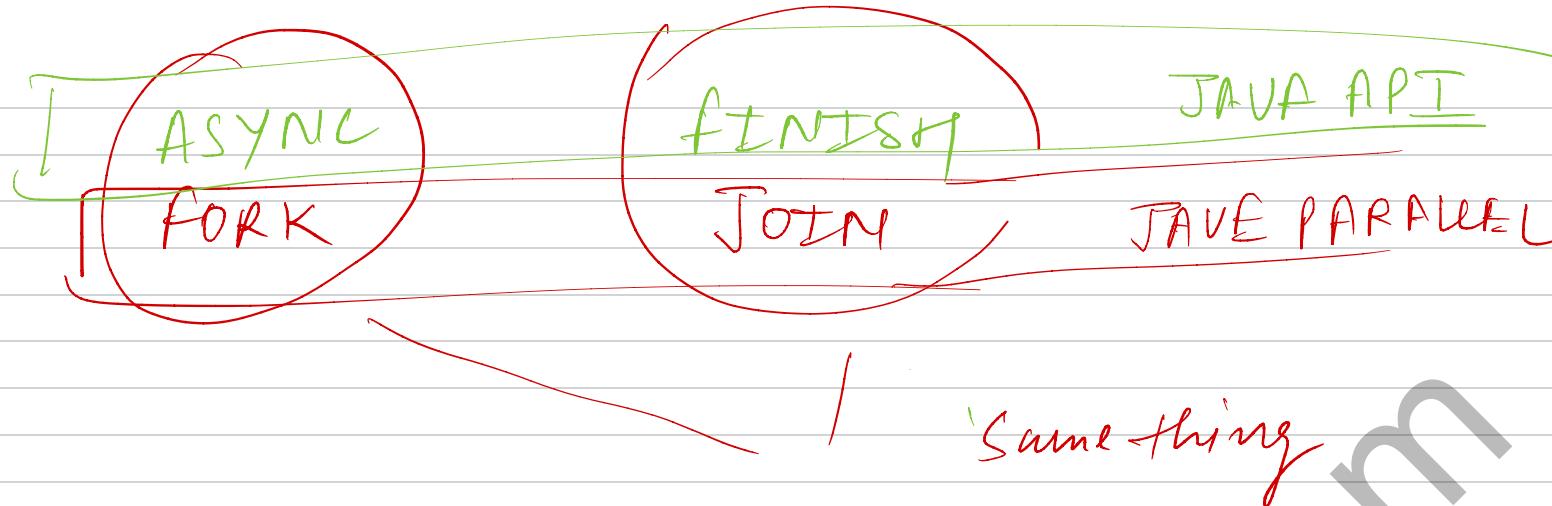
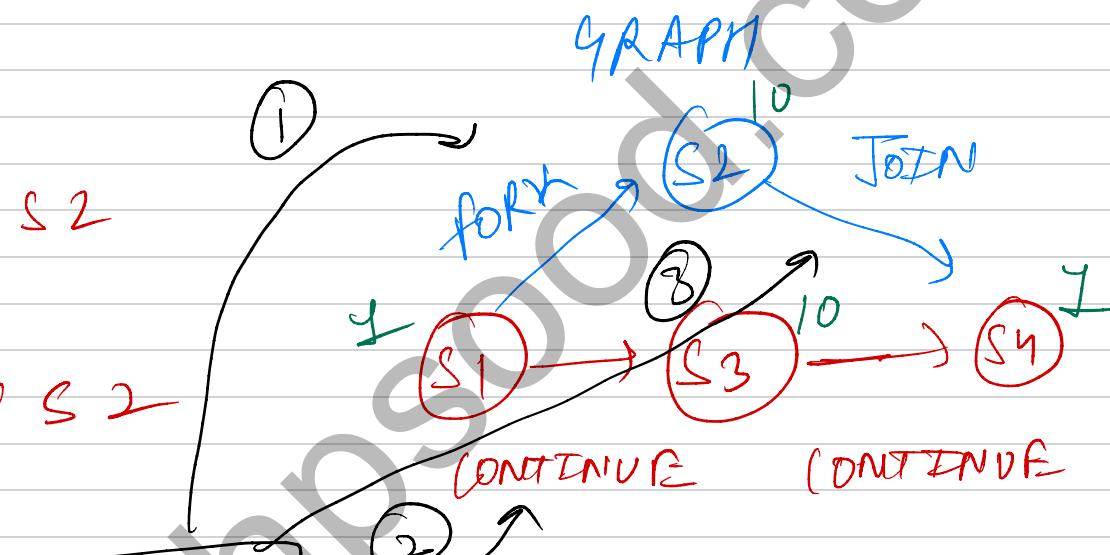


Parallel, concurrent &
distributed programming
in JAVA



Task List

S_1
FORK S_2
 S_3
JOIN S_2
 S_4



→ With these 3 edges we can model any parallel program

→ If there is a path directly b/w two nodes in the graph they cannot run in parallel for ex (S_2) & (S_3) can run in parallel, but (S_2) & (S_4) cannot

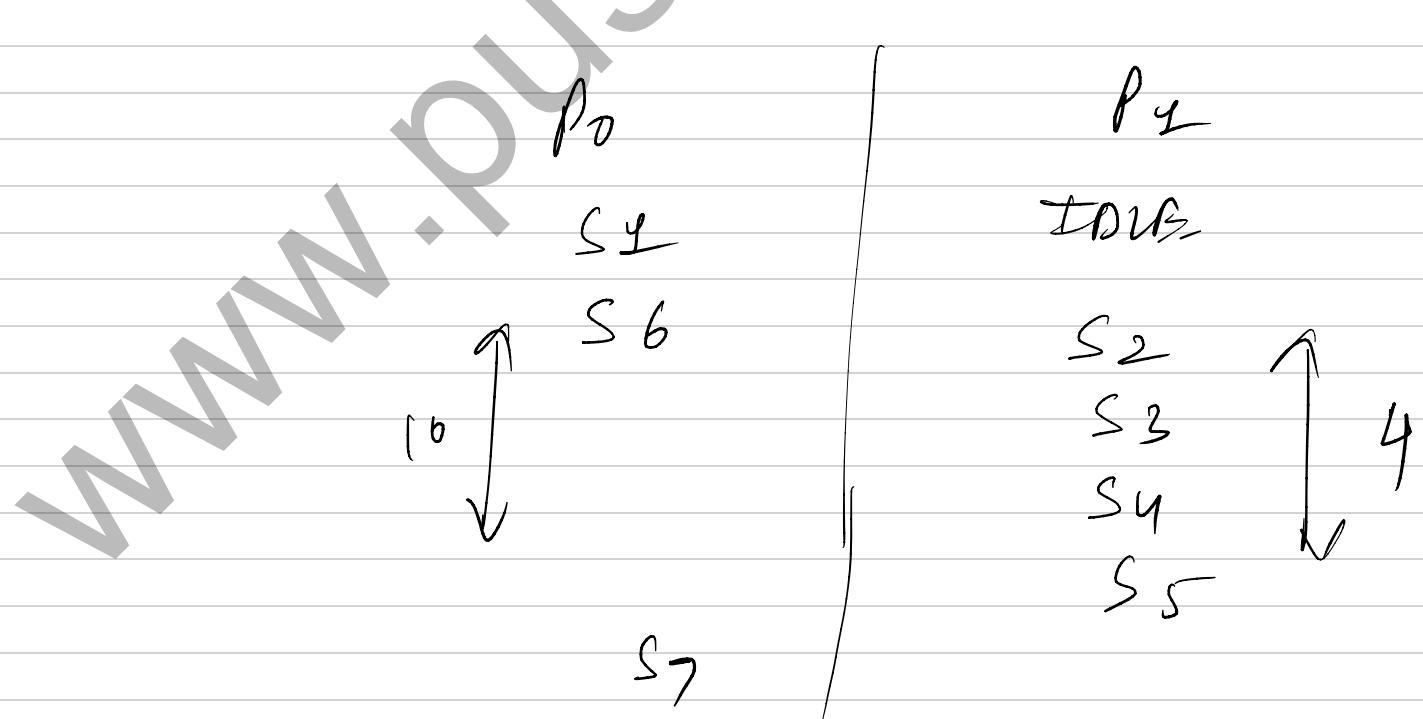
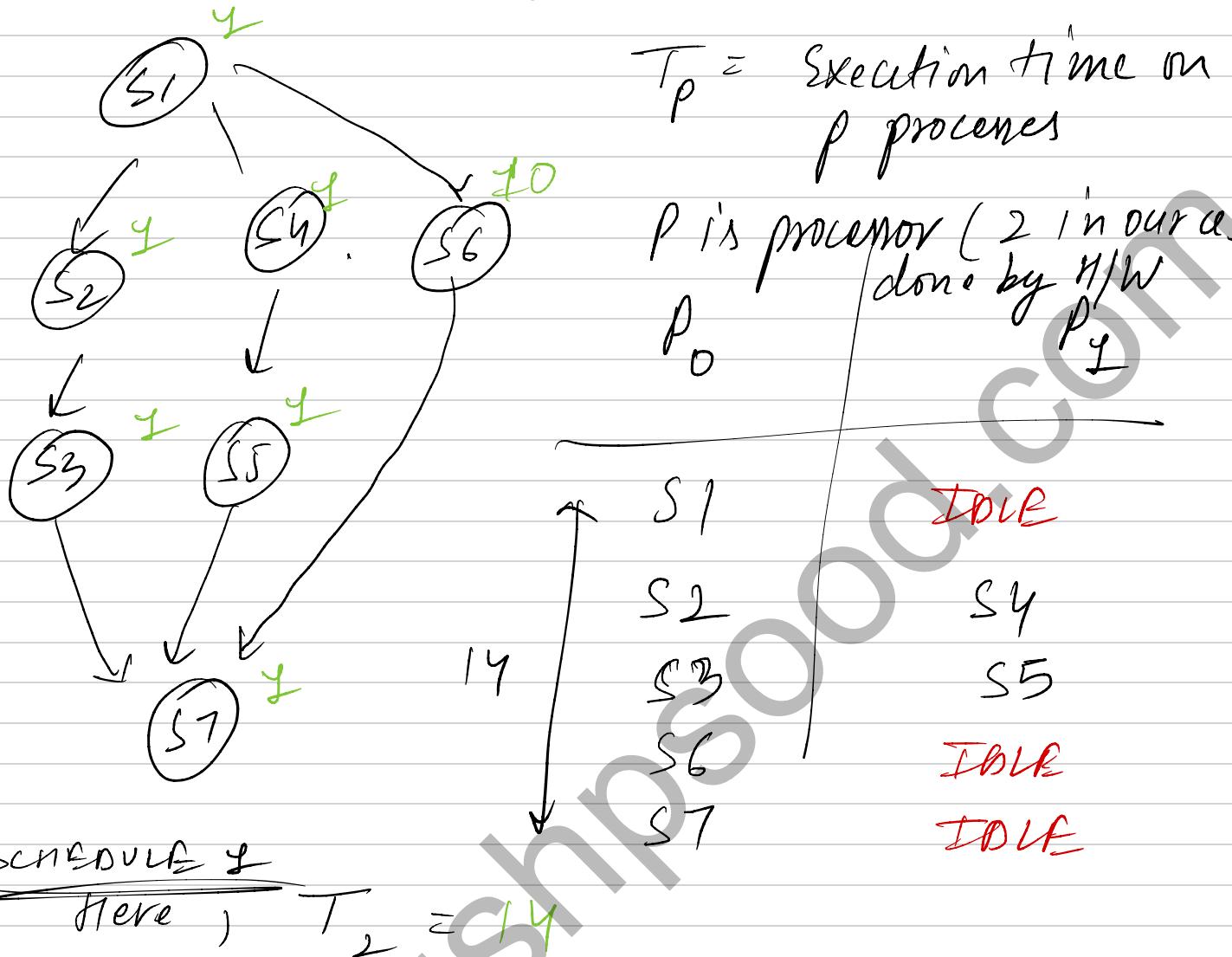
WORK $= 22$
(TOTAL sum of all nodes)

IDEAL
PARALLELISM

SPAN = 12
(sum of longest path)

$$\frac{\text{WORK}}{\text{SPAN}} = \frac{22}{12}$$

Multiprocessor Scheduling, parallel speedup



$$\text{Here } T_2 = 12$$

If $P = 1$, $P^T_1 = \text{WORK}$ $T_1 = 16$
(1 processor)

Here P is ∞

$T_\infty = \text{SPAN}$
(length of longest path) = 12

$$T_\infty \leq T_P \leq T_1$$

SPEEDUP = $\frac{T_1}{T_P}$ (How much faster the parallel version was able to run)

Bound on SPEEDUP

$$\text{SPEEDUP} \leq P$$

$$\text{SPEEDUP} \leq \frac{\text{WORK}}{\text{SPAN}}$$

= IDEAL PARALLELISM

Amdahl's law

(Interseting observation made by Amdahl 50 yrs back)

$$\text{SPEEDUP} \leq \frac{\text{WORK}}{\text{SPAN}}$$

$q = \text{FRACTION SEQUENTIAL}$

$$\text{SPEEDUP} \leq \frac{1}{q}$$

$$q = 0.5 \Rightarrow \text{SPEEDUP} \leq 2$$

No Need for large parallelism

$$q = 0.1 \Rightarrow \text{SPEEDUP} \leq 10$$

PROOF

get WORK must be included in the computational graph (Upper bound)

SPAN > q * WORK

$$\text{SPEEDUP} \leq \frac{\text{WORK}}{\text{SPAN}}$$

$$\text{SPEEDUP} \leq \frac{\text{WORK}}{q * \cancel{\text{WORK}}}$$

futures

$$A = f(B)$$

$$C = g(A)$$

$$D = h(A)$$

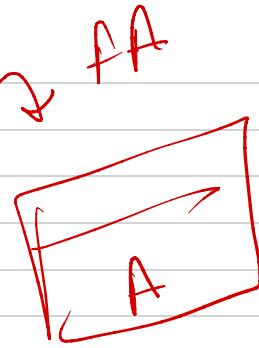
$$FA = \text{FUTURE}(f(B))$$

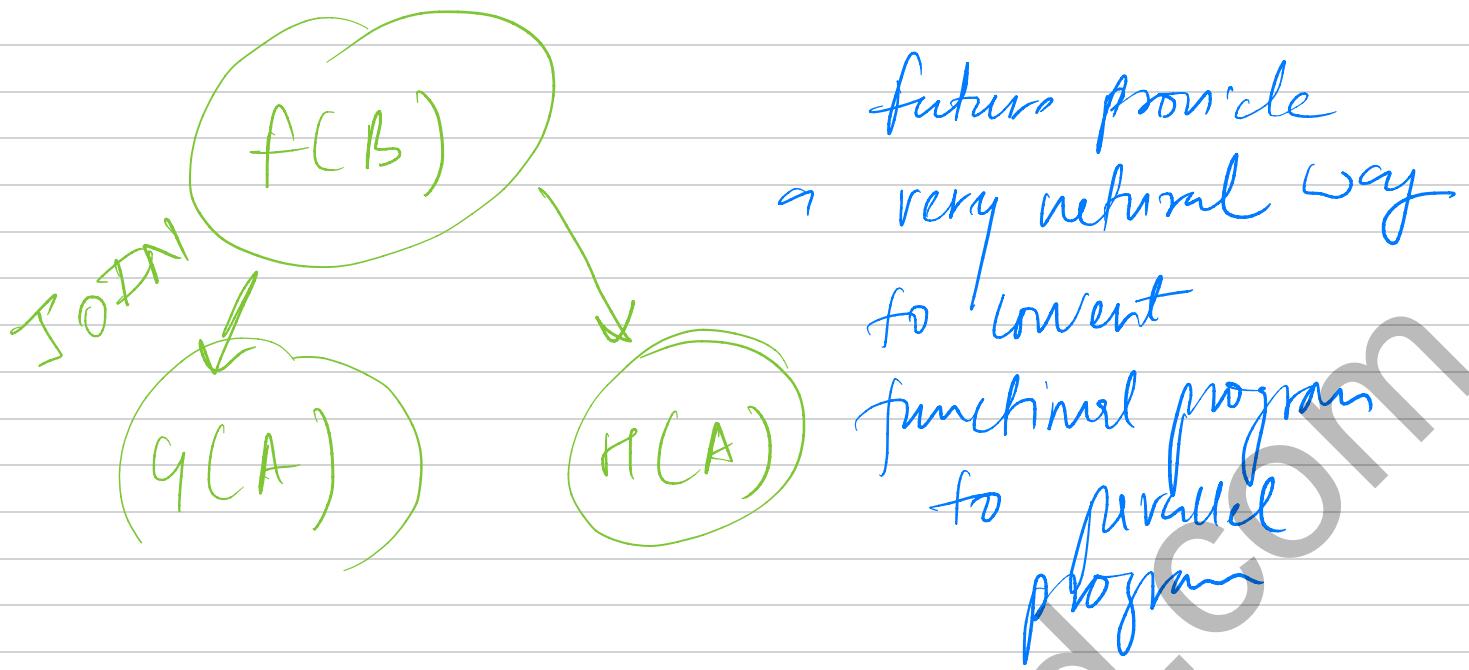
$$FC = \text{FUTURE}(g(f(A).get()))$$

$$FD = \text{FUTURE}$$

$$(h(FA).get())$$

BLOCK
WAIT
until A
is available





$fSUM1 = \text{FUTURE } \{ \text{LOWER } y \}$

$fSUM2 = \text{FUTURE } \{ \text{UPPER } y \}$

$\text{SUM} = fSUM1.get() + fSUM2.get()$

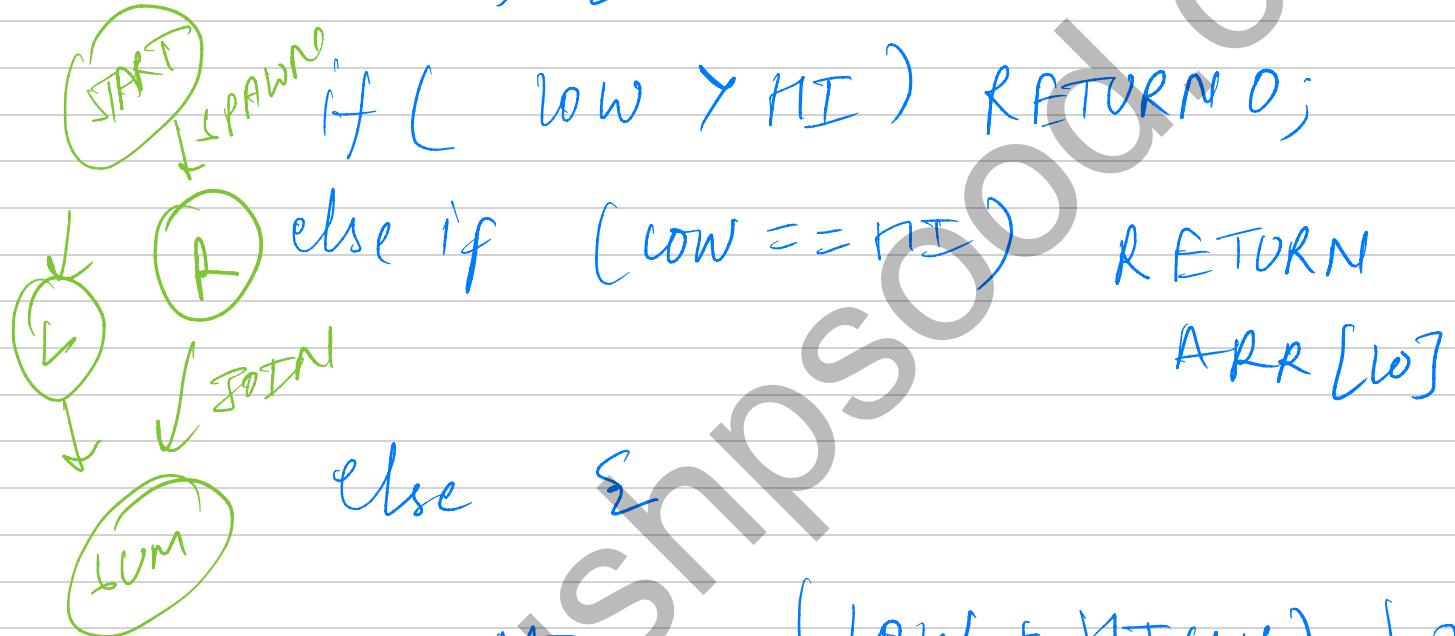
futures in JAVA FORK & JOIN

* NOT RECURSIVE ACTION

CLASS ASUM EXTENDS RECURSIVETASK

FIELDS : ARR, LO, HI

COMPUTE() {



Join for
future implementation

future get
operation

MID = $(\text{low} + \text{hi}) / 2$
 $L = \text{NEW ASUM}(\text{low}, \text{mid},$
 $\text{ARR})$

