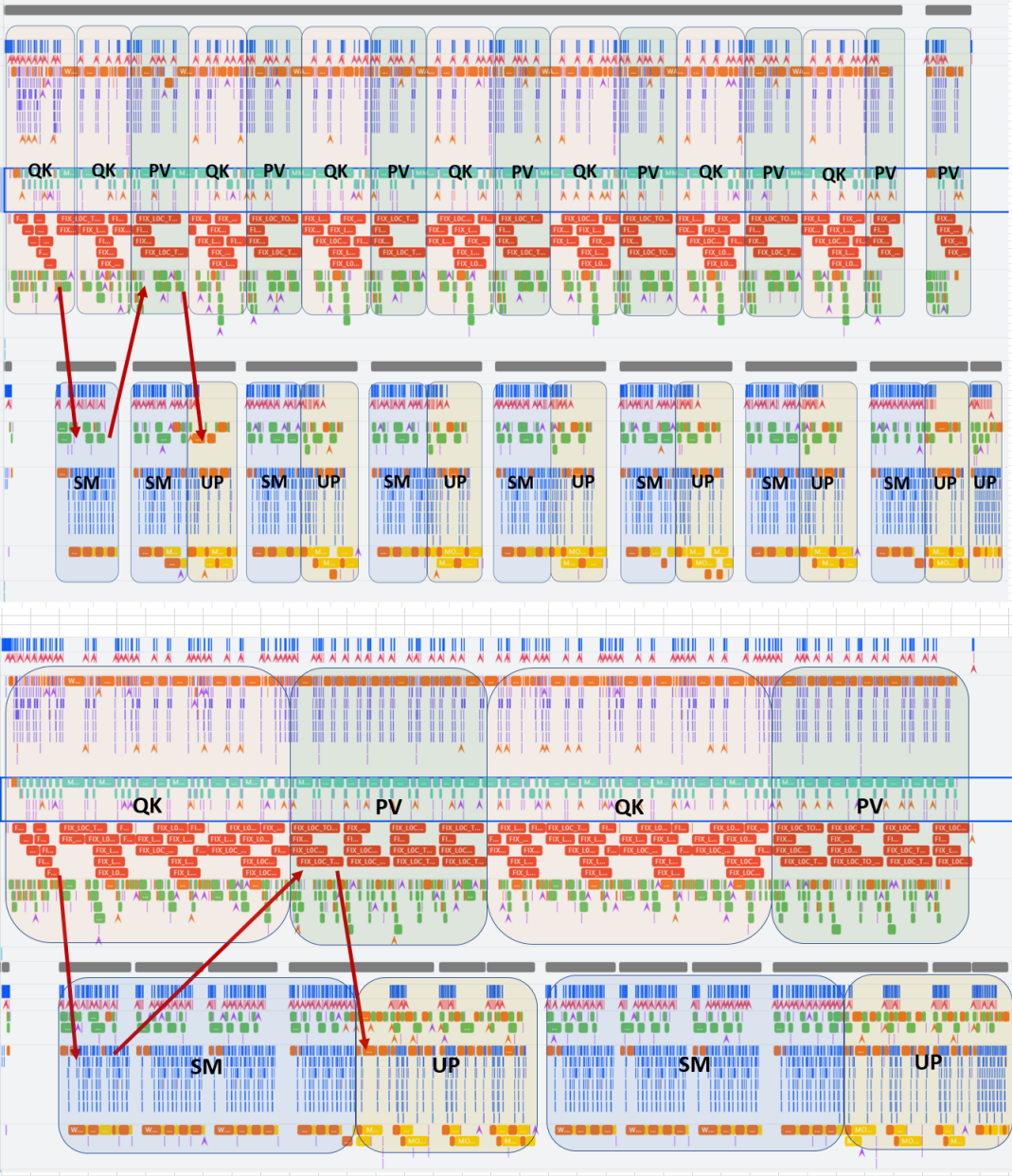


Int8 Optimizations: (Implemented from Zhentao)

(I)



Original Pipeline:

(1) chunk=1  
QK, QK, PV, QK, PV, QK, PV, PV, \_\_  
\_\_, SM, SM, UP, SM, UP, SM, UP, UP

New Pipeline:

(2) chunk=2  
QK, QK, PV, PV, QK, QK, PV, PV, \_\_  
\_\_, SM, SM, UP, UP, SM, SM, UP, UP

(3) chunk=4  
QK, QK, QK, QK, PV, PV, PV, PV ...  
\_\_, SM, SM, SM, SM, UP, UP, UP, UP...

SideNote: No need to change workspace.

SideNote: The gm workspace of S, P, Otmp would be doubled. But since in Int8, it still has the same workspace size compared to fp16.

- Observation:
- 1. The granularity of each stage is small, which requires perfect pipeline runtime to make the bubble disappear, which has less robustness.
  - 2. In the actual run(not the simulation), the SM duration is long, so there is a big bubble in the end for PV waiting for the SM.

Purpose of the New Pipeline with large chunksize: Enlarge the granularity of the pipeline between CV.

