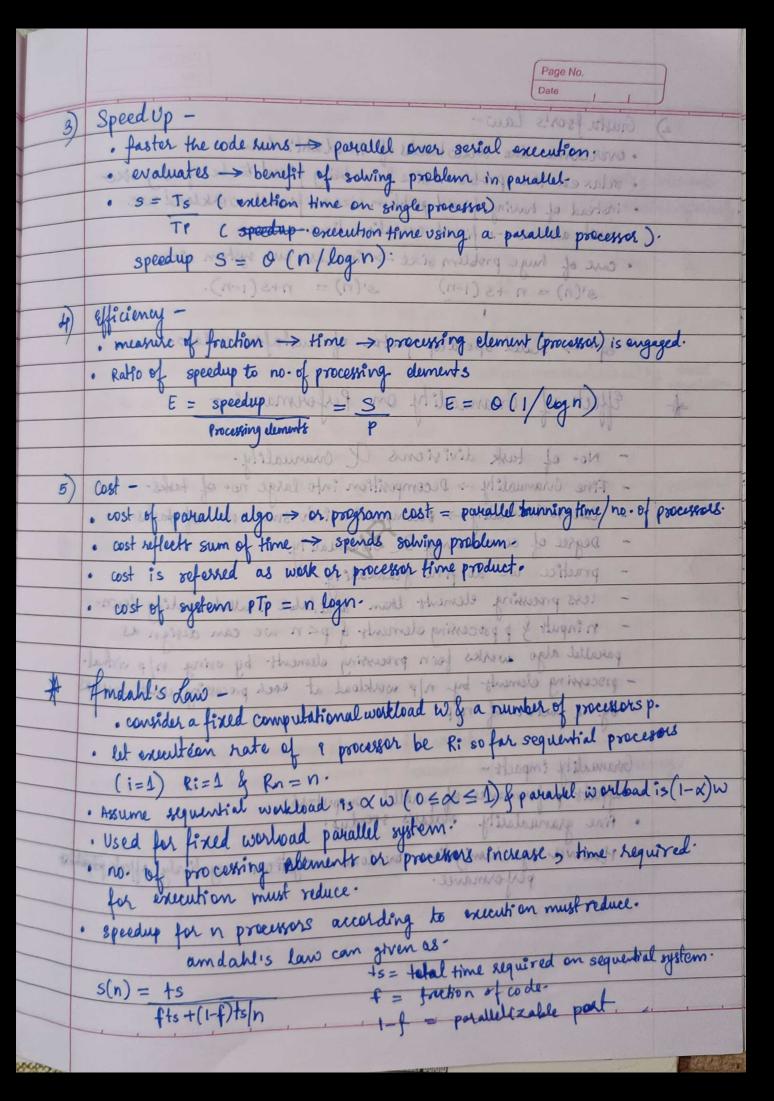
U-4(HPC)

Page No.

	Analyse and all a a a a a a a a
	Analytical Modeling of Parallel Programs.
0	
2	Sources of Overhead in Parallel programs. The land of the series
100	Lavor meter se se se se semente de position - able pins.
	· sequential algorithm -> its nuntime (evaluated by)
	- parallel nuntime > depends on input size margan los los
	> parameters of machine (communication)
	depends on Input no of processors.
	execution time of parallel algo > Inputsize t no of processing elements -
	· combination of parallel algo & parallel richtechile.
	Methics - 1) Efficiency 2) Speedlast 3) last of Parallel Brichard
*	Methics - 1) Efficiency 2) Specifiest 3) Cost of Parallel Erectical 5) Execution time -baldrovo file
· 10	· speed up - captures relative benefit of solving a problem in polved
	Interprocesses communication - made and a grand and a
	1) Processor working on non-trivial parallel problem will need to talk.
	@ source of parallel processing overhead -> time spent to interact & comm"
_1	@ data dependancy exists in problem & no of processing element
عالم الما	adata dependancy enistr in problem & no of processing element
	working on same date.
1	(4) Iwo processoring element are independent they can't work in parallel.
	starts to moment last processor finishes execution.
2)	Idling (Load Balance (im)) -
	process may become title because of many reason of
	· process may become title because of many reason. · Diff factors that contribute to dding of processing elements like
	load (mbalance, synchronization, preserve of serial component in
	dynamic substack generated
	number tack generation approcation if this very afficient to predict size of
	subtasks assigned to processor in warning
	to matchain unclosm workload among processors, ay namic subtasts go-detection
	· dunamic tage acculation apply -> subtasts assigned to processed in accuration
	· maentain uniform workeload among processors. Can't duide peoblem.
	· maintain uniform workeload among processors. Can't duide problem. statistically blue processing elements.

