

PROCESSAMENTO PARALELO

www.inf.pucrs.br/~linatural/corporas/processamento/txt/Processamento_07_JAN32.txt

Na área de processamento paralelo existem dois paradigmas principais de programação, Memória Compartilhada e Troca de Mensagens.

Cada um deles é adequado a uma arquitetura de hardware específica.

No entanto, existem arquiteturas de multiprocessadores para as quais o mapeamento para um desses paradigmas não é tão simples. Clusters de SMP (Symmetric multiprocessing), por exemplo, são construídos com máquinas de memória compartilhada, conectadas através de uma rede de interconexão.

Symmetric multiprocessing (SMP) involves a symmetric multiprocessor system hardware and software architecture where two or more identical processors connect to a single, shared main memory, have full access to all I/O devices, and are controlled by a single operating system instance that treats all processors equally, reserving none for special purposes. Most multiprocessor systems today use an SMP architecture. In the case of multi-core processors, the SMP architecture applies to the cores, treating them as separate processors.

Aplicações para clusters de SMP podem ser programadas para utilizar troca de mensagens entre todos os processadores.

Mas existe a possibilidade de um melhor desempenho se utilizado um modelo híbrido de comunicação com troca de informações por memória compartilhada dentro do nó SMP e troca de informações por mensagens entre os nós.

Nesse trabalho foi desenvolvido e avaliado um modelo híbrido de programação para uma aplicação na área de engenharia mecânica baseada no método dos elementos finitos.

O objetivo desse trabalho é avaliar esse modelo e comparar seu desempenho com uma versão pura, por troca de mensagens, da aplicação.

Na área de processamento paralelo, vários processadores são utilizados simultaneamente para executar uma única aplicação em menos tempo.

Do ponto de vista do programador, esse tipo de aplicação pode ser projetado através de dois modelos principais, memória compartilhada e troca de mensagens.

O modelo de programação em memória compartilhada é direcionado para arquiteturas nas quais múltiplos processadores compartilham um único espaço de memória.

A comunicação entre os processadores nesse modelo é realizada lendo-se e escrevendo-se dados nesse espaço de memória.

A noção do que é **privado** e **compartilhado** se torna importante nesse caso.

Dados compartilhados são visíveis para todos os processadores participando da execução paralela enquanto dados privados são locais para cada processador e não podem ser acessados por outros.

A comunicação entre os processadores ocorre através da leitura e escrita nesses dados compartilhados.

Um outro modelo de processamento paralelo é direcionado para arquiteturas de troca de mensagens. Nesse modelo, processadores não compartilham memória. Ao invés disso, eles enviam e recebem mensagens através da rede de interconexão. O melhor exemplo é o caso do MPI (MultiProcessing Interface).

Todos os dados são privados e a única forma de um processador obter uma informação que não está na sua memória local é requisitando-a ao processador que a possui.

Para que um programa seja executado em algum desses modelos é necessário algum tipo de construção de linguagem de programação.

Esse tipo de construção controla o compartilhamento de dados, a sincronização e assim por diante.

Cada um dos dois modelos de programação paralela oferece tipos diferentes de construções para alcançar esse fim.

No modelo de troca de mensagens, as construções geralmente se baseiam em bibliotecas de funções.

Essas bibliotecas incluem funções para enviar e receber mensagens, execução de sincronização, comunicação coletiva, etc.

Nesse modelo o programador precisa explicitamente particionar os dados, realizar a comunicação e a sincronização.

O MPI (Message Passing Interface) é o principal padrão desse tipo de biblioteca utilizado.

No padrão OpenMP para programação em memória compartilhada, o programador adiciona diretivas de compilação ao código fonte.

Essas diretivas não afetam a semântica do programa, apenas indicam como trabalho e dados devem ser compartilhados entre os processadores.

O código fonte é então compilado por um compilador (com suporte ao padrão) que gera código para criar threads que executam em paralelo nos diversos processadores do sistema.

Cada um desses padrões é adequado a seus respectivos modelos de arquitetura de hardware, mas existem arquiteturas de multiprocessadores para as quais o mapeamento para um deles não é tão simples.

Um exemplo importante são os clusters de SMPs (Symmetric Multi-Processors).

Clusters de SMPs são construídos a partir de diversos nós SMP (um tipo de máquina de memória compartilhada) conectados através de uma rede de interconexão.

De fato, existe uma tendência em computação de alto desempenho de construir computadores paralelos acrescentando gradativamente nós SMP interconectados por redes de interconexão simples.

Cada um desses nós consiste de um determinado número de processadores e uma grande quantidade de memória compartilhada.

Clusters de SMPs podem ser programados para utilizar troca de mensagens entre todos os processadores em todos os nós envolvidos no sistema.

No entanto, existe a possibilidade de um melhor desempenho ser for utilizado um modelo hierárquico de programação que seja mapeado diretamente para o modelo físico que combina memória distribuída e compartilhada.

Nesse contexto, um modelo de programação híbrido é definido como o modelo que utiliza multithread em memória compartilhada dentro do nó SMP e troca de mensagens entre os nós SMP.

Teoricamente, programas híbridos deveriam oferecer um desempenho melhor que programas puros de troca de mensagens por três razões. A troca de mensagens dentro do nó é substituída por um acesso de memória compartilhada mais rápido.

Há um volume menor de comunicação nos meios de transmissão de dados já que as mensagens intra-nó não são necessárias.

Há uma menor quantidade de processos envolvidos na comunicação o que deve levar a uma melhor escalabilidade.

Esse trabalho é um estudo sobre um modelo híbrido para aplicações de programação paralela.

O Departamento de Ciência de Computação da Universidade de Brasília possui como recurso, à sua disposição, um cluster de nós SMP.

Esse recurso foi utilizado em diversas pesquisas anteriores para desenvolvimento de aplicações científicas paralelas, mas nunca ficou claro se a arquitetura dessa máquina era utilizada da forma mais efetiva possível.

Decidiu-se então investigar se, com um estilo de programação que mapeasse diretamente para a arquitetura SMP de hardware, os resultados seriam melhores.

Na pesquisa realizada por foi feito um estudo comparativo (em modelos puros de troca de mensagens) do desempenho de códigos paralelos baseados no Método dos Elementos Finitos (MEF) aplicado à elasticidade linear para problemas estruturais que utilizam o método dos gradientes conjugados para solução de sistemas de equações.

Esse problema se apresenta como uma aplicação adequada para essa pesquisa por se tratar de um problema real no qual grandes benefícios podem ser obtidos pela paralelização do código.

De fato, problemas modelados pelo Método dos Elementos Finitos apresentam elevados custos computacionais em termos de tempos de execução e uso de memória, principalmente devido à grande quantidade de dados a serem processados durante a solução do sistema de equações.

O objetivo da pesquisa apresentada nessa dissertação é desenvolver e avaliar o uso de um modelo híbrido de programação para uma aplicação real de engenharia baseada no método dos elementos finitos.

Além disso, investigar quão efetivo é o seu uso através da análise do ganho de desempenho obtido e da comparação com o desempenho do modelo puro MPI.

Na introdução, é apresentada uma visão geral do problema a ser analisado, além da motivação e dos objetivos desse trabalho.

O restante desse documento é organizado da seguinte maneira, O capítulo 2 apresenta uma visão geral da área computação de alto desempenho e nesse capítulo são introduzidos vários conceitos utilizados na dissertação.

Segue-se o capítulo 3 sobre programação em memória compartilhada no padrão OpenMP.

O capítulo 4 aborda programação em memória distribuída e descreve em detalhes o padrão MPI para programação com troca de mensagens.

O capítulo 5 apresenta uma visão geral dos modelos híbridos de programação paralela, além de seus problemas, vantagens e desvantagens.

O capítulo 6 descreve a aplicação e em seguida no capítulo 7 são apresentados os resultados experimentais obtidos.

Finalmente são apresentas conclusões e sugestões de trabalhos futuros no Capítulo 8.

O tempo de execução, T , de um programa depende do número de instruções a serem executadas, do número médio de ciclos de clock consumidos por instrução e do tempo de ciclo de clock, A redução do tempo de um ciclo de clock é uma questão relacionada à engenharia e pode ser alcançada através do uso de materiais mais avançados e da construção de circuitos menores e mais eficientes.

Os outros dois fatores são melhorados através de estratégias de paralelismo, que replicam componentes básicos do sistema.

Paralelismo existe em uma grande variedade de máquinas e se apresenta de várias formas que podem ser classificadas em três níveis distintos, paralelismo de jobs, paralelismo de aplicação e paralelismo de instruções.

Paralelismo de jobs É o nível mais alto de paralelismo e é de interesse maior para administradores de sistema do que de usuários.

Nesse tipo de paralelismo, o mais importante é que um laboratório ou centro de comunicação execute uma maior quantidade de jobs possíveis em um período de tempo específico.

Isso pode ser alcançado, adquirindo-se mais máquinas de forma que uma maior quantidade de jobs seja executada ao mesmo tempo, apesar de para o usuário um job particular não executar mais rápido.

Nesse caso há uma diferenciação entre throughput (número de jobs em um período de tempo) e latência (tempo para executar uma aplicação).

Ocorre quando um programa simples é dividido em partes que são executadas ao mesmo tempo em múltiplos processadores ou múltiplas unidades funcionais.

Paralelismo em nível de programa geralmente se manifesta de duas formas, sobre seções independentes de um mesmo programa, ou sobre iterações individuais de um laço onde não há dependência entre de dados.

É invisível para os usuários e está no nível de organização de computadores.

Pipelines são a forma mais comum de alcançar esse tipo de paralelismo.

Nesse caso, instruções podem ser sobrepostas ou uma determinada instrução pode ser decomposta em suboperações e essas suboperações serem sobrepostas.

Em geral, programadores não precisam se preocupar com esse nível de paralelismo já que compiladores são capazes de organizar programas para explorá-lo.

Um conceito relacionado ao nível de paralelismo é o tamanho de grão das tarefas paralelas.

Em um sistema de grão grosso as tarefas que executam em paralelo representam trechos grandes da aplicação.

Em sistemas paralelos de grão fino, por outro lado, as tarefas representam trechos bem pequenos da aplicação constituídos de algumas poucas instruções.

No projeto de computadores de alto desempenho, deve-se tomar uma decisão entre um pequeno número de processadores poderosos ou um grande número de processadores simples para alcançar o desempenho necessário.

A vantagem de um pequeno número de processadores poderosos é a simplicidade de interconexão e a possibilidade de utilizar organizações de memória que facilitem a programação.

Por outro lado esses processadores são extremamente caros.

Utilizar um grande número de processadores simples oferece uma grande economia nesse sentido, em detrimento de uma maior complexidade nas estratégias de interconexão e de estratégias de organização de memória que dificultam a programação.

Com o rápido desenvolvimento dos microprocessadores, máquinas constituídas de algumas centenas de processadores alcançam o mesmo desempenho das máquinas de processadores especializadas mais poderosas e mais caras.

A principal terminologia para classificação de sistemas distribuídos foi proposta por Flyn e apesar de ser rudimentar, é muito útil para a classificação de computadores de alto desempenho.

Essa classificação se baseia na forma como fluxos de instruções e fluxos de dados são manipulados e inclui quatro classes de computadores, SISD (Fluxo de instruções único, fluxo de dados único).

Esses são os sistemas convencionais monoprocesados que possuem uma CPU e que acomodam um fluxo de instruções que é executado de forma serial.

SIMD (Fluxo de instruções único, vários fluxos de dados).

A mesma instrução é executada por diferentes fluxos de dados em diferentes processadores.

Cada processador tem sua própria memória de dados, mas existe uma única memória de instruções e uma única unidade de controle que busca e despacha as instruções.

As arquiteturas vetoriais são a classe mais ampla de processadores desse tipo.

MISD (Vários fluxos de instruções, fluxo de dados único).

Só existe um único fluxo de dados operando por sucessivas unidades funcionais.

Nenhum multiprocessador comercial desse tipo foi construído até hoje, mas pode ser elaborado no futuro.

Alguns processadores de fluxo de uso especial se aproximam de uma forma limitada dessa categoria.

MIMD (Vários fluxos de instruções, vários fluxos de dados).

Essas máquinas executam vários fluxos de instruções diferentes em paralelo em dados diferentes.

A diferença em relação às máquinas SISD reside no fato que as instruções e os dados são relacionados porque representam partes diferentes da mesma tarefa a ser executada.

Dessa forma, sistemas MIMD podem executar diversas subtarefas em paralelo para diminuir o tempo de solução da tarefa principal.

Existe uma grande variedade de sistemas MIMD e especialmente nessa classe, a classificação de Flynn se mostra inadequada, pois inclui arquiteturas totalmente distintas.

A categoria das máquinas MIMD pode ser subdividida em dois outros grupos, Sistemas de Memória Compartilhada e Sistemas de Memória Distribuída.

No caso de multiprocessadores com número pequeno de nós, é possível que os processadores compartilhem fisicamente uma única memória centralizada e que processadores e memória sejam interconectados por um barramento.

Com caches grandes, o barramento e a memória única podem satisfazer às demandas de memória de um número pequeno de processadores.

Substituindo-se o barramento único por vários barramentos ou até mesmo por um switch, um projeto de memória compartilhada pode ter sua escala aumentada até algumas dezenas de processadores.

Embora esse aumento de escala seja tecnicamente concebível, essa organização se torna menos atraente à medida que aumenta o número de processadores devido à contenção no acesso à memória.