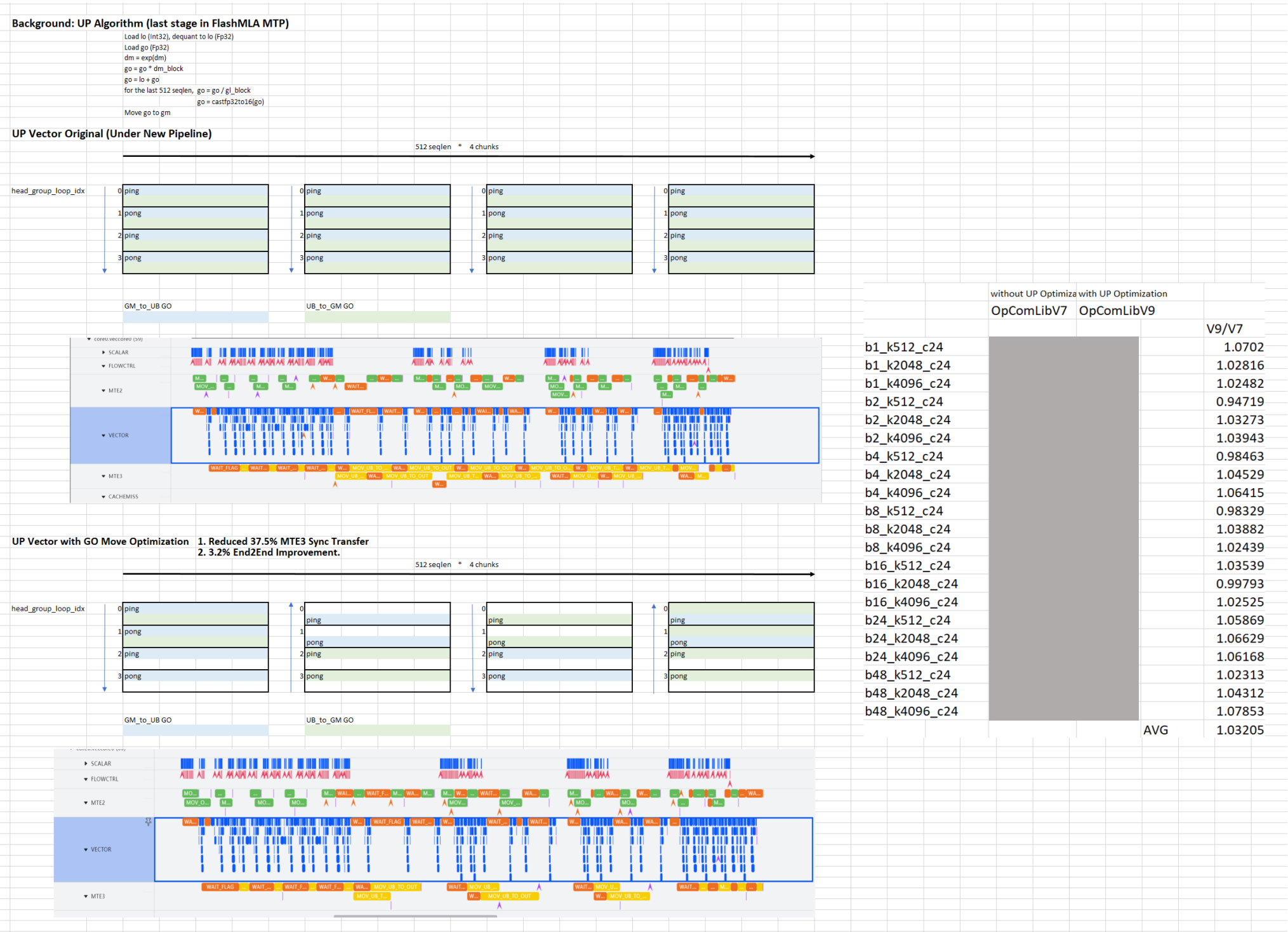
- Side Advantaged:
- 1. Utilized more L1 Buffer with pingpong
  - 2. leave the room for full cube utilization
  - 3. leave the room for UP memory transfer reduction.

- Side Disadvantaged:
- 1. Doubled the HBM usages on S, P, Otmp.

Result:

