

5: "Flavors" of Scalability in Manycore Processors

Seminars in *Scalable Computing*

Vittorio Scarano

Università di Salerno



Dottorato di Ricerca in Informatica

1/78



PLAN

- 1 PARALLELISM AND SCALABILITY FOR MULTICORE
 - Amdahl's Law
 - Gustafson-Barsis's Law
 - Amdahl's law and Multicore Processors
- 2 TIME- AND MEMORY-INDUCED SCALABILITY
 - Fixed-size scalability for multicore processors
 - Fixed-time scalability for multicore processors
 - Memory bound scalability for multicore processors
- 3 CONCURRENCY-INDUCED SCALABILITY
 - The model
 - Implications on design of multiprocessors
 - Conclusions

2/78



PLAN

- 1 PARALLELISM AND SCALABILITY FOR MULTICORE
 - Amdahl's Law
 - Gustafson-Barsis's Law
 - Amdahl's law and Multicore Processors
- 2 TIME- AND MEMORY-INDUCED SCALABILITY
 - Fixed-size scalability for multicore processors
 - Fixed-time scalability for multicore processors
 - Memory bound scalability for multicore processors
- 3 CONCURRENCY-INDUCED SCALABILITY
 - The model
 - Implications on design of multiprocessors
 - Conclusions

3/78



SCALABILITY

- **Scalability:** the property of a solution to a problem to maintain its *efficiency* as the *dimension* grows
- Some keywords to be addressed in the context of parallel programming:
 - Efficiency: speedup over the "corresponding" sequential solution
 - Dimension: processors number, type or interconnection; problem size (memory)
- Big-Oh notation for algorithms: scalability, but only in principle
 - what happens when you fill the current level of the memory hierarchy you are using
 - what happens when number of processors grows to infinity
 - ...

4/78



PLAN

1 PARALLELISM AND SCALABILITY FOR MULTICORE

- Amdahl's Law
- Gustafson-Barsis's Law
- Amdahl's law and Multicore Processors

2 TIME- AND MEMORY-INDUCED SCALABILITY

- Fixed-size scalability for multicore processors
- Fixed-time scalability for multicore processors
- Memory bound scalability for multicore processors

3 CONCURRENCY-INDUCED SCALABILITY

- The model
- Implications on design of multiprocessors
- Conclusions

A VERY "OLD" LAW: 1967

- Described (only informally!) by Gene Amdahl (IBM) in a 3 page papers of 1967 "Validity of the single processor approach to achieving large scale computing capabilities"
- The "validity of single processor" is described against the supporters of the parallel organization of computers (with parallel memories, connected by a bus or point-to-point, with parallel execution streams)
- The basic idea is that:
 - 1 the "data management housekeeping", that is inherently sequential, cannot be parallelized (and therefore improved) on a parallel computer (no matter how many resources it has)
 - 2 any "geometrically related" problem, given the irregularity of shapes/regions/etc. cannot be mapped onto a regular geometry of components

SPEEDUP ACCORDING TO AMDAHL

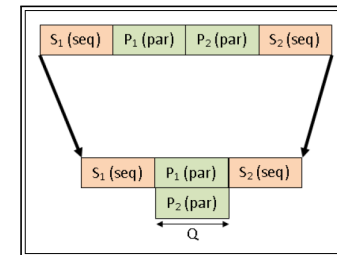
- The speedup S of X is the ratio between the sequential time to execute X and the parallel time (n processors) to execute X
- Let P be the part of X that can be parallelized
 - with n processors the parallel part takes time $\frac{P}{n}$ while the sequential takes time S
- Then the speedup is:

Amdahl's law

By executing X on n processors, the speedup is:

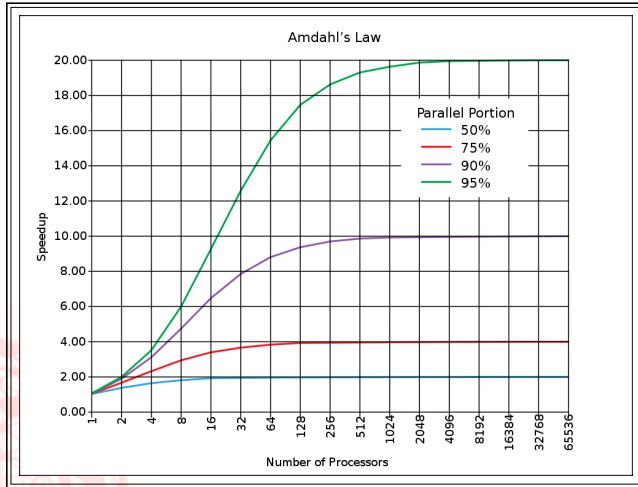
$$\text{Speedup}_A = \frac{S + P}{S + \frac{P}{n}}$$

AN EXAMPLE



- Assume a fixed program X with a sequential S and a parallelizable part P
- $n = 2$, $S = S_1 + S_2$, $P = P_1 + P_2$, $S_i = P_i$ fixed
- $\text{Speedup}_A(X) = \frac{S+P}{S+P/2} = 4/3$

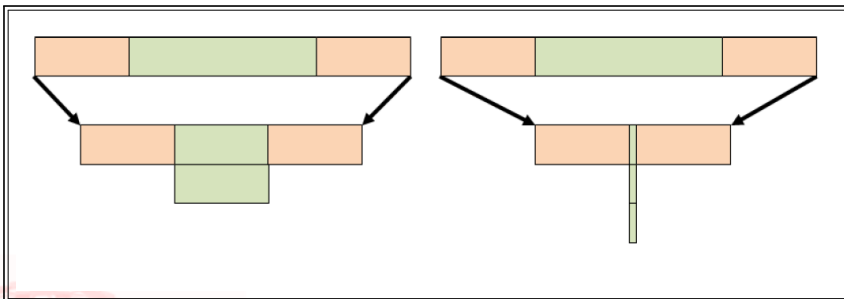
A DIAGRAM OF THE LAW



INTERPRETATION AND COMMENTS

- Amdahl's law indicates that the sequential part of a program will slow down any speedup that we can hope to get from parallelization
- $\lim_{n \rightarrow \infty} \text{Speedup}_A = \lim_{n \rightarrow \infty} \frac{S+P}{S+\frac{P}{n}} = \frac{S+P}{S}$
- If we set the Amdahl's coefficient $\alpha = S/(S + P)$, speedup is bounded by $1/\alpha$
- Law of diminishing return on investment
- \Rightarrow it is not enough to buy/invest in new hardware, but the sequential part must be negligible with respect to the parallel part (*good for us, Computer Scientists!*)

THE SITUATION AT THE LIMIT



PLAN

- PARALLELISM AND SCALABILITY FOR MULTICORE
 - Amdahl's Law
 - Gustafson-Barsis's Law
 - Amdahl's law and Multicore Processors
- TIME- AND MEMORY-INDUCED SCALABILITY
 - Fixed-size scalability for multicore processors
 - Fixed-time scalability for multicore processors
 - Memory bound scalability for multicore processors
- CONCURRENCY-INDUCED SCALABILITY
 - The model
 - Implications on design of multiprocessors
 - Conclusions

LAWS AND REALITY

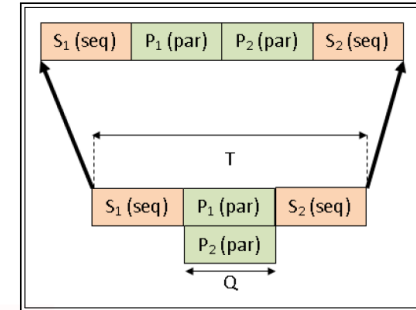
- 1988: Gustafson writes his (and his team) experience:

"Reevaluating Amdahl's Law" (Comm. of ACM, 1988)

The steepness of Amdahl's law graph when $S \rightarrow 0$ for $N = 1024$ implies that very few problems will experience even a 100-fold speedup. Yet, for 3 applications ($S = 0.4 - 0.8$ percent) we experience speedups between 1016 and 1021.

