



CS 165

Data Systems

Have fun learning to design and build modern data systems

class 15

joins

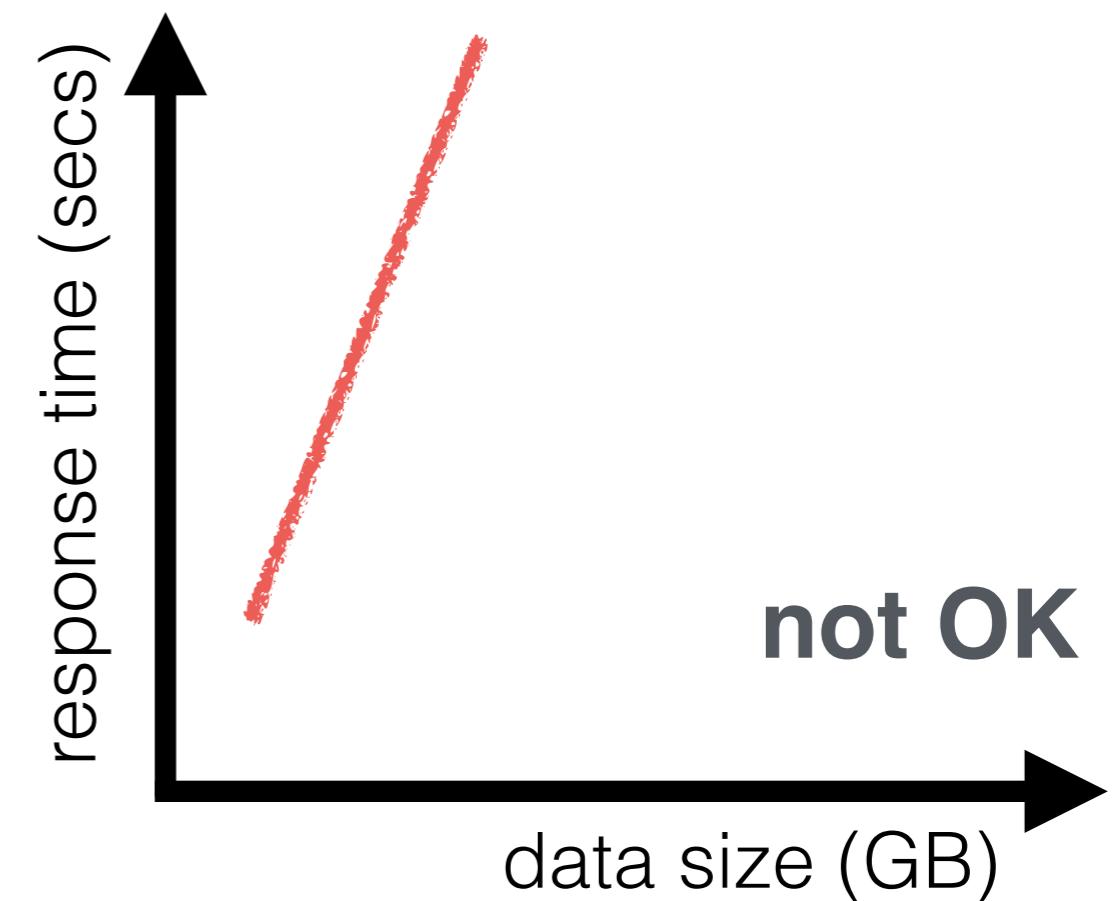
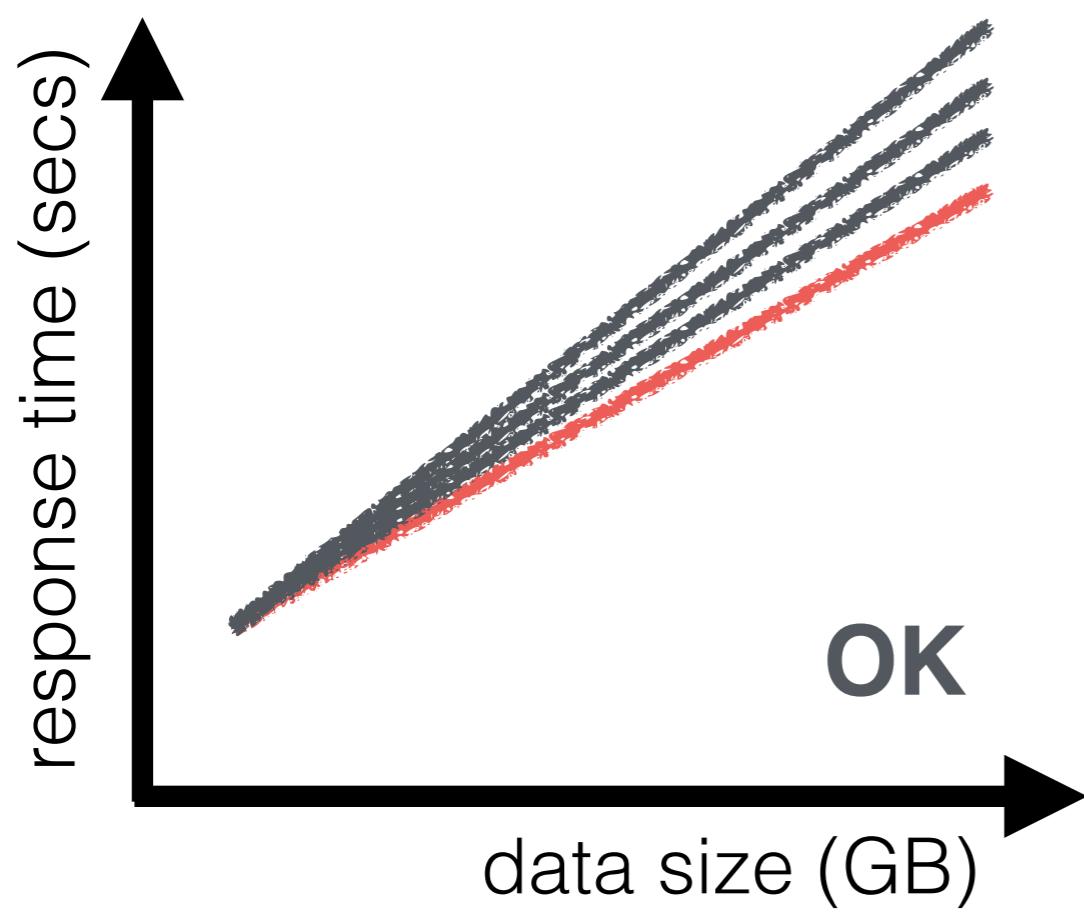
prof. Stratos Idreos

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



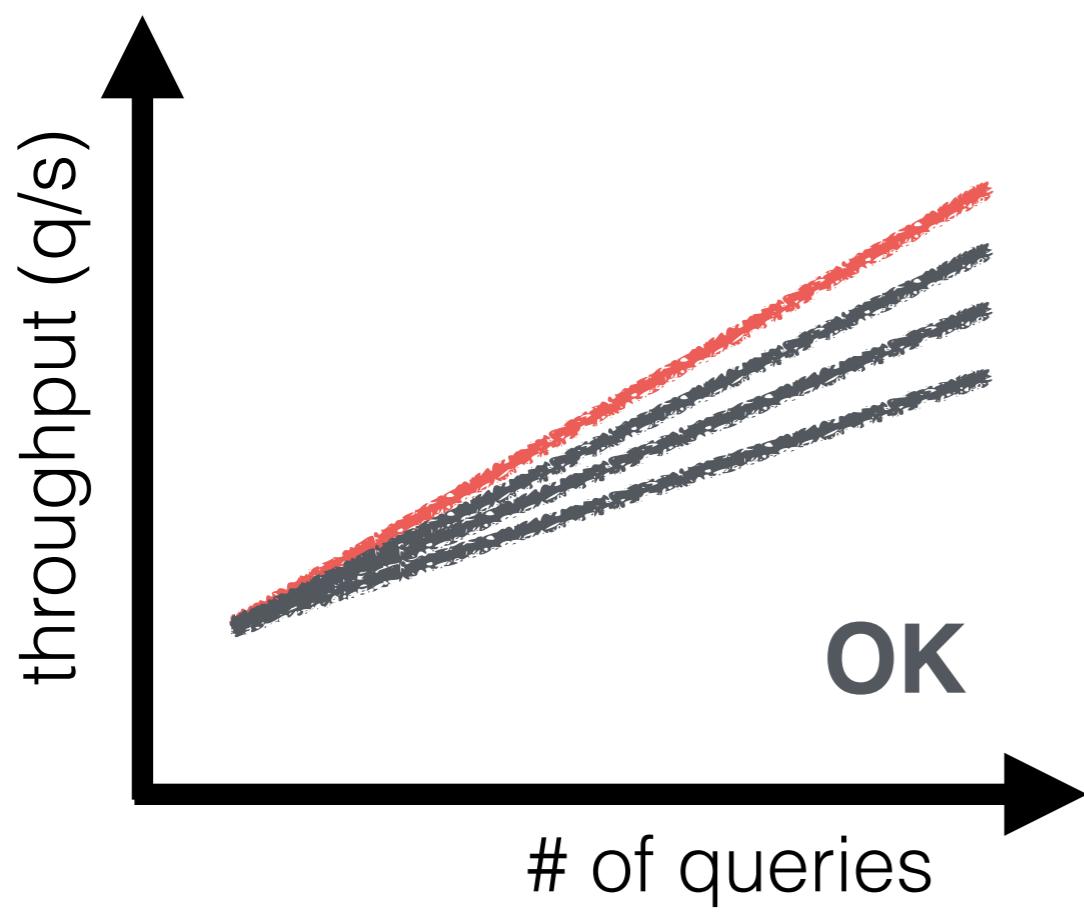


analysis/performance reports



e.g., to test the select operator

analysis/performance reports



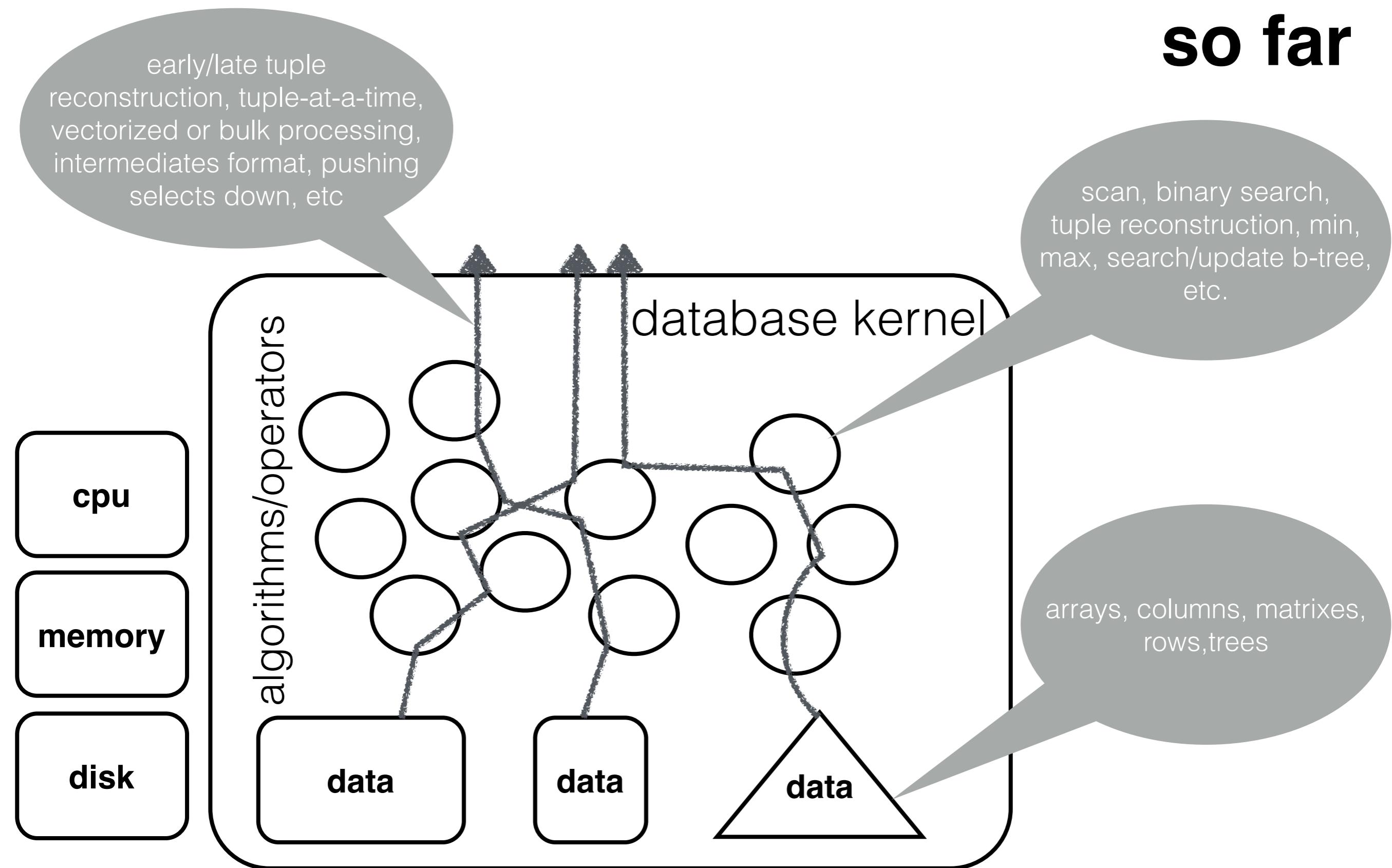
e.g., to test shared scans



hardware, data and query based optimizations
(project=m3)

