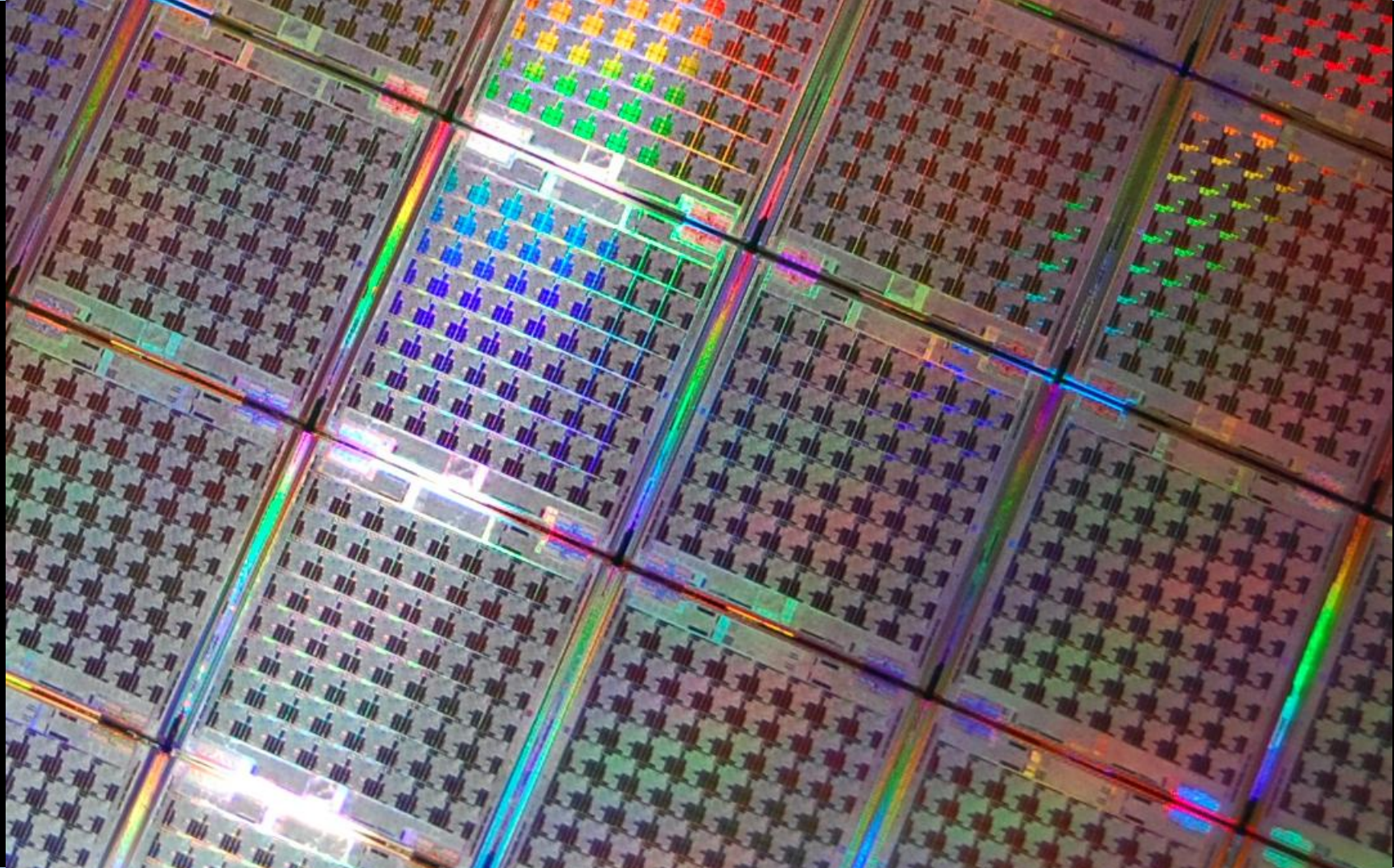


# ANÁLISE DE DESEMPENHO E ESCOLHA DINÂMICA DE ESCALONAMENTO PARA SISTEMAS MULTICORE

**EMILIO FRANCESQUINI, ALFREDO GOLDMAN**  
UNIVERSIDADE DE SÃO PAULO  
{EMILIO, GOLD}@IME.USP.BR

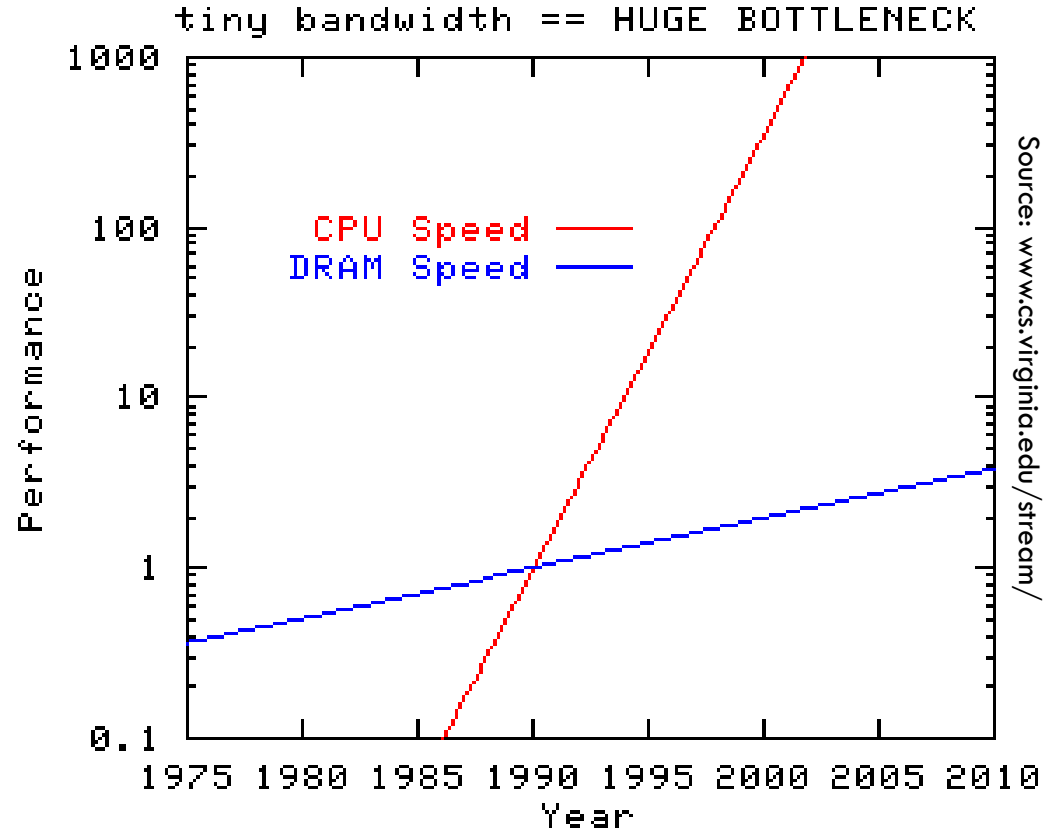
Julho de 2011



# Problemas

3

- Power Wall
- Memory Wall
- ILP Wall



NUMANode P#0 (16GB)

Socket P#0

L3 (24MB)

L2 (256KB) L2 (256KB) L2 (256KB) L2 (256KB) L2 (256KB) L2 (256KB) L2 (256KB) L2 (256KB)

L1 (32KB) L1 (32KB) L1 (32KB) L1 (32KB) L1 (32KB) L1 (32KB) L1 (32KB) L1 (32KB)

Core P#0 Core P#1 Core P#2 Core P#3 Core P#8 Core P#9 Core P#10 Core P#11

PU P#0 PU P#4 PU P#8 PU P#12 PU P#16 PU P#20 PU P#24 PU P#28

NUMANode P#1 (16GB)

Socket P#1

L3 (24MB)

L2 (256KB) L2 (256KB) L2 (256KB) L2 (256KB) L2 (256KB) L2 (256KB) L2 (256KB) L2 (256KB)

L1 (32KB) L1 (32KB) L1 (32KB) L1 (32KB) L1 (32KB) L1 (32KB) L1 (32KB) L1 (32KB)

Core P#0 Core P#1 Core P#2 Core P#3 Core P#8 Core P#9 Core P#10 Core P#11

PU P#1 PU P#5 PU P#9 PU P#13 PU P#17 PU P#21 PU P#25 PU P#29

idrouille  
-Intel(R) Xeon(R)  
Beckton X7560  
@ 2.27GHz

# Possíveis ações

5

- Definição de afinidade processo ↔ processador
- Colocação explícita de páginas de memória