- The criticism: "One does not take a fixed-size problem and run it on various numbers of processors (except in academic research)" 😊
- You should assume run time constant and not the problem size

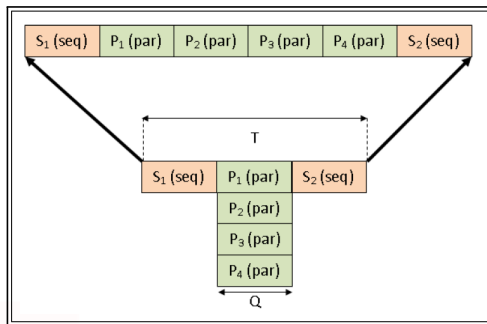
THE IDEA - 1



- Consider programs with a sequential part S , fixed, and a **fixed time frame**, $T = S + Q$.
- The speedup obtained by X is

$$\text{Speedup}_G(X) = \frac{S + 2Q}{S + Q} = 4/3$$

BUT, YOU CAN GET MORE "MILEAGE" . . .



- The speedup obtained by X is

$$\text{Speedup}_G(X) = \frac{S + 4Q}{S + Q} = 6/3$$

THE FORMULA

- Consider programs, with a sequential part S , fixed, and a **fixed time frame**, $T = S + Q$.
- Then the speedup by using n processors according to Gustafson (and Barsis) is:

Gustafson-Barsis' law

By executing X on n processors, the speedup is:

$$\text{Speedup}_G = \frac{S + nQ}{S + Q}$$

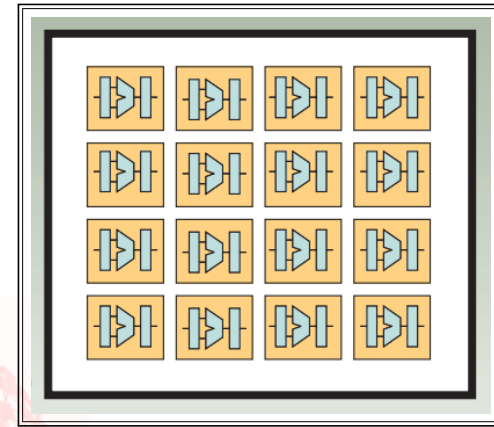
- With $n \rightarrow \infty$ the speedup is unbounded!

DESIGNER DILEMMA

- Assume that technological constraints bound the number of transistors on a multicore chip
- The dilemma: how to organize them? Many cores of small capacity or fewer cores of large capacity?
- The model: assume that on a chip you can place at most n Base Core Equivalents (BCE) of computational power 1
- Area of r BCEs can be used to obtain a processor with performances $perf(r)$
- In general, $perf(r)$ is sublinear; usually $perf(r) = \sqrt{r}$ (Pollack rule)
- "Amdahl's Law in the Multicore Era, M.D. Hill, M.R. Marty, IEEE Computer 41(7), Nov. 2008.



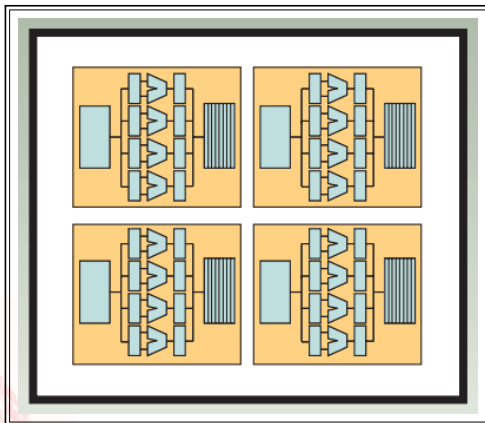
SYMMETRIC MULTICORE HIGHLY PARALLEL



PicoChip, Connex Machine, Tiler (TILE64)



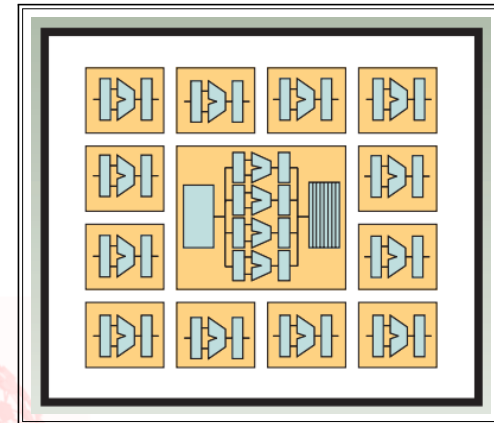
SYMMETRIC MULTICORE LOWLY PARALLEL



Intel, AMD



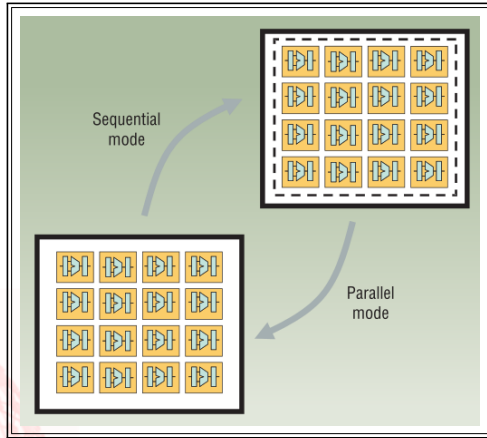
ASYMMETRIC MULTICORE



IBM/Sony Cell, Intel IXP



DYNAMIC MULTICORE



PLAN

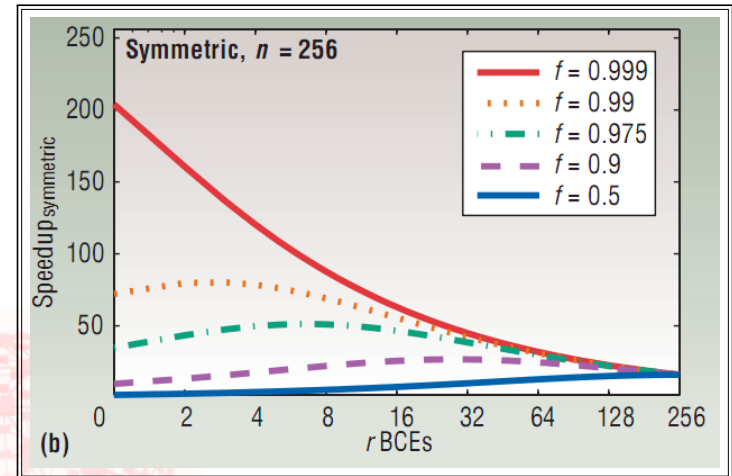
- 1 PARALLELISM AND SCALABILITY FOR MULTICORE
 - Amdahl's Law
 - Gustafson-Barsis's Law
 - Amdahl's law and Multicore Processors
- 2 TIME- AND MEMORY-INDUCED SCALABILITY
 - Fixed-size scalability for multicore processors
 - Fixed-time scalability for multicore processors
 - Memory bound scalability for multicore processors
- 3 CONCURRENCY-INDUCED SCALABILITY
 - The model
 - Implications on design of multiprocessors
 - Conclusions

