



CS 165

Data Systems

Have fun learning to design and build modern data systems

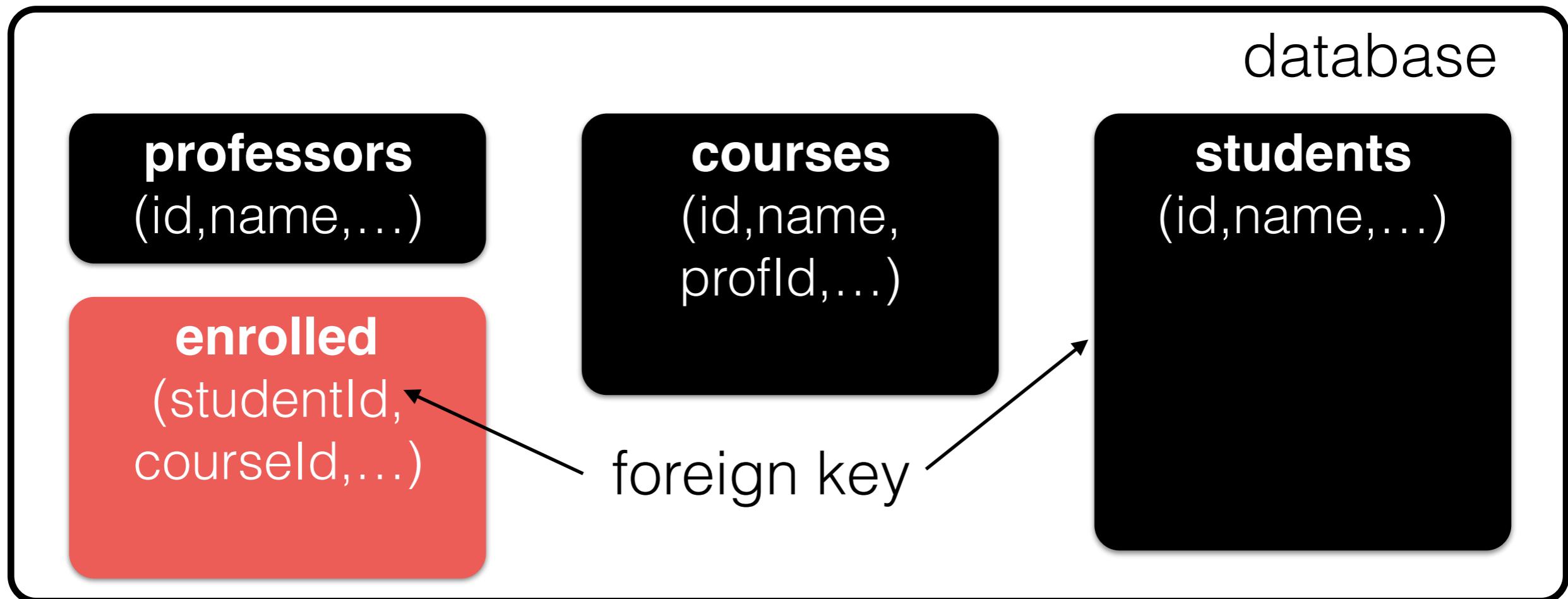
class 18

hashing

prof. Stratos Idreos

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

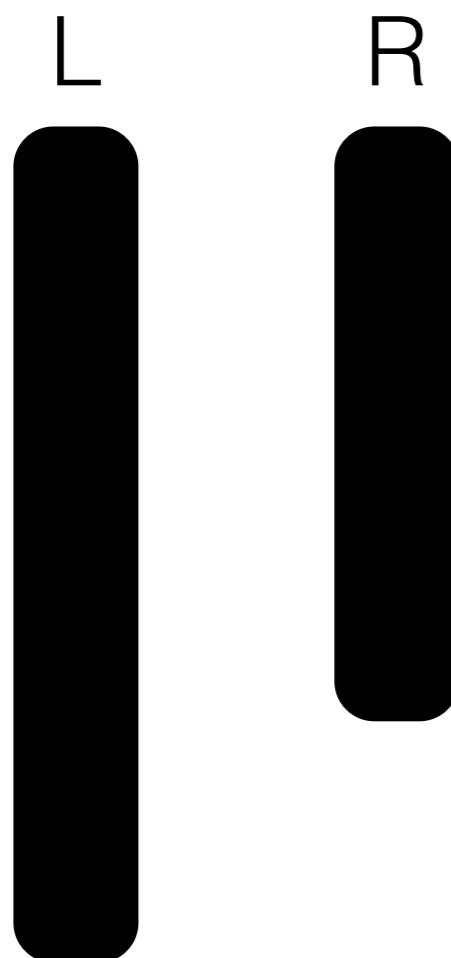




give me all students enrolled in cs165

select student.name **from** students, enrolled, courses **where**
 courses.name="cs165" and enrolled.courseId=course.id and
 student.id=enrolled.studentId **join**

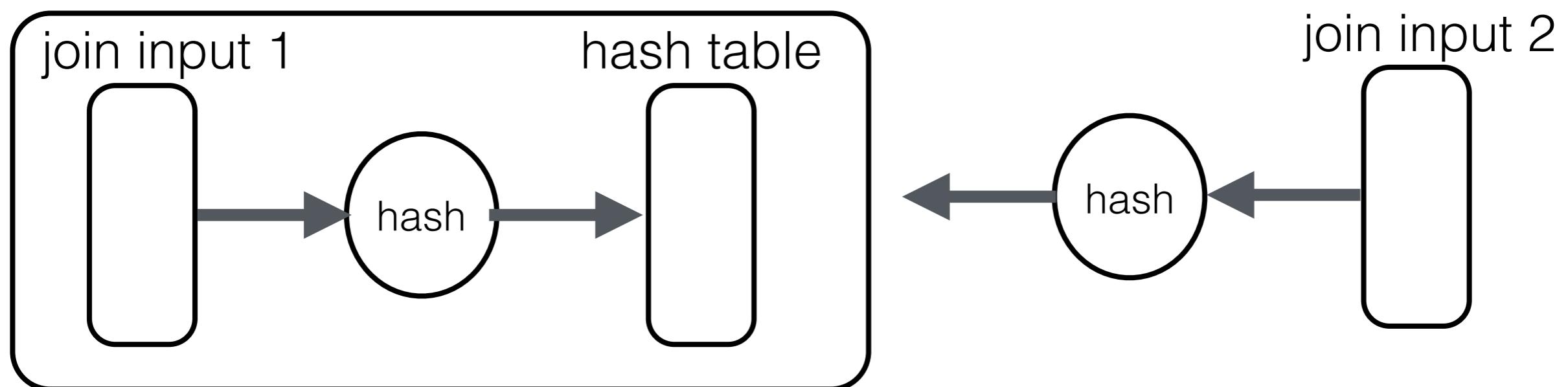


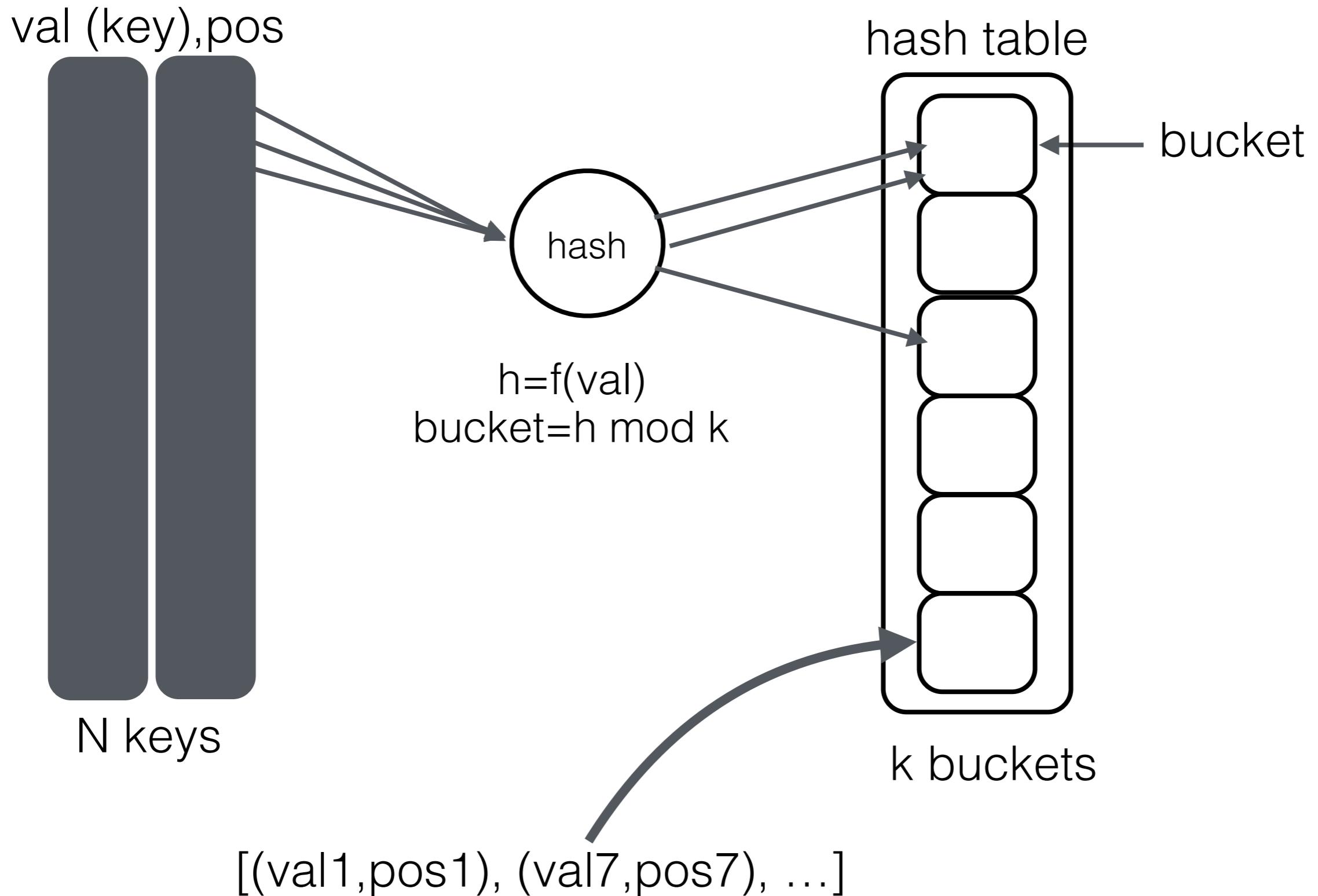


```
new resL[]; new resR[]; k=0
for (i=0;i<L.size;i++)
    for (j=0;j<R.size;j++)
        if L[i]==R[j]
            resL[k]=i
            resR[k++]=j
```



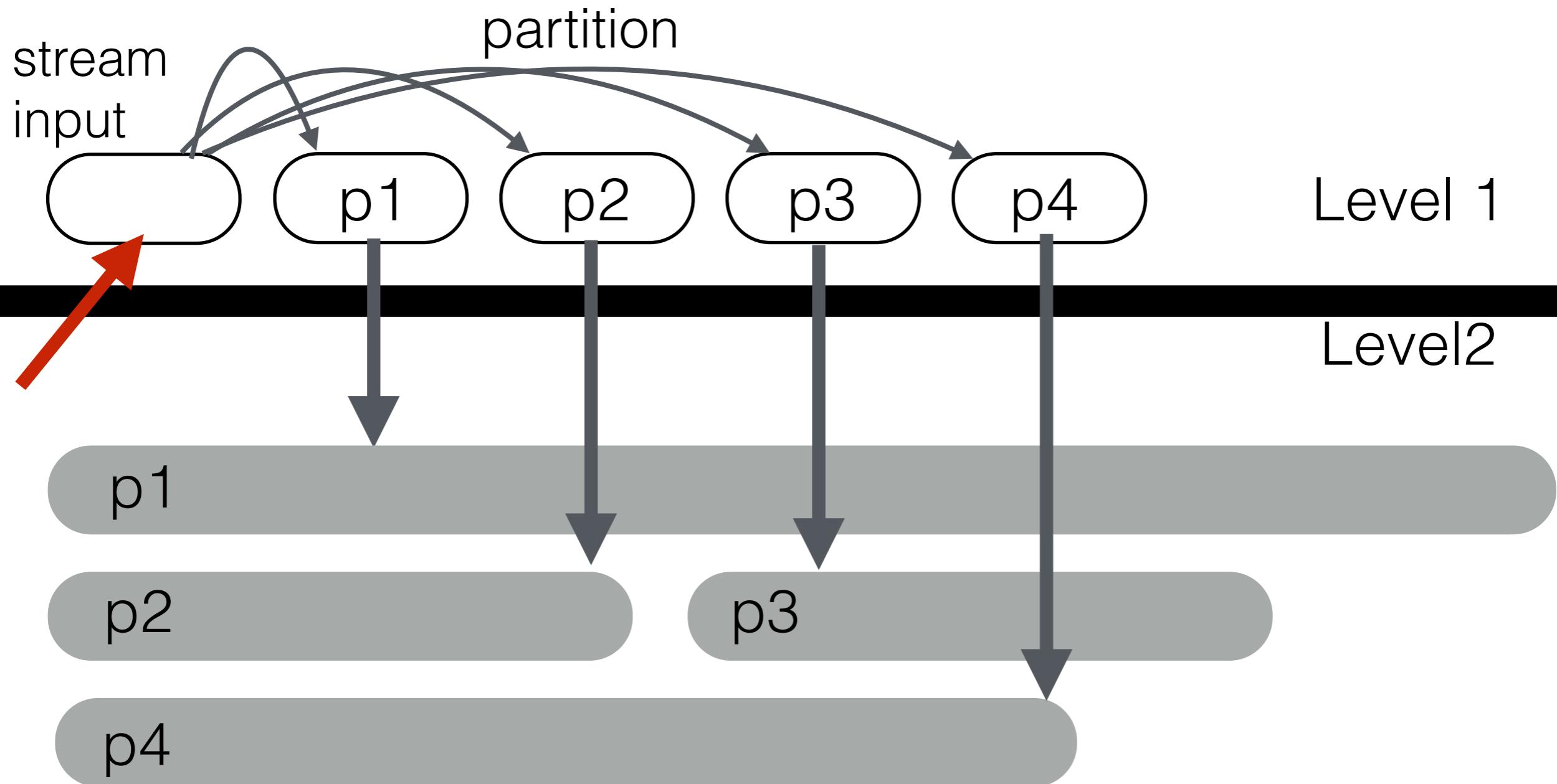
hash join



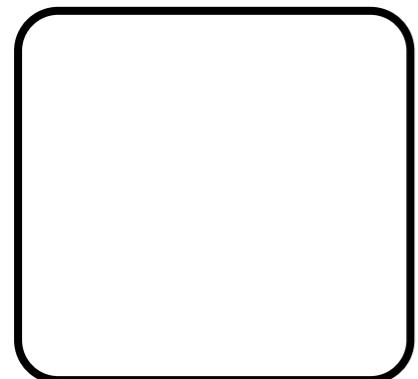


1. read input into stream buffer, hash and write to respective partition buffer
2. when input buffer is consumed, bring the next one
3. when a partition buffer is full, write to L2

we can partition into L1-1 pieces in one pass



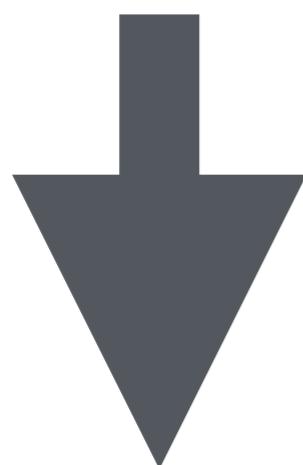
join input 1



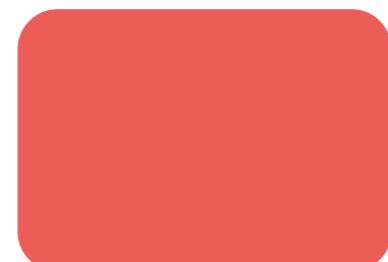
join input 2



grace hash join



hash partitioning



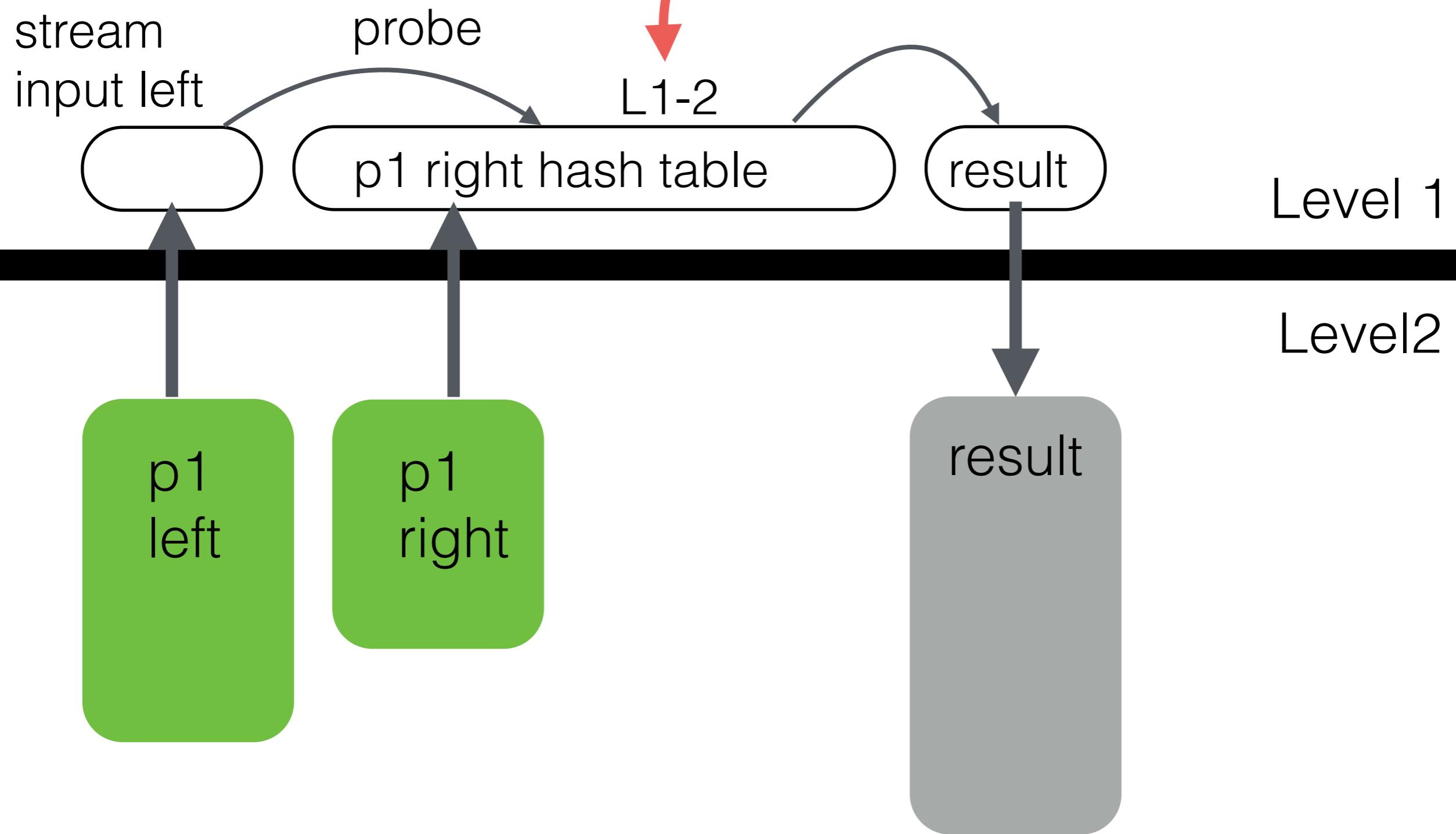
**then join each pair of partitions
independently in memory**



