

Performance Evaluation

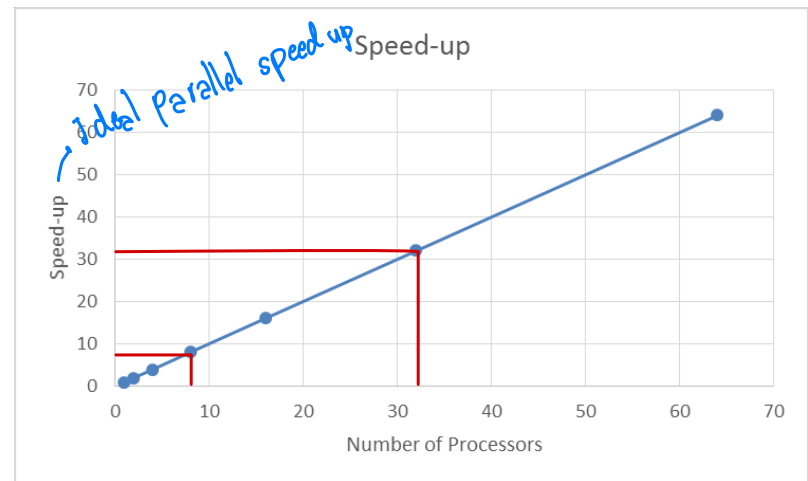
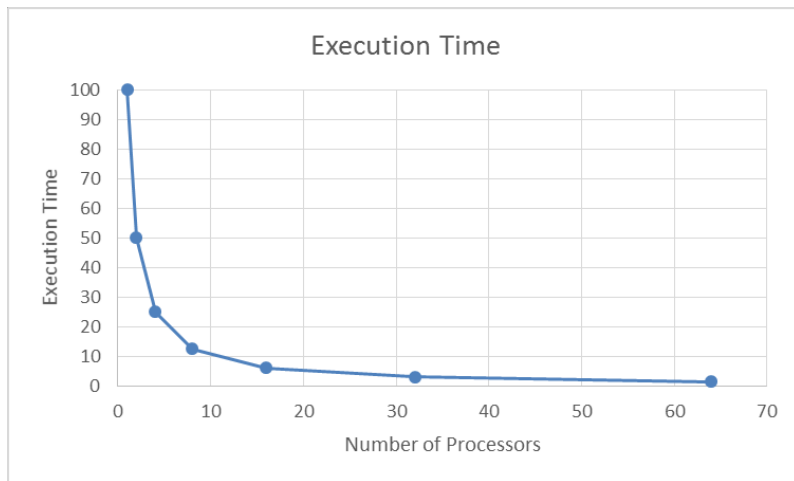
“The real problem is that programmers have spent far too much time worrying about efficiency in the wrong places and at the wrong times;

premature optimization is the root of all evil (or at least most of it) in programming” การพยายาม optimize เร็วเกินไป เป็นสาเหตุของความช้าร้ายในการเขียนโปรแกรม

--- Donald Knuth, 1974 Knuth บอกว่าเขียนโปรแกรมในลักษณะที่เร็วเกินไปจะทำให้โปรแกรมต้อง optimized กันแล้ว

Speed-up

- Speed-up determines how much faster parallel execution is versus serial execution. *ด.เร็วที่มากขึ้น เมื่อเทียบกับตอนยังไม่ปรับปรุงโปรแกรม*
- Speed-up = original run time / run time after change *เวลาของ code ที่ปรับปรุง / เวลา code เดิม*
- Parallel Speed-up = serial run time / parallel run time
- Ideal parallel speed-up = number of processors
maximum speed up ที่เป็นไปได้ของ Parallel computing = จ.น. processor เช่น ๕ 10 processors ด.เร็วขึ้น 10 เท่า



Note: When computing speed-up, the best and fastest serial algorithm & code must be used for measuring serial runtime. Do not use the parallel version with 1 thread.

