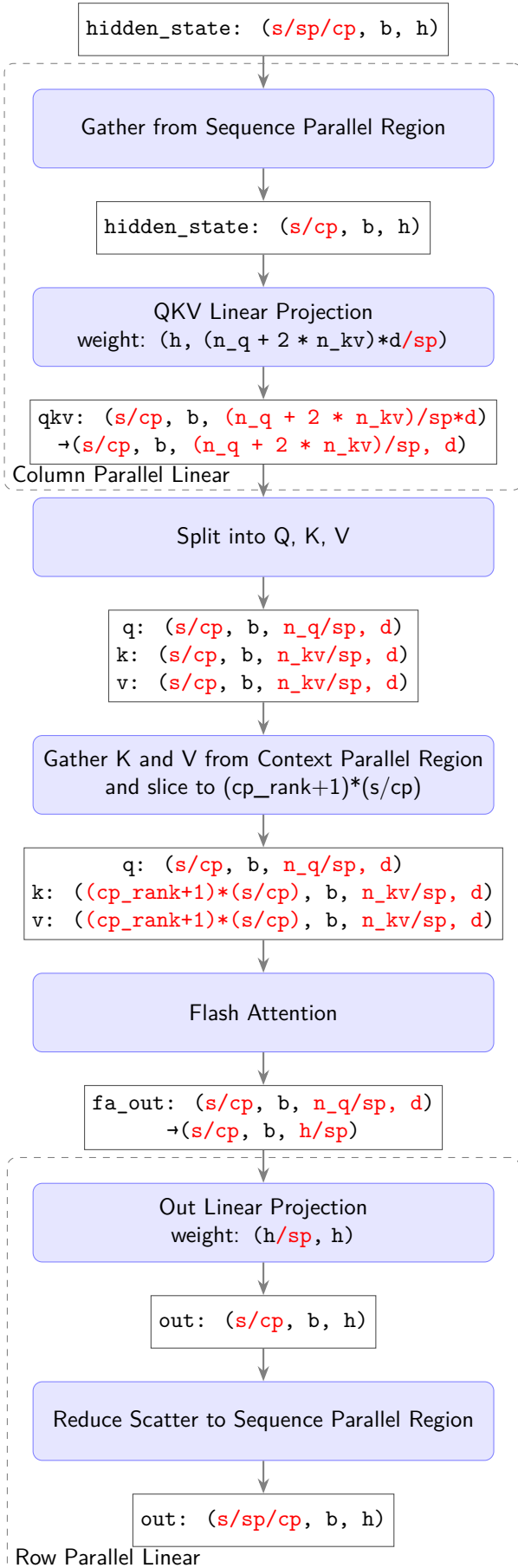


Llama Attention with Context Parallel



s: Sequence length

b: Batch size

h: Hidden size

d: Head dim

n_q: Number of Attention Heads for Q. $h = n_q * d$

n_{kv}: Number of Attention Heads for KV

sp: Size of Sequence parallel Process Group

cp: Size of Context parallel Process Group

cp_rank: Rank of Context parallel Process Group