$R = \text{NEW ASUM}(\text{mid} + 1, \text{hi},$
 $\text{ARR})$

R. FORK(); // FUTURE

RETURN L.COMPUTEC() +
R.JOIN();

Memoization

$$Y_1 = g(x_1)$$

INSERT (g, x_1, y_1)

$$Y_2 = g(x_2)$$

$$Y_3 = \cancel{g(x)} \rightarrow$$

LOOKUP (g, x_1)

FYI = FUTURE

$g(x_1) \rightarrow$

INSERT (g, x_1, f_{Y_1})

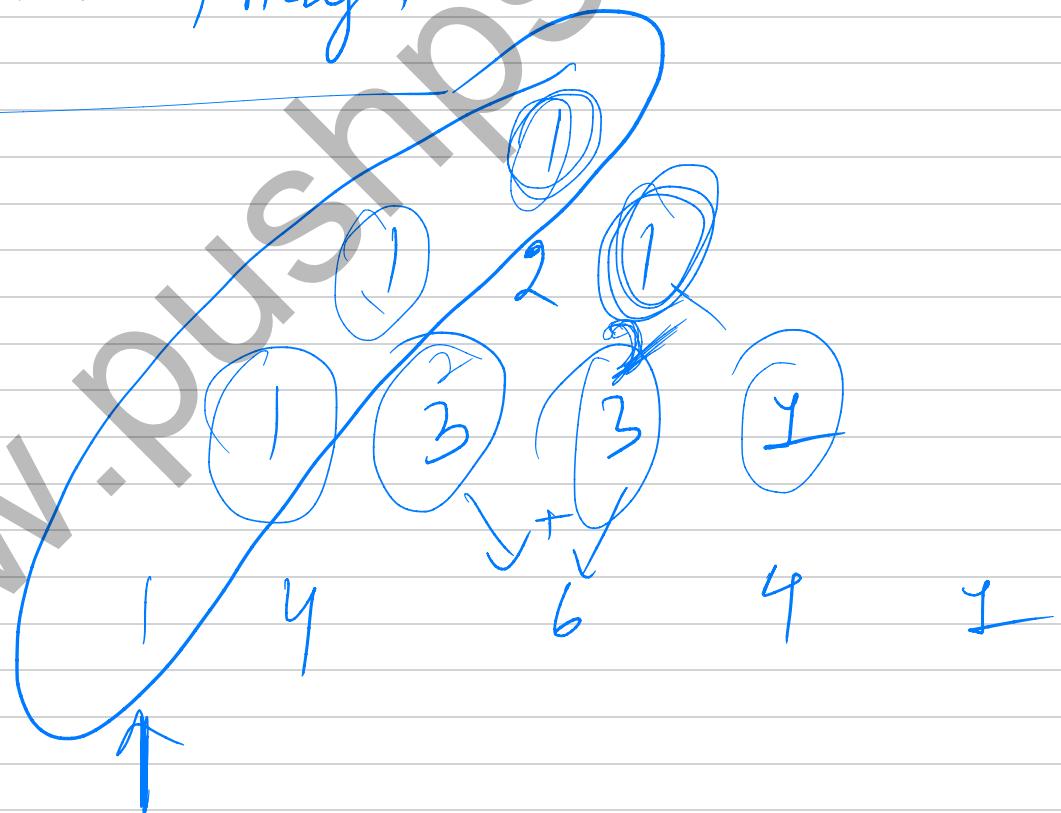
$y_2 = g(x_2)$

wanted work in

sequential

$y_3 = \text{LOOKUP}(g, x_1).$
get()

PASCAL'S Triangle



Java Streams

① `FOR(S : STUDENTS)`

`PRINT S`

`STUDENTS. STREAM() . (s → PRINTS)`

②

`FOR(S : STUDENTS)`

`if (s IS ACTIVE)`

`ACTIVE. APP(s)`

`FOR(A : ACTIVE)`

`AGRSUM+= A. AGE ;`

`Avg = AGRSUM / ACTIV. SIZE()`

USING STREAMS

STUDENTS - STREAM().

FILTER($s \rightarrow s \text{ IS ACTIVE}$).

MAP($s \rightarrow s.\text{AGE}$).

TAKAWAYS

AVERAGE();

- IT'S EASY TO CONVERT SEQUENTIAL STREAM TO PARALLEL
- USE PARALLEL STREAM instead OF STREAM
- AVG & FOR EACH ARE TERMINAL OPERATIONS
- FILTER & MAP ARE INTERMEDIATE OPERATIONS

DATA RACES &

2 TYPES DETERMINISM

FUNCTIONAL

DET

Same input will
always lead to
same output

FUTURE ASYNC S

Read

STRUCTURAL
DETERMINISM

We are interested in parallel programs that are both structurally & functionally plus no delta race but not always!

Benign Non Determinism
(Next Page)

$$\text{SUM 1} = \text{SUM OF LOWER HALF}$$

$$\text{SUM 2} = \text{SUM}$$

OF UPPER HALF.

$$\text{SUM} = \text{SUM 1} + \text{SUM 2}$$

SAME INPUT lead to

SAME COMPUTATIONAL

get()

GRAPH (structure of parallel program)

BIN DAY PARALLEL PROG

(DATA FENCE)
Read & Write

of WRITE & WRITE in parallel.

→ (vice versa if code is having DRF it is f₂S Deterministic)

DRF → Data Race Free

→ (can be achieved when we achieve functional and deterministic)

→ FUNCTIONAL + STRUCTURAL Deterministic

BENIGN NON DETERMINISM

Cases where we get diff output but they are all accepted

eg Search

PARALLEL

LOOP

FUNCTION 2

FOR [i : [0 : M-N])

FOUND = TRUE

ASYNC FOR [j : [0 : N-i]] {

IF [TEXT(i+j)] ≠ PAI[i]

FOUND = FALSE

If (FOUND = TRUE) Index = i;

If multiple both have

same string then

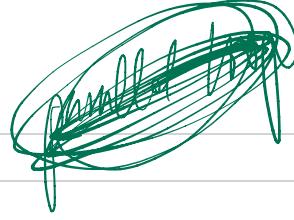
any index is the right

Example:

Value of
index is
initially
not acceptable
but acceptable
now the answer
is?

}

PARALLEL LOOPS



When the tasks set in being performed inside the loop is independent of each other we can freely use parallelism to compute them independently.

for {
 $p = \text{head}$; $p \neq \text{null}$;
 $p = p.\text{Next}$ }

{

$\text{async } p.\text{compute}$

y

Any for all

forall ($i : [0 - n - 1]$)

y

$$a[i] = b[i] + c[i]$$

This will be
converted to
a parallel form

A better way is by using the JAVA API's parallel stream

$q = IntStream.rangeClosed(0, n-1)$
parallel.

- to Array ($i \rightarrow b[i] + c[i]$),

Matrix multiplication This is used to measure speed of fast computation in the world (benchmarking)

$$C \leftarrow A \times B$$

$$C[i][j] = \sum_{k=0}^{n-1} A[i][k] * B[k][j]$$

PORALL \rightarrow parallelism

~~FOR([i], [j]) : [0:N-1, 0:N-1]~~

Can we run this in parallel

$C[i][j] = 0$ No! This will lead to data race

It will remove ~~FOR~~

$C[i][j] += A[i][k] *$
 $B[k][j];$

y

y

Barriers in parallel loop

OUTPUT

This order will be retained

FORALL ($I : [0 : N-1]$)

Hellow TWO
Hellow zero
BYE TWO

$S = \text{LOOKUP}(I);$

FRONT = "Hellow" + S

BELOW ON

FORALL? BARRIER

? FRONT = "BYE" + S

Different combination

choice 1
choice 2
choice 3
why not
choice 2
why so?

Can we print all
hallows before the
good byes?

Output with barrier
Hellow 2
n 0 D

n 1

< BARRIER >

good bye 1

n 2

n 0

We have to
repeat the
look up
if we have
more than
one parallel
loop
which can be
very
computationally
expensive

$$x_0 = 0 \quad x_N = 2$$

$$x_i^0 = (x_{i-1} + x_{i+1}) / 2$$

$$0 < i < N$$

$$N = 10$$

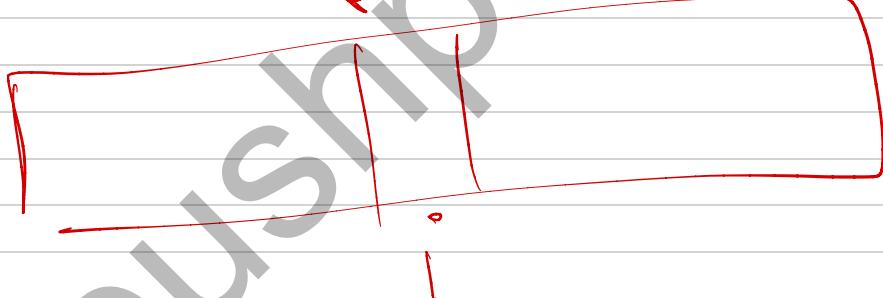
0

0	0.1	0.2	0.8	0.9
---	-----	-----	-------	-----	-----

old



New



ONE DIMENSIONAL ITERATIVE AVERAGING
APPROACH

sequential
FOR [ITER : [i : N STEPS]) {
 <sup>→ No of time we want
 to repeat</sup>



using barriers
FORALL (I : [i : N-1]) {
 ^{→ (N-1) TASKS}

 FORALL (I : [i : N-1]) {

$$\text{NEWX}[i] = \text{AVG} (\text{OLDX}[i-1], \\ \text{OLDX}[i+1])$$

 SWAP POINTERS

^{→ (N-1) TASKS}

^{(N-1) X N STEPS}

 FORALL (I : [i : N-1]) {
 ^{→ (N-1) TASKS}

 FOR [ITER : [i : N STEPS]) {

$$\text{NEWX}[i] = \text{AVG} (\text{OLDX}[i-1],$$

^{→ OLDX[i+1]}

 SWAP POINTERS

[→]

FOR ALL $(I : [0 : N-1])$

Too Many !
N TASKS

$$A[I] = B[I] + C[I]$$

Solution

FOR ALL $(g : [0 : Ng-1])$

No of groups
Total task will be Ng

for $i :$

How to determine?

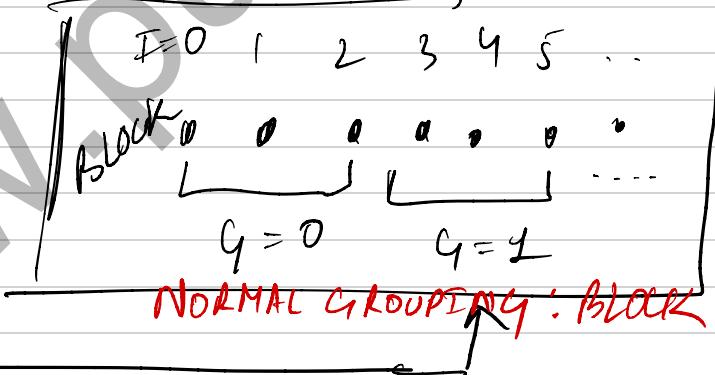
My group (g)

Not always

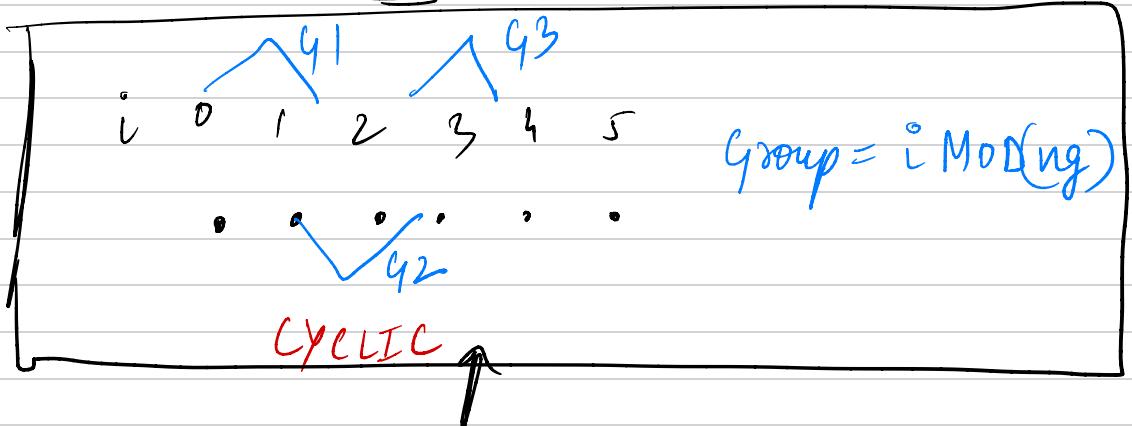
$Ng, [0 : N-1]$ create these many tasks

$$A[I] = B[I] + C[I]$$

CHUNKS OF EQUAL SIZE



fix to too many task
Is chunking vs iteration
grouping



$$\text{Group} = i \bmod(g)$$

Split phase barriers with JAVA PHASERS

CPLF SPAN = 200

FORALL { i : [1:N] } S ≈ 100 unit of time

Do local work

print "HELLO" + i;
ARRIVE
My id = LOOKUP(i);
AWAIT ADVANCE
print "GOODBYE" + NEXT + myid;

