

23/11/2011

UNIT - 2

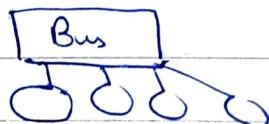
Communication Overhead is the time and depending upon the bandwidth the resources spent on transferring data and messages b/w the parallel components (processing elements or nodes)

- Depending upon the bandwidth, latency and how summary information is communicated among the processing elements, delays due to communication delay may vary from processor to processor

Sources of overhead① Interprocess Interaction

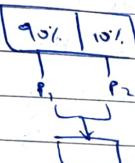
Any non-trivial parallel system requires its processing elements to interact and communicate data. The time spent in communicating data b/w processing elements is usually the most significant source of the processing overhead.

eg- a lot of data is getting exchanged but not that much o/p is produced



(2) Idling

Processing elements in parallel computing may become idle due to many reasons like load imbalance, synchronization, and presence of serial components in the program.

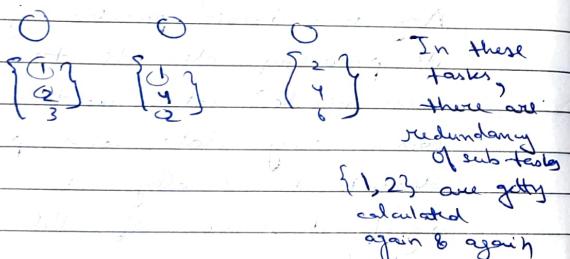


(3)

Excess / Essential Computation

The difference in computation performed by parallel program and the best serial program is the excess computation overhead.

eg - Fast Fourier Transform



In serial, this would get calculated once and used again.

#

Granularity

It is the measure of amount of work or computation performed by a task.

$$\text{Granularity} = \frac{\text{Computation time}}{\text{Communication time}}$$

Parallelism based on Granularity -

(1) Fine-grained parallelism

- program is divided into many smaller tasks
- it facilitates load balancing
- it is best exploited in architectures that support high speed communications

Drawback - If granularity is too fine, it is possible that the overhead required for communication and synchronisation b/w the tasks takes longer \Rightarrow time than the computation time.

eg - a) Connection machines (CM-2)

b) j-Machines

(2) Coarse-grained parallelism

- relatively large amounts of computational works are done b/w communication or synchronisation events.
- harder to load balance efficiently (bcz kuch task bohot jaldi ho rahi hai some are idle some are doing fast).
- message passing architectures are well suited for this parallelism.

eg - Cray Y-MP

(3) Medium-grained parallelism

eg - Intel iPSC

Impact of granularity on Performance

- (1) Synchronization overhead, scheduling strategies can negatively impact the performance of fine-grained tasks.
- (2) To reduce communication overhead, we can increase the granularity.

Definitions

- (1) Degree of concurrency - No. of tasks that can be executed in parallel.

(2) Critical path length - A directed path in the task dependency graph represents a sequence of tasks that must be processed one after the other. The longest path length in task dependency graph will be the critical path length.

Decomposition Techniques

The operation of dividing the computation into smaller parts; some of which may be executed parallelly

(1) Recursive Decomposition

based on divide and conquer algorithm.

(2) Data Decomposition

it divides the data in program into parts and then assigns them into instructions.

eg: Matrix Multiplication Problem,

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \quad B = \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix}$$

$$\text{Task 1 : } C_1 = [A_{11} \times B_{11} + A_{12} \times B_{21}]$$

$$\text{Task 2 : } C_2 = [A_{11} \times B_{12} + A_{12} \times B_{22}]$$

$$\text{Task 3 : } C_3 = [A_{21} \times B_{11} + A_{22} \times B_{21}]$$

$$\text{Task 4 : } C_4 = [A_{21} \times B_{12} + A_{22} \times B_{22}]$$

→ Tasks induced by data decomposition are performed entirely, and each task

performs useful computation towards the soln of the problem.

30/01/24

(3) Exploratory Decomposition

Decomposition are to a search of a state space of solutions

e.g. 15 puzzle problem

1	7	2	3
8	6	4	5
9	12	10	11
13	14	15	

This is the state space

→ Tasks induced by exploratory can be terminated before finishing as soon as the desired solution is found.