AMDAHL'S LAW FOR SYMMETRIC MULTICORE

- Speedup depends on the parallelizable part of the program, f , by resources n on chip (in BCEs) and by resources dedicated to each core (r BCE)
- There are n/r core, each with performance $perf(r) = \sqrt{r}$
- Amdahl's law, in this case, is

$$Speedup_{symm}(f, n, r) = \frac{1}{\frac{1-f}{perf(r)} + \frac{f}{perf(r) \cdot n/r}}$$

PERFORMANCES: SYMMETRIC MULTICORE



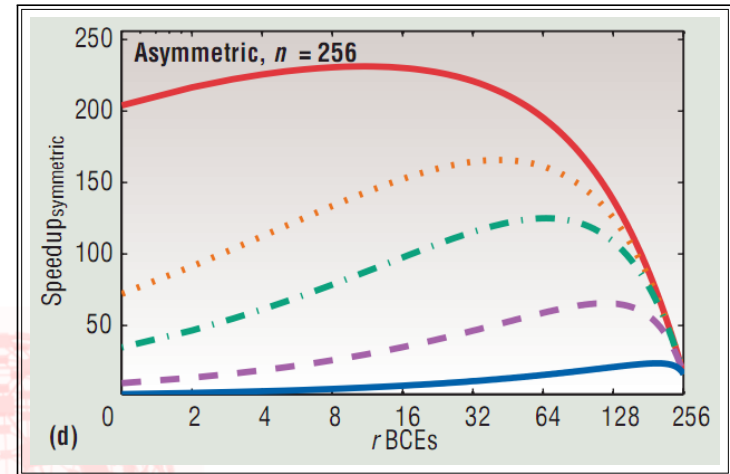
AMDAHL'S LAW FOR ASYMMETRIC MULTICORE

- Speedup depends on the parallelizable part of the program, f , by resources n on chip (in BCEs) and by resources dedicated to each core (r BCE)
- For the asymmetric multicore, a processor has more resources (r) and there are $n - r$ core with 1 BCE each
- In total $1 + n - r$ core, with different performances
- In the sequential part of the program, we can use the largest (r BCEs) core.
- In the parallel part, we can use all the cores (each with its performance)
- Amdahl's law, in this case, is:

$$\text{Speedup}_{\text{asymm}}(f, n, r) = \frac{1}{\frac{1-f}{\text{perf}(r)} + \frac{f}{\text{perf}(r)+n-r}}$$



PERFORMANCES: ASYMMETRIC MULTICORE



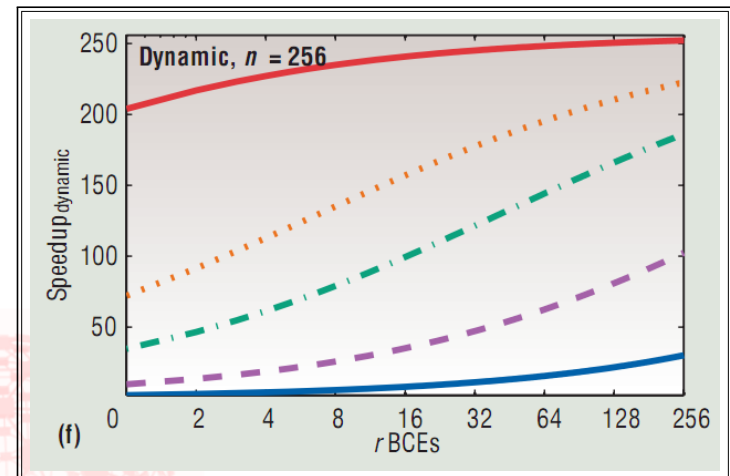
AMDAHL'S LAW IN THE DYNAMIC CASE

- Speedup depends on the parallelizable part of the program, f , by resources n on chip (in BCEs) and by resources dedicated to each core (r BCE)
- If it is possible to exploit each core (with multithread, for example) then each processor can be both a single processor (with r BCEs), in sequential, and n processors with 1 BCE of processing power
- In the sequential part, we can use the "largest" core (r BCE)
- In the parallel part, we can use all the n cores
- Amdahl's law, in this case, is:

$$\text{Speedup}_{\text{asymm}}(f, n, r) = \frac{1}{\frac{1-f}{\text{perf}(r)} + \frac{f}{n}}$$



PERFORMANCES: DYNAMIC MULTICORE



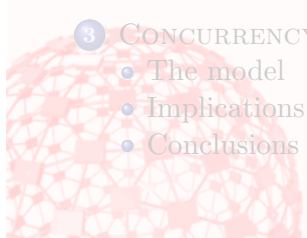
COMMENTS

- Software is not infinitely parallel/sequential
- Data movements and tasks add overhead
- Scheduling on asymmetric/dynamic can be more costly than on symmetric
- "Learning curve" for programmers
 - More costly to develop parallel software than sequential software
 - With asymmetric, double (at least) the number of platform to develop software on



PLAN

- 1 PARALLELISM AND SCALABILITY FOR MULTICORE
 - Amdahl's Law
 - Gustafson-Barsis's Law
 - Amdahl's law and Multicore Processors
- 2 TIME- AND MEMORY-INDUCED SCALABILITY
 - Fixed-size scalability for multicore processors
 - Fixed-time scalability for multicore processors
 - Memory bound scalability for multicore processors
- 3 CONCURRENCY-INDUCED SCALABILITY
 - The model
 - Implications on design of multiprocessors
 - Conclusions



THE CONCLUSIONS OF HILL-MARTHY PAPER

Pessimists will bemoan our model's simplicity and lament that much of the design space we explore can't be built with known techniques. We charge you, the reader, to develop better models, and, more importantly, to invent new software and hardware designs that realize the speedup potentials this article displays. Moreover, research leaders should temper the current pendulum swing from the past's underemphasis on parallel research to a future with too little sequential research. To help you get started, we provide slides