Pelo fato de existir uma única memória principal que tem um relacionamento simétrico com todos os processadores e um tempo de

acesso uniforme a partir de qualquer processador, esses multiprocessadores freqüentemente são chamados de multiprocessadores simétricos (de memória compartilhada) (SMP-Symetric Multiprocessors).

Mostra uma representação de um sistema de memória compartilhada.

O segundo grupo consiste em multiprocessadores com memória fisicamente distribuída.

Para dar suporte a quantidades maiores de processadores, a memória deve ser distribuída entre os processadores em vez de centralizada para atender à demanda de largura de banda sem incorrer em uma latência de acesso à memória longa demais.

A distribuição de memória entre os nós tem duas vantagens importantes.

Primeiro, é uma forma econômica de aumentar a escala da largura de banda de memória se a maior parte dos acessos se destina à memória local no nó.

Em segundo lugar, ela reduz a latência para acesso à memória local.

Estrutura básica de um multiprocessador de memória compartilhada centralizada.

A principal desvantagem é que a comunicação de dados entre os processadores se torna mais complexa e tem latência mais alta porque os processadores não compartilham mais uma única memória centralizada.

Mostra uma representação de um multiprocessador de memória distribuída.

Para os multiprocessadores de memória fisicamente compartilhada, a comunicação entre processadores é realizada simplesmente escrevendo-se e lendo-se informações do espaço de endereçamento comum.

É uma forma de comunicação bem eficiente, porém bastante restrita quanto ao número de processadores participantes, devido a problemas de acesso à memória.

Normalmente, para qualquer multiprocessador com uma quantidade muito grande de nós deve-se utilizar módulos de memória distribuídos fisicamente com os processadores.

Nesse caso, existem duas abordagens alternativas para troca de informações entre os processadores.

Na primeira abordagem, a comunicação ocorre por meio de um espaço de endereçamento virtual compartilhado.

Isto é, as memórias fisicamente separadas podem ser endereçadas com um único espaço de endereços compartilhado logicamente.

Cada referência à memória pode então ser realizada por qualquer processador e para qualquer posição no espaço de endereçamento global, supondo apenas que ele tenha os direitos de acesso corretos.

Esses multiprocessadores são chamados de arquiteturas de memória compartilhada distribuída (DSM distributed shared memory).

O termo memória compartilhada nesse caso não significa que existe um compartilhamento físico de memória e sim um compartilhamento lógico.

Na segunda abordagem, o espaço de endereços é constituído de vários módulos de memória disjuntos que não podem ser endereçados por um processador remoto.

Em tais máquinas, o mesmo endereço físico em dois processadores diferentes se refere a duas informações diferentes.

Cada módulo processador-memória é em essência um computador separado e por essa razão esses computadores paralelos são chamados de multicomputadores.

Um multicomputador pode consistir até mesmo de computadores completos, totalmente independentes e conectados apenas por uma rede local.

Esse tipo de multicomputador recebe hoje em dia a denominação popular de cluster.

A comunicação de dados em multiprocessadores é realizada através da troca explícita de mensagens.

Um cluster é uma coleção de computadores completos (nós) que são conectados fisicamente por uma rede de alta performance ou uma rede local.

Tipicamente, cada nó é uma estação de trabalho, um computador pessoal ou uma máquina SMP.

O mais importante é que todos os nós do cluster devem ser capazes de trabalhar juntos como um recurso computacional único e integrado.

Além disso, cada nó pode trabalhar como uma máquina individual.

O objetivo dos clusters é fornecer serviços de alta disponibilidade de alto desempenho.

A arquitetura conceitual de um cluster é apresentada.

Arquitetura de um Cluster de Computadores.

Cada nó é um computador completo.

Isso implica que cada nó possui seu próprio processador, cache, memória, disco e dispositivos de E/S.

Além disso, existe um sistema operacional completo e independente em cada nó.

Um nó pode possuir mais de um processador, mas apenas uma cópia do sistema operacional.

Um cluster é um recurso computacional único.

Isso o diferencia de um sistema distribuído típico cujos nós são usados como recursos individuais.

Um cluster realiza o conceito de recurso único através de uma das diversas técnicas de SSI (single system image, imagem única do sistema).

Conexão Entre Nós Os nós de um cluster geralmente são conectados através de uma rede padrão como Ethernet, FDDI, Fiber-Channel ou ATM.

Clusters representam uma forma de aumentar a disponibilidade de um sistema, ou seja a porcentagem de tempo que o sistema fica disponível para o usuário.

Aumento de Desempenho Um cluster deve oferecer alto desempenho para várias finalidades.

Uma finalidade é utilizá-lo como um superservidor.

Se cada nó de um cluster pode servir n clientes, então o cluster como um todo pode servir mn clientes simultaneamente.

Uma outra finalidade é utilizar o cluster para minimizar o tempo de execução de uma aplicação, distribuindo o trabalho em tarefas paralelas.

O tempo de execução é a medida de desempenho mais confiável e que melhor traduz o que se busca em termos de velocidade de processamento, tanto do hardware quanto do software.

Nesse sentido uma definição de desempenho pode ser, O tempo de execução pode ser visto de duas maneiras, tempo decorrido desde o início até o final da execução do programa ou o tempo de CPU, isto é, o tempo que efetivamente foi utilizado pela CPU para executar o programa, excluindo-se o tempo consumido pelo próprio sistema operacional durante a execução do programa.

O normal é a utilização do tempo decorrido ou do programa como um todo, ou parte dele.

Pois o tempo de CPU também não contabiliza o tempo de comunicação, que é adicionado para se utilizar as máquinas em paralelo.

O tempo total de execução pode ser decomposto em tempo de computação, tempo de comunicação e tempo de espera.

Tempo de Computação O tempo de computação de um algoritmo é o tempo consumido realizando algum trabalho.

O tempo de computação geralmente vai depender de alguma forma, do tamanho do problema.

Se o algoritmo paralelo replica computação, então o tempo de computação também dependerá do número de tarefas ou processadores.

Em sistemas heterogêneos, o tempo de computação também dependerá de em qual processador as operações foram realizadas.

Além disso, também dependerá das características do processador e de seu sistema de memória.

O tempo de comunicação de um algoritmo é o tempo que suas tarefas gastam enviando e recebendo mensagens.

Existem dois tipos básicos de comunicação, interprocessos e intraprocessos.

Na comunicação interprocessos, as duas tarefas comunicantes estão localizadas em processos diferentes.

Na comunicação intraprocessos, duas tarefas comunicantes estão localizadas no mesmo processador.

O custo de envio de uma mensagem entre duas tarefas localizadas em processadores diferentes pode ser representado por dois parâmetros, o tempo para início da comunicação e o tempo de transferência por palavra (tipicamente 4 bytes) que é determinado pela largura de banda física do link de comunicação.

Um processador pode estar em espera devida à falta de computação ou falta de dados.

No primeiro caso, o tempo espera pode ser reduzido usando técnicas de balanceamento de carga.

No segundo caso o processador está em espera enquanto computação e comunicação necessárias precisam ser realizadas.

Ambos, tempo de computação e tempo de comunicação são especificados explicitamente em um algoritmo paralelo.

Assim, é mais fácil determinar suas contribuições para o tempo de execução.

Tempo de espera pode ser mais difícil para determinar, já que muitas vezes depende da ordem na qual as operações são apresentadas.

O tempo de espera pode às vezes ser evitado estruturando o programa de forma que processadores realizem outras computações ou comunicações enquanto esperam por dados remotos.

Essa técnica é chamada de sobreposição de computação e comunicação já que computação local pode ser realizada concorrentemente com comunicação remota e computação.

Essa sobreposição pode ser alcançada de duas formas.

A abordagem mais simples é criar múltiplas tarefas em cada processador.

Quando uma tarefa está bloqueada esperando por dados remotos, a execução pode ser passada para outra tarefa para a qual dados já estão disponíveis.

Tratando-se de uma aplicação paralela, onde o objetivo é obter um aumento na velocidade de processamento pela subdivisão do problema em tarefas que podem ser executadas concorrentemente, é interessante verificar o desempenho do sistema através do speedup S_n , que é a razão entre o tempo gasto na execução seqüencial ou em uma única máquina, T_1 e o tempo de execução em paralelo, isto é, em mais de uma máquina, T_n .

O speedup evidencia o ganho de tempo obtido na execução paralela para um dado número de tarefas concorrentes.

A eficiência, E_n , é a razão entre o valor do speedup (S_n) e o número de tarefas n , utilizadas na execução paralela Assim, a eficiência mostra se os ganhos obtidos com a adição de máquinas estão sendo relevantes de forma a determinar a quantidade ótima de máquinas

necessárias para a execução paralela de um dado tipo de problema e volume de dados.

Um aspecto importante da análise de performance é o estudo de como o desempenho do algoritmo varia com relação a parâmetros como tamanho do problema, número de processadores, etc.

Para algoritmos paralelos em particular, é de interesse o comportamento com aumento no número de processadores.

Uma abordagem para quantificar escalabilidade é determinar como o tempo de execução e a eficiência variam como o aumento do número de processadores, para um tamanho de problema fixo.

Essa análise do problema fixo possibilita a resposta à questões como, o limite de velocidade para resolver um determinado problema em um computador específico.

O maior número de processadores que se deve utilizar para manter uma eficiência É importante considerar tanto eficiência quanto tempo quando a escalabilidade está sendo avaliada.

Enquanto a eficiência vai geralmente diminuir monotonamente com o número de processadores, o tempo de execução pode aumentar se o modelo de desempenho inclui um termo proporcional a uma potencia positiva do número de processadores.

Em alguns casos, não será produtivo mais do que algum número máximo de processadores para um tamanho de problema em particular ou escolha de parâmetros de máquina.

Gene Amdahl, formulou o que hoje é chamado de a Lei de Amdahl para caracterizar a maneira como uma aplicação poderia utilizar de forma eficiente processadores paralelos escaláveis.

Praticamente todo programa paralelo mescla partes que são seriais e partes paralelas.

Uma análise de engenharia é um bom exemplo.

A parte de inicialização provavelmente será completamente serial, os dados de entrada são lidos, a matriz é preenchida e os dados são particionados.

A fase de solução do problema, por outro lado, poderá se beneficiar da presença de múltiplos processadores e ser altamente paralela.

Com a lei da Amdahl pode-se verificar quanto de um programa pode executar em paralelo e quanto desse mesmo programa precisa executar usando apenas um processador.

Uma vez que a razão entre a porção paralela e a seqüencial é estabelecida, ela impõe um limite superior para o speedup possível da aplicação.

Seja uma aplicação que execute em 100 minutos onde essa aplicação pode executar em paralelo por 95 dos 100 minutos.

Nesse caso, mesmo se fossem utilizados tantos processadores que a porção paralela do código executasse em um piscar de olhos, o tempo

total de execução seria ainda de 5 minutos (mais um piscar de olhos) devido à parte serial do código.

Assim, mesmo adquirindo um número infinito de processadores, nós melhoramos o desempenho da aplicação apenas por um fator de 20.

Mostra um gráfico do speedup para o número de processadores e pode-se perceber que não há benefício em adicionar mais processadores a partir do momento que a parte serial do código se torna o fator dominante no tempo de execução.

Quando a lei de Amdahl começou a ser discutida parecia que máquinas paralelas em grande escala eram pouco vantajosas.

Pesquisas em códigos existentes colocam a porcentagem do código que pode-se executar em paralelo para aplicações típicas entre 60% e 95%.

Uma análise superficial concluiria que muito mais do que cerca de 8 processadores não seria muito vantajoso.

Lei de Amdahl para uma aplicação da qual 95% do código pode executar em paralelo.

A Lei de Amdahl é correta, mas existem algumas suposições erradas, feitas quando ela foi usada para concluir que paralelismo maciço não era vantajoso, Pelo fato de não haver processadores paralelos disponíveis no momento da concepção da maioria das aplicações da pesquisa citada, os programadores não se preocuparam em tentar tornar seus códigos mais adequados para serem paralelizados.

À medida que esses sistemas se tornam muito utilizados, muito esforço é realizado para repensar essas aplicações de modo a maximizar as operações que podem ser realizadas em paralelo.

Geralmente quando o tamanho do problema é duplicado, a parte serial leva o dobro do tempo para executar uma tarefa, enquanto a parte paralela leva de quatro a oito vezes mais tempo.

A disparidade no aumento relativo de tempo significa que a aplicação gasta mais tempo na parte paralela da aplicação.

À medida que as capacidades de memória e de processador aumentam, pesquisadores começaram a resolver problemas maiores.

Assim, tornando-se um problema maior, um problema com porção paralela de 95% se torna 99% paralelo e tornando-se ainda maior se torna 999% paralelo e assim por diante.

Devido ao fato do speedup potencial ser tão influenciado pelo tamanho do problema, surgiram algumas novas leis que foram criadas para capturar esse efeito.

Essas leis são chamadas Lei de Gustafson.

A lei de Gustafson analisa como o aumento do tamanho de um problema afeta a escalabilidade.

Já a lei de Ni analisa a relação entre o crescimento do tamanho do problema e a habilidade de executá-lo em paralelo.

Ocorreram também avanços em pesquisas de ferramentas de análise de fluxo de dados para detectar e extrair o paralelismo de códigos onde essa tarefa realizada apenas pela análise manual do código era muito difícil.

O progresso na construção e no uso de processadores paralelos efetivos e eficientes é lento.