apply to all
algo/data structures

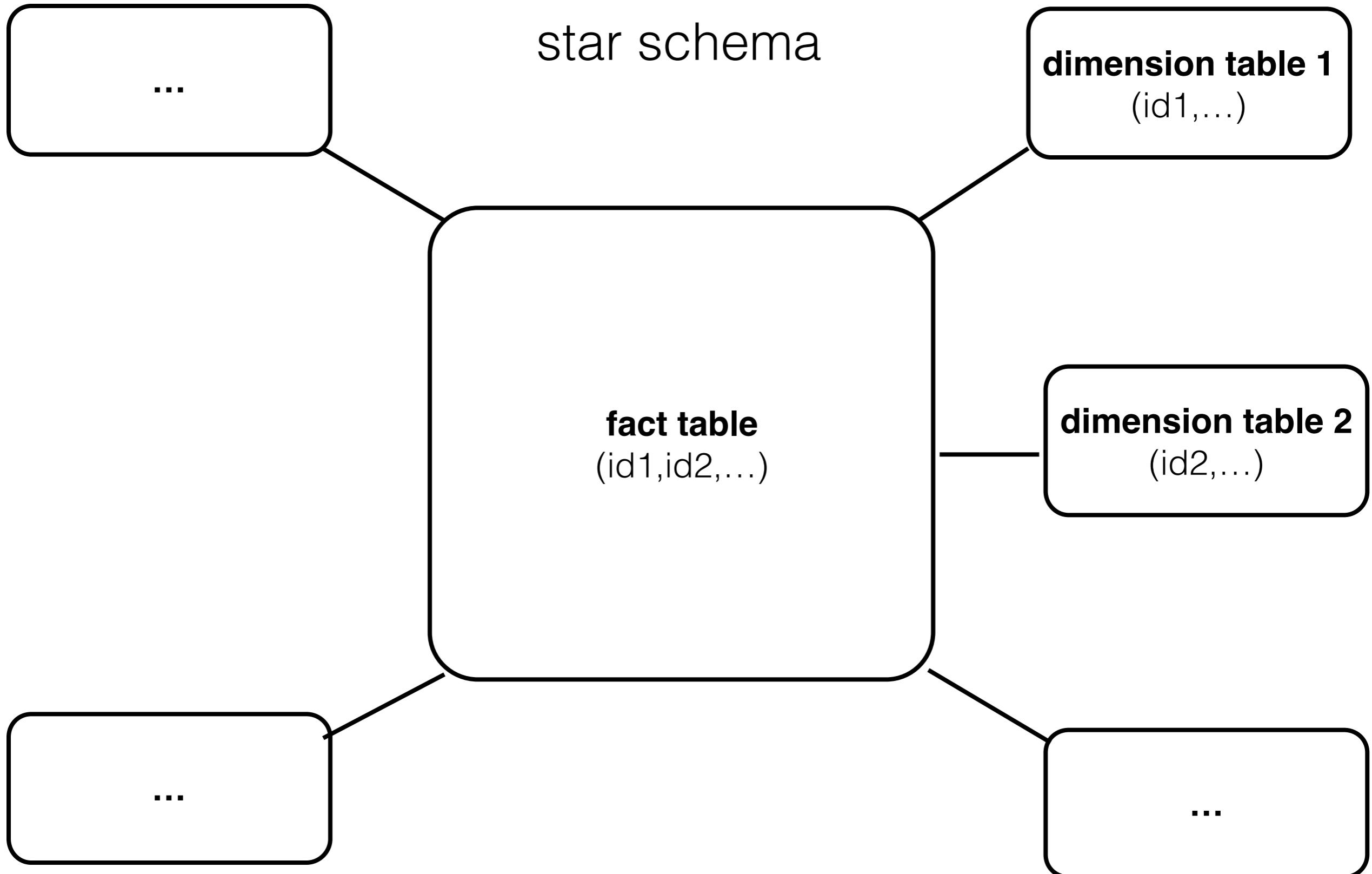
so far



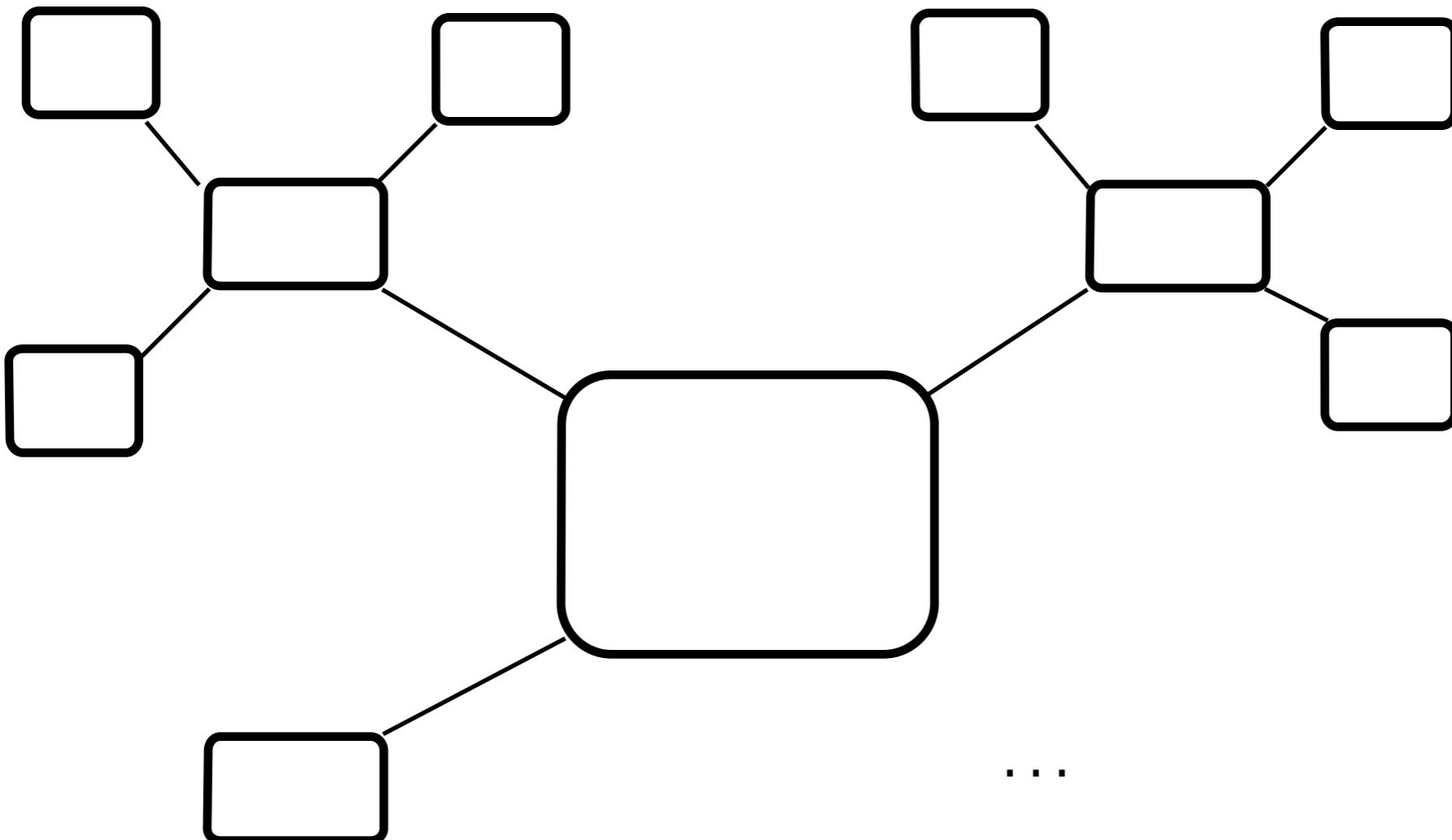


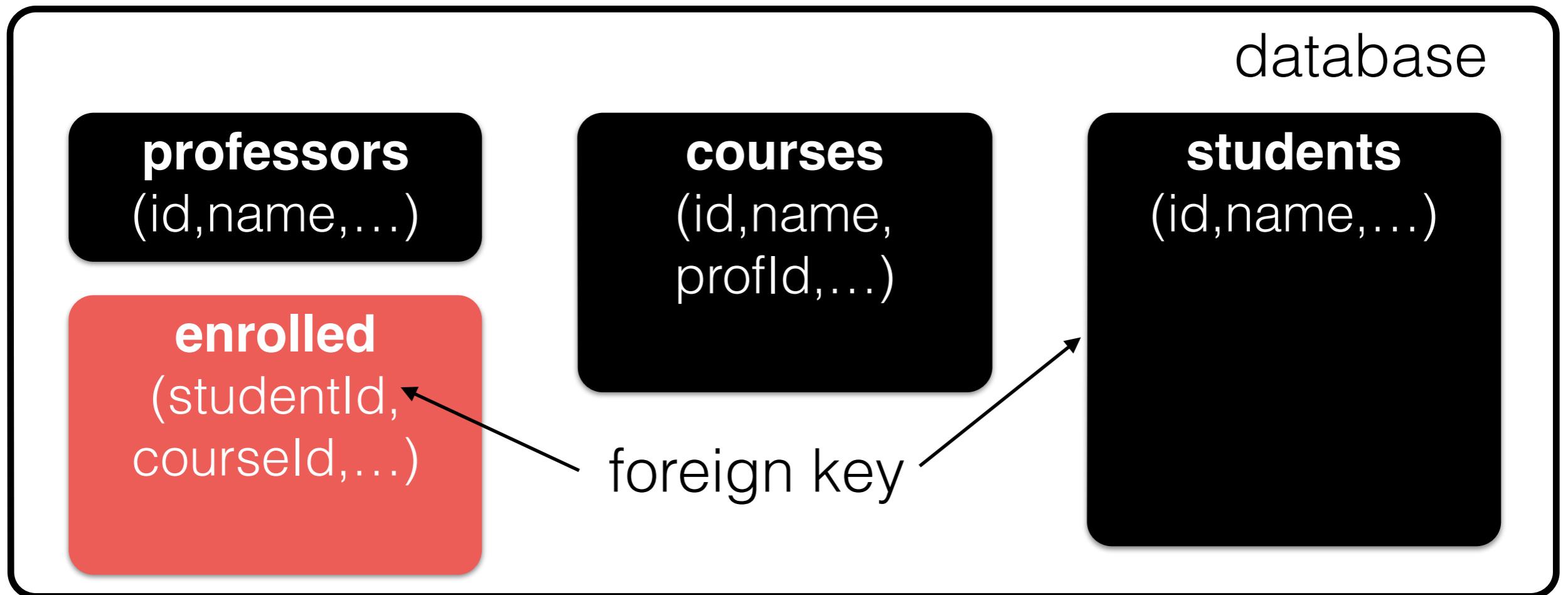
joins
(project=m4)

