

ITERATORS & JOINS:

`hasNext()`: Is there more data?

`next()`: Fetch the next piece of data

`[R = Pgs * P/Pgs]`

cardinality, total no records $|R|$, data pgs.

= total no of pgs * no of records/pg
 $|R|$ P_f

DIFFERENT JOINS

compute both ways for R & S & check which one gives smaller cost.

SMJ: [Sort Merge Join]

cost: cost of sorting R + cost of sorting S + cost of merging
 cost of sorting a relation w N pgs: $2N * \text{ceil}(1 + [\log(N/B)])$
 $B=1$

cost of merging:

$WC: |R| + |R| * S$ Avg C: $|R| + |S|$

GHTJ:

COST: cost to hash [R] & [S] + cost of Native Hash Join

hash until $\leq B-2$

write: $\text{ceil}(\text{no of pgs to write}/B)$

total I/Os: $\text{ceil}(\text{no of pgs to write}/B) * B-1 + \text{Read I/Os}$

Native: [R] + [S] I/Os - **last pass's I/O write**

BLNJ: $\min(RMS, SMR)$

COST: $[R] + \text{ceil}[R]/(B-2)] * [S]$

SNLJ: $\min(RMS, SMR)$

COST: $[R] + P_R[R][S]$

PNLJ: keep smaller relation b/w R + S

COST: $[R] + [R][S]$ as the outside loop * need at least 2 buffer pgs

INLJ: [Index Nested]

COST: $[R] + |R| * (\text{cost to lookup matching record})$

Clustered: 1 I/O / pg of matching record

Unclustered: 1 I/O / matching record.

NOTE:

INLJ, SMJ joins based on data being & GHTJ - ordered a certain way (equality predicate)

BLNJ, SNLJ - handles all conditions.

Alt 1 has I/O cost of \Rightarrow default sorted.

(height)

cost to reach above leaf + num of leaves read.

Alt 2/3 has I/O cost of

cost to reach above leaf + num of leaves read + num data pages read.

Cost to reach above leaf is generally

Height of tree, h,

but when given matching tuples its h+1.

EXAMPLE: Grace Hash Join:

Grace Hash Join:

Justice League = 100 pgs $\frac{1}{2}$ of records = earth: — equality predicate.

T = 200 pgs

B = 12

J

perI read 100 pgs = 100 I/Os

perI write: $\text{ceil}(\text{no of pgs to write}/B)$ since only half $\text{ceil}(200/11) = 19$

$\therefore \text{ceil}(50/11) = 5$

$\therefore \text{total I/Os: } 5 \times B-1 = 5 \times 11 = 55$

$\therefore \text{total cost: } 100 + 55 + 200 + 20 \cdot 9 + 55 + 20 \cdot 9 = \underline{\underline{828}} \text{ I/Os}$

T

200 pgs = 200 I/Os

$\text{ceil}(200/11) = 19$

total I/Os: 19×11

$= 204$

QUERY OPTIMIZATION:

Given a plan some things we can do to cheapen a query:

① push selections / projections down a tree

② materialize immediate relations (write to a temp file)

→ results in additional write I/Os but better

③ use indices INLJ. in the long run

SELECTIVITY: % of tuples selected by a predicate.

STEPS: selectivity * # tuples \downarrow * if not stated use 1/10

Predicate	Selectivity	Assumption
c = v	1 / (number of distinct values of c in index)	We know c .
c = v	1 / 10	We don't know c .
c1 = c2	1 / MAX(number of distinct values of c1, number of distinct values of c2)	We know c1 and c2 .
c1 = c2	1 / (number of distinct values of c1)	We know c1 but not [other column].
c1 = c2	1 / 10	We don't know c1 or c2 .

INEQUALITIES OF INTEGERS

Predicate	Selectivity	Assumption
c < v	$(v - \min(c)) / (\max(c) - \min(c) + 1)$	We know max(c) and min(c). c is an integer.
c > v	$(\max(c) - v) / (\max(c) - \min(c) + 1)$	We don't know max(c) and min(c). c is an integer.
c <= v	$(v - \min(c)) / (\max(c) - \min(c) + 1) + (1 / c)$	We know max(c) and min(c). c is an integer.
c >= v	$(\max(c) - v) / (\max(c) - \min(c) + 1) + (1 / c)$	We don't know max(c) and min(c). c is an integer.
c < v	1 / 10	
c > v	1 / 10	

INEQUALITIES OF FLOATS:

Predicate	Selectivity	Assumption
c >= v	$(\max(c) - v) / (\max(c) - \min(c))$	We know max(c) and min(c). c is a float.
c >= v	1 / 10	We don't know max(c) and min(c). c is a float.
c <= v	$(v - \min(c)) / (\max(c) - \min(c))$	We know max(c) and min(c). c is a float.
c <= v	1 / 10	We don't know max(c) and min(c). c is a float.