Parallel Efficiency

- Parallel efficiency indicates how well software utilizes the computational resources of the system.
- $\text{Parallel efficiency} = \text{speed-up} / \text{number of processors}$
ประสิทธิภาพการใช้ทรัพยากรประมวลผลแบบขนาน อาจไม่ได้เท่ากับ Ideal parallel speed up เมื่อ parallel computing เติบโตขึ้น?
- Maximum efficiency = 1 or 100% เช่น 10 processor & speed up = 10 $= \frac{10}{10} = 1$
- Ex. 12x speedup on 16 cores $\rightarrow 75\%$ efficiency $\frac{12}{16} \times 100 = 75\%$

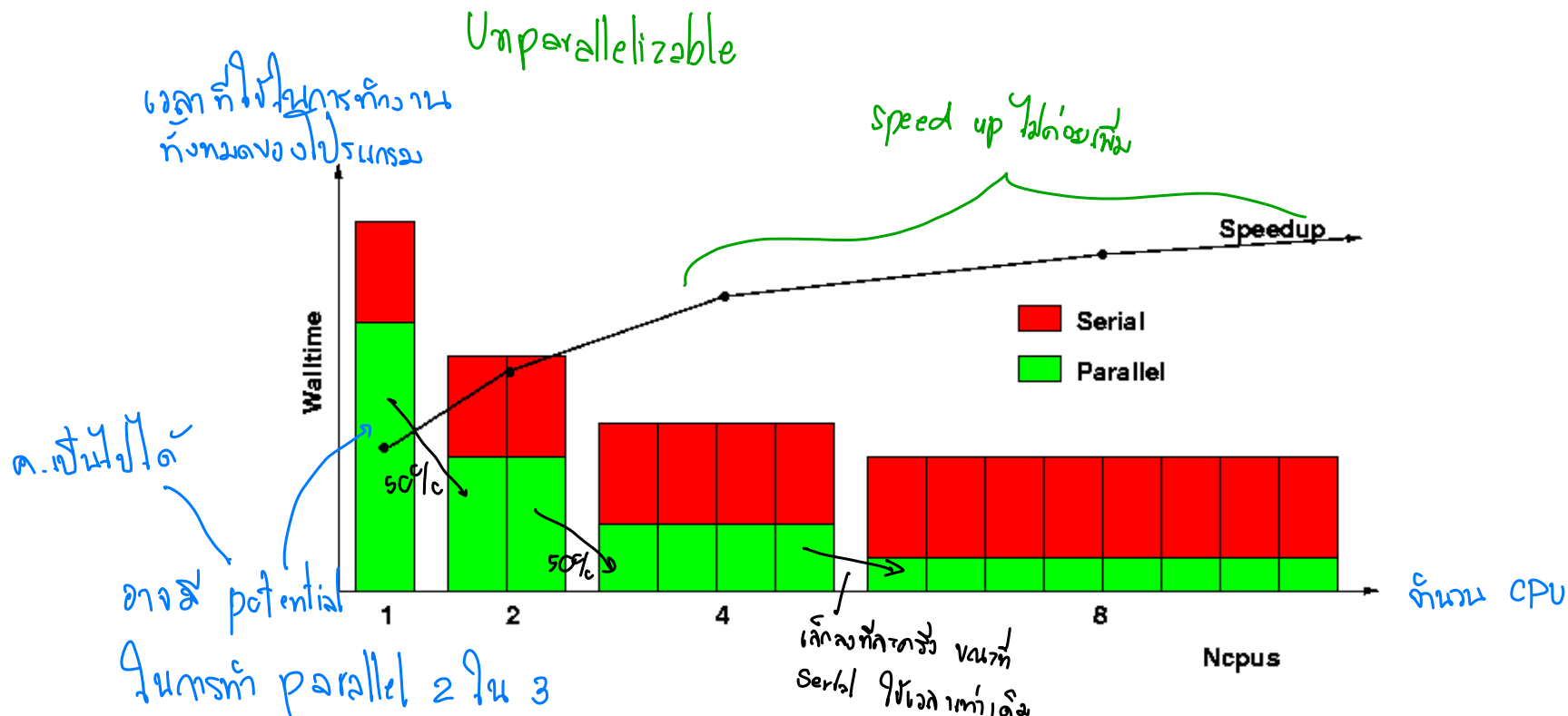
This means that, on average, each core is idle about 25% of the time.