star schema



snowflake schema





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**



table 1

foreign key
referencing table 2

table 2

primary key

join

find all tuples where FK=PK

join: glue the data back together



key,payload

1,d1,e1,f1
1,d2,e2,f2
2,d3,e3,f3
2,d4,e4,f4
2,d5,e5,f5
3,d6,e6,f6

join

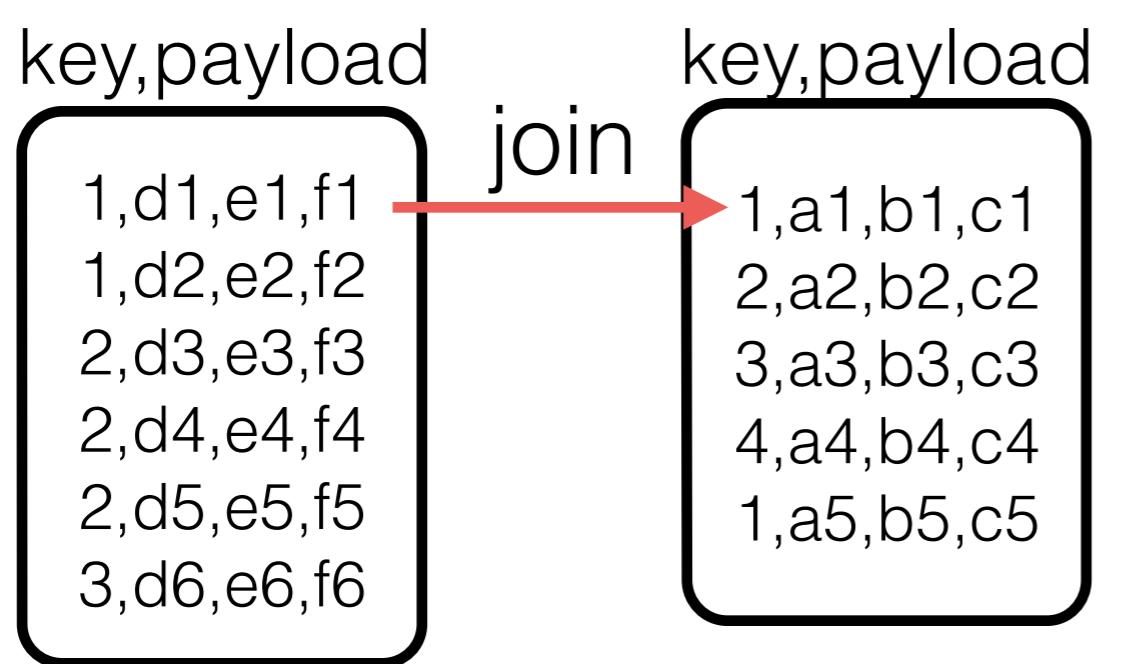
key,payload

1,a1,b1,c1
2,a2,b2,c2
3,a3,b3,c3
4,a4,b4,c4
1,a5,b5,c5

join result

1,d1,e1,f1,a1,b1,c1
1,d1,e1,f1,a5,b5,c5
1,d2,e2,f2,a1,b1,c1
1,d2,e2,f2,a5,b5,c5
2,d3,e3,f3,a2,b2,c2
2,d4,e4,f4,a2,b2,c2
2,d5,e5,f5,a2,b2,c2
3,d6,e6,f6,a3,b3,c3

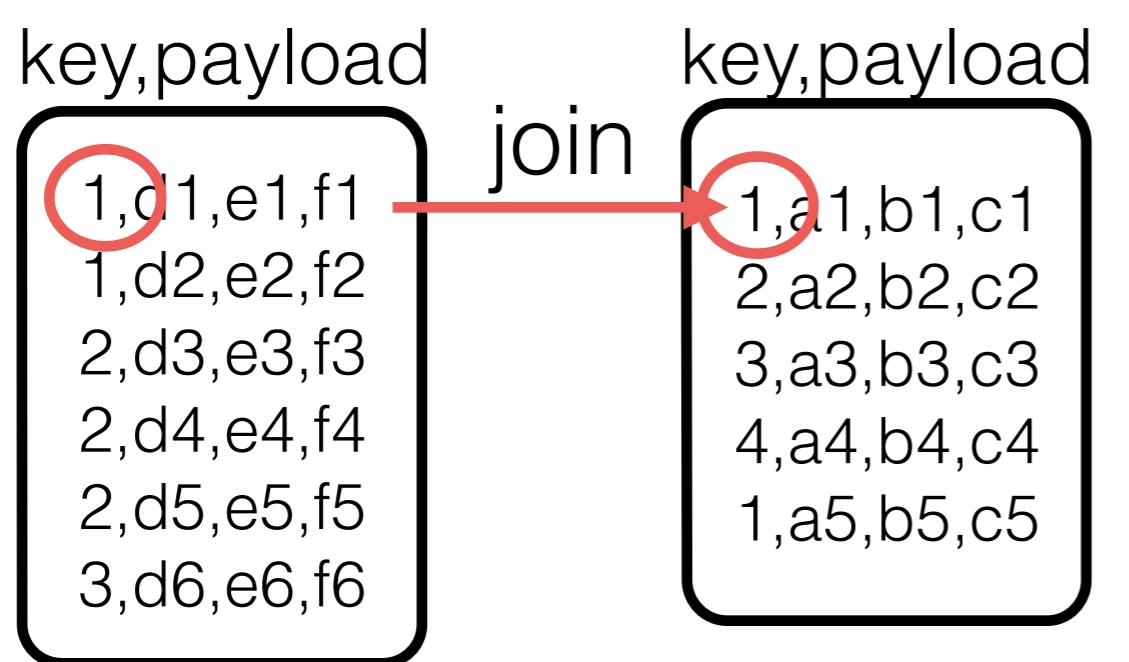




join result

1,d1,e1,f1,a1,b1,c1
1,d1,e1,f1,a5,b5,c5
1,d2,e2,f2,a1,b1,c1
1,d2,e2,f2,a5,b5,c5
2,d3,e3,f3,a2,b2,c2
2,d4,e4,f4,a2,b2,c2
2,d5,e5,f5,a2,b2,c2
3,d6,e6,f6,a3,b3,c3

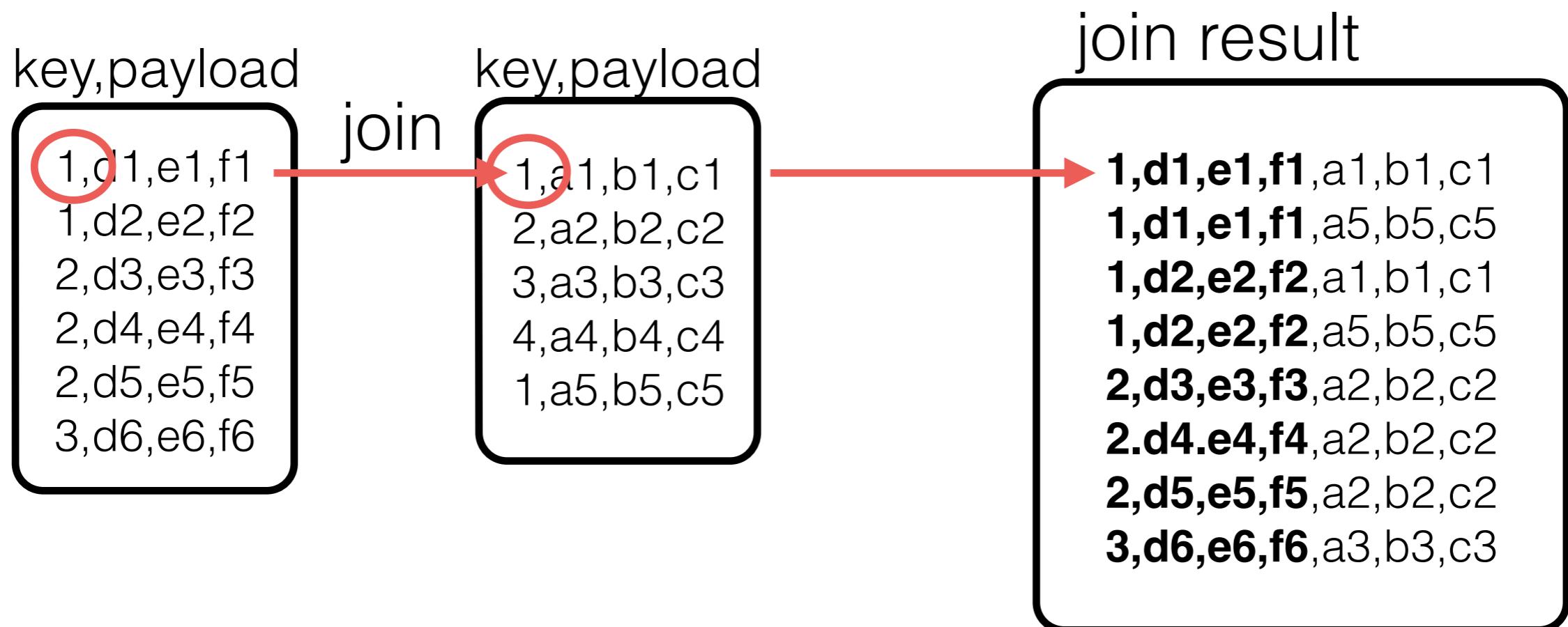




join result

1,d1,e1,f1,a1,b1,c1
1,d1,e1,f1,a5,b5,c5
1,d2,e2,f2,a1,b1,c1
1,d2,e2,f2,a5,b5,c5
2,d3,e3,f3,a2,b2,c2
2,d4,e4,f4,a2,b2,c2
2,d5,e5,f5,a2,b2,c2
3,d6,e6,f6,a3,b3,c3





key,payload

1,d1,e1,f1
1,d2,e2,f2
2,d3,e3,f3
2,d4,e4,f4
2,d5,e5,f5
3,d6,e6,f6

join

key,payload

1,a1,b1,c1
2,a2,b2,c2
3,a3,b3,c3
4,a4,b4,c4
1,a5,b5,c5

join result

1,d1,e1,f1,a1,b1,c1
1,d1,e1,f1,a5,b5,c5
1,d2,e2,f2,a1,b1,c1
1,d2,e2,f2,a5,b5,c5
2,d3,e3,f3,a2,b2,c2
2,d4,e4,f4,a2,b2,c2
2,d5,e5,f5,a2,b2,c2
3,d6,e6,f6,a3,b3,c3



key,payload

1,d1,e1,f1
1,d2,e2,f2
2,d3,e3,f3
2,d4,e4,f4
2,d5,e5,f5
3,d6,e6,f6

join

key,payload

1,a1,b1,c1
2,a2,b2,c2
3,a3,b3,c3
4,a4,b4,c4
1,a5,b5,c5

join result

1,d1,e1,f1,a1,b1,c1
1,d1,e1,f1,a5,b5,c5
1,d2,e2,f2,a1,b1,c1
1,d2,e2,f2,a5,b5,c5
2,d3,e3,f3,a2,b2,c2
2,d4,e4,f4,a2,b2,c2
2,d5,e5,f5,a2,b2,c2
3,d6,e6,f6,a3,b3,c3



key,payload

1,d1,e1,f1
1,d2,e2,f2
2,d3,e3,f3
2,d4,e4,f4
2,d5,e5,f5
3,d6,e6,f6

join

key,payload

1,a1,b1,c1
2,a2,b2,c2
3,a3,b3,c3
4,a4,b4,c4
1,a5,b5,c5

join result

1,d1,e1,f1,a1,b1,c1
1,d1,e1,f1,a5,b5,c5
1,d2,e2,f2,a1,b1,c1
1,d2,e2,f2,a5,b5,c5
2,d3,e3,f3,a2,b2,c2
2,d4,e4,f4,a2,b2,c2
2,d5,e5,f5,a2,b2,c2
3,d6,e6,f6,a3,b3,c3