CONNECTIVES:

Predicate	Selectivity	Assumption
p1 AND p2	S(p1)*S(p2)	Independent predicates
p1 OR p2	S(p1) + S(p2) - S(p1)*S(p2)	Independent predicates
NOT p	1 - S(p)	

2 cases: if predicates dependent on each other.

$P(A \cup B) = P(A) + P(B) - P(A \cap B)$ - dependent.

$P(A \cup B) = P(A) + P(B)$ - total exclusivity

EXAMPLES:

Select * FROM R INNER JOIN S

ON R.col = S.col WHERE R.col < 5

$B=12$ [R] = 20 pgs [S] = 10 pgs Full scan for R:

What's the cost of R
 BNLJ S?
 $20 + \text{ceil}(5/12-1) * 10 = 20 + 1 * 10 = 30$ I/Os

but :: we push down only 5 is used.

$= 20 + 5 * 10 = 70$ I/Os

TRANSACTIONS & CONCURRENCY

A : **Atomicity**: All ops OR No ops [commit or abort]

C : **Consistency**: Data is constant & maintained

I : **Isolation**: 1 transaction at a time

D : **Durability**: transaction persists after committing

Throughput: no of transactions / unit time

Latency : response time / transaction

Serial Schedule: each transaction runs from start to finish w/o intervening actions from other transactions.

We use serializability to provide isolation

Check for Serializability.

- same initial reads - starts the same

- same winning writes - ends the same

CS186 MT2

Table A: takes 50 I/Os to perform a scan

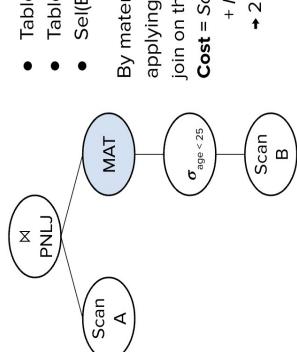
Table B: takes 100 I/Os to perform a scan

$\sigma_{age < 25} = 0.5, [B] = 100$

By materializing the intermediate relation, we're applying $\sigma_{age < 25}$ before PNLU, and performing the join on the result of the selection.

$$\text{Cost} = \text{Scan A (50)} + \text{Scan B (100)} + \text{Materialize (100 * 0.5)} + \text{PNLU (50 * 100)}$$

→ 2,700 I/Os in total



2 schedules in conflict:

- diff transacn

- same resource

- 1 write at least

conflict equivalent:

- ops within each transaction is same b/w both schedules

- every pair of conflict ops are ordered same way.

S is conflict = some serial schedule

⇒ S is serializable.

- no cycles in dependency graph

View Serializable:

- same initial read

- same & dt reads

- same winning writes.

NOTE:

Blind Writes:

writes w/o RS b/w them.

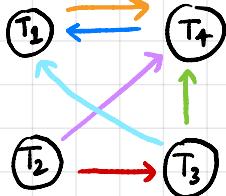
Conflict Serializable cards abt

order of BW!

T1	R(A)		W(C)	R(B)
T2	R(A)	R(B)		
T3	R(A)		W(B)	
T4		W(A)		R(C)

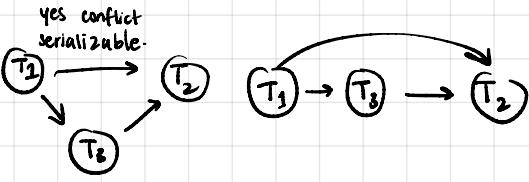
Reads + writes go together.

- i) R(A) before W(A)
- ii) R(A) before W(C)
- iii) W(A) before R(A)
- iv) R(B) before W(B)
- v) W(C) before R(C)
- vi) W(B) before R(C)



- cycle in graph so not conflict serializable.