# Fixação de processos

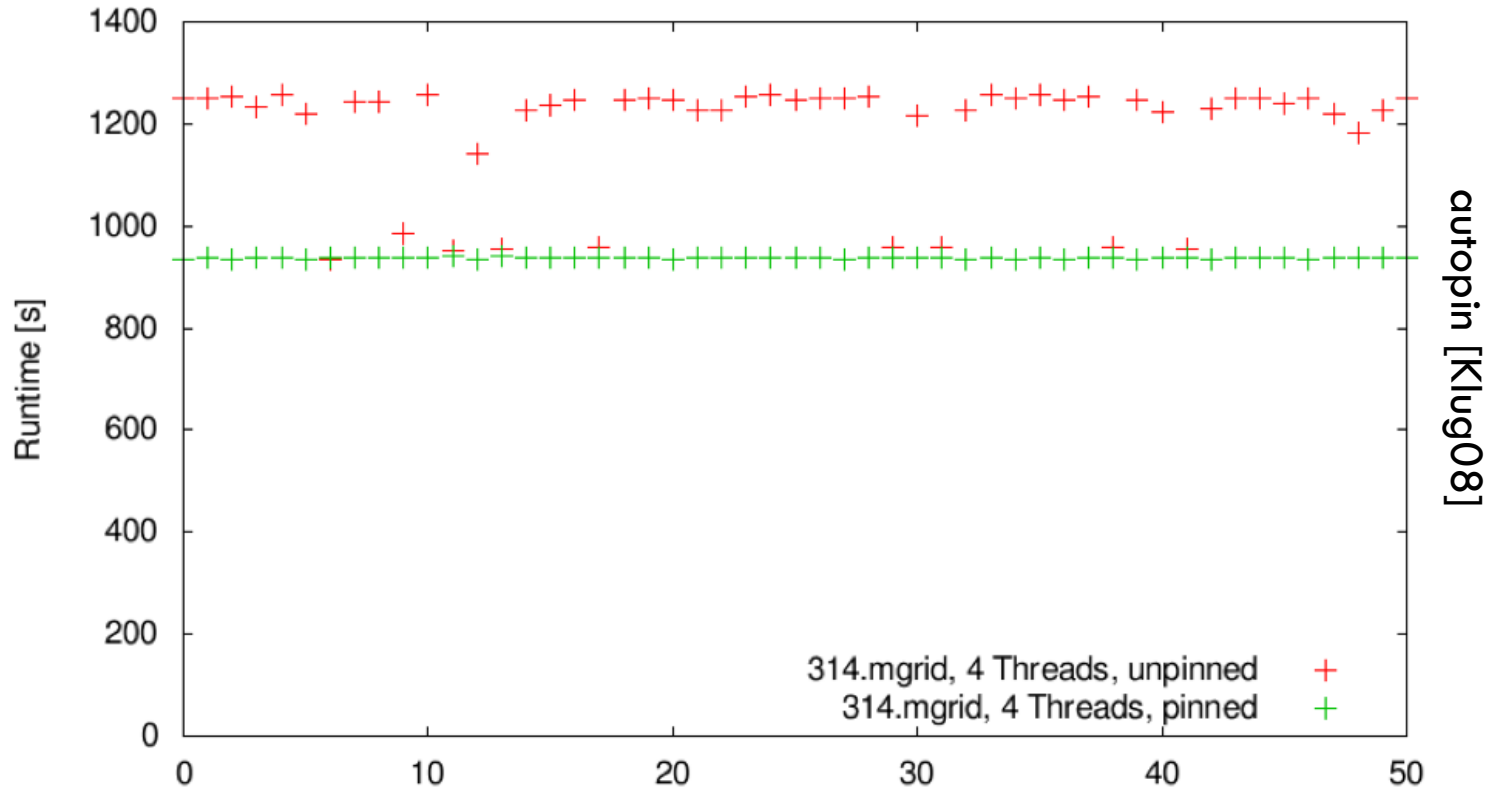
6

- Linpack em 16 cores

| Configuration                       | CPUs | Time (sec) | % Speedup |
|-------------------------------------|------|------------|-----------|
| No affinity control used            | 16   | 467.16     |           |
| taskset 0xf                         | 16   | 481.83     | -3.04%    |
| taskset 0x1; 0x2; 0x4; 0x8          | 16   | 430.44     | 8.53%     |
| -cpu_bind=map_cpu:0,1,2,3<br>-B 1:1 | 16   | 430.36     | 8.55%     |