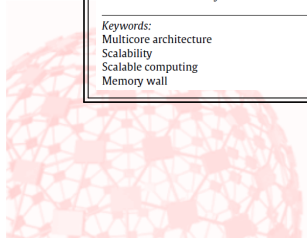


THE ANSWER 😊

Reevaluating Amdahl's law in the multicore era

Xian-He Sun*, Yong Chen
Computer Science Department, Illinois Institute of Technology, United States

<p>ARTICLE INFO</p> <p><small>Article history:</small> Received 5 December 2008 Received in revised form 22 March 2009 Accepted 9 May 2009 Available online 22 May 2009</p> <p><small>Keywords:</small> Multicore architecture Scalability Scalable computing Memory wall</p>	<p>ABSTRACT</p> <p>Microprocessor architecture has entered the multicore era. Recently, Hill and Marty presented a pessimistic view of multicore scalability. Their analysis was based on Amdahl's law (i.e. fixed-workload condition) and challenged readers to develop better models. In this study, we analyze multicore scalability under fixed-time and memory-bound conditions and from the data access (memory wall) perspective. We use the same hardware cost model of multicore chips used by Hill and Marty, but achieve very different and more optimistic performance models. These models show that there is no inherent, immovable upper bound on the scalability of multicore architectures. These results complement existing studies and demonstrate that multicore architectures are capable of extensive scalability.</p> <p style="text-align: right;"><small>© 2009 Elsevier Inc. All rights reserved.</small></p>
--	---



HISTORY SEEMS REPEATING ITSELF

- Worried of the pessimistic implications of Amdahl's law, computer vendors were concentrated thirty years ago on machines with few powerful processors
- After Gustafson-Barsi's law, many vendors built successfully large parallel machines (nCUBE, Thinking Machines) ...
- Actually, parallel machines do enroll tens of thousands of processors
- The view of the authors: the Hill-Marthy analysis of manycore processors has the same pitfalls of Amdahl's law (fixed size vs. fixed time) and, more important, there is the *memory wall* problem

33/78



SPEEDUP ON FIXED-TIME - 1

- Amdahl's law assumes the the problem size (the *workload*) is fixed when run on the unenhanced (sequential system) (fixed-size)
- Gustafson-Barsi's law argues that powerful machines are designer for larger problems and the workload should be scaled up (fixed-time)

$$\text{Speedup}_{FT} = \frac{\text{Seq time of solving Scaled Workload}}{\text{Par time of solving Scaled Workload}}$$

- Let w = original workload and let w' = the scaled workload be such that they finish in the same time with 1 processor (sequential) and with m processors

34/78



SPEEDUP ON FIXED-TIME - 2

- Assume scaling only occurs on the parallel part
- Then:

$$w' = (1 - f)w + fmw$$

$$\text{Speedup}_{FT} = \frac{\text{Seq time of solving } w'}{\text{Par time of solving } w'}$$

$$\text{Speedup}_{FT} = \frac{\text{Seq time of solving } w'}{\text{Seq time of solving } w}$$

$$\text{Speedup}_{FT} = \frac{w'}{w} = \frac{(1 - f)w + fmw}{w} = (1 - f) + mf$$

- ... that is the Gustafson-Barsi's law: assuming that the workload can increase, building large scale parallel systems is beneficial

35/78



THE MEMORY WALL - 1

- Not always the workload can increase without bounds: the physical constraint of memory is often "hit"
- Let w^* be the scaled workload under memory space constraints

$$\text{Speedup}_{MB} = \frac{\text{Seq time of solving } w^*}{\text{Par time of solving } w^*}$$

- Assume each node is a processor-memory pair: each node has a memory of size M (L1-cache)
- Let $g(x)$ be the increase of the parallel part of the workload as the memory capacity increases by x
- That is: $w = g(M)$ and the parallel part is $g(m \cdot M) = g(m \cdot g^{-1}(w))$

36/78



THE MEMORY WALL - 2


- Let w = original workload and let w^* = scaled workload be such that they finish in the same time resp. with 1 processor (sequential) and with m processors

$$\text{Speedup}_{MB} = \frac{\text{Seq time of solving } w^*}{\text{Par time of solving } w^*}$$

$$\text{Speedup}_{MB} = \frac{(1-f)w + f \cdot g(m \cdot g^{-1}(w))}{(1-f)w + \frac{f \cdot g(m \cdot g^{-1}(w))}{m}}$$

- If $g(x)$ is a power function $g(x) = ax^b$ we can write:

$$g(mx) = a(mx)^b = m^b \cdot ax^b = m^b g(x) = \bar{g}(m)g(x)$$
 where $\bar{g}(m)$ is the power function with coefficient 1.
- By assuming $g(x)$ polynomial, approximated by its highest degree term, we have:

37/78 

THE MEMORY WALL - 3


- We have:

$$\text{Speedup}_{MB} = \frac{(1-f)w + f \cdot \bar{g}(m)w}{(1-f)w + \frac{f \cdot \bar{g}(m)w}{m}}$$

- By simplifying the w we get the final version:

$$\text{Speedup}_{MB} = \frac{(1-f) + f \cdot \bar{g}(m)}{(1-f) + \frac{f \cdot \bar{g}(m)}{m}}$$

- This is called the *Sun-Ni's law* for memory bounded scalability.
- When $\bar{g}(m) = 1$ it becomes the Amdahl's law ...
- When $\bar{g}(m) = m$ it becomes the Gustafson-Barsi law

38/78 

AN EXAMPLE: MATRIX MULTIPLICATIONS


- Two matrices $N \times N$ to multiply
- Computation needed: $2N^3$
- Memory needed: $x = 3N^2 \Rightarrow x/3 = N^2$

$$g(x) = 2 \left(\sqrt{x/3} \right) 3 = cx^{3/2}$$

$$\Rightarrow \bar{g}(x) = x^{3/2}$$


- The speedup for matrix multiplication is:

$$\text{Speedup}_{MB} = \frac{(1-f) + f \cdot \bar{g}(m)}{(1-f) + \frac{f \cdot \bar{g}(m)}{m}} = \frac{(1-f) + f \cdot m^{3/2}}{(1-f) + f \cdot m^{1/2}}$$

39/78 

PLAN

- PARALLELISM AND SCALABILITY FOR MULTICORE
 - Amdahl's Law
 - Gustafson-Barsi's Law
 - Amdahl's law and Multicore Processors
- TIME- AND MEMORY-INDUCED SCALABILITY
 - Fixed-size scalability for multicore processors
 - Fixed-time scalability for multicore processors
 - Memory bound scalability for multicore processors
- CONCURRENCY-INDUCED SCALABILITY
 - The model
 - Implications on design of multiprocessors
 - Conclusions

40/78 

A VIEW OF HILL-MARTHY ANALYSIS - 1

- Fixed-size workload assumption (like Amdahl's law)
- The speedup is the ratio of the performances (enhanced over the original)
 - since the performances are the reciprocal of execution time

$$\text{Speedup} = \frac{T_{original}}{T_{enhanced}}$$

- If w is the workload, the $T_{original} = w/perf(1) = w$
- The performances over an n -BCE multicore is:

$$T_{enhanced} = \frac{(1-f)w}{perf(r)} + \frac{fw}{n/r \cdot perf(r)}$$

A VIEW OF HILL-MARTHY ANALYSIS - 2

- Then the speedup is:

