

# CAD - Lista 03

Pedro Antonio Maciel Soraiva - 202303358

① - SMP apresenta processadores idênticos com memória compartilhada; acesso uniforme à memória; sistema operacional único e apresenta uma escalabilidade limite pelo nº de processadores. Em comparação com um uniprocessador, o SMP apresenta maior desempenho, escalabilidade incremental, mais simples que usar um cluster, alta disponibilidade em caso de um processador falhar e menor latência de comunicação.

② - Coerência de cache é como se ~~garantem~~ mantém a consistência de dados em múltiplos processadores que possuem cópias de dados em caches locais. A diferença de esquemas de hardware e software é no tempo de execução (hardware) e compilação (software); a abordagem de software é mais simples - evita dados compartilhados no cache, entretanto é menos eficiente no uso do cache. No esquema de hardware, há protocolos de diretório (controle de escrita e leitura), protocolos de monitoração (mantem a coesão de cache entre processadores), Protocolo MESI (Modified, Exclusive, Shared, Invalid)

③ Clusters são um sistema de computação baseado em nós independentes mas conectados por um middleware, seria como se nos juntássemos vários PCs em um só, pela rede. A vantagem é tolerância a falhas, fácil de escalar, balanceamento de carga e você consegue agregar recursos computacionais que seria mal utilizado.



④ - UMA: Unified Memory Access: Tempo de acesso uniforme como o nome diz, mais simples, menos escalável, geralmente usado SMP.

NUMA (Non UMA): Tempo variado de acesso a memória, depende da distância do processador, mais complexo, maior escalabilidade, geralmente usado em cluster.

CC-NUMA (Cache Coherence NUMA): Protocolos de software e hardware garante coerência de informações, escalabilidade de desempenho em níveis alta de paralelismo.

⑤ - SISD (Single Instruction, Single Data): Ele executa uma instrução em um dado; não apresenta paralelismo; usado em computadores de Von Neumann tradicionais.

SIMD (Single Instruction, Multiple Data): Um fluxo de instrução para múltiplos fluxos de dados; apresenta paralelismo (como a GPU), ideal para operações singulares em conjunto grande, como a soma de constante em grande array.

MISD (Multiple Instruction, Single Data): Vários fluxos de instrução para um fluxo de dados; não apresenta aplicações práticas.

~~(Multi)~~ MIMD (Multiple Instruction, Multiple Data): Vários fluxos de instruções e de dados; apresenta paralelismo; flexível, usado em clusters, sistemas distribuídos e multiprocessadores.

⑥ - Memória compartilhada: programação intuitiva e fácil de aplicar, com balanceamento dinâmico; problemas de escalabilidade limitada, problemas com sincronização e coerência de cache.

Memória distribuída: Mais complexo; mais escalável, tolerante a falha e com maior custo benefício; problema com



sobrecarga de comunicação e duplicação de dados.

⑦: Particionamento: Dividimos o problema em tarefas menores e executáveis para paralelismo

- Comunicação: Dois tarefas comunicam e conectam entre si, para minimizar a latência

- Aglomeração: Agrupamos tarefas que comunicam em "clusters" para serem processadas, a ideia é equilibrar o paralelismo e as comunicações entre atividades

- Mapeamento: Atribuir as aglomerações para processadores, considerando as comunicações entre processadores e topologias das máquinas, garantido uso total do hardware.

⑧: O particionamento pode ser decomposto em funcional (tarefas) baseado nas tarefas que devem ser executadas e funções que devem ser coordenadas entre si; pode ser definido um pipeline ou fluxo lógico de como deve ser executadas cada parte do problema. No paralelismo de dados, como deve ser executadas e comunicadas os dados em paralelismos massivos e escaláveis, como deve ser feita operações em cada dado e como serão sincronizadas no final

- Paralelismo de dados → divide os dados → quais devem ser processados e como → sincronização final como necessário

- Paralelismo de instrução → processo dividido em partes → como associadas a cada instrução → coordenação da execução dessas partes

9) Expressão:  $(X/K - Y * j * g + Z * g * h)$

constantes:  $K, j, g$  e  $h$

vetores:  $X, Y, Z$  (tamanho  $m$ )

#### • Paralelismo de Tarefas

Tarefa 1.1: calcula  $X/K$

Tarefa 1.2: Calcula  $Y * j * g$

Tarefa 1.3: Calcula  $Z * g * h$

Tarefa 2: soma os resultados das tarefas 1.1, 1.2 e 1.3

#### • Paralelismo de Dados:

1- Calculamos duas novas constantes  $a = j * g$  e  $b = g * h$

Vamos considerar os vetores, então podemos fazer operações pelo índice, da seguinte maneira

$$\text{Resultado}_i = X[i]/K + Y[i] * a + Z[i] * b$$

Que ao ser executado, teríamos algo como

$$\text{resultado} = \sum_{i=0}^{m-1} X[i]/K + Y[i] * a + Z[i] * b$$

E para paralelizar isso, cada Thread de uma GPU poderia executar uma operação com os elementos  $i$ . Se houver  $m-1$  threads, haveria uma resolução quase instantânea.