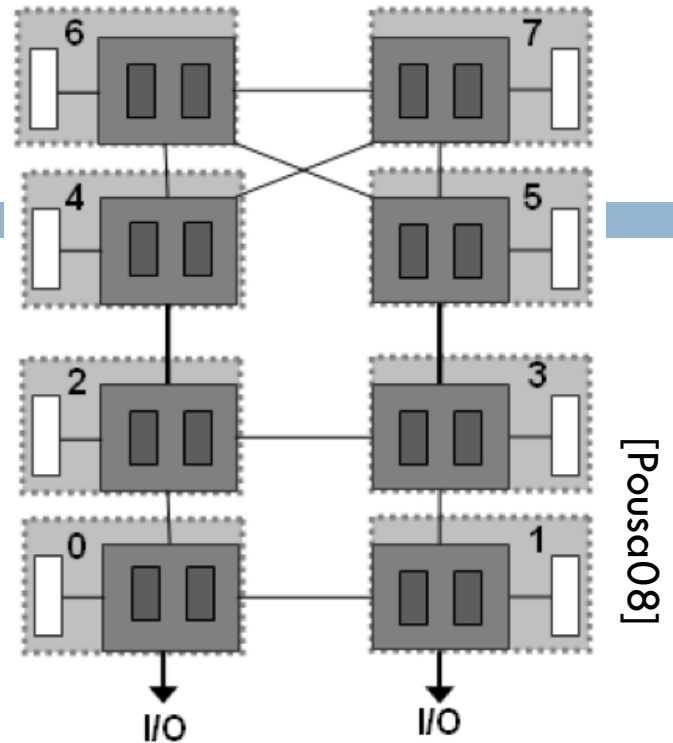
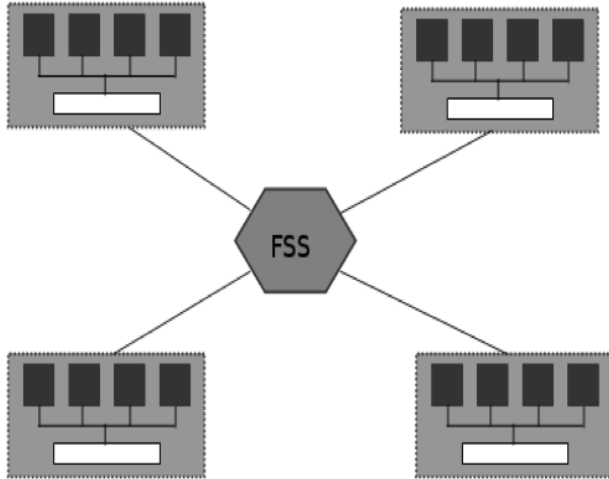
# Fixação de processos

7



# NUMA

8

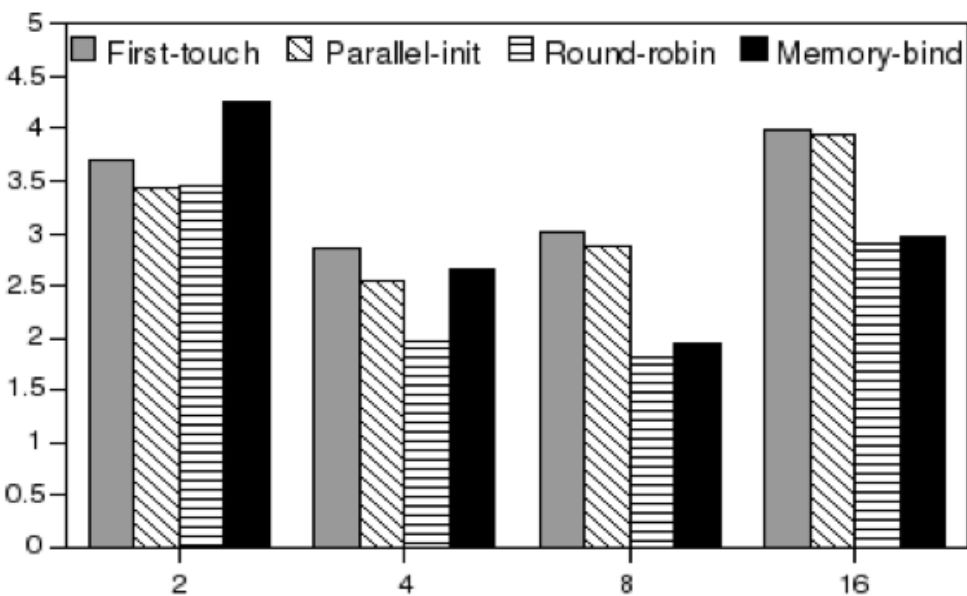


- ❑ Fator NUMA: 2  $\rightarrow$  2.5
- ❑ 16 Itanium2 @1.6 GHz
- ❑ 64 GBytes RAM

- ❑ NUMA factor: 1.2  $\rightarrow$  1.5
- ❑ 8 dual-core Opteron @ 2.2 GHz
- ❑ 32 Gbytes RAM

