



CS 165

Data Systems

Have fun learning to design and build modern data systems

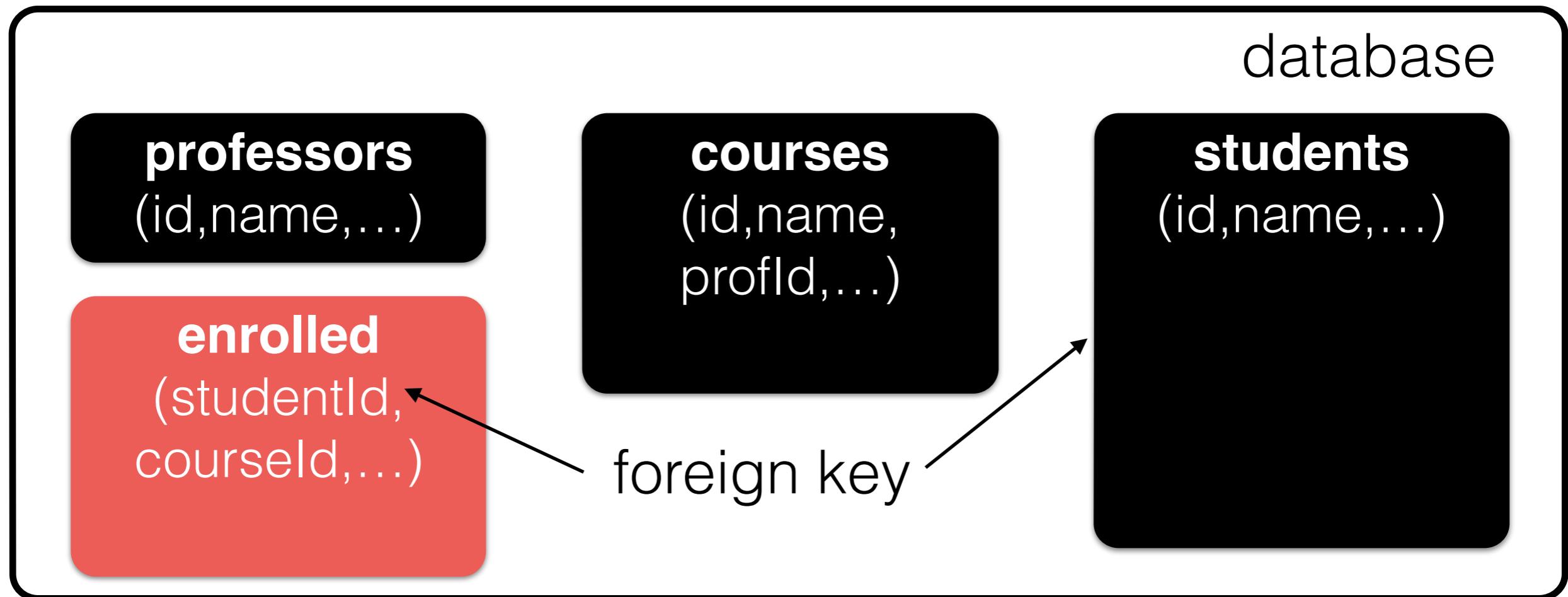
class 16

hash joins

prof. Stratos Idreos

[HTTP://DASLAB.SEAS.HARVARD.EDU/CLASSES/CS165/](http://DASLAB.SEAS.HARVARD.EDU/CLASSES/CS165/)

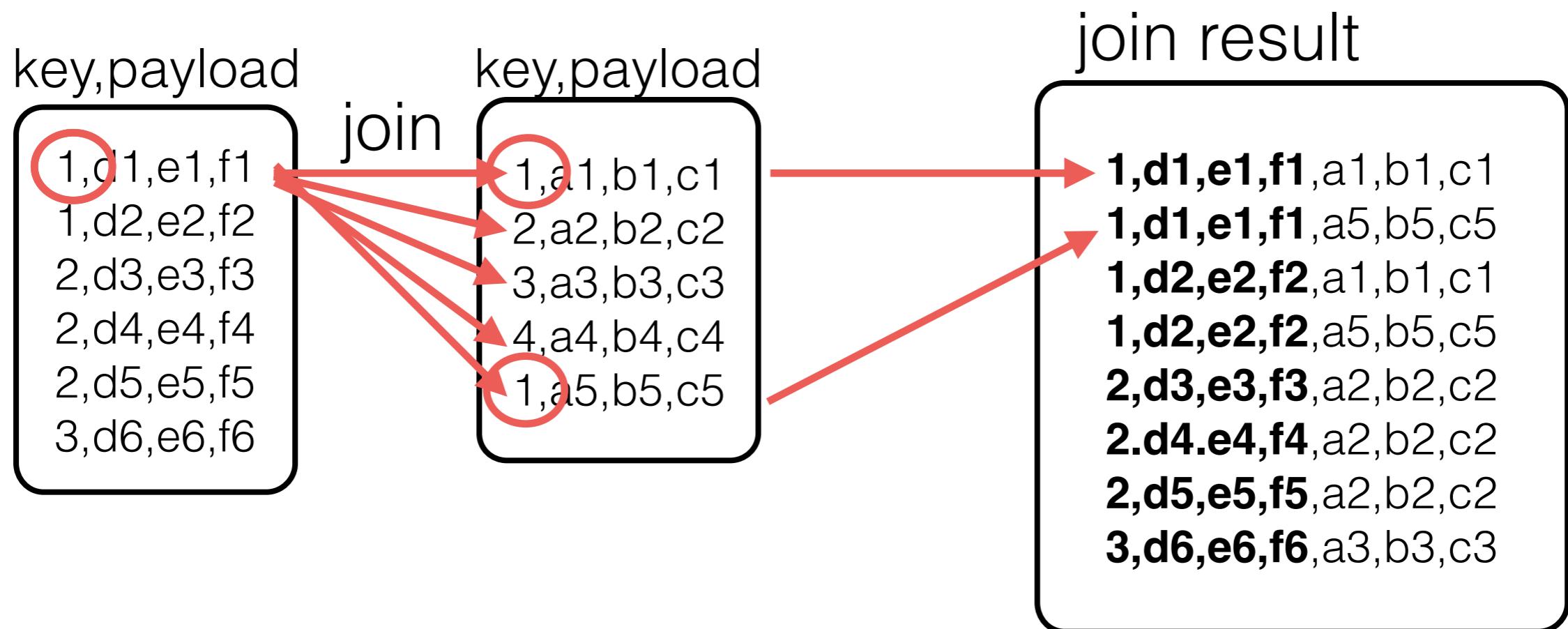


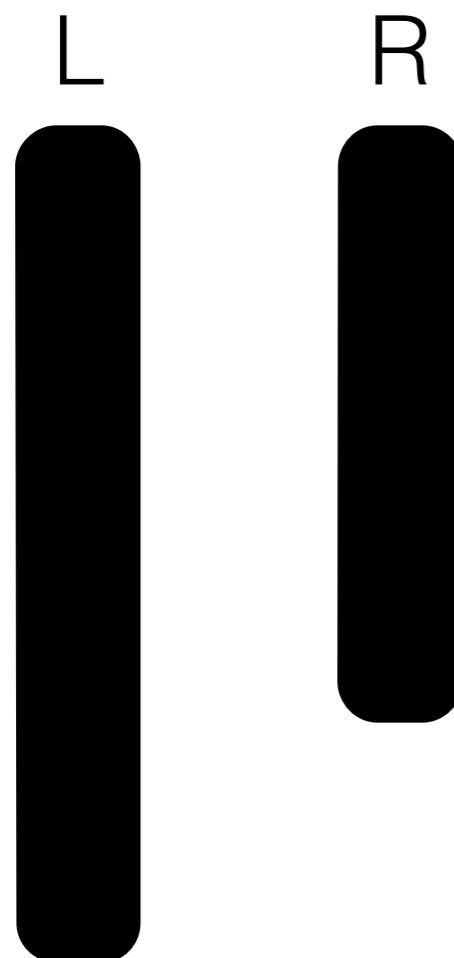


give me all students enrolled in cs165

select student.name **from** students, enroled, courses **where**
courses.name=“cs165” and enroled.courseId=course.id and
student.id=enroled.studentId **join**

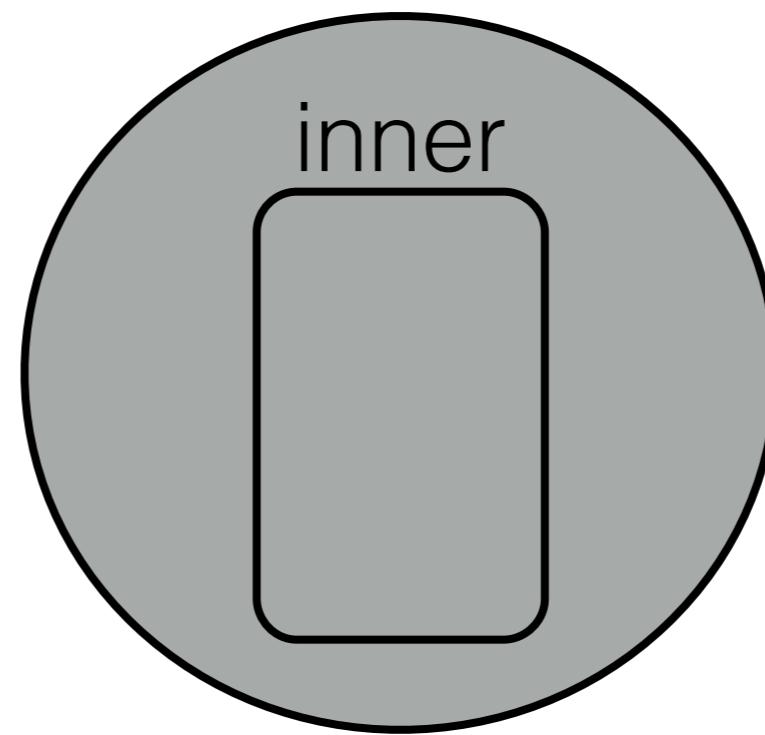
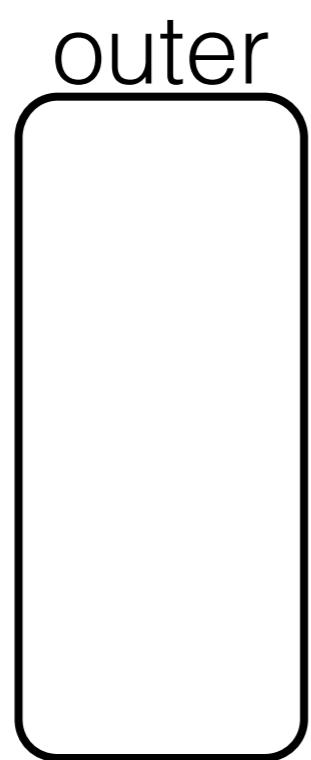






```
new resL[]; new resR[]; k=0
for (i=0;i<L.size;i++)
    for (j=0;j<R.size;j++)
        if L[i]==R[j]
            resL[k]=i
            resR[k++]=j
```