Essa taxa de progresso foi limitada por sérios problemas de software, bem como por um longo processo de evolução da arquitetura de multiprocessadores para aumentar a facilidade de uso e melhorar a eficiência.

A grande variedade de abordagens arquitetônicas e o sucesso limitado, além da curta direção de muitas arquiteturas, representam as principais dificuldades no que diz respeito ao software.

Entretanto, o progresso realizado apresenta alguns motivos para otimismo quanto ao futuro do processamento paralelo e dos multiprocessadores.

Em primeiro lugar, o uso de processamento paralelo em alguns domínios começa a ser compreendido.

Dentre eles talvez o principal seja o da computação científica e da engenharia.

Esse domínio de aplicações tem uma ânsia quase ilimitada por maior capacidade de computação e nele existem muitas aplicações que têm uma grande porção de paralelismo natural.

Outra área de aplicações importante e muito maior é a dos sistemas de bancos de dados e processamento de transações em larga escala.

Esse domínio de aplicações também tem muito paralelismo natural disponível através do processamento paralelo de solicitações independentes, mas sua necessidade de computação em larga escala em oposição ao simples acesso a sistemas de armazenamento são menos compreendidas.

Além disso, existe hoje uma grande crença de que o modo mais efetivo de construir um computador que ofereça maior desempenho do que pode ser alcançado com um microprocessador de um único chip, é construir um cluster que amplie as vantagens significativas de preço-desempenho dos microprocessadores produzidos em massa.

Um terceiro motivo é que os multiprocessadores são altamente eficientes para cargas de trabalho multiprogramadas, que são com freqüência o uso dominante de mainframes e servidores de grande porte, como também de servidores de arquivo e servidores Web.

Essas aplicações constituem efetivamente um tipo restrito de carga de trabalho paralela.

No futuro, essas cargas de trabalho poderão representar um grande mercado para multiprocessadores de alto desempenho.

Quando uma carga de trabalho quiser compartilhar recursos como o armazenamento de arquivos ou puder compartilhar de forma eficiente um recurso como uma memória extensa, um multiprocessador poderá ser uma alternativa muito eficiente.

Os multiprocessadores se mostram muito eficazes para certas cargas de trabalho comerciais intensivas e aplicações de pesquisa na Web em larga escala.

No caso de aplicações comerciais que não exigem alto desempenho de comunicação, que tenham pequena necessidade de memória ou demanda limitada por computação, é provável que os clusters sejam mais econômicos que os multiprocessadores.

Atualmente o espaço comercial é uma mistura de clusters de PCs básicos, SMPs e clusters de SMPs, com diferentes estilos arquitetônicos.

E finalmente, o multiprocessamento no chip cresceu em importância por duas razões.

Primeiro, no mercado embutido no qual o paralelismo natural existe com frequência, tais abordagens representam uma alternativa óbvia para processadores mais rápidos e possivelmente menos eficientes no uso do silício.

Em segundo lugar, a diminuição dos rendimentos nos projetos de microprocessadores de ponta incentiva os projetistas a buscar o multiprocessamento no chip como uma solução potencialmente mais econômica.

O mercado de computação de alto desempenho sempre se caracterizou pela rápida mudança de fornecedores, arquiteturas e tecnologias.

Apesar de todas essas mudanças, a evolução do desempenho em larga escala parece ser um processo contínuo.

A lei de Moore ("O número de transistores em um chip dobra a cada dois anos") geralmente é citada nesse contexto.

Mostra o gráfico dos computadores que alcançaram o pico de desempenho em suas respectivas épocas nas últimas seis décadas.

Na média houve um aumento de desempenho de duas magnitudes a cada década.

Fica claro que a lei de Moore foi verdadeira praticamente em todo o desenvolvimento da computação moderna.

No final da década de 90, os clusters eram comuns no ambiente acadêmico, mas principalmente como objeto de pesquisa e não como uma plataforma para aplicações reais.

A maioria desses clusters era pouco poderosa e, como resultado, a edição de Novembro de 1999 do TOP500 listava apenas sete sistemas baseados em cluster.

Desempenho dos computadores mais rápidos nas últimas seis décadas.

Isso mudou a partir do momento que os clusters chamaram a atenção de desenvolvedores de aplicações comerciais e industriais e que aplicações com requisitos de comunicação menos restritivos tornaram o custo benefício dos clusters baseados em componentes de prateleira interessante.

Em pouco tempo, a maioria dos fornecedores no mercado de computação de alto desempenho comercializava esse tipo de sistema.

Mostra que em Novembro de 2006 os clusters eram a arquitetura dominante no Top500, com 361 sistemas.

Há, no entanto, ainda uma grande diferença entre a utilização principal de clusters e o uso das outras arquiteturas mais integradas.

Os grandes supercomputadores são usados principalmente para turnaround computing onde o poder de processamento máximo é aplicado para um único problema.

O objetivo é resolver um grande problema que não pode ser resolvido em um tempo razoável seqüencialmente ou resolver um problema simples em um intervalo de tempo menor.

Esse tipo de computação permite a solução de problemas com restrições de tempo real.

O objetivo principal nesse caso é o tempo para solução do problema.

Os clusters por outro lado, geralmente rodam diversos jobs simultaneamente e o objetivo principal é obter o máximo desempenho por custo.

Principais arquiteturas no Top 500 na última década.

A comunidade da computação de alto desempenho já utilizava componentes de prateleira maciçamente nos anos 90.

MPPs e Constellations (Clusters de SMPs) tipicamente usavam processadores padrões de estações de trabalho apesar de usarem redes de interconexão customizadas.

Havia, porém uma grande exceção, praticamente ninguém usava processadores Intel.

Baixo desempenho e a limitação do projeto de 32 bits eram as principais razões para isso.

Isso mudou com o surgimento do Pentium III e especialmente em 2001 com o surgimento do Pentium 4, que trazia grandes avanços no desempenho de memória graças ao novo modelo de barramento e suportava pontos flutuantes de 64 bits.

O número de processadores no Top500 com processadores Intel passou de apenas seis em Novembro de 2000 para trezentos e dezoito em Novembro de 2004.

O interesse em arquiteturas inovadoras sempre foi grande na comunidade de computação de alto desempenho.

Isso não chega a ser surpreendente já que essa área nasceu como e sustenta seu crescimento nas inovações tecnológicas.

Uma das preocupações atuais diz respeito à necessidade crescente de espaço e poder computacional nos clusters modernos.

Processadores no TOP 500 na última década.

No desenvolvimento do BlueGene/L, a IBM atacou esse problema projetando um sistema muito eficiente tanto do ponto de vista computacional como do ponto de vista de O BlueGene/L, não utiliza os processadores mais poderosos disponíveis no mercado e sim processadores projetados especialmente para sistemas embarcados.

Além de uma redução grande na memória principal total disponível, uma característica marcante nesse sistema é ele ser extremamente denso.

Para alcançar o desempenho esperado um número enorme desses processadores (mais de 128000) foi combinado, usando mecanismos de interconexão especializados.

Havia muitas dúvidas quanto à capacidade desse sistema em alcançar o desempenho prometido e quanto à sua viabilidade como um sistema de uso geral.

Os primeiros resultados da versão beta foram encorajadores e uma versão quatro vezes menor figurou na primeira posição da edição de Novembro de 2004 da lista dos top500.

Ao contrário do progresso no desenvolvimento de hardware, houve pouco, talvez uma regressão, na facilidade de programação de sistemas escaláveis.

Algumas tentativas iniciadas na década de 90 na direção de software mais acessível foram completamente abandonadas.

O movimento para o modelo de memória distribuída forçou mudanças no paradigma de programação.

O alto custo de comunicação e sincronização da interação processador-processador exige novos algoritmos que minimizem essas operações.

O uso de sistemas de memória distribuída provocou o crescimento de novos modelos de programação, principalmente o paradigma de troca de mensagens.

No entanto, os progressos na área de debuggers e ferramentas de desempenho é lento e a maioria dos usuários considera as ferramentas de programação para supercomputadores paralelos extremamente inadequadas.

Baseando-se nos dados atuais da lista dos Top500 que cobre os últimos 20 anos e assumindo que o desenvolvimento atual de desempenho se mantenha nos próximos anos, o desempenho observado atualmente pode ser extrapolado para um provável cenário futuro.

Essa projeção foi feita utilizando regressão linear sobre a escala logarítmica dos níveis de desempenho do Top500.

Projeções para o Top500 na próxima década.

Olhando quatro anos a frente espera-se que o primeiro sistema PetaFlop apareça na lista por volta de 2009.

Olhando ainda mais longe, pode-se especular que baseado na duplicação de desempenho todo ano, o primeiro sistema excedendo 100 Petaflops/s deve estar disponível por volta de 2015.

No entanto, devido à rápida mudança nas tecnologias usadas em computação de alto desempenho não há como imaginar ao certo como será a arquitetura desses sistemas e como alcançarão esses níveis.

O fim da lei de Moore como nós a conhecemos já foi previsto muitas vezes e talvez um dia aconteça.

Novas tecnologias como a computação quântica, possivelmente permitirão no futuro aumentar capacidades computacionais muito além daquelas previstas nessas projeções.

No entanto, apesar da área de alto desempenho ter mudado radicalmente diversas vezes desde o surgimento do Cray1 quarenta anos atrás, não há perspectiva para uma mudança nesse ciclo de constantes redefinições.

Em um sistema de memória compartilhada, cada processador tem acesso direto à memória dos demais, isto é, pode recuperar ou gravar informações em qualquer endereço que faz parte da memória compartilhada.

O programador pode ainda declarar certas partes da memória como exclusivas a um processador o que constitui um simples, mas poderoso modelo para expressar e gerenciar o paralelismo em uma aplicação.

Apesar de sua simplicidade e da escalabilidade conseguida com avanços tecnológicos recentes, muitos desenvolvedores de aplicações paralelas resistem ao modelo e uma das razões diz respeito à portabilidade.

Ao longo do tempo cada fornecedor de hardware para o modelo criou suas próprias extensões de C ou Fortran para programação paralela em memória compartilhada de forma que uma aplicação a ser portada de uma plataforma para outra precisasse ser reescrita.

Esse capítulo apresenta uma alternativa portátil para memória compartilhada, o OpenMP.

Pelo lado da programação com troca de mensagens, o MPI praticamente padronizou o modelo.

O MPI é portátil, disponível para diversas plataformas e é aceito como um padrão para esse tipo de aplicações.

No entanto, troca de mensagens é uma técnica de programação difícil.

Ela exige que os dados do programa sejam divididos explicitamente e então toda a aplicação deve ser paralelizada para trabalhar com os dados particionados.

Não há uma forma incremental de paralelizar uma aplicação.

Além disso, arquiteturas de multiprocessadores modernas têm gradualmente fornecido suporte a técnicas para memória compartilhada (hardware para coerência de cachê, por exemplo), o que torna a troca de mensagens uma opção às vezes desnecessária e muito restritiva.

Existe ainda o padrão Pthreads.

Esse é um padrão algumas vezes utilizado para memória compartilhada em sistemas de baixo nível.

No entanto não está focado no campo de computação de alto desempenho.

Há muito pouco suporte a pthreads em Fortran além de não ser uma estratégia escalável.

Até mesmo para aplicações escritas em C, o modelo com pthreads é de nível baixo demais, muito mais do que o necessário para a maioria das aplicações científicas.

Além disso, é um modelo voltado mais para paralelismo de instruções e não paralelismo de dados além de possuir portabilidade limitada.

Um outro ponto é que desenvolvedores de software para aplicações científicas e laboratórios governamentais possuem um grande volume de código que precisa ser paralelizado de forma portátil.

Os desenvolvedores precisam paralelizar tais códigos sem a necessidade de reescrevê-los completamente, mas isso não é possível com a maioria dos padrões de programação paralela.

O OpenMP possibilita isso.

O OpenMP é o modelo ideal para programadores que precisam paralelizar rapidamente aplicações científicas já existentes, mas é flexível o suficiente para atender a um conjunto muito maior de finalidades.

O OpenMP fornece um caminho incremental para a conversão de software existente para execução em paralelo.

Também fornece escalabilidade e performance para reescrever uma aplicação completamente ou para desenvolver uma nova aplicação.

Basicamente, OpenMP é um conjunto de diretivas de compilação e uma biblioteca de rotinas que estendem o Fortran e C/C++ para expressar paralelismo em termos de memória compartilhada.

Ele não especifica a linguagem base e pode ser implementado em qualquer compilador.

Vários fornecedores possuem produtos que suportam OpenMP, incluindo compiladores, ferramentas de desenvolvimento e ferramentas de análise de desempenho.

O "OpenMP Review Board" inclui membros como Digita, HP, Intel, IBM.

Todas essas companhias estão ativamente desenvolvendo compiladores e ferramentas para o OpenMP.

Um compilador com suporte a OpenMP, transforma o código original em uma versão paralela para memória compartilhada, se guiando pelas anotações em forma de diretivas.

O OpenMP define assim, um processo de paralelização automático, porém guiado pelo usuário.

Isto significa que o compilador não precisa realizar uma vasta análise de código, já que se baseia apenas em informações fornecidas pelo programador.

Isto dá ao usuário total controle sobre o que deve ser paralelizado e como, ao mesmo tempo que diminui sensivelmente a complexidade do

compilador O OpenMP foi projetado para ser um padrão flexível e de fácil implementação em diferentes plataformas.

O padrão tem quatro partes distintas.

Estrutura de controle, O OpenMP possui um conjunto minimalista de estruturas de controle.