Parallel Performance Limitation

สถานะที่ efficiency
ไม่ได้สมบูรณ์แบบ ideal

- Unparallelizable computation (inherently serial computation)
การคำนวณที่ทำงานไม่ได้แบบ parallel เพราะมี dependencies
- Parallel computing overhead
 - Thread creation & scheduling การสร้างและสลับการทำงานของ thread มี overhead
 - Communication & synchronization overhead การสื่อสารส่งข้อมูล
- Load balancing การกระจายงานที่ไม่เท่ากัน งานย่อยๆ มีขนาดเล็กลงไปอยู่ไม่เท่ากัน
CPU ที่รับงานเล็ก ๆ ทำเสร็จเร็วก็ไม่ได้ทำอะไร (idle)
เวลาที่ไถ่จริง ๆ ทำทำงานเสร็จเร็ว (processor ที่สุดทำงานเสร็จ งานในที่สุด)
เวลาแทนที่เวลาไถ่จริง มันดันไปทำที่ไม่เกี่ยวข้องกับการคำนวณ = overhead

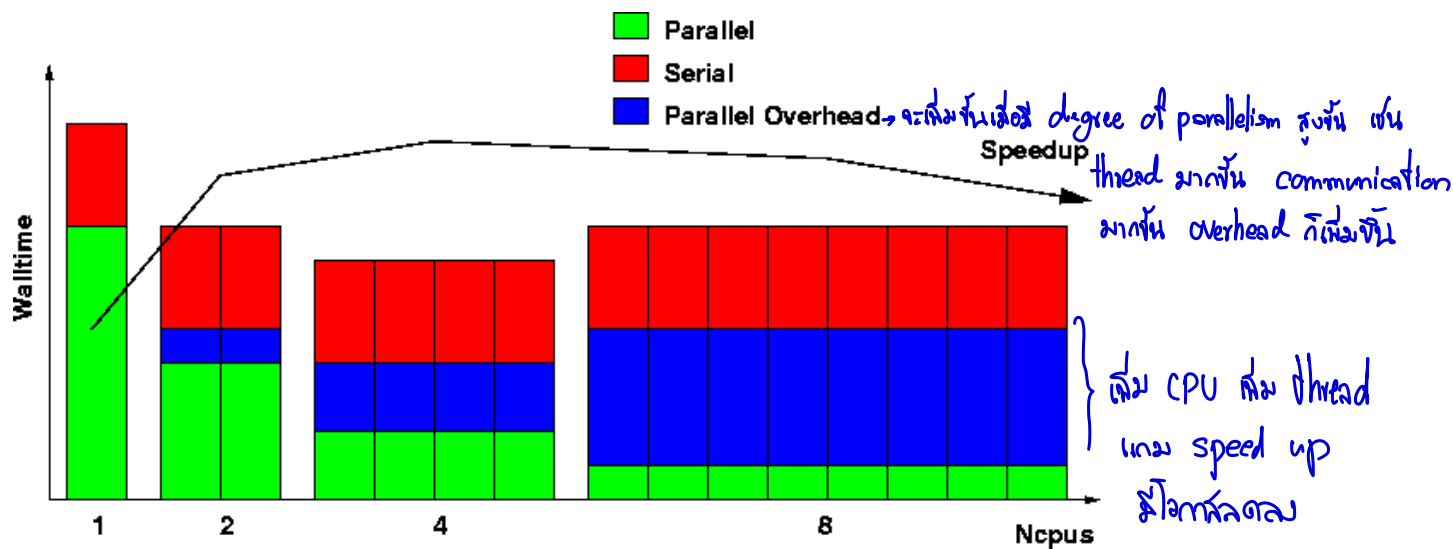
Realistic Speedup



- Inherently serial computation limits speedup gain.