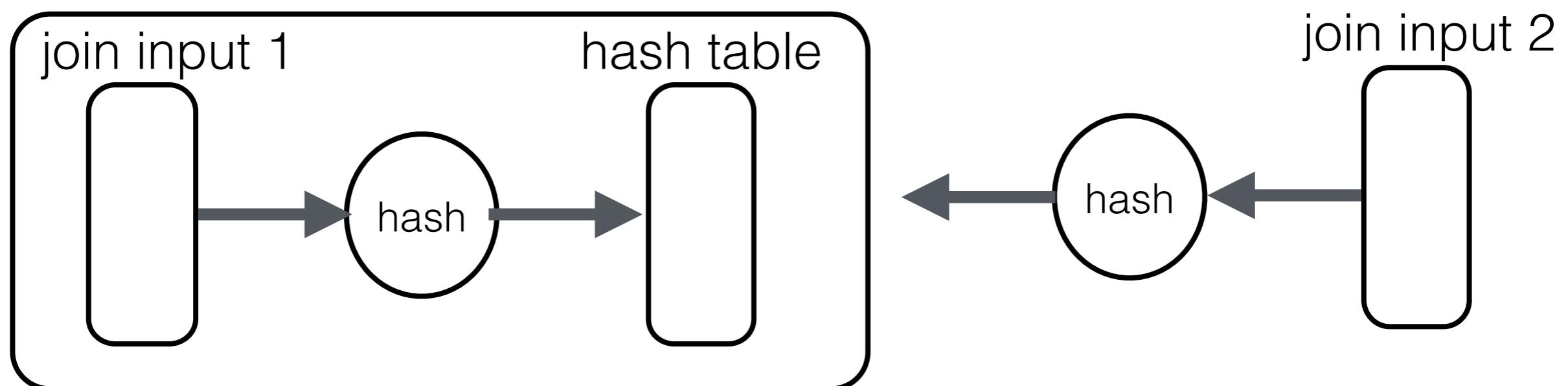




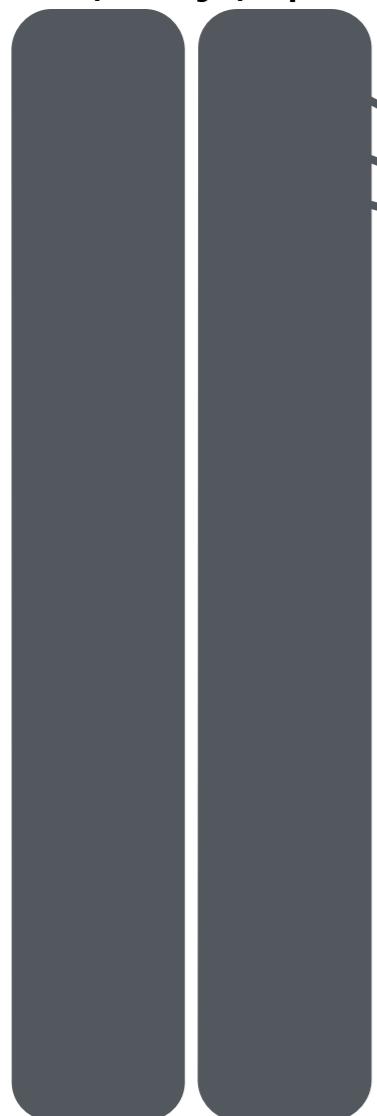
oracle
search in $O(1)$



hash join



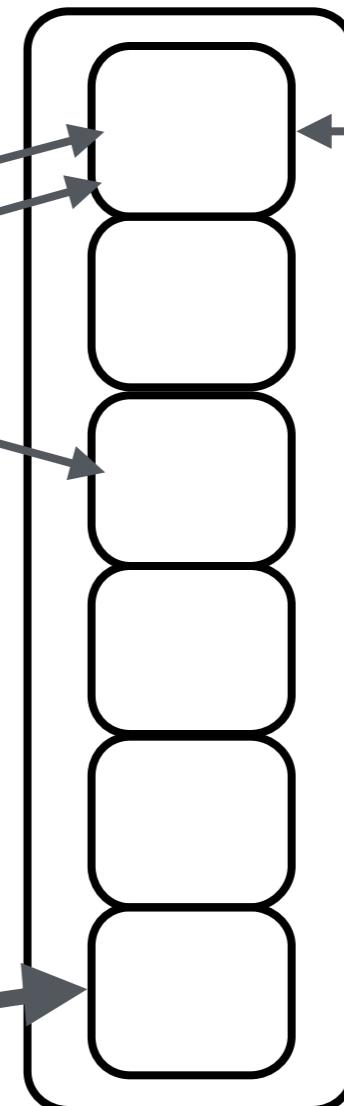
val (key),pos



hash

$$h=f(\text{val})$$
$$\text{bucket}=h \bmod k$$

hash table



bucket

$[(\text{val}1, \text{pos}1), (\text{val}7, \text{pos}7), \dots]$

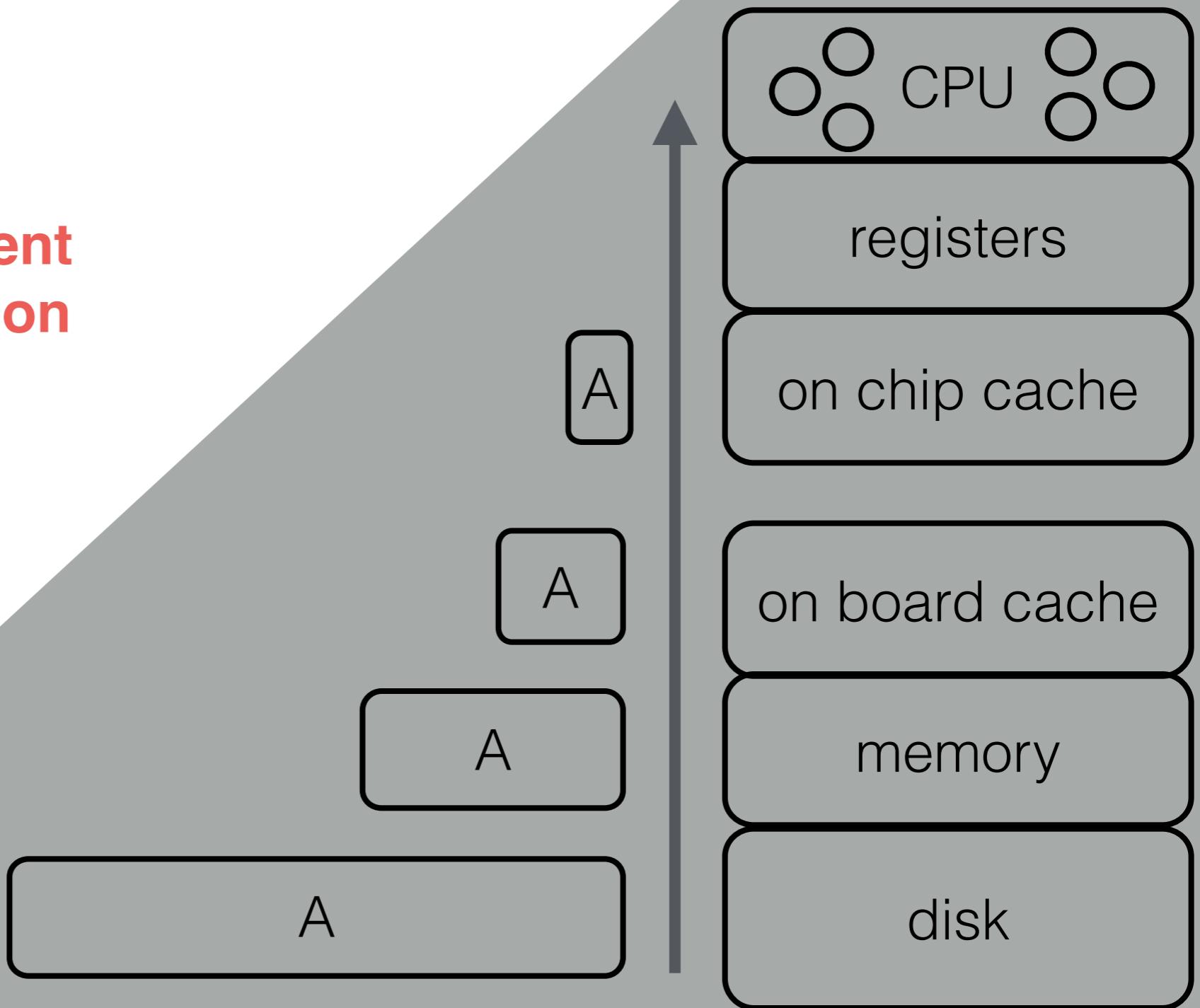
bucket size



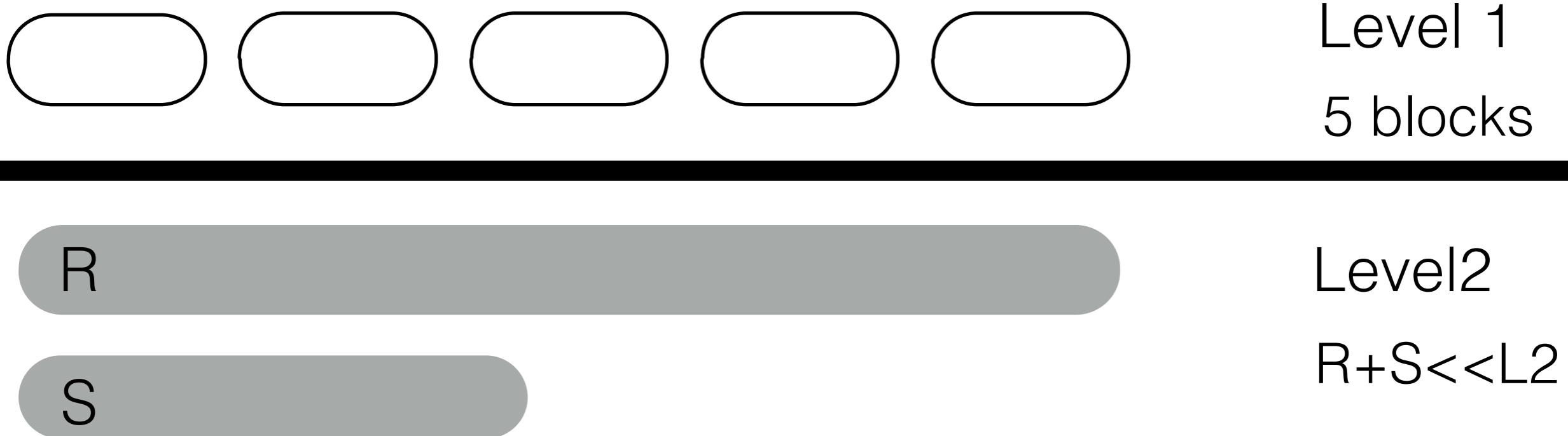