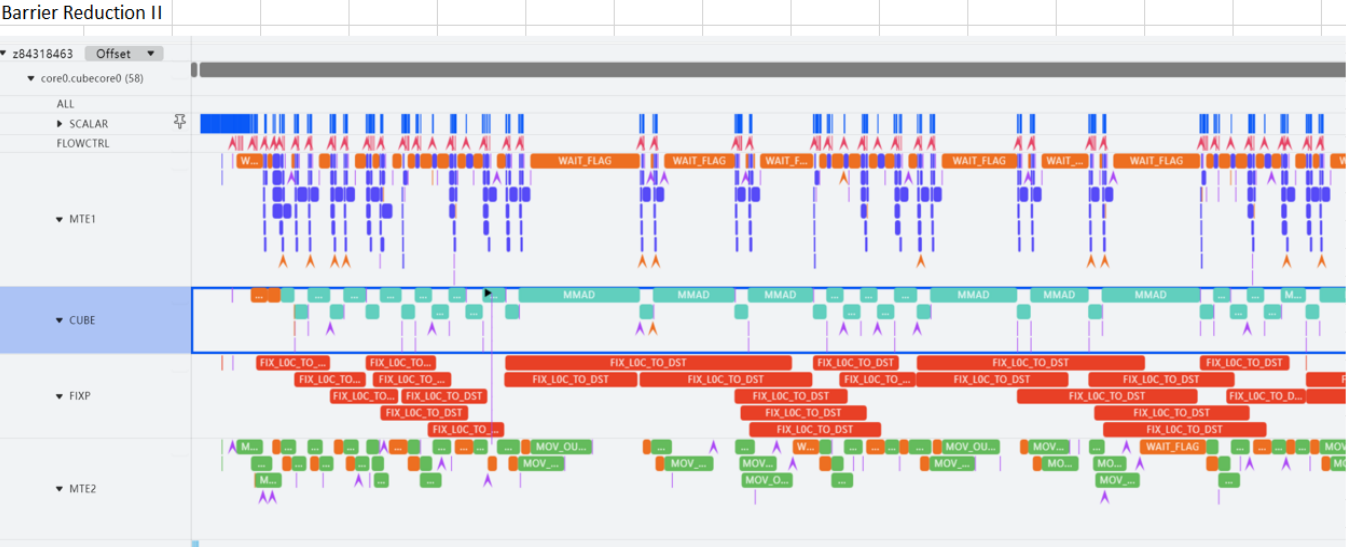
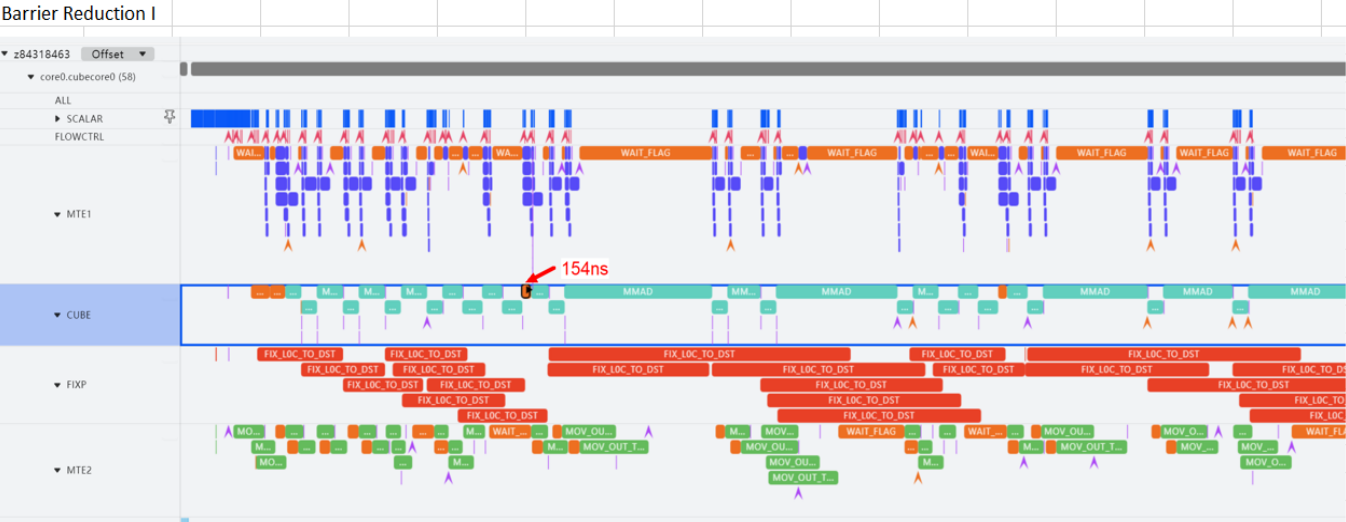
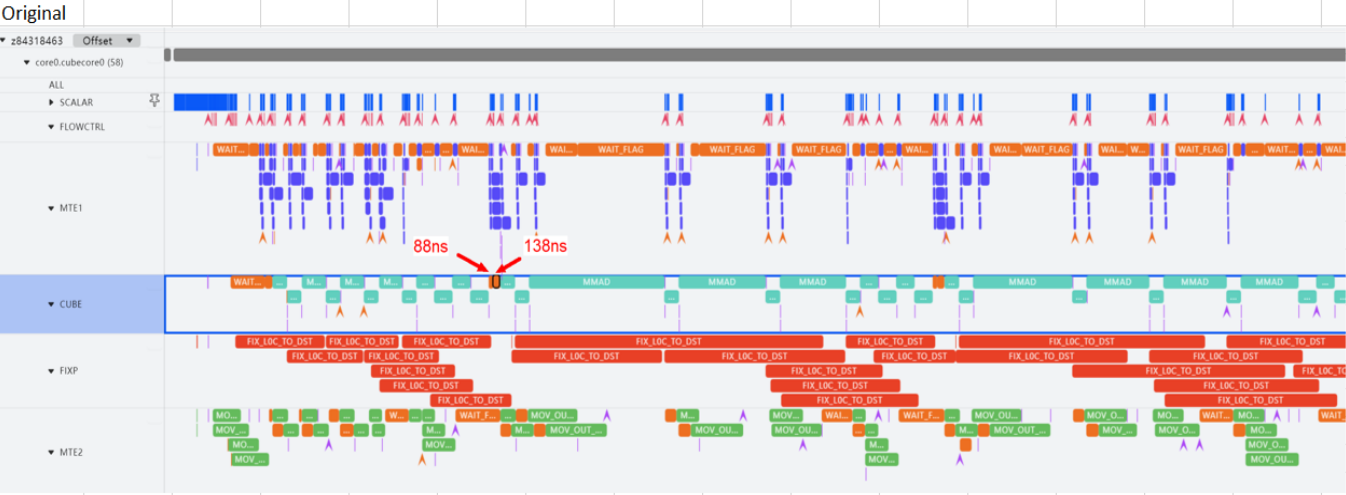
	Bs1_K4096_C24	Bs1_K2048_C24
Original Pipeline:	us	us
New Pipeline:	us	us
Improvement:	4.5%	9.2%