T1	R(A)			R(C)	W(B)	COMMIT
T2		R(A)	W(A)			COMMIT
T3	R(A)	R(C)			W(B)	COMMIT



SIMPLE LOCKING:

S: lets transaction read (R) [Shared]

- many transac hold S at once

X: lets transaction modify (W)[Exclusive]

	NL	S	X
NL	✓	✓	✓
S	✓	✓	✗
X	✓	✗	✗

DEADLOCK: $T_2 \xrightarrow{\text{WAITS}} T_2$

Waits - for graph

if T_i holds lock that conflicts w lock T_j wants

" T_j waits for T_i " $T_j \longrightarrow T_i$

PRIORITY: current time - start time

(older ↑ & Priority ↑)

high low

wait-die: T_i wants a lock but T_j holds conflict

T_i high? waits for T_j

T_i low? aborts (die)

high low

wound-wait: T_i wants a lock but T_j holds conflict

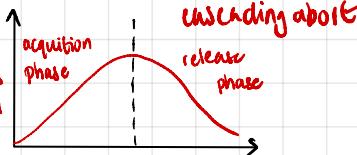
T_i high? ABORT T_j (wounds)

T_i low? waits for T_j

2 PHASE LOCKING (2PL)

problem

PROBLEM:
2PL lets new locks held
read new vals
but trans commits



STRICT 2PL:

- allows release at end of transaction
acquire b/t abort/commit.



Consider the following schedule:

T1	R(A)				R(C)	W(B)	COMMIT
T2			R(A)	W(A)			COMMIT
T3	R(A)	R(C)				W(B)	COMMIT

Which of the following changes to the schedule above, will result in a schedule that is possible using strict 2PL?

- A) Make T1 abort instead of commit
- B) Make T2 abort instead of commit
- C) Remove R(A) from T1 and T3
- D) Remove T2 from the schedule
- E) Swap W(B) with W(D) for T3
- F) No change needed, it is already strict 2PL
- G) None of these changes provide strict 2PL

Consider the following schedule:

T1	R(A)	W(A)		R(B)	W(B)	W(A)	COM
T2			R(B)				COM
T3				R(C)		W(C)	COM

Does the schedule follow 2PL?

Yes, no transaction acquires a lock after releasing a lock.
Consider T2: it needs to release its locks until it is ready to commit. This means T1 should be unable to acquire X(B) until T2 commits.

Does the schedule follow strict 2PL?

No, under strict 2PL, T2 would not release its locks until it is ready to commit. This means T1 should be unable to acquire X(B) until T2 commits.

MULTI-GRANULARITY LOCKING:

Compatibility matrix:

Parent-child lock relationship:

Parent	Pass child locks				
NL	None				
S	None				
X	None				
IS	IS, S				
I X	I X, X, I S, S				
SIX	I X, X				

* IS: Intent to acquire Shared lock at a lower level.

* IX: Intent to acquire exclusive lock at a lower level

* SIX: Shared Intent to acquire exclusive lock at a lower level

→ can read entire table + acquire X later
→ prevents writes but allows reads!

We have a table, T, with pages A and B. Page A has tuples A1, A2 and Page B has tuples B1, B2.

1) List the locks T1 has at the end of timestep 5 in the order they were acquired.

2) What locks are on T2's waiting queue at the end of timestep 5?

time	1	2	3	4	5
T1	R(A1)		R(A)		W(A)
T2		W(B1)		W(A2)	

1)

After timestep 1: T1 has IS(T), IS(A), S(A1)

After timestep 3: T1 has IS(T), S(A) (T1 escalates on the page-level lock for A)

After timestep 5: T1 has IX(T), X(A) (T1 promotes S(A) to X(A) - note: promote does not affect acquisition order)

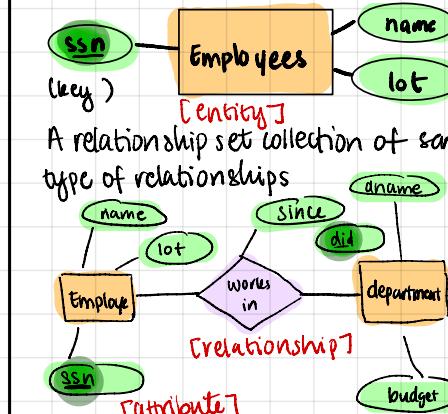
2)

T2's waiting queue: acquire IX(A)

Note that T2 will want to acquire X(A2) eventually, but it gets blocked on trying to acquire IX(A) before it gets the chance to try to acquire X(A2).