assumption for this class:

we know how many buckets we need
(more in next class about this)

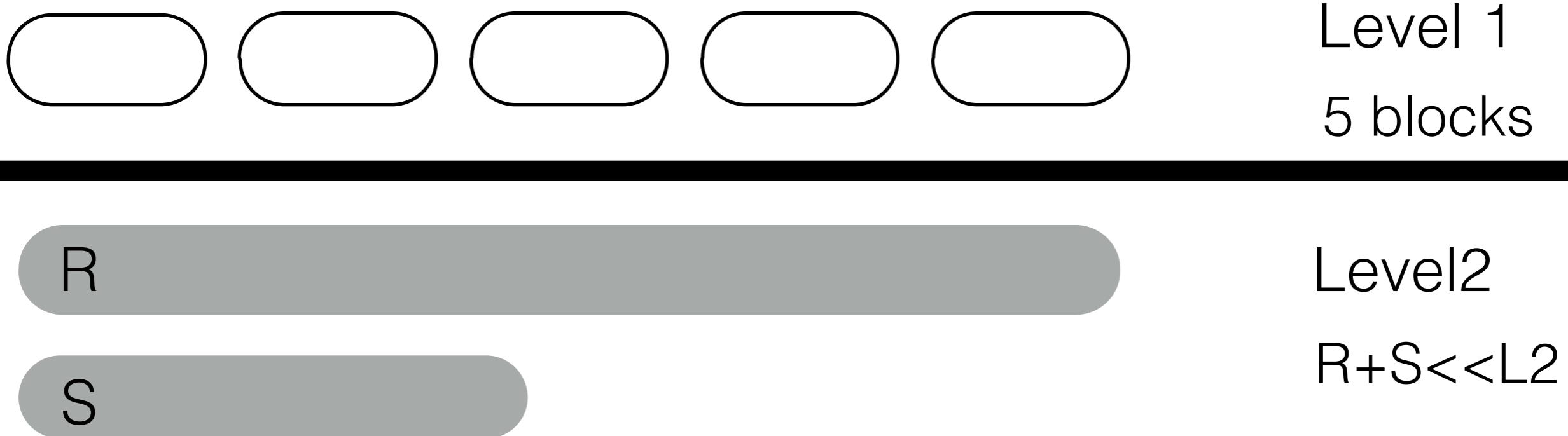
**minimize data movement
maximize CPU utilization**



goal: join R(val,pos) and S(val,pos) = Res(posR,posS)

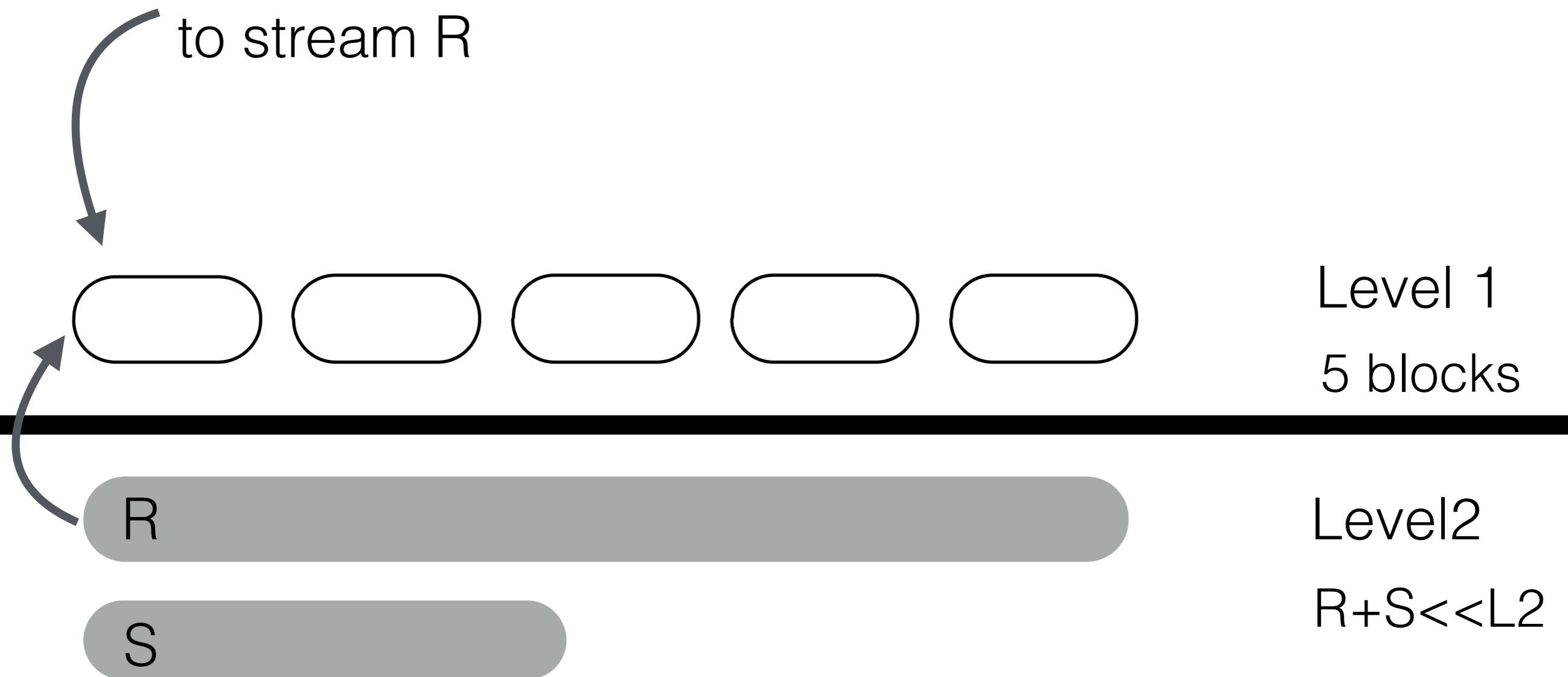


goal: join R(val,pos) and S(val,pos) = Res(posR,posS)



goal: join R(val,pos) and S(val,pos) = Res(posR,posS)