(II)



(III)

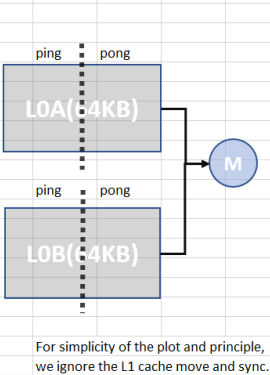
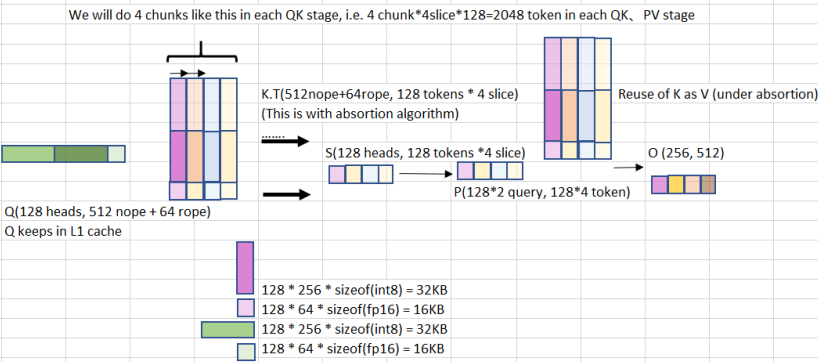
Int8 Kernel Optimization of Cube Bubble Under New Pipeline of chunksize 4



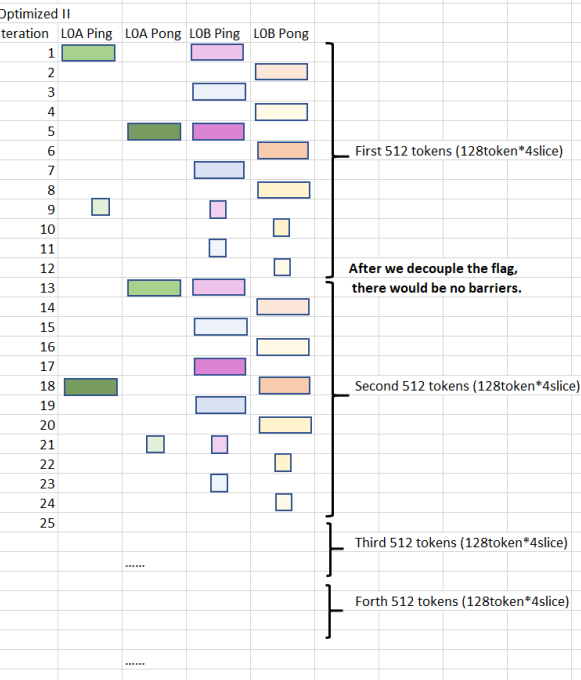
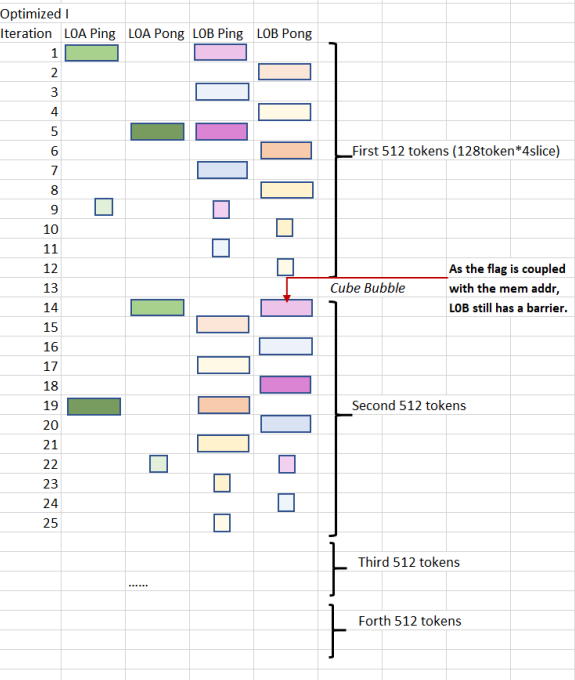
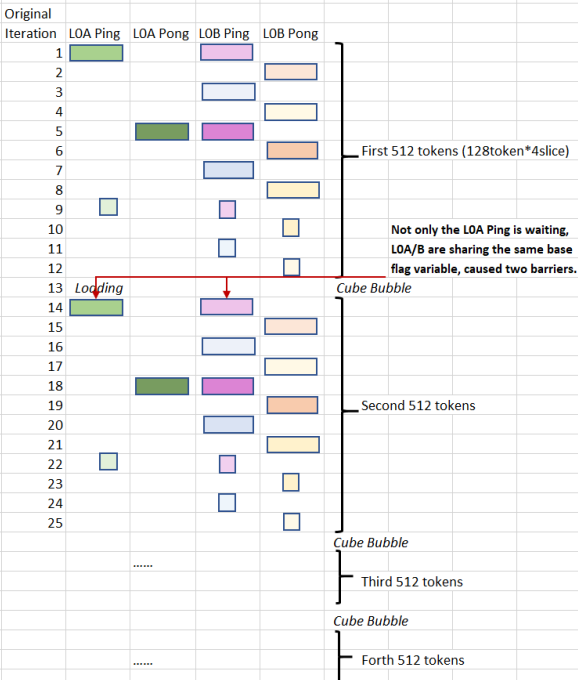
In the end, there is zero bubbles in the whole pipeline.