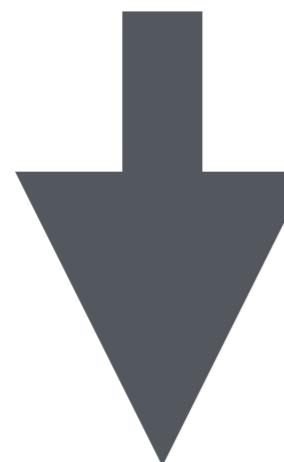
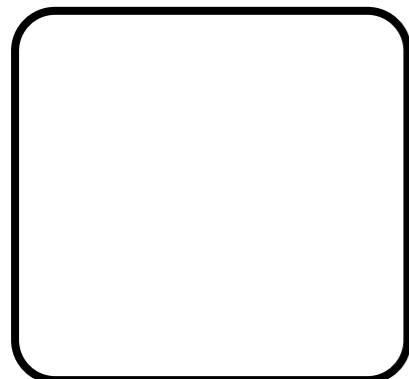
HARVARD
School of Engineering
and Applied Sciences

CS165, Fall 2015
Stratos Idreos

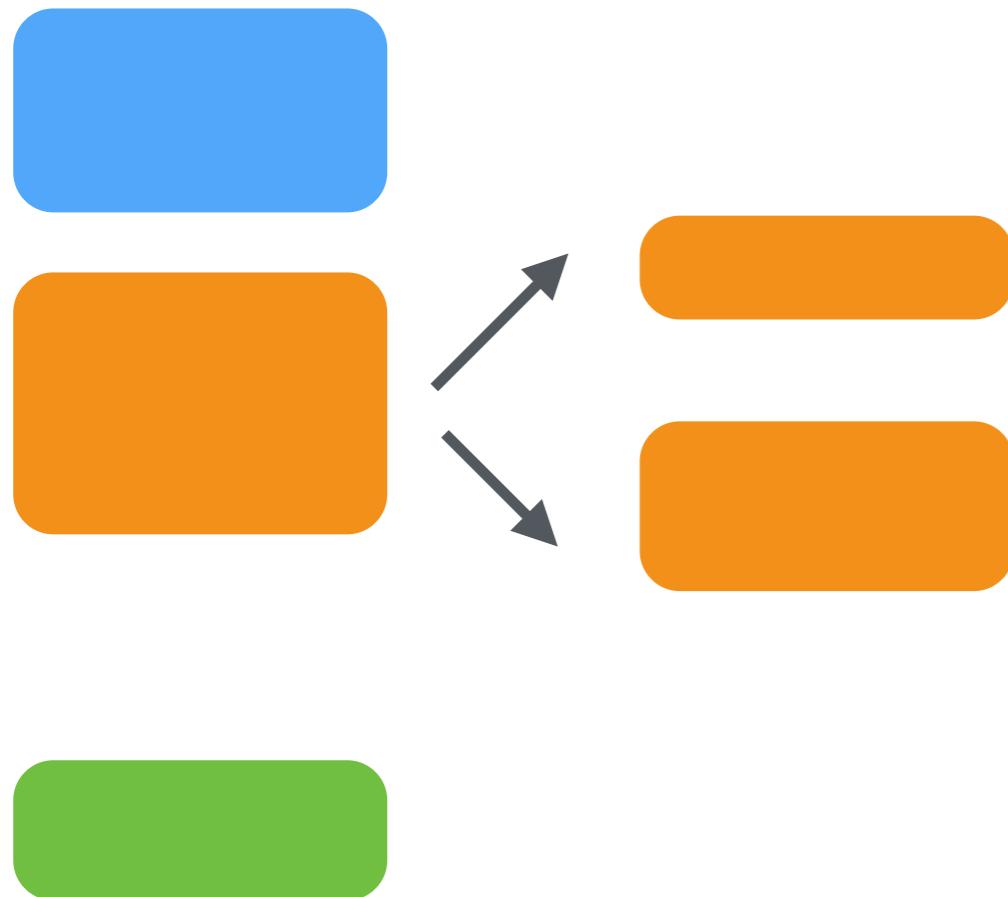
as long as at least one
of the pieces $\leq L1-2$



grace hash join



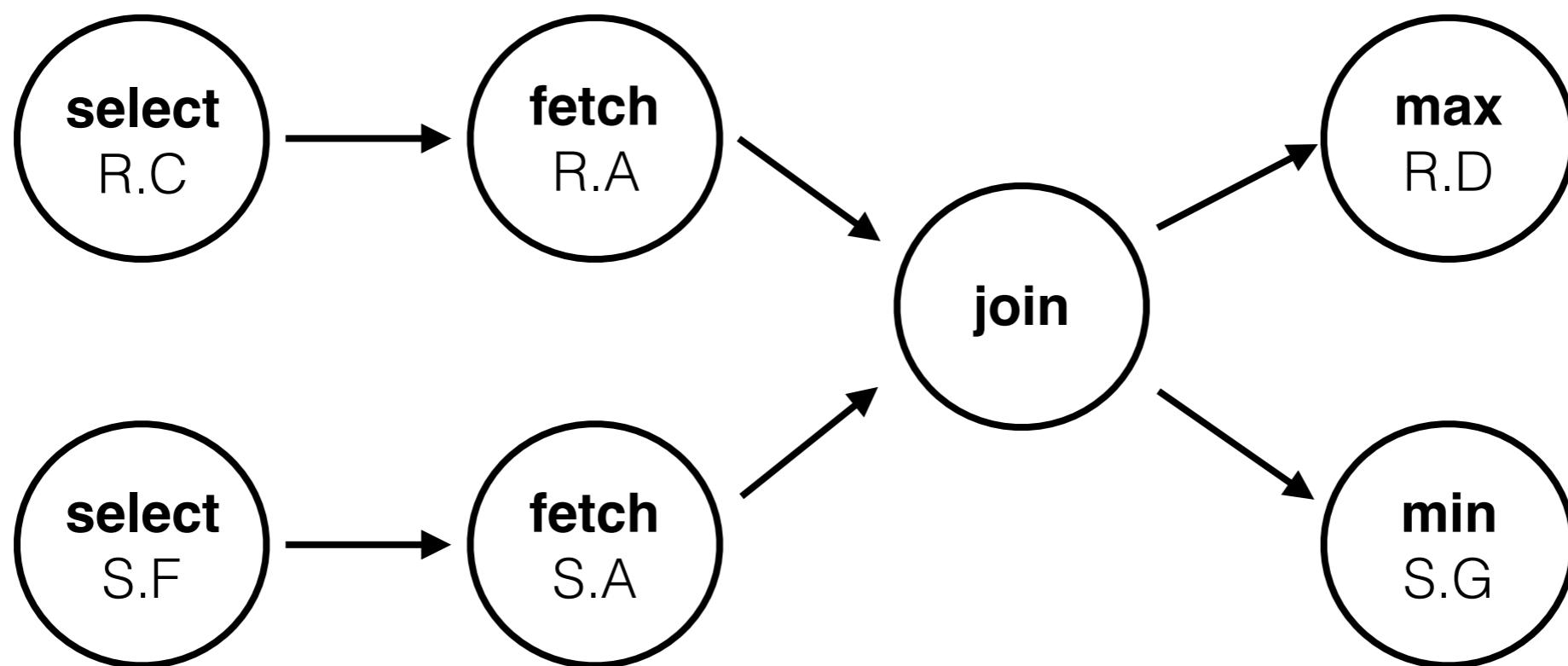
hash partitioning



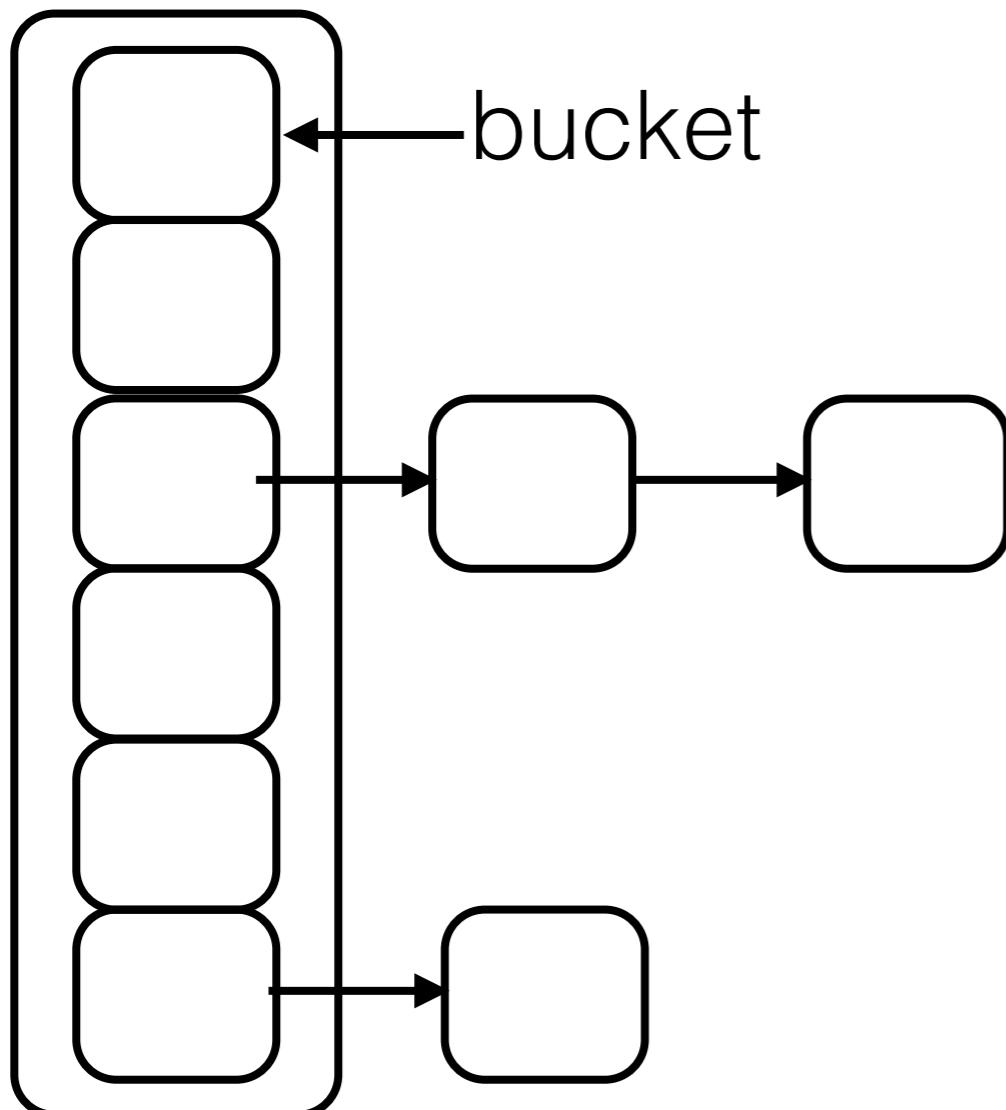
apply recursively if a partition does not fit in memory (L1-2)



today: how to create the hash table

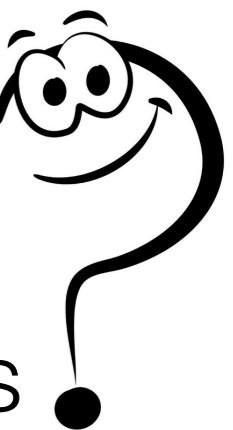


hash table

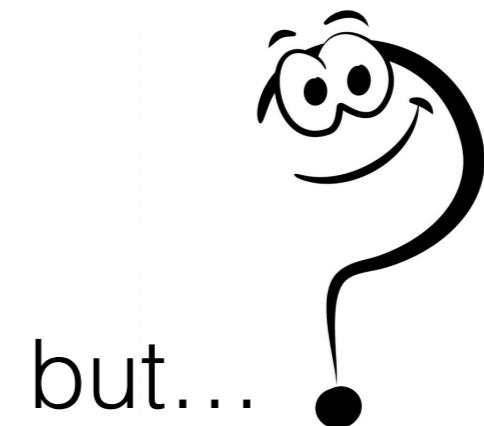


updates may create long lists

can this happen anyway,
i.e., even with no updates



simple solution would be to
double the hash table and rehash everything...

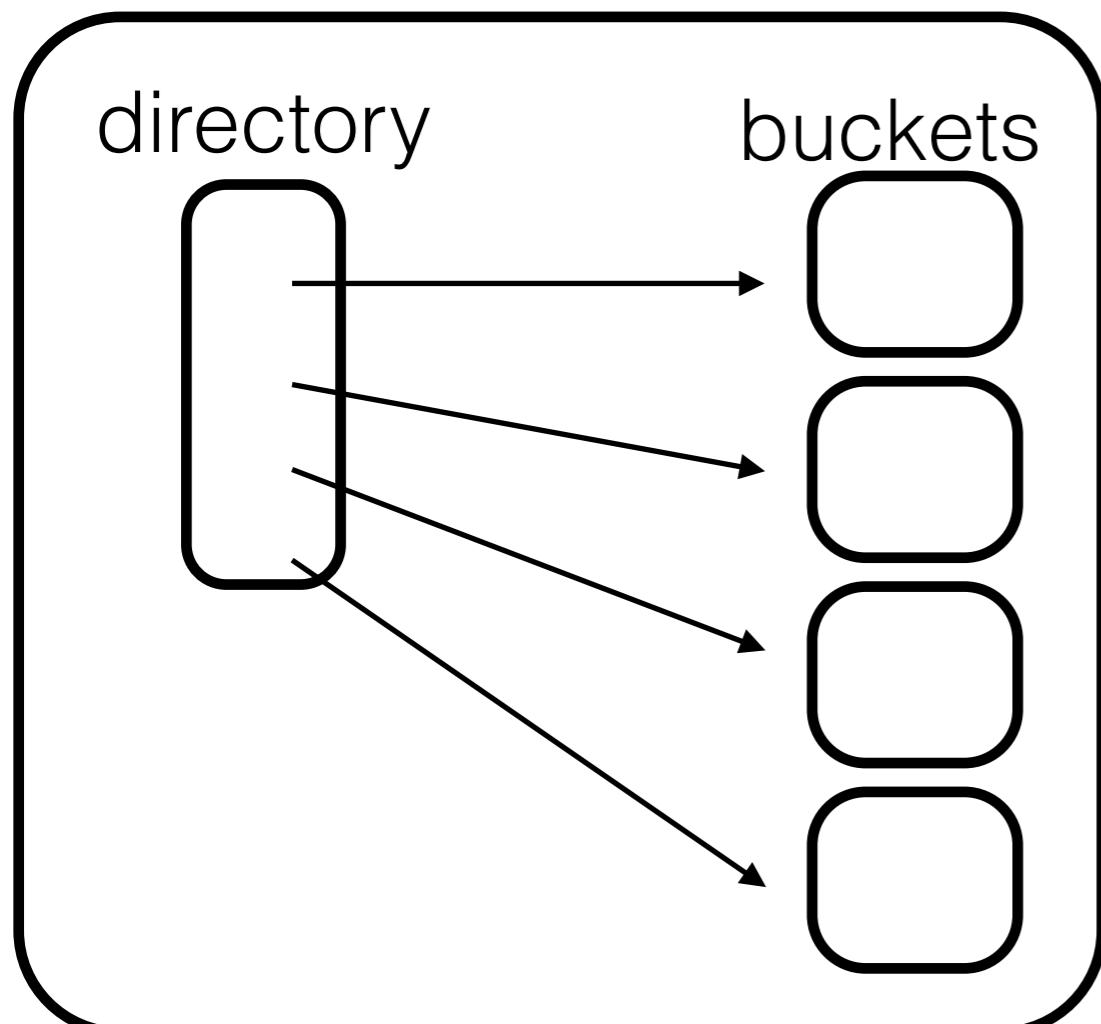


dynamic hashing

extendible hashing - linear hashing

extendible hashing

hash table

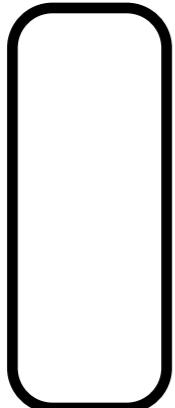


when there is no more space
split individual buckets
and if needed **double the directory**