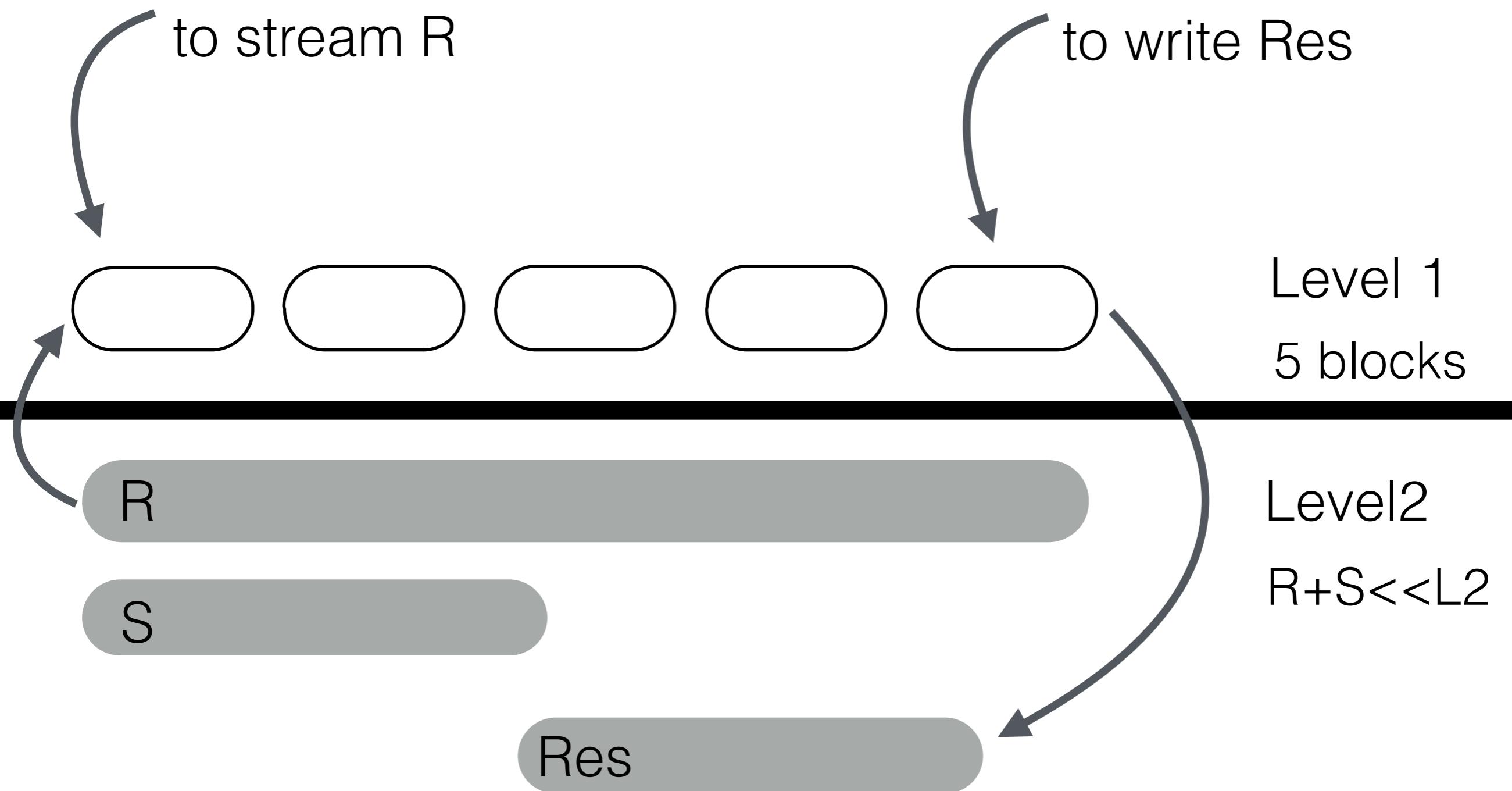
need one block
to stream R



goal: join $R(\text{val}, \text{pos})$ and $S(\text{val}, \text{pos}) = \text{Res}(\text{posR}, \text{posS})$

need one block
to stream R

need one block
to write Res



goal: join $R(\text{val}, \text{pos})$ and $S(\text{val}, \text{pos}) = \text{Res}(\text{posR}, \text{posS})$

need one block
to stream R

need one block
to write Res

hash table, $S \leq L_1 - 2$

Level 1
5 blocks

R

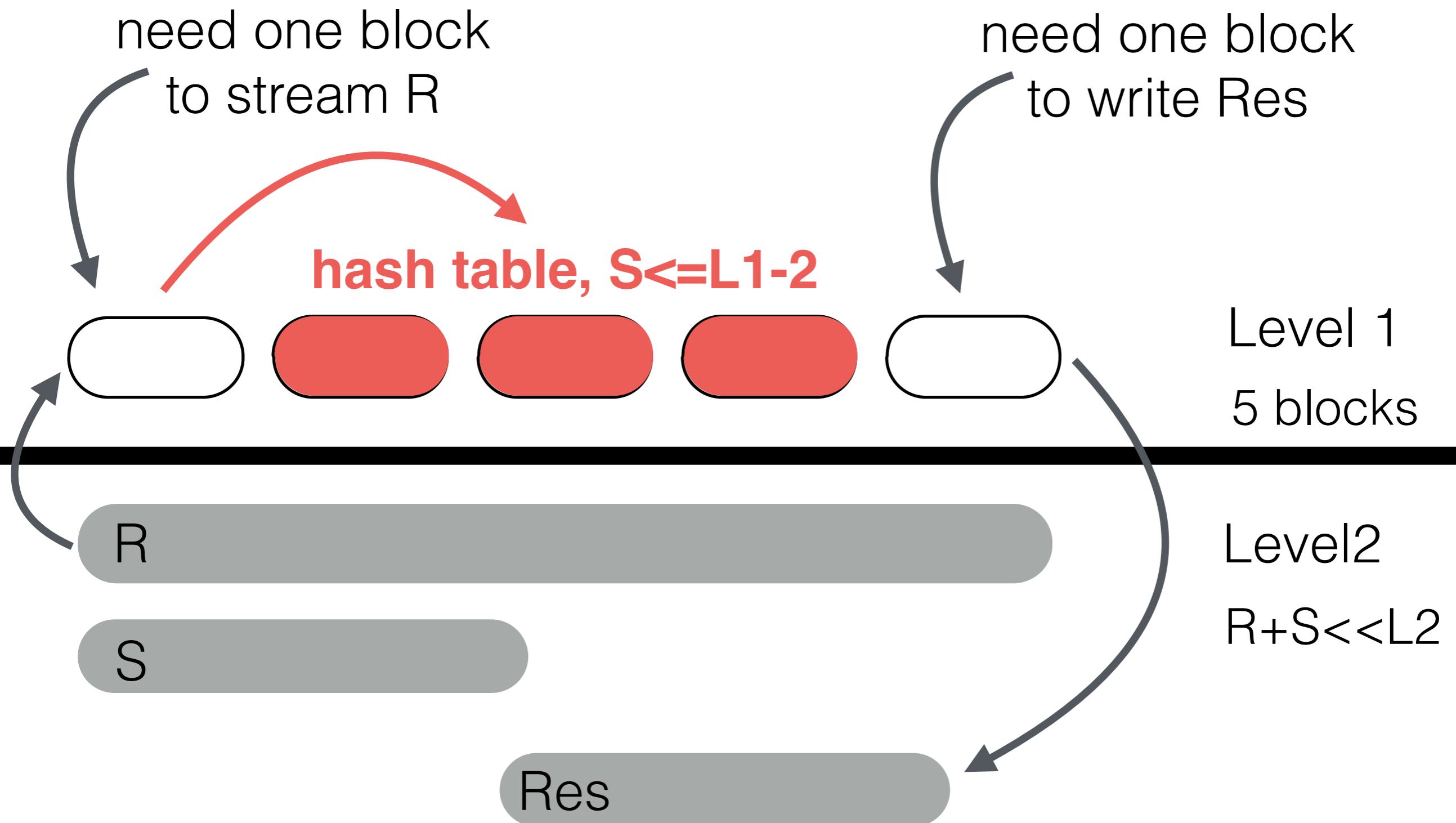
S

Res

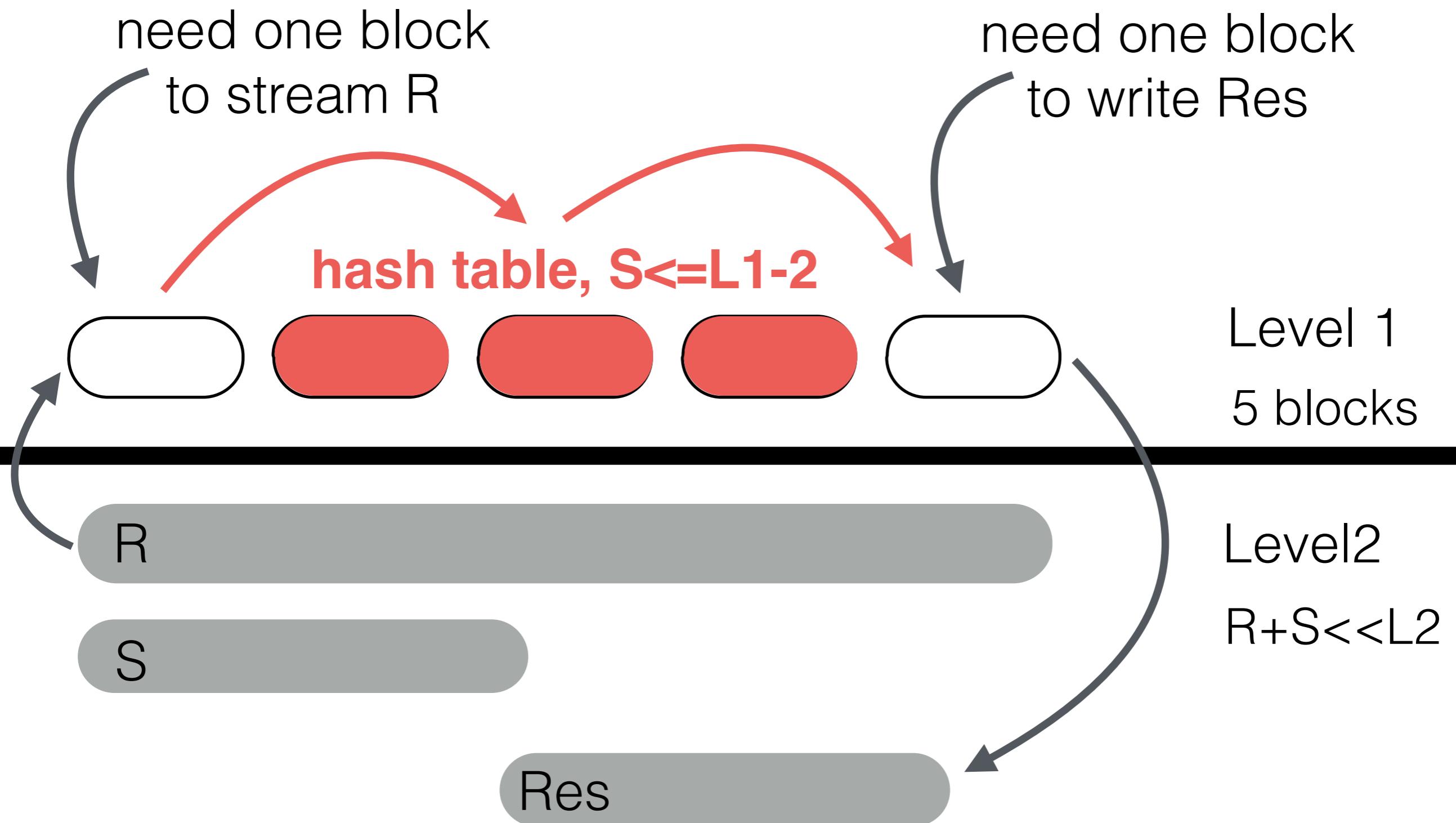
Level2

$R + S \ll L_2$

goal: join $R(\text{val}, \text{pos})$ and $S(\text{val}, \text{pos}) = \text{Res}(\text{posR}, \text{posS})$



goal: join $R(\text{val}, \text{pos})$ and $S(\text{val}, \text{pos}) = \text{Res}(\text{posR}, \text{posS})$



goal: join R(val,pos) and S(val,pos) = Res(posR, posS)

need one block
to stream S

create hash table on S

unused

Level 1
5 blocks

R

S

Level2

R+S<<L2

goal: join R(val,pos) and S(val,pos) = Res(posR, posS)

need one block
to stream S

create hash table on S

unused

Level 1
5 blocks

R

S

Level2

$R+S \ll L2$

Total Cost=R+S+Res (blocks)

what if $R > L1-2$ & $S > L1-2$?

