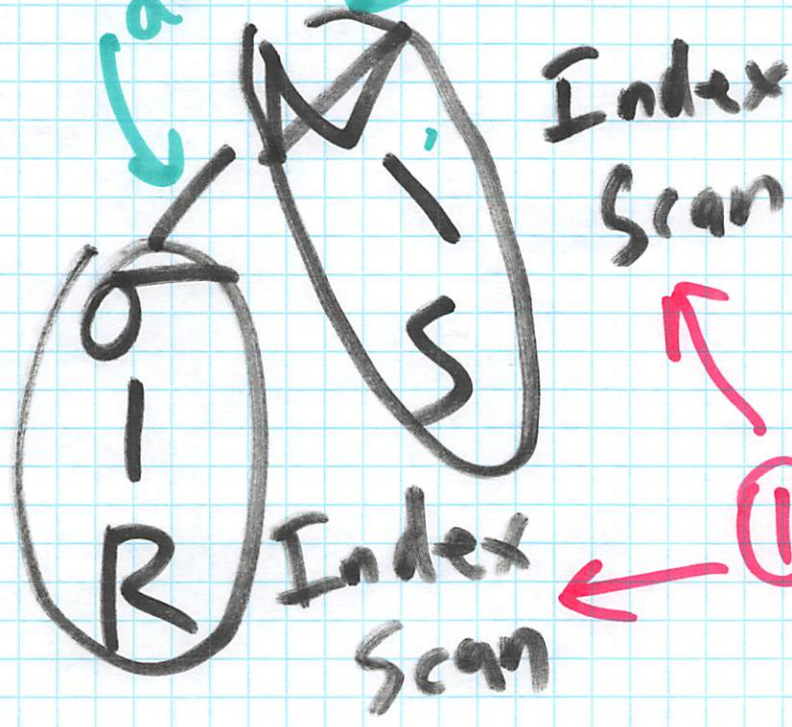


② Cardinality estimation
 How much data? / How much data?



RMS Costs

- 1P H Join
- 3P I Join
- Sort Merge Join
- INL Join
- agg → 2P H Join

IOs

$$\frac{\text{merge}}{0} + \frac{\text{sort}}{|R| \log |R| + |S| \log |S|}$$

can skip if R (or S) is already sorted

either #records R or #records R * log |S|

$$2(|S| + |R|)$$

1p HT

HashTable t

for s in S :

add s to t

for r in R :

find matches for r in t :

emit(match r)

1p TT

In Mem
tree

2p HT

for s in S :

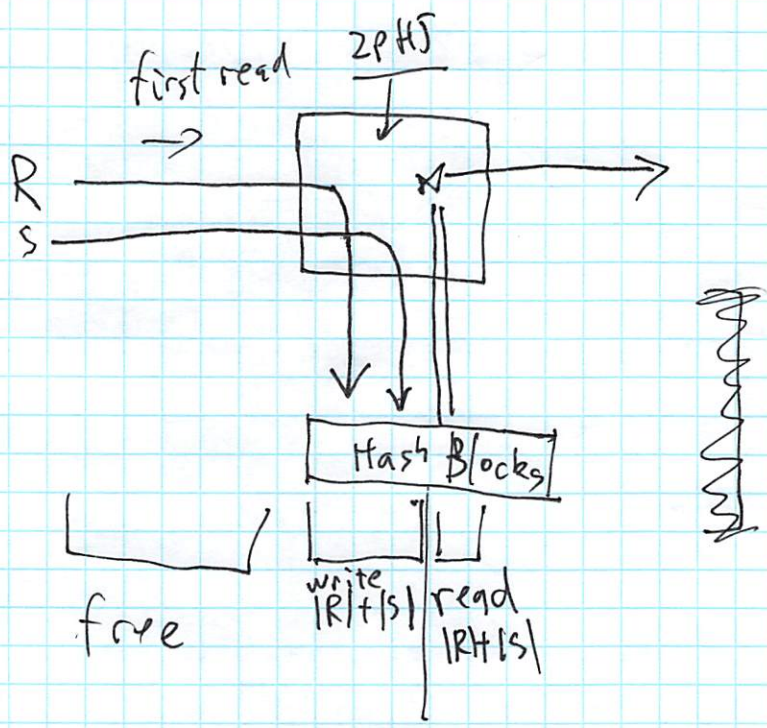
write s to on-disk bucket $\text{hash}(s)$

for r in R :