(4) Speculative Decomposition

When a program may take one of many possible branches depending on the result from computations preceding the choice.

e.g. Nested if

if ()

{

if ()

In a program if and if statement

contains another if and the program is already decomposed, the lines of code are assigned to multiple processors without even checking the conditions.

- different processor execute both if's.
- The processors will run both the if statements and its inner if's.

Inner if is discarded if 1st if becomes false.

(5) Functional Decomposition

Parallel Processing in Memory

3 types -

- ① Shared Memory Architecture
- ② Distributed Memory Architecture
- ③ Hybrid " " (Combination of 2)

(1) Shared Memory Architecture

- a) Multiple processors operate independently but share the common memory as global address space

b) Shared memory systems are tightly coupled systems

- c) Only one processor can access the shared mem location at a time.
- d) Changes in a memory location affected by one processor are visible to all other processors.

They are able to use the updated values.

→ divided in 2 categories based on the Memory Access Time.

- (1) Uniform Memory Access (UMA)
 - (2) Non-Uniform " " (NUMA)
- special category → Cache only Mem. Architecture (COMA)

Uniform Mem Access

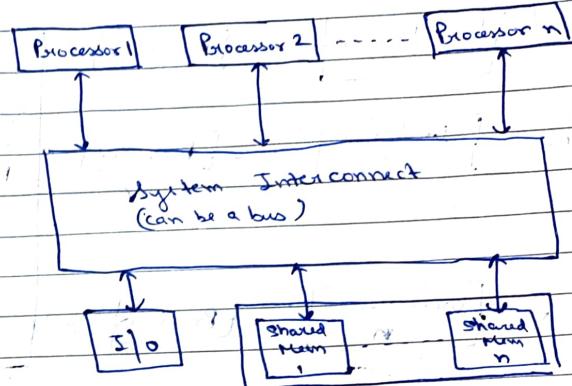
- (1) The physical mem is shared by all the processors.
- (2) All processors have equal access time to all memory locations.
- (3) used in symmetric multiprocessor machines (SMP)
- (4) Each processor may have private ~~access~~ cache is called cc-UMA

aka cc-UMA
Cache coherent

cache coherency = if processor no update
koi baki bhi

access kar paye update ho

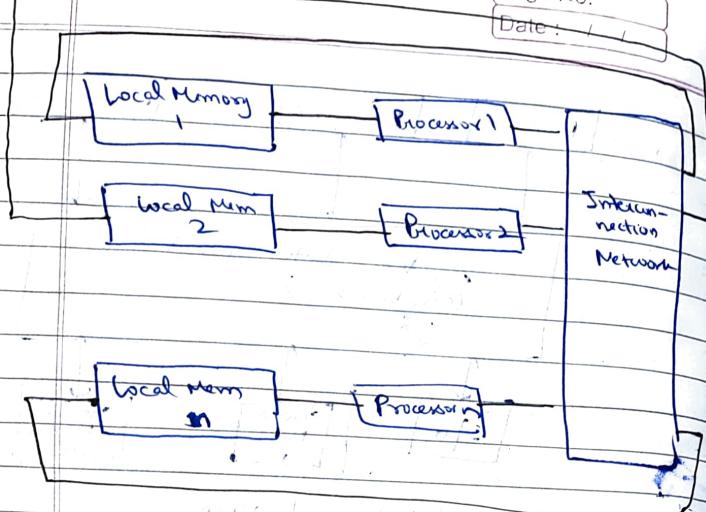
- if one processor updates a location in shared memory all the other processors know about the update.



Architecture of UMA

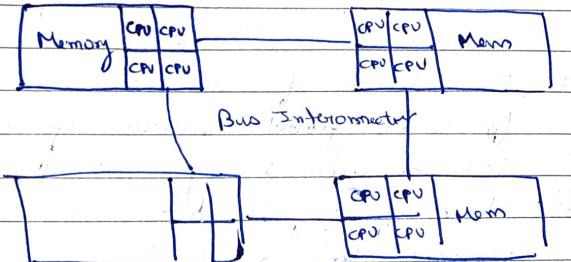
Non-Uniform Memory Access

- (1) access time varies with the location of the memory
- (2) the shared memory is physically distributed among all the processors called 'local memories'
- (3) The collection of all local memories forms the global address space which can be accessed by all the processors.
- (4) It is faster to access a local memory ^{with} a local processor.
- (5) The access of remote memory attached to the processors is slower due to the added delay through the interconnected network.