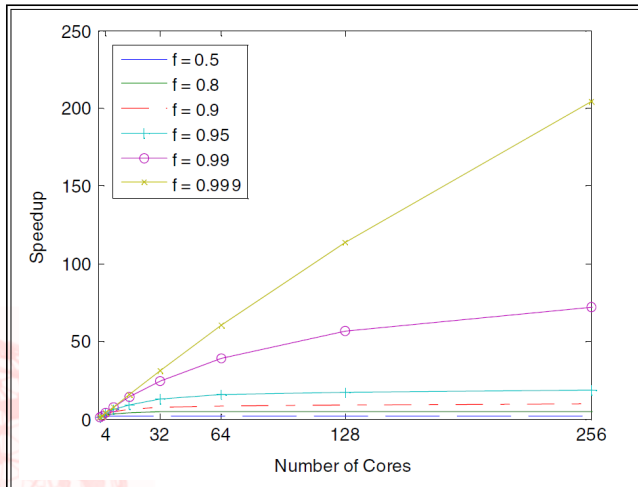
$$\text{Speedup} = \frac{T_{original}}{T_{enhanced}} = \frac{w}{\frac{(1-f)w}{perf(r)} + \frac{fw}{n/r \cdot perf(r)}}$$

$$\text{Speedup} = \frac{1}{\frac{(1-f)}{perf(r)} + \frac{f}{n/r \cdot perf(r)}}$$

- Since $perf(r)$ is a constant c (given a fixed design) and $m = n/r$ we can write "Amdahl's style"

$$\text{Speedup} = \frac{c}{(1-f) + \frac{f}{m}}$$

THE RESULT OF FIXED-SIZE SCALABILITY



PLAN

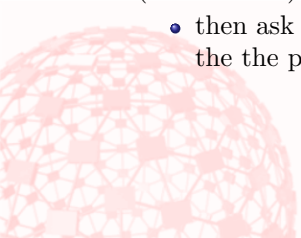
- 1 PARALLELISM AND SCALABILITY FOR MULTICORE
 - Amdahl's Law
 - Gustafson-Barsis's Law
 - Amdahl's law and Multicore Processors
- 2 TIME- AND MEMORY-INDUCED SCALABILITY
 - Fixed-size scalability for multicore processors
 - Fixed-time scalability for multicore processors
 - Memory bound scalability for multicore processors
- 3 CONCURRENCY-INDUCED SCALABILITY
 - The model
 - Implications on design of multiprocessors
 - Conclusions

HILL-MARTHY ANALYSIS IN FIXED-TIME - 1

- Assume symmetric architecture, each core with its L1 cache
- The Hill-Marthy law is:

$$\text{Speedup} = \frac{1}{\frac{1-f}{\text{perf}(r)} + \frac{f \cdot r}{\text{perf}(r) \cdot n}}$$

- Now, to measure the fixed-time scalability, we need to fix the *initial* point (let $n = r$) and the *scaled* number of cores (let $n = mr$)
 - then ask that the sequential time at initial point be equal to the the parallel time at the scaled point



HILL-MARTHY ANALYSIS IN FIXED-TIME - 2

- $$\frac{(1-f)w}{\text{perf}(r)} + \frac{fw}{\text{perf}(r)} = \frac{(1-f)w}{\text{perf}(r)} + \frac{fw'}{\text{perf}(r) \cdot m}$$

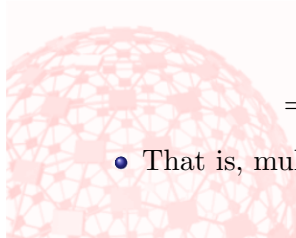
- Hence, we obtain that $w' = mw$

$$\text{Speedup}_{FT} = \frac{\text{Seq time of solving } w'}{\text{Par time of solving } w'} = \frac{\text{Seq time of solving } w'}{\text{Seq time of solving } w}$$

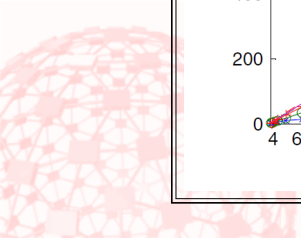
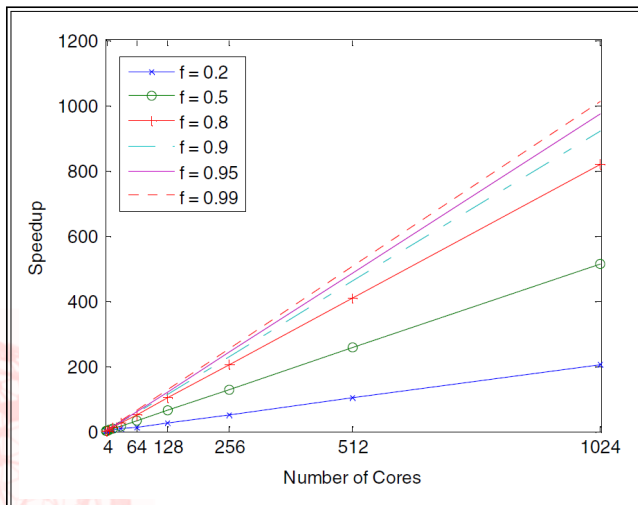
$$\text{Speedup}_{FT} = \frac{\frac{(1-f)w}{\text{perf}(r)} + \frac{fw'}{\text{perf}(r)}}{w/\text{perf}(r)} =$$

$$= \frac{(1-f)w + f(mw)}{w} = (1-f) + fm$$

- That is, multicore are scalable

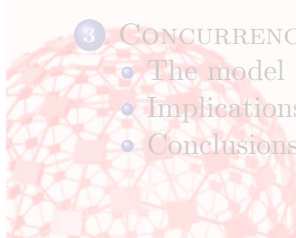


THE RESULT OF FIXED-TIME SCALABILITY



PLAN

- 1 PARALLELISM AND SCALABILITY FOR MULTICORE
 - Amdahl's Law
 - Gustafson-Barsis's Law
 - Amdahl's law and Multicore Processors
- 2 TIME- AND MEMORY-INDUCED SCALABILITY
 - Fixed-size scalability for multicore processors
 - Fixed-time scalability for multicore processors
 - Memory bound scalability for multicore processors
- 3 CONCURRENCY-INDUCED SCALABILITY
 - The model
 - Implications on design of multiprocessors
 - Conclusions



