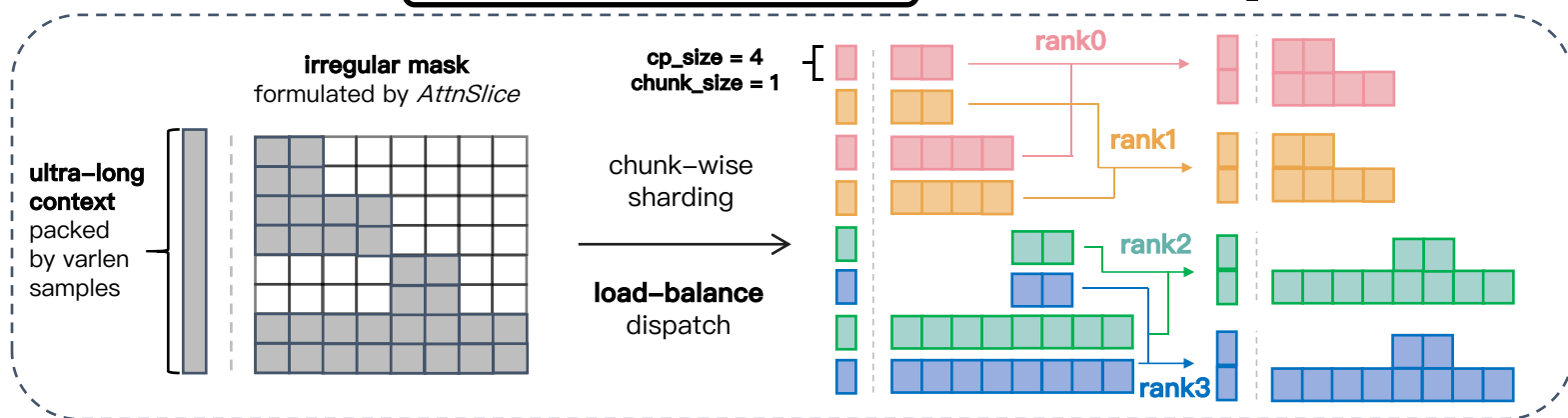


Dispatch Solver for

computation load-balance

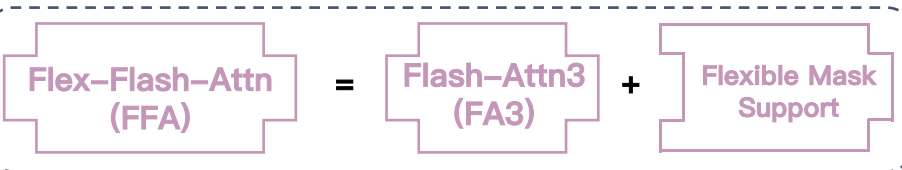
chunk_partitions = [(0,2),(1,3),(4,6),(5,7)] (2)

minimax_loads = 10



(1)

flexible and efficient attention kernels



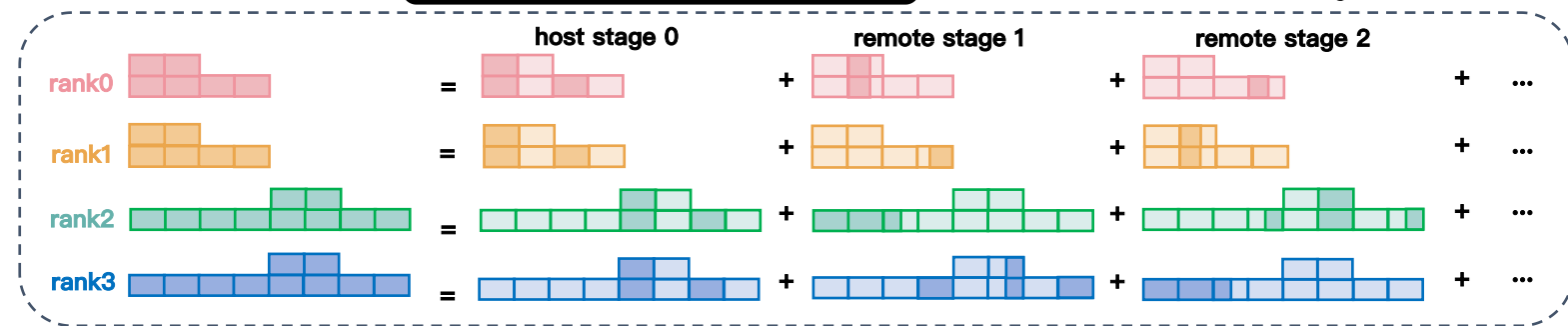
Overlap Solver for

adaptive multi-stage overlap

num_stages_fwd = 4

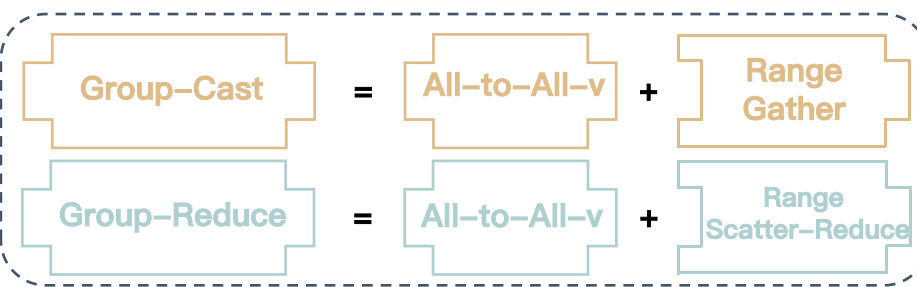
(4)

num_stages_bwd = 3



(3)

zero-redundant communication primitives



MagiAttention forward and backward timelines

(5)