A experiência indica que apenas algumas poucas são necessárias para escrever a maioria das aplicações paralelas.

Dessa forma, o OpenMP inclui estruturas de controles apenas para aquelas situações onde o compilador pode oferecer funcionalidade e desempenho superiores ao que poderia ser desenvolvido pelo próprio programador.

Ambiente de dados, Associado a cada tarefa existe um único ambiente de dados fornecendo um contexto para execução.

A tarefa inicial possui um ambiente de dados que existe por tanto tempo quanto dure o programa.

Ela constrói novos ambientes de dados apenas para aquelas tarefas criadas durante a execução do programa.

Os objetos que fazem parte de um ambiente de dados podem ter um dos três atributos básicos, shared, private, reduction.

Sincronização, Há dois tipos de sincronização, explícita ou implícita.

A sincronização implícita ocorre no começo e no fim de blocos parallel e blocos de diretivas de controle.

O usuário especifica sincronizações explícitas para gerenciar ordem ou dependência de dados.

Sincronização é um tipo de comunicação entre processos e como tal pode afetar significativamente o desempenho do programa.

Assim, em geral, quando a sincronização é minimizada obtêm-se melhor desempenho.

Por essa razão, o OpenMP fornece um rico conjunto de funcionalidades de sincronização para que os programadores possam ajustar adequadamente a sincronização em suas aplicações.

Biblioteca de Rotinas, O OpenMP fornece ainda uma biblioteca de rotinas e um conjunto de variáveis de ambiente.

A biblioteca de rotinas inclui rotinas de query e lock.

Além de permitir à uma aplicação especificar o modo como ela deve ser executada.

As variáveis de ambiente por sua vez, ajudam os programadores a criar além de aplicações portáteis, ambientes de execução também portáteis.

O paralelismo no OpenMP é baseado em threads através do padrão fork-join.

Todos os programas OpenMp iniciam com uma única thread, a thread mestre.

A thread mestre executa seqüencialmente até que a primeira região paralela é encontrada.

Ela então realiza um fork criando um conjunto de threads e assim as instruções originalmente delimitadas no código pela diretiva, são executadas em paralelo nas diversas threads do conjunto.

Quando as threads filhas completam sua execução das instruções da região paralela elas sincronizam e terminam, restando por fim apenas a thread mestre novamente.

Modelo de Execução do OpenMP.

O OpenMP é de fácil utilização e consiste basicamente de dois tipos de construções, pragmas e rotinas de biblioteca.

Os pragmas do OpenMP instruem o compilador para paralelizar seções de código.

Todos os pragmas do OpenMP começam com #pragma omp.

Como toda diretiva pragma, essas também são ignoradas pelo compilador se ele não suporta a funcionalidade, nesse caso OpenMP.

A finalidade principal das rotinas do OpenMP é de atribuir e recuperar informações sobre o ambiente.

Há também rotinas para alguns tipos de sincronização.

Para utilizar as rotinas de biblioteca do OpenMP, o programa deve incluir o arquivo de header omp h'.

Se a aplicação utilizar apenas pragmas, o arquivo de header pode ser ignorado.

Formato de pragmas do OpenMP.

As diretivas incluem o seguinte, parallel, for, parallel for, section, sections, single, master, critical, flush, ordered, e atomic.

Essas diretivas especificam compartilhamento de trabalho entre threads ou instruções de sincronização.

As cláusulas são opcionais e alteram o comportamento das diretivas.

Cada diretiva possui um conjunto diferente de cláusulas disponível e cinco diretivas (master, critical, flush, ordered e atomic) nunca aceitam cláusulas.

Apesar de existirem muitas diretivas, é possível escrever aplicações relevantes utilizando apenas algumas delas.

A diretiva mais comum e importante é a diretiva parallel.

Essa diretiva cria uma região paralela para o bloco estruturado que segue a diretiva.

Essa diretiva diz ao compilador que o bloco estruturado de código deve ser executado em paralelo em múltiplas threads.

Cada thread irá executar o mesmo fluxo de instruções, no entanto não necessariamente o mesmo conjunto de instruções, devido a possíveis instruções de controle de fluxo como if-else.

Segue-se, um exemplo que calcula a média de dois valores em um vetor e armazena o resultado em outro vetor.

Aqui é introduzida uma nova diretiva OpenMP, `#pragma omp for`.

Essa é uma diretiva de compartilhamento de trabalho que instrui o compilador para dividir as iterações do loop que se segue entre as threads do time.

Exemplo de um trecho de código utilizando a diretiva `#pragma omp for`.

Nesse caso, se `size` tem o valor 100 e o loop é executado em uma máquina com 4 processadores, as iterações do loop são alocadas de forma que o processador 1 fica responsável pelas iterações de 1 a 25, o processador 2 pelas iterações de 26 a 50, o processador 3 de 51 a 75 e o processador 4 pelas iterações de 76 a 99.

Nesse caso é utilizada uma política estática de escalonamento, mas outras políticas podem ainda ser selecionadas.

Se o trecho de código não utilizasse o `pragma for`, então cada thread utilizaria o loop completo e realizariam trabalho redundante.

Em todos os loops precedentes não existia dependência entre as iterações de loop.

No exemplo que se segue, existem duas dependências de loop diferentes, Dependências entre iterações de loop.

A paralelização do loop 1 é problemática porque, para executar a iteração i do loop, é necessário o resultado da iteração $i-1$ e portanto há uma dependência entre a iteração i e a iteração $i-1$.

A paralelização do loop 2 é problemática também, mas por uma razão diferente.

Nesse caso é possível calcular o valor de x , mas fazendo isso não é possível calcular o valor de x .

Há uma dependência da iteração $i-1$ para a i .

Quando for realizada paralelização de loops, o programador deve se certificar de que não existem dependências entre as iterações.

Quando não há dependências, o compilador pode executar o loop em qualquer ordem, inclusive em paralelo.

Esse é um requisito importante porque o compilador não é capaz de fazer esse tipo de verificação.

Ao escrever programas paralelos, é muito importante entender quais dados são privados e quais dados são compartilhados para uma execução correta e para um desempenho adequado.

No OpenMP essa distinção é bem clara.

Se variáveis são compartilhadas por todas as threads no time, uma mudança em uma dessas variáveis realizada por uma thread é vista por todas as outras.

Por outro lado, as variáveis privadas, possuem cópias exclusivas para cada uma das threads, de forma que mudanças realizadas por uma delas não são visíveis para as outras.

Por padrão, todas as variáveis são compartilhadas exceto em três exceções.

Em primeiro lugar, para loops `parallel for`, o índice é privado.

Em segundo, variáveis declaradas dentro do bloco `parallel` também são privadas.

E ainda em terceiro lugar, quaisquer variáveis listadas em uma das cláusulas de controle de escopo, `private`, `firstprivate`, `lastprivate`, ou `reduction` são privadas.

Cada uma das cláusulas de controle de escopo de variáveis recebe uma lista de variáveis, mas cada uma possui uma semântica diferente.

A cláusula `private` diz que cada variável em sua lista deve possuir uma cópia privada em cada uma das threads.

A cópia será inicializada com o valor padrão do tipo (por exemplo, 0 para o tipo inteiro).

As cláusulas `firstprivate` e `lastprivate` possuem a mesma semântica da cláusula `private`, exceto que para `firstprivate` o valor da variável imediatamente antes da região paralela é copiado para cada uma das cópias e que para `lastprivate` o valor da variável é copiado para a cópia principal na thread mestre na última iteração.

A cláusula `reduction` possui uma semântica similar à da cláusula `private`, mas recebe além de uma variável, um operador.

O conjunto de operadores é limitado (`+`, `*`, `&`, `|`, `&&`, `||`) e as variáveis de redução devem ser do tipo escalar (`float`, `int` ou `long`).

Ao final do bloco de código, o operador de redução é aplicado às cópias privadas das variáveis e ao seu valor original.

O valor de `sum` é implicitamente inicializado em cada thread com o valor 0.

Assim que o bloco `#pragma omp for` é completado, as threads aplicam o operador `+` para todas as cópias privadas e para o valor original e o resultado é atribuído para a variável `sum`, original da thread mestre.

Redução no OpenMP.

O OpenMP é tipicamente utilizado para paralelismo de loop, mas também suporta paralelismo em nível de funções.

Esse mecanismo é chamado de sessões OpenMP e é útil em várias situações.

Segue-se um trecho de código de uma rotina de QuickSort que utiliza seções, QuickSort utilizando Sections no OpenMP.

Nesse exemplo, o primeiro `#pragma` cria uma região paralela com sessões.

Cada sessão é precedida por uma diretiva `#pragma omp section`.

Cada sessão na região paralela é atribuída a uma única thread do time e todas elas podem executar concorrentemente.

Com várias threads executando concorrentemente, muitas vezes é necessário sincronizá-las.

O OpenMP suporta vários tipos de sincronização.

Um tipo de sincronização importante é a barreira implícita ao final de uma região `parallel`.

Uma sincronização de barreira exige que todas as threads alcancem aquele ponto antes que qualquer uma possa continuar.

Há ainda uma barreira implícita ao final de cada bloco `#pragma omp for` e `#pragma omp single`.

Para remover essa barreira implícita basta incluir a cláusula `nowait`, Remoção de barreira implícita utilizando a cláusula `nowait`.

Um outro tipo de sincronização é realizado com barreiras explícitas.

Em algumas situações, é necessário incluir uma barreira em algum ponto da região paralela.

Isso é incluído no código com a diretiva `#pragma omp barrier` Além disso, em uma região paralela, às vezes pode ser necessário limitar o acesso a uma única thread, como por exemplo, se for necessário escrever em um arquivo no meio de uma região paralela.

Em muitos desses casos, não importa qual thread executa o código contanto que apenas uma o faça.

O OpenMP possui a diretiva `#pragma omp single` para essa finalidade.

Há ainda uma diretiva muito parecida `#pragma omp master` que especifica que a thread única que utilizará o trecho de código será a thread mestre.

Apresenta um exemplo de uso da diretiva `#pragma omp single`.

Uso da diretiva `#pragma omp single`.

Além das diretivas de compilação, O OpenMP possui também um conjunto de rotinas muito útil para o desenvolvimento de aplicação.

Existem três classes de rotinas disponíveis, rotinas de ambiente, rotinas de sincronização e rotinas de tomada de tempo.

Todas essas rotinas começam com `omp_` e são definidas no arquivo de header `omp.h`.

As rotinas de ambiente permitem ao programador recuperar e alterar vários aspectos do ambiente operacional onde o OpenMP está sendo executado.

Funções que começam com `omp_set_` devem, de uma forma geral, ser executadas apenas fora da região paralela.

Todas as outras podem ser executadas em trechos de códigos paralelos ou não paralelos.

Algumas das funcionalidades da biblioteca do OpenMP incluem, por exemplo, especificar ou recuperar o número de threads no time atual (`omp_set_num_threads`, `openmp_get_num_threads`) ou recuperar o número de processadores disponíveis.

Outra construção de sincronização utilizada é `lock`, que controla o acesso a seções críticas de código.

Existem dois tipos de locks, `simple` e `nestable` e cada lock pode existir em um dos três estados, não-inicializado, bloqueado e desbloqueado.

Locks simples (`omp_lock_t`) não podem ser obtidos mais de uma vez, mesmo pela mesma thread.

Locks nestable (`omp_nest_lock_t`) são idênticos aos locks simples, exceto quando uma thread tenta obter o lock que ela já possui, ela não será bloqueada.

Além disso, locks nestable são contadores e mantêm registro de quantas operações de set foram realizadas sobre eles.

Existem rotinas de sincronização que atuam sobre esses locks.

Para cada rotina existe uma variante simples em uma nestable.

Existem cinco ações que podem ser realizadas em um lock, `initialize` (inicializa), `set` (adquire um lock), `unset` (libera um lock), `test` (verifica se o lock está livre) e `destroy` (destrói o lock).

Os programadores podem escolher tanto as rotinas de biblioteca como os pragmas para sincronização.

A vantagem dos pragmas é que eles são extremamente estruturados.

Isso torna mais fácil a compreensão dos programas já que olhando para o código é fácil ver o local de entrada e saída das regiões críticas.

Já as rotinas de biblioteca têm como grande vantagem a flexibilidade.

É possível, por exemplo, passar um lock como parâmetro para uma outra função e dentro dessa função realizar uma operação de set ou unset.

Isso não é possível com os pragmas.

De uma forma geral, faz mais sentido usar os pragmas de sincronização a menos que seja necessário um grau maior de flexibilidade.

Uma primeira desvantagem do OpenMP é que com ele uma aplicação só pode ser executada em um espaço de endereçamento único, ou seja, não é possível executar uma aplicação totalmente OpenMP em um cluster.

Além disso, OpenMP é construído em cima de um modelo nativo de threads e por isso adiciona overhead à aplicação.

Um outro ponto importante é que, pelo fato da paralelização ser gerada pelo compilador, é muito fácil escrever trechos de código incorretos em relação a condições de corridas e deadlocks entre outros problemas.

Quanto às vantagens, talvez a principal seja a simplicidade e o pouco esforço de programação exigido, se comparado a outros modelos.

Além disso, a paralelização de uma aplicação seqüencial já existente pode ser feita de forma incremental pela anotação do código com diretivas.

Isso elimina a necessidade de reescrever o código.

O código original nesse caso ainda é preservado já que se o programa for compilado sem a opção de suporte a OpenMP, todas as diretivas incluídas são ignoradas.

