

Parallelisme

- Race condition, starvation, dead lock, live lock.
→ solució: ordre de lock predeterminat (ex. alfabètic)

T_1 = temps lineal

T_{∞} = cost mínim

Parallelism = T_1 / T_{∞}

Parallel slackness = Parallelism / P_{min}

Efficiency $P = T_1 / (T_P \cdot P)$

T_P (temps amb P processadors)

Speedup $P = T_1 / T_P$

Strong scaling (canvi temps, canvi P , problema fixat)

Weak scaling (canvi temps, canvi P , problema canvia però se'n fa un per P fixat)

[Data sharing]

$T_{\text{access}} = t_s + m \cdot t_w$
↳ bytes

SMP (symmetric multi-processor)

- tots els processadors tenen accés igual a tota la memòria (UMA)

- un sol SO

- bottleneck: bus de memòria

- coherència de memòria

↳ snooping-based coherence protocol

- write-update (↑ tràfic al bus)

- write-invalidate (dirty = exclusive ownership)

MSI / MESI

↳ False sharing: dos processadors escriuen freqüentment a la mateixa línia de cache

- Consistència entre escriptures a diferents variables

NUMA

- "hub" entre cpu i memòria per assegurar coherència

- directory-based cache coherence protocol

Distributed-memory architectures

- OS level (on page fault)
- user level (message-passing library)

[Amdahl's Law]

ϵ - graus de temps que el programa s'executa en paral·lel

$T_P = (1 - \epsilon)T_1 + (\epsilon \cdot \frac{T_1}{P}) + \text{overhead}$

S_P amb $P \rightarrow \infty$: $\frac{1}{1 - \epsilon}$

• Pipelining, ILP (instruction level P)

TLP (thread level parallelism) (hyperthreading)

DLP (data level paral.) (SIMD)

• Memory cache

- Temporal locality (eg. loop)

- Spatial locality (eg. array)

network topology: direct / indirect (arbre) etc.
- latència (lineal/mesh)
- ample de banda