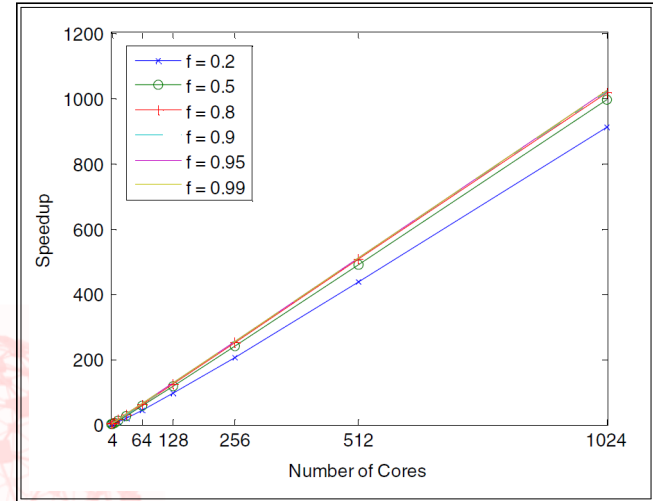
MEMORY BOUND SCALABILITY

- Let w^* be the scaled workload under memory capacity constraint
- Assume we scale from r cores to mr cores
-

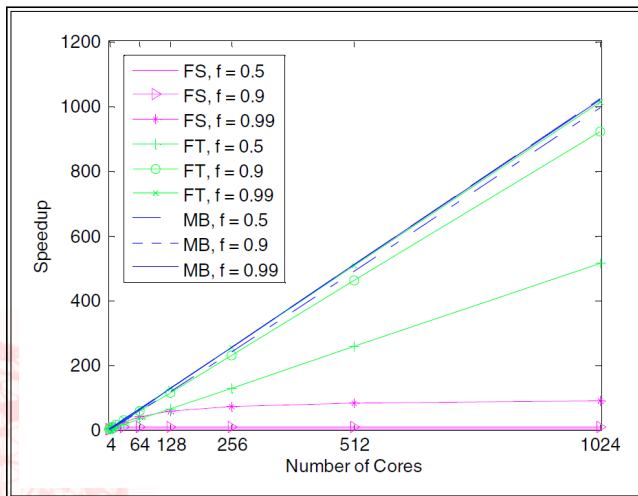
$$\begin{aligned} \text{Speedup}_{MB} &= \frac{\text{Seq time of solving } w'}{\text{Par time of solving } w^*} \\ &= \frac{\frac{(1-f)w}{\text{perf}(r)} + \frac{fw^*}{\text{perf}(r)}}{\frac{(1-f)w}{\text{perf}(r)} + \frac{fw^*}{m \cdot \text{perf}(r)}} = \\ &= \frac{(1-f)w + fw^*}{(1-f)w + \frac{fw^*}{m}} \\ &= \frac{(1-f)w + f \cdot \bar{g}(m)w}{(1-f)w + \frac{f \cdot \bar{g}(m)w}{m}} = \frac{(1-f) + f \cdot \bar{g}(m)}{(1-f) + \frac{f \cdot \bar{g}(m)}{m}} \end{aligned}$$



THE RESULT OF MEMORY-BOUND SCALABILITY



SUMMARY OF SCALABILITY



CONCLUSIONS OF THE AUTHORS

- Hill-Marthy's work is a corollary of Amdahl's law and applies only if users don't increase their computing demands when given more computing power
- Multicore processors **are** scalable, with fixed-time and memory-bound scalability models



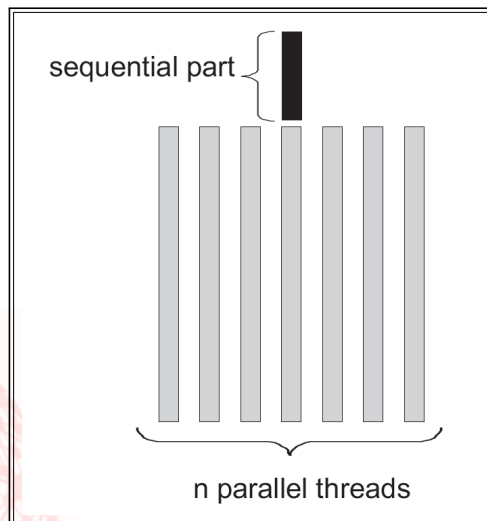
PLAN

- 1 PARALLELISM AND SCALABILITY FOR MULTICORE
 - Amdahl's Law
 - Gustafson-Barsis's Law
 - Amdahl's law and Multicore Processors
- 2 TIME- AND MEMORY-INDUCED SCALABILITY
 - Fixed-size scalability for multicore processors
 - Fixed-time scalability for multicore processors
 - Memory bound scalability for multicore processors
- 3 CONCURRENCY-INDUCED SCALABILITY
 - The model
 - Implications on design of multiprocessors
 - Conclusions

MOTIVATIONS

- As we have seen, Amdahl's law (and its derivations) have a significant impact on the design of multicore processors
 - "How many/large should be the core(s) in a multicore?"
- The software model in Amdahl's law is simple: either sequential or completely parallel code
- Something is missing in this view: the cost of synchronization through critical sections in the parallel part of the code
- By adding this factor, the analysis will show that the parallelism is also fundamentally limited by the synchronization

AMDAHL'S LAW SOFTWARE MODEL



└ The model

PLAN

- 1 PARALLELISM AND SCALABILITY FOR MULTICORE
 - Amdahl's Law
 - Gustafson-Barsis's Law
 - Amdahl's law and Multicore Processors
- 2 TIME- AND MEMORY-INDUCED SCALABILITY
 - Fixed-size scalability for multicore processors
 - Fixed-time scalability for multicore processors
 - Memory bound scalability for multicore processors
- 3 CONCURRENCY-INDUCED SCALABILITY
 - The model
 - Implications on design of multiprocessors
 - Conclusions

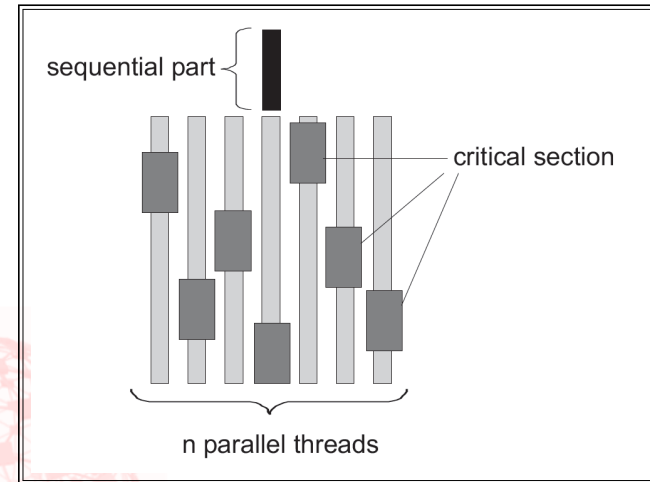
MODELING CRITICAL SECTIONS

- Let us split the program execution in three parts:
 - f_{seq} is the sequential part: it cannot be parallelized (thread spawning, data distribution, merging data from threads, etc.)
 - $f_{par,ncs}$ is the parallelizable part: infinitely parallelizable since it does not need synchronization (i.e. outside critical sections)
 - $f_{par,cs}$ is the parallelizable part that requires synchronization: if two or more threads compete for entering, their execution is serialized
- Of course, $f_{seq} + f_{par,ncs} + f_{par,cs} = 1$