Ainda como vantagens pode-se citar a portabilidade, já que a maioria dos fornecedores de hardware possui compiladores com suporte ao padrão e ainda o mapeamento natural do modelo OpenMP em arquiteturas SMP.

Saber quando usar o OpenMP é quase tão importante quanto saber como usá-lo.

Os seguintes pontos são úteis nessa decisão, Plataforma alvo é multiprocessador ou multicore, Nesse caso, se a aplicação estiver saturando um núcleo ou processador, transformá-la em uma aplicação multithread com o OpenMP irá melhorar o desempenho da aplicação.

Aplicações Multi-Plataforma, O OpenMP é multiplataforma e uma API largamente utilizada e, como usa pragmas, a aplicação pode ser compilada mesmo se o compilador não suportar o padrão OpenMP Paralelização de Loops, O OpenMP é ideal para a paralelização de loops.

Se a aplicação possui loops caros computacionalmente e que não possuem dependências entre as iterações, usar OpenMP é a escolha ideal.

Otimização de última hora necessária.

Pelo fato do OpenMP não exigir a reconstrução da aplicação, é uma ferramenta perfeita para realizar pequenas mudanças para melhorar o desempenho.

OpenMP não é para todo problema multithread.

Ele foi desenvolvido para ser usado em computação de alto desempenho e funciona melhor em estilos de programação que possuem código com muitos loops e arrays de dados compartilhados.

A criação das threads para as regiões paralelas do OpenMP implica em algum overhead.

Para que ocorra um ganho de desempenho, o speedup obtido pela região paralelizada deve se sobrepôr ao overhead de inicialização das threads.

As pragmas do OpenMP, apesar da facilidade de programação, não oferecem um bom feedback quando ocorrem erros.

Se uma aplicação crítica precisa ser desenvolvida e esta precisa detectar falhas e se recuperar delas, então o OpenMP provavelmente não é uma boa escolha (pelo menos em sua versão atual).

Uma das melhorias previstas para as próximas versões é um mecanismo de tratamento de erros mais consistente.

Troca de mensagens é a principal estratégia para comunicação de dados em multiprocessadores.

Um sistema de troca de mensagens tipicamente combina uma memória local e um processador em cada nó, além de uma rede de interconexão.

Nesse modelo não há memória compartilhada global e para movimentar dados de uma memória local para outra é necessário trocar mensagens através da rede de interconexão.

Isso é geralmente realizado através de pares de comandos send/receive que devem ser explicitamente chamados no código da aplicação.

Isso elimina a necessidade de uma grande memória global, assim como os seus requisitos de sincronização.

Dois fatores importantes devem ser considerados no projeto de redes de interconexão para sistemas de troca de mensagens, a largura banda do link e a latência de rede.

A largura de banda do link é definida como o número de bits que podem ser transmitidos por unidade de tempo (bits/s).

Já a latência de rede é definida como o tempo para completar a transferência de uma mensagem.

Quando uma determinada aplicação é executada nesse modelo, o programa é dividido em tarefas concorrentes e cada uma pode ser executada em um processador diferente.

Se o número de tarefas é maior que o número de processadores, então algumas delas terão que dividir um único processador em um esquema de compartilhamento de tempo.

Tarefas executadas em um único processador utilizam canais internos para troca de mensagens enquanto tarefas em processadores diferentes usam canais externos.

Uma vantagem importante desse tipo de troca de dados é a eliminação da necessidade de construções de sincronização como locks e semáforos, o que resulta em um melhor desempenho.

Além disso, um esquema de troca de mensagens possui uma melhor escalabilidade.

Uma mensagem é definida como uma unidade lógica para comunicação interprocessos.

Cada mensagem é considerada como um conjunto de informações relacionadas enviadas como uma entidade.

Uma mensagem pode ser uma instrução, dados, sincronização ou um sinal de interrupção.

Um sistema de troca de mensagens interage como o mundo exterior recebendo mensagens de entrada e enviando mensagens de saída.

É essencial que o mundo exterior perceba um comportamento consistente em um sistema desse tipo.

O MPI é um padrão de uma interface de troca de mensagens para programação paralela em memória distribuída.

O MPI inclui rotinas de comunicação ponto a ponto, operações coletivas, assim como suporte a grupos de tarefas, contextos de comunicação e topologias de aplicação.

O principal objetivo do MPI é fornecer portabilidade entre diferentes plataformas.

O mesmo código fonte de troca de mensagens pode ser executado em diversas arquiteturas de hardware e sistemas operacionais desde que a biblioteca MPI esteja disponível.

Apesar de troca de mensagens ser geralmente associada a computadores paralelos com memória distribuída, o mesmo código pode também ser executado em computadores de memória compartilhada.

Ele pode ainda ser executado em uma rede de estações de trabalho ou até mesmo em uma única máquina multiprocessada.

O conhecimento de que existem implementações eficientes para o MPI para uma grande variedade de plataformas possibilita uma grande flexibilidade no desenvolvimento e na escolha do ambiente a ser utilizado.

Um outro tipo de compatibilidade possibilitada pelo MPI é a habilidade de executar aplicações de forma transparente em sistemas heterogêneos, isto é, coleções de processadores com arquiteturas de hardware distintas.

Isso, graças a um modelo de máquinas virtuais que esconde as diferenças de arquitetura.

O usuário não precisa se preocupar se o código está enviando mensagens entre processadores de uma arquitetura semelhante ou diferente.

Portabilidade é fundamental, mas um padrão não seria largamente utilizado se isso fosse obtido em detrimento do desempenho.

Um ponto crucial é que o MPI foi cuidadosamente projetado para possibilitar o desenvolvimento de implementações eficientes.

As decisões de projeto aparentemente foram corretas, já que muitas plataformas executam aplicações MPI com excelentes desempenhos.

Outro objetivo importante para processamento paralelo diz respeito à escalabilidade.

O MPI possibilita uma boa escalabilidade através de várias características de seu projeto.

Por exemplo, uma aplicação pode criar subgrupos de tarefas, que por sua vez, possibilitam operações de comunicação controladas e restritas a alguns delas, possivelmente envolvidas em um escopo específico.

Finalmente, MPI, define um comportamento bem conhecido e mínimo para implementações de troca de mensagens.

Isso tira das costas do programador a preocupação com certas questões de nível mais baixo e o deixa livre para se concentrar na lógica do problema a ser resolvido.

Por exemplo, o MPI garante que a transmissão de mensagens é confiável liberando o programador do trabalho de verificar se cada mensagem foi devidamente recebida.

O maior atrativo do paradigma de troca de mensagens é a sua portabilidade.

Programas escritos dessa forma podem ser executados em multicomputadores de memória distribuída, multiprocessadores de memória compartilhada, redes de workstations e combinações de todos esses.

O paradigma não se torna obsoleto em arquiteturas que combinam as visões de memória distribuída e compartilhada nem por aumentos expressivos na velocidade das tecnologias de comunicação.

Assim é possível e eficiente a utilização desse padrão nas mais diversas variedades de máquinas, incluindo aquelas constituídas de uma coleção de outras máquinas, paralelas ou não, conectadas por uma rede de comunicação.

O padrão também é adequado para programas escritos em um estilo mais restrito de SPMD (Single Program Multiple Data), onde todos os processadores executam o mesmo fluxo de instruções.

Apesar de não ser fornecido suporte explícito para threads, o projeto do MPI não prejudica o seu uso.

O padrão possui ainda uma série de características que melhorariam o desempenho e a escalabilidade de hardwares de comunicação interprocessos especializados.

Assim, espera-se que implementações da interface para essas máquinas sejam criadas em breve.

Enquanto isso, as implementações do MPI para protocolos de comunicação interprocessos padrão Unix possibilitam portabilidade e comunicação eficiente para aplicações em clusters e redes de workstations.

Uma aplicação MPI pode ser visualizada como uma coleção de tarefas concorrentes comunicantes.

Um programa inclui código escrito pelo programador da aplicação que é ligado à biblioteca de funções do MPI.

A cada tarefa é associado um rank.

Esses ranks são utilizados pelas tarefas MPI para identificar umas às outras em cooperações que venham a ser realizadas.

Tarefas MPI podem ser executadas no mesmo processador ou em processadores diferentes de forma concorrente como ilustrado.

Enviar uma mensagem para uma tarefa na mesma ou em uma máquina diferente é transparente para a aplicação.

O MPI automaticamente seleciona o mecanismo de comunicação disponível mais eficiente em uma máquina ou entre máquinas.

O uso de ranks torna todas as operações de cooperação independentes da localização física dos participantes.

No decorrer do texto tarefas MPI também serão chamadas de processos MPI.

Aplicação MPI.

Um grupo de tarefas é uma coleção ordenada na qual cada tarefa é identificada pelo seu rank dentro da ordem.

Para um grupo de n tarefas os ranks variam de 0 até $n-1$.

Grupos de tarefas são usados com duas finalidades importantes.

Em primeiro lugar, eles são utilizados para especificar quais tarefas estão envolvidos em uma operação coletiva de comunicação, como um broadcast, por exemplo.

Em segundo lugar, eles são utilizados para introduzir paralelismo no código de forma que grupos diferentes resolvam partes diferentes da aplicação.

Se isso é feito carregando-se códigos executáveis diferentes para cada grupo, então temos um paralelismo de tarefas MIMD.

Por outro lado, se cada grupo executa um desvio condicional dentro do mesmo executável, temos um paralelismo de tarefas SIMD (paralelismo de controle).

A especificação inicial do MPI adota um modelo estático de tarefas de forma que existe um número fixo delas do início até o fim da execução do programa.

Apesar do modelo de tarefas do MPI ser estático, os grupos são dinâmicos já que eles podem ser criados e destruídos durante a execução e cada tarefa pode pertencer a vários grupos simultaneamente.

No entanto, os membros de um grupo não podem ser alterados e, portanto, é necessário criar um novo grupo para alterar os membros de um já existente.

Em MPI um grupo é um objeto referenciado através de um handle.

O MPI possui rotinas para criar novos grupos especificando os ranks ou particionando um grupo existente através de uma chave.

Outras rotinas consultam o rank de uma determinada tarefa dentro de um grupo específico, testam se uma tarefa faz parte de um determinado grupo, realizam sincronização de barreira em um grupo e consultam seu tamanho e os membros que a ele pertencem.

Contextos de comunicação foram propostos inicialmente para permitir a criação de fluxos de mensagens distintos e bem delimitados entre tarefas, com cada fluxo pertencendo a um único contexto.

Um uso importante dos contextos é garantir que mensagens enviadas em uma fase da aplicação não sejam incorretamente interceptadas em outra fase.

Os contextos representam um critério adicional para seleção de mensagens e possibilitam a criação de espaços de tags de mensagens independentes.

O usuário realiza operações explícitas sobre contextos, já que não existe um tipo de dados visível para eles.

No entanto, os contextos estão associados aos comunicadores de forma transparente ao usuário de maneira que mensagens enviadas através de um determinado comunicador só podem ser recebidas através do comunicador correto correspondente.

O escopo de uma operação de comunicação é especificado pelo contexto de comunicação utilizado e o grupo ou grupos envolvidos.

Em uma operação coletiva ou ponto a ponto entre os membros de um mesmo grupo, apenas este precisa ser especificado e as tarefas fonte e destino são identificadas por seus ranks dentro do grupo.

Em uma comunicação ponto a ponto entre tarefas de grupos diferentes, os dois grupos precisam ser especificados e nesse caso tarefas fonte e destino são identificadas pelo seu rank em seu respectivo grupo.

No MPI um objeto abstrato chamado Comunicador é usado para definir o escopo de uma operação de comunicação.

Comunicadores usados em comunicação intra-grupos ou inter-grupos são chamados de intra ou inter-comunicadores respectivamente.

Um intra-comunicador pode ser considerado como uma ligação entre um contexto e um grupo enquanto um inter-comunicador liga um contexto e dois grupos, um com a tarefa fonte e outro com a tarefa destino.

Os comunicadores são passados como parâmetro para todas as rotinas ponto a ponto ou coletivas para especificar o contexto e o grupo ou grupos envolvidos na operação.

Em muitas aplicações, as tarefas são organizados em uma topologia particular como um grid bidimensional ou tridimensional, por exemplo.

O MPI suporta diversas topologias, especificadas por grafos nos quais tarefas que se comunicam são ligadas por um arco.

Por conveniência, o padrão já inclui suporte explícito para grids cartesianos n-dimensionais.

No MPI um grupo possui uma topologia cartesiana, uma topologia de grafo ou não possui topologia.

Apresenta uma topologia em grafo para cinco tarefas.

Topologia em Grafo para uma aplicação MPI.

O MPI utiliza um rico conjunto de mensagens de envio e recebimento de mensagens e a comunicação entre tarefas envolve os seguintes componentes.

Remetente, geralmente identificado por seu rank.

Receptor, geralmente identificado por seu rank.

Dados da mensagem.

Tag da mensagem.

Ajuda a distinguir múltiplas mensagens entre duas tarefas que devem ser tratadas de forma diferente ou em uma ordem pré-estabelecida.

O comunicador que fornece um contexto para a comunicação.

O MPI, para comunicação ponto a ponto, seleciona mensagens explicitamente através de suas tarefas fonte, tag de mensagens e contextos de comunicação.

A tarefa fonte e a tag podem ser ignoradas na seleção de mensagem, o contexto de comunicação não.

As tarefas fonte e destino são especificadas através de um grupo e um rank.