key,payload

1,d1,e1,f1
1,d2,e2,f2
2,d3,e3,f3
2,d4,e4,f4
2,d5,e5,f5
3,d6,e6,f6

join

key,payload

1,a1,b1,c1
2,a2,b2,c2
3,a3,b3,c3
4,a4,b4,c4
1,a5,b5,c5

join result

1,d1,e1,f1,a1,b1,c1
1,d1,e1,f1,a5,b5,c5
1,d2,e2,f2,a1,b1,c1
1,d2,e2,f2,a5,b5,c5
2,d3,e3,f3,a2,b2,c2
2,d4,e4,f4,a2,b2,c2
2,d5,e5,f5,a2,b2,c2
3,d6,e6,f6,a3,b3,c3



key,payload

1,d1,e1,f1
1,d2,e2,f2
2,d3,e3,f3
2,d4,e4,f4
2,d5,e5,f5
3,d6,e6,f6

join

key,payload

1,a1,b1,c1
2,a2,b2,c2
3,a3,b3,c3
4,a4,b4,c4
1,a5,b5,c5

join result

1,d1,e1,f1,a1,b1,c1
1,d1,e1,f1,a5,b5,c5
1,d2,e2,f2,a1,b1,c1
1,d2,e2,f2,a5,b5,c5
2,d3,e3,f3,a2,b2,c2
2,d4,e4,f4,a2,b2,c2
2,d5,e5,f5,a2,b2,c2
3,d6,e6,f6,a3,b3,c3



key,payload

1,d1,e1,f1
1,d2,e2,f2
2,d3,e3,f3
2,d4,e4,f4
2,d5,e5,f5
3,d6,e6,f6

join

key,payload

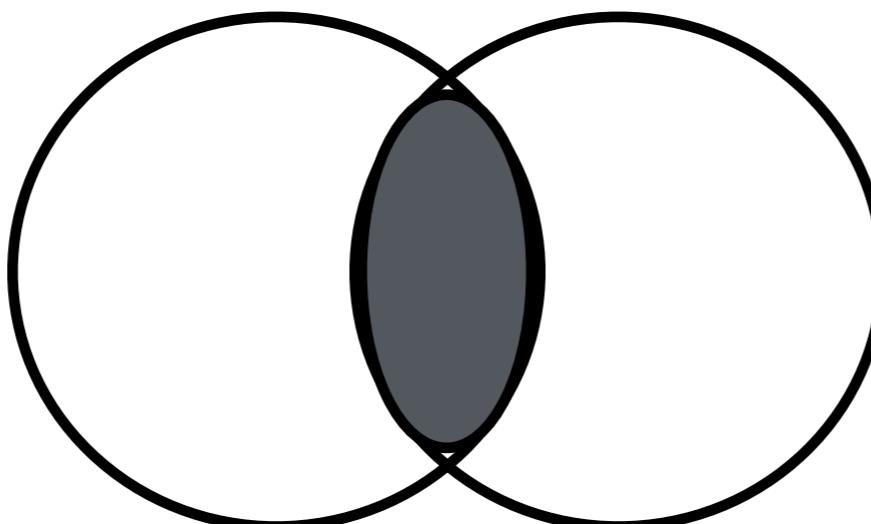
1,a1,b1,c1
2,a2,b2,c2
3,a3,b3,c3
4,a4,b4,c4
1,a5,b5,c5

join result

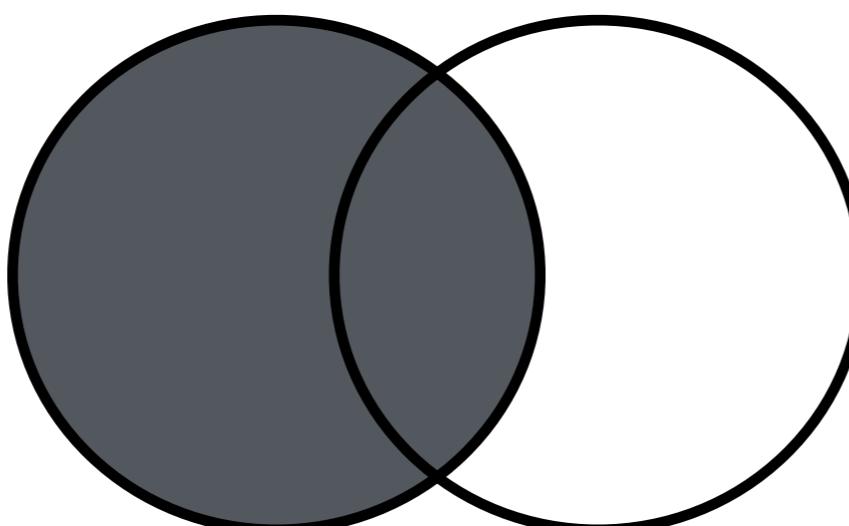
1,d1,e1,f1,a1,b1,c1
1,d1,e1,f1,a5,b5,c5
1,d2,e2,f2,a1,b1,c1
1,d2,e2,f2,a5,b5,c5
2,d3,e3,f3,a2,b2,c2
2,d4,e4,f4,a2,b2,c2
2,d5,e5,f5,a2,b2,c2
3,d6,e6,f6,a3,b3,c3