57/78



INCLUDING THE CRITICAL SECTIONS



58/78



EXECUTION TIME: TWO CASES

- The evaluation is distinguished in two cases:
 - the total time can be approximated by the average per-thread execution time
 - the threads all execute equally long: the average makes "sense" as a measure
 - the total time is determined by a slowest thread
 - one of the thread is consistently slower than others

59/78



(1) AVG THREAD: THE MODEL - 1

- The sequential part does not scale with parallelism: $T_{seq} \propto f_{seq}$
- The parallel fraction that lies outside the critical section is parallelizable: $T_{par,ncs} \propto \frac{f_{par,ncs}}{n}$
- We need to compute $T_{par,cs}$
- Define the *probability for a critical section* during the parallel execution P_{cs}

$$P_{cs} = Pr\{\text{critical section}|\text{parallel}\} = \frac{f_{par,cs}}{f_{par,ncs} + f_{par,cs}}$$

60/78



(1) AVG THREAD: THE MODEL - 2

- The probability for i threads out of n threads to be in CS is binomially distributed (random entering into CS) hence:

$$Pr\{i \text{ of } n \text{ threads in critical section}\} = \binom{n}{i} P_{cs}^i (1 - P_{cs})^{n-i}$$

- We define the *Contention probability*, P_{ctn} as the probability that 2 (or more) critical sections will contend the same resource (lock, transaction, etc.)
- Now, assuming that the CSs contend randomly, the probability that j threads contend for the same CS is binomially distributed:

$Pr\{j \text{ of } i \text{ threads contend for the same CS}\} =$

$$\binom{i}{j} P_{ctn}^j (1 - P_{ctn})^{i-j}$$

61/78



(1) AVG THREAD: THE MODEL - 3

- We need to compute the average time spent in CS
- Let one thread T be in the critical section, the probability of other i out of remaining $n - 1$ to be in CS is given by:

$$\binom{n-1}{i} P_{cs}^i (1 - P_{cs})^{n-1-i}$$

- The probability that j out of these i will contend is given by:

$$\binom{i}{j} P_{ctn}^j (1 - P_{ctn})^{i-j}$$

- If other j thread contend the CS with T , $j + 1$ thread will be serialized ...
- ... and it will take $\frac{(j+1)f_{par,cs}}{n}$ to execute the critical section
 - $\frac{jf_{par,cs}}{n}$ to execute the j critical sections (and T is waiting)
 - $\frac{f_{par,cs}}{n}$ to execute the CS of T .

62/78



(1) AVG THREAD: THE MODEL - 4

In summary, the average time spent in a critical section is:

- by averaging over all possible i
- of the probability that i other threads are in a CS ...
- times the probability that j of the i threads contend with T
- times the time spent in to execute the (delayed because of synchronization) Critical Section

$$T_{par,cs} \propto \sum_{i=0}^{n-1} \binom{n-1}{i} P_{cs}^i (1 - P_{cs})^{n-1-i} \cdot \sum_{j=0}^i \binom{i}{j} P_{ctn}^j (1 - P_{ctn})^{i-j} \cdot \frac{j+1}{n} f_{par,cs}$$

63/78



(1) AVG THREAD: THE MODEL - 5

$$T_{par,cs} \propto \sum_{i=0}^{n-1} \binom{n-1}{i} P_{cs}^i (1 - P_{cs})^{n-1-i} \cdot \sum_{j=0}^i \binom{i}{j} P_{ctn}^j (1 - P_{ctn})^{i-j} \cdot \frac{j+1}{n} f_{par,cs}$$

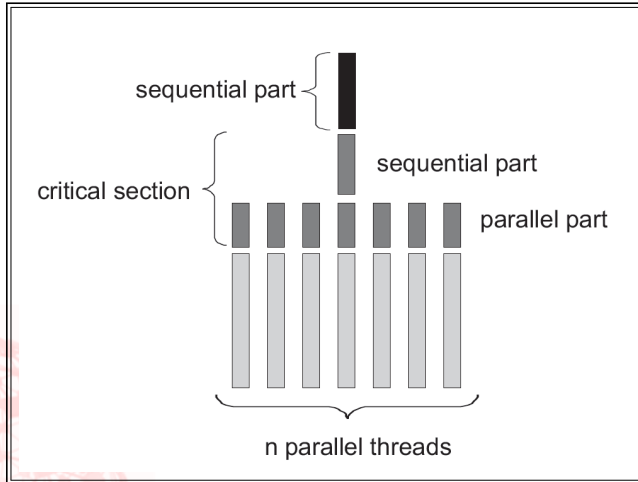
- By exploiting the property that $\sum_{i=0}^n \binom{n}{i} P^i (1 + P)^{n-i} = nP$ (avg of a binomial distribution), we obtain:

$$T_{par,cs} \propto f_{par,cs} \cdot \left(P_{cs} P_{ctn} + (1 - P_{cs} P_{ctn}) \cdot \frac{1}{n} \right)$$

64/78



MODELING AS A SEQUENTIAL PLUS A PARALLEL PART

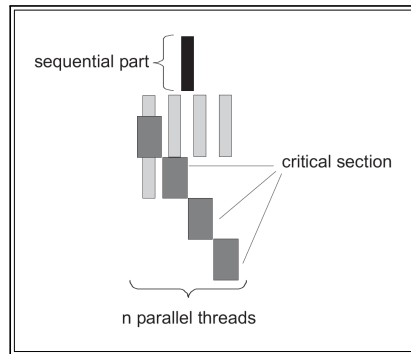


EXECUTION TIME: TWO CASES

- The evaluation is distinguished in two cases:
 - the total time can be approximated by the average per-thread execution time
 - the threads all execute equally long: the average makes "sense" as a measure
 - the total time is determined by a slowest thread
 - one of the thread is consistently slower than others

(2) SLOWEST THREAD: THE MODEL - 1

- The assumption that the time of execution of parallel part is dominated by the average per-thread execution time does not always hold
- In case the contention is high, the total time of execution may be dominated by one slowest thread.