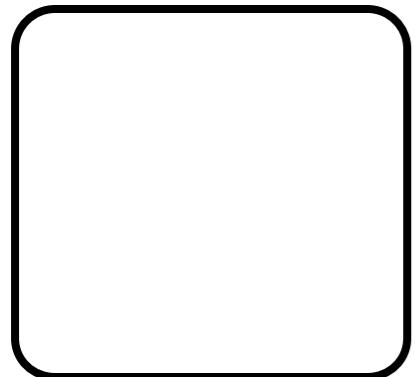
```
new resL; new resS; k=0;r=0;  
for (i=0,i<S.size;i++,k++)  
    addToHash(ht,S[i], S[i].val mod L1-2)          hash first L1-2  
    if (k==L1-2 || i==S.size-1)  
        k=0  
        for (j=0;j<R.size;j++)  
            res=probe(ht,R[i].val);  
            if (res!=null)  
                resR[r]=R[j].[pos]  
                resS[r++]=res.pos  
empty(ht);
```

scan all R
and probe

repeat until
we hash everything in S

Total Cost= $R*S/(L1-2)+S+Res$ (blocks)

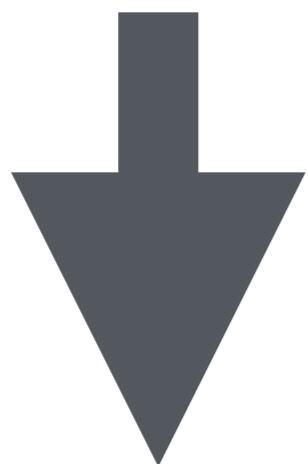
join input 1



join input 2

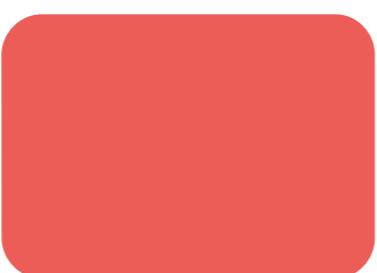
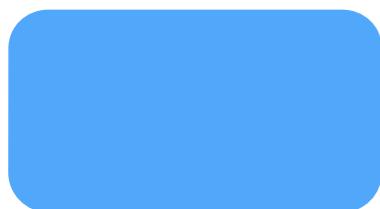


grace hash join



hash partitioning

one pass to partition



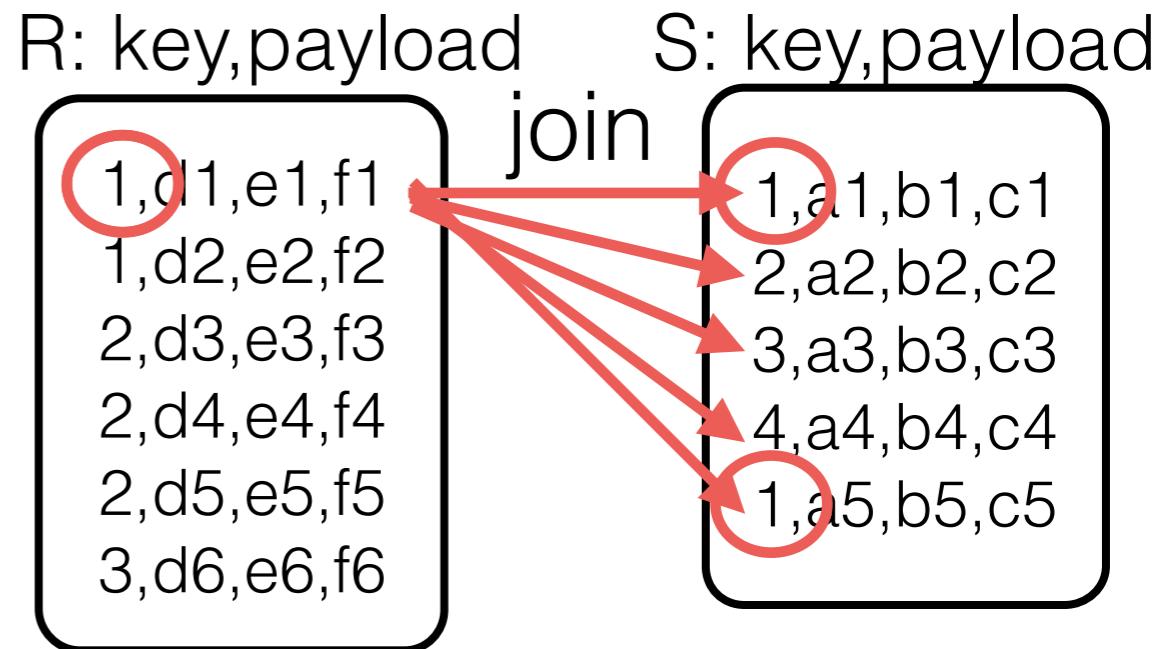
then one pass to
join each pair of partitions
independently in memory





for every
L.key = R.key pair
return [L.pos,R.pos]

data/results
one column-at-a-time



CPU

level 1

level 2

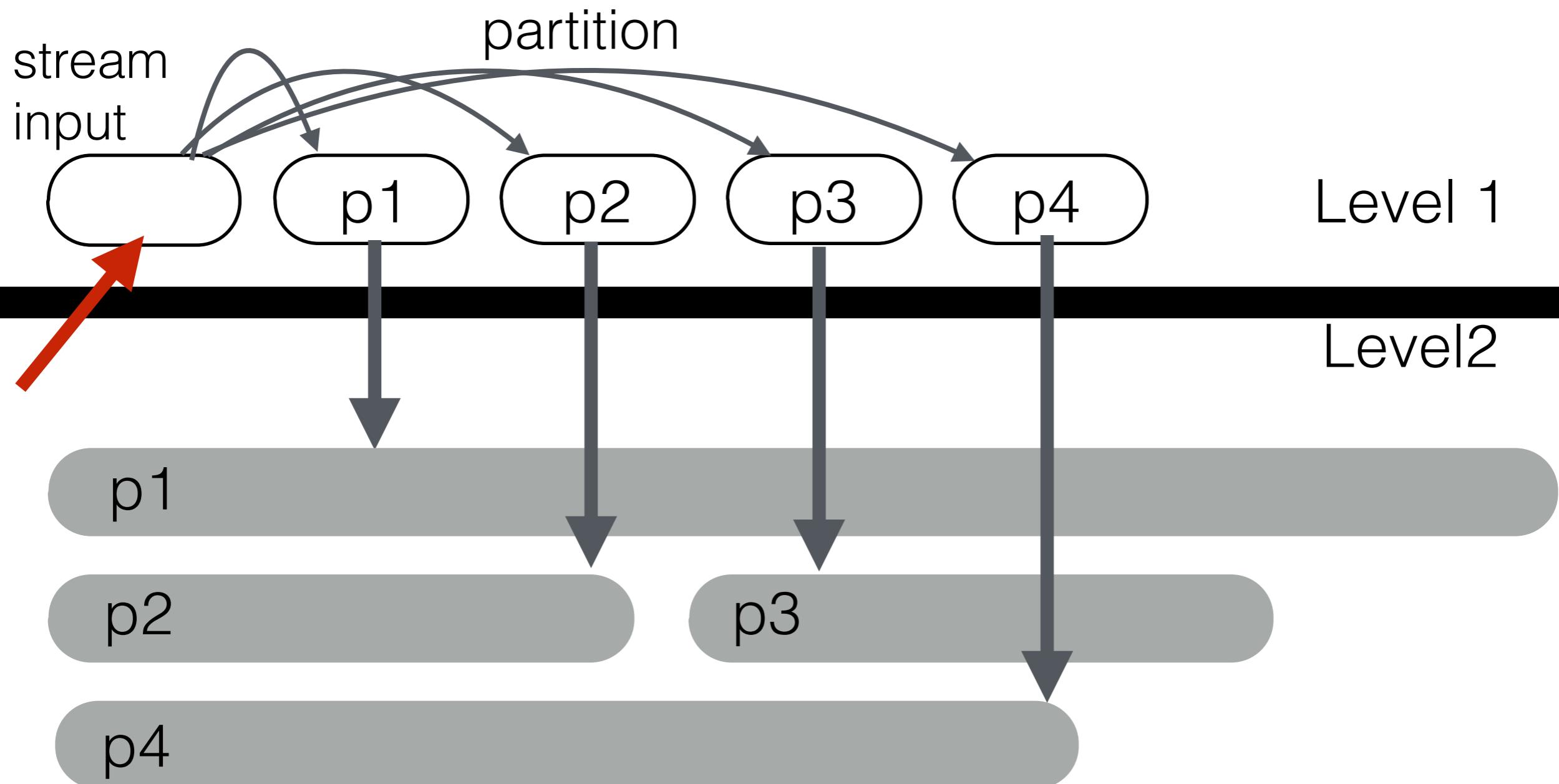
- 1) design a grace join algorithm and give its cost
- 2) which table do we start from?
- 3) what if partitions end up being >L1-2?
- 4) when can we do grace join in 2 passes max?