ER DIAGRAMS:

entities are real world obj described with attributes



many-to-many .



many-to-one :



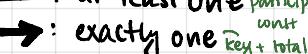
one-particip



one-to-one



one-to-many



Weak entities: can be identified uniquely only w the key of another entity.

entity when combined w partial key (---) identifies other entity's key

must be 1-to-many rel. w total partic.

Inference Rules:

$$\begin{aligned} y \leq x &\Rightarrow x \rightarrow y & x \rightarrow y \text{ and } y \rightarrow z &\Rightarrow x \rightarrow z \\ x \rightarrow y &\Rightarrow xz \rightarrow yz & x \rightarrow y \text{ and } x \rightarrow z &\Rightarrow x \rightarrow yz \\ x \rightarrow yz &\Rightarrow x \rightarrow y \text{ and } x \rightarrow z \end{aligned}$$

* BCNF always lossless.

BCNF (Boyce-Codd Normal Form): no redundancy R in BCNF if for every FD $X \rightarrow A$ that holds over R, either ① $A \subseteq X$ or ② X is a superkey.

also: for each FD $X \rightarrow Y$ in F⁺:

if $X \rightarrow Y$ violates BCNF → decompose R into (R-X⁺) U X and X⁺.

* don't want lossiness!
Lossless iff F⁺ contains
④ (X intersects Y) → Y or ⑤ (X intersect Y) → Y.

FUNCTIONAL DEPENDENCIES

functional dependency

$$X \rightarrow Y \quad (X \text{ determines } Y)$$

superkey X is superkey of R if $X \rightarrow$ all attributes of R

candidate key: at set of R

keys that determines all columns in a relation

attribute closure: of an attribute X given a set of FDs is X⁺

Relation R is BCNF if:

for every FD $X \rightarrow A$ that holds over R either $A \subseteq X$ or X is a superkey.

never lossless.

- no redundancy in BCNF!

lossless - NOT able to constrain original relation.

Closure: $F = \{A \rightarrow B, AB \rightarrow AC, BC \rightarrow PD, DA \rightarrow C\}$

$$A^+ = A \rightarrow B \quad (AB) \quad AB^+ = AB \rightarrow AC \quad ABCD \Rightarrow$$

$$BC \rightarrow BD \quad (ABCD)$$

superkeys b/c A is candidate key.

$$BC^+ = BC \rightarrow BD \quad (BCD)$$

$A \rightarrow$ candidate key b/c $A^+ = \text{minimal set of keys to cover all symbols in } F$.

RECOVERY:

NO STEAL:

if pg not flushed to disk yet
transn cannot steal
SLOW: other transn wait for
pgs.
few uncommitted transacs

STEAL:

evict/flush pg w/ uncommitted
updates
FAST but on crash, need to UNDO
unwanted updates (for atomicity)

NO FORCE:

- commit w/o forcing dirty pg to
disk
FAST but on crash need to KEDO
for durability.
- guaranteed during buffer
manager's MRU pg replacement
policy.
- less I/Os Durability.

FORCE:

- force dirty pgs to disk on
COMMIT.
slow: yes to write pg to disk
Atomicity

No steal

Steal

No Steal

Steal

no undo

undo

redo

redo

No F

F

no undo

undo

no redo

no redo

F

Performance

Logging / Recovery

Steal, No Force: Write Ahead Logging

- log records written to disk before
data pg to disk.
- all log records written to dis when
transacn commits.

Undo Logging:

we want to undo if not
committed.

4 types: start commit abort update.

- if T modifies element X, update log record
written to disk before dirty page. [STEAL]
- if T committed, WRITE to disk before commit
scan log from end to find if T completed. [FORCE]
if T not completed write X to disk.

Redo Logging:

Analysis

Redo, Undo-

Analysis: rebuild the transn table, Dirty Page
table. (DPT)

if not END, add to transn table, set last
log sequence number. (LSN)

if COMMIT or ABORT, transn status change
if UPDATE, not in DPT, add to DPT, recLSN to
if END, remove from transn table. LSN.