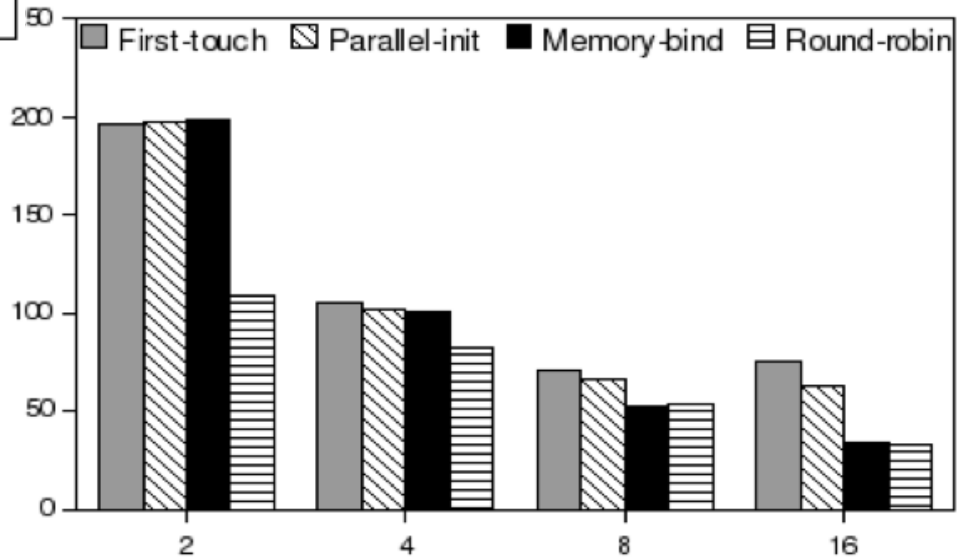


# NAS Benchmark



[Pousca08]

Itanium



Opteron

CG Kernel – CFD

Alto uso de memória

Padrão de acesso irregular

# Resumindo

10

- Ambientes heterogêneos
  - ▣ Alocações distintas levam a tempos de execução diferentes
    - Máquinas heterogêneas
      - Cada core pode possuir diferentes tempos de acesso à memória
    - Variações no tempo de comunicação entre cada um dos cores
      - Propriedades da interconexão são dinâmicas
        - Localização
        - Carga
        - Perfil da aplicação
        - ...
  - ▣ Tencionamos minimizar o tempo total de execução da aplicação

# Problemas com as soluções atuais

11

- Falta de escalabilidade
  - ▣ Teste e escolha manual de alocações
    - Para cada aplicação
    - Para cada máquina
- Falta de adaptabilidade dinâmica
  - ▣ Ao ambiente
  - ▣ Ao comportamento da aplicação