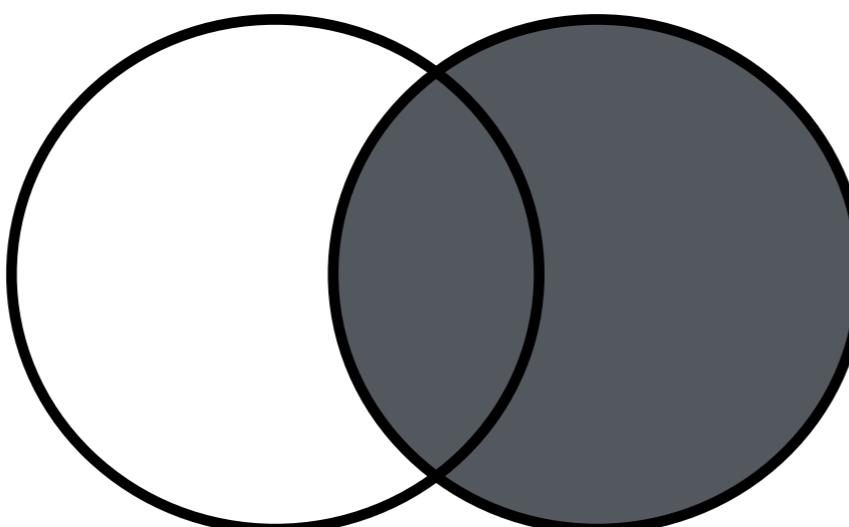
inner join



left outer

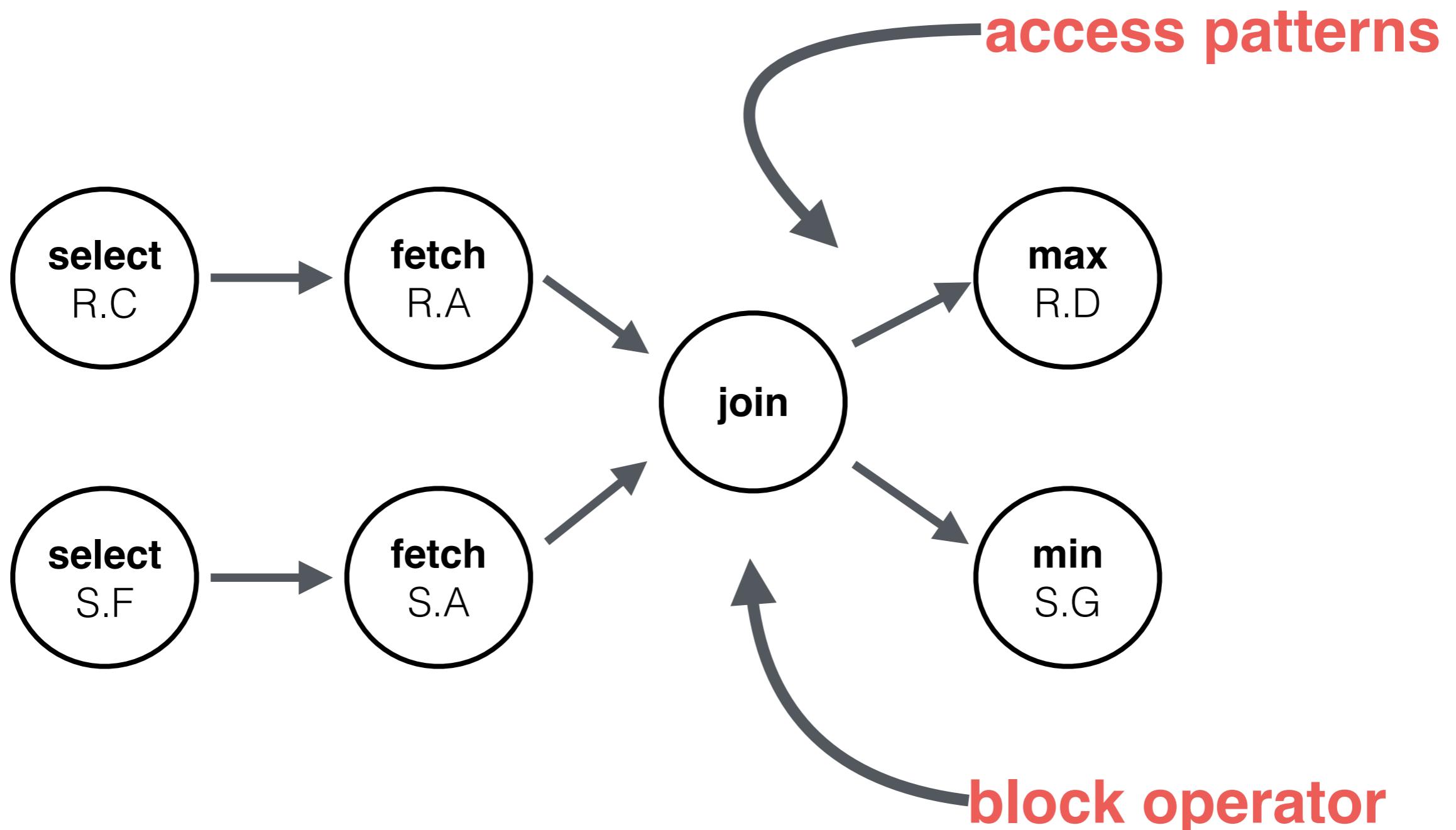


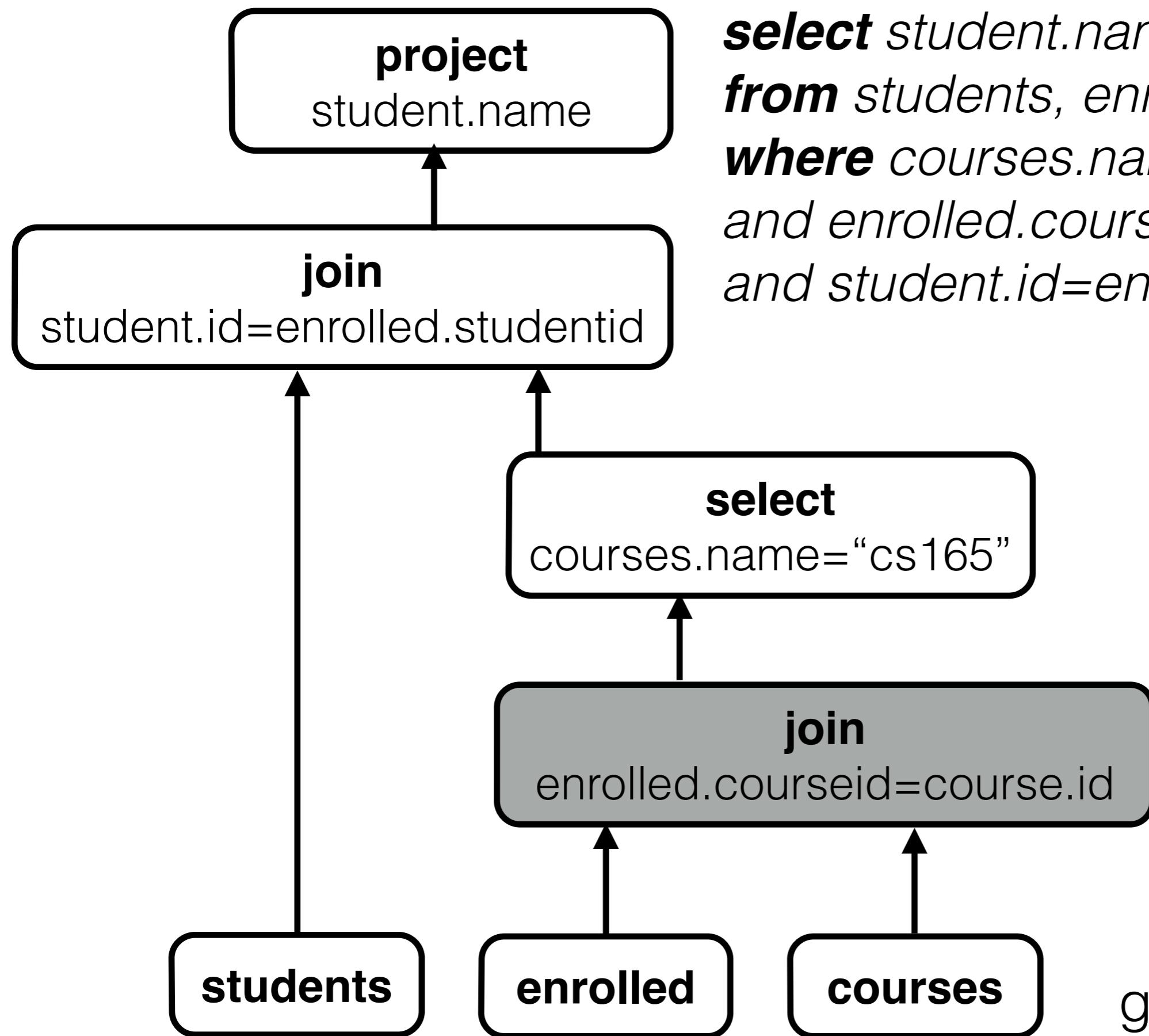
right outer



$R(A,B,C,D) - S(A,E,F,G)$

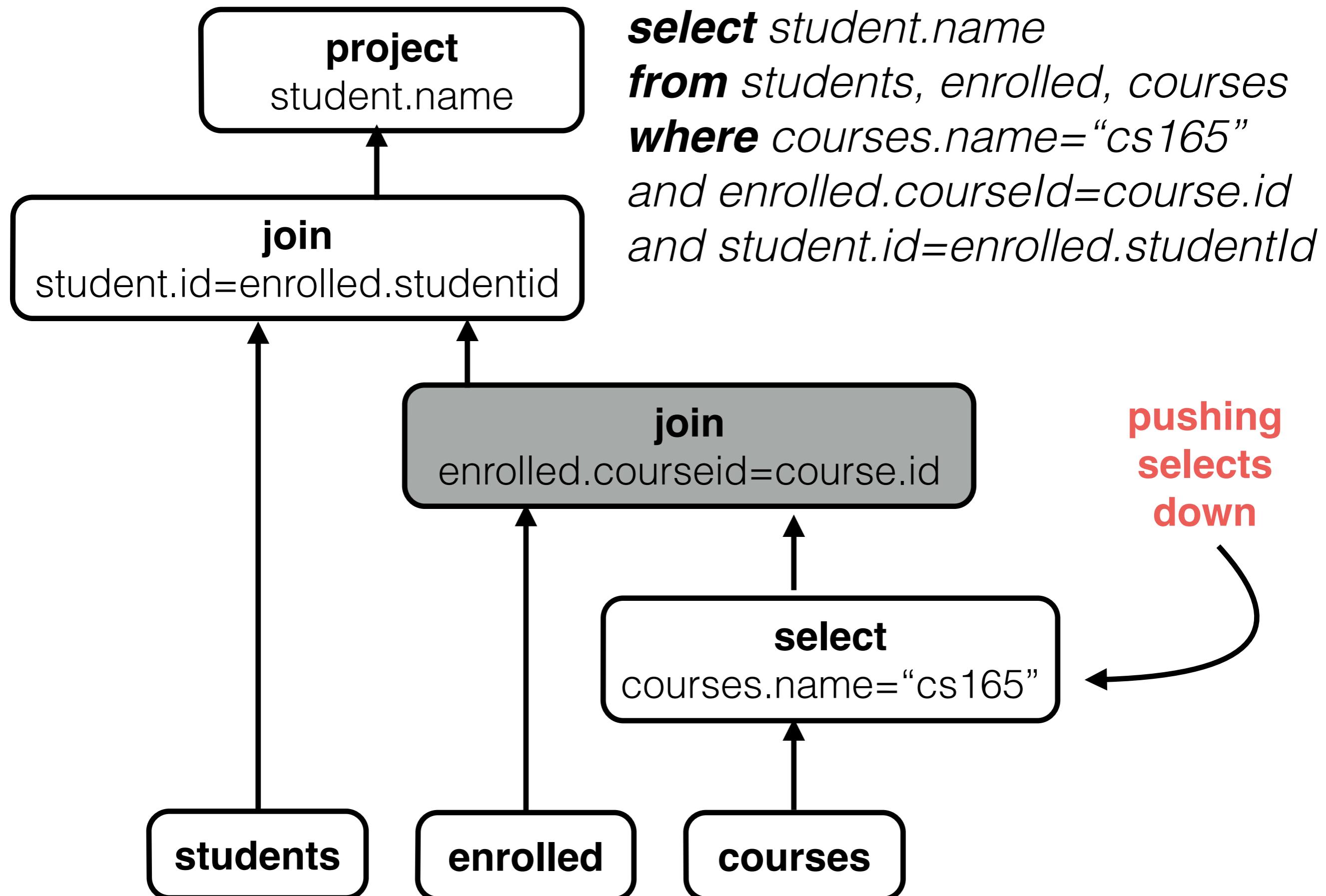
select max(R.D),min(S.G)
from R,S
where R.A=S.A and R.C<10 and S.F>30





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





Initial Status			Query and Query Plan (MAL Algebra)		
Relation R			Relation S		
Ra	Rb	Rc	Sa	Sb	
3	12	12	17	11	
16	34	34	49	35	
56	75	53	58	62	
9	45	23	99	44	
11	49	78	64	29	
27	58	65	37	78	
8	97	33	53	19	
41	75	21	61	81	
19	42	29	32	26	
35	55	0	50	23	

select sum(R.a) from R, S where R.c = S.b and
5 < R.a < 20 and 40 < R.b < 50 and 30 < S.a < 40

1. inter1 = **select**(Ra,5,20)
2. inter2 = **reconstruct**(Rb.inter1)
3. inter3 = **select**(inter2,40,50)
4. join_input_R = **reconstruct**(Rc,inter3)
5. inter4 = **select**(Sa,50,65)
6. inter5 = **reconstruct**(Sb,inter4)
7. join_input_S = **reverse**(inter5)
8. join_res_R_S = **join**(join_input_R,join_input_S)
9. inter6 = **voidTail**(join_res_R_S)
10. inter7 = **reconstruct**(Ra,inter6)
11. result = **sum**(inter7)

Initial Status

Relation R			Relation S	
Ra	Rb	Rc	Sa	Sb
3	12	12	17	11
16	34	34	49	35
56	75	53	58	62
9	45	23	99	44
11	49	78	64	29
27	58	65	37	78
8	97	33	53	19
41	75	21	61	81
19	42	29	32	26
35	55	0	50	23

Query and Query Plan (MAL Algebra)

select sum(R.a) from R, S where R.c = S.b and
5 < R.a < 20 and 40 < R.b < 50 and 30 < S.a < 40

1. inter1 = **select(Ra,5,20)**
2. inter2 = **reconstruct(Rb.inter1)**
3. inter3 = **select(inter2,40,50)**
4. join_input_R = **reconstruct(Rc,inter3)**
5. inter4 = **select(Sa,50,65)**
6. inter5 = **reconstruct(Sb,inter4)**
7. join_input_S = **reverse(inter5)**
8. join_res_R_S = **join(join_input_R,join_input_S)**
9. inter6 = **voidTail(join_res_R_S)**
10. inter7 = **reconstruct(Ra,inter6)**
11. result = **sum(inter7)**

select(Ra,5,20)

Ra
3
16
56
9
11
27
8
41
19
35

reconstruct(Rb,inter1)

inter1
2
4
5
7
9

Rb
12
34
75
45
49
58
97
75
42
55

select(inter2,40,50)

inter2
2
4
5
7
9

inter3
4
5
9

reconstruct(Rc,inter3)

Rc
12
34
53
23
78
65
33
21
29
0

join_input_R
4
5
9
23
78
65
33
21
29

(1)

(2)

(3)

(4)



Initial Status			Query and Query Plan (MAL Algebra)		
Relation R			Relation S		
Ra	Rb	Rc	Sa	Sb	
3	12	12	17	11	
16	34	34	49	35	
56	75	53	58	62	
9	45	23	99	44	
11	49	78	64	29	
27	58	65	37	78	
8	97	33	53	19	
41	75	21	61	81	
19	42	29	32	26	
35	55	0	50	23	

select sum(R.a) from R, S where R.c = S.b and
 5 < R.a < 20 and 40 < R.b < 50 and 30 < S.a < 40

1. inter1 = **select**(Ra,5,20)
2. inter2 = **reconstruct**(Rb.inter1)
3. inter3 = **select**(inter2,40,50)
4. join_input_R = **reconstruct**(Rc,inter3)
5. inter4 = **select**(Sa,50,65)
6. inter5 = **reconstruct**(Sb,inter4)
7. join_input_S = **reverse**(inter5)
8. join_res_R_S = **join**(join_input_R,join_input_S)
9. inter6 = **voidTail**(join_res_R_S)
10. inter7 = **reconstruct**(Ra,inter6)
11. result = **sum**(inter7)

Initial Status

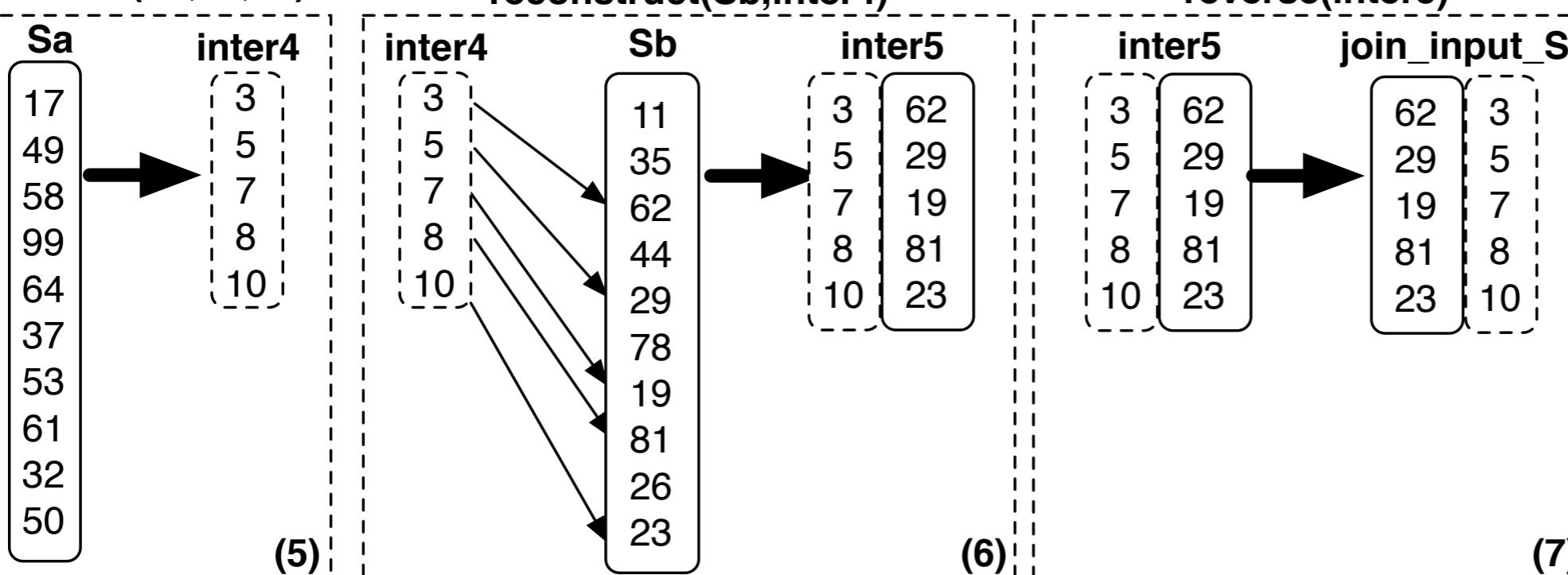
Relation R			Relation S	
Ra	Rb	Rc	Sa	Sb
3	12	12	17	11
16	34	34	49	35
56	75	53	58	62
9	45	23	99	44
11	49	78	64	29
27	58	65	37	78
8	97	33	53	19
41	75	21	61	81
19	42	29	32	26
35	55	0	50	23