Ideal for
I/O performance

ARIES: Recovery system → steal, no force

- uncommitted data can be flushed (\uparrow mem)
- data pgs don't need to be flushed b4 xact
can commit (\downarrow latency of xact).

The Log: (WAL) → Log recs are assigned an
LSN (log seq. #); prev LSN - LSN of prev
log rec for xact. Compensation Log Rec (CLR)
log entry undoing another non-CLR log entry.

Memory

- Xact Table
- last LSN
- status
- DPT → recLSN
- Flushed LSN
- Buffer Pool
- Log Tail

Log

- LSN
- prev LSN
- XID
- Type
- PgID
- len
- offset
- LA-IM
- ACT-IM

DB

- Data Pgs
- Loc w/ Pg LSN
- Master Rec.

Memory

- Xact Table → tracks active xactions.
- DPT: tracks dirty pgs in buffer.
- rec LSN - recovery

LSN indicates first log rec that caused pg to ~~be dirty~~
be dirty. Flushed LSN: max LSN successfully flushed

Log Recs: update, commit, abort, end, CLR,
begin checkpoint, end checkpoint.

* ea rec associated w/ a txn.

Pages: Pg LSN = LSN of last rec modifying a pg.
→ stored on pg itself, flushed w/ pg, can differ
in mem v. on disk pg.

ARIES: Commit: ① write commit rec to log
② flush log up to commit rec ③ write end rec

ARIES: ABORT: ① write abort rec to log ② undo
all changes made by txn. (start from last LSN
go back via prev LSN) (undo all log rec that can
be undone by adding CLR to log) ③ write end rec.

ARIES: Recovery: → analysis, redo, undo →

Step 1: Analysis - start from begin-chkpt.

↳ update txrn table + DPT → end all committing
txns and abort all running txns. * smallest LSN
for ea. dirty pg.

Step 2: Redo - start from smallest/oldest rec LSN
(in DPT) → redo all update and CLR recs UNLESS:

- a) pg not in OPT → disk up to date already.
- b) rec LSN of pg > LSN → no need to redo;
change was already flushed.

c) pg LSN ≥ LSN ⇒ pg LSN is authority to do flush.

Step 3: Undo - start from largest t/most recent
last LSN (in txrn table) → undo all recs of aborting
txns, skipping over CLR → write CLR as do undos.
Abort all txns in xact table, not already committing.

ARIES Checkpointing: ~~don't need to start from beginning every time!~~

↳ fuzzy checkpoints → taken while other
txns are running + making changes.

→ Write a Begin Chkpt rec (latest rec can
start analysis from) ; write end chkpt rec
later (contains DPT/xact table valid @ the time
of begin checkpoint).

DISTRIBUTED TRANSACTIONS

2 PC: 2 PHASE COMMIT

2 PC: Handling Failures

participants

coordinator

coordinator thinks participants went down

↳ if p didn't vote ABORT

↳ if waiting for ACK, RECOVER

p thinks c went down:

↳ if prep not logged, vote NO

↳ if prep logged, RECOVER

prep* or abort*

vote Yes or no

commit* or
abort*

commit / abort

ACK

end

* = flush to disk • optimized (if commit)

2PC Recovery:

• if have commit/abort log rec: do that action

• if have prep log rec, but no commit/abort

→ send log INQUIRIES about status

• if no prep, commit, abort → abort

Optimized: assume can abort if vote NO, don't

need ack, abort not flushed, abort as soon as

abort is known

NOTES:

→ 2PC cannot handle indefinite failures
if coordinator doesn't recover no participant gets msg

→ By partitioning/centralizing locks database can still use 2PL

- 2PC w/ premmmed abort does not result in fewer records.



NO SQL STUFF

online transaction processing (OLTP): simple lookups

online analytical processing (OLAP): read only queries w/ lots of joins, agg

scaling: ① **Partitioning (sharding)** →

split data among mult mach. To ↑ parallelism [faster W, slower R]

② **Replication** → copy data among mult mach for fault tolerance (slower writes, faster reads)

Distributed Systems Desired Properties: ① **consistency**:

2 clients same view ② **availability** every req must get a response ③ **partition tolerance** even if mach disconnects delay or drop packets etc.

CAP Theorem: only can get 2/3 properties (① v ②)

BASE Semantics: eventual consistency systems gaurantee:

① Basic Atail → R/W available ATM but not always consistent

② soft state → DB can change even w/o inputs.