Realistic Speedup

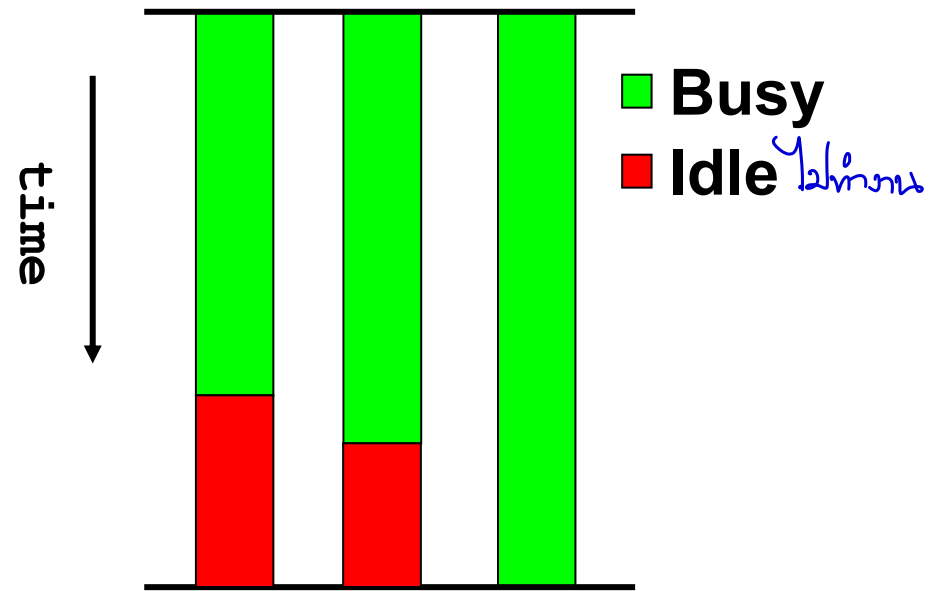
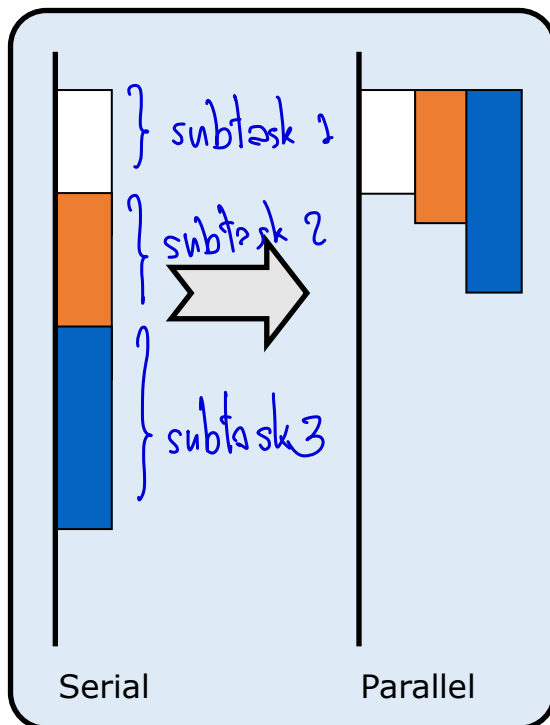
Overhead



- Overhead can even reduce speedup when adding more processors.

Load Imbalance

- Unequal work loads lead to idle processors and wasted time.



Amdahl's Law

พิจารณาการคำนวณ Unparallelizable code

Let s be fraction of sequential code, p be fraction of parallelizable code ($s + p = 1$)
 พิจารณาการคำนวณ speed up จากสัดส่วนที่โปรแกรมสามารถทำ parallel ได้มากน้อยแค่ไหน