Para comunicações intra-grupos, o grupo e o contexto são ligados em um intra-comunicador enquanto para uma comunicação inter-grupos a fonte e o destino são ligados em um inter-comunicador.

Assim, uma rotina de send recebe como parâmetros um comunicador, o rank da tarefa destino e o tipo de mensagem para especificar o contexto e o destino de uma mensagem.

Uma rotina de receive recebe os mesmos três parâmetros para selecionar a mensagem que deve receber.

Comunicação ponto a ponto envolve a transmissão de uma mensagem entre um par de tarefas.

O MPI possui uma grande variedade de rotinas de comunicação ponto a ponto, ao contrário de outras bibliotecas de troca de mensagens que em geral possuem apenas um método de comunicação desse tipo.

Isso dá ao programador muito mais controle sobre como as mensagens serão tratadas.

De uma forma geral as rotinas de comunicação ponto a ponto podem ser classificadas quanto ao bloqueio (bloqueantes ou não bloqueantes) e quanto ao modo de comunicação (padrão, síncrono, bufferizado e ready).

Quanto ao bloqueio, as rotinas bloqueantes garantem que a tarefa transmissora ou receptora ficará bloqueada até que a transmissão da mensagem seja completada.

Caso contrário, tem-se uma rotina não bloqueante.

Quanto ao modo de comunicação existem quatro modos possíveis, modo síncrono, modo bufferizado, modo ready e modo padrão.

No modo síncrono, o receptor deve enviar uma confirmação de recebimento de mensagem, de maneira que o transmissor possa ter certeza de que a mensagem foi recebida.

No modo bufferizado, a transmissão de uma mensagem utilizando buffers permite que esta se complete rapidamente já que depois de copiada para o buffer, fica a cargo do sistema transmitir a mensagem quando possível.

No modo ready, a operação é finalizada rapidamente sem a utilização de buffers ou de confirmações da tarefa receptora, objetivando um desempenho melhor para ambientes computacionais específicos.

No modo padrão, o término de uma operação de comunicação pode ou não significar que o receptor foi ativado.

Construções de sincronização são utilizadas para forçar uma determinada ordem de execução entre atividades de tarefas paralelas.

Em alguns casos, em algum ponto da execução, é necessário que algumas tarefas paralelas realizem sincronização com outras tarefas.

No MPI as operações de sincronização ocorrem devido a operações de troca de mensagens bloqueantes e através de operações de barreira.

Quando é utilizada, uma operação de receive bloqueante, força a tarefa que está recebendo a esperar até que a mensagem seja recebida.

Se for utilizado o modo síncrono de comunicação, as duas tarefas devem se encontrar em um ponto de sincronização.

Se tanto a rotina de envio quanto a de recebimento são bloqueantes, então a rotina de comunicação não se completará até que ambos, o remetente e o receptor se encontrem.

Tarefas em um grupo podem também ser sincronizadas em um ponto da execução através de uma barreira.

Nenhuma tarefa pode continuar sua execução além de uma barreira, até que todas as outras tenham chegado até ela.

O grupo pode incluir todas as tarefas ou apenas um subconjunto das tarefas, dependendo do comunicador.

Na prática, quando uma tarefa chama a rotina de barreira ele se bloqueia e permanece assim até que todas as tarefas do subconjunto em questão tenham chamado a mesma rotina.

Sincronização de barreira, para todas as tarefas do comunicador.

Essas rotinas possibilitam comunicação coordenada entre um grupo de tarefas.

O grupo de tarefas e o contexto são especificados pelo intra-comunicador passado como parâmetro para a rotina.

As operações coletivas do MPI forma projetadas de tal forma que sua sintaxe e semântica são consistentes com as das rotinas ponto a ponto.

As operações coletivas do MPI não possuem um argumento tag e assim devem ser chamadas por todos os membros do grupo.

Assim que uma tarefa realiza seu papel na comunicação coletiva, ela pode continuar com outras instruções.

Existem três tipos dessas operações, controle de tarefas, computação global, e movimento de dados.

A função de barreira, pode ser classificada como uma operação coletiva de controle de tarefas.

Serão apresentadas as rotinas de computação global e de movimentação de dados.

O MPI não inclui operações de comunicação coletiva não bloqueantes.

Existem três tipos básicos de rotinas coletivas de movimentação de dados, broadcast, scatter e gather.

Existem duas versões para cada uma delas, Um para todos e todos para todos.

A versão um para todos do broadcast transmite dados de uma tarefa para todas as outras no grupo.

Já na versão todos para todos os dados são transmitidos de cada tarefa para todas as outras.

Assim, cada uma das tarefas termina a operação com o mesmo buffer de saída que nada mais é do que a concatenação dos buffers de envio de cada tarefa na ordem dos ranks.

Apresenta uma operação de broadcast um para todos.

A versão um para todos da rotina scatter envia dados distintos de uma tarefa para todas as outras no grupo.

Essa operação também é chamada de comunicação um para todos personalizada.

No caso da rotina scatter todos para todos cada uma das tarefas envia dados para todas as outras no grupo sendo que esses dados diferem de uma tarefa destino para outra.

É também conhecida como comunicação todos para todos personalizada.

Apresenta a uma operação de scatter.

Os padrões de comunicação na rotina gather são os mesmos da rotina scatter, exceto que a direção do fluxo de dados é reversa.

Na versão um para todos do gather, uma tarefa recebe dados de todas as tarefas do grupo.

Broadcast um para todos.

A tarefa raiz recebe a concatenação do buffer de entrada de todas as tarefas na ordem do rank.

A versão todos para todos do gather é idêntica à versão todos para todos do scatter.

Existem duas rotinas básicas de computação global no MPI, `reduce` e `scan`.

Ambas requerem a especificação de uma função para a computação.

Existe uma versão na qual o usuário seleciona uma função de uma lista pré-definida e outra na qual o usuário fornece um ponteiro para função.

A operação `reduce` combina os elementos no buffer de entrada de cada tarefa no grupo usando a operação especificada e retorna o valor combinado no buffer de saída da tarefa cujo rank foi especificado como o root da operação.

Apresenta uma operação de `reduce`.

Operação de `reduce`.

Existe uma variação da rotina `reduce`, chamada `allReduce` que ao invés de retornar o valor da computação global apenas para a tarefa root, retorna para o buffer de todas as tarefas.

Apresenta uma operação de `allReduce`.

A Existem dois tipos de operações de `scan`, préfixada e posfixada.

O resultado de uma operação de `scan` é diferente em cada tarefa, de acordo com o rank da tarefa.

Por exemplo, sejam T_0, T_1, \dots, T_{n-1} os membros de um grupo contendo os dados d_0, d_1, \dots, d_{n-1} respectivamente e um operador.

O resultado em uma tarefa T_i será $d_0 d_1 d_2 \dots d_i$ para o `scan` prefixado.

No `scan` pós-fixado, o resultado para a tarefa T_i será $d_i d_{i+1} \dots d_{n-1}$.

As operações pré-definidas suportadas pelas operações de `reduce` e `scan` são, máximo, mínimo, soma, produto, AND lógico, AND bitwise, OR lógico, OR bitwise, XOR lógico, XOR bitwise.

Avanços tecnológicos recentes tornaram possível diversos processadores acessarem um único espaço de memória de forma eficiente e recolocaram as arquiteturas de memória compartilhada em foco na área de computação de alto desempenho.

Cresce também nessa área a tendência de agrupar sistemas SMP em clusters em busca de um tempo de processamento ainda mais rápido.

Essa tendência torna desejável que aplicações desenvolvidas para clusters de SMPs sejam portáteis e eficientes.

Aplicações baseadas em troca de mensagens, escritas com bibliotecas como o MPI, possuem naturalmente excelente portabilidade e podem ser utilizadas em clusters de SMPs sem maiores problemas.

No entanto, apesar de ser claro que para a comunicação entre as máquinas SMP, a troca de mensagens é uma boa estratégia, o mesmo não se pode dizer para a comunicação realizada dentro das máquinas SMP.

Teoricamente um modelo de programação específico para memória compartilhada como o OpenMP seria uma estratégia mais eficiente para a comunicação dentro de uma máquina SMP.

Portanto, uma combinação dos paradigmas de memória compartilhada e troca de mensagens na mesma aplicação tem potencial de oferecer um melhor desempenho do que uma abordagem MPI pura.

Em arquiteturas de memória compartilhada (SMP Multiprocessador simétrico), diversos processadores compartilham fisicamente um único espaço de endereçamento de memória.

Isso pode limitar a escalabilidade do sistema porque o acesso à memória se torna um gargalo à medida que o número de processadores aumenta.

Apesar do problema de contenção no acesso à memória, é mais simples programar no modelo de memória compartilhada do que no modelo com troca de mensagens já que o programador não precisa se preocupar com questões como partição e distribuição dos dados entre os processadores.

Representação Memória Compartilhada (SMP).

Um multiprocessador simétrico, no entanto, não atende aos requisitos de desempenho de diversas aplicações muito complexas.

Unir vários sistemas SMP é uma forma de aumentar o número total de processadores disponíveis e conseqüentemente fornecer poder computacional suficiente para esse tipo de aplicação.

Clusters de SMPs podem ser descritos como uma união de memória compartilhada e memória distribuída em uma mesma arquitetura de hardware.

Eles são constituídos de alguns nós SMP conectados através de uma rede, onde cada um deles contém dois ou mais processadores que compartilham memória fisicamente.

Alguns sistemas possuem suporte de hardware e de software para que um processador acesse a memória de um nó remoto diretamente.

No entanto, na maioria dos casos, clusters de SMPs necessitam de trocas de mensagens explícitas para a comunicação entre os nós do cluster.

Clusters de SMPs, introduzem um nível adicional na hierarquia de memória constituído por um conjunto de sistemas de memória compartilhada.

Infelizmente, esse nível adicional de memória torna o comportamento desses sistemas menos previsível e a sua programação mais difícil.

A utilização de um modelo híbrido de programação se beneficia das vantagens dos dois modelos.

Por exemplo, a programação híbrida nos permite o uso das políticas de particionamento explícito de dados, característico do paradigma de troca de mensagens junto com o paralelismo de grão fino, característico do paradigma de memória compartilhada.

Nesse trabalho utiliza-se a biblioteca MPI para a programação da parte de troca de mensagens e as diretivas de compilação e bibliotecas do padrão OpenMP para a programação da parte de memória compartilhada.

Esses são os dois principais padrões de programação paralela e conseqüentemente, são os mais utilizados para estratégias híbridas de programação.

Por isso, o restante do trabalho se concentra mais especificamente nesses dois modelos.

A maioria das aplicações híbridas apresenta um modelo hierárquico com a paralelização MPI em um nível superior e a paralelização OpenMP em um nível abaixo.

Mostra um array bidimensional que foi dividido entre quatro processos MPI.

Cada um desses subarrays foi então subdividido entre threads OpenMP.

Esse modelo hierárquico é um mapeamento muito próximo à arquitetura física de um cluster de SMPs.

Em um programa híbrido típico, o MPI é iniciado e finalizado de forma usual, com as rotinas MPI_INIT e MPI_FINALIZE.

Uma região paralela do OpenMP ocorre entre essas chamadas, criando uma ou mais threads adicionais em cada processo.

Se, o programa for executado utilizando quatro processos MPI e duas threads OpenMP então o fluxo de execução seria observado.

Para garantir a portabilidade da aplicação, o ideal é realizar as chamadas de comunicação do MPI fora de regiões paralelas do OpenMP.

Em geral, as chamadas de comunicação são realizadas naturalmente fora dessas regiões, já que a paralelização OpenMP costuma ser transparente para o processo que participa de uma execução MPI.

Entretanto, quando é inevitável realizá-las dentro da região paralelas é imprescindível usar alguma rotina de sincronização que garanta que a comunicação seja realizada por apenas uma das threads.

No OpenMP pode-se utilizar as rotinas de sincronização CRITICAL, MASTER ou SINGLE para esse propósito.

Fluxo de execução com quatro processos MPI e duas threads OpenMP.

Ao escrever uma aplicação híbrida OpenMP/MPI é importante considerar também como cada paradigma pode paralelizar o problema e como combinar os dois para alcançar o melhor resultado.

O problema do grid bidimensional envolve uma decomposição do grid em uma dimensão com MPI e em outra com OpenMP pois a conversão de problemas com estrutura hierárquica é simples de ser utilizada.

Outros tipos de problema podem exigir soluções mais complexas.

Além de combinar explicitamente troca de mensagens e multithreading em um programa, outros estilos híbridos de programação incluem o uso de versões de bibliotecas de troca de mensagens que utilizam a memória

compartilhada para comunicação dentro do SMP e ambientes de memória virtual compartilhada.

Para executar em cluster de SMPs um sistema pode utilizar simplesmente MPI para comunicação dentro dos nós SMP.

No entanto, um grande número de aplicações não escala bem com MPI, do qual grande parte é de aplicações que envolvem um balanceamento de carga complexo.

Com um código híbrido, o MPI é utilizado apenas para comunicação entre nós e, uma vez que o OpenMP é menos suscetível a problemas de balanceamento, há uma grande possibilidade de um melhor desempenho.

O OpenMP geralmente apresenta melhor desempenho em aplicações de grão fino, onde a quantidade de comunicação MPI pode ser superada pela quantidade de comunicação OpenMP.

Quando uma aplicação exige boa escalabilidade, com uma paralelização de grão fino, um código híbrido pode ser mais eficiente.