NUMA architecture

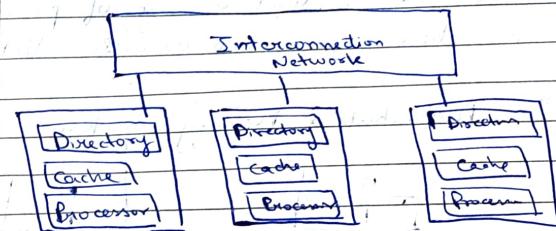
- ⑥ It is often made by physically linking 2 or ~~more~~ more SMPs.
(symm. multiprocessor)



Connections when multiple SMPs

special card
→ COHMA

- ① The distributed / main mems are converted to caches.
- ② All caches form a global address space and there is no mem hierarchy at each processor / node.
- ③ Remote cache access is mediated by the distributed cache directory



Directory - it helps in finding which data is located at which position and how to access that position.

- # Advantages of Shared Mem Architecture
- ① Global address space provides a user friendly programming perspective to memory
 - ② Data sharing b/w tasks is both fast and uniform due to the proximity of memory to CPUs.

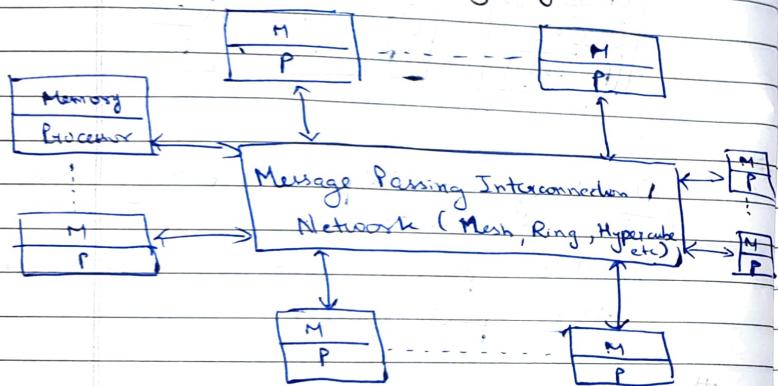
* Disadv

- ① Lack of scalability b/w mem and CPU.
- ② Programmer responsibility for synchronisation constructs that ensure correct access of global memory.
- ③ Expensive - as more processors are involved

There is a limit on no. of processors in shared mem ^{can} so we use distributed

31/01/24

Distributed Memory System



- ① It consists of multiple computers called nodes, interconnected by a message passing network.
- ② Each node is an autonomous computer consisting of a processor, local memory, and sometimes attached disk or input-output peripherals.

- ③ Distributed memory are loosely coupled system (shared-area tightly-coupled)
- ④ Memory addresses in one processor don't map to another processor. So there is no concept of global address space.
- ⑤ Bcz each processor has its own local memory, it operates independently.
- ⑥ Changes made to the local memory have no effect on the memory of other processors. Hence, the concept of cache-coherency doesn't imply.
- ⑦ When a processor needs to access the data from another processor, it is done by passing message b/w the processors, which is usually the task of the programmer.

Advantages

- ① Memory is scalable with no. of processors
- ② Increasing the no. of processors, the size of the memory also increases independently.
- ③ Each processor can rapidly access its own memory without interference and without the overhead incurred with trying to maintain cache-coherency.
- ④ Cost effectiveness [because we can use other topologies like mesh ring]

commodity, off the shelf processors and networking

Disadvantages

- ① The programmer is responsible for many of the details associated with data communication b/w the processor.
- ② NUMA takes a longer time

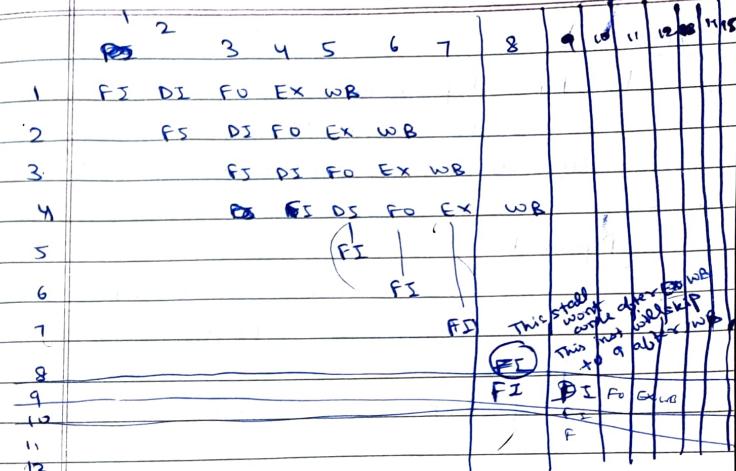
Hybrid - combo of all adv & disadv

Numerical (Instruction Pipeline)