• Speedup

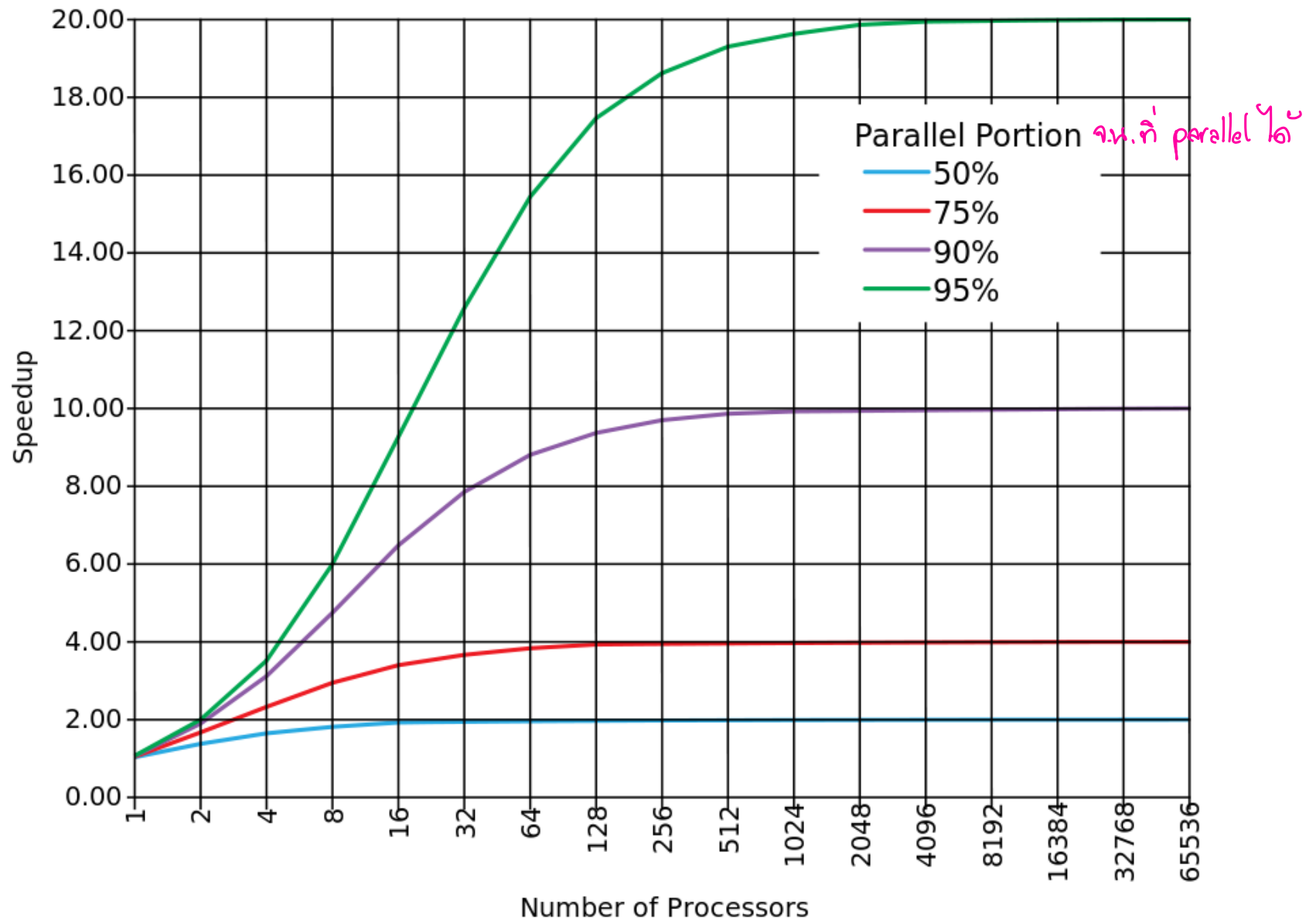
$$= \frac{s+p}{s+\frac{p}{N}}$$

sequential time
parallel time
a.n. processor

$$= \frac{1}{s+\frac{(1-s)}{N}}$$

- If N approaches infinite then speed-up approaches $1/s$ and efficiency approaches 0. จำนวน processor เพิ่มขึ้น ทำให้ speed up เข้าใกล้ $\frac{1}{s}$ และประสิทธิภาพเข้าใกล้ 0
- Amdahl's Law implies that parallel computing is only useful with small number of processors or the problem can be perfectly parallelized.

Amdahl's Law



Fix Time model Gustafson's Law

ธรรมชาติการเพิ่มงาน ส่วนที่ limit จริง ๆ ไม่ใช่เพิ่มงาน แต่เป็นเวล

Usually, we add more processors to solve larger problems.