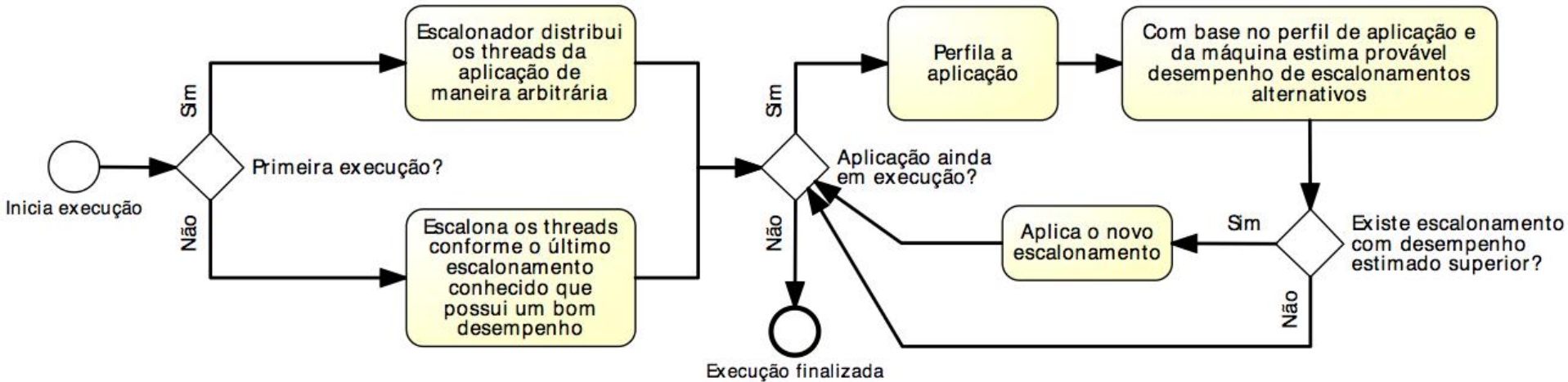
# Nossa proposta

12

- Perfilamento online e offline tanto das aplicações quanto do ambiente
- Escalonamento dinâmico de tarefas
- Colocação dinâmica de páginas de memória
- Proceder gradualmente para uma solução abrangente

# Nossa proposta (cont.)

13



# ○ Perfilamento

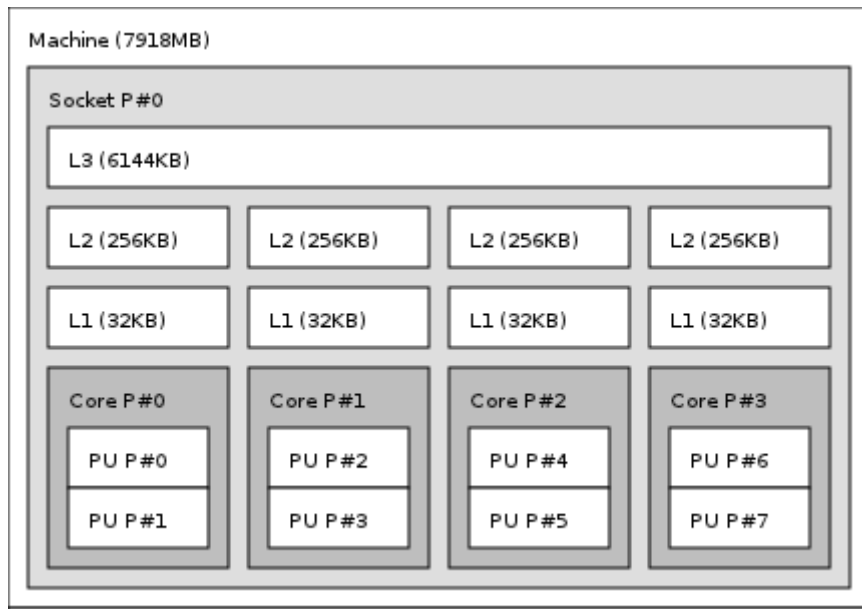
14

## PAPI

- Autopin [Klug08]
  - Retired instructions
- Exemplos
  - Level 1/2/3 data/instruction cache misses
  - Cache Line Invalidation (SMP)
  - Data/Instruction translation lookaside buffer misses
  - Integer/FP instructions executed, FLOPS



hwloc



# Outras ferramentas úteis

15

- Netperf
- OProfile 
- STREAM
- NAS Parallel Benchmarks

- SPEC Benchmarks



# Desafios de pesquisa

16

- Escalonamento de processos leves
  - ▣ Relacionamento entre os processos
    - Aplicação de BubbleSched
  - ▣ Fixação de processos aos cores
  - ▣ Modelo de atores
- Construção de um escalonador dinâmico e automático



# Referências

- [Balle2007] Enhancing an Open Source Resource Manager with Multi-Core/Multi-threaded Support, S. M. Balle and D. Palermo, Job Scheduling Strategies for Parallel Processing, 2007.
- [Pousa08] Christiane Pousa Ribeiro, Jean-Francois Mehaut, Alexandre Carissimi, Marcio Castro, and Luiz Gustavo Fernandes. Memory affinity for hierarchical shared memory multiprocessors. Computer Architecture and High Performance Computing, Symposium on, 0:59–66, 2009.
- [Klug08] Tobias Klug, Michael Ott, Josef Weidendorfer, and Carsten Trinitis: "autopin - Automated Optimization of Thread-to-Core Pinning on Multicore Systems" Transactions on High-Performance Embedded Architectures and Compilers, 3(4), 2008

Obrigado!

# Perguntas