③ Eventual Consist → given enuf time all reads consistent after updates propagate.

NO SQL Data Models:

• Data model: (key, val) pairs • ops: get(k), put(k, v)

• Distr / Partit (w/ hash func) → no replicate / multi-way rep → propagate update w/ eventual consistency.

Document Stores: (like JSON) → store collec of docs (data)

JSON: * object (key: value); arrays []; atomics support

* self-describing → schema elems in data itself.

	T (row ← enf or read)	R (relation ← enf. for write)
flexibility	flexible, can rep	Less flexible
schema enf.	complex + nested data	Schema fixed
rep	text-based	B in rep (good 4 disk)

(ex) `{ "student": [{ "id": "0001", "name": "me", "classes": ["61A", "61B", "33"], "id": "0002" }] }`

Mongo DB: store field: value.

i.e. doc has special `"_id"` as primary key.

↳ syntax on last pg!

Distributed File System (DFS): for v large file; ea chunk replicated several times ≥ 3 for fault tol.

MAP Reduce: high-level prog model & important for large scale parallel processing.

① Map → extract s/t you care abt from each record.

* MR slower if complex + need many MRs.

② Reduce → agg, summarize, filter, transform

④ write results.

MAP: user provides MAP funct w/ input: (in key, val) and output: bag of (intermediate key, val)

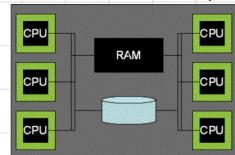
* sys applies map func in parallel to all (in key, val) in the input file

* fault tol by writing intermed files to disk.

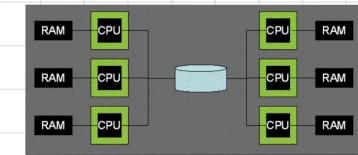
REDUCE: sys groups all pairs w/ same intermed key & passes bag of vals to reduce func user provides reduce func input (intermed key, bag of vals); o/p: bag of o/p vals.

PARALLEL QUERY PROCESSING

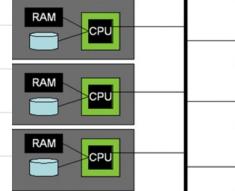
[modern computer] Shared memory



[network file system] shared Disk



Shared nothing [cluster]



Intraquery

Parallelism:

spread work of one query to multiple machines.

Intraquery Parallelism

↳ **Intraoperator**: make one operator run as (ex sorting multiple) quickly as possible.

↳ **Interoperator**: running operators in parallel.

↳ **Pipeline parallelism**: records passed to parent as soon as done

↳ **Bushy tree parallelism**: diff branch of tree run in parallel.

Interquery Parallelism: gives each one diff queries for ↑ throughput & finish more queries

Sharding: each data pg stored on 1 machine

Replication: each data pg on multiple machine

Network Cost: how much data to send over network for op

Parallel sorting / hashing: range partition table, local sort / hash on each machine.

Passes: 1 (partition across machines) + ceil [$1 + \log_{\text{base}} [N/\text{MB}]$] [no of passes to sort table]

SMJ Passes: 2 (1 pass/table to partition) + ceil [$1 + \log_{\text{base}} [\frac{S}{R/\text{MB}}]$] [passes to sort R + S]

Partitioning Data:

Range: key →'s data based on which range key belongs to.

pro: get(key); sorting; SMJ, etc...

con: distr of data may be skewed.

Hash: → data based on hash func.

pro: get(key); hashing; GHJ, etc...

con: distr may be skewed.

Round Robin: cycles thru partitions as data comes in.

pro: even distribution summing all data across mach.

con: not ordered based on key.

Parallel Sorting:

* partition data over mach w/ range partit.

* perform ext. sort on ea. mach. indep.

Parallel SMJ:

* partition data for both rel over mach w/ range partit. (same range for ea. rel).

* perform SMJ on ea. mach. indep.

Parallel Hashing:

* use hash func to partition data over mach.

then, run ext hashing on ea. mach. indep.

Parallel Hash Join:

* use hash partit on both rel, then perf

a normal GHJ on ea. mach. indep.

Symmetric Hash Joins: → streaming hash join algo → doesn't req all tuples of 1 rel to be available b4 starting.

ASYMMETRIC SHUFFLE

if R is partitioned well just partition J then local join at every node & union results.