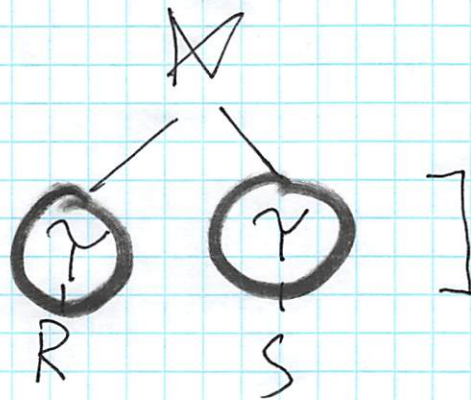
write r to on-disk bucket $\text{hash}(r)$

for b in bucket:

in mem join on all tuples in b



SMS

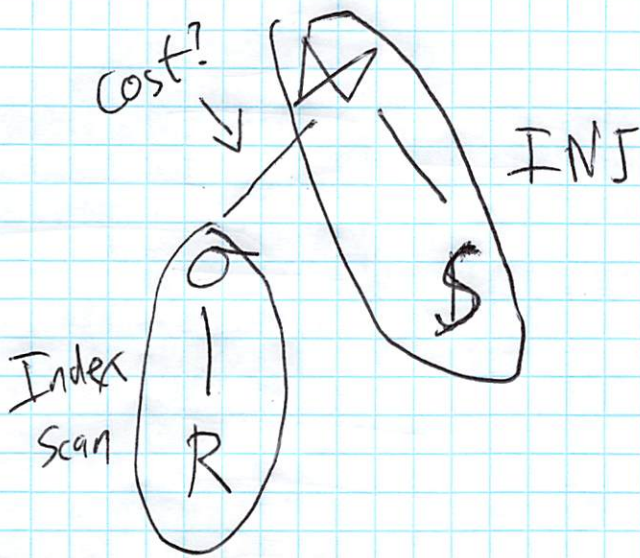


INL

```
for r in R:
  for s in S:
    index lookup(r)
    emit(r.s)
```

} loop for # of records in R

cost depends on type of index
Hash: $O(1) \rightarrow |$
Tree: $\log |S|$



Π, γ : no change

$R \cup S$: sum of $|R|, |S|$

$R \times S$: product of $|R|, |S|$

$\sigma_c R$

$A = \text{const}$

$A < \text{const}, \text{const} < A < \text{const}$

$A \neq \text{const}$

$A_1 = A_2$

A = Const

Idea 1

SELECT A, COUNT(*)
FROM R
GROUP BY A

} statistic

CTMI

Idea 2

Store ^{only} Max # of rows with a given A

↑ upper bound on $|\sigma_{C=A} R|$

Avg : avg of $|\sigma_{C=A} R|$

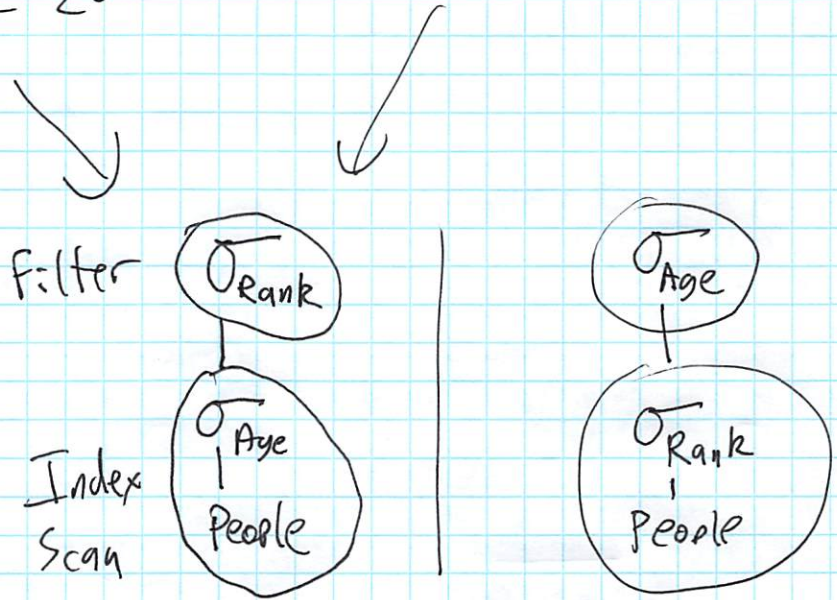
Min : lower bound

Idea 3

Store min, max of A]

SELECT Name
FROM People
WHERE Rank = 3
AND Age = 20

People
Index on Rank
Index on Age



Optimal for (3, 20)

Rank	#	Age	#
1	3	19	1
2	2	20	4
3	3	21	2
		22	1

	Rank	age
Count	max 3	max 4
	avg 2.6	avg 2
	min 2	min 1

	Rank	age
Values	min 1	min 19
	max 3	max 22

overall size 8

8

(σ_{Rank} R) vs (σ_{Age} R)

Using Statistics

Idea 2 → no more than 3 tuples in $\sigma_r R$ ✓
no more than 4 tuples in $\sigma_a R$

Idea 2.1 → avg of 2.6 tuples in σ_r
2 tuples in σ_a ✓

Idea 3 → avg of ~~$\frac{8}{3}$~~ tuples in σ_r
avg of $\frac{8}{4}$ tuples in σ_a ✓