT^2$	MM TD
	1...(Tg-1)th forward pass: q.kt_cached: $BK(T-1)(T-1)$ y_cached: $BK(T-1)\frac{H}{K}$					
	q.kt zeros + $[:, :, T-1, T-1]$ $(B \times T)$		$BT^2$	ZSC		
	forward_cached: $q.kt = q \cdot k^{T_{c+1}}$ $(\frac{q.kt}{\sqrt{d_k}}) \in O(1)$ $(B \times T), (H \times \frac{T}{K}), (B \times T), (H \times \frac{T}{K})$		$BK \frac{T}{K} T = BT^2$ —	MM		
	self.mask[:, :, T, T] + att.masked_fill + softmax		$BT^2$	MFS		
	y_new = new_att.v $(B \times T), (H \times \frac{T}{K}), (B \times T), (H \times \frac{T}{K})$ cat(y_cached, y_new) to_device		$BK T \frac{T}{K} = BT^2$ $BT^2$ $BKT^2$	MM ZSC TD		

Can be sped up: Mem-selfattn:  $M = 2BT_c^3 \cdot MM + BT^2 \cdot MFS + \sum_{i=1}^{T_g-1} [2B(T_c+i)^2 \cdot MM + 2B(T_c+i)^2 \cdot ZSC + B(T_c+i)^2 \cdot MFS]$

Naive-selfattn:  $N = \sum_{i=0}^{T_g-1} [2B(T_c+i)^3 \cdot MM + B(T_c+i)^2]$

Cannot be sped up:  $C = \sum_{i=0}^{T_g-1} BK(T_c+i)^2 \cdot TD$   
(communication cost)

$$P = \frac{N}{N+C}, S = \frac{M}{N}$$

$$Slotting = \frac{1}{(1-P) + \frac{P}{S}} = \frac{1}{\frac{C}{N+C} + \frac{\frac{N}{N+C}}{\frac{M}{N}}} = \frac{1}{\frac{C}{N+C} + \frac{N^2}{N+C}} = \frac{N+C}{C + \frac{N^2}{M}}$$

Ballpark latency numbers  
MM/MFS/ZSC TD  
L1 cache : comm-cost  $\approx 10:1$   
1 TB/s (900 GB/s for Tesla V100)