É claro que uma aplicação apenas OpenMP teria um desempenho melhor, mas para clusters de SMPs o MPI ainda é necessário para comunicação entre os nós.

Reduzindo o número de processos MPI a escalabilidade, na maioria das situações, é melhorada.

Códigos que usam uma estratégia de dados replicados, geralmente sofrem de limitações de memória e de baixa escalabilidade devido a comunicações globais.

Utilizando um modelo de programação híbrido em um cluster de SMPs, o problema pode ser limitado à memória do SMP e não à memória de dois processadores como é o caso da paralelização MPI pura.

Essa é uma vantagem clara que permite o estudo de problemas de tamanhos maiores. Algumas aplicações MPI exigem um número específico de processos para serem executadas.

Por exemplo, um código que limita o número de processos MPI para algumas combinações ou um número fixo ou que escalam apenas em potências de 2.

Isso pode criar problemas de duas formas.

Em primeiro lugar, o número de processos necessários pode não ser igual ao número de máquinas.

Se for muito grande, o número de processos pode tornar a execução inviável e se for muito pequeno, torna a utilização dos recursos de hardware ineficiente com máquinas não utilizadas.

Se um código híbrido é utilizado, a estratégia de decomposição natural do MPI pode ser usada, executando o número desejado de processos e threads OpenMP usadas para distribuir o trabalho entre processadores de forma a permitir uma utilização eficiente dos recursos.

A técnica chamada de balanceamento de poder computacional dinamicamente ajusta o número de processos trabalhando em uma determinada computação.

A aplicação é escrita em modo híbrido com diretivas OpenMP dentro do código MPI.

Inicialmente o trabalho é distribuído entre processos MPI, mas quando a carga em um processador dobra, o código usa diretivas MPI para criar uma nova thread em outro processador.

Assim, sempre que a carga de processamento de um processo MPI se torna excessiva, o trabalho pode ser redistribuído.

A primeira escolha a ser feita no desenvolvimento de uma aplicação para um cluster de SMPs é entre um modelo unificado e um modelo híbrido.

No modelo de programação unificado, o programador usa uma única API para descrever tanto a comunicação dentro do nó, quanto a comunicação entre nós diferentes.

Todos os sistemas de troca de mensagem ou DSM pertencem a essa categoria.

Os modelos híbridos, por sua vez, misturam memória compartilhada dentro do multiprocessador e troca de mensagens entre nós.

MPI + OpenMP, MPI+threads são dois exemplos de modelos híbridos.

O desempenho de tais modelos depende pelo menos de três fatores, o compartilhamento de suporte de comunicação (memória do sistema e interface de rede) entre processadores, isto é, como a latência por processo e a largura de banda evoluem quando vários processadores utilizam a mesma interface de rede.

O grau de paralelismo de memória compartilhada que pode ser alcançado com o modelo híbrido e o speed-up obtido na seção paralela.

A partir de um código MPI já existente, a abordagem mais simples consiste em paralelizar com o OpenMP os loops dentro do código MPI.

Essa abordagem é chamada de paralelização de grão fino, em nível de loop ou ainda de paralelização incremental.

Várias abordagens diferentes podem ser utilizadas.

A primeira possibilidade consiste em paralelizar os loops na parte procedural do MPI sem quaisquer otimizações.

Apenas a corretude da versão paralela contra a versão procedural é verificada.

Mas a abordagem incremental pode ser melhorada significativamente aplicando várias otimizações manuais (permutação de loops, troca de loops, uso de variáveis temporárias).

Essas otimizações são necessárias por exemplo, para transformar um loop que não pode ser paralelizado em um que pode ser, ou então para melhorar a eficiência do loop paralelo, evitando falso compartilhamento e diminuindo o número de pontos de sincronização.

Uma outra questão é a escolha dos loops a serem paralelizados.

Uma opção é paralelizar todos os loops.

Essa claramente não é uma boa opção pois pode incluir loops que não contribuem de forma significativa para o tempo total de execução.

Nesse caso a paralelização pode tornar essas porções de código mais lentas que a versão serial.

A melhor alternativa consiste em selecionar através de uma análise detalhada os loops que influem significativamente para o tempo total de execução.

Mostra um fluxograma de um processo ideal para paralelização de grão fino.

Fluxograma de paralelização de Grão Fino.

Uma outra abordagem utilizada é a paralelização de grão grosso (ou SPMD).

Nessa abordagem, o OpenMP ainda é utilizado para obter vantagem sobre a memória compartilhada dentro do nó SMP, mas um estilo de programação SPMD é utilizado.

Nesse caso, o OpenMP é usado para criar N threads no início do programa principal, onde essas threads se comportam exatamente como um dos processos MPI original, mas se comunicam por memória compartilhada.

Da mesma forma que em um modelo SPMD de troca de mensagens, o programador deve tomar cuidado com algumas questões, distribuição de dados entre as threads, distribuição de trabalho e coordenação entre elas.

Já que a distribuição de dados é através da memória compartilhada, ela implica apenas a atribuição de diferentes regiões da mesma estrutura de dados para diferentes threads em execução.

Geralmente, o programador calcula a região atribuída a cada thread com base em alguma função de hash sobre o número (identificador) da thread.

A coordenação das threads envolve a gerência de seções críticas e barreiras, além do uso de ou a diretiva MASTER do OpenMP ou da chamada de biblioteca `omp_get_threads_num` para instruções condicionais (exclusivas de uma ou grupo de threads).

Existem poucas publicações sobre resultados em paralelização OpenMP+MPI de grão grosso.

Normalmente, quando as chamadas MPI são realizadas fora da região paralela do OpenMP a única thread que realiza comunicação é thread deste.

Enquanto a comunicação está sendo realizada por uma thread todas as outras estão bloqueadas esperando a conclusão da operação, de forma que o tempo de CPU não é aproveitado da melhor maneira possível.

Existe a possibilidade de realizar a sobreposição de comunicação e computação através de threads concorrentes.

Nesse caso, enquanto a comunicação é realizada pela thread mestre, outras threads não comunicantes, podem executar algum código da aplicação.

Essa categoria exige que o código da aplicação seja separado em duas partes, código que pode ser sobreposto com comunicação e código que deve esperar até que as operações de comunicação sejam finalizadas.

A sobreposição de comunicação e computação é útil para obter o máximo aproveitamento do hardware.

Porém, há algumas desvantagens.

Em primeiro lugar, a maioria das aplicações não está preparada para distinguir entre instruções que podem ser realizadas antes do término das operações de envio e recebimento de dados e as que podem.

Em segundo lugar, um modelo de programação de grão grosso é necessário, o que geralmente exige que a aplicação seja reescrita para que distribuição de trabalho entre as threads seja realizada manualmente baseada nos ranks de cada uma.

E finalmente, é preciso implementar alguma funcionalidade de balanceamento de carga para compensar as diferentes cargas de comunicação e computação entre as threads.

A paralelização OpenMP dentro de processos MPI pode causar overhead adicional.

A criação das regiões paralelas e a sincronização ao seu final induzem algum trabalho adicional, principalmente se uma abordagem de grão fino é utilizada.

Esse overhead pode ser reduzido com uma estratégia de grão grosso para a paralelização, a região paralela é inicializada apenas uma vez no início da aplicação e diretivas OMP MASTER e BARRIER são utilizadas para sincronização antes e depois da comunicação MPI.

Se não for possível paralelizar todo o trabalho do MPI ou se a paralelização OpenMP não puder ser balanceada satisfatoriamente, então o speedup também será reduzido em virtude da Lei de Amdahl.

Existem muitos trabalhos que estudam o desempenho de modelos de programação paralela.

Em, por exemplo, foi realizado um estudo de desempenho sobre arquitetura de troca de mensagens em clusters SMP.

Em foi realizada uma comparação entre os modelos de troca de mensagens e de memória compartilhada.

Ambos são focados em modelos tradicionais e não abordam modelos híbridos.

A respeito desses modelos existem alguns trabalhos como, onde foram realizados testes para 6 diferentes benchmarks.

Nesse trabalho compara-se o desempenho de versões puras MPI dos programas com versões híbridas modificadas com diretivas OpenMP para diferentes tamanhos de dados e arquiteturas de cluster.

Em um modelo híbrido de programação com MPI e OpenMP foi aplicado para a solução de um sistema linear de equações em um cluster SMP.

Dois trabalhos interessantes são onde os autores discutem diversos aspectos relevantes no desenvolvimento de aplicações em um modelo híbrido de programação, além de realizarem uma grande variedade de testes de desempenho.

Particularmente em foi desenvolvida uma aplicação híbrida para um código baseado no método dos elementos finitos para simulações climáticas e geológicas no computador Earth Simulator.

Alguns trabalhos como apresentam modelos híbridos baseados em ambientes de memória distribuída compartilhada (DSM) como uma alternativa para unificar os modelos de memória.

A aplicação escolhida para este estudo de arquiteturas híbridas de programação paralela trata-se da simulação de um problema de elasticidade linear baseada no método dos elementos finitos (MEF).

Essa simulação utiliza o método dos gradientes conjugados para a solução de um sistema de equações.

Pelo método dos elementos finitos um modelo matemático descrito por equações diferenciais parciais em um domínio contínuo, é convertido em um modelo discreto de elementos finitos, com um número finito de graus de liberdade (GLs).

O trabalho realizado em apresenta resultados de execuções de uma versão 3 D do método dos gradientes conjugados (MGC) na solução das equações de equilíbrio de um problema de elasticidade tridimensional, utilizando programação paralela por troca de mensagens.

Esse trabalho foi baseado no algoritmo desenvolvido em.

A aplicação, originalmente desenvolvida para o paradigma de troca de mensagens, foi transformada em uma aplicação híbrida para se beneficiar da comunicação via memória compartilhada dentro do nó SMP, com a comunicação entre os nós realizada por troca de mensagens.

Para essa transformação foi utilizada uma abordagem de grão fino com a paralelização dos loops responsáveis pela solução do sistema de equações através de threads OpenMP.

Estas threads compartilham a mesma memória física.

Isto possibilita que o número de tarefas MPI seja reduzido e conseqüentemente seja diminuído o volume de comunicação inter-processos.

No caso um cluster com 8 máquinas SMP de dois processadores cada uma, utilizando-se todos os processadores em uma execução da aplicação, tem-se 16 tarefas MPI (duas por nó SMP) que trocam mensagens durante seu ciclo de vida.

No programa modificado, pode-se iniciar a execução da aplicação com apenas 8 tarefas MPI (uma por nó SMP) e à medida que regiões paralelas são encontradas uma nova thread t2 é criada no processador ocioso do nó e executará algum trecho de código em paralelo com a thread t1 que corresponde ao programa principal da tarefa MPI.

A comunicação entre t_1 e t_2 ocorre através da memória, compartilhada pelos processadores do mesmo nó SMP.

A consequência prática é que a comunicação MPI entre duas tarefas é substituída por uma mais rápida na memória compartilhada.

Um material é elástico se ele deforma quando uma carga é aplicada a ele, mantém uma deformação constante enquanto a carga é mantida constante e retorna ao seu formato original, não deformado, quando a carga é removida.

As propriedades mecânicas de um material podem ser verificadas por um teste de tensão no qual uma barra ou cilindro de um material com comprimento L e seção A é fixada em uma de suas extremidades e sujeita a uma carga F .

Enquanto a carga é gradualmente aumentada, o corpo de prova irá sofrer deformação até se romper em dois pedaços.

Normalmente deseja-se compreender o comportamento desses materiais para diferentes tamanhos e formas de objetos, especialmente como a carga F aplicada se relaciona com a deformação ocorrida e que cargas o material suporta sem que ocorra uma fratura.

Teste de Tensão.

Sólidos constituídos de materiais dúcteis, como aço e ferro, possuem uma relação entre tensão e deformação bastante complexa, envolvendo comportamentos mecânicos diversos.

Tais sólidos, quando sujeitos a tensões suficientemente pequenas exibem comportamento elástico linear e apresentam valores de deformações e deslocamentos diretamente proporcionais às forças aplicadas.

Apresenta o gráfico dos regimes de deformação de materiais dúcteis.

No gráfico pode-se ver que quando o sólido sai do regime elástico linear e entra no regime plástico a relação entre tensão e deformação deixa de ser constante.

Além disso, no regime plástico o sólido não retorna à sua forma original quando são retirados os esforços mecânicos.

Este trabalho está restrito a problemas no regime elástico linear.

Gráfico de tensão-deformação de materiais dúcteis.

Para materiais elásticos lineares, a forma geral de uma equação constitutiva é, A constante de proporcionalidade é determinada experimentalmente e diz respeito às propriedades que constituem o material em questão.

Problemas de elasticidade linear podem ser resolvidos através do método dos gradientes conjugados aplicado a um sistema de equações gerado pelo método dos elementos finitos.

O método dos elementos finitos, aplicado a um problema estático linear, gera um sistema de equações algébricas, onde K é uma matriz simétrica e positiva definida constante e representa a matriz de rigidez do material, F é um vetor de forças externas aplicadas ao

sólido e u é um vetor de incógnitas e representa os deslocamentos resultantes da aplicação das forças externas.

Os vetores u e F possuem dimensão igual ao número de graus de liberdade do sólido.

K é uma matriz quadrada de dimensões também iguais ao número de graus de liberdade do sólido.

No campo da engenharia, a análise de problemas complexos requer a modelagem matemática do sistema físico e muitas vezes a solução analítica é complexa e geralmente impossível.

Nesse caso é necessário recorrer ao uso de técnicas numéricas.