with directory in memory
we need 1 I/O to fetch the bucket

not true from memory to caches in general

directory



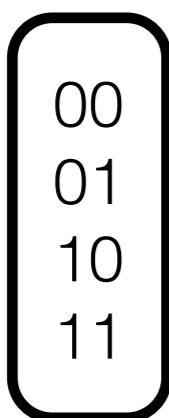
maps hash values to buckets

$h = \text{hash}(\text{key})$

use binary form of h

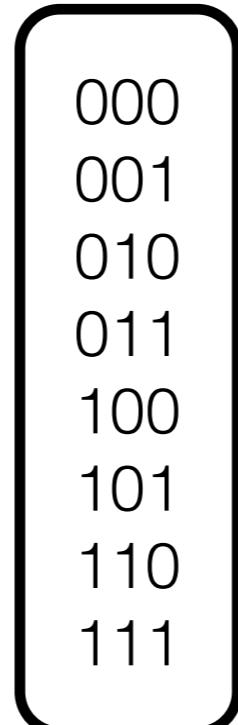
use X last bits to map 2^X buckets

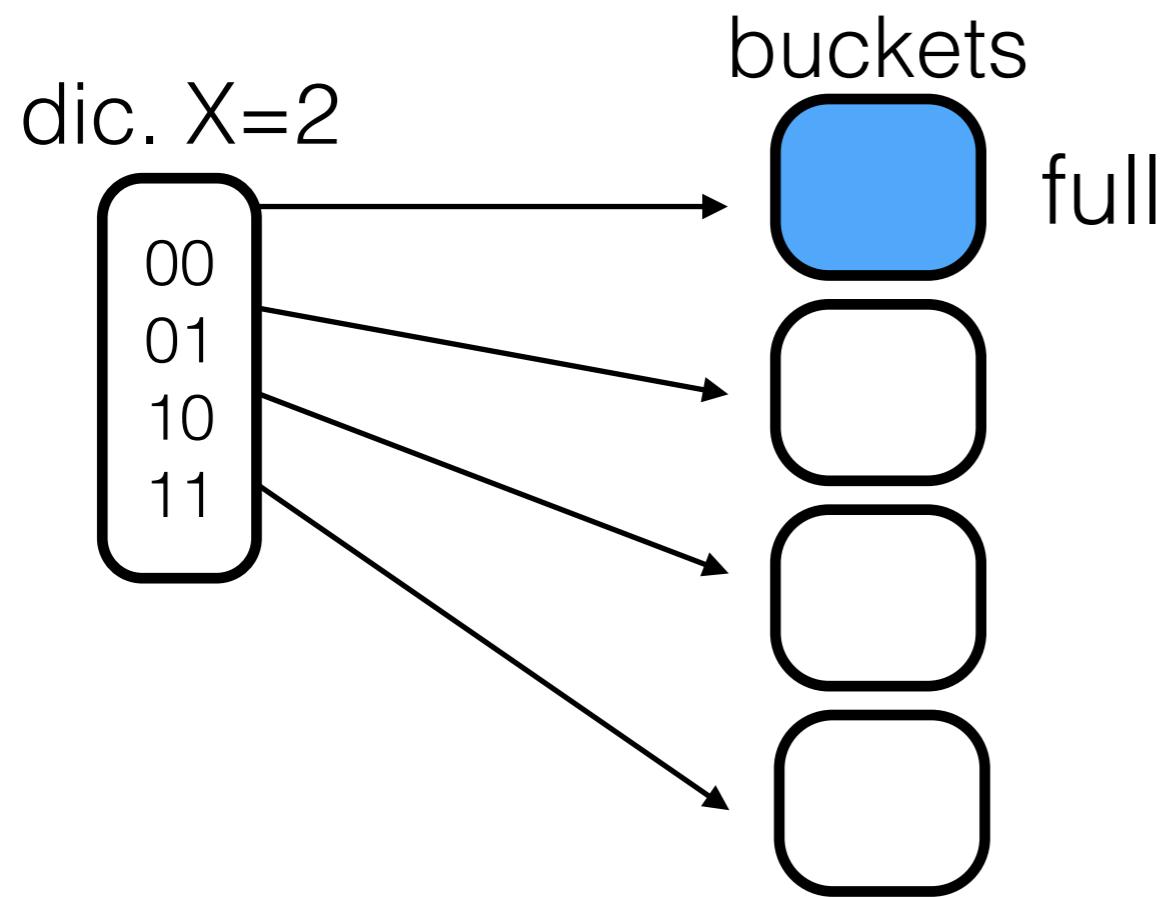
for $X=2$

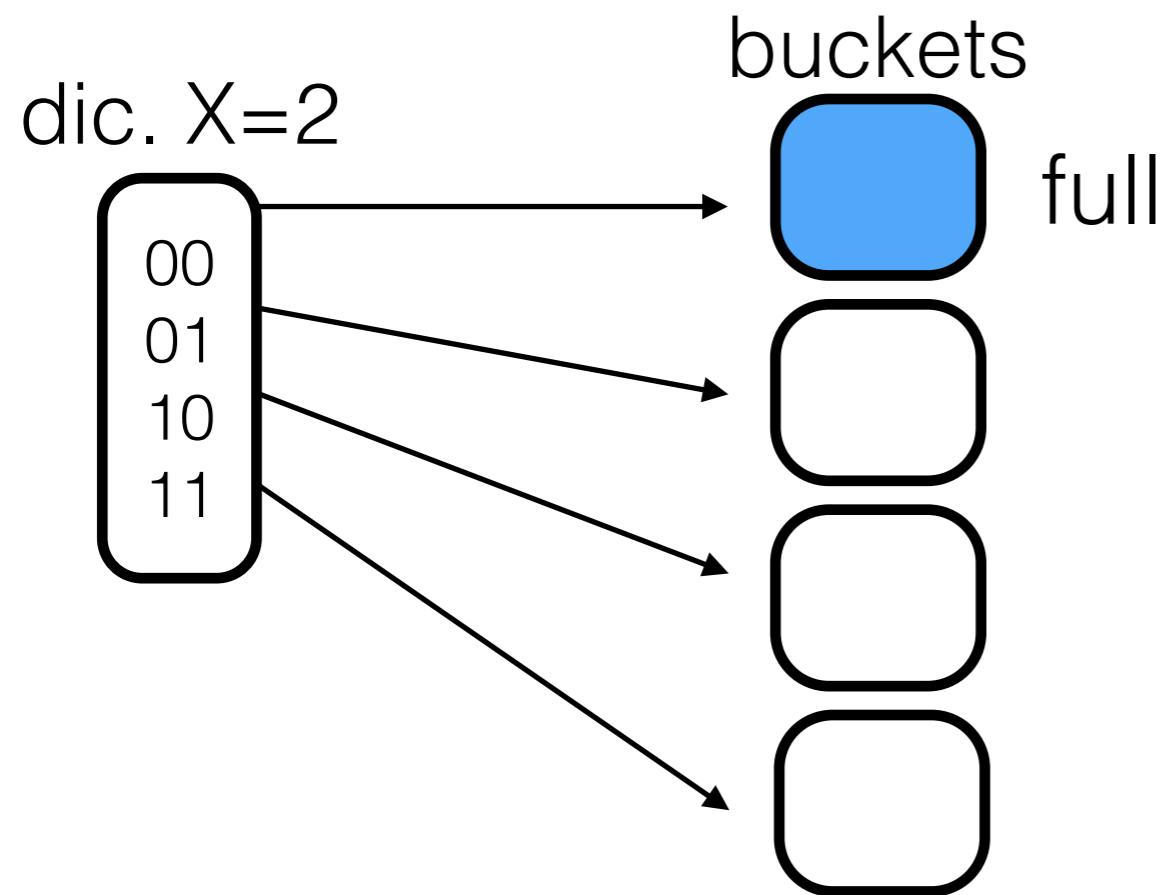


extending
the directory →

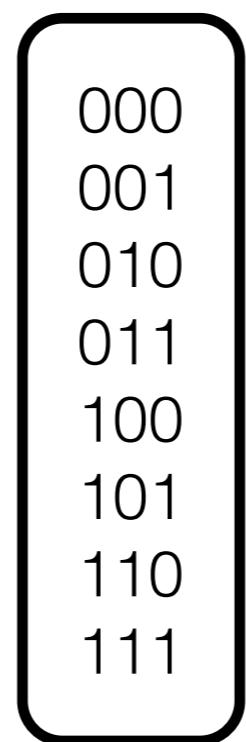
for $X=3$

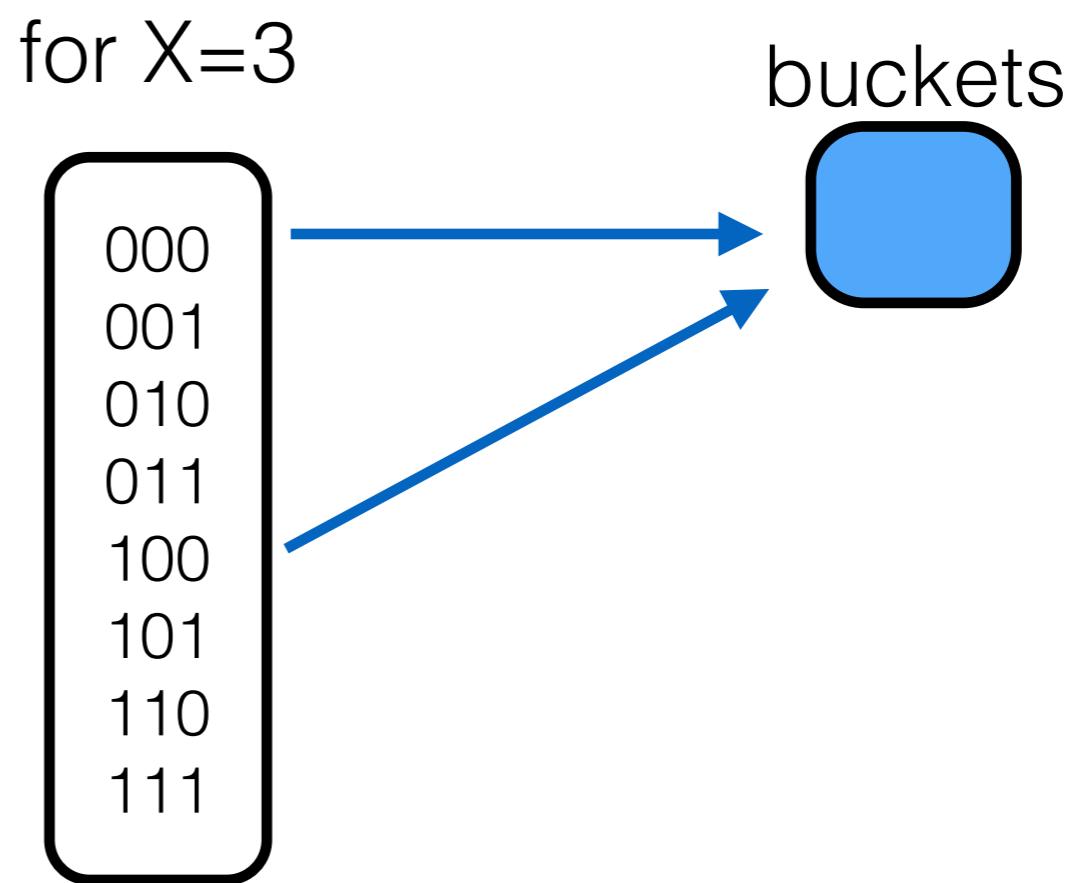
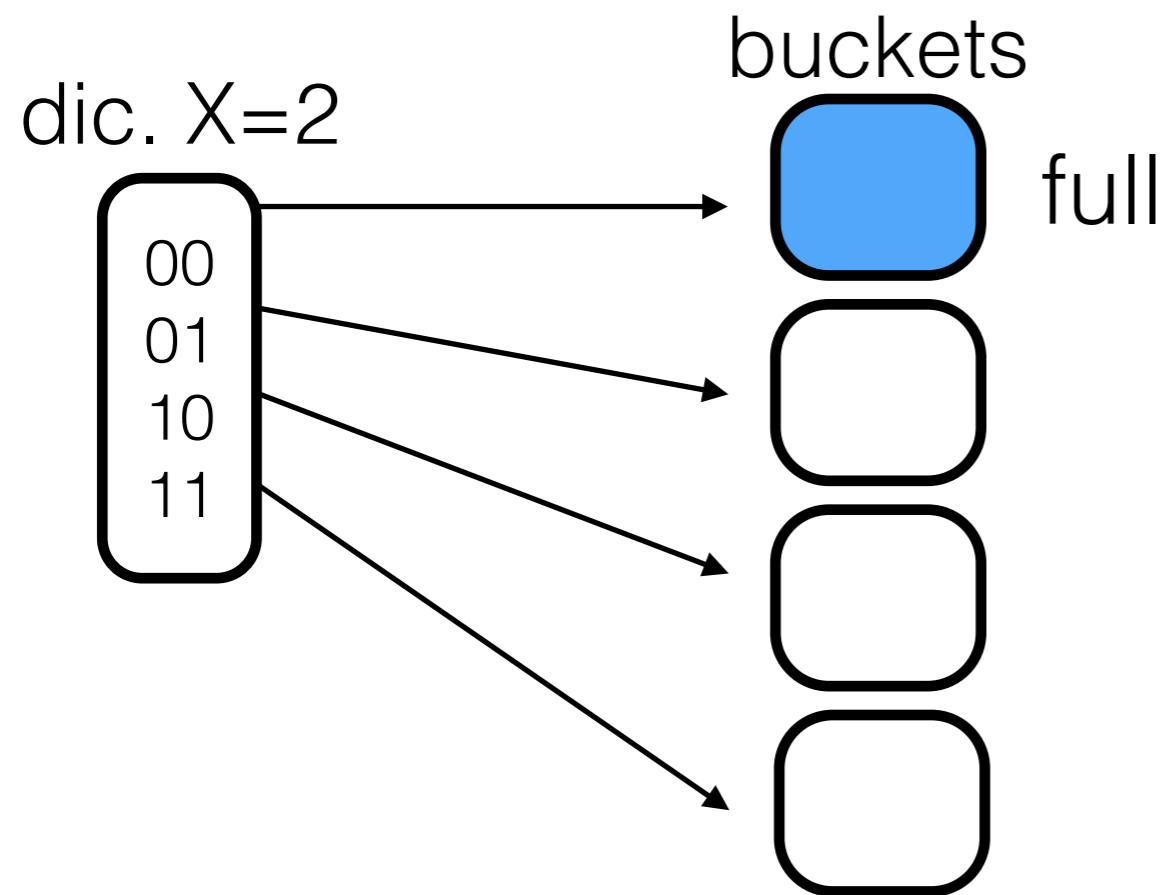


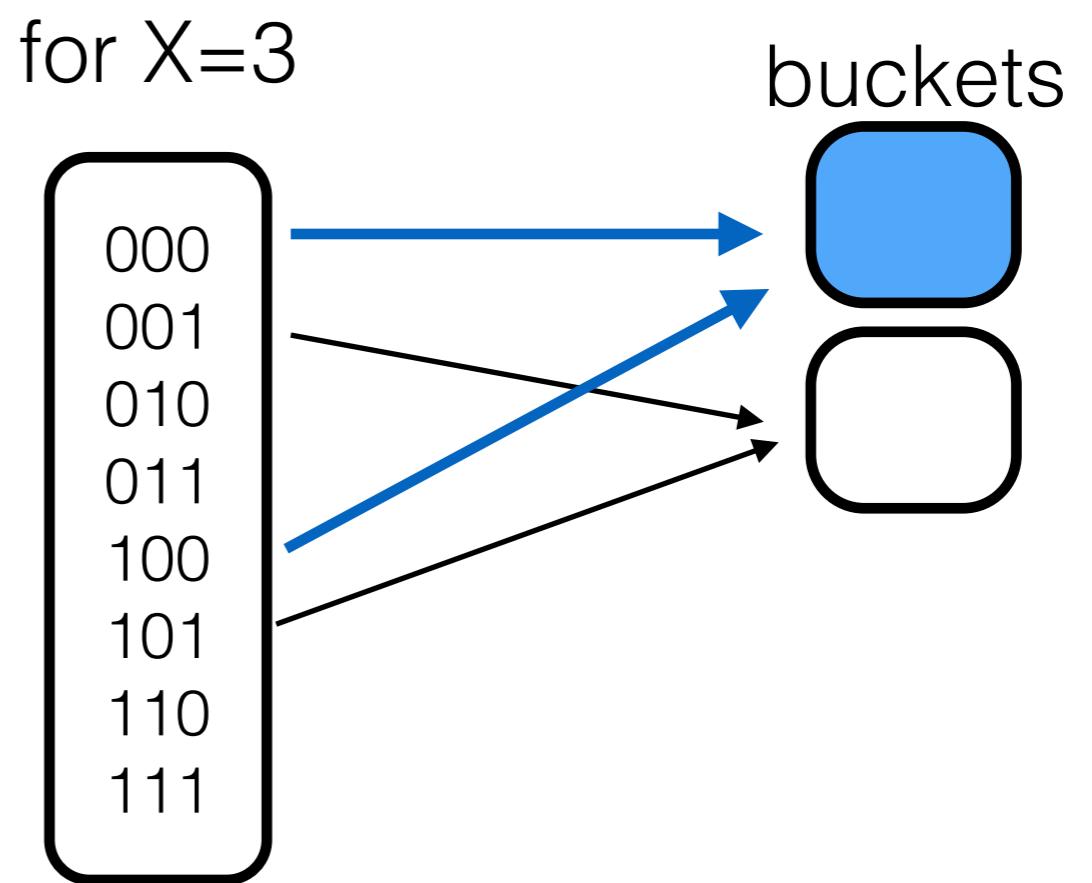
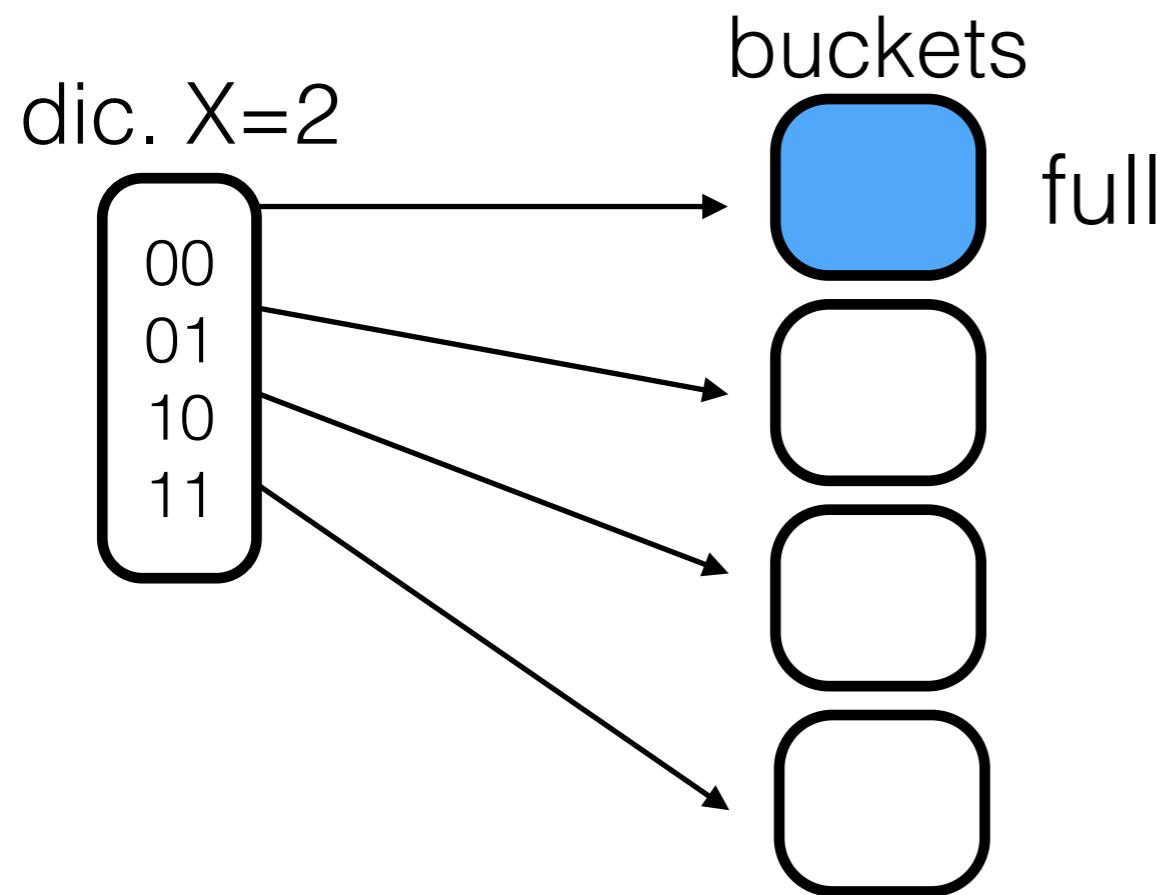