g

W/ PHASER CPLF = 200
SPAN

- ARRIVE AND AWAIT ADVANCE

it is actually a barrier

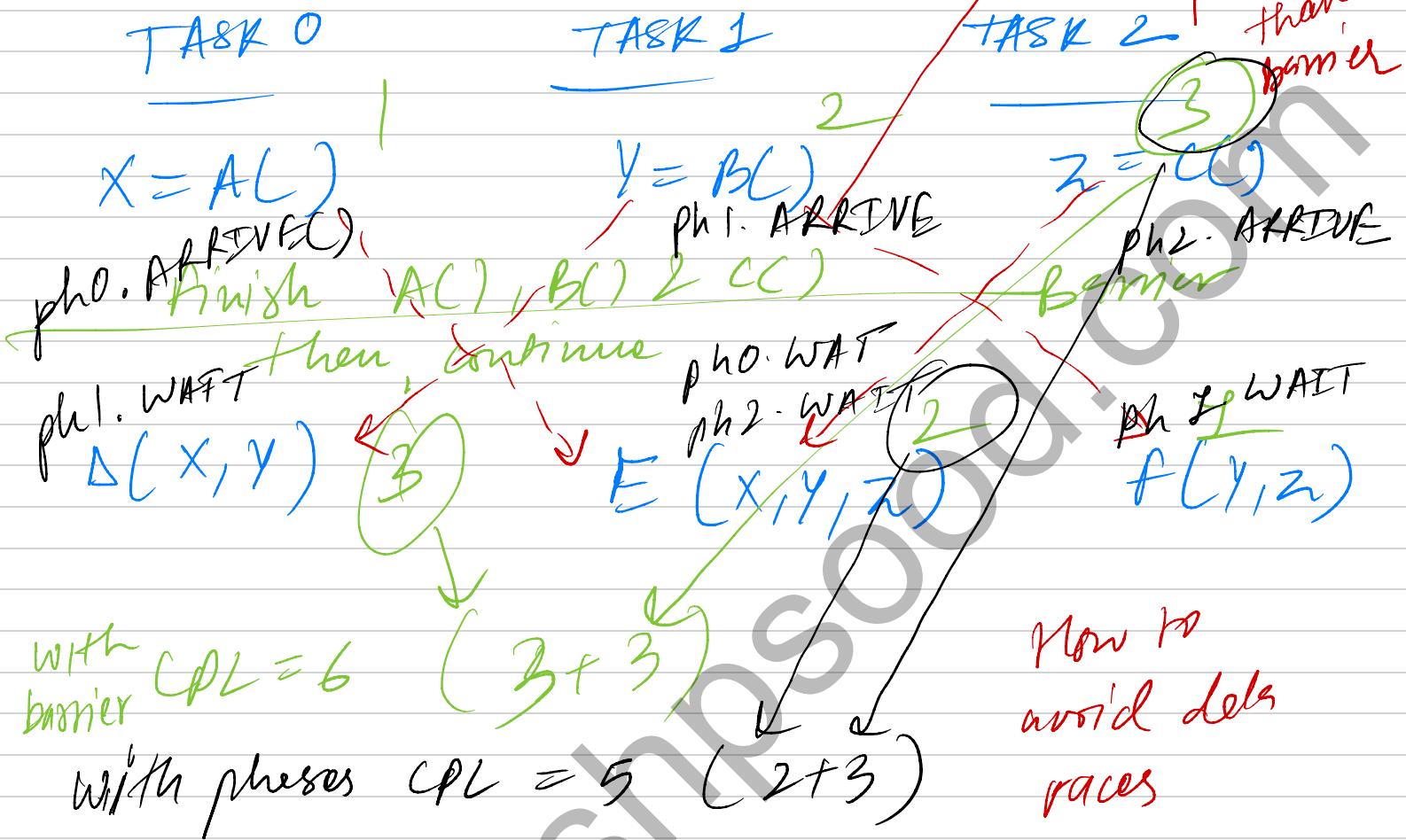
- ARRIVE

LOCAL WORK

- AWAIT ADVANCE

Not associated
or can be executed

POINT - TO - POINT SYNC



1-D Iterative Averaging with Processors

$\text{ph}[i] = \text{new_phaver}[N+2];$
FORALL ($i : [0 : N-1]$) {

FOR (ITER: $[0 : \text{N_STEPS}-1]$) {

$\text{NEWX}[i] = \text{AVG}[\text{OLDX}[i-1]]$

$\text{ph}[i].\text{ARRIVE}();$

$\text{ph}[i-1].\text{WAIT}([\text{OLDX}[i+1]]);$

$\text{ph}[i+1].\text{WAIT}()$ Barrier

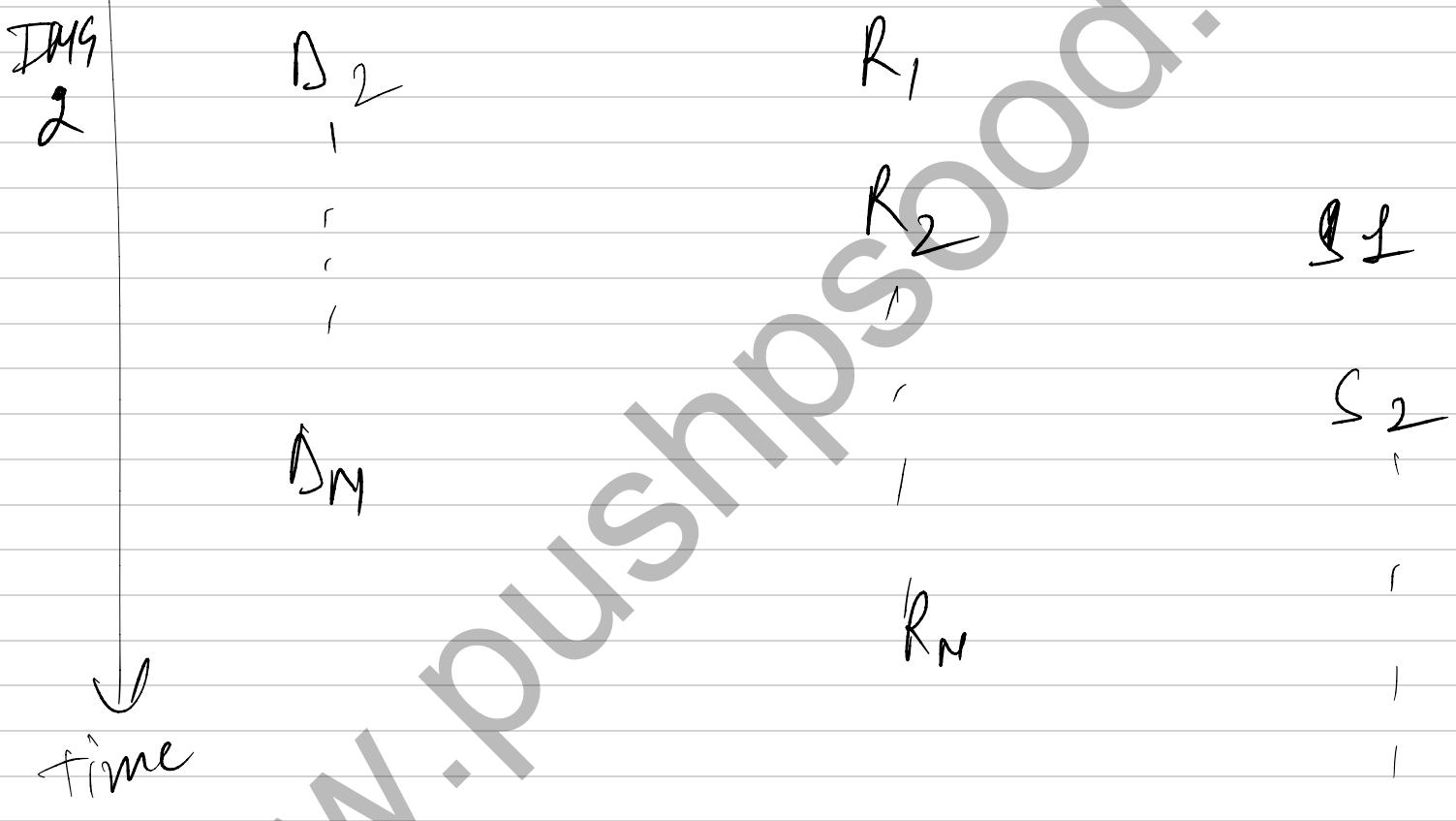
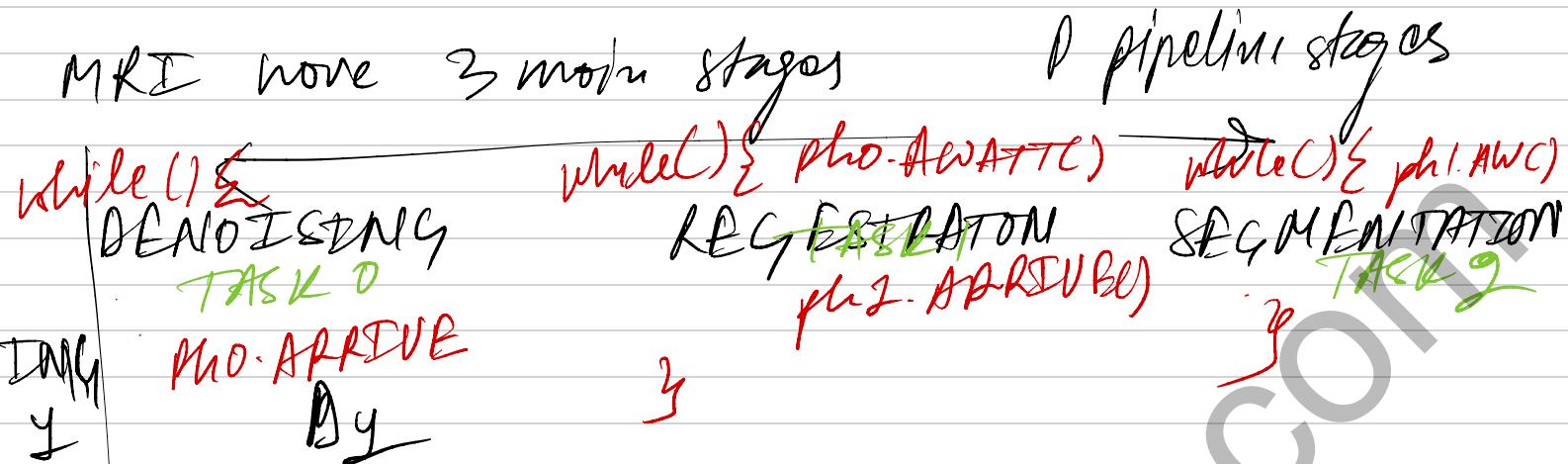
fixed
earlier

3

SWAP

NEWX, OLDX

Pipelining

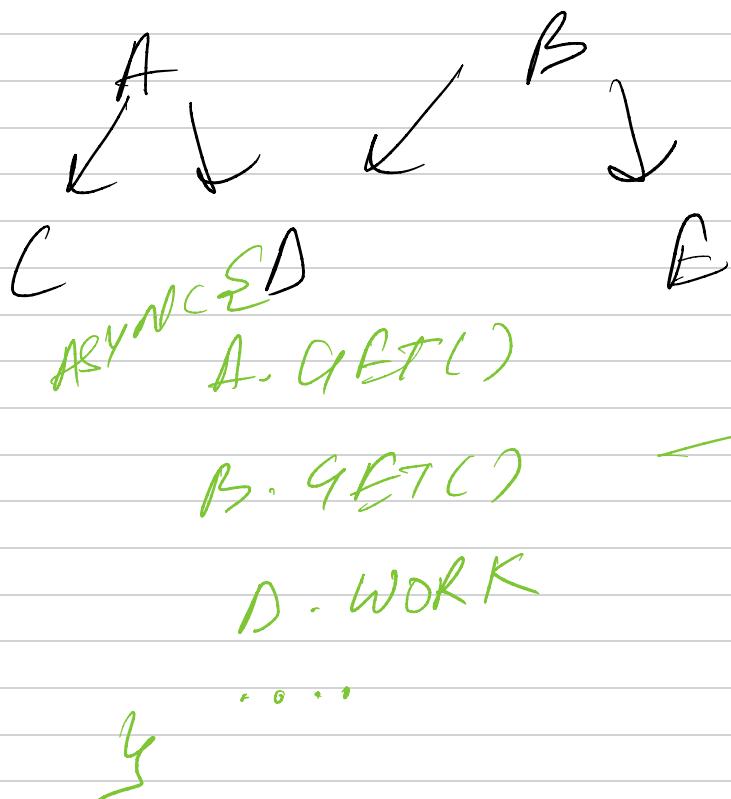


$$WORK = N + P$$

$$CL = N + P - 1$$

$$\text{PARALLELISM} = \frac{N+P}{(N+P-1)} \approx P$$

DATA Flow



possibility 1

works but implicit dependency!

ASYNC { A. WORK ; A.PUT() }
ASYNC { b. WORK ; b.PUT() } explicit

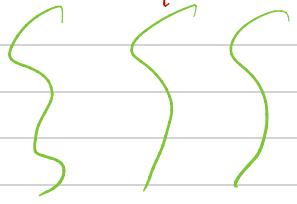
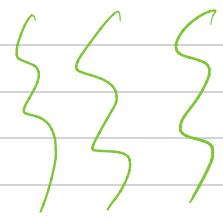
ASYNC AWAIT(A) {
 C. WORK
 precondition

ASYNC AWAIT(A, B) {
 D. WORK
 postcondition

ASYNC AWAIT(B) {
 E. WORK
 postcondition

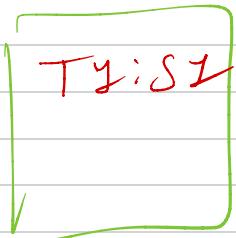
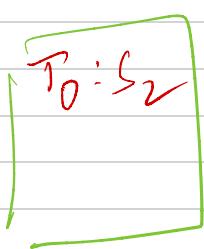
possibility 2

Threads



→ Thread
* Main program is itself a thread

OS will schedule the threads for processing on processor



P₀

P₁

Here, P₁ is processor

MAIN THREAD: T₀

CREATE: t₁ = newThread(S₁) $x = 0$
 (S₁ uses the var x)

t₂ = newThread(S₂)

START = t₁.START()

t₂.START()
S₂:

JOIN(WAIT = t₁.JOIN());

because of this

(S₃ uses the var x)

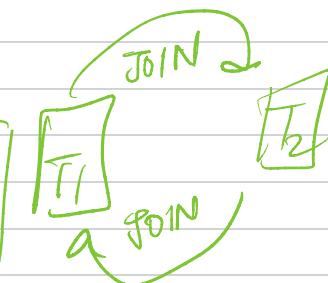
S₃ || Requires x

if S₁ is performing t₂.JOIN()

leads to Dead lock

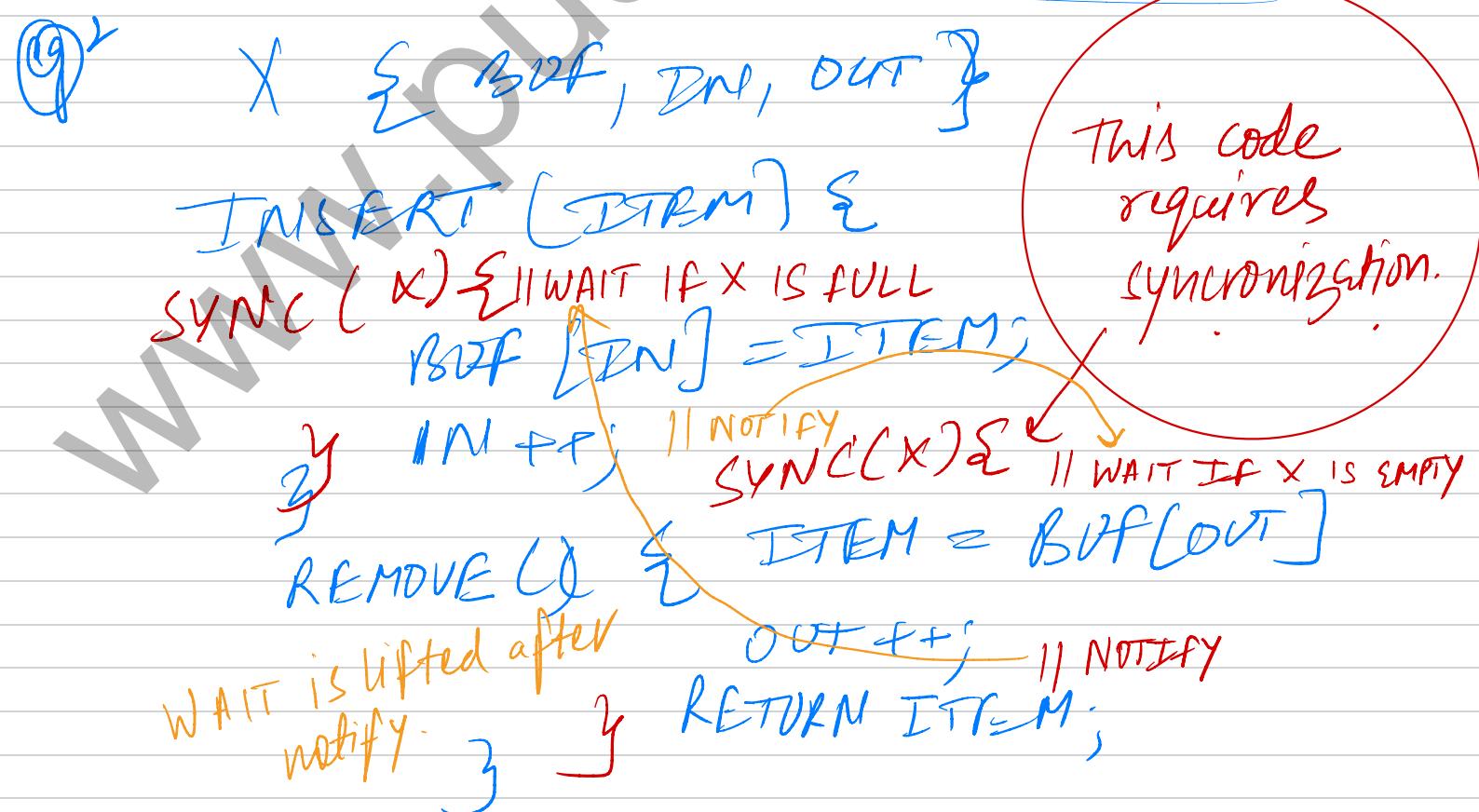
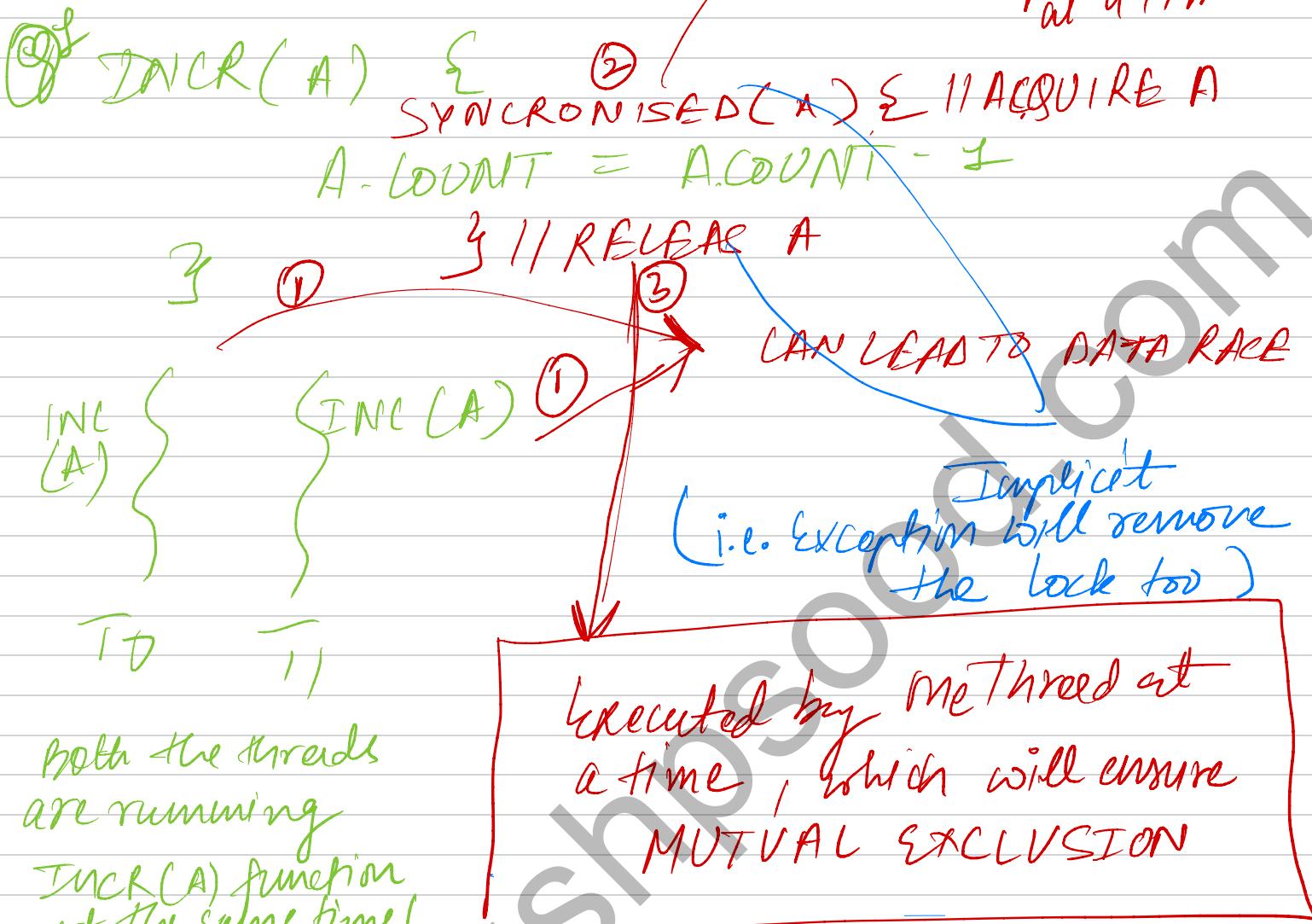
and if S₂ is performing t₁.JOIN()