	Date Date
	(39H) N-U (Date)
3)	Excess computation - believed for probable has explained
1	
	· computation -> not by social version:
	· social alen > difficult to varallelize or are repeated across proun
	· serial algo -> difficult to parallelize or are repeated across procus · excess = parallel - Best serial to min recommendor
	a : Book sowing moran -> diff factors parallel code -> excess comp
	· Best serial program -> diff factors parallel code -> excess comp
	I depends on Input no of process ers.
*	Production measure of Analysis - ale Jelling to ent materials
N.	Performance measure of Analysis - also tallaced to another idence a continuous a
	Mobiles -> 1) Efficiency 2) Speed(ost 3) (ost 4) Parallel Overhease
	5) Execution time - bashow file
	speed up - captures relative benefit of solving a problem in parallel.
	S = speedup = To Chime taken by best sequential algorm / by singleprocess
-4	S = speedup = To Chime taken by best sequential algorm / by singleprocess TP [time taken by p processors].
Conso	Execution Time -> Jime to complete a task & denoted by Jesse.
1	1) sequential or social runtime (Ts) - time chapse -> b/10 start of execution.
	1) Sequential or social huntime (Ts) - time elapse - blue start of execution. end of execution - sequential computer
	Mas serial runtinge in the
- 14	ii) Parallel Runtime (Tr) - time that clapse from moment first processor
	starts to moment last processor finishes execution.
1	2) Idling (Load Balance (im)) -
2) to total Overlacad - mum to surround all amound mem seeing.
)	denoted by To. Privallel processing time sport by and I among
	time spent by parallel program - time spont by serial program-
har	ts semention time on a single processor
6.00	Tp - emultion time using a parallel precures.
برادا	· Speed to Ps geven by S = O(n/logn)
	dynamic task gonderent part I subject amount to proceed in reference
1	. much in uniform weekle at among processed. Can't dide problem.
1	take walk blue processing elements.



		Paga No.
		Page No.
2	Grusta fsor's Law-	Speedup -
	· overcome the draw backs of Amdahl's law.	
	· relaxed the problem size from being fixed to	
	· instead of having fixed problem size or fixed we	
	· Case of huge problem size - increase our enter	
	· case of huge problem sixe - increase our syste	
	s'(n) = n + s(1-n) $s'(n) = n + s(1-n)$	
yest.	pan - exted speed up tacket at fourtheleans	1 chicona
	eqn -> scaled speed up factor of Grustafson's	
4	Ellet of Granusting dements	· Karlo of Species
	Effect of Granuality on Performance.	a = 3
	- No. of task divisions & Granuality.	19
	- Fine Granuality - Decomposition into lases	1 1.12
Processed	- Fine Granuality -> Decomposition into large no.	of tasks 100)
	- Carrie Granuality -> occomposition into small no	of tasks?
	- Degree of concurrency of Granuality-	06 LV201-34 L500
	- practice we do fine gramatity.	101 12 80-87
	- less processing elements than available is called	scating down-
	- n inputs & p processing elements of p< n we can parallel also works for a processing elements I	aurign as
	- processing elements by n/p workload at each process	ing hip virtual-
	ised computational workland tig no for workland and best	and thetaks
We	water of a processor be R: so far sequential proces	· conviour a
	Chauman in those in the control of t	
(h-1) 20	affects performance of parallel computers:	=18 (L=j)
	· affects performance of parallel computers: · fine granualarity invests speedup.	· Assume sque
· Jorg	MONO SALISMAN TO THE MANUAL THE MANUAL TO TH	a hazalu a leak para
	performance	allos for
	in processors according to execution must reduce.	tot brown
	The state of the s	and drymods .
Steps -	relables law can given as" the telak time sequenced on sequential	
	2000 12 118000 = 1	5(n) = 15
	They should be a single of the	1, 614