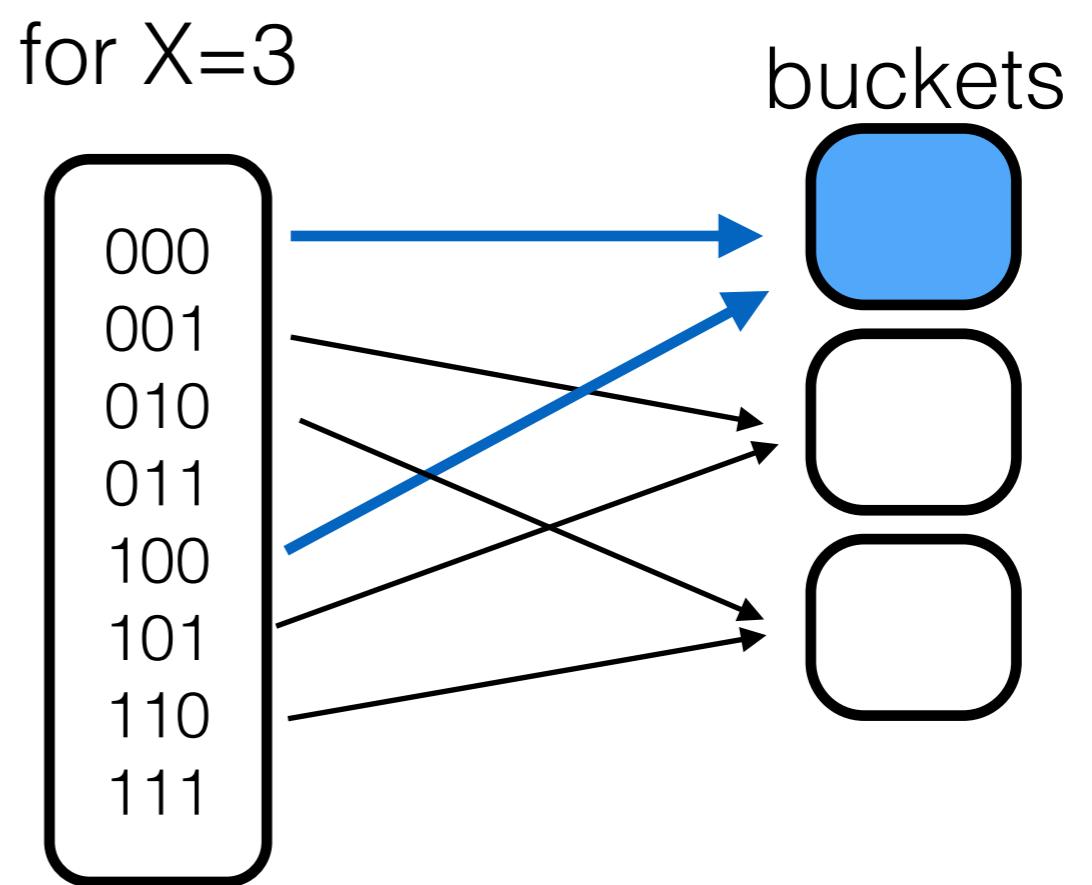
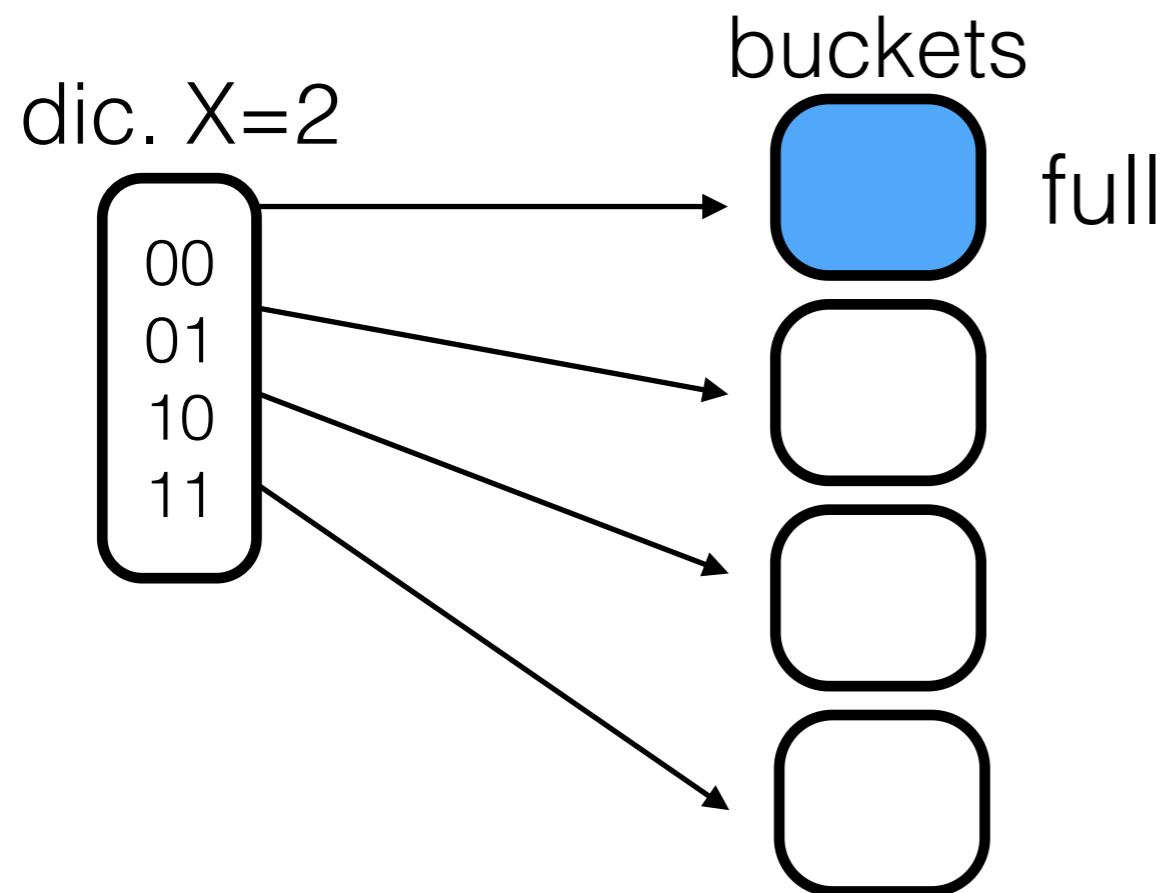


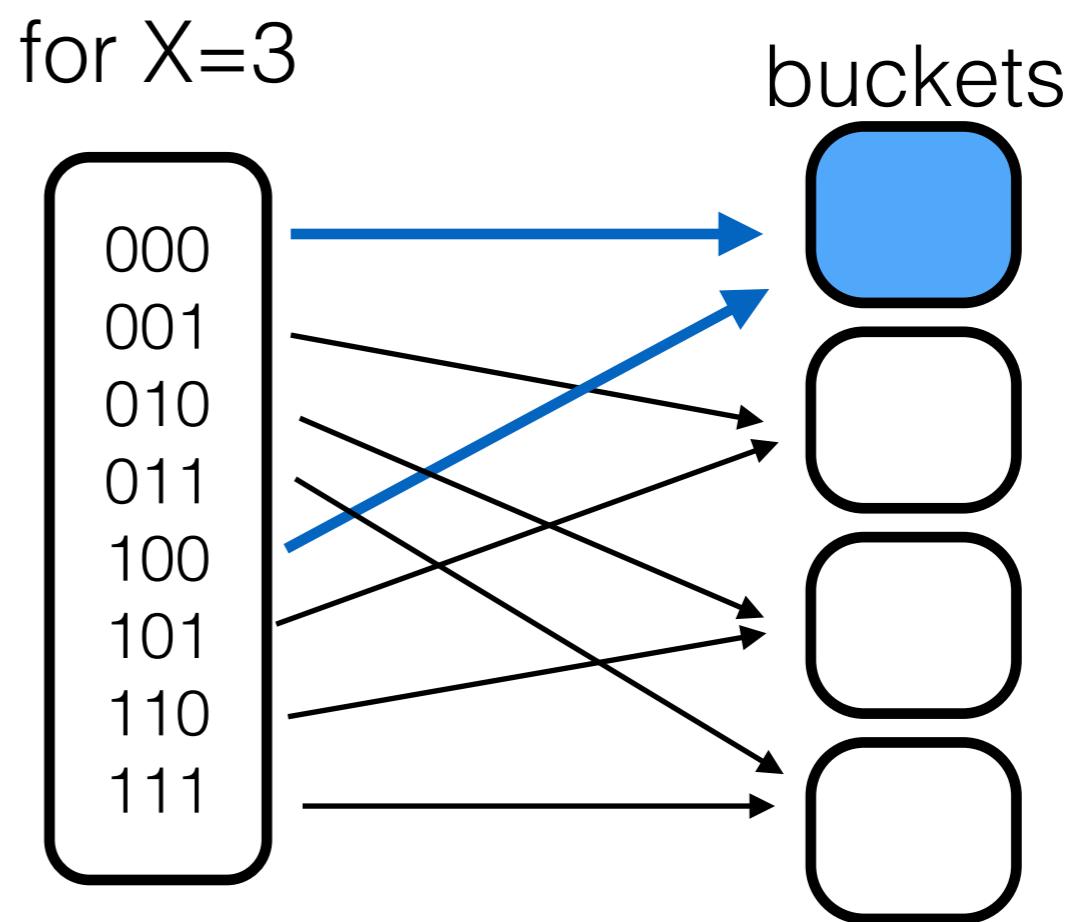
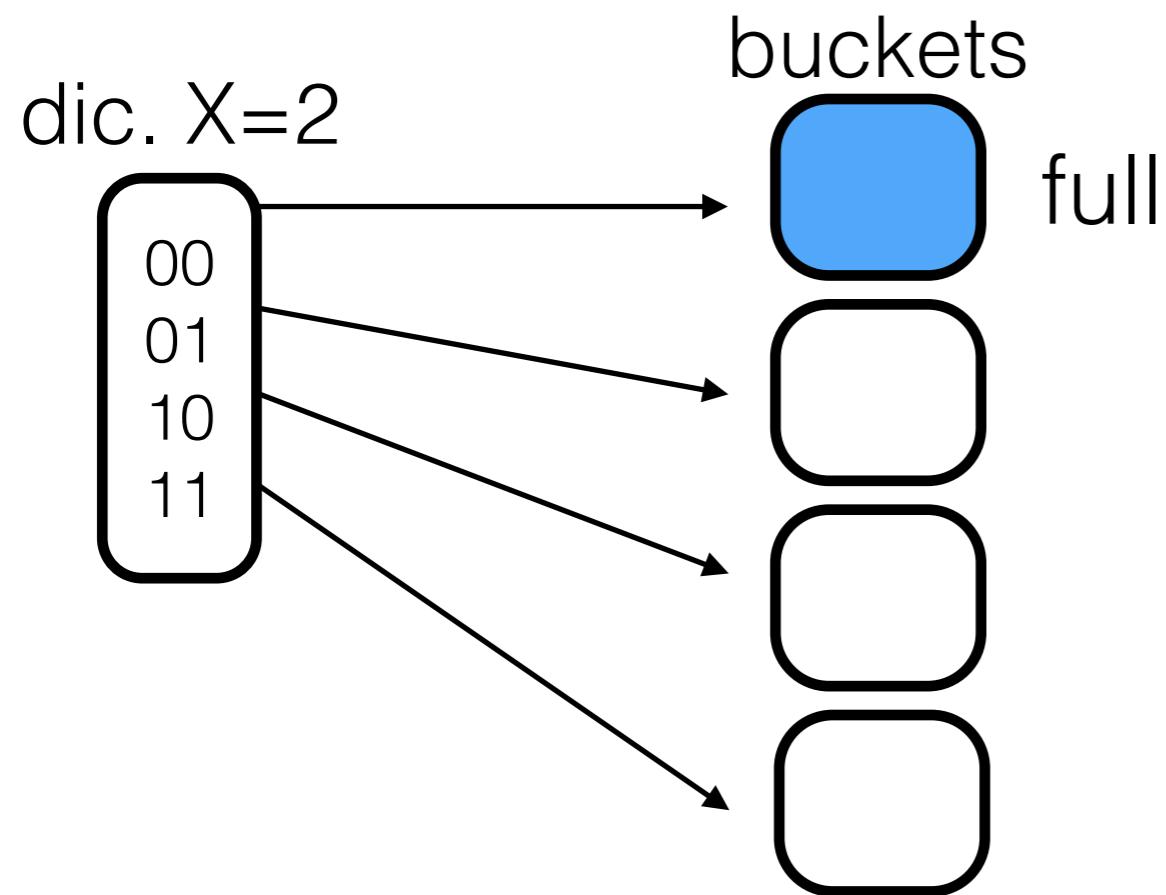
for X=3

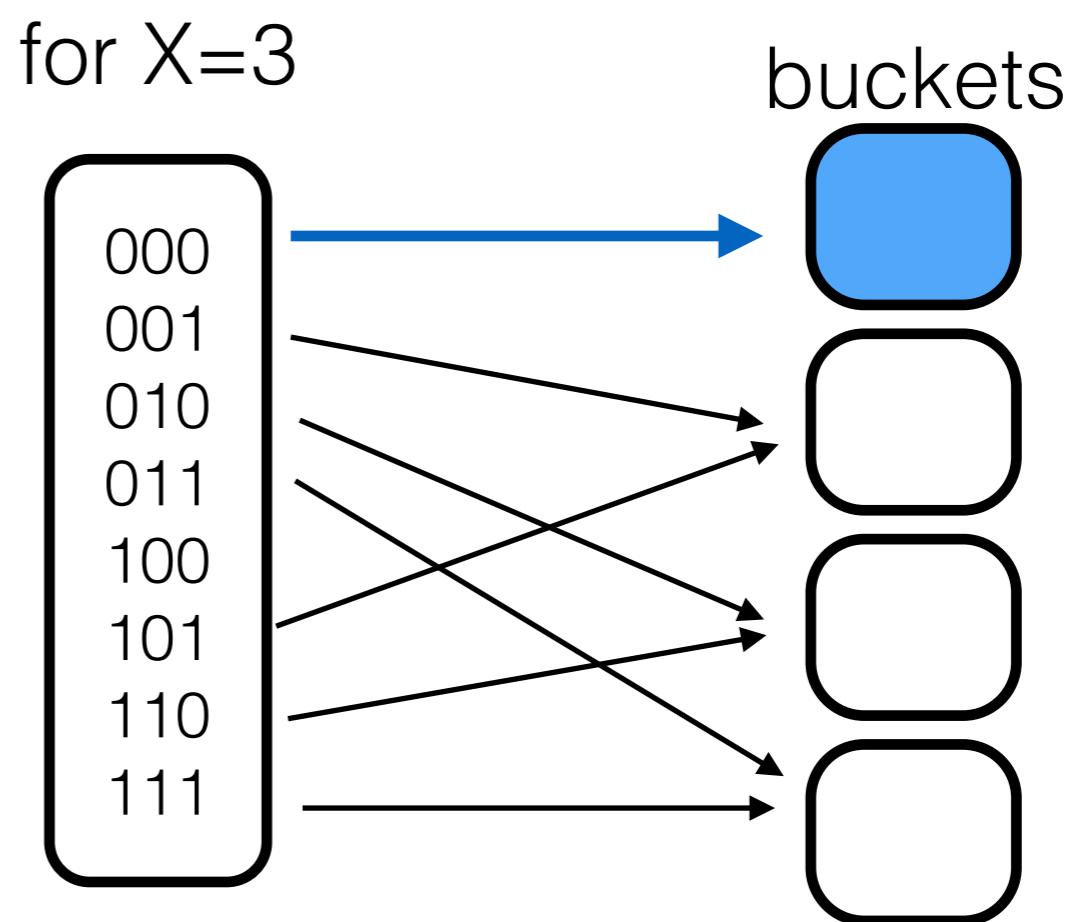
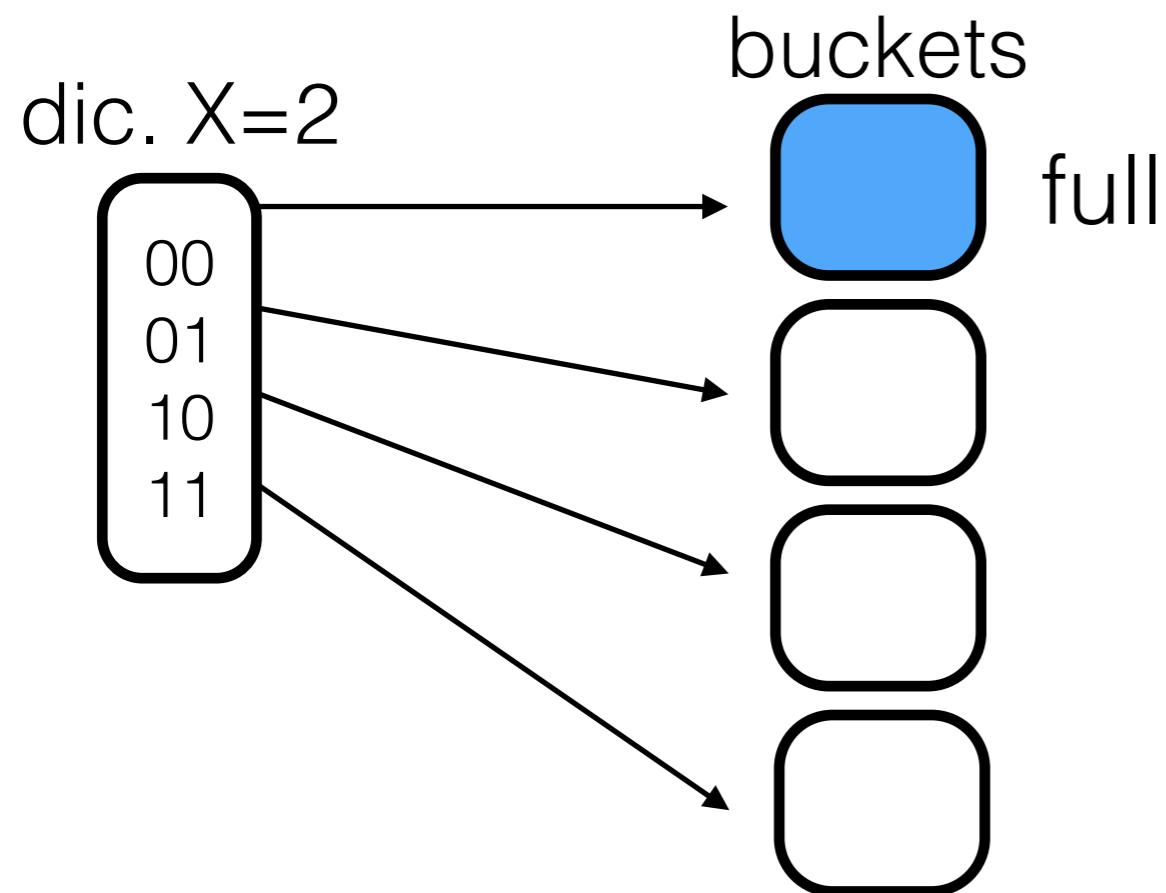


