### **Self Introduction**

Dr Tay Seng Chuan Tel: 65168752 Email: <u>scitaysc@nus.edu.sq</u> Office: S16-02, Dean's Office at Level 2 URL: http://www.physics.nus.edu.sg/~phytaysc

I have been working in NUS since 1990, and I teach mainly computer and computing subjects. (Deputy Director at Centre for Remote Imaging Sensing and Processing: 2001 to 2004)

I am (was) a hands-on person on High Performance Computing. Todate I have made my hands dirty on PC-transputer System, network of workstations, and GRID System. Participated in IBM-iHPC's Blue Challenge in 2002.

Current Duties in NUS:

Day: SM2/SM3 Programme (MOE), Faculty IT Unit (Science), Physics Dept (Senior Lecturer). I am also a mentor of MOE Gifted Education Programme.

Night: Assistant Master of Temasek Hall, Resident Fellow (Block E).























# High Performance Computing : *The Motivation*

The distance calculated is reduced by a factor of 2 to 3 in many materials used to build computer. As such the distance traveled is reduced to 15 cm or less (dimension of PDA). Can we find a PDA (Personal Digital Assistant) of this size and yet is able to achieve the supercomputer performance?











### **Performance Metrics**



#### <u>Speedup</u>

Speedup (S) is the ratio between the time needed for the most efficient sequential algorithm  $(T_1)$  to perform a computation, and the time needed to perform the computation on a parallel machine incorporating parallelism  $(T_p)$  where p processing elements, p > 1, are used.





<b>Example:</b> <i>A</i> In the followin number. What number?	<b>AS</b> ngg tis	rid of the	r <b>ch</b> cells time	<b>Pr</b> the rec	<b>oblem</b> re is one quired to	e negati o find t	ve he
	4	8	7	5			
	$\frac{2}{1}$	12	10	9			
	3	11	13	-7			
							22

#### Search Problem (cont'd)

4	8	7	5
2	12	10	6
1	16	15	9
3	11	13	-7

Let c be the time required to process a cell, and a top down search on the columns is adopted.

When 1 processor is used,  $T_1 = 16c$ .

When 4 processors are used and assume no overhead,  $T_4 = 4c$ .



### **Speculation :**

Can speedup be greater than p (or can efficiency be greater than 100%)?

Given a sequential algorithm (G1) that incurs the least runtime  $(T_1)$  for a problem. Suppose its parallel version is able to achieve a superlinear speedup. We have

$$S(p) = \frac{T_1}{T_p} > p$$
$$T_p < \frac{T_1}{p}$$











# Overheads incurred in Parallel Computing (cont'd)

Give a square grid filled with random data. The heap property can be established by applying row sort and column sort.







# Overheads incurred in Parallel Computing (cont'd) $T_{p} = (\alpha + n\beta) + s + (\alpha + n\beta) + (\alpha + n\beta) + s + (\alpha + n\beta)$ $= 2s + 4 (\alpha + n\beta)$

### Overheads incurred in Parallel Computing (cont'd)

The communication overhead should be treated as a relative term with respect to computation granularity. Why?

35

36

$$S(p) = \frac{8s}{2s + 4(\alpha + n\beta)}$$

If  $s >> \alpha + n\beta$ ,  $S(p) \approx 4$ . (What is \$2000 if you already have \$1,000,000!!)

Otherwise S(p) is sub-linear, and in the worst case it can become a <u>slowdown</u>.

## Load Balancing

Should workload be evenly (equally or fairly) distributed?

- Common belief : Yes.
- My answer : Not always yes if our objective is to minimize the runtime of a parallel program.
- It has been analytically and experimentally established that an even workload distribution may aggravate communication overhead, resulting in a longer program execution time.

























Para	lle	I B	loc	ck	So	rt (	(co	ont'	d)		
	6	26	25	18	3	45	49	5			
	32	2	63	56	36	59	15	7			
	41	29	54	22	46	12	10	60			
	61	52	28	34	62	23	30	51			
	33	42	53	11	8	38	48	9			
	17	57	21	58	47	64	55	35			
	16	27	43	37	1	20	44	4			
	13	39	14	24	19	31	40	50			
			Ra	and	om	Dat	a				
										5	50









### Conclusions

- ② The speed of the traditional computer cannot be increased indefinitely. A parallel platform offers the potential to reduce program runtime.
- Parallel computing also incurs communication overhead and parallelization penalty.
- An even workload distribution scheme may not result in the least runtime.

55

Conclusions (cont'd)
To achieve a good speedup, the interacting effects of process granularity, communication overhead, load balancing, and parallelization penalty must be considered. You got to know how to do it, and I will teach you in this course.
Good speedup has been achieved in many HPC applications. That's why CZ4102 is not a fairy tale.