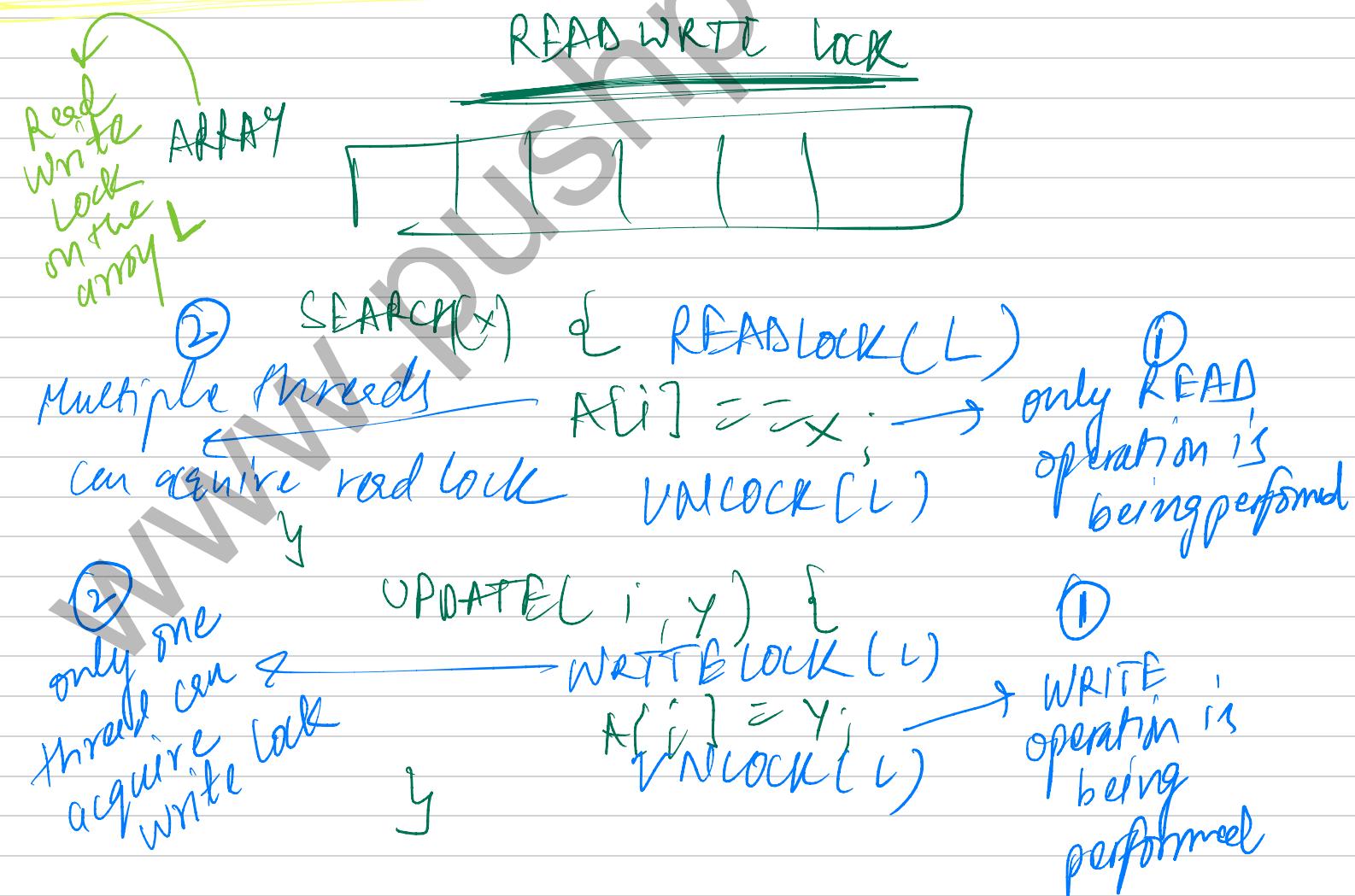
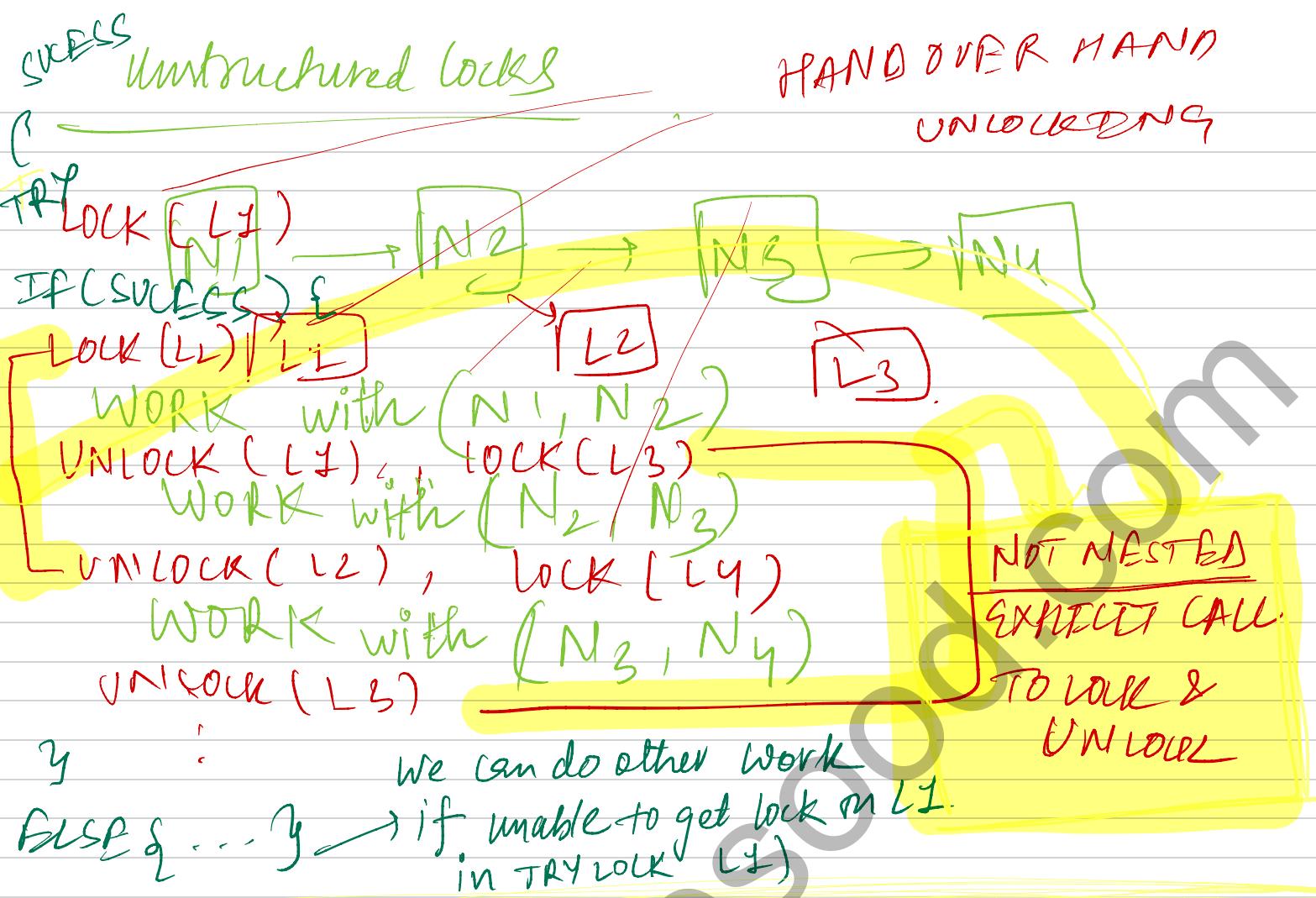
we should use threads cautiously and avoid such deadlock condition



STRUCTURED LOCKS

Only one resource can acquire lock at a time.





LIVENESS → fundamental correctness

property of concurrent programming

DEADLOCK ↓

LIVENESS COMPROMISED

① DEADLOCK

T1

SYNC(A){

 SYNC(B)

y

T2

SYNC(B){

 SYNC(A)

→ Here they are not blocked but they are not making any progress

② LIVELock

MUTABLE INT

T1

DO {

 SYNC(R) {
 X-INCR()

 R = X.GET()

DO {

 SYNC(R) {
 X-DECR()

Even adding synchronise wont fix the issue and lead to live lock loop (infinite loop)

Y Y

WHILE(R < 2)

Diagram

 R = X.GET()

Y Y
 R = X.GET()

WHILE(R > 2)

③

STARVATION

$S_1 \dots S_{100}$

T_1

do {

$IN[1] = S_1.\text{READLINE}()$

PRINT IN[1]

}
while (...)

T_2

:

T_W

do {

STARVES AS THE
OTHER THREADS ARE
TAKING UP ALL THE
THREADS

$IN[1] = S_{100}.\text{READLINE}()$

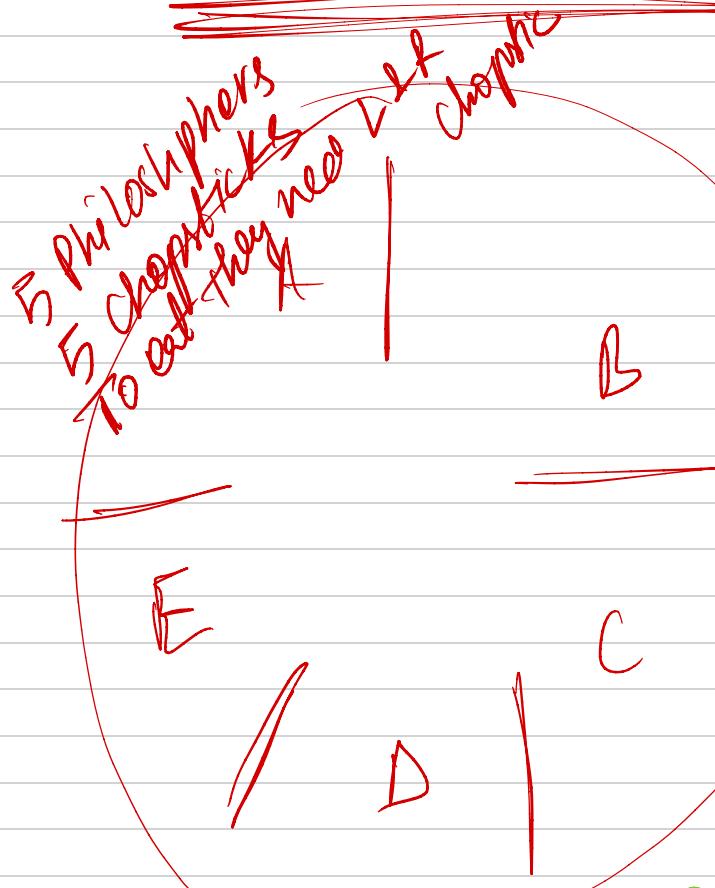
PRINT IN[1]

}

while (...)

DINING PHILOSOPHERS

(LIVELOCK &
DEADLOCK CLASSICAL
EXAMPLE)



(DEADLOCK, NO LIVELOCK)
STRUCTURED LOCKS

WHILE (...) {

THINK

SYNC (LEFT) {

SYNC (RIGHT) {

EAT

y

yy

- THINK
- PICK UP CHOPSTICK
- EAT
- PUT DOWN CHOPSTICK
repeat

code

Never able to eat

(NO deadlock, but livelock)
UNSTRUCTURED Locks

WHILE (...) {

THINK

S1 = TRYLOCK (LEFT)

if (! S1) CONTINUE;

S2 = TRYLOCK (RIGHT)

if (! S2) {
UNLOCK (LEFT)
CONTINUE
y y

POSSIBLE DEADLOCK IN STRUCTURED LOCK

All philosopher have ranking and pick the chopstick in their LEFT and keep waiting for the right one, which will never happen

DEADLOCK

UNSTRUCTURE LOCKS (SAME EXAMPLE IN PREVIOUS PAGE)

while (. . .) {

 THINK

$S_1 = \text{TRYLOCK(LEFT)}$;

 if ($|S_1|$) CONTINUE;

$S_2 = \text{TRYLOCK(RIGHT)}$

 if ($|S_2|$) {

UNLOCK(LEFT) ;

 CONTINUE;

 }

MODIFICATION

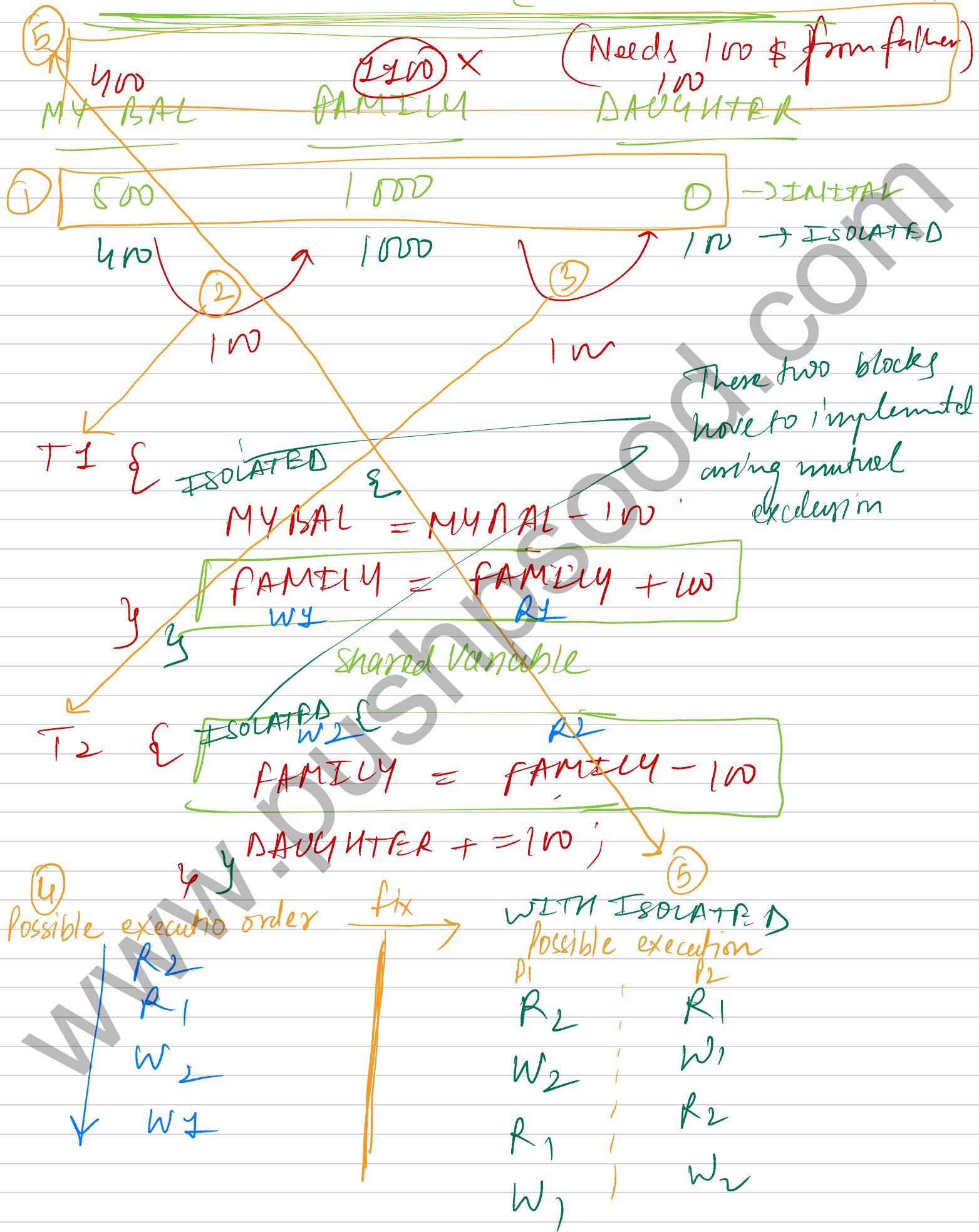
1 philosopher → PICK RIGHT

1 philosopher → PICK LEFT

Change order
of one
philosopher

CAN LEAD TO STARVATION which is known issue,
but this can be fixed using SEMAPHORES

Critical Section (Also known as ISOLATE construct)



Object based Isolation (MONITORS)



DELETE (CUR)

ISOLATED CUR, CUR.PREV, CUR.NEXT
ISOLATED E → WORKS! But
this becomes sequential

$$\text{CUR.PREV.NEXT} = \text{CUR.NEXT}$$

$$\text{CUR.NEXT.PREV} = \text{CUR.PREV};$$

y

can lead to wrong answer

Execution of Delete function without ISOLATED

T₁ : DELETE(B) → (A, C)

(A, B, C)

T₂ : D(C) → (B, D)

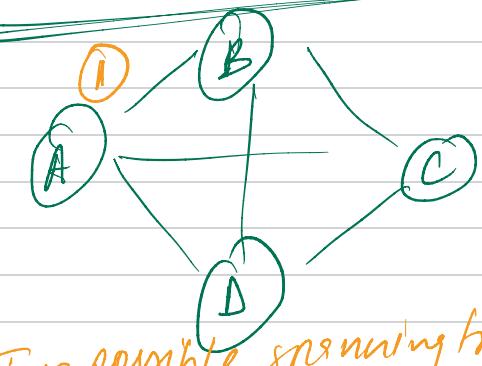
can be run parallelly

T₃ : D(E) → (D, F)

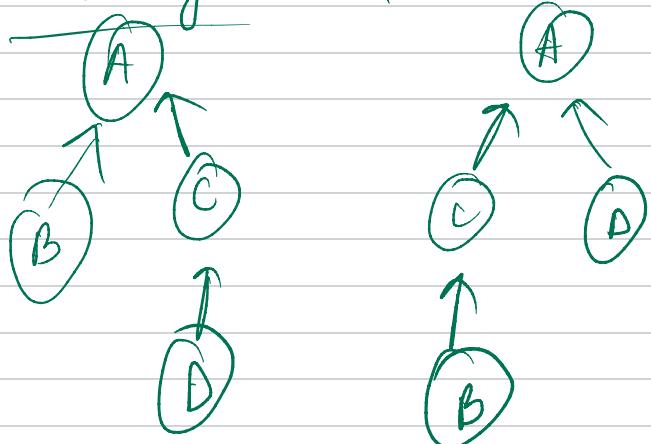
(D, E, F)

No objects in common so they will run parallel

SPANNING TREE OF AN UNDIRECTED GRAPH



② Two possible spanning tree
Possibility 1 Possibility 2



③ COMPUTE(v) \leftarrow boolean success flag

FOR EACH NEIGHBOR C OF v
CAN LEAD TO Data Race

```

    S ← MAKE PARENT(v, c)
    if (S)
        can be
        computed in
        COMPUTE(c) a new
        thread.
        i.e. parallel
    }
  
```

MAKE PARENT(v, c) {

ISOLATED { if (c.parent == null) {

c.parent = v;

s = TRUE

else {

s = FALSE

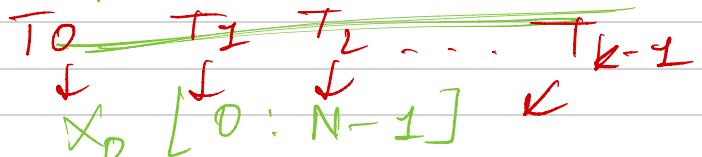
3g
RETURN s;

This is a
famous pattern
compare & set

Here multiple threads
can be trying to make
parent

Atomic Var & References are higher level than lower level locks

ATOMIC VARIABLES



Very standard problem

```
DO {  
    ISOLATED(CUR) {  
        J = CUR;  
        CUR = CUR + 1;  
        IF (J > N)  
            BREAK;  
    }  
    PROCESS ITEM(X(J))  
    WRITE(TRUE)  
}
```

CUR is ATOMIC

Replace with atomic integer

INTEGER Now

special HW can handle this really well

CUR.GET AND ADD(1)

Other use cases:

ATOMIC REFERENCES

COMPARE AND SET (Expected Val, New Val)

ISOLATED(THIS) {

IF (THIS.VAL == ^{expected} Val)
 THIS.VAL = newVal

RETURN TRUE

ELSE { RETURN FALSE }

y y

Refinement of object based isolation

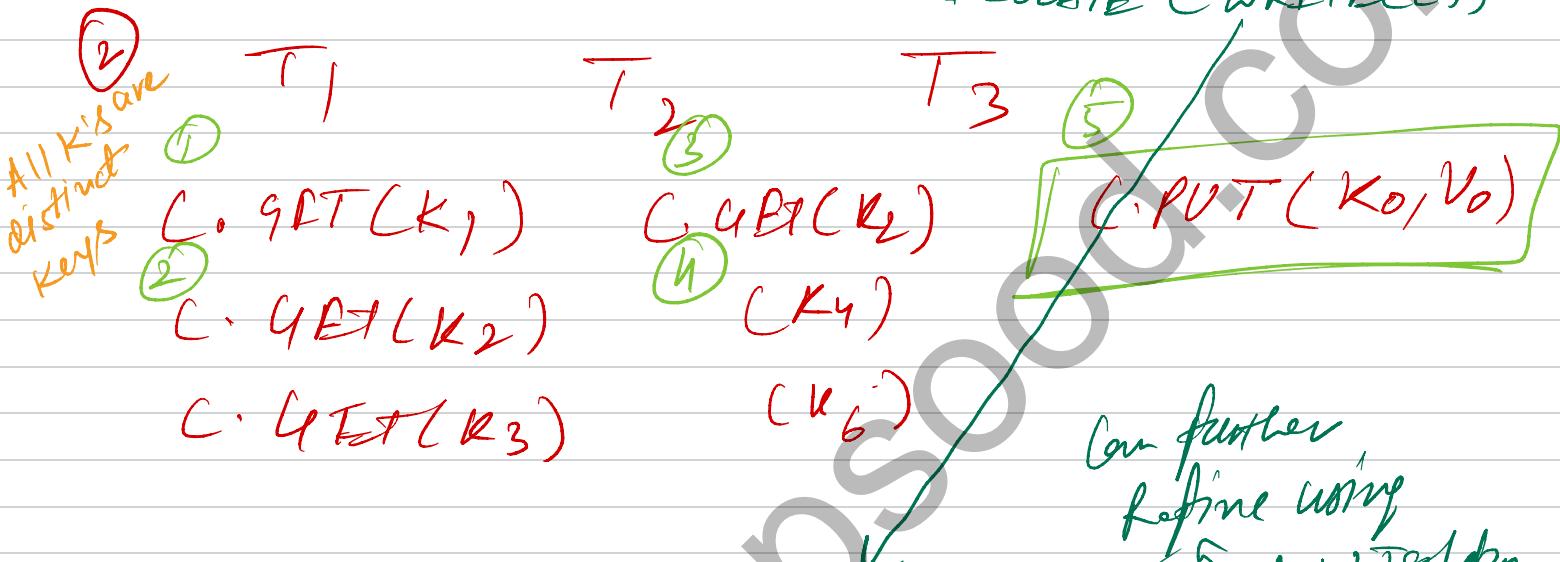
READ-WRITE ISOLATION

① Collection C

$V_1 \leftarrow C.\text{Get}(K_1)$ ISOLATED (Option 1)
 $C.\text{PUT}(K_2, V_2)$ ISOLATED (Option 2)