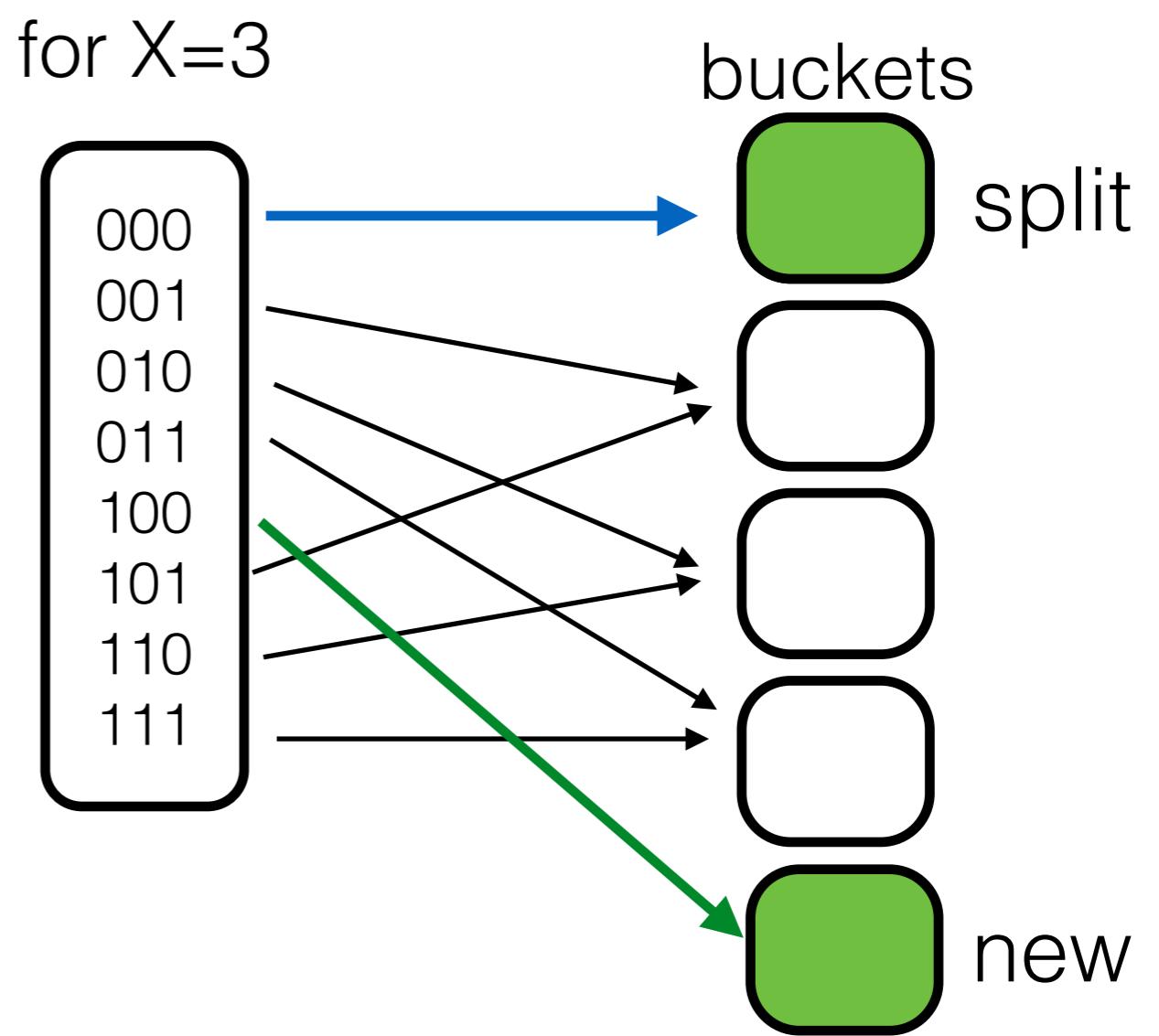
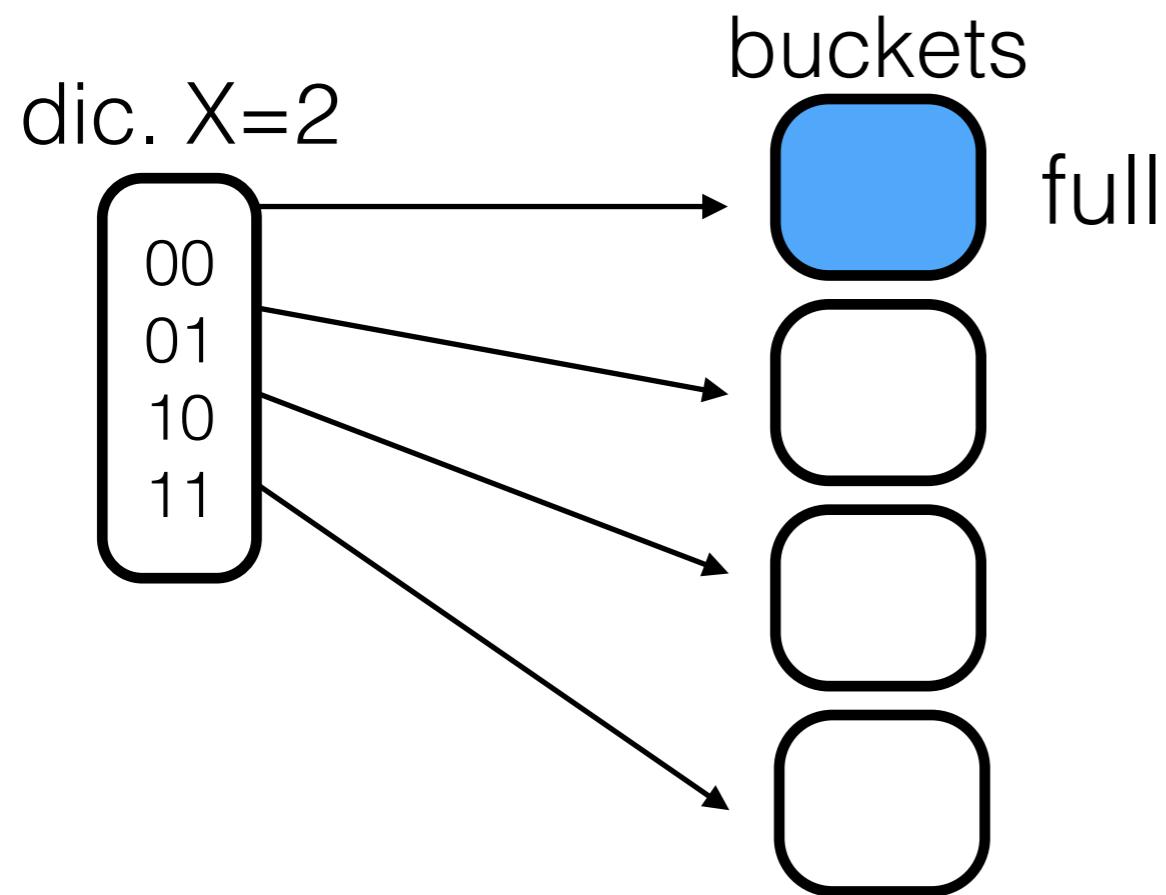


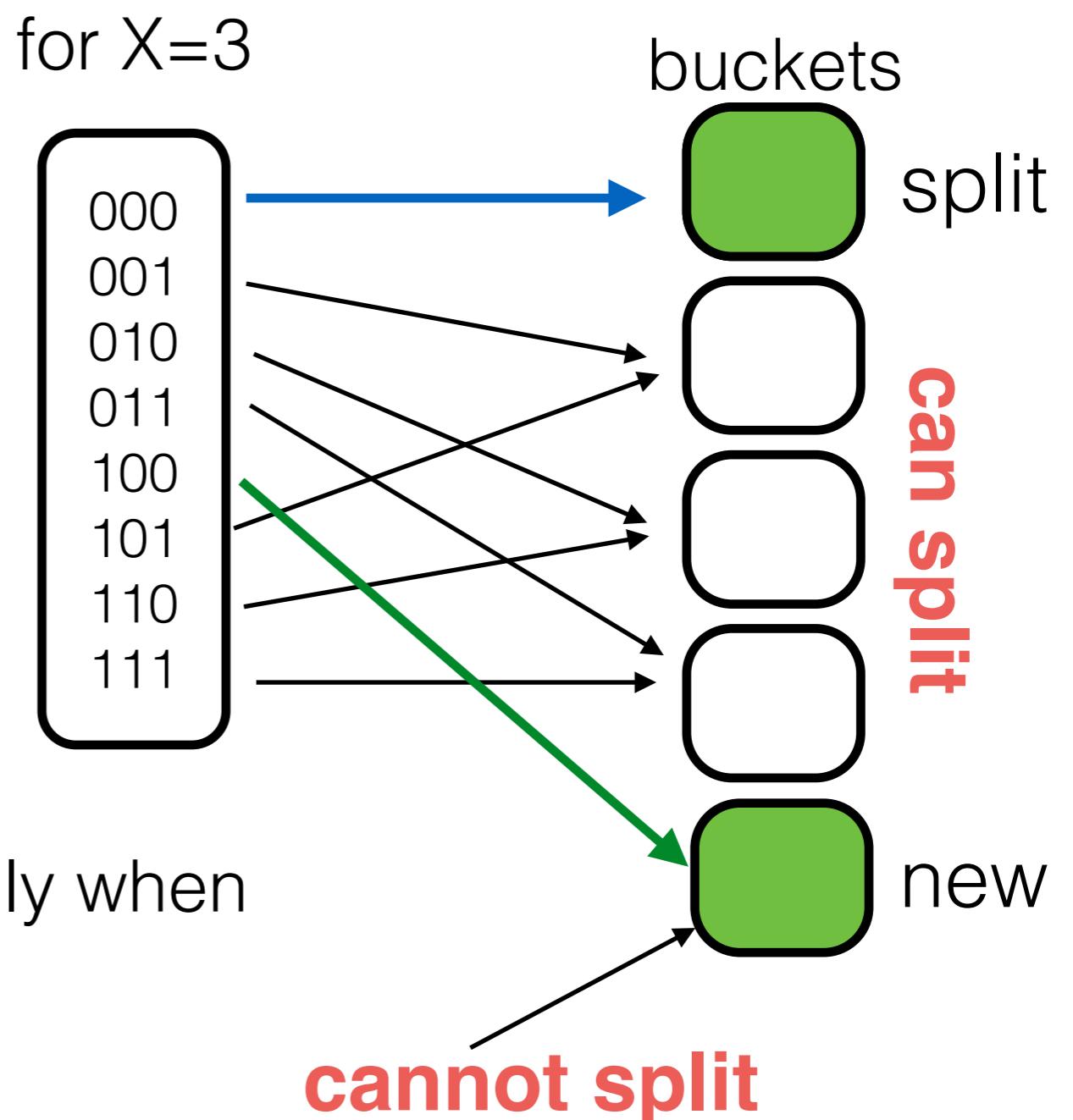
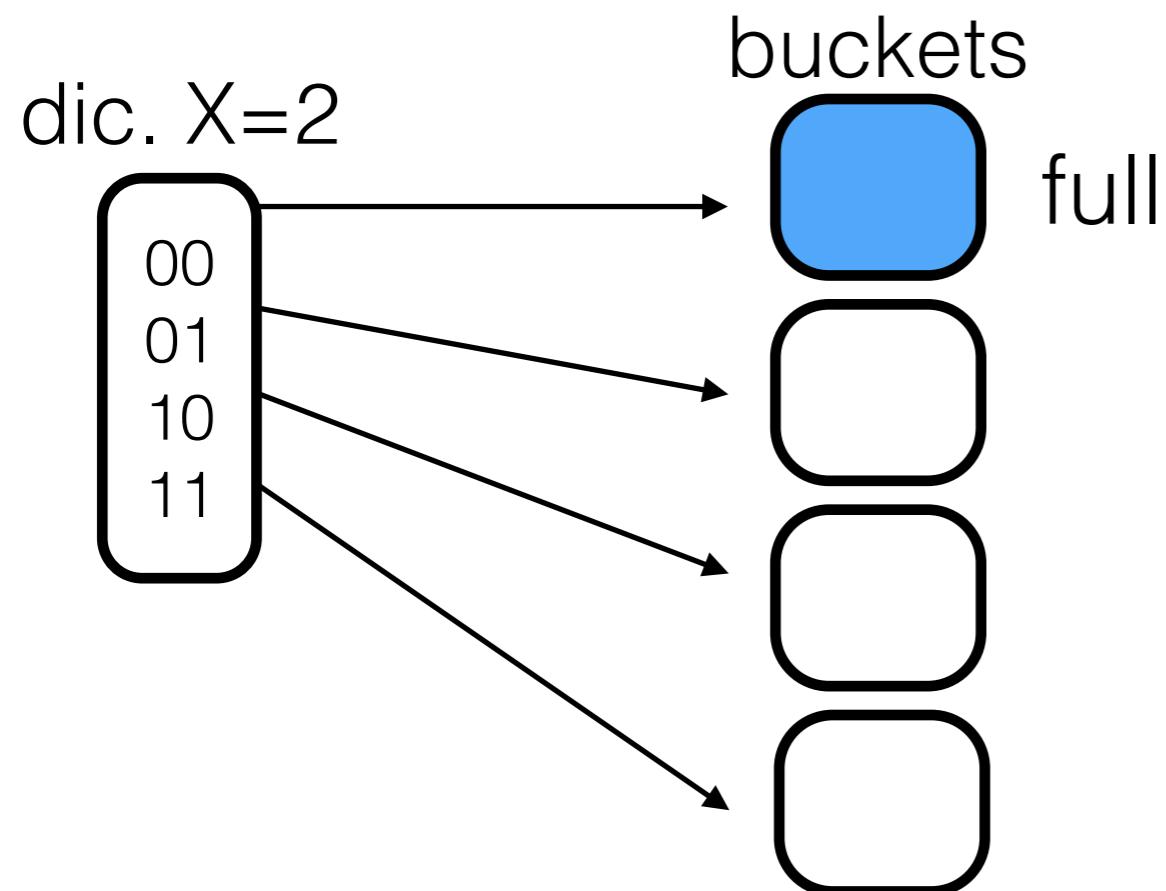












extend directory only when
no split is possible

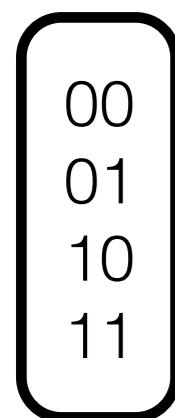
cannot split

linear hashing

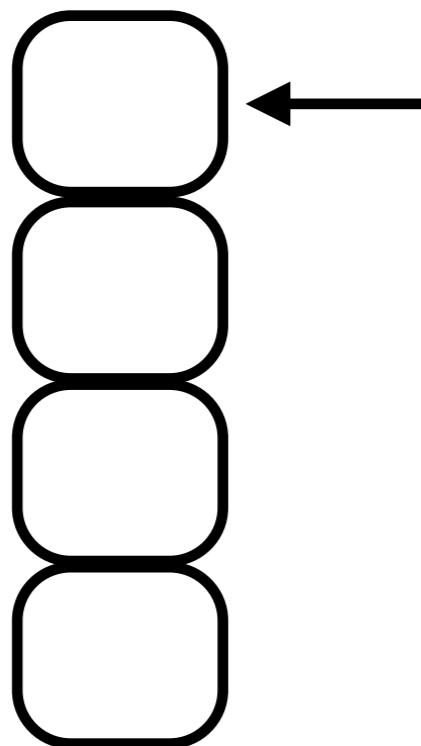
do not use a directory (indirection cost)

instead incrementally extend
the hash table one bucket at a time

e.g., use 2 bits of hash value (h_0)



dictionary is implicit

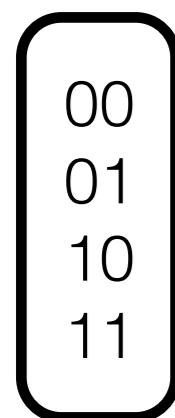


next bucket
to split

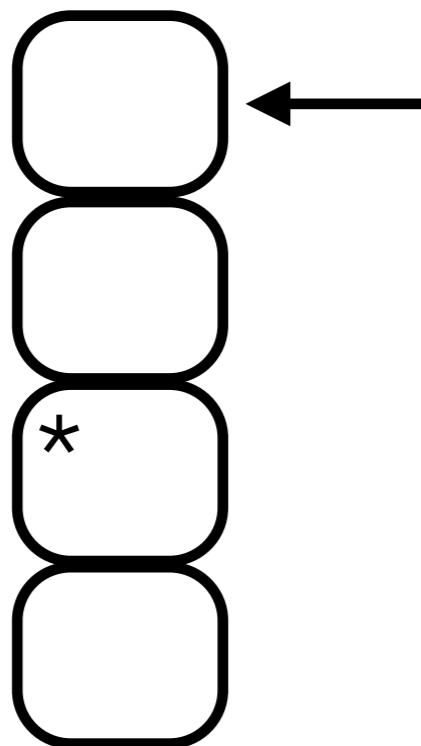
when insert causes overflow somewhere split *next bucket*

use +1 bit of hash value

e.g., use 2 bits of hash value (h_0)



dictionary is implicit



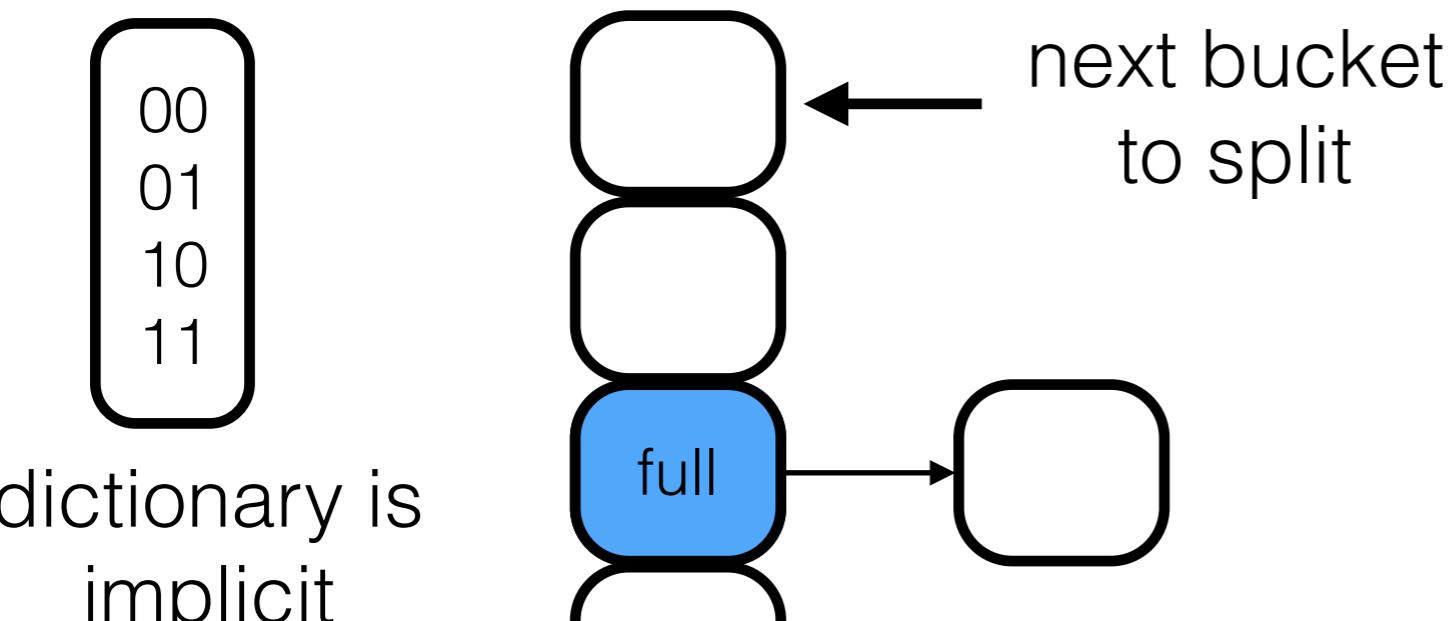
next bucket
to split

when insert causes overflow somewhere split *next bucket*

use +1 bit of hash value



e.g., use 2 bits of hash value (h_0)