Query and Query Plan (MAL Algebra)

select sum(R.a) from R, S where R.c = S.b and
5 < R.a < 20 and 40 < R.b < 50 and 30 < S.a < 40

1. inter1 = **select(Ra,5,20)**
2. inter2 = **reconstruct(Rb.inter1)**
3. inter3 = **select(inter2,40,50)**
4. join_input_R = **reconstruct(Rc,inter3)**
5. inter4 = **select(Sa,50,65)**
6. inter5 = **reconstruct(Sb,inter4)**
7. join_input_S = **reverse(inter5)**
8. join_res_R_S = **join(join_input_R,join_input_S)**
9. inter6 = **voidTail(join_res_R_S)**
10. inter7 = **reconstruct(Ra,inter6)**
11. result = **sum(inter7)**

select(Sa,50,65)



Initial Status			Query and Query Plan (MAL Algebra)		
Relation R			Relation S		
Ra	Rb	Rc	Sa	Sb	
3	12	12	17	11	
16	34	34	49	35	
56	75	53	58	62	
9	45	23	99	44	
11	49	78	64	29	
27	58	65	37	78	
8	97	33	53	19	
41	75	21	61	81	
19	42	29	32	26	
35	55	0	50	23	

select sum(R.a) from R, S where R.c = S.b and
 $5 < R.a < 20$ and $40 < R.b < 50$ and $30 < S.a < 40$

1. inter1 = **select**(Ra,5,20)
2. inter2 = **reconstruct**(Rb.inter1)
3. inter3 = **select**(inter2,40,50)
4. join_input_R = **reconstruct**(Rc,inter3)
5. inter4 = **select**(Sa,50,65)
6. inter5 = **reconstruct**(Sb,inter4)
7. join_input_S = **reverse**(inter5)
8. join_res_R_S = **join**(join_input_R,join_input_S)
9. inter6 = **voidTail**(join_res_R_S)
10. inter7 = **reconstruct**(Ra,inter6)
11. result = **sum**(inter7)

Initial Status

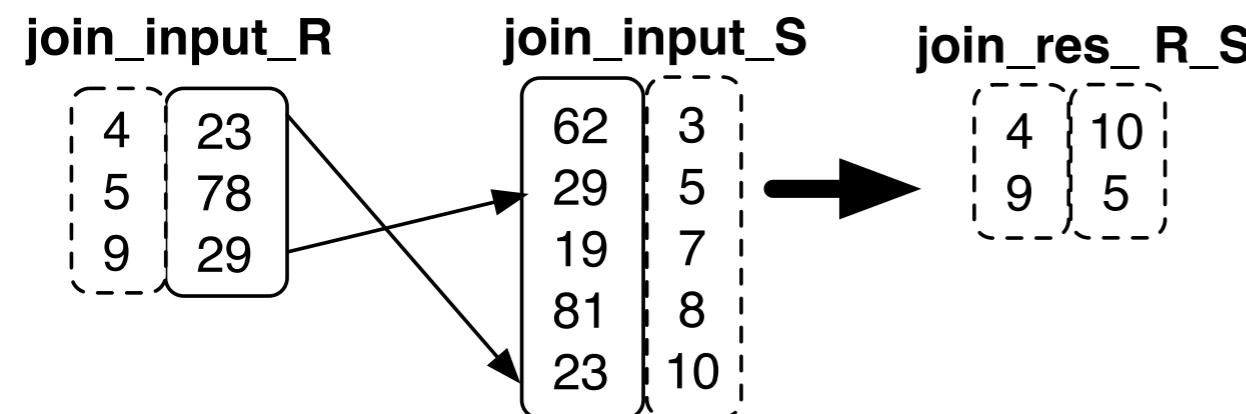
Relation R			Relation S	
Ra	Rb	Rc	Sa	Sb
3	12	12	17	11
16	34	34	49	35
56	75	53	58	62
9	45	23	99	44
11	49	78	64	29
27	58	65	37	78
8	97	33	53	19
41	75	21	61	81
19	42	29	32	26
35	55	0	50	23

Query and Query Plan (MAL Algebra)

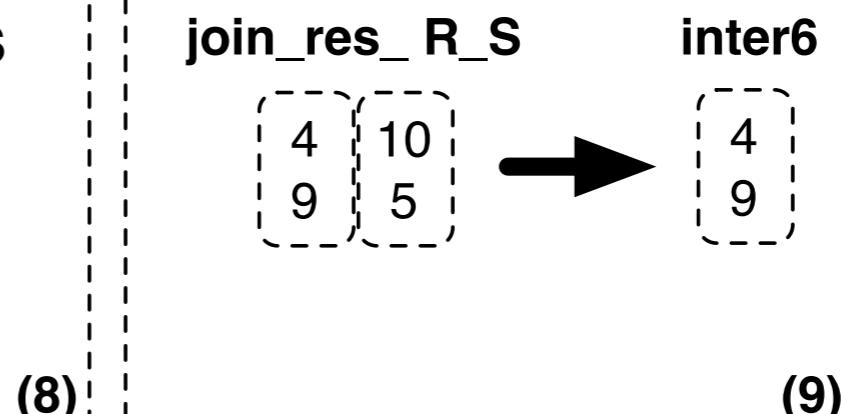
select sum(R.a) from R, S where R.c = S.b and
5 < R.a < 20 and 40 < R.b < 50 and 30 < S.a < 40

1. inter1 = **select(Ra,5,20)**
2. inter2 = **reconstruct(Rb.inter1)**
3. inter3 = **select(inter2,40,50)**
4. join_input_R = **reconstruct(Rc,inter3)**
5. inter4 = **select(Sa,50,65)**
6. inter5 = **reconstruct(Sb,inter4)**
7. join_input_S = **reverse(inter5)**
8. join_res_R_S = **join(join_input_R,join_input_S)**
9. inter6 = **voidTail(join_res_R_S)**
10. inter7 = **reconstruct(Ra,inter6)**
11. result = **sum(inter7)**

join(join_input_R,join_input_S)



voidTail(join_res_R_S)



project m4



Initial Status			Query and Query Plan (MAL Algebra)		
Relation R			Relation S		
Ra	Rb	Rc	Sa	Sb	
3	12	12	17	11	
16	34	34	49	35	
56	75	53	58	62	
9	45	23	99	44	
11	49	78	64	29	
27	58	65	37	78	
8	97	33	53	19	
41	75	21	61	81	
19	42	29	32	26	
35	55	0	50	23	

select sum(R.a) from R, S where R.c = S.b and
 $5 < R.a < 20$ and $40 < R.b < 50$ and $30 < S.a < 40$

1. inter1 = **select**(Ra,5,20)
2. inter2 = **reconstruct**(Rb.inter1)
3. inter3 = **select**(inter2,40,50)
4. join_input_R = **reconstruct**(Rc,inter3)
5. inter4 = **select**(Sa,50,65)
6. inter5 = **reconstruct**(Sb,inter4)
7. join_input_S = **reverse**(inter5)
8. join_res_R_S = **join**(join_input_R,join_input_S)
9. inter6 = **voidTail**(join_res_R_S)
10. inter7 = **reconstruct**(Ra,inter6)
11. result = **sum**(inter7)

Initial Status

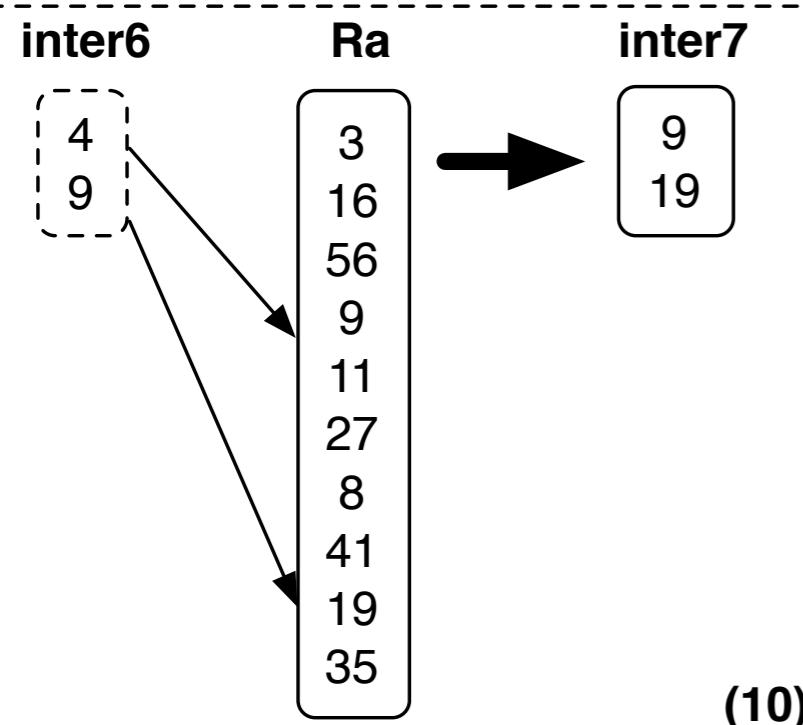
Relation R			Relation S	
Ra	Rb	Rc	Sa	Sb
3	12	12	17	11
16	34	34	49	35
56	75	53	58	62
9	45	23	99	44
11	49	78	64	29
27	58	65	37	78
8	97	33	53	19
41	75	21	61	81
19	42	29	32	26
35	55	0	50	23

Query and Query Plan (MAL Algebra)

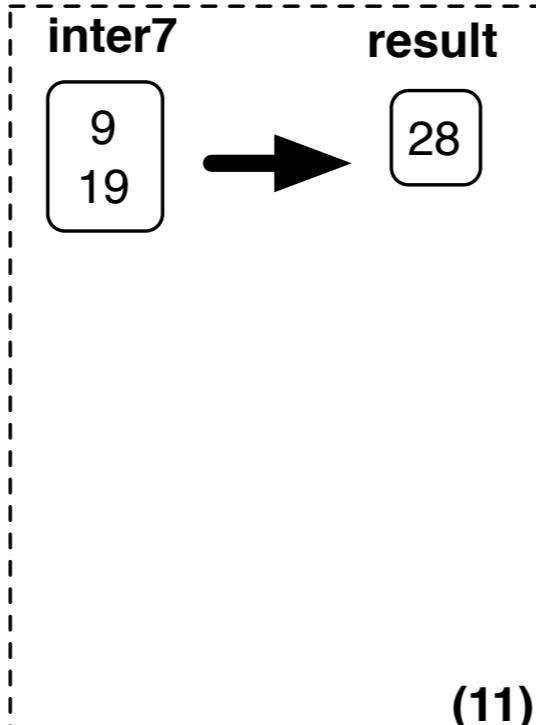
select sum(R.a) from R, S where R.c = S.b and
5 < R.a < 20 and 40 < R.b < 50 and 30 < S.a < 40