Consider an instruction pipeline with 5 stages (F1, D1, F0, Ex, WB) without any branch prediction. The stage delays are 5ns, 7ns, 10ns, 8ns & 6ns.

There are intermediate storage buffers after each stage and delay of each buffer is 1ns. A program

consisting of 12 inst is executed in this pipelined processor. I9 is the only branch instructions and its branch target is I9. If the branch is taken during exe of the program then what is the time needed to complete the program.



$$\text{Max delay} = 10$$

~~$$\text{Cycle time} = 11 \text{ ns}$$~~

~~$$\text{Execution time} = (\text{Inst} + 1) \times \text{IP}$$~~

$$\text{Total cycles} = 15$$

$$w=2$$

$$1 \times 1 \text{ ns}$$

3 stall

$$\text{no. of cycles} \times \text{Cycle time}$$

$$\text{Total time} = 15 \times 11 \text{ ns}$$

$$= 165 \text{ ns}$$

Q2- There is inst pipeline with 4 stages. Stage delays 5, 6, 11ns, 8ns. Delay of an interstage register is 1ns. Find speedup of pipeline in steady state under ideal conditions as compared to corresponding non pipeline implementation.

$$(B+1) \times 11 \\ \approx 12$$

Ans -
Non-Pipelined processor

Time to execute N instructions =
 $\text{cycle time} = (5+6+11+8) \times N = 30N$

Pipelined

$$\text{Cycle Time} = \max(5, 6, 11, 8) + 1 = 12 \text{ ns}$$

For N inst = 12 N

$$\text{Speedup} = \frac{30}{12} = 2.5$$

06/02/24

(5) Functional Decomposition

In this approach, the focus is on the computation that is to be performed rather than the data manipulated by the computation

Constructs for Parallelism

(1) Creating concurrent processes

- uses **FORK-JOIN**

in the original construct, a FORK statement generates one new path for a concurrent process and the concurrent process uses the JOIN statement at their ends.

(2) UNIX Heavy weight Processes

The UNIX system call `fork()` creates a new process. The new process is an exact copy of the calling process except that it has a unique process id. It has its own copy of the parent's variable.

They are assigned to the same values as the original variables initially.

On successfully executing `fork()` returns zero to the child process and returns the process ID of the child process to the parent process.

Processes are joined with system calls `wait()` & `exit()`.

Main program

FORK

FORK

These processes
are aka
spawned
processes

FORK

JOIN

JOIN

JOIN

Definition

→ Heavy weight process - It completely separates program with its own variables, stack and memory location.

→ light-weight processes = Shares the same memory space and global variables b/w the processes.

→ POSIX Thread (pthread) -

→ Portable Operating System Interface
pthread is an execution model that exists independently from a language as well as a parallel executor model.

- It allows a program to control multiple different flows of work that overlap in time.

Thread

header used < pthread.h >

- There are around 100 thread procedures all ~~are~~ prefixed 'pthread-' and they can be characterised into 4 groups -

① Thread management
(creating threads)

② ~~Mutexes~~

③ Condition Variables

④ Synchronization b/w threads using read or write locks and barriers.

a) Thread management - Thread operations include thread creation, termination, synchronisation (JOIN, blocking), scheduling, data management, and process interaction

Syntax:

```
int pthread_create(pthread_t * thread,  
const pthread_attr_t * attr,  
void * (* start_routine)(void *),  
void * arg);
```

① thread - gives pointer to an unsigned integer value that returns the thread id of the thread created.

② attribute - pointer to a structure that is used to define thread attributes like detached state, scheduling policy, stack address, etc

(3) set to null for default thread attributes
 start routine - pointer to the self-routine
 that is executed by the
 thread.
 return type is void.

→ The function has a single attribute but if multiple values need to be passed to the fn, a struct must be used.