- 5) utilize multi-cores
- 6) can you keep all cores at 100% all the time?

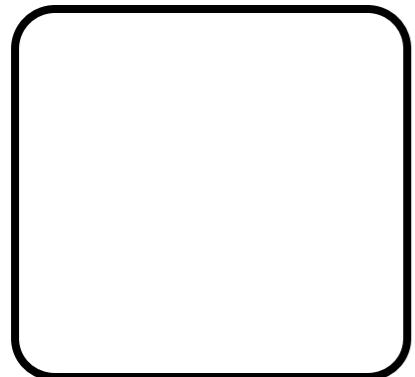


1. read input into stream buffer, hash and write to respective partition buffer
2. when input buffer is consumed, bring the next one
3. when a partition buffer is full, write to L2

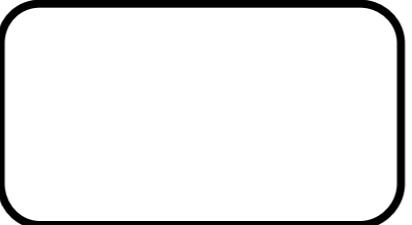
we can partition into L1-1 pieces in one pass



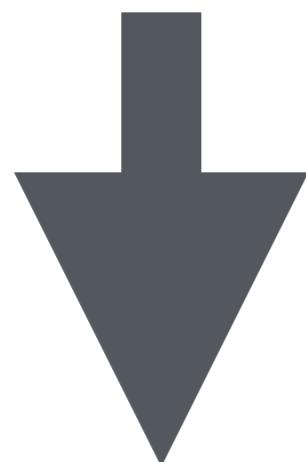
join input 1



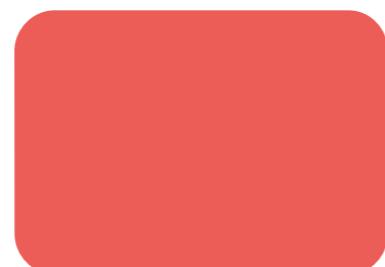
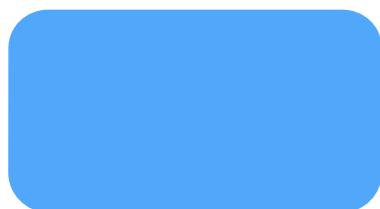
join input 2



grace hash join



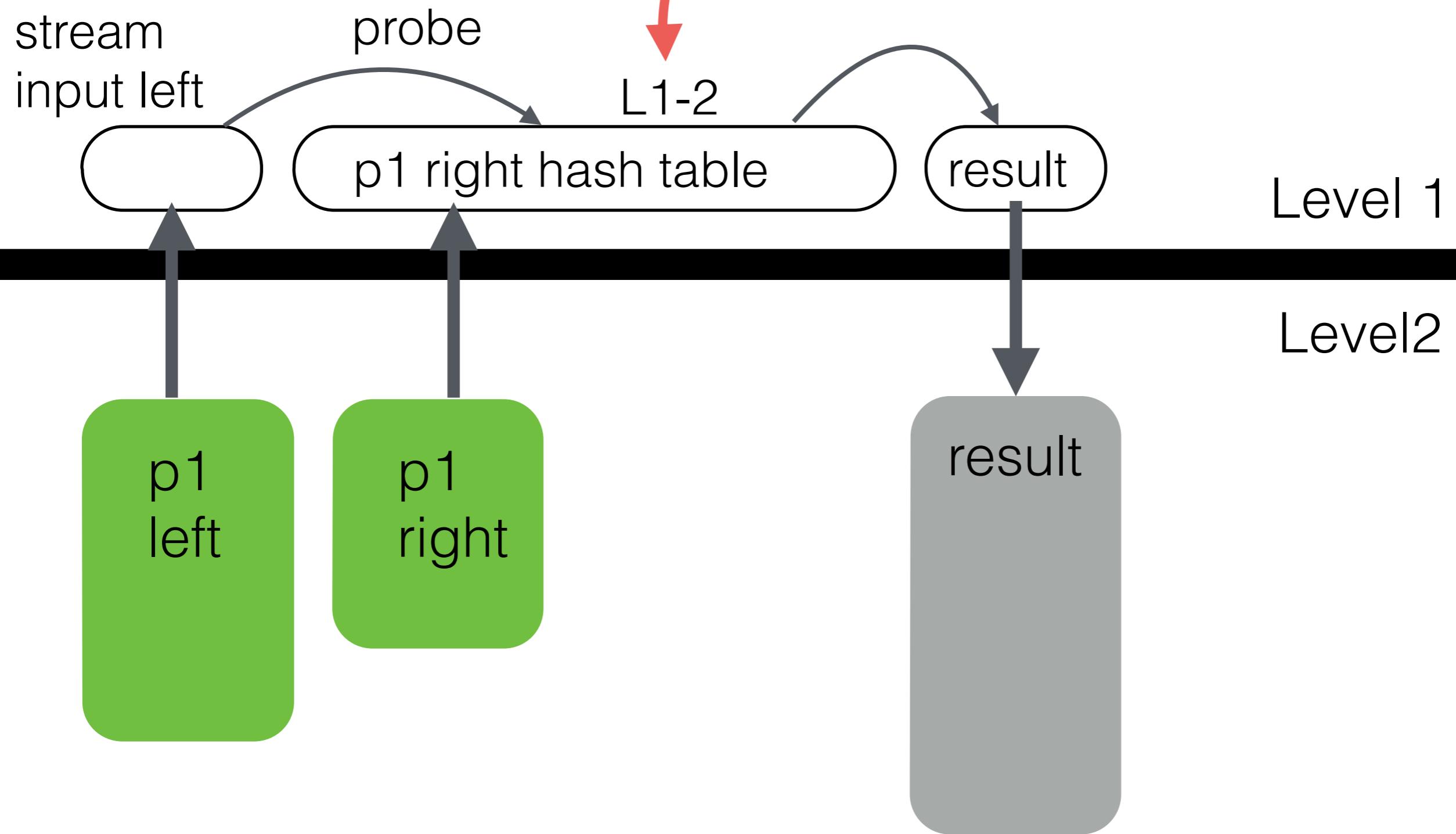
hash partitioning



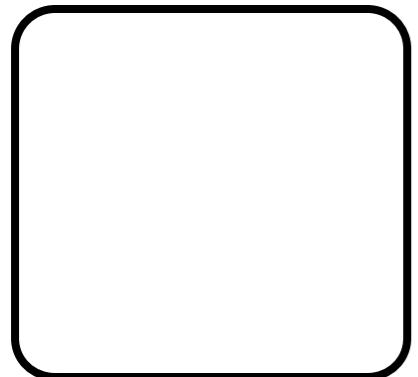
**then join each pair of partitions
independently in memory**



as long as at least one
of the pieces $\leq L1-2$



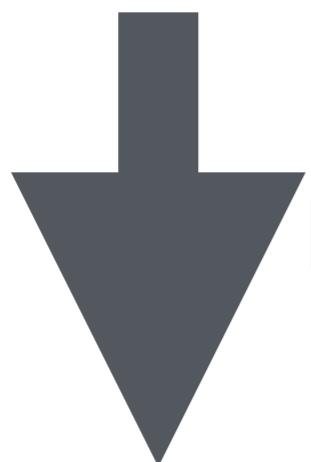
join input 1



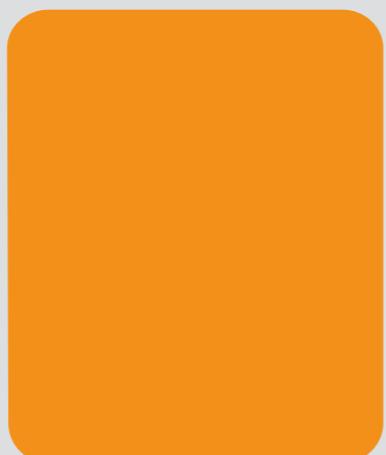
join input 2



grace hash join



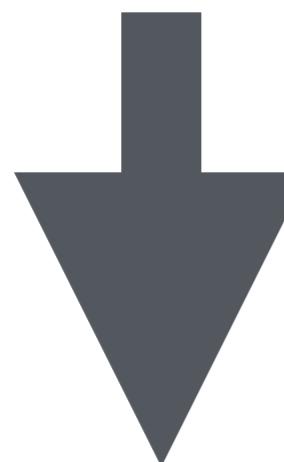
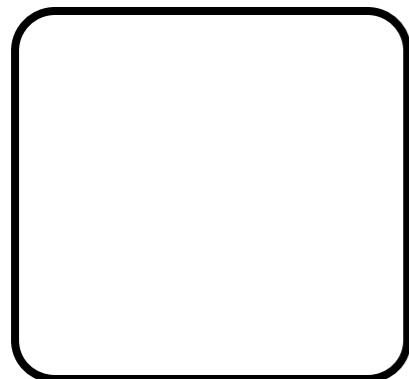
hash partitioning