Fixed the issue
but gets can run concurrently

Concurreny
ISOLATE (WRITELC))



Can further refine using R-W Isolation

T1, T2, T3, T4 are parallel

5 will make sure it's isolated

And then other gets 4 will get executed

DELETE(CUR) { }

writeMode

Access to CUR is in read mode

ISOLATED(CUR.NEXT, CUR.PREV, CUR.F)
 $\text{CUR}.\text{PREV}.\text{NEXT} = \text{CUR}.\text{NEXT}$

$$\text{CUR}.\text{NEXT}.\text{PREV} = \text{CUR}.\text{PREV};$$

y

y

ACTORS

(Even higher level of concurrency)

Implementation of atomic int get&Add to show issue

GETANDADD(int DELTA) {

J = CUR

CUR += DELTA

RETURN J;

isolated

with CC &
Isolation!

ensures
concurrency

if we
forgot to
put isolation
in CUR in
a function
leads to bugs.

fro C {

CUR = ...

if isolation is
mixed this can
put bugs, how to
put this at a global level

GETANDADD()

S

①

every request is sent to
the private thread
but below

②

CUR

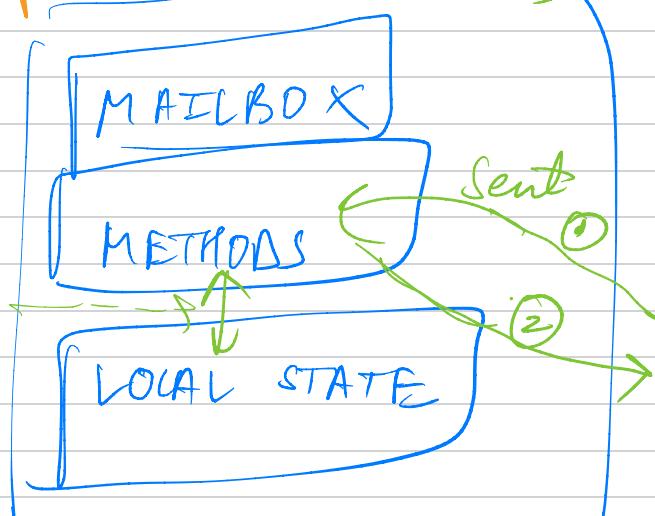
GAA()

③

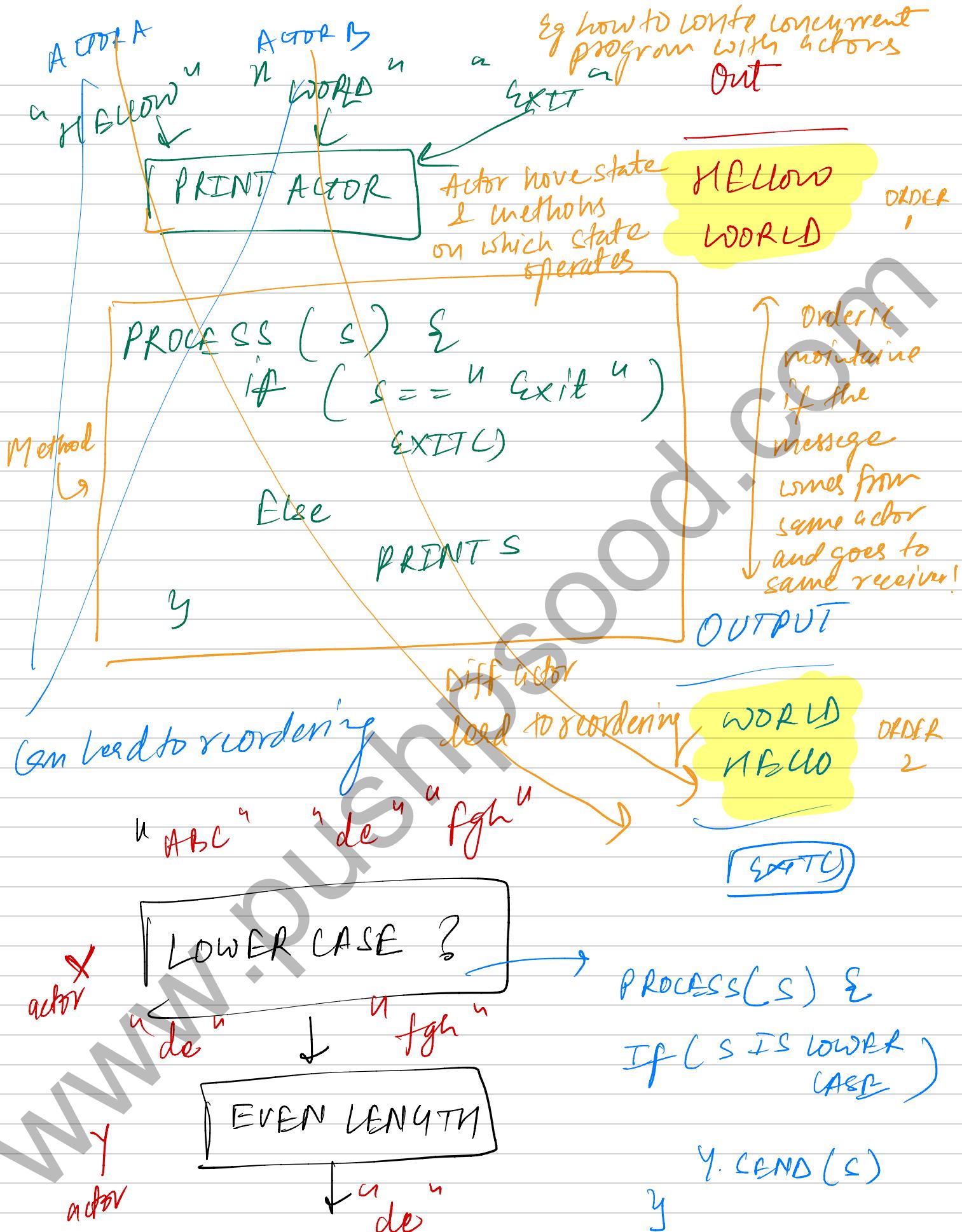
6
No way any other
thread can access
the val directly if
can interact
with owner of CUR
value throw sending
messages

updates

ACTOR (Reactive)



This combination
is called actor

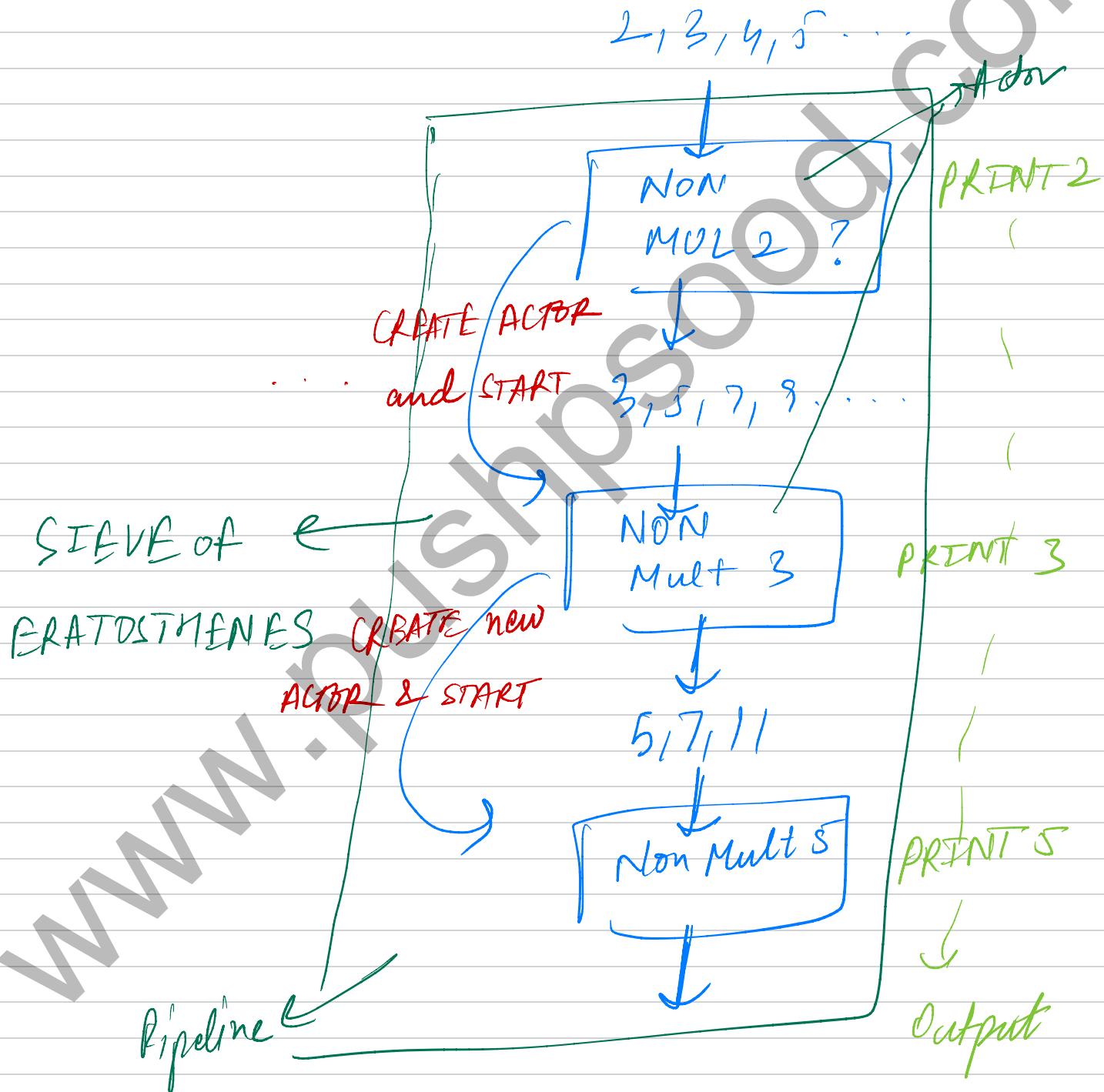


Generating prime numbers $\leq n$ using Actors

Normal Version:

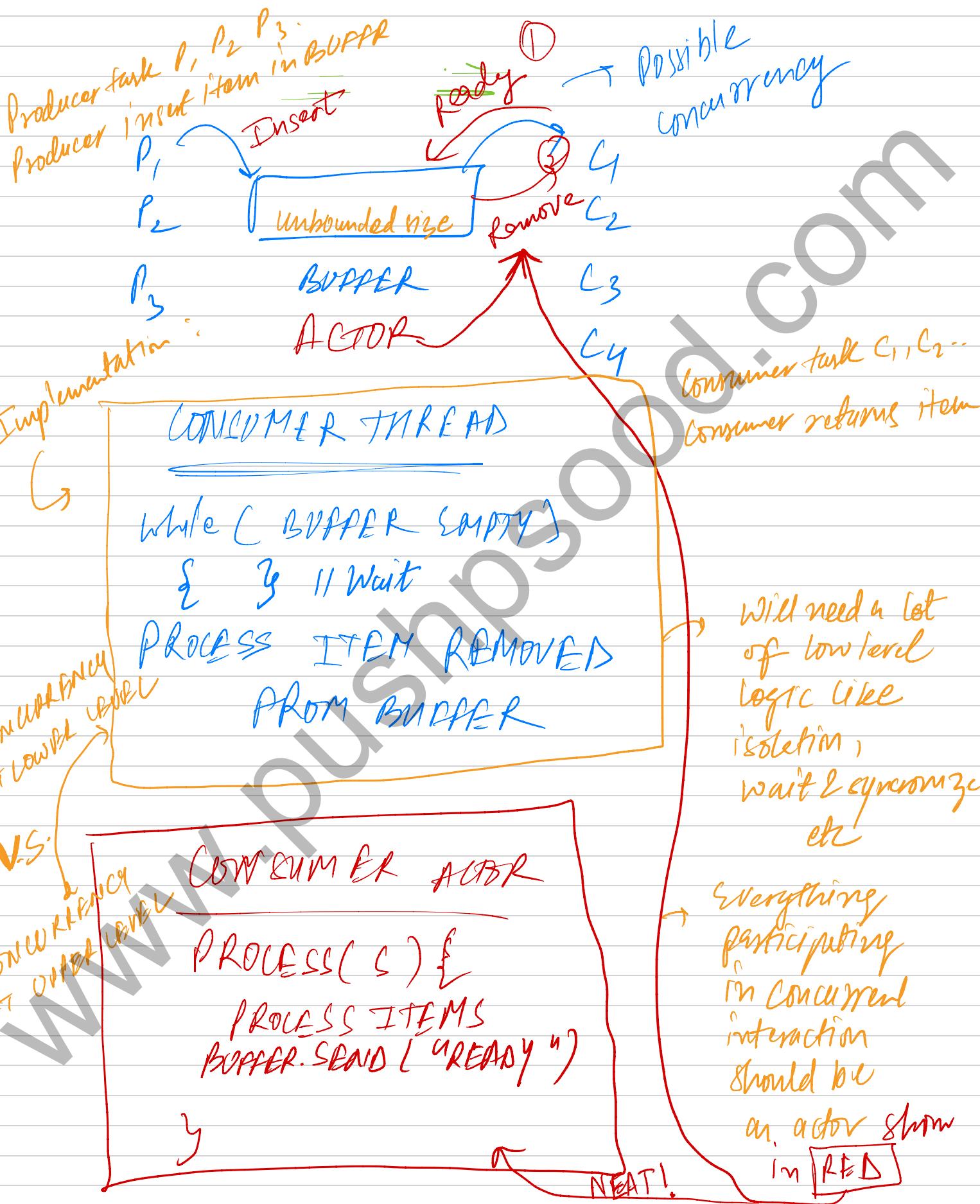


Optimized Greek algorithm : SIEVE OF ERATOSTHENES



PRODUCER CONSUMER PROBLEM

(Classic concurrent programming problem)



BOUNDED BUFFER PRODUCER CONSUMER

(Realistic version of producer consumer example)

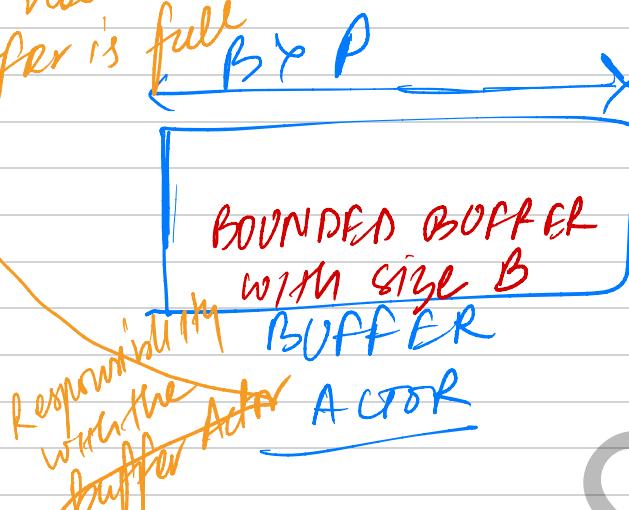
Main challenge with bounded buffer is
producer should not add item to
buffer if buffer is full

PRODUCER

ACTORS

CONSUMERS

ACTORS.



P₁

P₂

:

:

:

REQUEST

DATA

SEND DATA

INSERT)

READY

G

C₂

:

:

:

REMOVE

Though this looks so simple this
problem has been studied for decades
in advance OS classes and there
we use primitives such as
semaphores

OPTIMISTIC CONCURRENCY

Normal Implementation

INTEGER Σ

get()
set()

getAndAdd(delta) Σ

```

CUR = this.get();
NEXT = CUR + delta
this.set(NEXT);
return CUR;

```

Advance concept
to create concurrent
ds by experts

Multiple threads
calling get and
add will lead
to data race.

Expert Implementation

ATOMIC INTEGER Σ

GET()
ATOMIC INTEGER Σ
SET()
GET()

COMPARABLE AND SET()

COMPARE AND ADD(delta) Σ

GET AND ADD(delta) Σ
CUR = this.get();

WHILE(TRUE) {
 CUR = this.get();

CUR = CUR + delta;

IF (this.compareAndSet(CUR, CUR + delta)) {

IF (this.compareAndSet(CUR, CUR + delta)) {

RETURN CUR;

① If multiple threads
were updating the value
of another thread can update the
integer value in between
then the cur value will be stale

② Returns
a boolean value
to indicate
success

③ Each thread
will repeatedly
run this loop till
successful

The code ensures no DEADLOCK
and LIVELOCK

CONCURRENT QUEUE

(Popular concurrent OS)

- flow block first
- IN BETWEEN

Implementation:

Queue loop helps in concurrency

Multiple thread access

ENQ (x) {

 TAIL.NEXT = x; SSS

 TAIL = x; SSS

DEQUEUE() {

 if (EMPTY) ...

 R = HEAD; SSS

 HEAD = HEAD.NEXT; SSS

 Return R

ENQ (x) {

 WNLDE(true) {

 TF (!TAIL.NEXT.COMPARISON AND

 SET (NULL, x))

 } CONTINUE;

Only succeed when it has correct
val of TAIL.NEXT

①

In concurrent prog
every shared var is potential hazard

SSS
multiple set
can be enqing

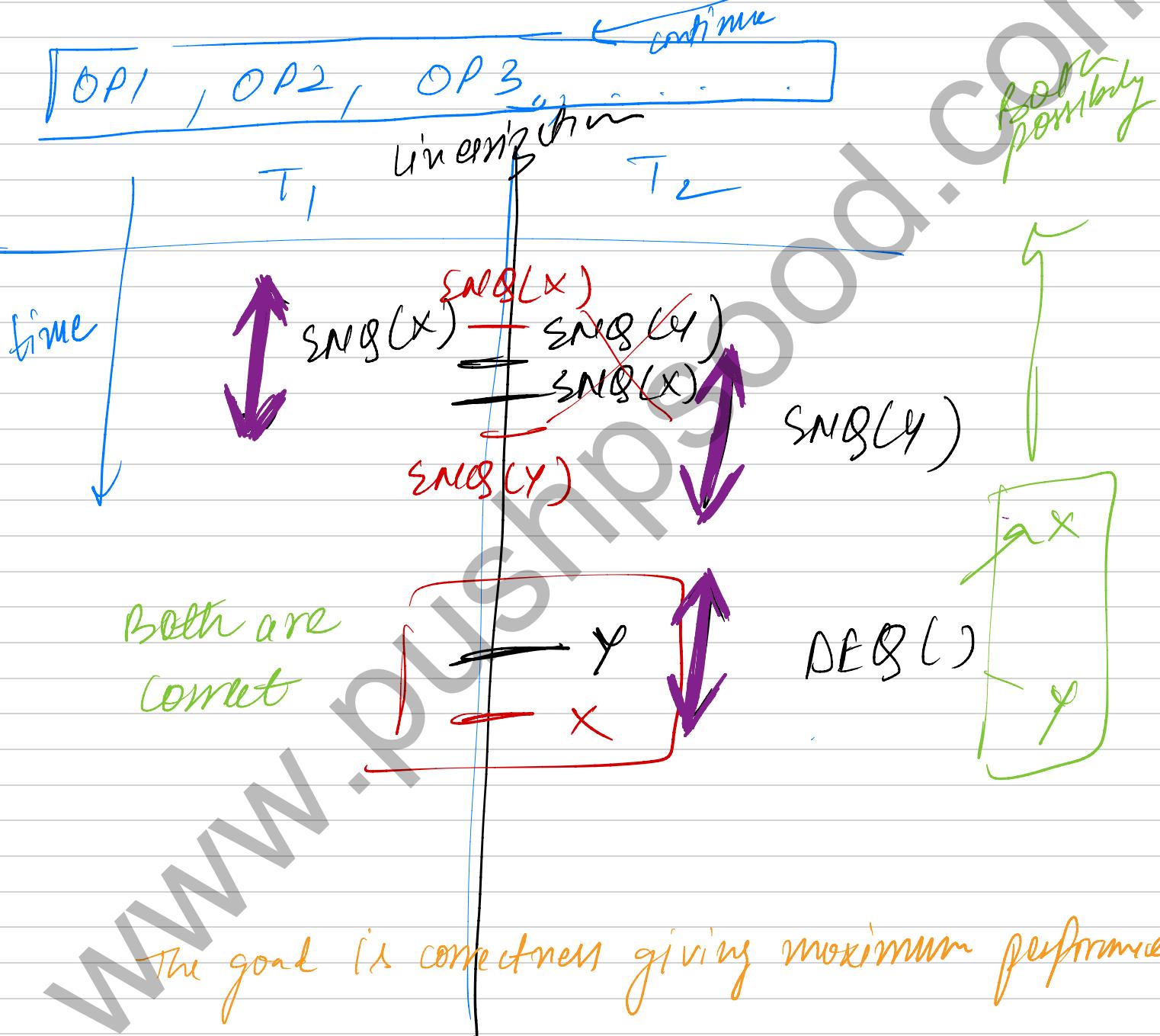
②

③

usual way of
dealing with
concurrency is
locks & isolation
but optimistic
concurrency allows
more concurrency