	Page No.	
	Date	
1	Scalaboility of parallel systems - silving to stribut stalgrapes	4
4	Example of adding no rues with in preceding demarks -	**
	. capacity of increase speedup in proportion to no of processing el	ements.
	. reflects parallel system's ability to utilize increasing processing resources	
	· effectional is given as - E = 5 = Ts = 1	• (
	P PTP 11+TE	1
	. evaluates behaviour of parallel algo - input size approaches infailly	1
	. total overheard to increasing for of pools or such as	1
1	overy parallel program always has a sequential that runs for time to	1
2.07	· overload increases superlinearly due to communication overhead idealing;	excell and him
	e As p1 Et goes down Eco/pie de de la selection de la selectio	
	. scale more than efficiency shall decrease. The most of the	
	· efficiency may be more than efficiency shall decrease	
	efficiency may increase if profesen size is increased keeping no of p	nemous
AL	· facilitatio comparistion by aiff parallel algo & system config.	August .
A k	Minimum execution time & min cost -	
*	TO CAR TO LIVE JOHN A PONCE	*
	· no- of processors in crease parallel runtime continues to descrease & asym	phoally
	D. V I I I A SALIII A A LIVILLA VIV. III LAVINII TA III	
	The state of the s	
	· give mine execution me as of the safet to age of the proof of age of a safet ag	
	- sommers of the prime some state of age .	
	also sign sin some to for which This mintelland .	
	gives no of processing elements for which Tp is min- gives no of processing elements is bounded by degree of concurrency max. no of processing elements is bounded by	.1
	HUN. TO OF DECOME	
	San politing beach ((with a three)	
-	min = with to (with c [w]) lead prilling was .	
Laid	. Verely specialized compulses decise as	
and It	and what a second of the secon	
8218	· Epinolastion techniques title lang search	
	* approximation techniques tile loop revealing, date bleeking synchriter	
	ephinicalism techniques like loop remeding data blocking sprechization fuller temporare effectioned of possibile makes weeks multiplication algo.	
	epitonication techniques like loop remelling, data blocking synchrischer fullen temprene effectioned of possible makes weeter materialises	

	Page No. Date
*	Asymptotic Analysis of parallel and programs -
	Example of adding n. nos. with in processing elements.
- the land	$T_p = O(\log n)$. $S = O(n)(\log n)$
us efficiely	after partitle system's abiling to will be indeased, processing resource
	· Offsetency is given is - E = 5 = Ts - 1 Stable
X	Maria Very Muliplication
	· evaluates behaviour of parallel algo - input size approaches infinity
0	· char actorizing also complexity & performance scaling
10	examines how algorithmic officiency changes with including problem
James 1	common asymptotic notation -> Big O, Theta & Big Omega-
	· provides insight into scalability & efficiency of parallel algo.
	· identify dominant computational growts.
	· gui des algo selection & captimization
- PHENETY -	with lawer asymptotic complexities are preffered for large only.
1. Marking	· facilitates comparision b/w diff parallel algo & system config.
	or wheremen orwers on time & min with -
*	Matrix Vector Multiplication -
Mealgraph	one of processions in crease parallel runtime continues to descrease of a
	· computes product of a matrix of rector, resulting in new rector
	· fundamental operation in linear algebra & computational science.
	· can be parallelized to organite concurring & improve performance.
	· algo for include sono-wise, columnise & black wise approaches
	· parallelized aim to distribute computational workload acress
· fran	muliple processing unets officiently.
	· optimize date weating of minimize commen ovalhead.
	· key building block in scientific & engineering applications & img processing.
	· executial for accelerating computations in large scale numerical simulations
11001	o Using specialized computing devices like GPUs enhaces performances notionization techniques like loop unabling the plants - veeks multiplates
	eptimization techniques like loop unalling, data blocking & vector multiplication. further empreve efficiency of parallel matrix vector multiplication.
1	algo-
15	

Matrix Multiplication Algorithm -· computes product of two matrices using parallel processing. · divides matrices using parallel processing into smaller blacks of distributes among multiple processing. . each process computes a partion of resulting matrix independently. · reduces overall time & leverage consultent processing. . utilizes techquies like data processing partioning & common among processes scales well with no. of available processing units, improving parformence. · implement ed using parallel algo like Cannon's algo. optimizes resource whilization by distributing workload across muliples . everheads menimized by efficiently exchanging data b/w processors o 61905 one: no mediate priority todoenof multiples (1900) Will Execution Rate & Redundarry = x000 x 12/129 1 along meanues inclase in required. computation when using more processing units to defined on per component basis. & letterny person see parts arrighed a sample time based on sample times of components it is concerted to-· COOA - · Compute Unitied Device Architecture · launched by NYEDER in SECT. As an preparating against interface providing popular comprehasion using art WOR and thetre -· Befor e: CLIDA and: without & pixel shadens were used for pesselve , pixel shades is a core compared that can execute per proces-. winter shade is also - a cope compenent este press & it is assembly language specific, used for geometrice operations. . com entension ef. C++ traduce to support 6.40 eccorations. compensate - , purallel compute empires inside NVIDIA ELVI . 05 haved level ouggest for hardware-. Veer made dained. . durice level At I for devolu

i prox translation 's et actualectures for how