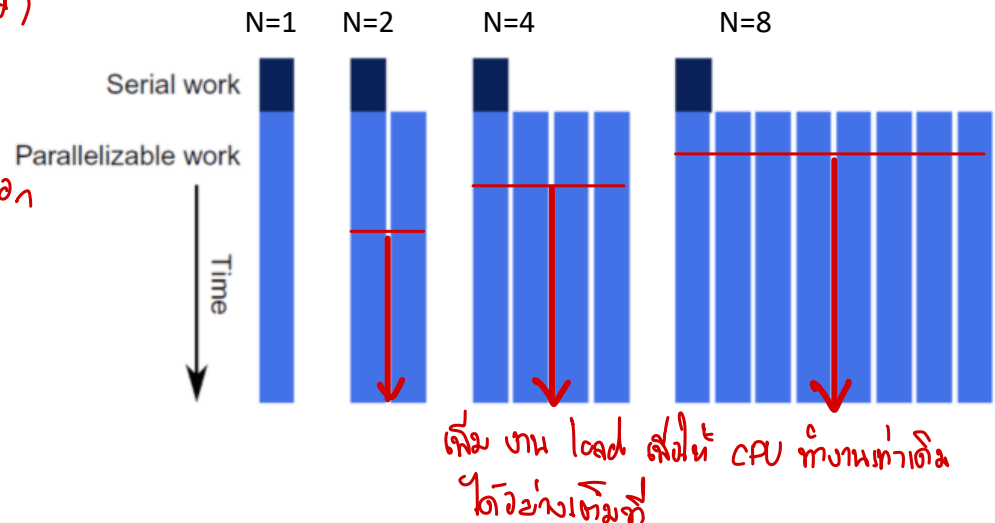
So, let the problem size scale with the number of processors and fix the time. แล้วเพิ่มงานให้ fit กับเวลาที่

Gustafson's speed-up measures how much more work can be done on a parallel machine compared to a sequential machine.

$$\text{Speedup} = \frac{s + \underbrace{pN}_{\substack{\text{เพิ่มงาน (vary load)} \\ \text{parallel ทุก CPU}}}}{\underbrace{s + p}_{= 1 \text{ ต่อตัวส่วนออก}}}$$