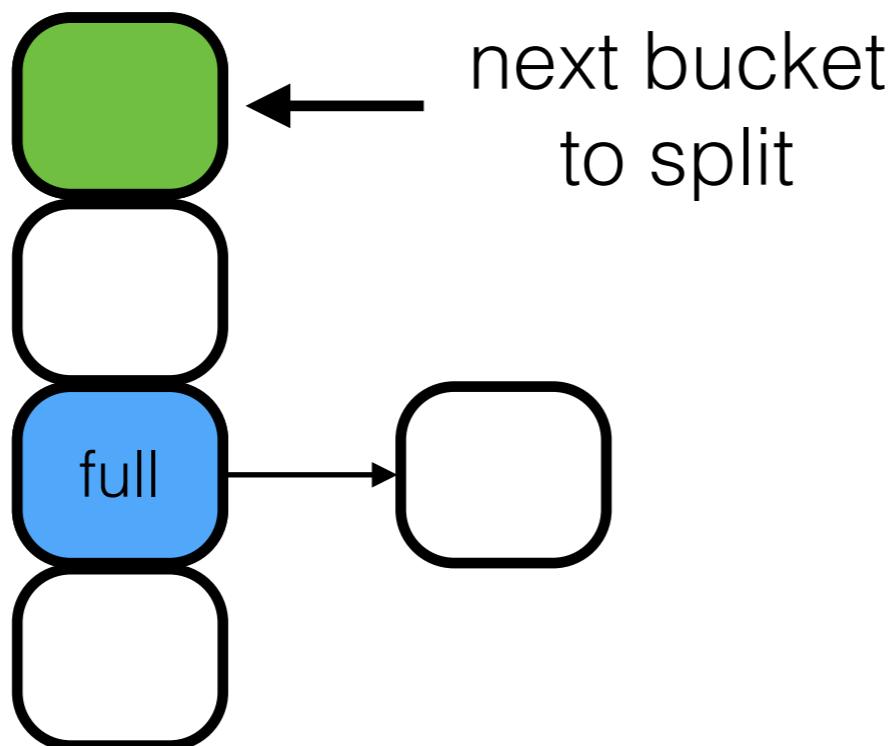
when insert causes overflow somewhere split *next bucket*

use +1 bit of hash value

e.g., use 2 bits of hash value (h_0)

00
01
10
11

dictionary is implicit



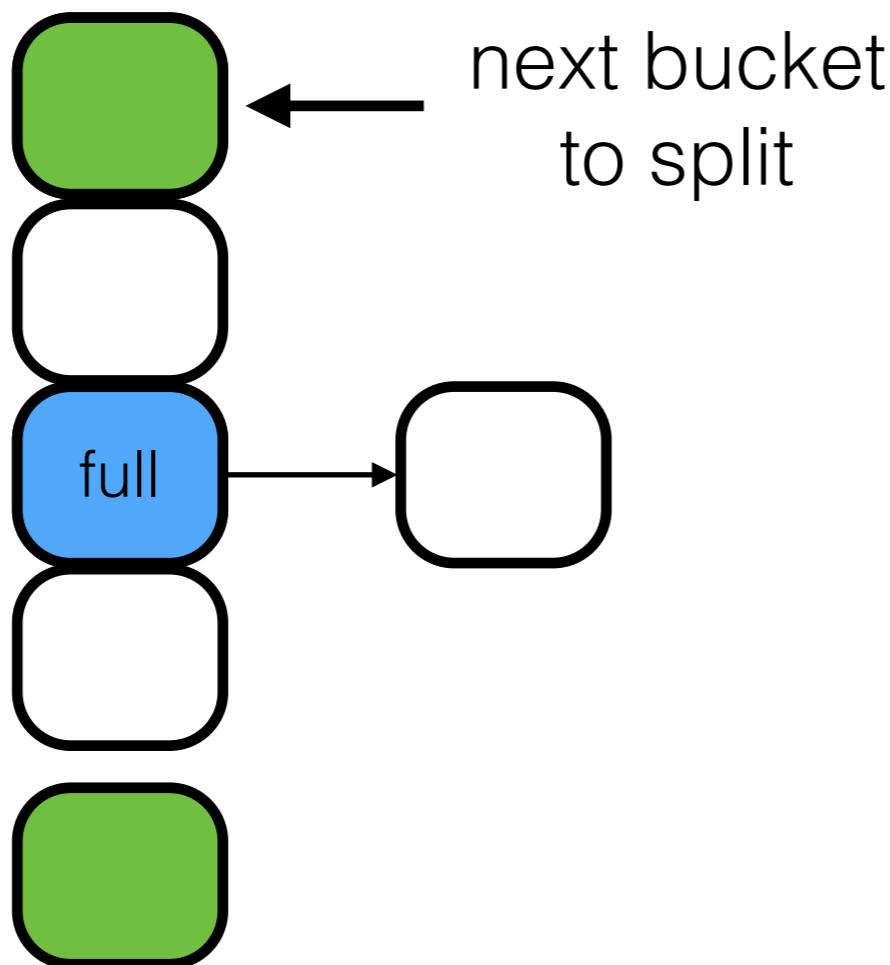
when insert causes overflow somewhere split *next bucket*

use +1 bit of hash value

e.g., use 2 bits of hash value (h_0)

00
01
10
11

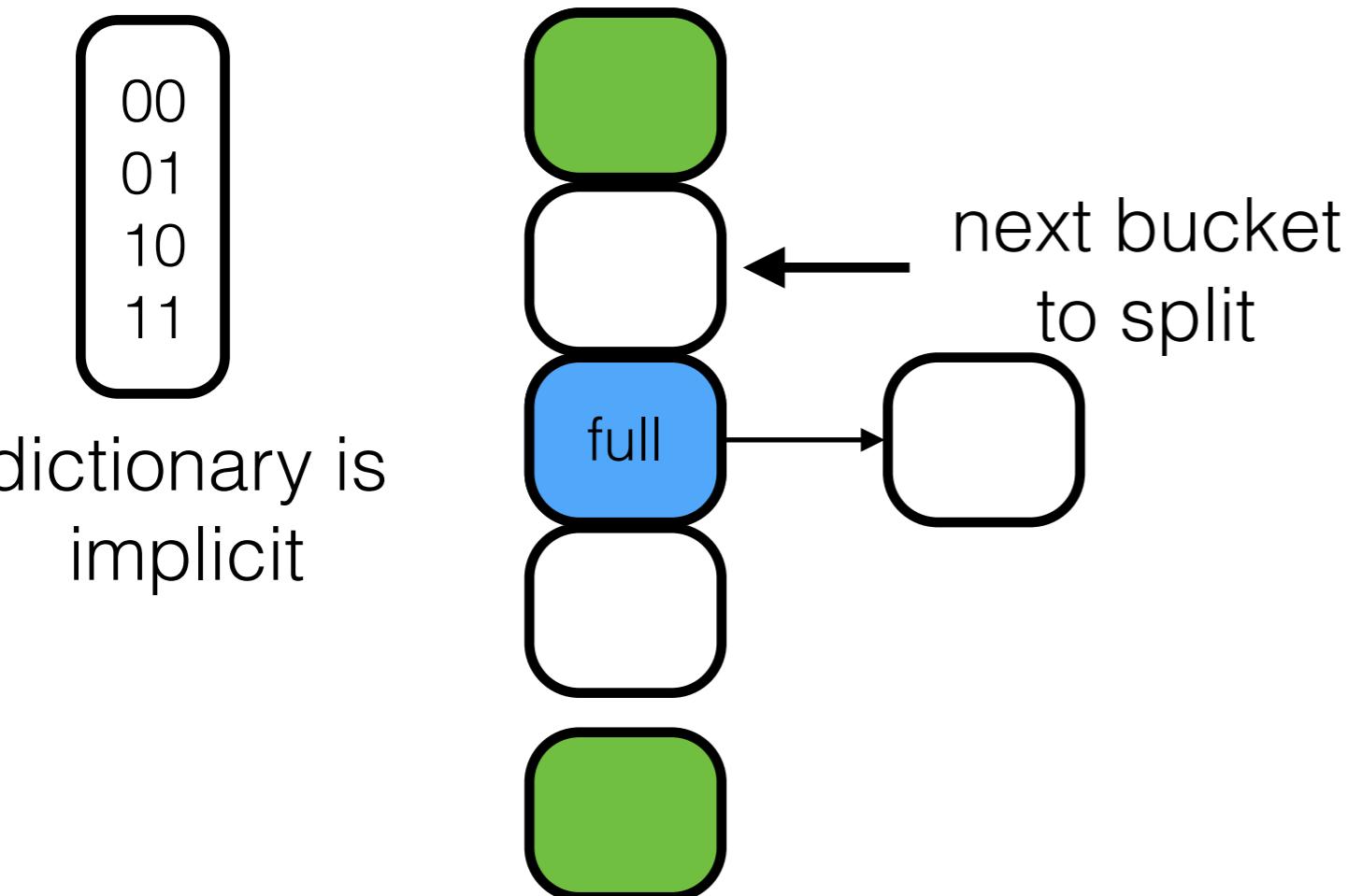
dictionary is implicit



when insert causes overflow somewhere split *next bucket*

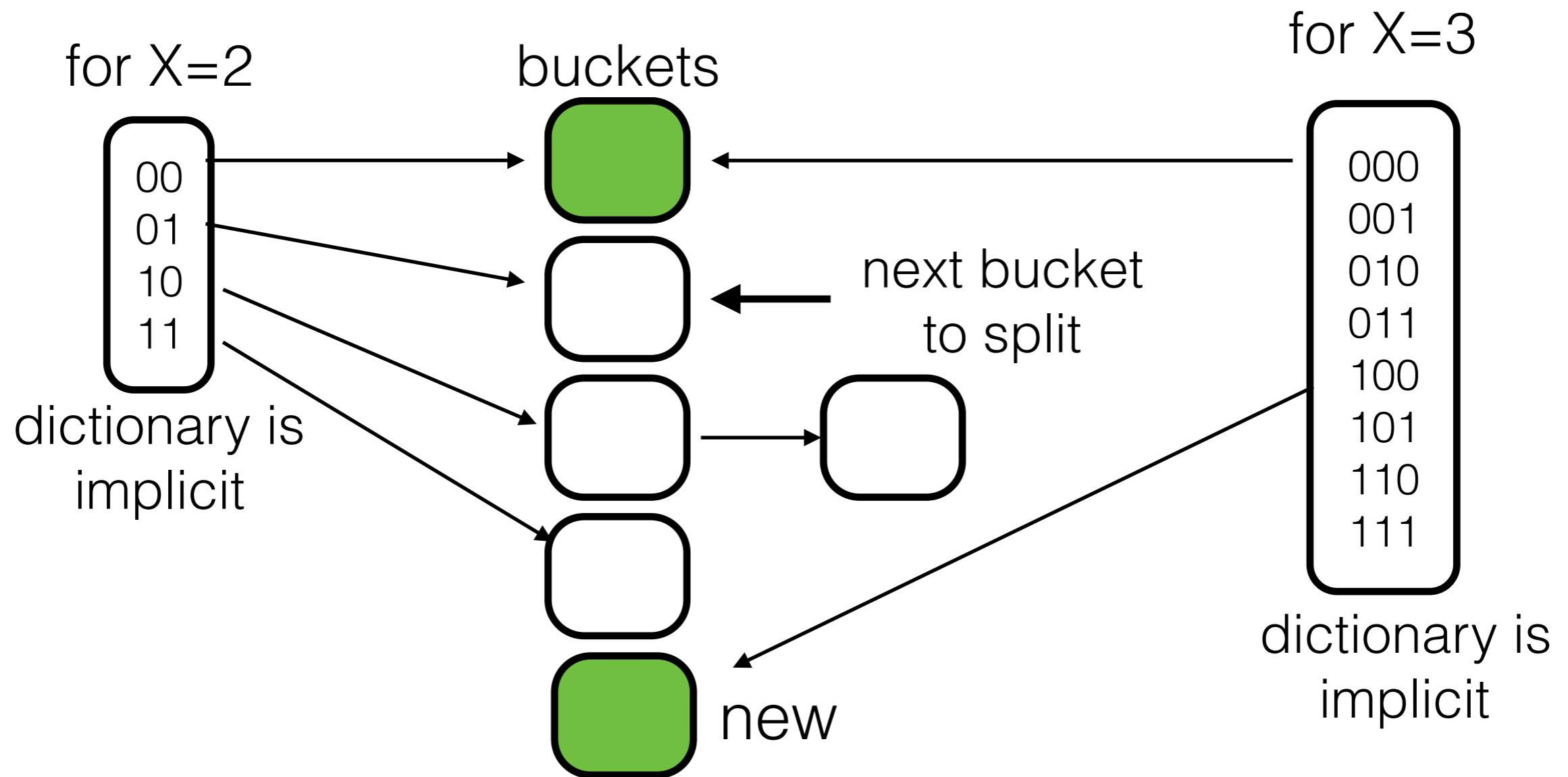
use +1 bit of hash value

e.g., use 2 bits of hash value (h_0)

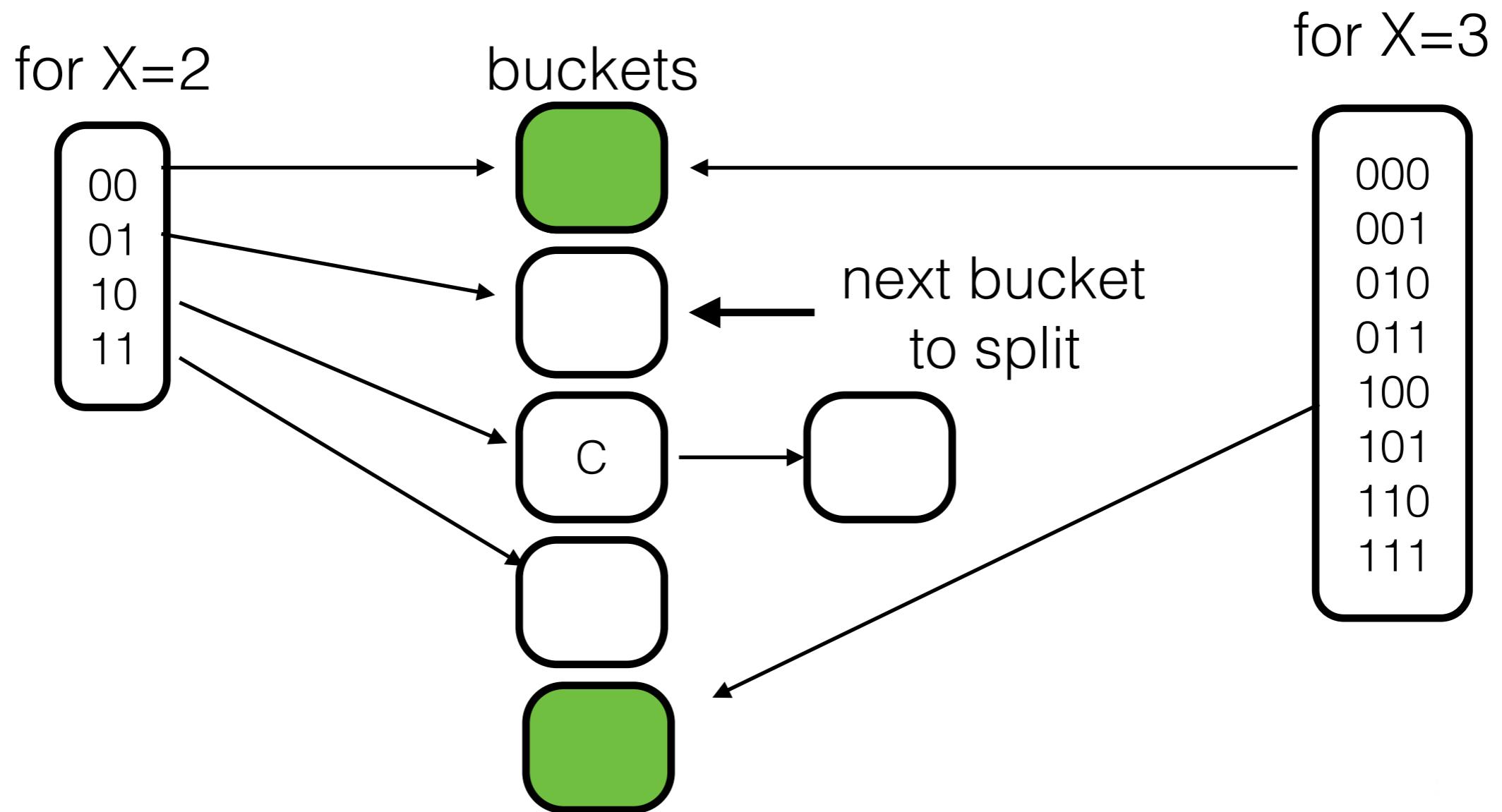


when insert causes overflow somewhere split *next bucket*

use +1 bit of hash value



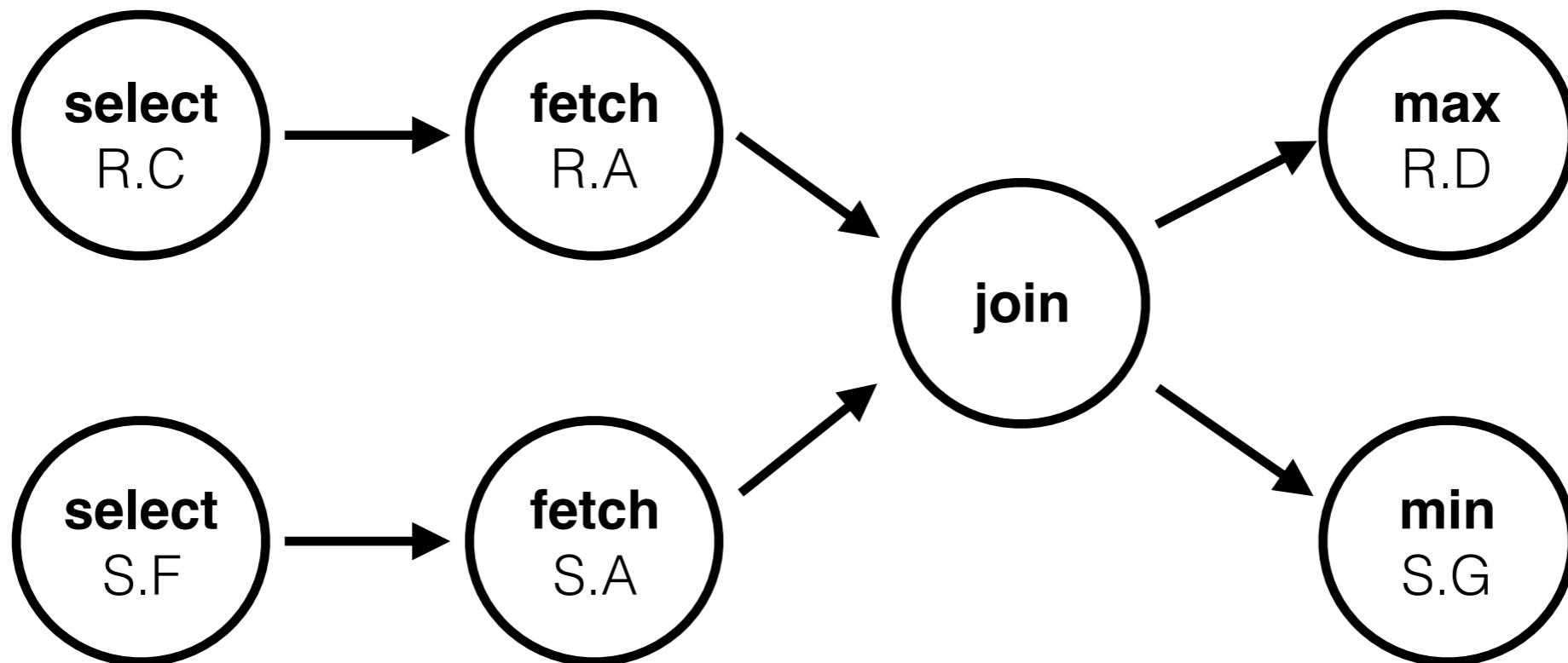
search: use current hash function (bits)
 if after split bucket ok
 else use next hash function (bits+1)