LINEARIZABILITY (Important property of concurrent object)

Concurrent object is an object which is made "thread safe", i.e. safe for multiple threads to operate on it e.g. Atomic integer, concurrent queue, concurrent hashmap etc. It has some sets of operations $\{OP_1, OP_2, OP_3\}$



The goal is correctness giving maximum performance

CONCURRENCE HASHMAP (one of many implementation of Java Concurrent DS)

CON HASHMAP

CON LINKED LIST

CON SIGHT LIST SET

Implements operations of normal map with some additional functions!

linearizable

GET (key)

PUT (key, value)

PUTIFABSENT (k, v)

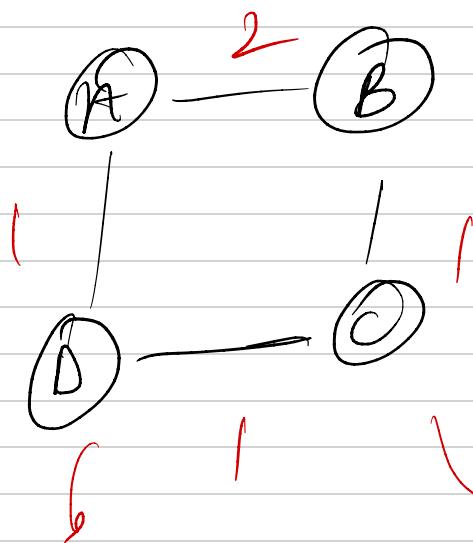
Clear()

PutAll()

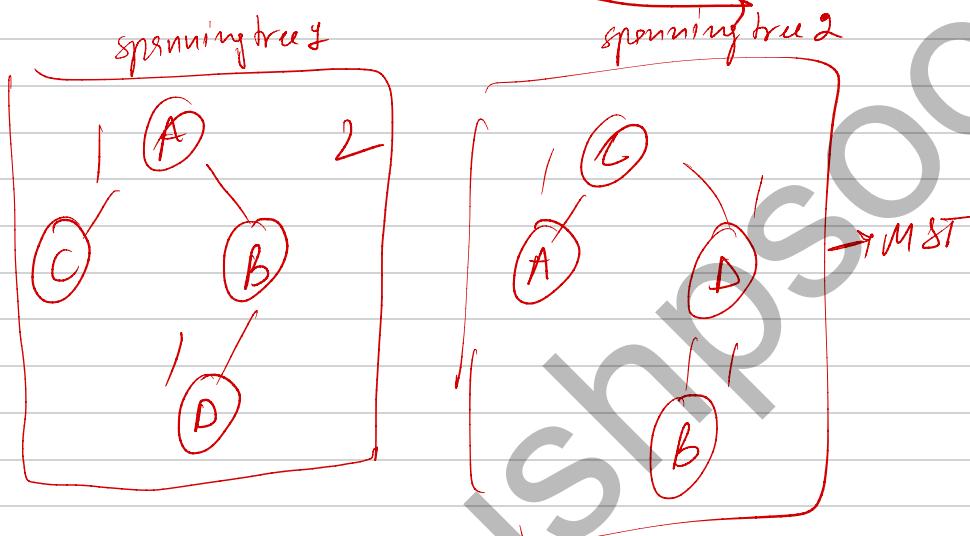
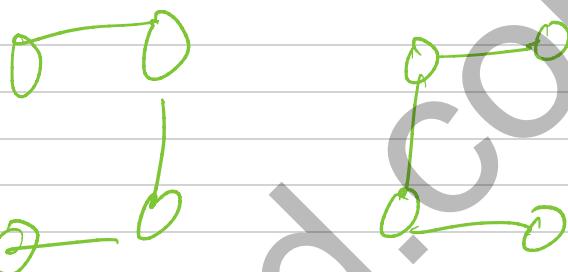
Not Linearizable

Don't call
Clear() in
Parallel
with a put()

Minimum Spanning Tree



Tree with subset of edges that connects all the nodes



minimum spanning tree Algo

$L \leftarrow$ Vertices

concurrent linked Queue

WHILE ($\text{size}(L) > 1$) {

if (! TRYLOCK(N_1))



$N_1 \leftarrow \text{REMOVE}(L)$

if (! trylock(N_1)) $E \leftarrow \text{GETMINEDGE}(M)$

{ fixup
continue;
} $N_2 \leftarrow \text{GETNEIGHBOR}(N_1, E)$

$N_3 \leftarrow \text{MERGE}(N_1, N_2)$

$\text{REMOVE}(N_2)$, $\text{INSERT}(N_3)$

Distributed
Programming in
Java

www.pushpsood.com

Map Reduce

Big Data

How BIG DATA WORKS
Abstract the data into
Key Value Pairs (K, V)

MAP

grouping
goes through
all the
 (K, V) pairs
from map function
and puts together
all the same keys

for($KEY, VALUE$) pair (K_A, V_A)

Map f |
[K_A, V_A] |
KV01 |
KV02 |
 K_B, V_B

Output

(K_A, V_A)

(K_A, V_A)

Group with same key

if $K_{A1} = K_{B1} = K$

($K, (V_{A1}, V_{B1})$)

REDUCE

REDUCE

($K, g(V_{A1}, V_{B1})$)

EXAMPLE

Input Data

$(W_R, 10)$ $(H, 11)$ $(W, 12)$ $(B, 13)$

Map (enumerate the factors in the Val)
(factors)

$(W_R, 2)$ $(H, 11)$ $(W, 2)$ $(B, 13)$

$(W_R, 5)$

$(W_R, 10)$

$(W, 4)$

$(W, 3)$

$(W, 6)$

$(W, 12)$

prime factors
of 13

prime factors of
10

Group (Same)

sum of prime factor
of college WR

Reduce
(sum)

$(WR, 17)$

$(H, 11)$

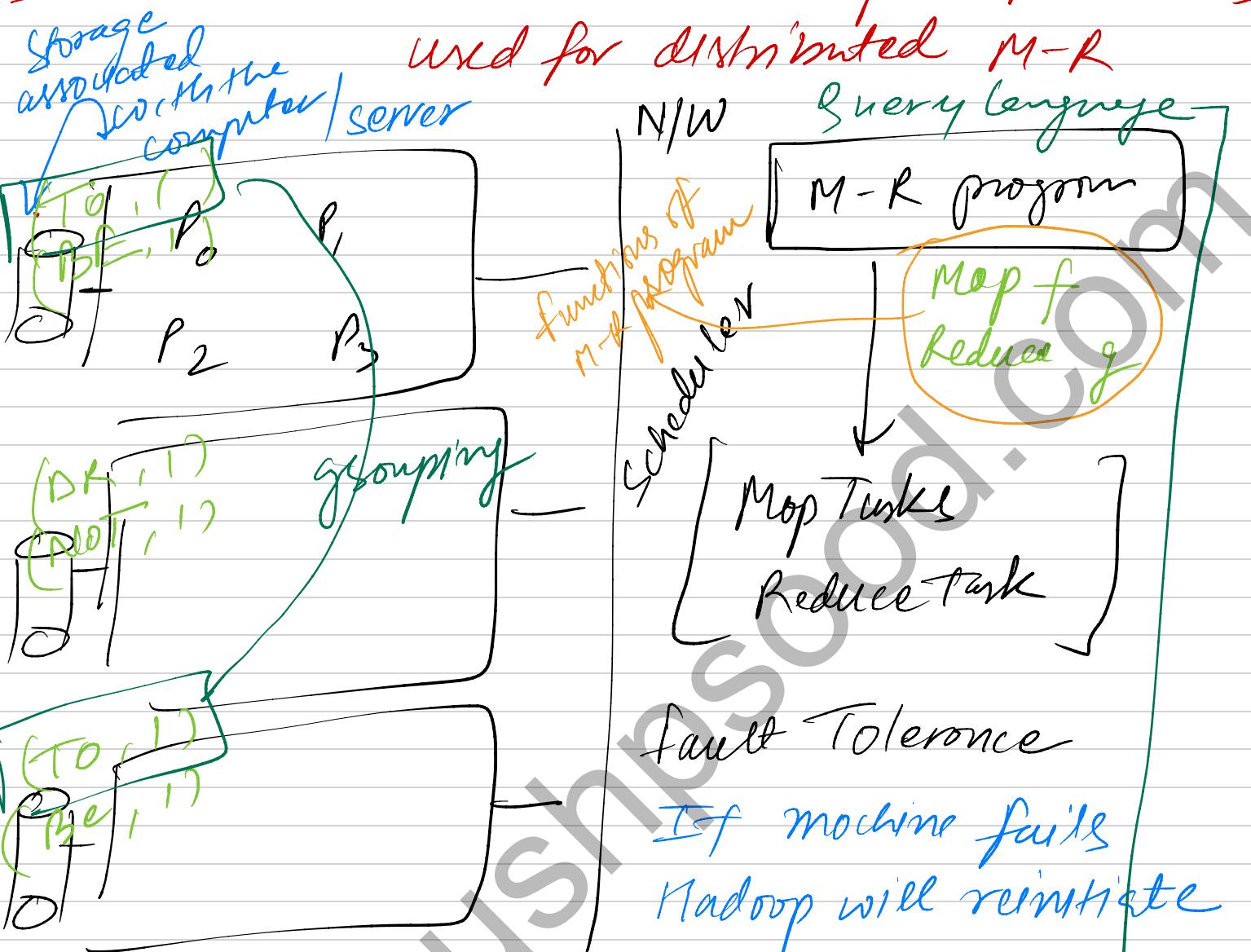
$(W, 27)$

$(B, 13)$

Simple operation M & R

These are keys basically
Colleges inside RICE University

Apache Hadoop Framework (Implementation of Map & Reduce)



word count

(TO, 1)
(Be, 1)
(OR, 1)
(NOT, 1)

INPUT K-V pairs

Map, f

Intermediate (K,V) pairs

grouping

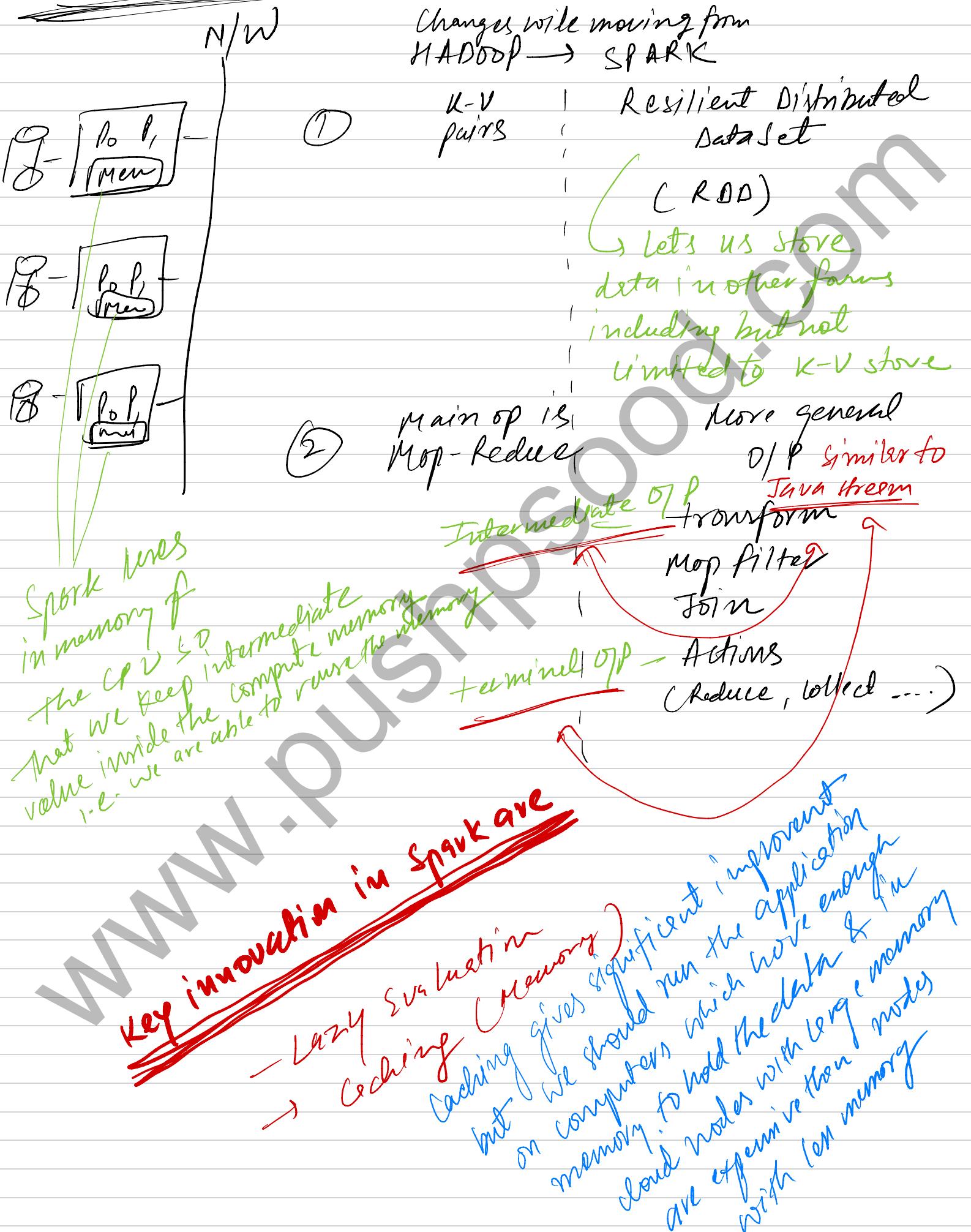
Reduce, g

Output (K,V pairs)

Hive,
Pig etc

SQL type
to automatically
generate
Map Reduce
queries

Apache Spark Project (Modern version of Hadoop)



EXAMPLE OF SPARK M/R FRAMEWORK IN JAVA

WORD COUNT (using spark API in Java)

- ① SparkContext (sc) = new JavaSparkContext ();
- ② Input = sc.Textfile ("...")
- ③ words = Input.flatMap (line →
line.split (" "));
- ④ pairs = words.mapToPair (x → new Tuple2 (x, 1))
- ⑤ counts = pairs.reduceByKey ((x, y) → x + y)

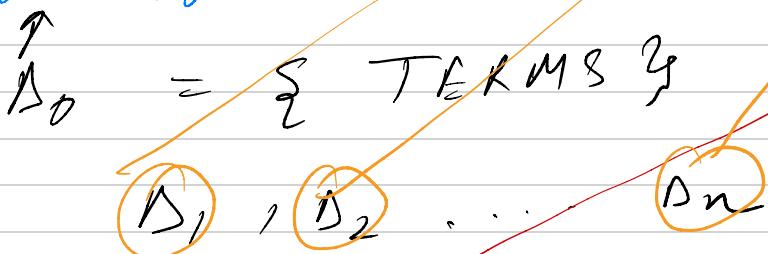
Eg of Hadoop M/R Framework

Document 1, 2..n

(Find the most similar doc across large no of doc)

TF-IDF Weights for info retrieval and document

Document Do



→ No of time term 1 occurs in Document 2

DF

→ No of documents the term x occurs in

Term 1 Tf_{11}, Tf_{12}, \dots

DA
DF_x

Term 2 Tf_{21}, Tf_{22}, \dots

Words like "the" appears a lot of time

IDF - Inverse Document freq = $\frac{N}{DF}$

($\frac{N}{DF}$ for commonly occurring word "the" will be 1)

Weight (Term_i, Doc_j)

$$= Tf_{(i,j)} \times \log \left(\frac{N}{DF_x} \right)$$

eg

The idea is to use this as weights to give more prominence to word that occur infrequently.

$TF(i,j)$ → Computing total frequency

Input : (DOC, TERMS Σ)

↙ Map

((DOC, TERM), Σ)

↙ Reduce

((DOC, TERM), TF)

↙ Map

when TF $\neq 0$

(TERM, Σ)

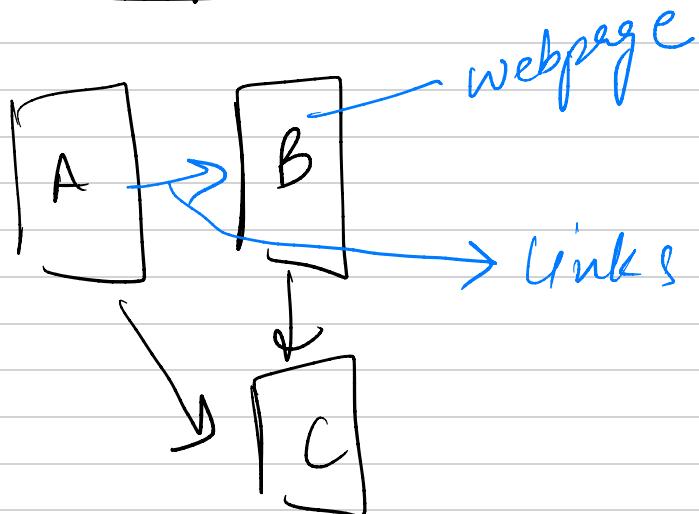
↙ Reduce

(TERM, DF)

Total frequency
the term appears
in the doc

for every occurrence
of term in the
document when
 $TF \neq 0$

Page Rank Example



$$\text{rank}(B) =$$

$$\sum_{A \text{ ESRC}(B)} \frac{\text{RANK}(A)}{\text{DEST_COUNT}(A)}$$

FOR (ITER - ...) Σ

1) FOR EACH PAGE $A \rightarrow B$

$$\text{CONTRIBS}(B) += \frac{\text{RANK}(A)}{\text{DEST_COUNT}(A)}$$

weighting
factor

2) FOR EACH PAGE B

$$\text{RANK}(B) = 0.15 + 0.85 \times \text{CONTRIBS}(B)$$

BETTER SUITED FOR
SPARK?

Spark is good for iterative programming
because of Cache!