(2) SLOWEST THREAD: THE MODEL - 2

- Assume that the execution of the slowest thread is determined by the serialized execution of all the contending critical sections
- The average time of the slowest thread is influenced by:
 - sequential part: $T_{seq} \propto f_{seq}$
 - time spent in contending: $T_{par,cs} \propto f_{par,cs} P_{ctn}$
 - $T_{par} \propto \frac{f_{par,cs}(1-P_{ctn}) + f_{par,ncs}}{2n}$
- Then, the average execution of the slowest thread is:

$$T \propto f_{seq} + f_{par,cs} P_{ctn} + \frac{f_{par,cs}(1 - P_{ctn}) + f_{par,ncs}}{2n}$$

INTEGRATION TO AMDAHL'S LAW

- Since the time of a parallel program is:

$$T \propto f_{seq} + \max \left\{ f_{par,cs} P_{cs} P_{ctn} + \frac{f_{par,cs}(1 - P_{cs} P_{ctn}) + f_{par,ncs}}{n}, \right. \\ \left. f_{par,cs} P_{ctn} + \frac{f_{par,cs}(1 - P_{ctn}) + f_{par,ncs}}{2n} \right\}$$

- Then, the Amdahl's law can be, asymptotically, reduced to:

$$\lim_{n \rightarrow \infty} S = \frac{1}{f_{seq} + \max \{ f_{par,cs} P_{cs} P_{ctn}, f_{par,cs} P_{ctn} \}} \\ = \frac{1}{f_{seq} + f_{par,cs} P_{ctn}}$$

PLAN

- 1 PARALLELISM AND SCALABILITY FOR MULTICORE
 - Amdahl's Law
 - Gustafson-Barsis's Law
 - Amdahl's law and Multicore Processors
- 2 TIME- AND MEMORY-INDUCED SCALABILITY
 - Fixed-size scalability for multicore processors
 - Fixed-time scalability for multicore processors
 - Memory bound scalability for multicore processors
- 3 CONCURRENCY-INDUCED SCALABILITY
 - The model
 - Implications on design of multiprocessors
 - Conclusions

ASYMMETRIC MULTICORE - 1

- One of the arguments in favour of asymmetric multicore processors is that the sequential part gets quickly executed by the (only) large core ...
- ... while the parallel part is spedup by the many (small) cores
- Critical sections do limit the potential benefit of a single large core: when contention occurs, sequential execution due to queuing on a resource is performed on the small cores, thereby inefficiently.

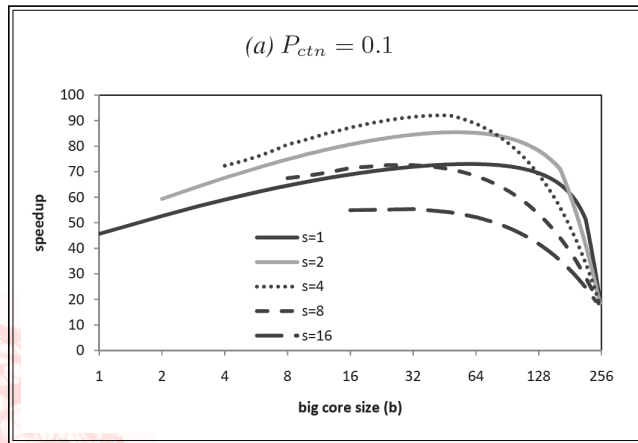
ASYMMETRIC MULTICORE - 2

- Assuming n cores with performances of 1, and one single core with performance p the time is:

$$T \propto \frac{f_{seq}}{p} + \max \left\{ f_{par,cs} P_{cs} P_{ctn} + \frac{f_{par,cs}(1 - P_{cs} P_{ctn}) + f_{par,ncs}}{n + p}, \right. \\ \left. f_{par,cs} P_{ctn} + \frac{f_{par,cs}(1 - P_{ctn}) + f_{par,ncs}}{2(n + p)} \right\}$$

- Then, small cores should be sufficiently large as not to impact negatively on the execution of the sequential part due to contention to CSs

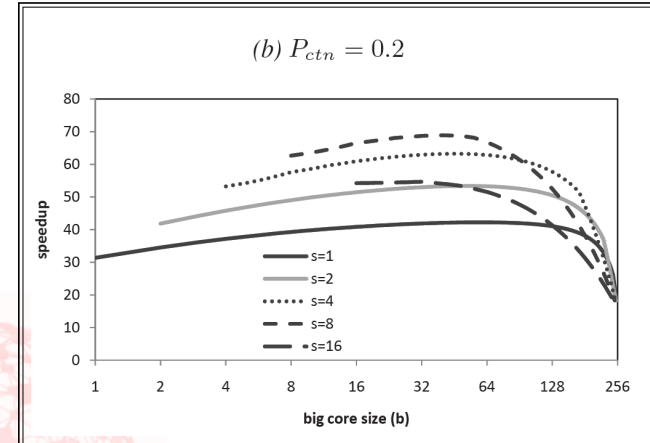
BIG (B) VS. SMALL (S) CORES - 1



73/78



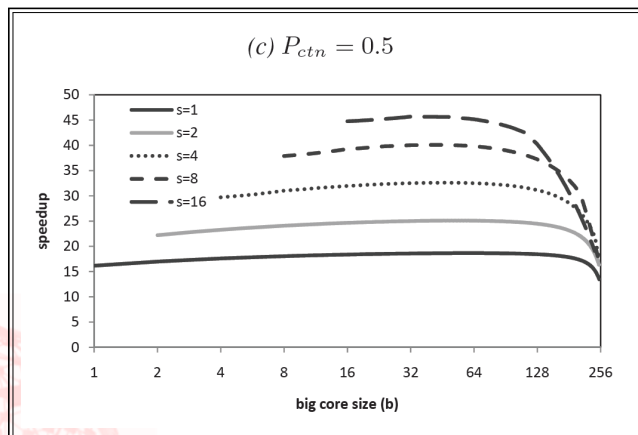
BIG (B) VS. SMALL (S) CORES - 2



74/78



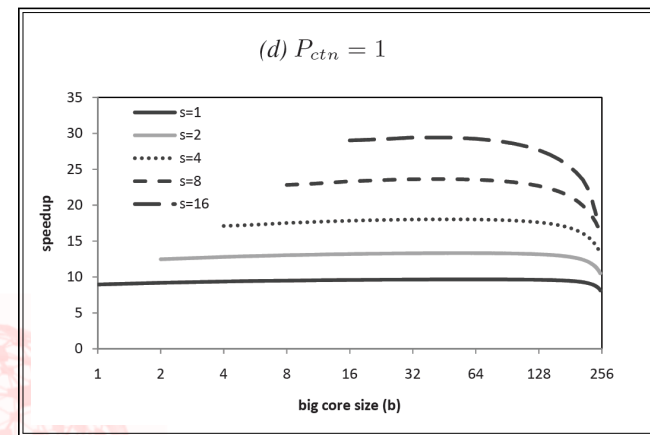
BIG (B) VS. SMALL (S) CORES - 3



75/78



BIG (B) VS. SMALL (S) CORES - 4



76/78



PLAN

- 1 PARALLELISM AND SCALABILITY FOR MULTICORE
 - Amdahl's Law
 - Gustafson-Barsis's Law
 - Amdahl's law and Multicore Processors
- 2 TIME- AND MEMORY-INDUCED SCALABILITY
 - Fixed-size scalability for multicore processors
 - Fixed-time scalability for multicore processors
 - Memory bound scalability for multicore processors
- 3 CONCURRENTLY-INDUCED SCALABILITY
 - The model
 - Implications on design of multiprocessors
 - Conclusions



BY CONSIDERING CONCURRENCY . . .

- . . . asymmetric multicore offer smaller performance benefits than expected
 - contention is executed on one of the small cores
- . . . many small cores in an asymmetric multicore may significantly bound the performances, when high contention rates are expected
- . . . techniques for accelerating critical sections are needed (such as executing CSs on the big core, or predicting which CSs will compete, . . .)
- Research still working on this . . .