1. inter1 = **select(Ra,5,20)**
2. inter2 = **reconstruct(Rb.inter1)**
3. inter3 = **select(inter2,40,50)**
4. join_input_R = **reconstruct(Rc,inter3)**
5. inter4 = **select(Sa,50,65)**
6. inter5 = **reconstruct(Sb,inter4)**
7. join_input_S = **reverse(inter5)**
8. join_res_R_S = **join(join_input_R,join_input_S)**
9. inter6 = **voidTail(join_res_R_S)**
10. inter7 = **reconstruct(Ra,inter6)**
11. result = **sum(inter7)**

reconstruct(Ra,inter6)



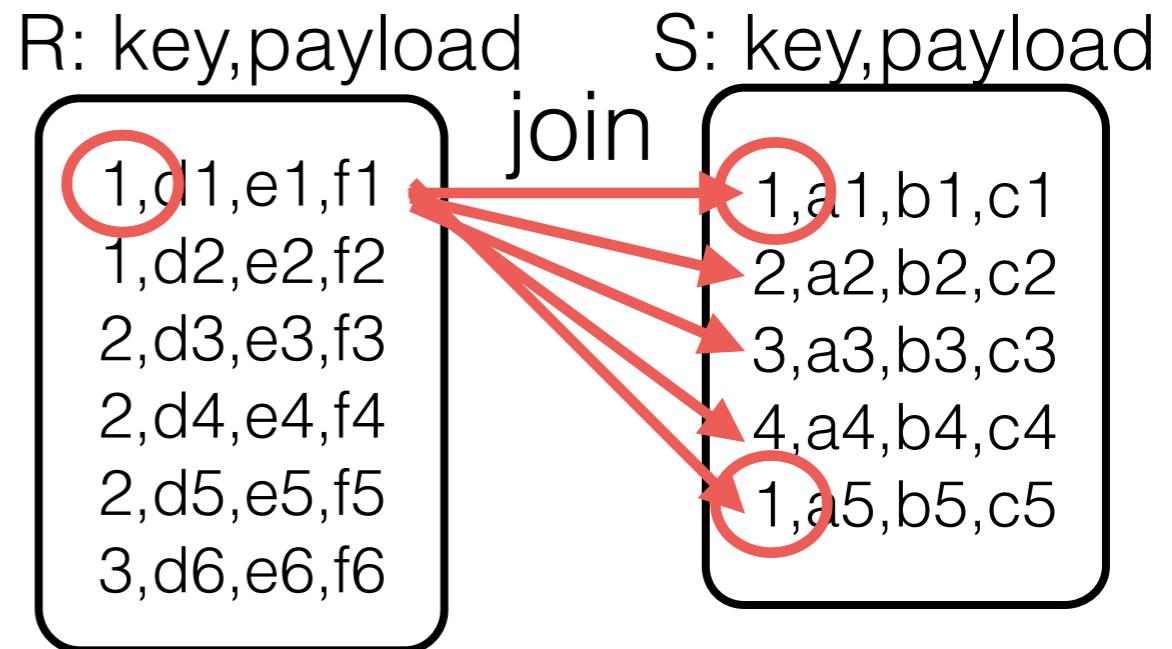
sum(inter7)





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

data/results
one column-at-a-time



CPU

level 1

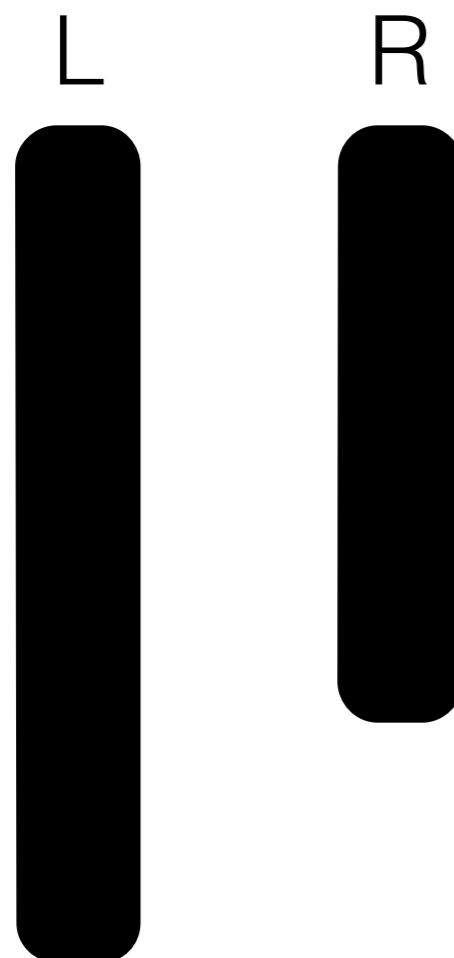
level 2

- 1) design a nested loops join algorithm and give its cost
- 2) which table do we start from?
- 3) think of optimizations to minimize L1 misses
- 4) can we use sorting?

Assume we have a b-tree on one of the join columns:

- 5) how would you use the b-tree to do a join?
- 6) would you use the b-tree or the nested loops join?

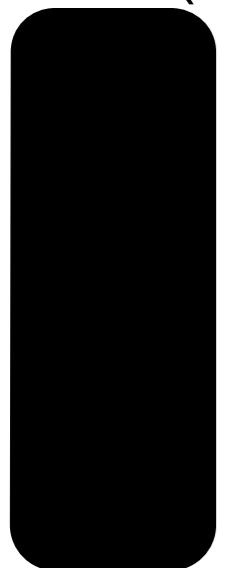
all R columns > L1, all S columns > L1, R+S << L2, L1 block = L2 block
CPU can read directly from L1 only



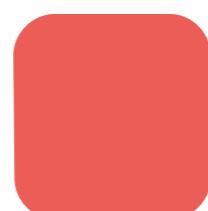
```
new resL[]; new resR[]; k=0
for (i=0;i<L.size;i++)
    for (j=0;j<R.size;j++)
        if L[i]==R[j]
            resL[k]=i
            resR[k++]=j
```



outer(r)



inner(s)