Reduction Strategy:

Int8 Kernel Optimization of Cube Bubble Under New Pipeline of chunksize 4



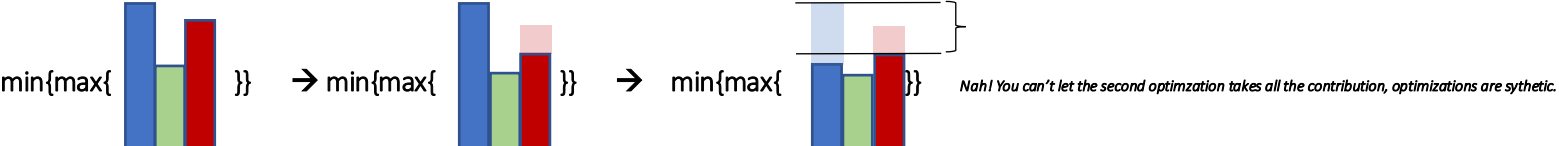
For simplicity of the plot and principle, we ignore the L1 cache move and sync.



## Int8 Overall Improvement (includes other optimization from teammates): 19%

		OpComLibV0 (us)	OpComLibV11 (us)	V11/V0
b1_k512_c24				1.07568
b1_k2048_c24				1.39753
b1_k4096_c24				1.22115
b2_k512_c24				1.10849
b2_k2048_c24				1.20156
b2_k4096_c24				1.17252
b4_k512_c24				1.125
b4_k2048_c24				1.24277
b4_k4096_c24				1.19364
b8_k512_c24				1.13337
b8_k2048_c24				1.25688
b8_k4096_c24				1.22447
b16_k512_c24				1.04
b16_k2048_c24				1.23246
b16_k4096_c24				1.2045
b24_k512_c24				1.10615
b24_k2048_c24				1.23452
b24_k4096_c24				1.22962
b48_k512_c24				1.13204
b48_k2048_c24				1.2534
b48_k4096_c24				1.25591
			AVG	1.1924

### Side Explanation:



The real case is more complicated since each component can not be separated independently

## Other Optimizations:

Some solid optimization without End2End Improvement (but definitely no harm):