(4) argument - pointer to void that contains the arguments to the function defined in the earlier argument

eg - Main program

```
pthread_create(&thread1, NULL, proc1, &arg),
}
pthread_join(thread1, &status);
```

Thread1

```
proc1(&arg)
```

```
}
```

```
void pthread_exit(void * retval);
Pointer to an integer that stores the
```

return status of the thread terminated.
 The scope of this variable must be global so that any thread waiting to 'join' this thread may read the return status.

```
int pthread_join(pthread_t th, void ** thread_return);
```

- ① th → thread id of the thread for which current thread waits.
- ② thread_return → pointer to the location where the exit status of the thread mentioned in th is stored.

`pthread_t pthread_getselfid();`
 self - used to get the thread id of the current thread.

`int pthread_equal(pthread_t t1, pthread_t2);`
 if two threads are equal then it returns a non-zero value (here p1, p2 are threadids).

```
int pthread_cancel(pthread_t thread);
```

b) Mutexes

- Used to prevent data inconsistencies due to race conditions.

Race condition often occurs when 2 or more threads need to perform operations on the same memory area. ~~address~~ but the results of ~~operati~~ computations depends on the order in which these operations are performed.

c) Condition variables

is a variable of type `pthread_cond_t` and is used with the appropriate functions for waiting and later process continuation.

A condition variable must always be associated with the mutex ~~without~~ to avoid the race conditions.

13/02/24

~~Processor~~ Parallelism

2 Types of HSM

- ① Hardware Parallelism → It is defined by the machine architecture and H/w multiplicity. It displays the resource utilization, patterns of simultaneously executable operations.

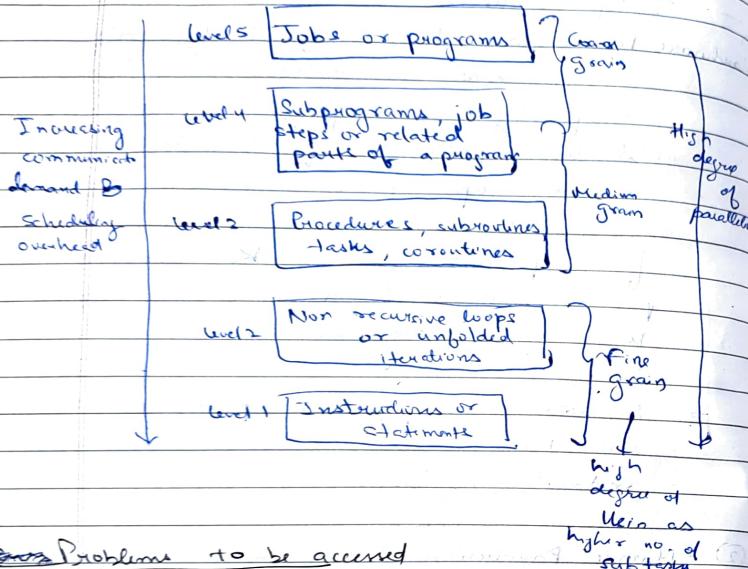
If a processor issues K instructions per machine cycle then it is called K -issued processor.

→ A conventional pipelined processor takes one M/c cycle to issue a single instruction cycle.

- ② Software Parallelism → This type of HSM is revealed in the program profile or in the program flow graph. (flow chart)

* threads (ekhudi dhikna)
and
processes

Grain Packing



~~Problems to be addressed~~

Q1 - How can we partition a program into 16 branches, program modules, micro-tasks or grains to yield the shortest possible execution time.

Q2 - What is the optimal size of the concurrent grains in a computation.

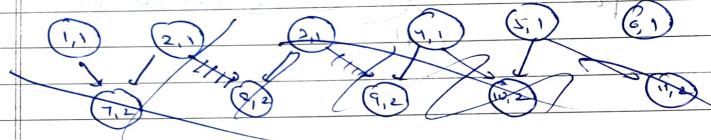
Solution - Grain Packing - The idea of grain packing solves the above ~~above~~ challenges.

→ The idea of GP is to apply fine

grain first in order to achieve higher degree of parallelism. Then one combines multiple fine grain nodes into a coarse grained node if it can eliminate unnecessary communication delays or reduce the ~~overhead~~ all scheduling overhead.