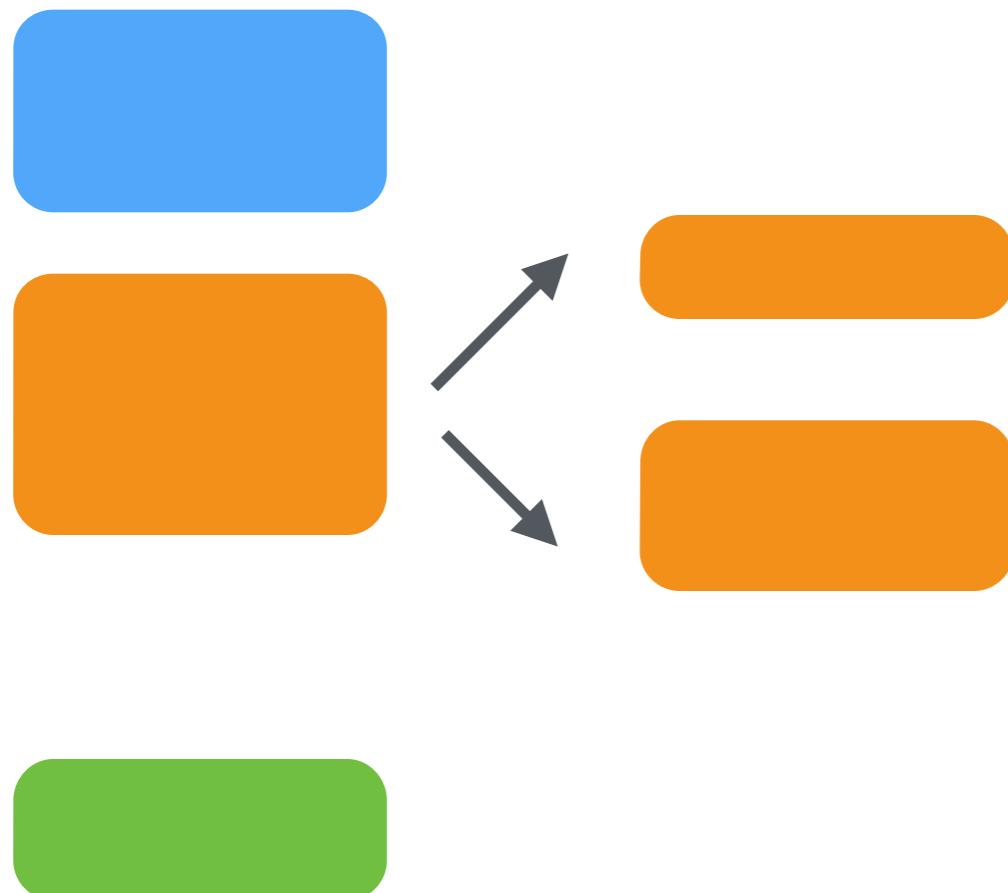
**both left and
right side >L1-2**



grace hash join



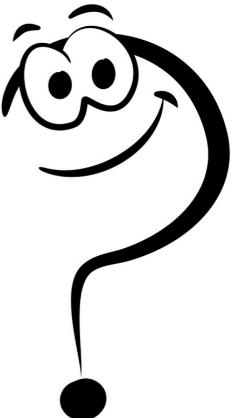
hash partitioning



apply recursively if a partition does not fit in memory

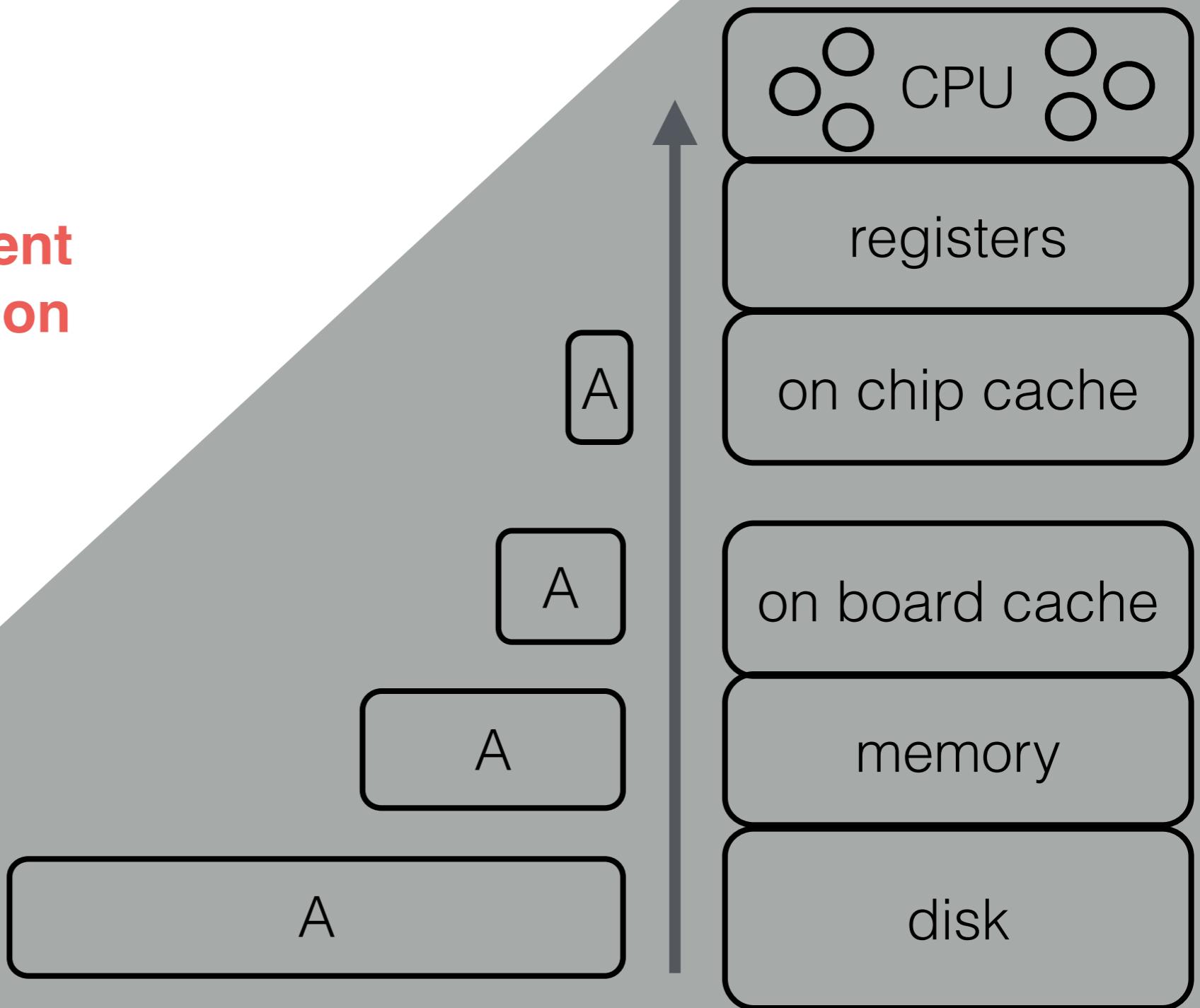


when can we do grace join
in exactly 2 passes



- 1) at least one side should fit in L1-2
- 2) simplify problem by considering one side only:
if all partitions we create fit in L1-2 we are ok
- 3) the maximum number of partitions we can create in one pass is L1-1
- 4) so if $R/L1-1 \leq L1-2 \rightarrow R \leq (L1-1)(L1-2)$
we will not need to repartition any pieces

**minimize data movement
maximize CPU utilization**



how to partition in parallel?