keep splitting next bucket with each overflow
until all buckets are split (for $x=2$)
then restart for $x=3$

any problems
 what would happen
 when we split bucket C



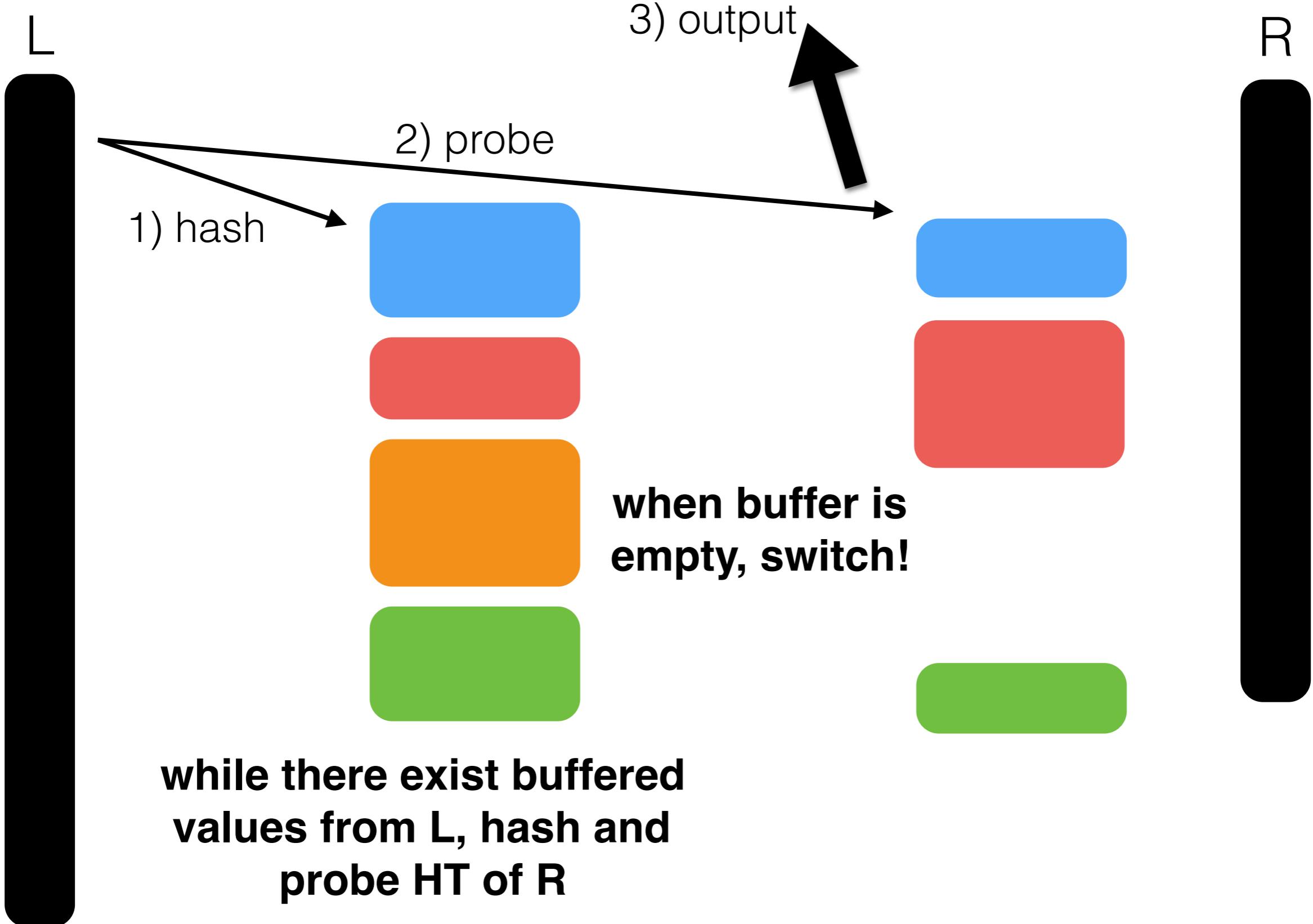


What can we do to start working immediately?
(hint: vectorized processing)

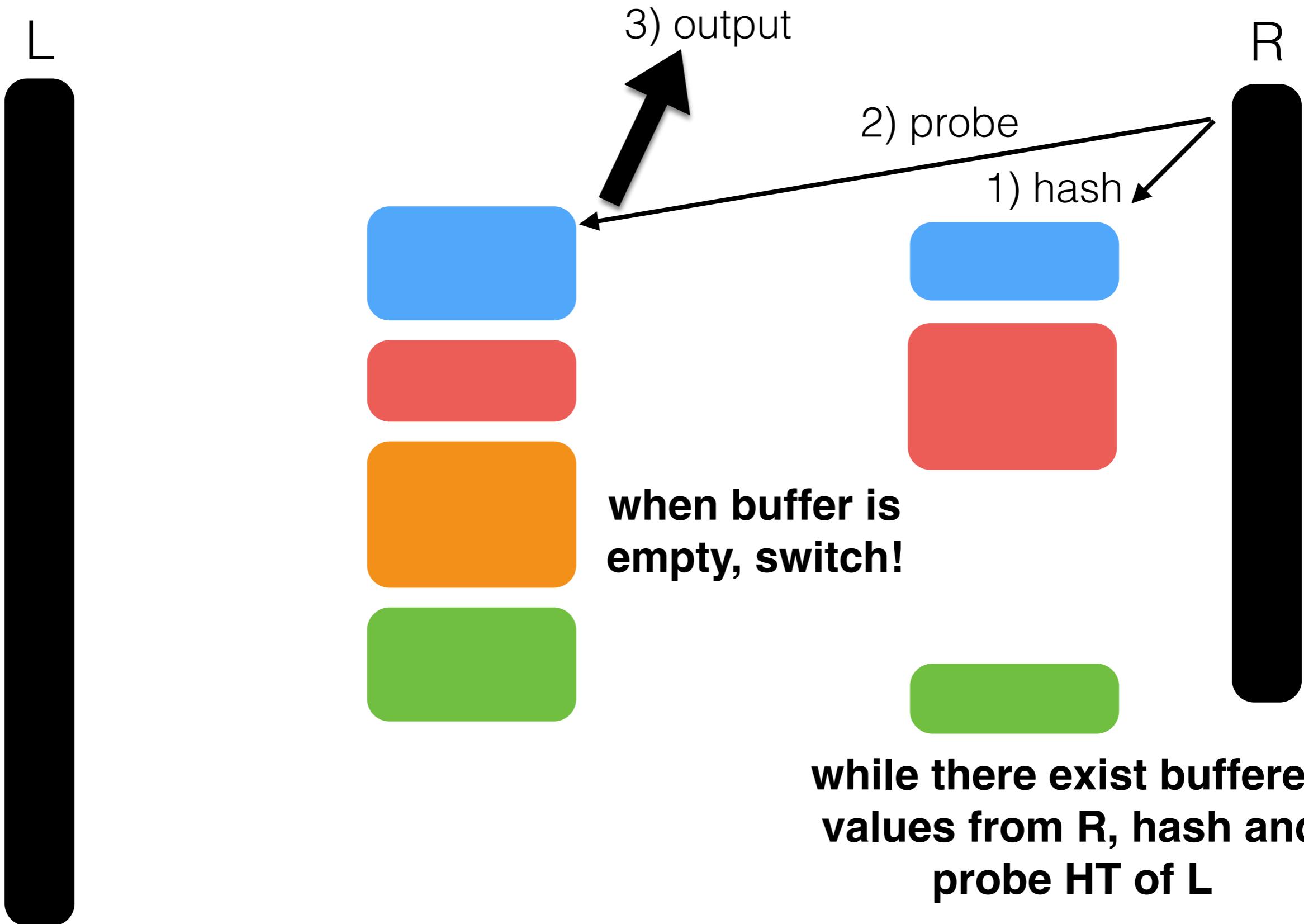
What can we do if we wait for all data to arrive?
(hint: bulk processing)