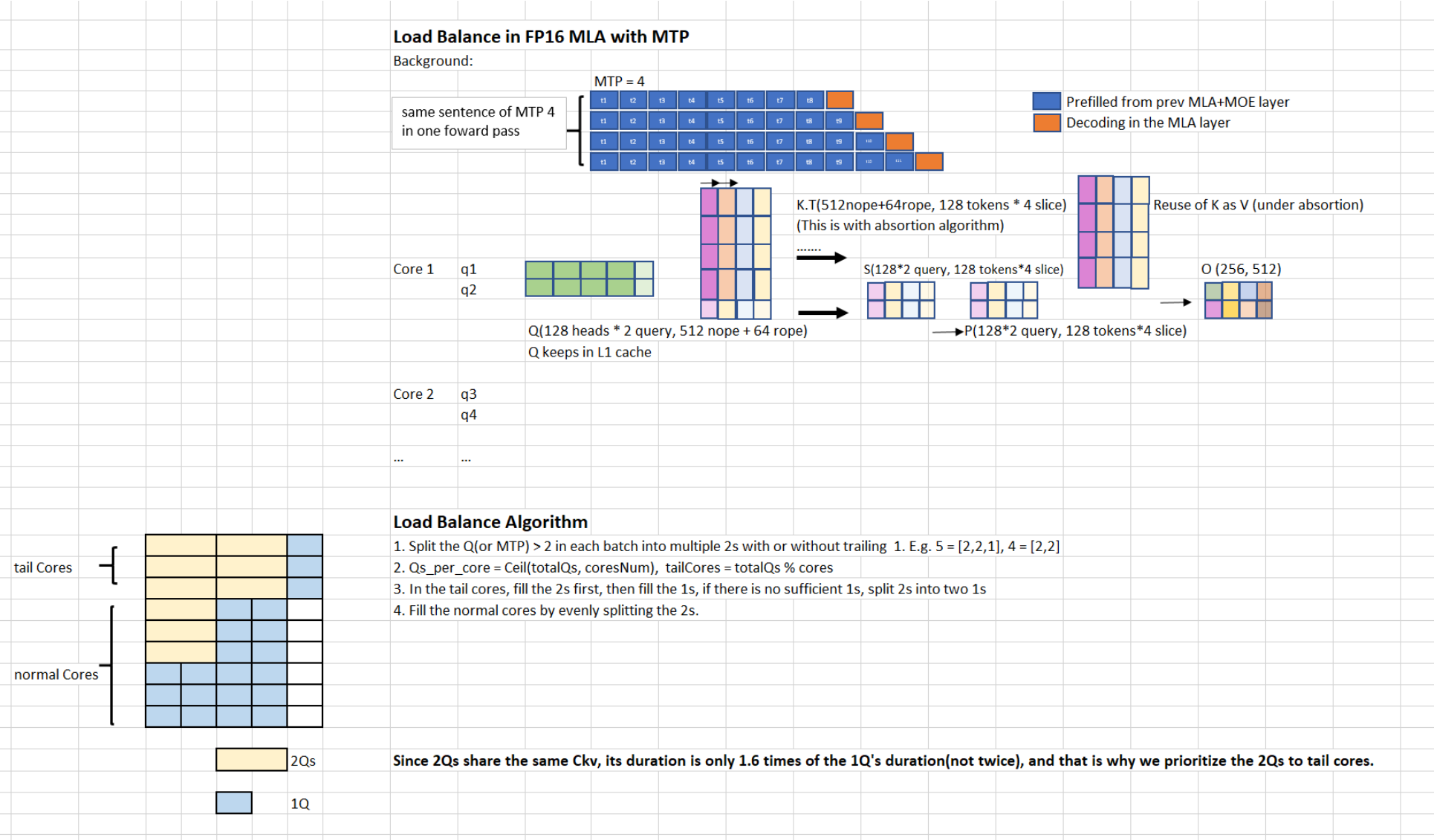
1. S/P Layout optimization in Int8 case.
2. Setting two stages of cross-core synchronization for QK stage, one for nope(S\_nope int32), one for rope(S\_rope fp32). So once the S\_nope is calculated, it can move to vector core and do dequantization.

Other concerns: multi-platform; code vulnerability.

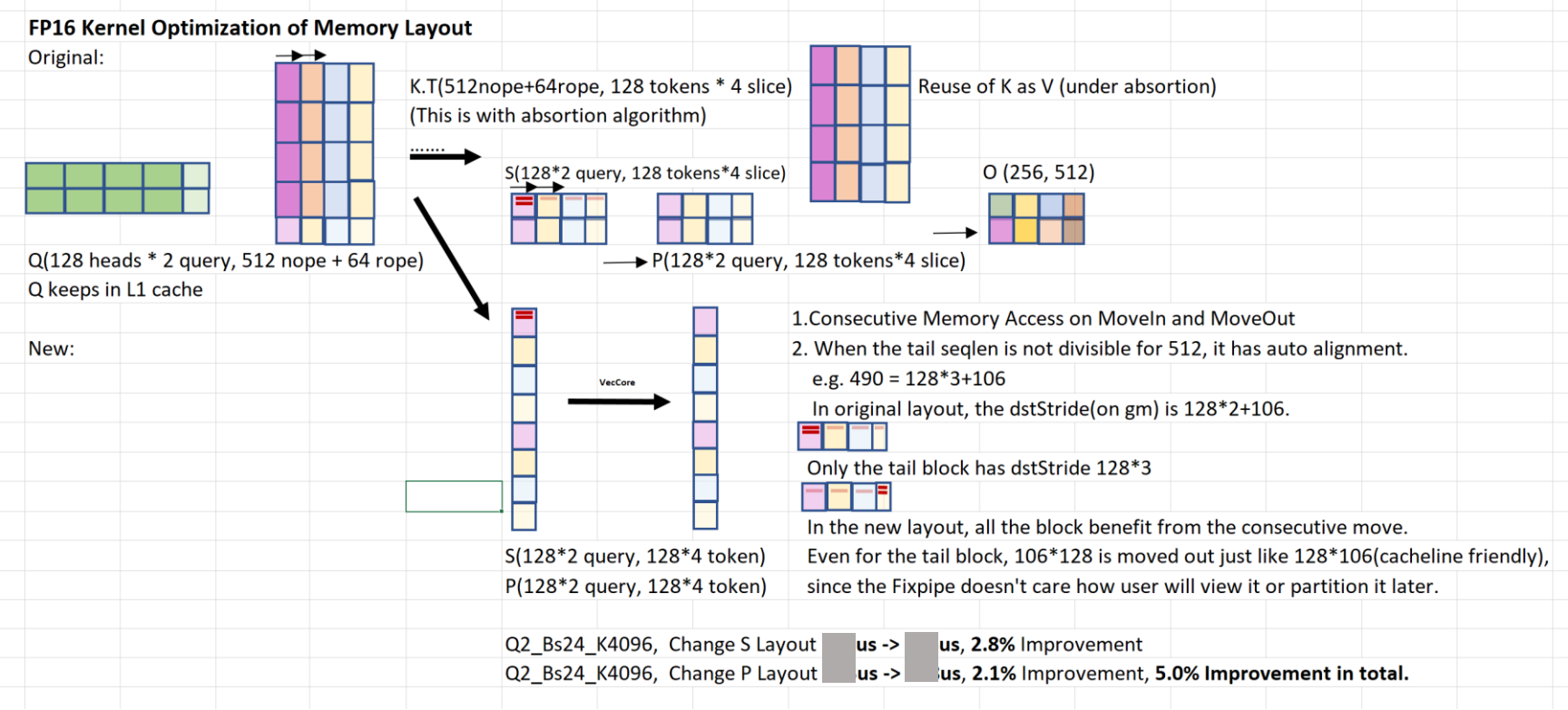
**Teammates' optimizations:** tail rowsum/max; cube/LOC unitflag; nope/rope separation; UP stage pingpong (me participated in this as well); enlarge the headsgroup size in SM stage; synchronization reduction, etc.