$$= s + (1 - s)N$$

$$\text{Efficiency} = \frac{s + (1 - s)N}{N}$$

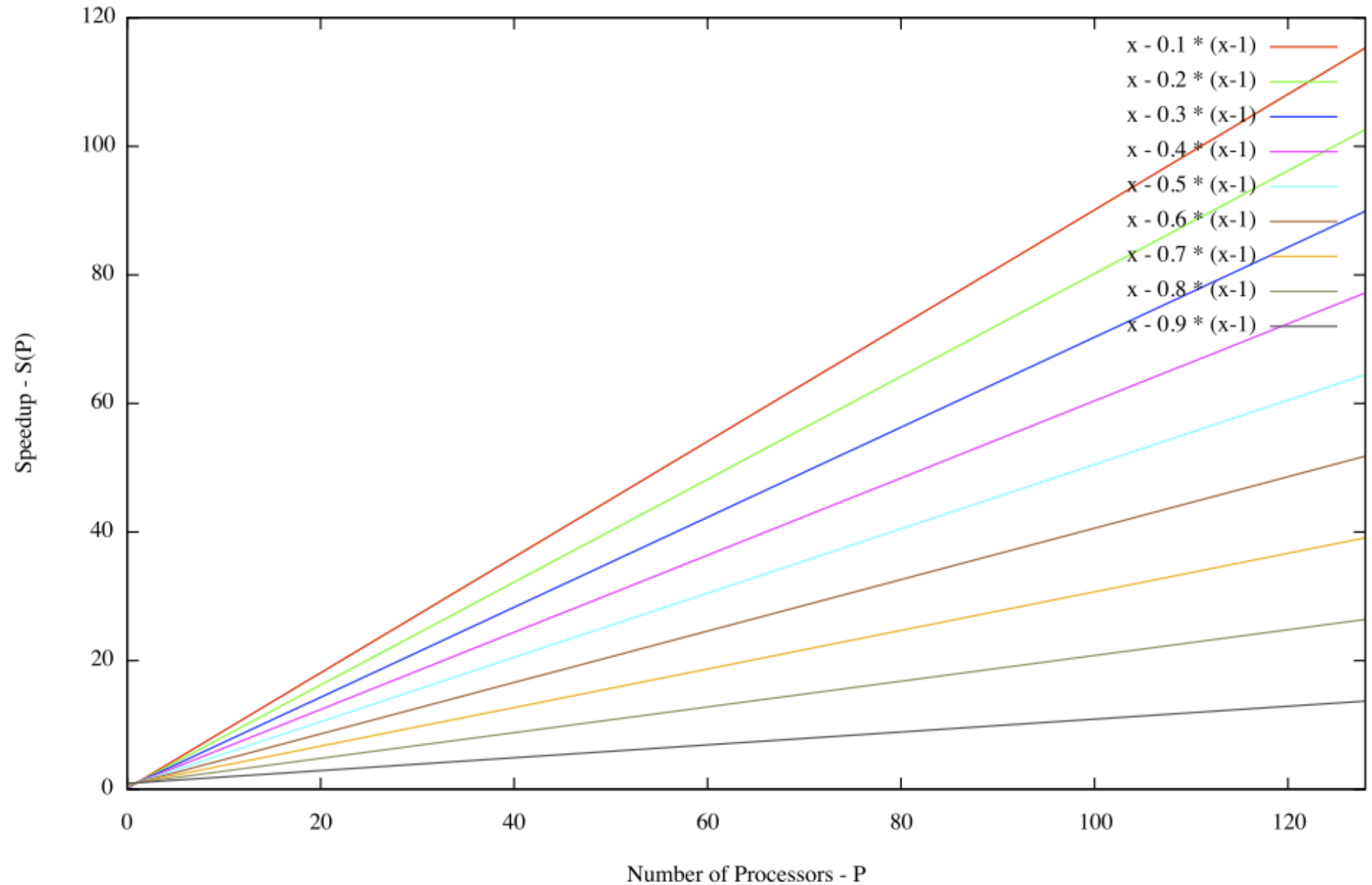


Scaled speed-up has no upper-bound.

If N is very large then efficiency approaches $(1 - s)$ or fraction of parallel code.

Gustafson's Law: $S(P) = P - a \cdot (P - 1)$

ไม่ limit เป็นเส้นตรงที่ความชันลดลง



Super-linear Speedup speed up > processor

- When speedup is greater than number of processors.
- Possible cases: สาเหตุ
 - Smaller subtasks fit better in cache lines/memory pages and therefore reduce misses. ขนาดของงานที่เล็กลงบางตัวใส่ลงใน cache ได้พอดี (ลด cache miss)
 - Parallel execution has more cache memory to use and share. ใช้ cache มากกว่า
 - Parallel algorithm reduce amount of computation, e.g. parallel search finds solution faster. algorithm บางตัวทำในโปรแกรมจบลงเร็วขึ้น
- Caution: Data sets that are too small – smaller than typical data set size – can give a false sense of performance improvement. ชุดข้อมูล เล็กเกินไป บางทีจะได้ผลลัพธ์ที่ mislead (เข้าใจผิด)

parallel search แบ่งข้อสอบ 10 ส่วนให้ 10 คนไป search แล้วปรากฏว่าคนที่ต้องกรอกอยู่ในหน้าเรา
ของข้อที่ 10 process ที่ทำงานข้อที่ 10 ก็จะสามารถทำงานเร็วกว่า 10 เท่า
(speed up > 10)

Performance Testing

- Load test
 - Test how system performs under expected load.
- Stress test
 - Test how system reacts under extreme load. At what point the system break and where it breaks.
- Endurance test
 - Test how system behaves when it is used at expected load for a long period of time.
 - Some problems appear only after long running time, e.g. memory leak, unclosed resources/connections.