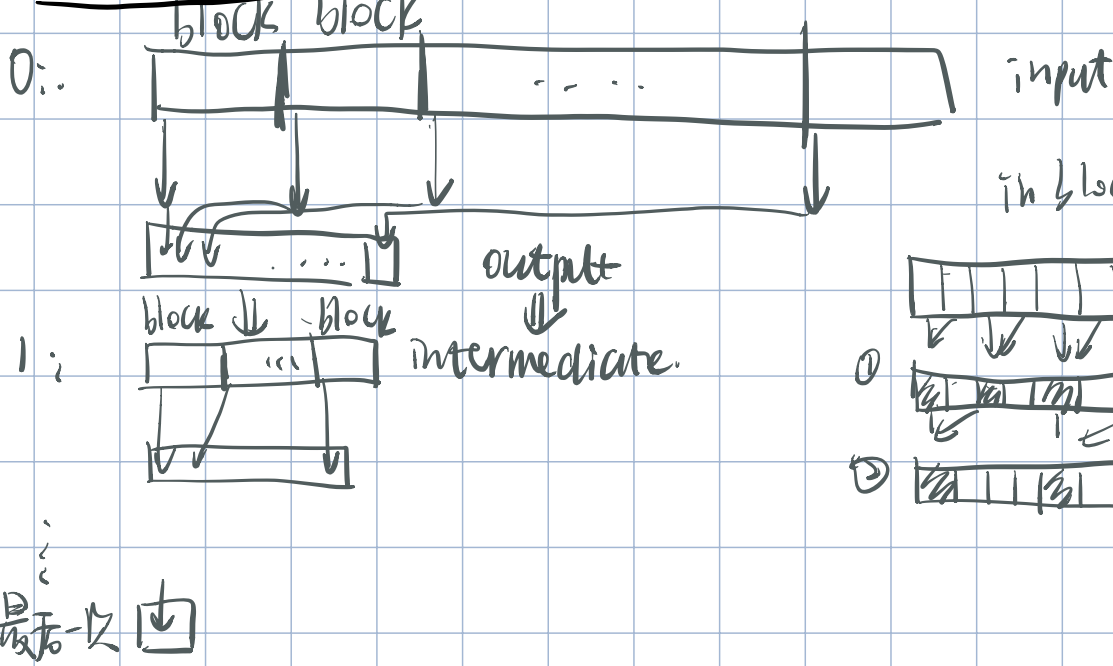
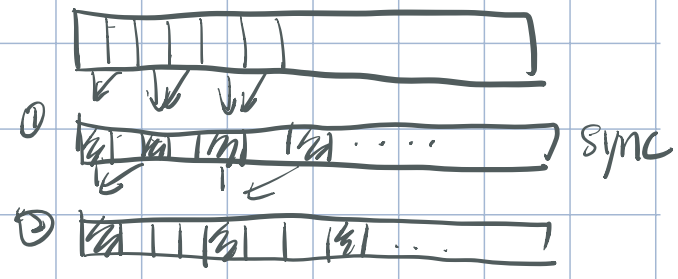


⚠ 这里所有减100都不对. 注意后期看如何改一下逻辑

Reduction 0&1.

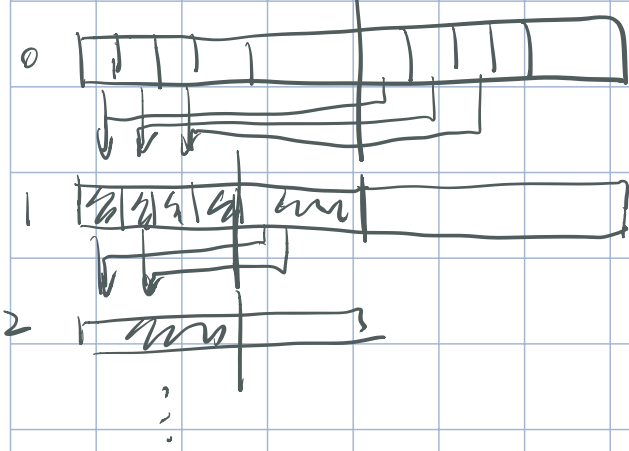


in block:

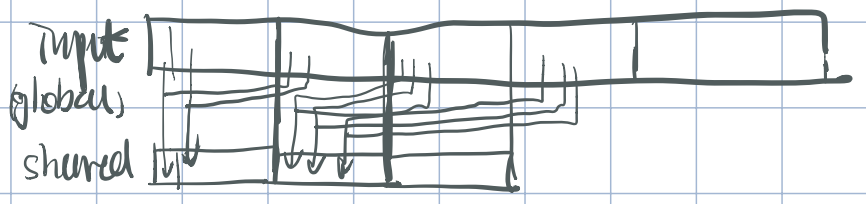


Reduction 2.

within block



Reduction 3

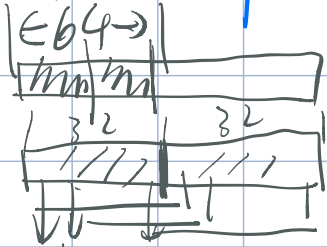


每个 block 里 shared mem 接受 2 个
block 的 global mem 里的内容.
所以 kernel 可以减少一半.

Reduction 4

原理: 同一 warp 用的是同一指令.

可以不用 sync. 展开循环.



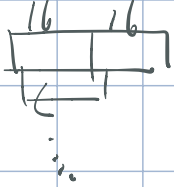
可以保证一定
按顺序执行.
因为 SLMT

Reduction 5

用模板去教. 告知 block 大小.

让 compiler 完成 loop unrolling.

block size 的半最大是 512 threads,
所以可以有 2 个循环展开.



Reduction 6

应该是每个 grid 间聚合起来!