FP16 Optimizations: (Implemented from Zhentao)

(I) Yicai He did the first version but Zhentao Fan rewrote it with a different algorithm.



(II)



FP16 Overall Improvement (includes other optimization from teammates):

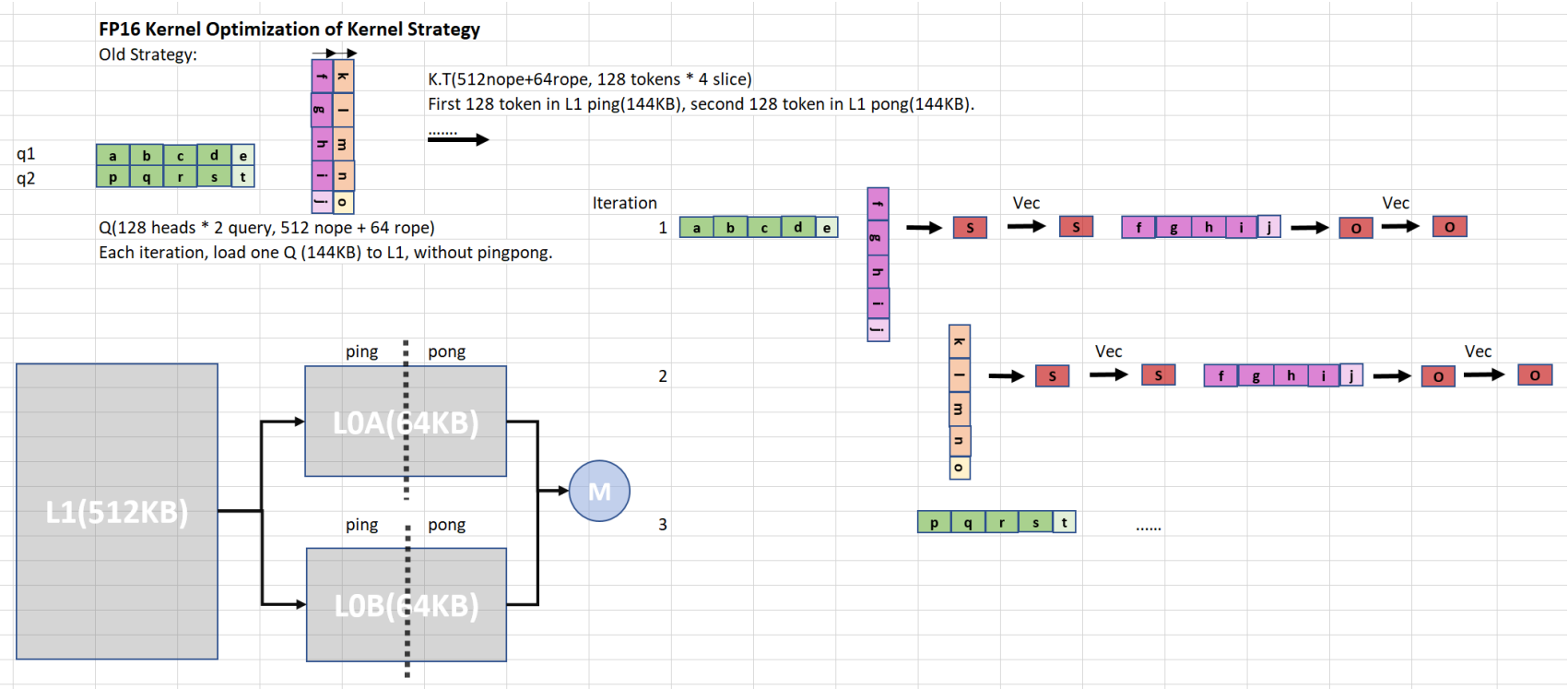
On XX cluster, with optimized operator, the **full model**(not the single operator) inference has **5.XX%** speed up on MTP1\_Bs32/64\_Prefill2048\_Decode2048 case, and **1.XX%** speed up on MTP1\_Bs16\_Prefill2048\_Decode2048. (Results from the model team.) The single **operator** improvement on MTP1 case is also **around 20%**.