join

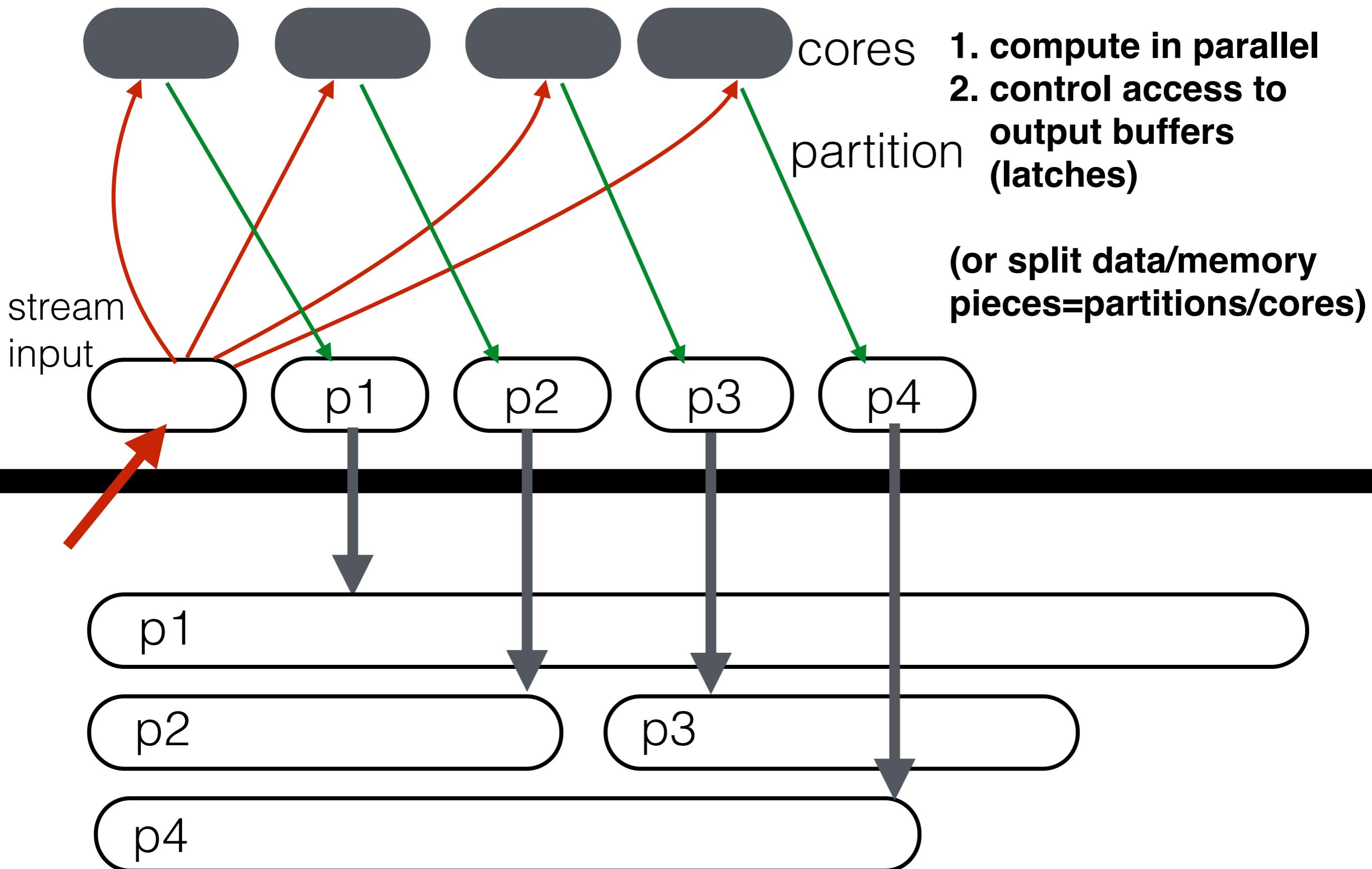
partition each input

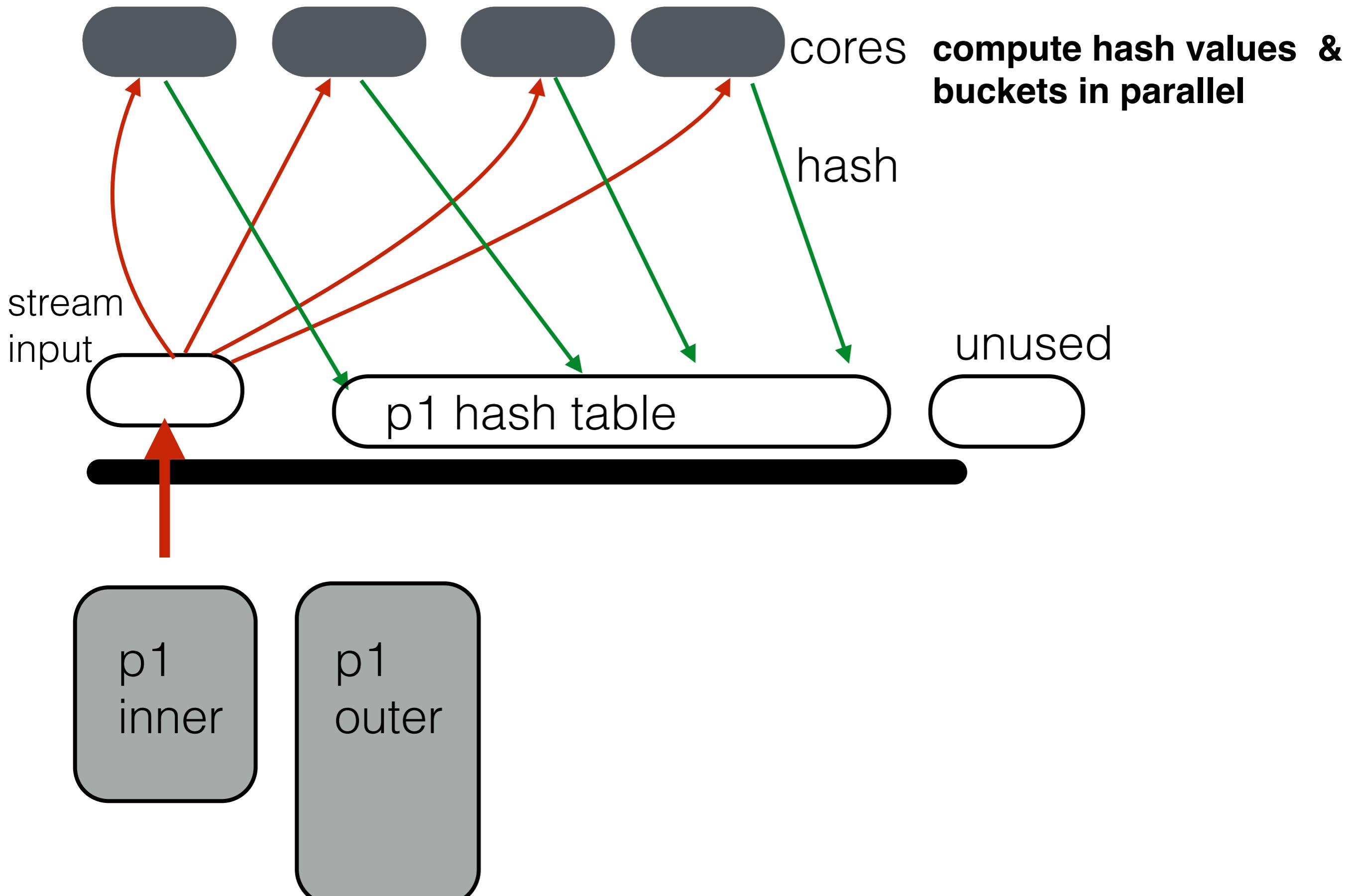
for each pair of partitions
create hash table
probe hash table

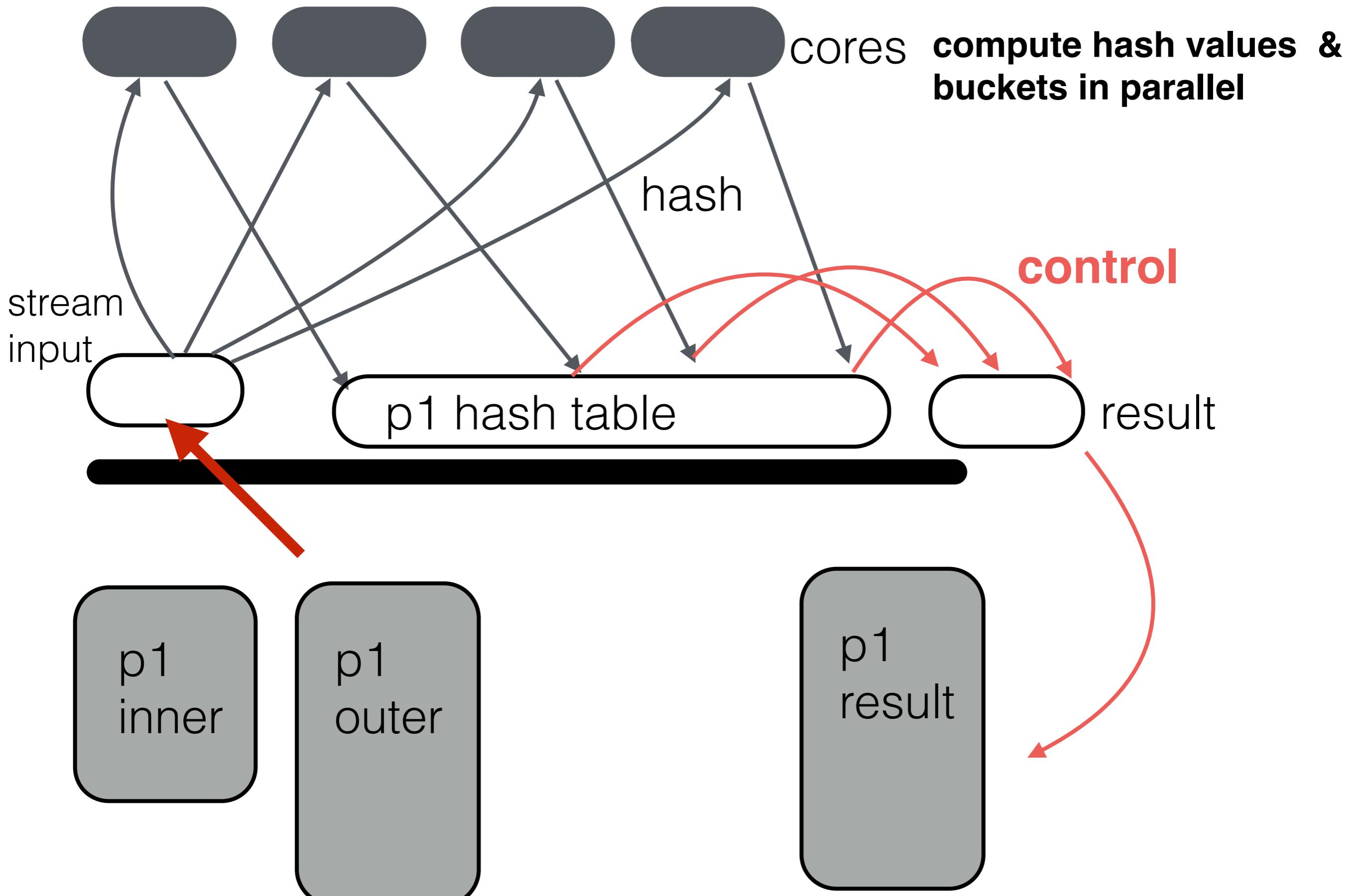
how to create HT in parallel?

how to probe HT in parallel?









Translation Look aside Buffer (TLB)

CPU

L3

memory

disk

we need to translate virtual memory addresses to physical memory

virtual memory



CPU

cache conscious grace join

L1

L2

L3

memory

disk

partition R such as each hash table fits in L3 (L3-2)
as much as $K=L3-1$ partitions in one pass

every time a buffer/cache line is full spill to memory
we may be writing to K memory areas

if $K > TLB$ then we incur TLB misses



radix join

recursively partition with maximum **outputs \leq TLB**

may do more passes but sequential access for reads
and random access for writes $<$ TLB

no TLB misses





joins

(project=m4)

what to do in m4?

nested loops and hash joins

cache conscious and multi-core (within reason)



Join Processing in Databases with Large Main Memories

L. Shapiro

ACM Transactions on Database Systems. 11(3), 1986

Cache Conscious Algorithms for Relational Query Processing

A. Shatdal, C. Kant and J. Naughton

Very Large Databases Conference, 1984

Database Architecture Optimized for the new Bottleneck: Memory Access

P. Boncz, S. Manegold and M. Kersten

Very Large Databases Conference, 1999

Sort vs. Hash Revisited: Fast Join Implementation on Modern Multi-Core CPUs

Changkyu Kim, et all.

International Conference on Very Large Databases, 2009

hash joins

DATA SYSTEMS

prof. Stratos Idreos