1. $a=1$ 7. $g = arb$ 13. $n = 4 \times l$
2. $b=2$ 8. $h = cxd$ 14. $m = 3 \times n$
3. $c=3$ 9. $i = dxe$ 15. $o = nai$
4. $d=4$ 10. $j = exf$ 16. $p = oxh$
5. $e=5$ 11. $K = dxh$ 17. $q = pag$
6. $f=6$ 12. $l = j \times k$

Grain cycle :



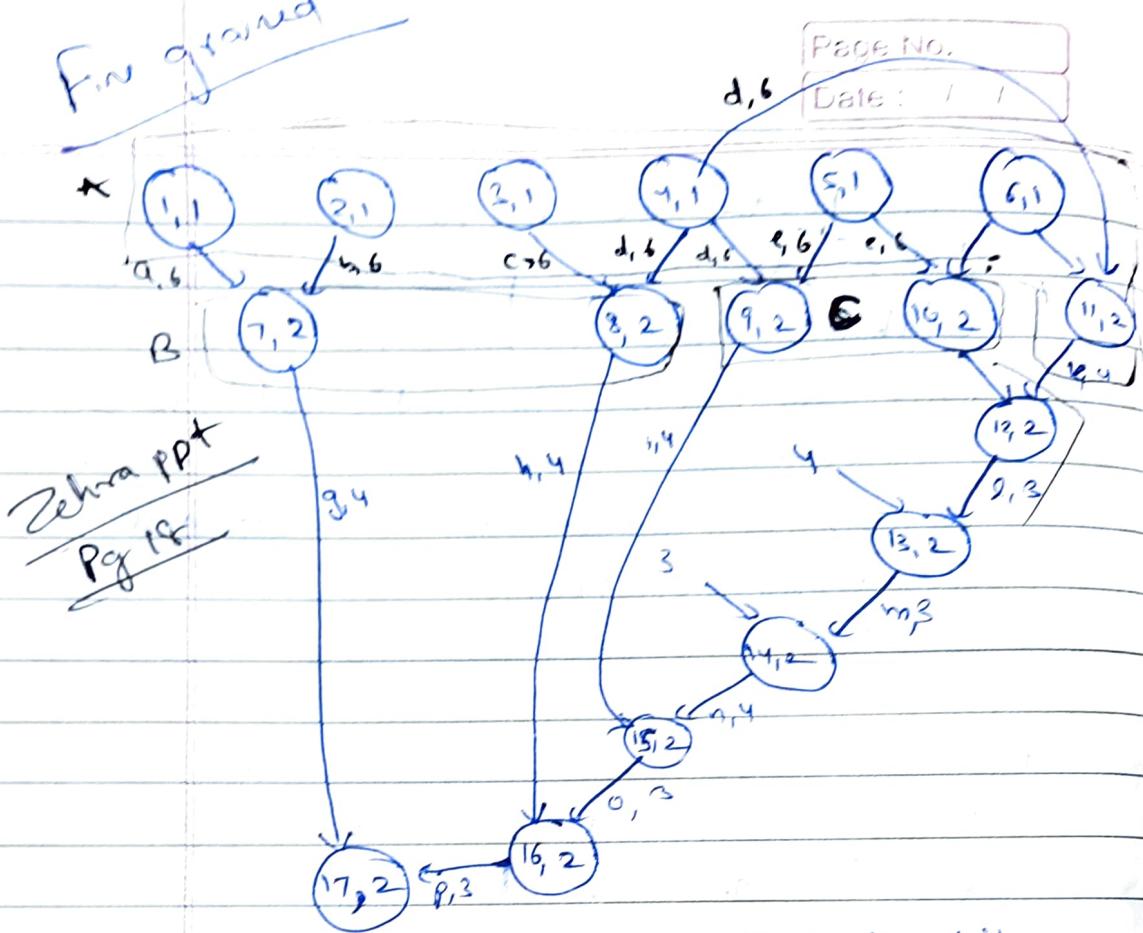
P.T.J

FIN around

Page No.

Date: / /

d,6



Zebra PPT
Pg 18

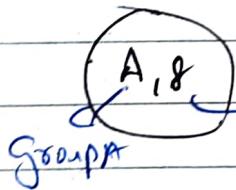
(n, i) = (node, grain size)

(x, i) = (input, delay)

(u, v) = (output, delay)

use like kitne cyclagenge

Affter
Packing



add all the
grain