CODE

FOR (ITER . . .) {

FOR EACH PAGE $A \rightarrow B$ {

CONTRIBS += LINKS

. JOIN (RANKS)

. VALUE

. FLATMAP TO PAR C

. DEST \rightarrow (DEST,

RANK (SRC) / DESTCount);

}

FOR EACH PAGE B {

RANKS = CONTRIBS,

REDUCE BY KEY (NEW SUM())

MAPVALUES (

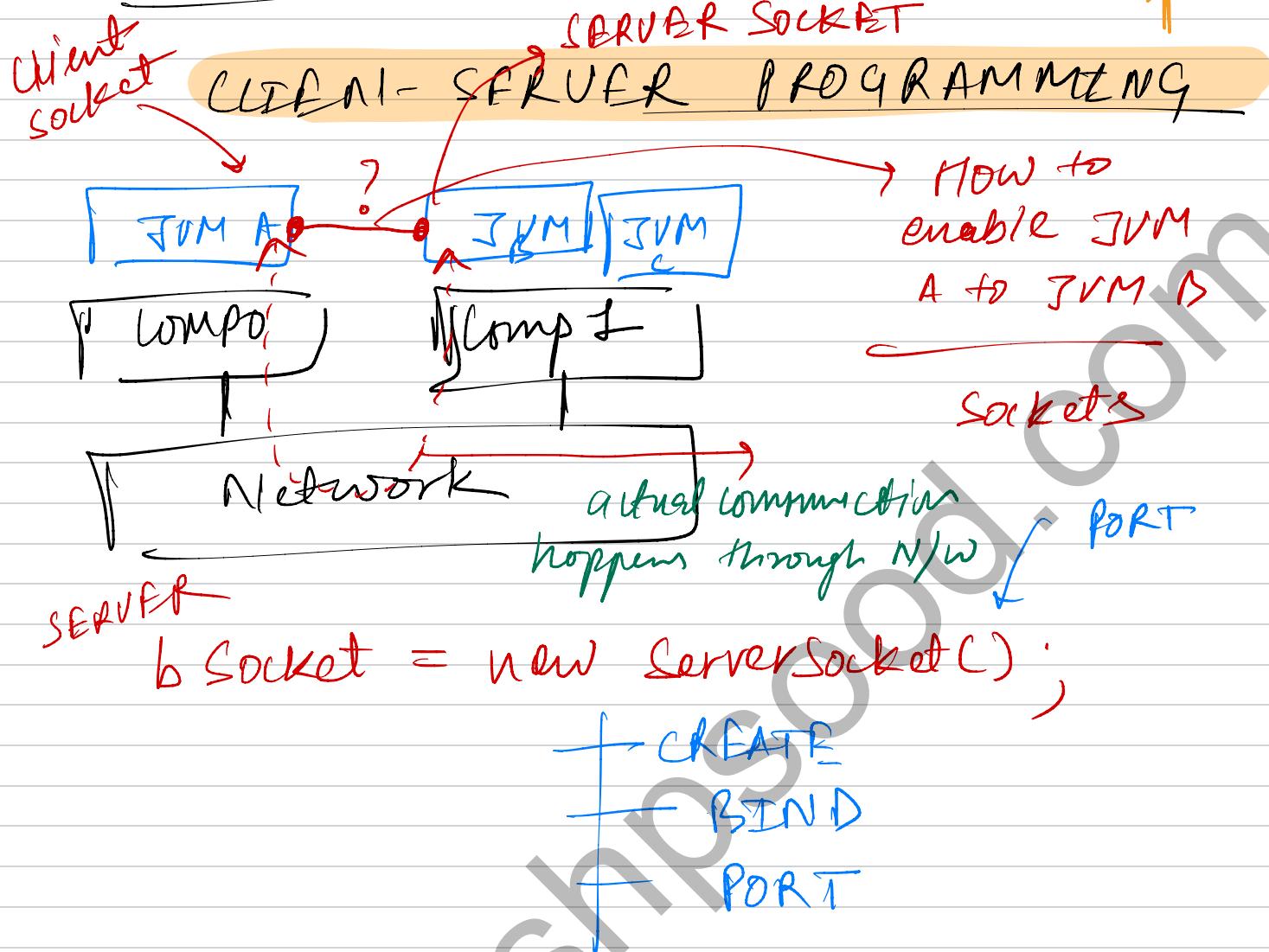
g

SUM $\rightarrow 0.15 + 0.85 * \text{SUM};$

g

INTRODUCTION TO SOCKETS

lower level concept
than distributed comp



a = bSocket.accept()

a.getInputStream() // Reads

a.getOutputStream() // Writes

CLIENT

aSocket = new Socket()

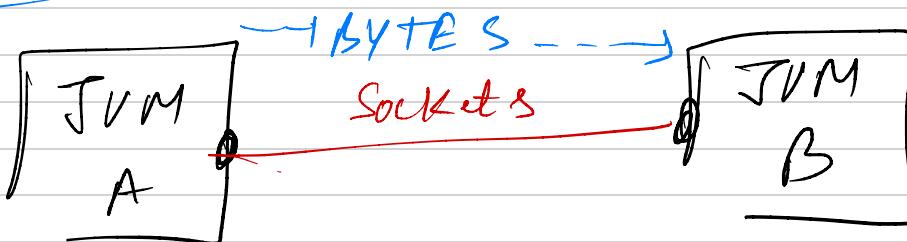
aSocket.getInputStream() // Read

aSocket.getOutputStream() // Write

SERIALIZATION / DE SERIALIZATION

Objects
A₀, A₁, A₂

B₀, B₁, B₂



Objects to Bytes

Serialization

Bytes to Objects

Deserialization

Choices for DB | Serialization

① complicated!
when network
objects

Custom

SER | DESR

② Use XML (Heavy weight)

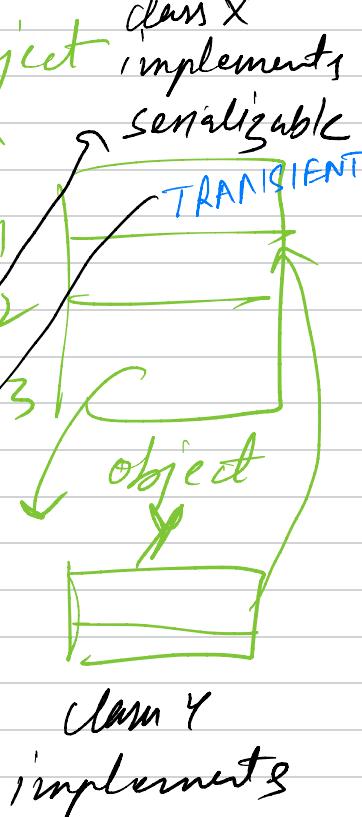
③ Using Java Serialization
& Deserialization
(It is robust)

④ Interface
Definition
language (IDL)

(C++, JAVA
etc) ⑤ PROTOCOL BUFFER

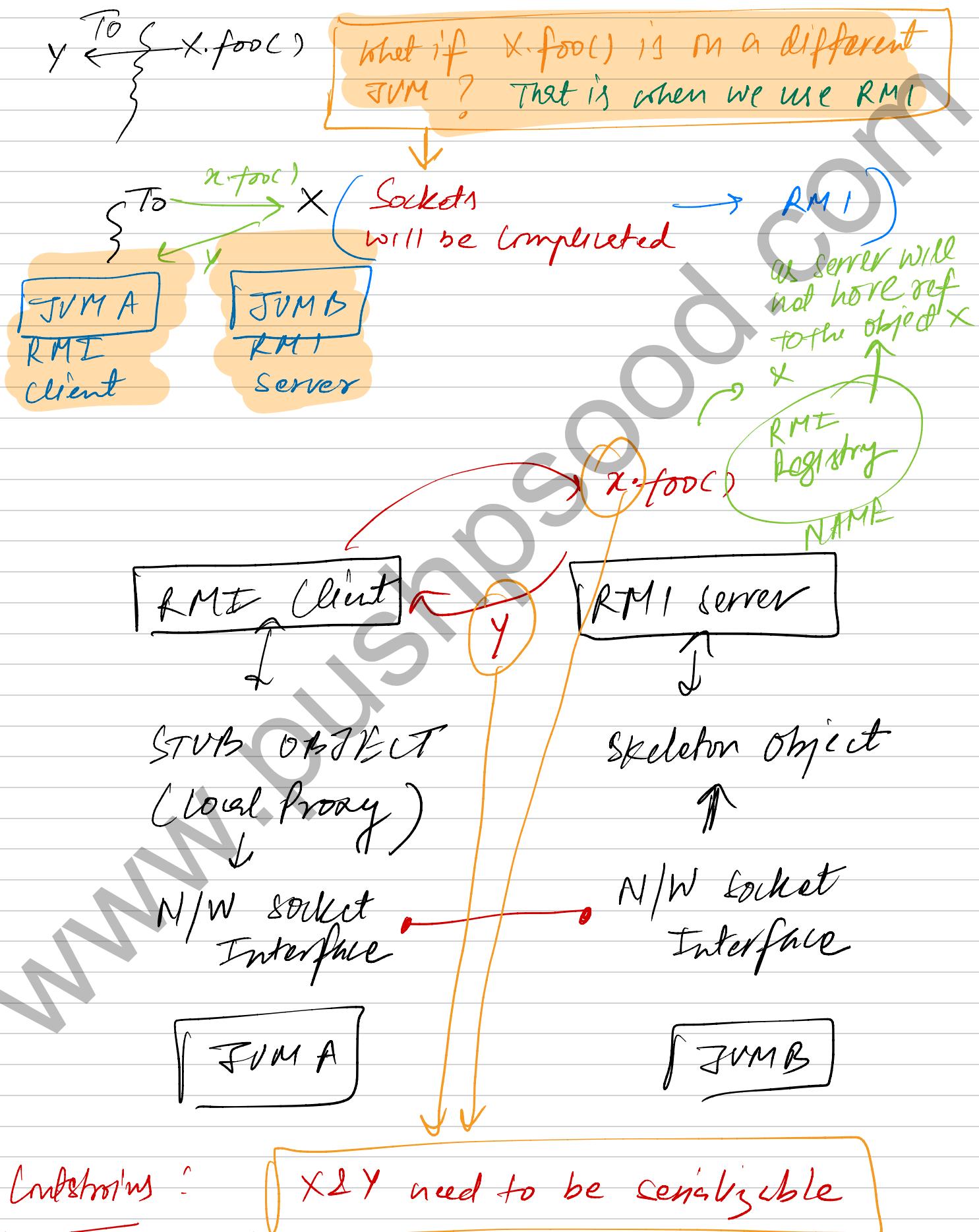
Language Agnostic PROTO FILE needs to be created!

if f₃ is
Transient its
content will
not be copied
Collection
overkill!



REMOTE METHODS INVOCATION (RMI)

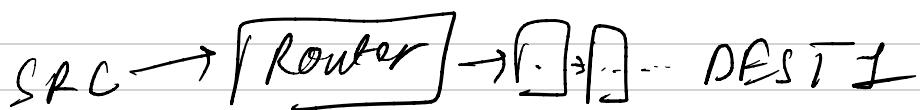
(Higher level concept build on Sockets and serialization)



MULTICAST SOCKETS (since JDK 1.1)

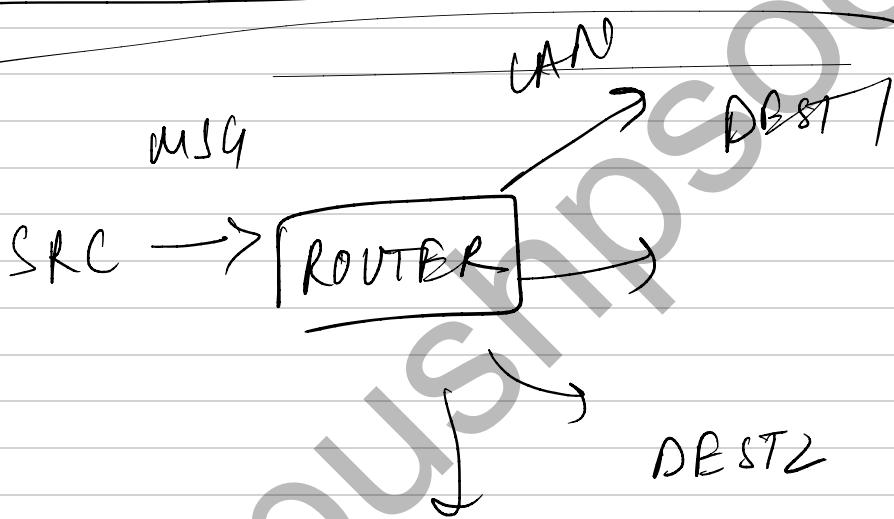
Usage of multicast: newsfeeds, video conference, gaming etc

Earlier we studied UNICAST



This is unicast Point to Point

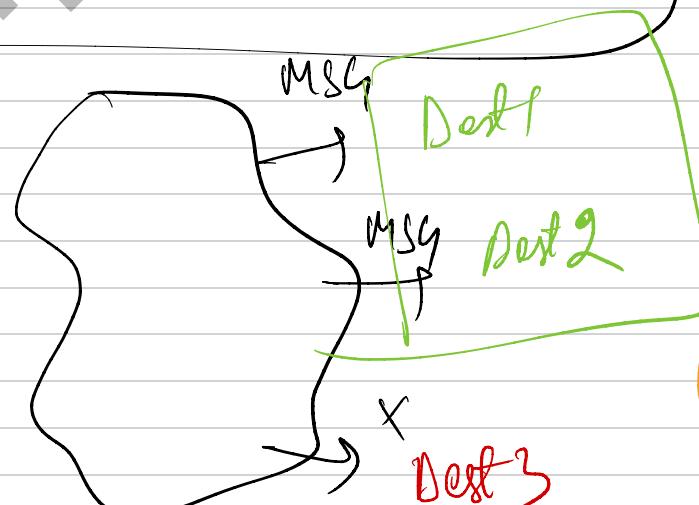
Unicast
(one message goes to one destination)



BROADCAST

(one message goes to all destination but local area)

only send to T1,2
MSG
SRC →



MULTICAST

(one message goes to selected destinations)

Multicast Socket Application

S = new multicast socket (PORT)

S. JOINGROUP (G);

MS4 = NEW DATAGRAM PACKET (...);

S. SEND(MS4);

MSG2 = NEW DATAGRAM PACKET (...);

S. RECEIVE(MSG2);

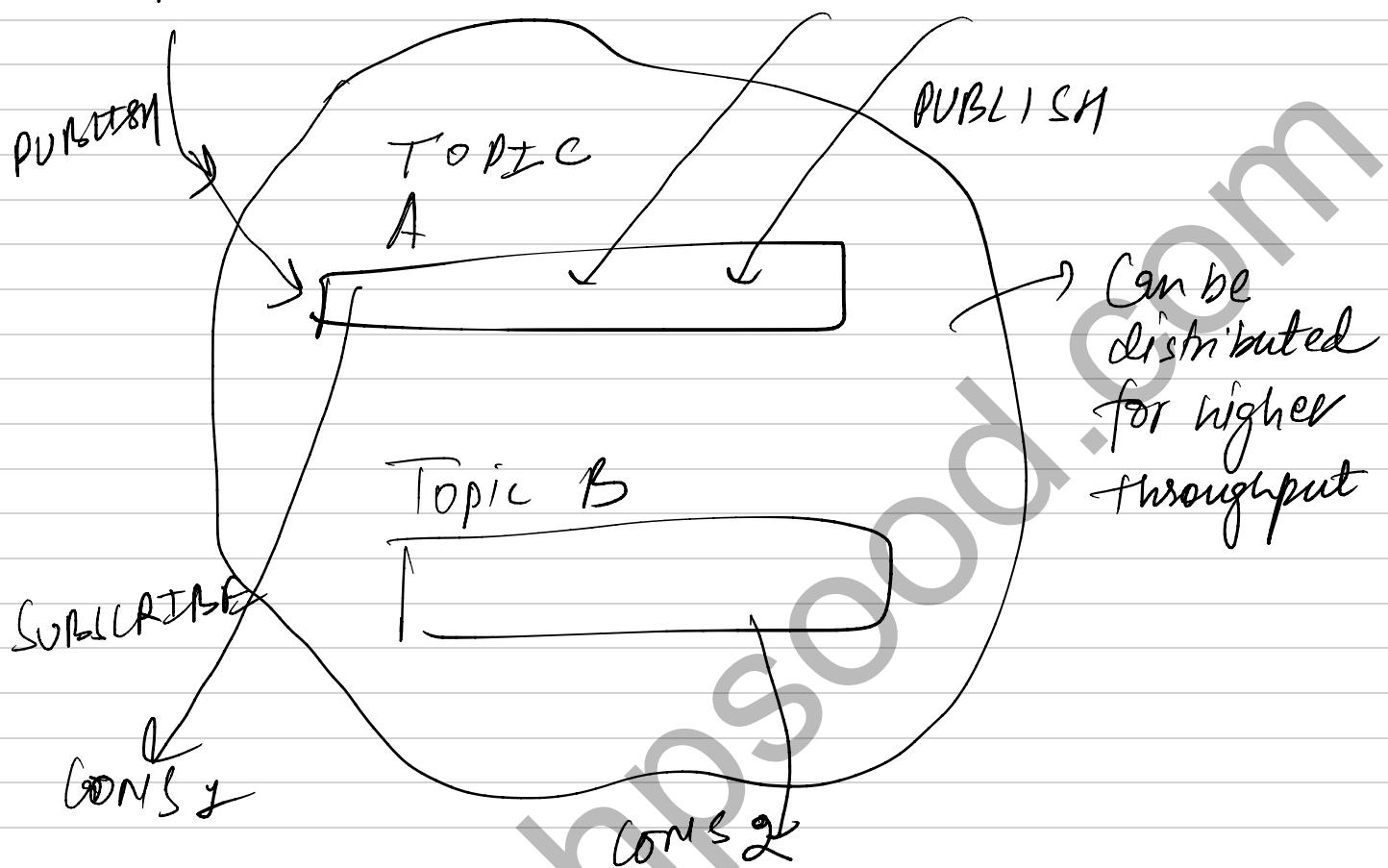
S. LEAVEGROUP (G);

PUBLISHER & SUBSCRIBER

(Higher level pattern to group communication than MULTICAST)

PROD 1

PROD 2



PRODUCER

P = new KafkaProducer(...)

P.send(
 how producer record
(topic, key, value))

P.close()

CONSUMER

C = new KafkaConsumer(...)

C.subscribe(
 (TOPIC, ...))