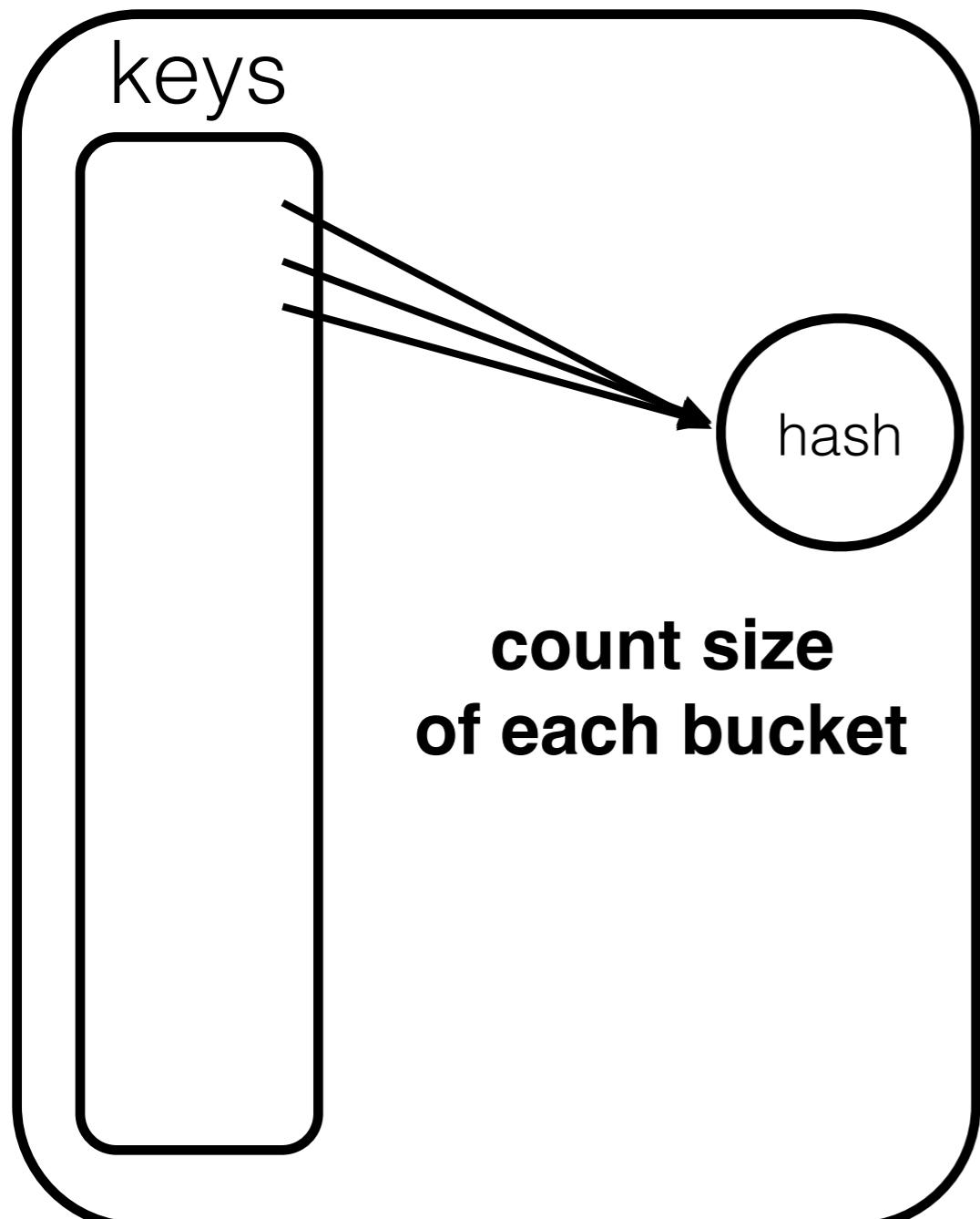
symmetric hash join



symmetric hash join

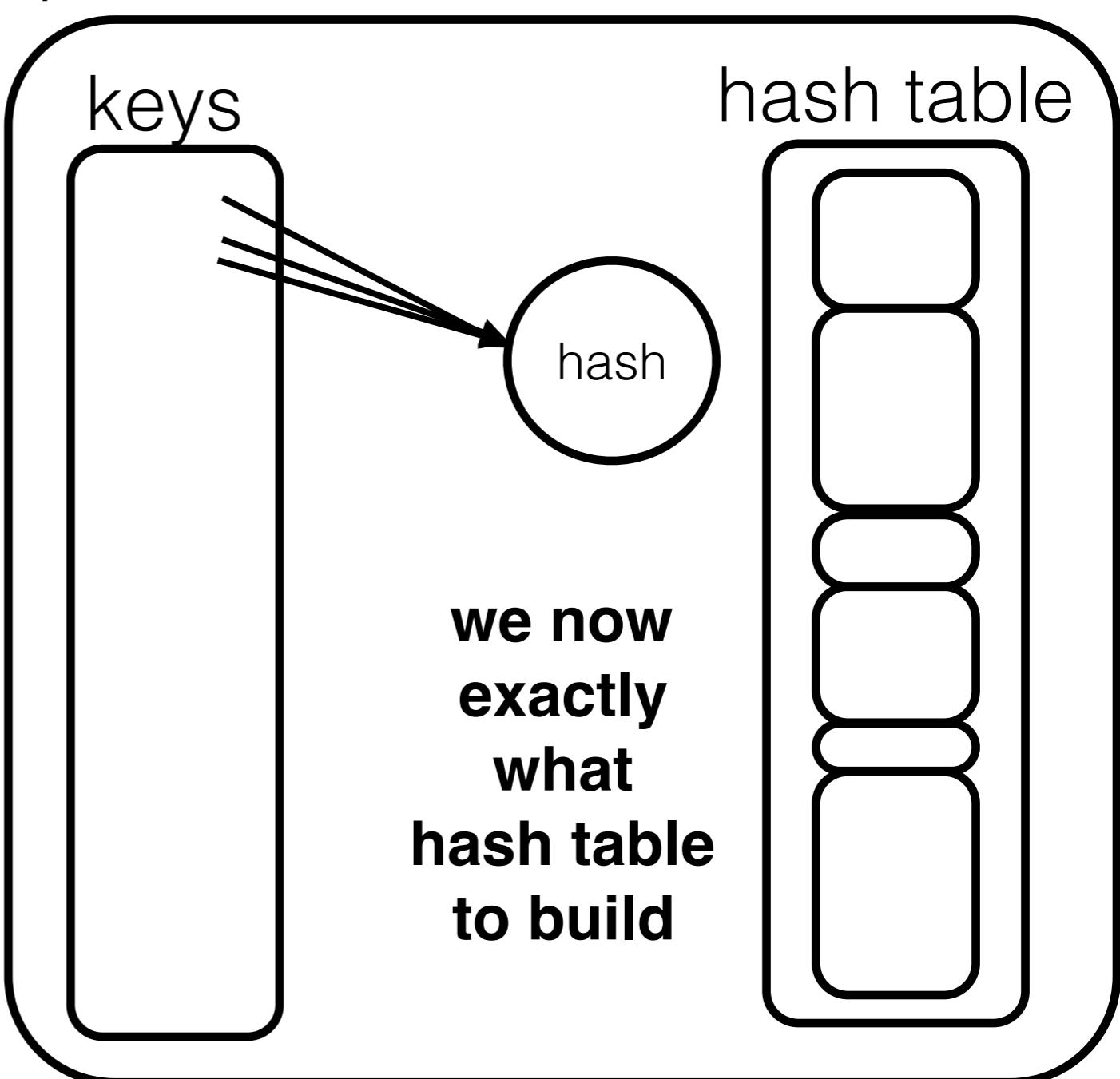


phase 1

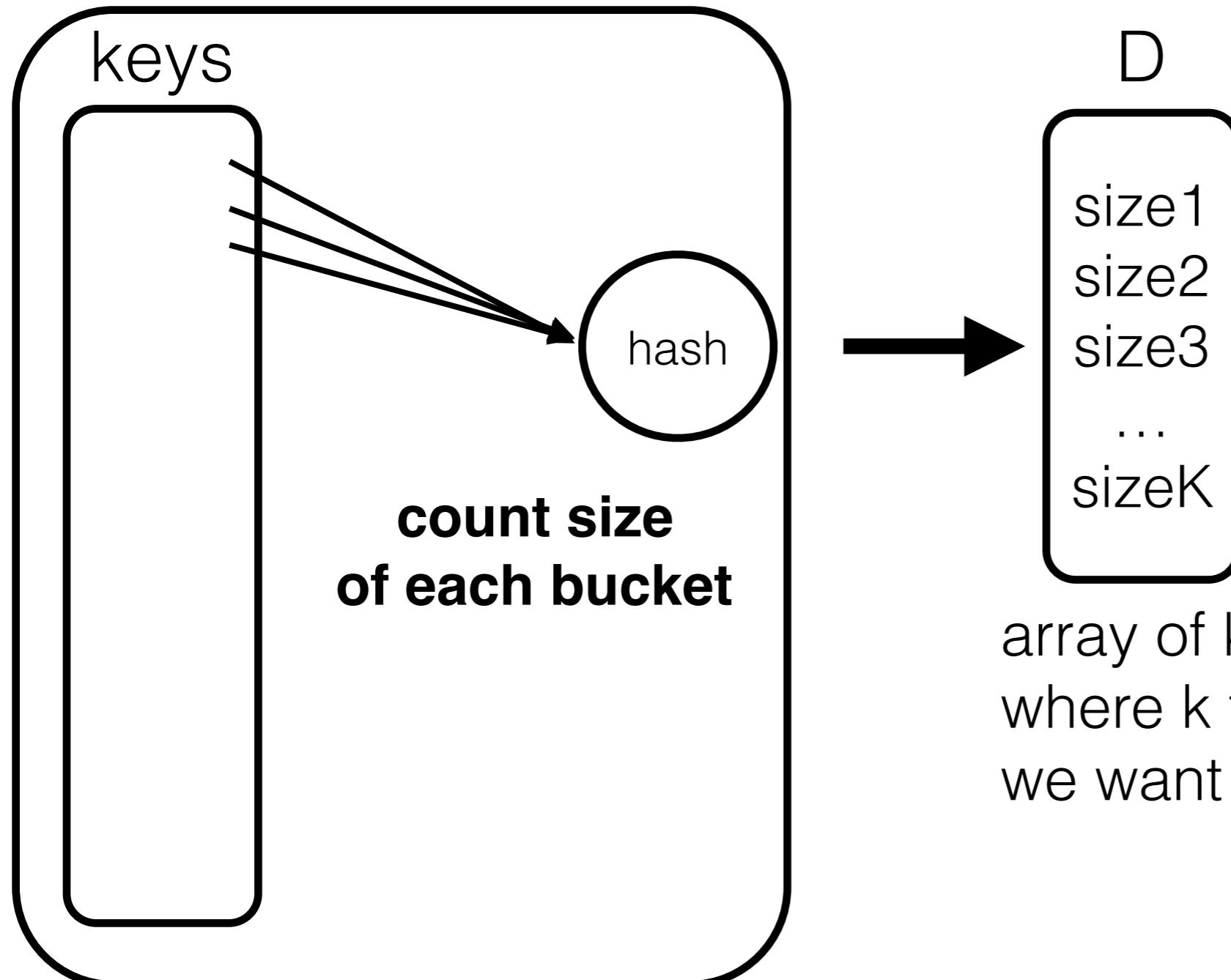


static hashing with 2 passes

phase 2



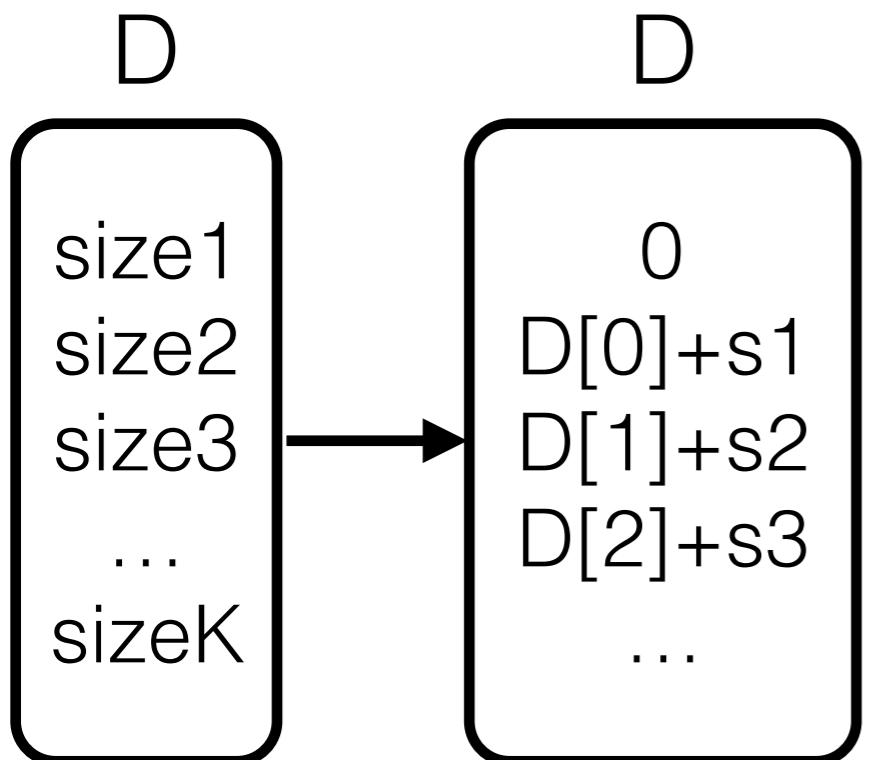
phase 1



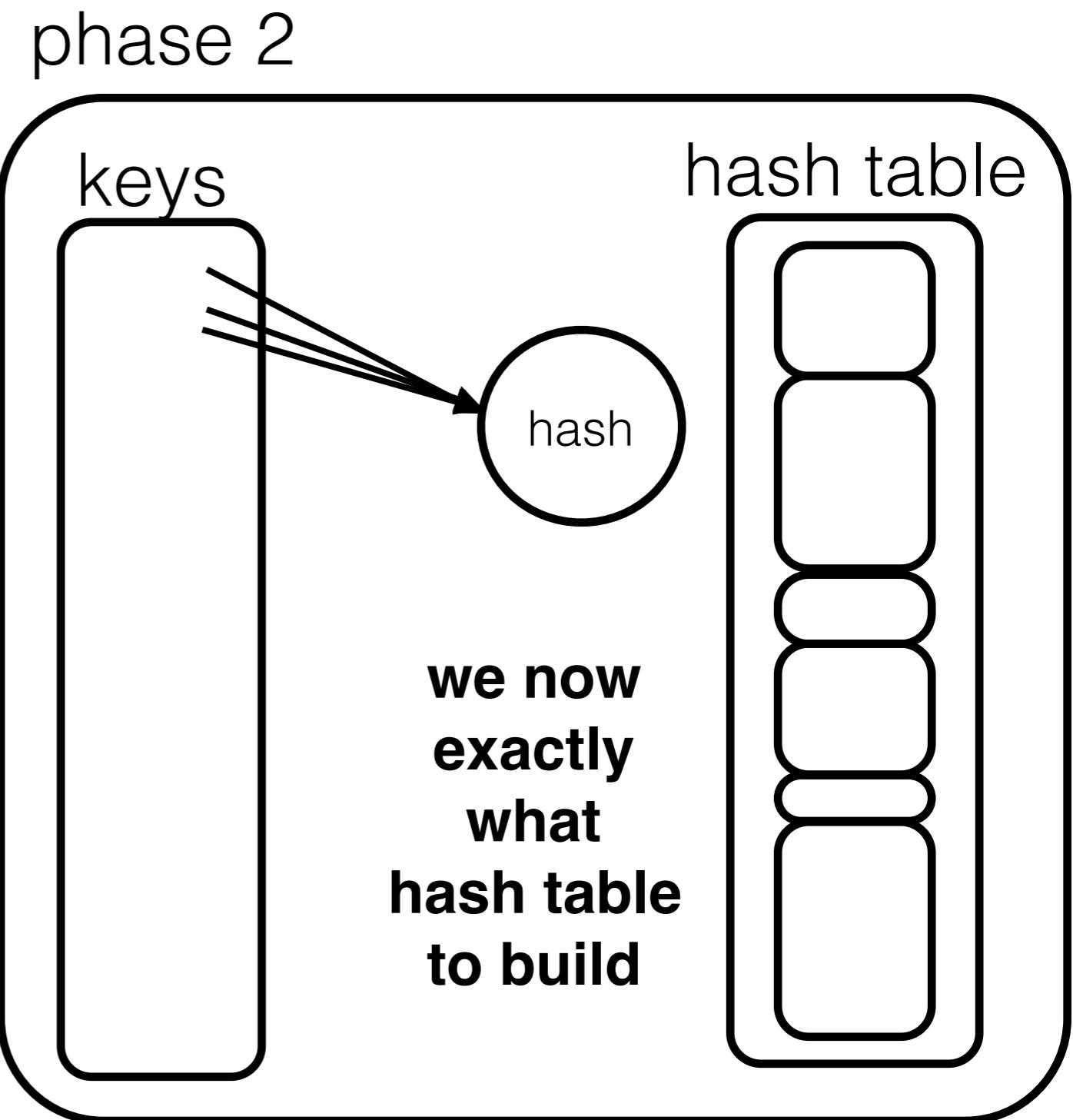
for each bucket
we know its size

array of k slots
where k the buckets
we want to have





pass D and sum all counts
to get offsets in a sequentially
stored hash table



which hash function should we use?

$$h = a * \text{key} + b$$

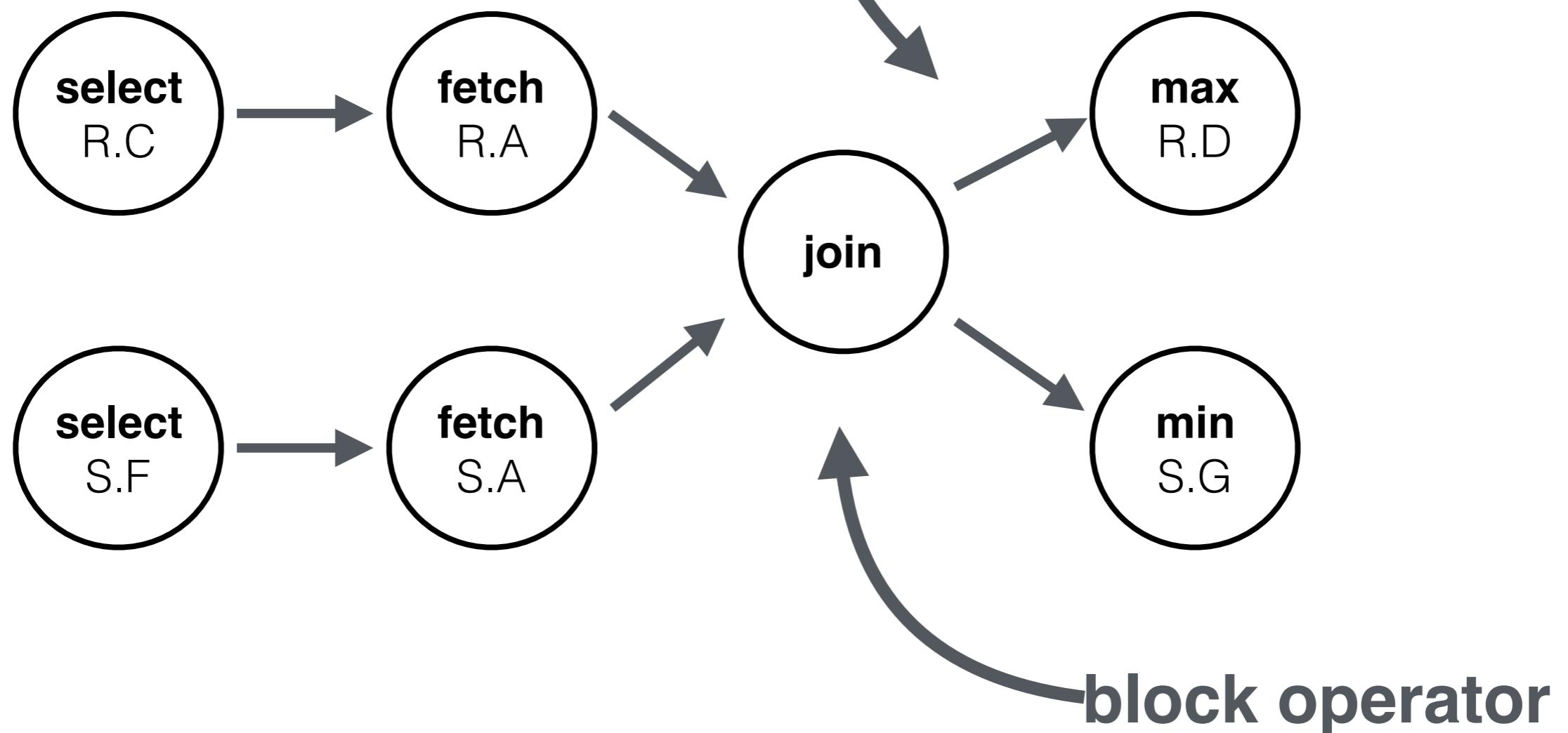
$$b = h \bmod K$$

**can be done with
binary ops only
use X LSBs**



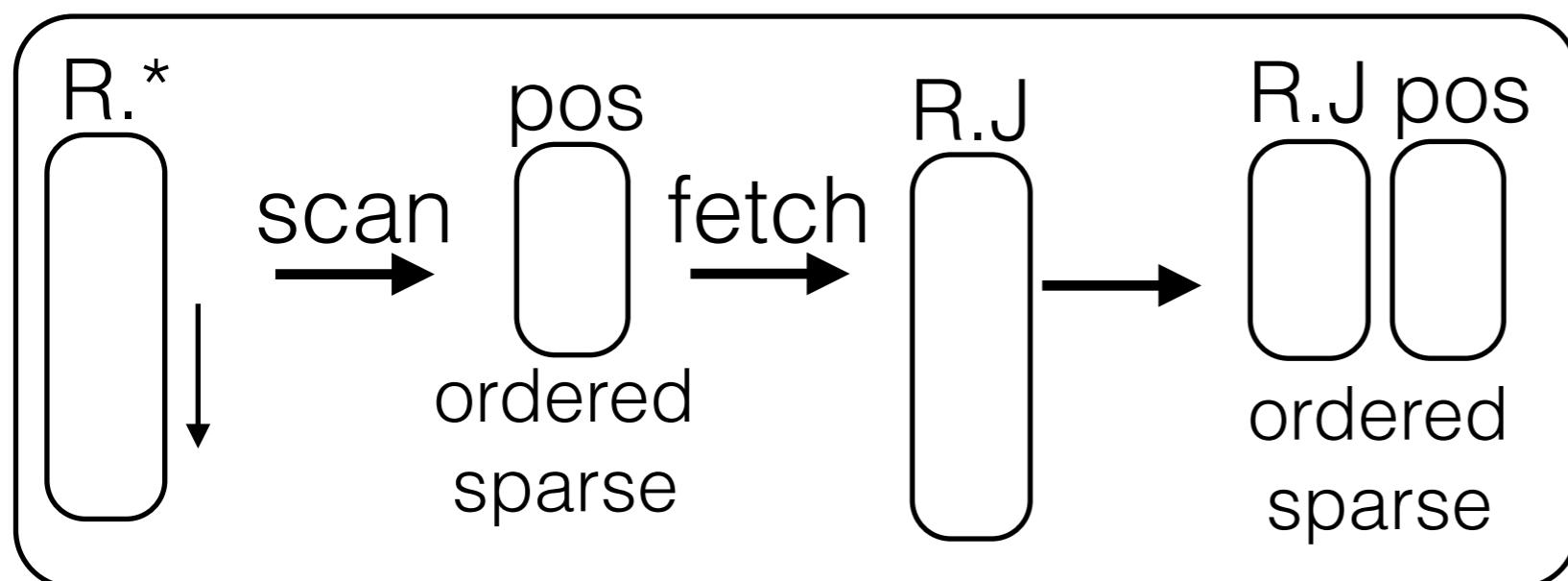
```
select max(R.D),min(S.G)  
from R,S  
where R.A=S.A and R.C<10 and S.F>30
```

what happens after the join?
access patterns



```
select R.A, R.B, R.C, S.A, S.B, S.C  
from R, S  
where R.J=S.J and ...
```

preparing the R join input



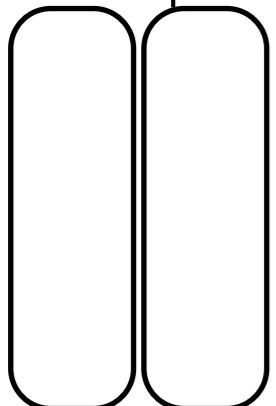
same for the S join input

```
select R.A, R.B, R.C, S.A, S.B, S.C  
from R, S  
where R.J=S.J and ...
```

join

join input R.J

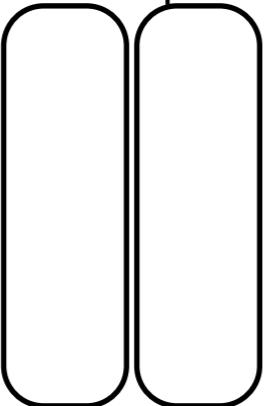
R.J + posR



ordered
sparse

join input S.J

S.J + posS



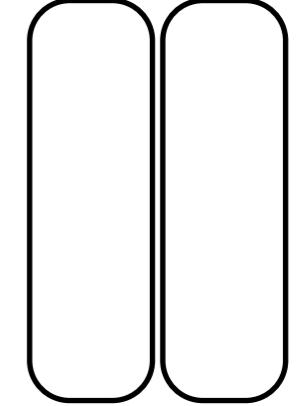
ordered
sparse

partition (=reorder)
both join inputs to join



join result

posR + posS

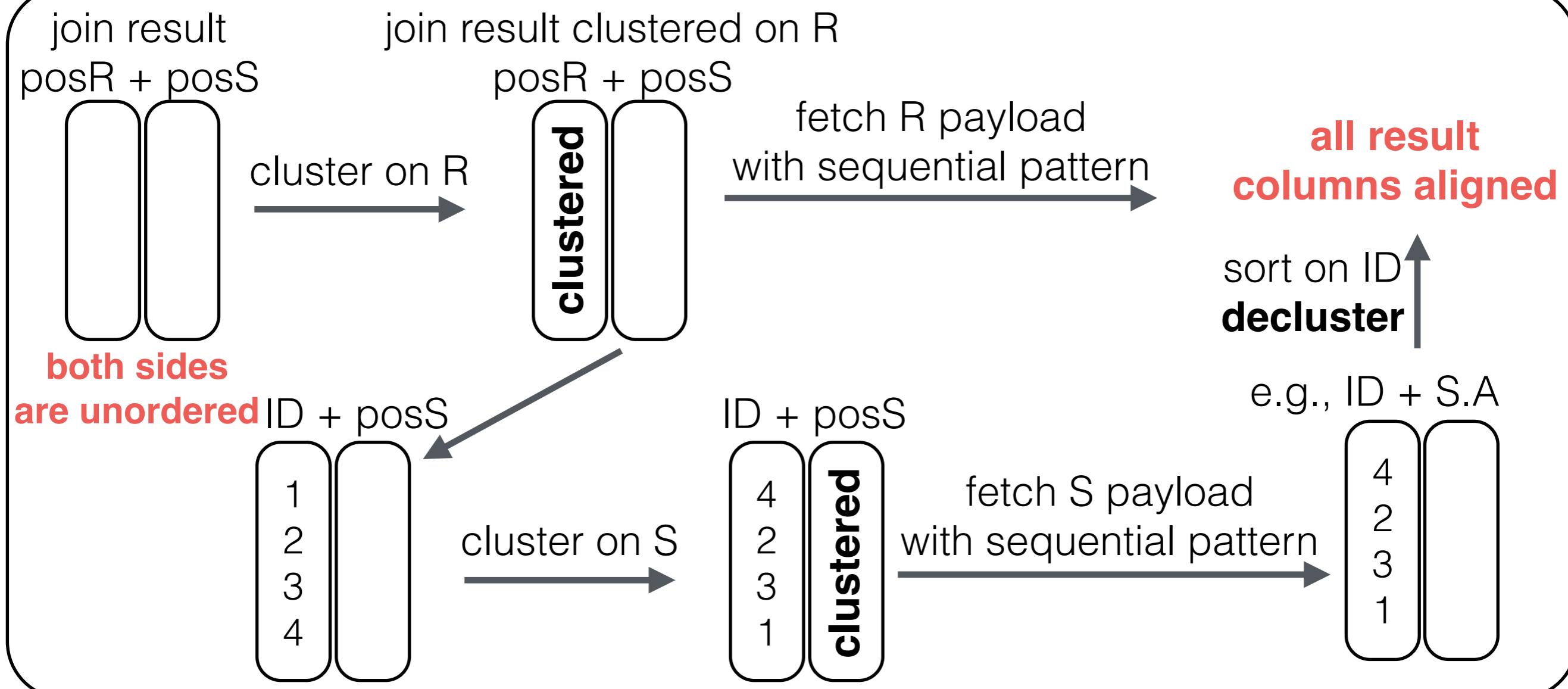


**both sides
are unordered**

```
select R.A, R.B, R.C, S.A, S.B, S.C
from R, S
where R.J=S.J and ...
```

Cache-Conscious Radix Decluster Projections
 By S. Manegold, P. Boncz, N. Nes, and M. Kersten
 Very Large Databases Conference, 2004

radix declustering





Textbook Chapter 11

Cache-Conscious Radix Decluster Projections
S. Manegold, P. Boncz, N. Nes, and M. Kersten
Very Large Databases Conference, 2004



HARVARD
School of Engineering
and Applied Sciences

CS165, Fall 2015
Stratos Idreos

hashing

DATA SYSTEMS

prof. Stratos Idreos