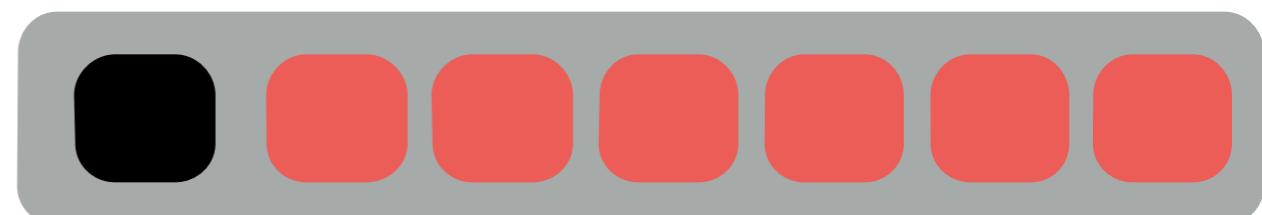
join

probe the small side
for (possibly) better **data locality**

stream outer pages
(try to) hold inner pages

$r/p+s/p, s < L_1$

Level1



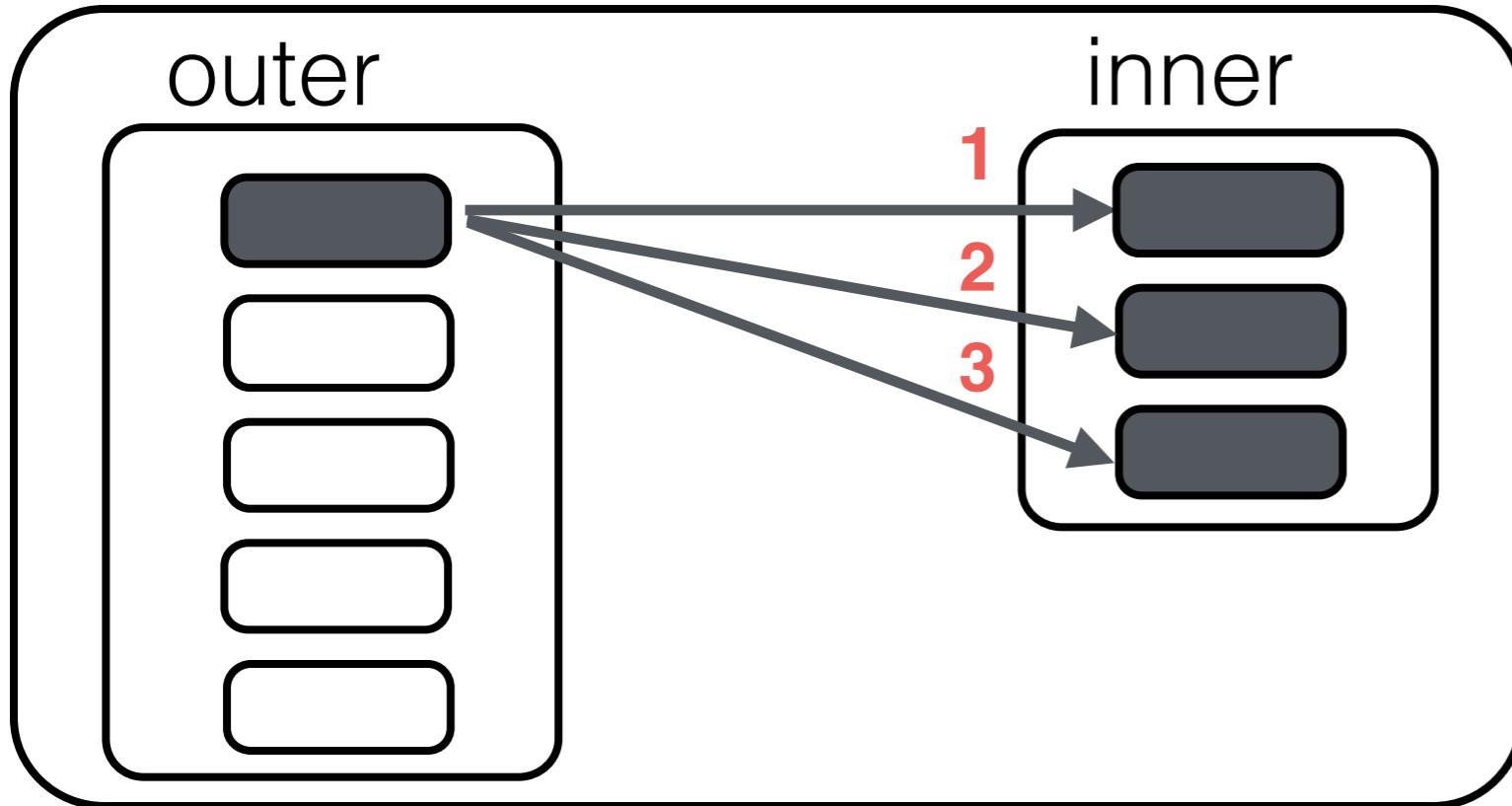
Level2



res

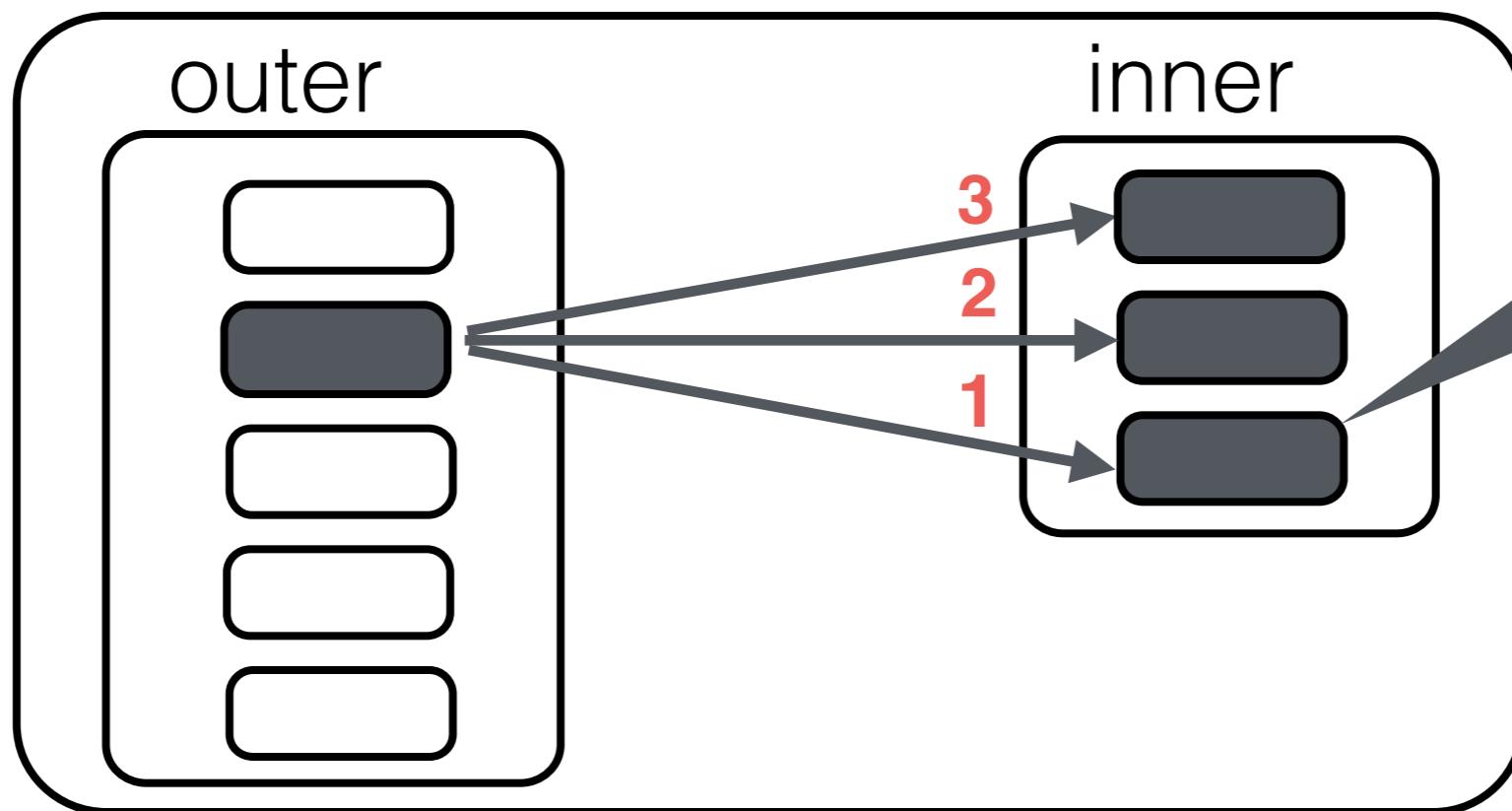


phase 1



zig zag

phase 2



will still be in LLC

...



blocked nested loops

```
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
```

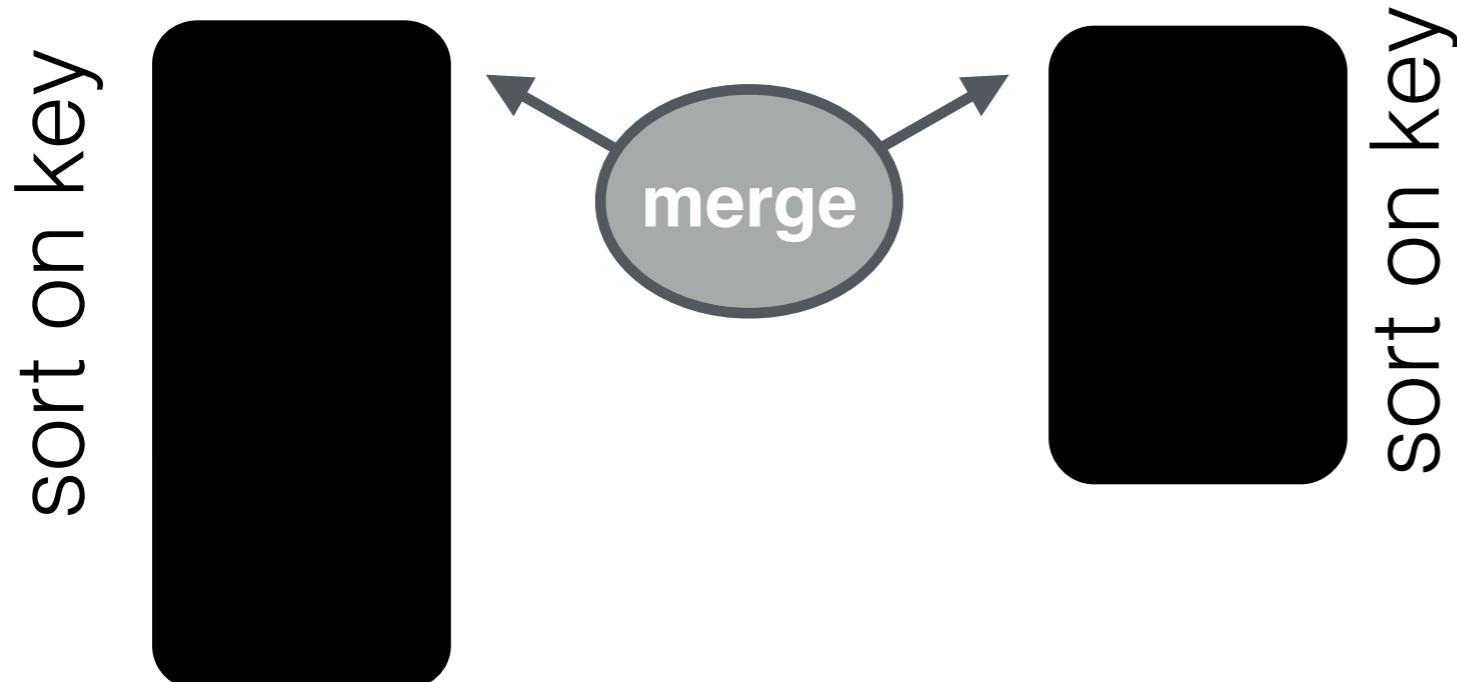
comp $O(L \times R)$
I/O $O(L/p + L \times (R/p))$

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

comp $O(L \times R)$
I/O $O(L/p + (L/p) \times (R/p))$



sort merge join

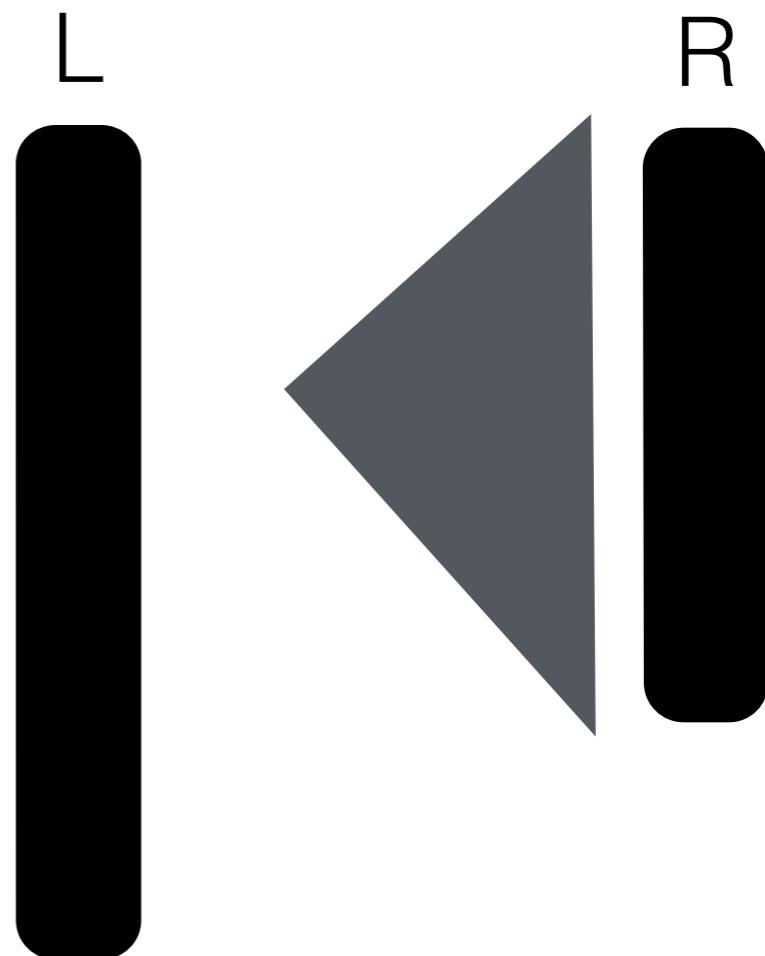


perfect if data is sorted or
needs to be sorted anyway
for the rest of the plan

```
while left and right still have tuples
  if left.val < right.val left++
  else if left.val > right.val right++
  else
    add to result, left++, right++
```

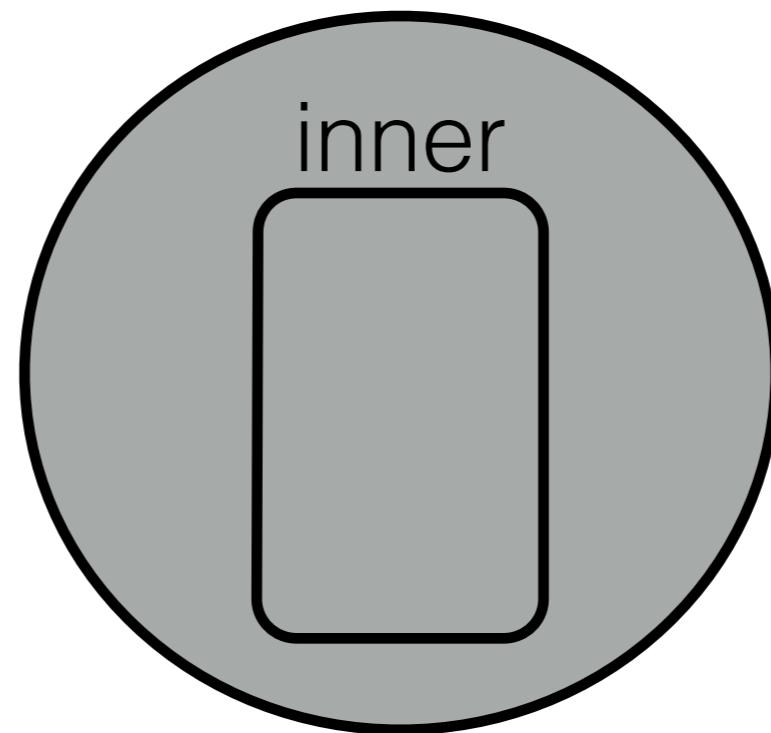
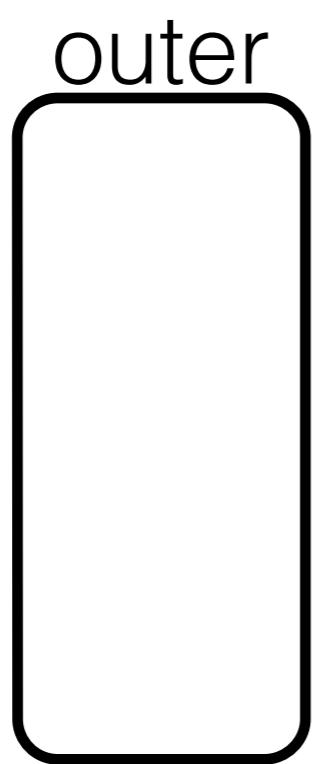


```
new resL[]; new resR[]; k=0
for (i=0;i<L.size;i++)
    jk=R.btree.probe(L[i])
    if (jk!=null)
        resL[k]=i
        resR[k++]=jk.pos
```



even with buffering, we might have to read several (random) pages per probe (to traverse the tree)





oracle
search in $O(1)$



what if not all data fits in **main-memory**?

what if not all data fits in **L3 cache**?

what if not all data fits in **L2 cache**?

what if not all data fits in **L1 cache**?

can we exploit the **>1 cores** in modern **CPUs**?





joins

(project=m4)

what to do in m4?

nested loops and hash joins

cache conscious and multi-core (within reason)



textbook: chapters 4, 14

joins
DATA SYSTEMS

prof. Stratos Idreos