For FP16 kernel, the main speed up comes from the load balancing(shown above) and a new Q2 kernel strategy(shown below). Most code of this new kernel strategy (below) is developed by my teammates, and I was mainly for the bugs they stuck with, solved quite lot critical bottlenecks during that development.

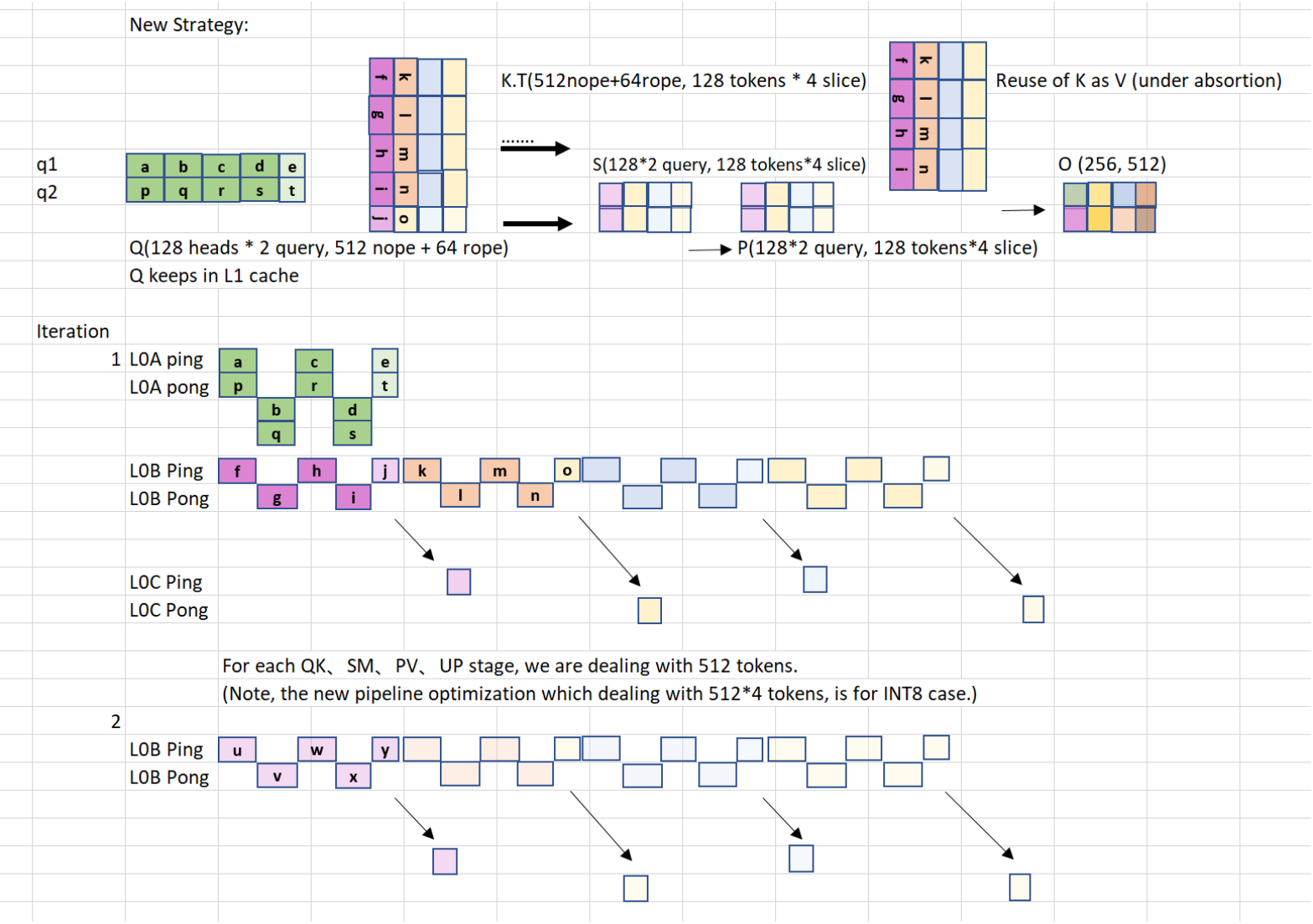


Appendix:

1.Old FP/BF16 kernel strategy:



2.New FP/BF16 Kernel Strategy:



3. MLA Inference Operators on Ascend (Execute Sequentially):

- 1. MLAPO(i.e. Preprocess Operator)[rope, Wk absortion, data preparation for paged attention]
- 2. MLAPA(i.e. Paged Attention) [flashMLA] ←This is what we are optimizing at
- 3. Einsum [Wv absortion]