RECORD = Producer
 (time)

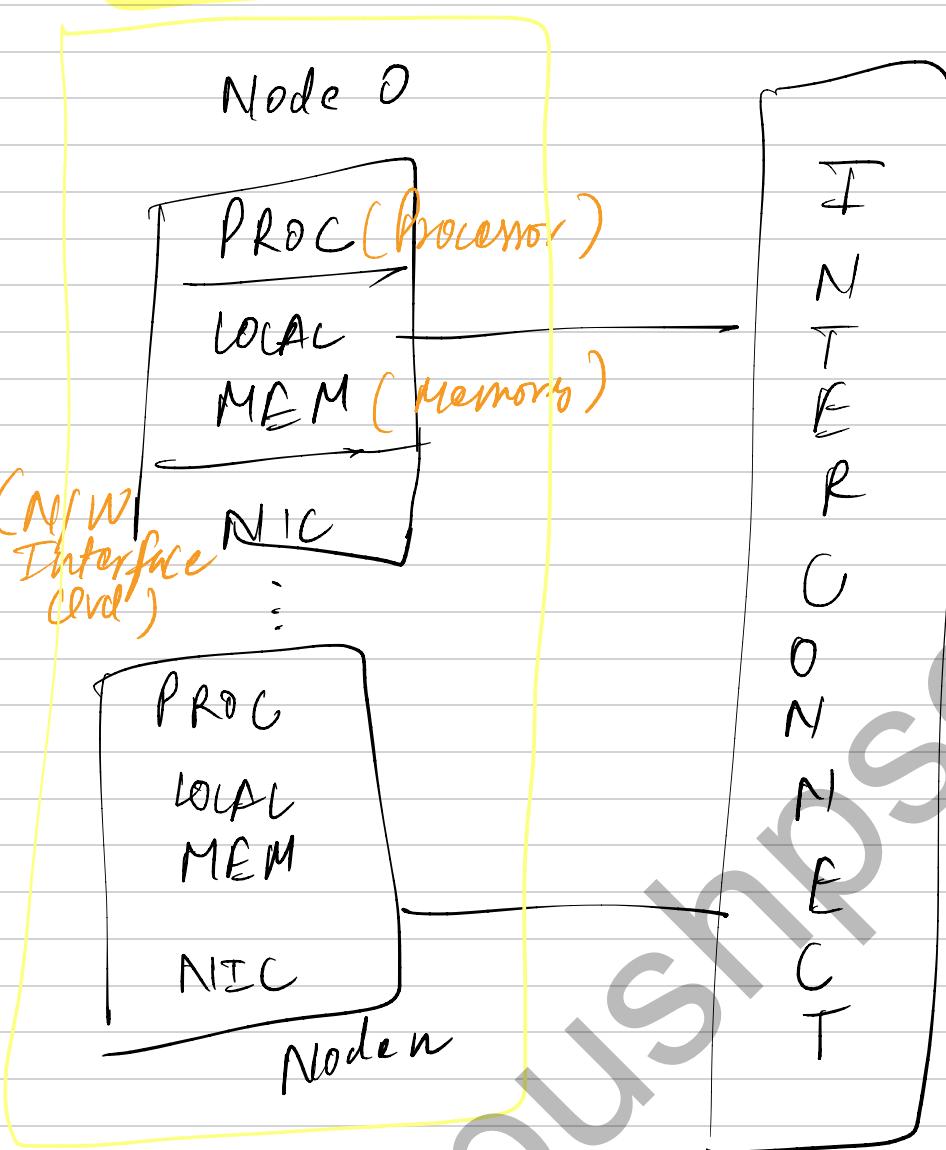
for (R : Records)

R.offset, key, value

How to use collection of computer as single parallel computer

DISTRIBUTED PARALLELISM

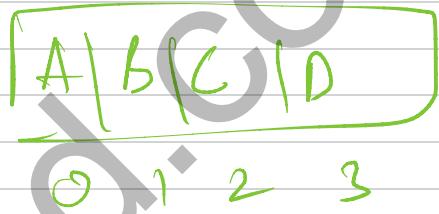
CLUSTER (Collection of Nodes)



DATA

DISTRIBUTION

GLOBAL VIEW: X_{global}



LOCAL VIEW

Node 0 : X_{local}

$$X_L = \boxed{A, B}$$

Node 1 : X_{local}

$$X_L = \boxed{C, D}$$

Single Program Multiple Data Abstraction (SPMD)

Interface : MPI

Main()
MPI - INIT
for (I = 0; I < LENGTH; I++)
 X_L[I] = X_L · LENGTH + RANK() * I

OUTPUT :

NODE 0 : X_L[0] = 0
 X_L[1] = 1
 X_L[2] = 2
 X_L[3] = 3

NODE 1 : X_L[0] = 1
 X_L[1] = 2

POINT TO POINT COMMUNICATION

Core primitives (Send & Receive)

(Basic primitive in distributed parallelism)

RANK 0

$s = "ABCD"$

SEND S TO RANK 1

using MPI program (different things using same program using conditions)

RANK 1

ALLOCATE BUF

RECEIVE INTO BUF

FROM RANK 0

PRINT BUF;

SEND

10

4

RECEIVE

PRINT BUF

MPI.FINALIZE();

We can serialization &
deserialization for objects

MAIN() {

 MPI. INIT();

 if (RANK == 0) {

 S = "ABCD";

 MPI. send (S, 0,

 S.LENGTH, 1,

 MPI. CHAR, 1, 99);

 DEST = TAG;

 } else {

 BUF = new CHAR[4];

 MPI. receive (

 BUF, 0, BUF.LENGTH,

 MPI. CHAR, 0, 99);

 Sender

 TAG

 } PRINT BUF;

Message Ordering & Deadlock

MPI is

blocking

operation

R₀

BLOCKED

send X to

R₁

BLOCKED

SEND Y

② → R₂

① receive X

Receive

Y from R₁, from R₀

leads to
deadlock!

DEADLOCK FIX

APPROACH 1

SEND → R₁

RECV ← R₀

RECV ← R₁

SEND → R₀

SAME ORDERING
FOR TWO MSGS

— SAME sender

— SAME receiver

— SAME type

— SAME tag

APPROACH 2

PART IS matched
allows 2 way
exchange

SENDRECV()

SEND
RECV()

MPI has all
API that clubs
Send & Receive
and takes care
of addressing
deadlock

(Extension to the message passing paradigm)

Non-blocking Communication

Immediate

operations
as opposed to
send & receive

(ISEND,IRECV)

fixes
Deadlock
Condition

R₀

R₁

S₁: Send

①

IDLE
Time

②

Ack

③

ACK

S₂:

S₃:

acknowledgment
that the
send has
been matched

Time
write

fix

N
E
+
W
B
R
K

S₄: RECV

S₅:

S₆:

Using IRECV
ISEND fixed
below deadlock &
improved parallelism

DEADLOCK

R₀

R₁

ISEND

ISEND

-X TO R₁

Y TO R₀

REQ0 =

REQ1 =

I RECEIVE

I RECEIVE

Y FROM R₁

X FROM R₁

;

;

WAIT (REQ0)

WAIT (REQ1)

USEY

USEX

WAIT (Req0)

WAIT (Req1)

S₃:

OVERWRITE

Send Value

Value

Non
blocking
communication

WaitAll ()

WaitAny ()

API waiting
for waiting
for one group
for one group

Collective Communication

(includes multiple ranks as opposed to P2P communication)

(multiple ranks)

R₀ : Send X to All Ranks

R₁ :

R₂ : Sum of all Y's into Z

R₃ :

Implementation:

main() { MPI. INIT();

INIT X[] = ...;

if (RANK == 0) X[0] = 99; // Root

BCAST (X, 0, 1, MPI. INT, 0)

PRINT (RANK, X[0]);

REDUCE [Y, 0, z / 0, 1,
MPI. INT, 2]

If (RANK == 2) {

PRINT Z[0];

}

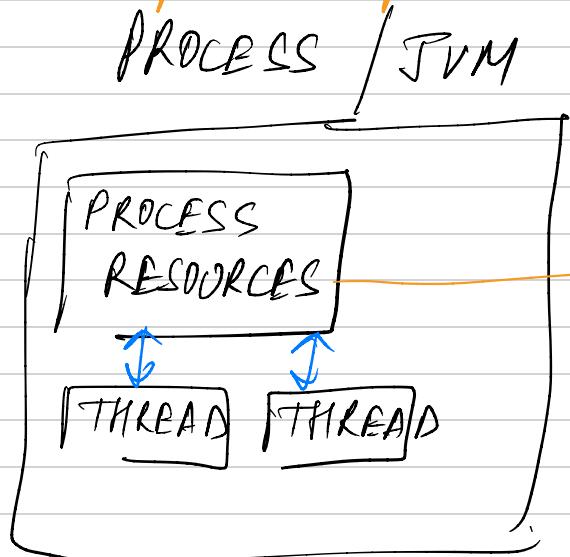
uses full parallelism
full consistency
writing of the CPU

FUNDAMENTAL BUILDING BLOCK OF DISTRIBUTED COMPUTING



PROCESSES & THREAD

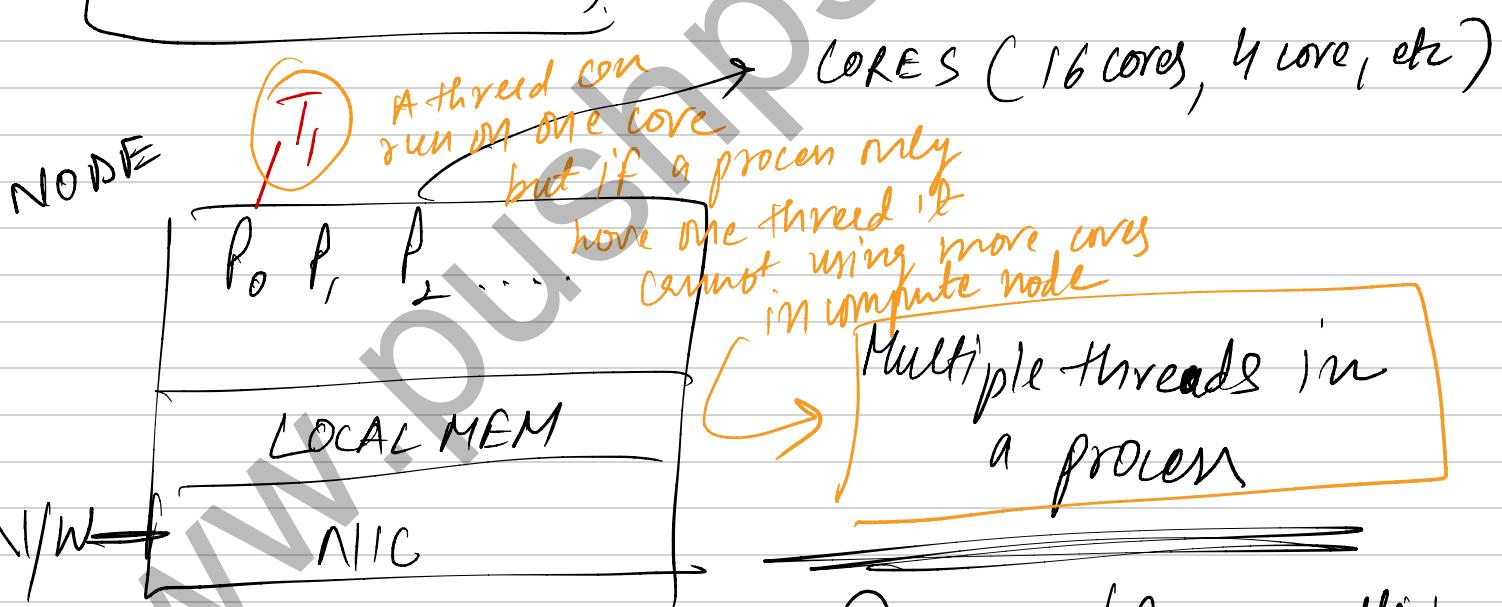
In context of JAVA, process is JVM



THREAD CAN RUN ON ONE CORE

Within a process we can create multiple threads

Threads share the process resources



- ① Memory/Resource efficiency due to sharing.
- ② Responsiveness (specially in N/W delay)
- ③

Adv Multiple Processors in a Node

① Responsiveness (Turn Delays)

example is Garbage Collection

② Scalability

Even though we can use multiple threads to increase throughput there is a limitation on throughput to flatten out or even decrease after a certain point

③

Availability / Tolerance

Using processes & thread to build something useful

↳ Multithreaded Servers

PORT, ... etc

Implementing file server:

① LISTENER = new SERVERSOCKET()

2) WHILE (TRUE) {

 2A) S = LISTENER.ACCEPT

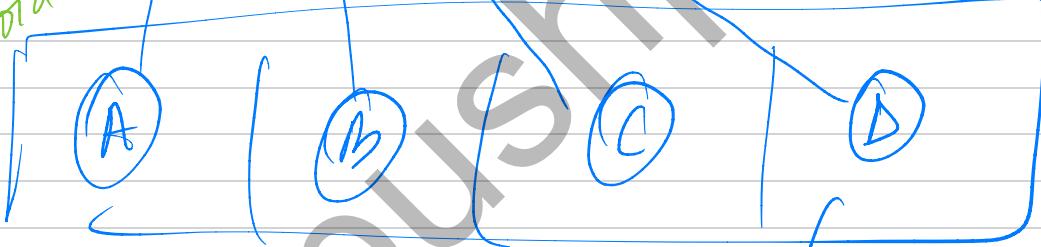
 2B) READ FILE REQUEST INFO
 FROM S. GETINPUTSTREAM();

 2C) ACCESS THE FILE

 2D) SEND FILE TO S. GETOUTPUT
 STREAM()

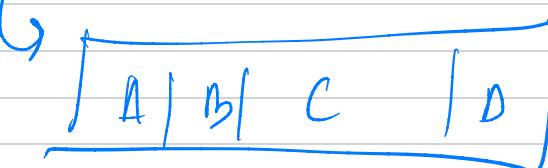
T =
new Thread
T.start();

We should
use thread pool
or tank
or avoid overused
(X)



delay because
of waiting!

Second
request



1st req



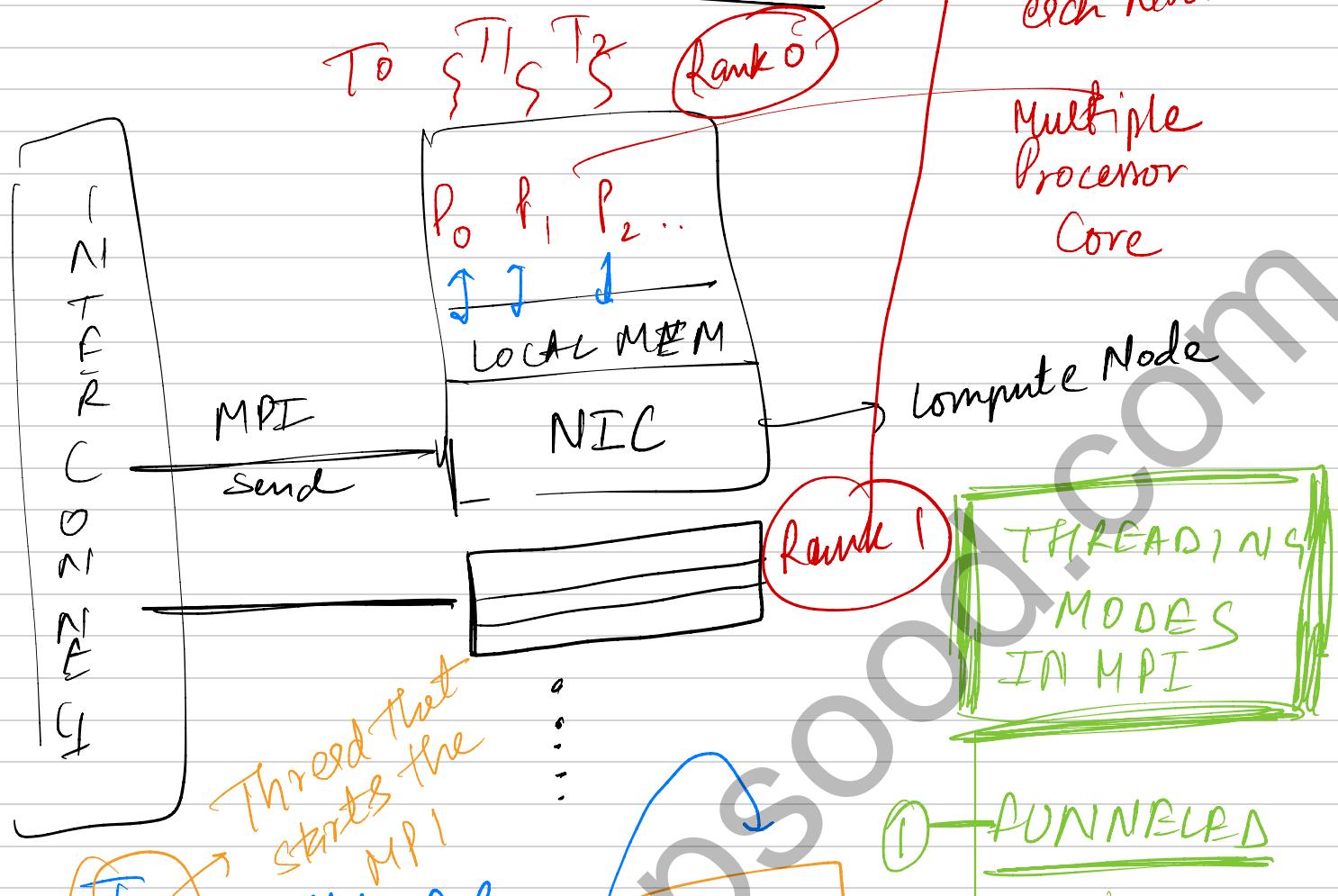
No
delay



Next req

Creating a thread everytime will be an overhead
so, we should be using a thread pool

MPI + THREADING



MPI::INIT()

CREATE WORKERS THREADS

YPI · SPND

MPT - RECEIVE

MPI : REDUCE

`MPI_FINALIZE()`

Can't have
2 threads
Waiting on
Same MPICell

MPI process in
each Node

Multiple Processor Core

Compute Node

Rank

THREADING MODES IN MPI

① PUNNELERD

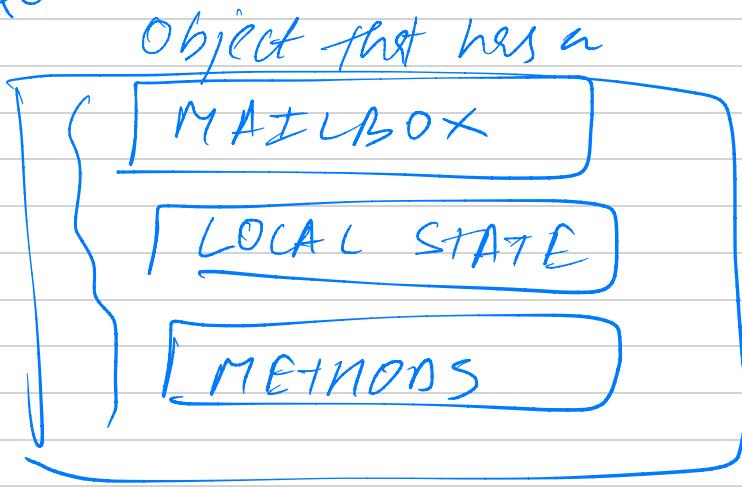
All MPI calls
are performed
by one thread

② - SERIALIZED
one MI
call at a
time

MULTIPLE
Multiple
MRI cells at
the same
time

DISTRIBUTED ACTORS

ACTOR



(2, 3, 4, 5, ...)



Non Multiples of 2

↓ (3, 5, ...)

Non Multiples of 3

↓ 3

↓ (5, ...)

↓ 5

↓ (7, ...)

SIEVE OF ERATOSTHENES

① CREATE CONFIG FILE
ACOR-MOST

Node 0

② CREATE REMOTE ACTOR

Node 1

③ LOOKUP REMOTE ACTOR BY LOGICAL NAME

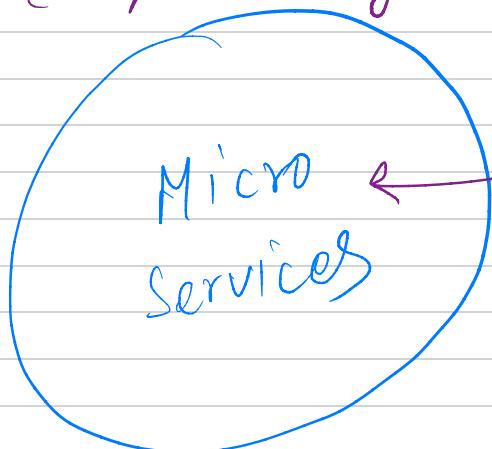
Non Multiples of 7

(but extra work
has to be
done in case
of multiple
nodes)

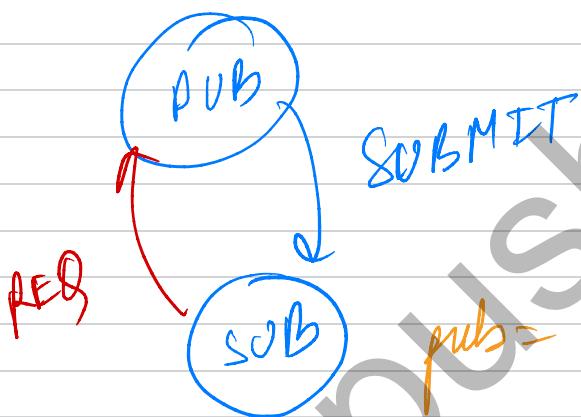
④ MESSAGES NEEDED TO BE SERIALIZED

REACTIVE PROGRAMMING

(Implementing distributed service oriented architecture)



ASYNCRONOUS EVENTS



REACTIVE STREAM
SPECIFICATION

FLOW. PUBLISHER
FLOW. SUBSCRIBER
PROCESS
SUBSCRIPTION

pub = new SubmittablePublisher();

pub.subscribe(s1);

pub.subscribe(s2);

pub.submit(...)

req item
at a time
f

SUBSCRIBER

onSubscribe(s1) {

SN.REQUEST(100); }

ONNEXT(ITEM) { PROCESS ITEM;

if (...) SN.REQ(100); }