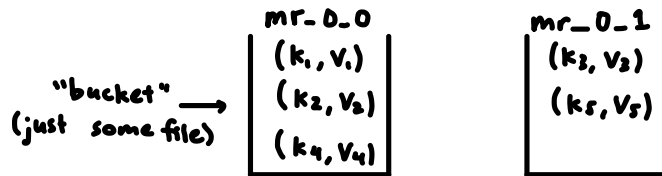


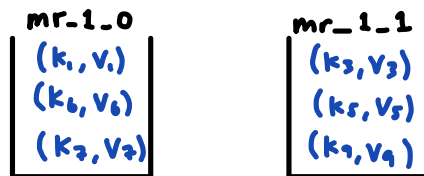
map task 0 (p.txt)



output of map function: n_{out} (k,v) pairs

Let $n_{out} = 5$.

map task 1 (q.txt)



Example: p.txt contains "apple" 5 times. We're running WC.

So $(k_1, v_1) = ("apple", 5)$. q.txt also contains apple, but 3 times.

So $(k_1, v_1) = ("apple", 3)$. Note that the k_1 from map

task 0 is the same as the one from map task 0.

Worker 1 now works on reduce task 0.

mr-0-0	mr-1-0
(k_1, v_1)	(k_1, v_1)
(k_2, v_2)	(k_6, v_6)
(k_4, v_4)	(k_7, v_7)

Concat all 6 (k,v) pairs and sort.

(k_1, v_1)
(k_1, v_1)
(k_7, v_7)
(k_2, v_2)
(k_6, v_6)
(k_4, v_4)

Reduce k_1 : Take in (k_1, v_1) and (k_1, v_1) and produce (k_1, v_1) .

Reduce k_7 : Take in (k_7, v_7) and produce (k_7, v_7) .

etc.

Note that the following scenario where a reduce task doesn't "fully aggregate" all the vals for a key is impossible.

Ex:

mr-0-0	mr-0-1
(k_1, v_1)	(k_3, v_3)
(k_2, v_2)	(k_5, v_5)
(k_4, v_4)	

mr-1-0	mr-1-1
(k_5, v_5)	(k_3, v_3)
(k_6, v_6)	(k_1, v_1)
(k_7, v_7)	(k_9, v_9)

Here, reduce task 0 won't be aware of (k_1, v_1) .

This is impossible b/c the hash function and modulo will put (k_1, v_1) and (k_1, v_1) into the "same bucket" (i.e. if (k_1, v_1) goes into mr-0-i, (k_1, v_1) goes into mr-1-i).