O método dos elementos finitos é uma técnica numérica extremamente poderosa para a solução de problemas nas áreas de Estruturas, Transferência de Calor, Mecânica dos Fluidos, etc.

O procedimento geral do método é que cada estrutura ou corpo é dividido em elementos menores de dimensões finitas, chamados elementos finitos.

O corpo original é então considerado como uma composição desses elementos.

Esses elementos são conectados entre si através de junções chamadas nós ou pontos nodais para formar a estrutura completa.

As propriedades desses elementos individuais são então formuladas e a partir dessas, as propriedades do corpo como um todo, são obtidas.

O método dos elementos finitos é uma técnica numérica e as respostas obtidas com ele não são soluções exatas e sim soluções aproximadas.

Entretanto, utilizando procedimentos apropriados e tendo à disposição recursos computacionais poderosos, é possível obter um alto grau de precisão.

As seguintes informações são associados a um elemento finito individual.

Esses dados são usados por aplicações de elementos finitos para realizar os cálculos, Dimensionalidade, Os elementos podem ter uma, duas ou três dimensões espaciais.

Cada elemento possui um conjunto de pontos distintos chamados pontos nodais ou apenas nós.

Nós servem a dois propósitos, definir a geometria do elemento e os graus de liberdade.

Eles se localizam nos cantos ou nas terminações dos elementos.

A geometria do elemento é definida pelo posicionamento dos pontos nodais.

Na prática, a maioria dos elementos usados possui geometrias simples.

Em uma dimensão, elementos são geralmente linhas retas ou segmentos curvos.

Em duas dimensões eles costumam ter formato triangular ou quadrilateral.

Em três dimensões costumam ser tetraedros, pentaedros ou hexaedros.

Mostra algumas geometrias típicas em uma, duas e três dimensões.

Os graus de liberdade determinam a forma do elemento.

Eles também funcionam como junções através das quais elementos adjacentes são conectados.

Graus de liberdade são definidos como valores de variáveis primárias nos pontos nodais.

No caso de aplicações na área de estruturas, essas variáveis primárias especificam, por exemplo, o deslocamento ocorrido em cada ponto nodal.

Há sempre um conjunto de forças nodais em uma relação de um para um com graus de liberdade.

Em elementos mecânicos a correspondência é estabelecida em termos de energia.

Para um elemento mecânico, são as relações que especificam as propriedades do material.

Em uma barra elástica linear, considerando-se a análise estrutural, é suficiente especificar o módulo elástico E .

Geometrias típicas de elementos finitos em uma, duas e três dimensões.

O Método dos Gradientes Conjugados (MGC) é um dos mais populares métodos iterativos para a resolução de grandes sistemas de equações lineares.

O MGC é efetivo para sistemas, onde x é um vetor desconhecido, b é um vetor conhecido e A é uma matriz simétrica, positiva definida e quadrada.

Tais sistemas de equações aparecem em diversos problemas importantes da Matemática, da Física e da Engenharia.

Métodos iterativos como o MGC, são adequados para uso com matrizes esparsas.

Se A é densa, a melhor alternativa é fatorá-la e resolvê-la por retrosubstituição.

Resolver as equações combinadas da Equação 64 é equivalente a encontrar o mínimo da forma quadrática correspondente ao sistema, Devido ao relacionamento entre a matriz A e a função escalar $f(x)$, é possível ilustrar algumas fórmulas da álgebra linear com imagens, mais intuitivas.

Por exemplo, a matriz A é chamada positiva definida se a seguinte propriedade é verdadeira para qualquer vetor x .

Mostra as formas quadráticas para matrizes do tipo positiva definida, negativa definida, positiva indefinida e indefinida.

A implicação é que se a matriz é positiva definida como em 64 então, ao invés de resolver o sistema da Equação 63, pode-se obter o mesmo resultado encontrando-se o mínimo da sua função quadrática.

Com o método dos gradientes conjugados isso pode ser realizado em n ou mais passos, de forma iterativa.

Formas quadráticas para tipos diferentes de matrizes.

Positiva definida.

Negativa definida.

Positiva indefinida.

Indefinida.

A busca do mínimo pode ser ilustrada através do método dos gradientes e basicamente inclui os seguintes passos, o gradiente é calculado no ponto inicial $x(0)$ e o movimento continua, na direção de busca, na linha do antigradiente enquanto a função objetivo continua decrescendo.

No ponto onde a função pára de decrescer, o gradiente é calculado de novo e o movimento continua em outra direção.

O processo continua até que o ponto mínimo é encontrado.

Busca de mínimo no método dos gradientes conjugados.

A seguir serão apresentados os funcionamentos da versão pura MPI e híbrida do código de elementos finitos e explicados seus funcionamentos.

A entrada para a aplicação é uma malha constituída de valores nodais dos elementos finitos.

Essa malha é gerada pela aplicação GID para a geometria selecionada com as restrições e forças impostas a ela.

O GID é uma ferramenta para modelagem geométrica para simulações numéricas que utilizam o método dos elementos finitos.

Dentre as informações que descrevem uma malha têm-se as propriedades do material (módulo de elasticidade, coeficiente de Poisson), informações de forças e restrições, coordenadas de cada nó e conectividades.

Inicialmente, essa malha é particionada através da biblioteca METIS.

O número de partições é igual ao número de tarefas MPI.

Assim, cada tarefa manipula aproximadamente o mesmo volume de dados.

Cada tarefa então monta sua matriz de rigidez e algumas estruturas de dados auxiliares baseado na partição recebida.

E então, o método dos gradientes conjugados paralelo é aplicado para resolver o sistema de equações.

Na implementação paralela do método dos gradientes conjugados a matriz de rigidez A deve ser subdividida e distribuída entre as várias tarefas.

Cada tarefa possui também uma subdivisão do vetor x e uma subdivisão do vetor u .

Quando ocorre divisão alguns graus de liberdade ficam compartilhados por mais de uma tarefa.

Esses graus de liberdade são chamados de graus de liberdade de fronteira (ou compartilhados), enquanto todos os demais são chamados de graus de liberdade internos.

Para fazer essa distinção, a matriz A de cada processador é subdividida em quatro A_p , A_s , B_p e B_p .

$T A_p$ é uma matriz quadrada com dimensão igual ao número de graus de liberdade internos de cada subdomínio e nela estão armazenados os valores de rigidez referentes aos graus de liberdade internos.

A matriz A_s também é quadrada e tem dimensão igual ao número de graus de liberdade compartilhados e nela estão os valores da matriz de rigidez para os graus de liberdade de fronteira.

A dimensão da matriz B_p é função tanto dos graus de liberdade internos quanto dos de fronteira.

Nesta, estão os valores de rigidez que relacionam os graus de liberdade privados com os de fronteira.

B_p é a matriz transposta de B_p .

Fluxo da Aplicação.

Devido à reestruturação da matriz A , os vetores x e b também são reformulados.

Cada um deles é separado em outros dois.

O vetor b é desmembrado em b_p e b_s , sendo que o primeiro tem dimensão igual aos graus de liberdade internos e o segundo aos de fronteira.

O mesmo procedimento se aplica ao vetor x .

O MGC é um método iterativo que busca a convergência para a solução de um sistema de equações.

Esse método, não calcula um valor exato e sim uma solução aproximada para o sistema.

Para isso especifica-se um erro máximo que representa o critério de parada do algoritmo.

Quando é obtida uma solução com um erro menor ou igual a, assume-se que a solução está próxima o suficiente do valor real e ela pode ser considerada a solução para o sistema de equações.

Inicialmente arbitra-se um valor aproximado para o vetor de resposta x .

O valor arbitrado é então introduzido no sistema de equações para obter o resíduo.

Se o resíduo é menor ou igual ao erro, o algoritmo pára e o valor de x arbitrado inicialmente é considerado a resposta para o método.

Caso contrário o método entra no loop de resolução.

Em cada iteração, o resíduo é recalculado e se ele é menor ou igual ao erro, o algoritmo pára.

Não é possível prever quantas iterações o método realizará para chegar à resposta final.

Eventualmente, o método pode falhar em convergir para uma solução e entrar em um loop infinito e por isso o número de iterações é limitado a um valor máximo.

A computação do método dos gradientes conjugados envolve um produto de matriz por vetor, duas operações de produto interno, três ajustes de vetores e uma operação de atualização do vetor x .

O algoritmo é de ordem $O(n)$ e quando ele é executado por p processadores, cada um deles resolve um subproblema em tempo $O(n/p)$.² As tarefas MPI realizam comunicação apenas no método dos gradientes conjugados.

Se uma malha é dividida entre alguns processadores, eles compartilham graus de liberdade em suas fronteiras e devem realizar comunicação intensiva para trocar informações durante os cálculos.

Além disso, para calcular o erro, é necessário utilizar rotinas de computação global nos dados espalhados no diversos processadores.

A comunicação interprocessos nessa aplicação é restrita a dois procedimentos, o produto interno e a atualização dos graus de liberdade compartilhados.

O procedimento de produto interno é realizado em vetores localizados em tarefas diferentes.

Para isso é utilizada uma rotina de computação global de soma (All-Reduce) para calcular um escalar.

No procedimento de atualização dos graus de liberdade, cada tarefa envia para cada um das tarefas com os quais compartilha graus de liberdade de fronteira o seu valor para esses graus.

Cada um deles então, soma o valor que possui com os valores que recebeu.

Para obter a versão híbrida da aplicação foi utilizada uma abordagem incremental para obter uma versão de grão fino da aplicação MPI original.

Inicialmente, foi realizada uma análise dos trechos de código que consumiam mais tempo para a execução.

A aplicação híbrida foi então desenvolvida através da modificação do código com diretivas OpenMP.

A partir de alguns testes iniciais, a paralelização foi gradualmente refinada para garantir que as threads fossem usadas para computação apenas nas partes do programa onde o benefício da execução fosse maior que o custo de criação das threads.

Isso naturalmente exclui seções de execução muito rápida e coloca o foco em trechos de código de longa duração e que podem ser executados concorrentemente (sem dependência de dados).

Essa análise levou à paralelização da parte do programa que executa o cálculo do método dos gradientes conjugados já que essa consome 90% do tempo total da aplicação.

Quanto à comunicação, não ocorre dentro das regiões paralelas de modo que a comunicação MPI da aplicação não é alterada.

Nas regiões modificadas com OpenMP as diversas threads que realizam operações em paralelo se comunicam através da memória compartilhada dentro do nó SMP.

O balanceamento de carga entre os nós SMP na aplicação depende diretamente da distribuição realizada nas tarefas MPI.

Entretanto durante a execução, em pontos onde a computação é intensiva, a carga é subdividida entre as threads dentro do nó SMP, de acordo com as diretivas OpenMP inseridas no código.

Para um cluster com 16 processadores (2 por nó SMP) e que utiliza todos os processadores em uma execução da aplicação, na versão pura MPI o domínio será dividido em 16 partes (para 16 tarefas MPI, duas por nó SMP).

Na implementação MPI/OpenMP, no entanto, o domínio é dividido em 8 partes (usando uma tarefa MPI por nó SMP e duas threads nas regiões paralelas, uma por processador).

A paralelização OpenMP acontecerá então em nível de loop, com cada thread realizando metade dos cálculos em paralelo.

Nesse caso, se a paralelização OpenMP é eficiente para trechos de código que consomem grande quantidade de tempo, assumindo que a comunicação dentro do nó SMP é melhorada com a comunicação de memória compartilhada, então a versão híbrida deverá apresentar um desempenho melhor.

A questão que se apresenta é se a paralelização OpenMP implementada é eficiente o bastante para tirar vantagem desse efeito na prática.

Essa eficiência também deve ser tal que não seja prejudicada pelo overhead introduzido pela criação das threads nas regiões paralelas.

Apresenta a comparação entre a decomposição de dados na versão pura MPI (68) e a versão híbrida MPI/OpenMP (68).

Decomposição de dados, versão pura, versão híbrida.

Este capítulo apresenta os resultados experimentais obtidos a partir dos testes realizados para a versão híbrida da aplicação, os resultados obtidos para a versão original e a comparação entre eles.

Fazem parte desses dados informações como speedup, tempo de execução e tempo de comunicação.

São apresentadas as configurações de hardware e é apresentada a metodologia utilizada.

Seguem-se a essas, as medidas de speedup relativo, a comparação dos tempos de execução, é apresentada a comparação entre os tempos de execução divididos em comunicação e computação e finalmente é apresentada a interpretação dos resultados.

Os resultados foram obtidos em um cluster de SMPs constituído de 8 máquinas dual-processor AMD rodando o Sistema Operacional Linux Red Hat 9.

Segue-se a configuração detalhada do cluster.

Especificação do Cluster.

Aplicamos inicialmente a abordagem incremental ao código original MPI para obter a versão híbrida da aplicação.

O código foi então instrumentado para distinguir entre tempo de comunicação e tempo de computação.

O tempo de comunicação inclui as chamadas de sincronização como sincronizações de barreira, por exemplo.

Essa separação é importante para investigar nossa hipótese inicial de que com a solução híbrida, o tempo total consumido com comunicação seria reduzido.

Neste capítulo, por simplificação de notação, na comparação dos dois modelos chamamos processos MPI e threads OpenMP de tarefas.

Assim uma execução pura, com dois processos MPI é considerada como possuindo duas tarefas e uma execução híbrida com duas threads OpenMP também.

As medidas foram realizadas para tamanhos diferentes de malhas de uma mesma geometria.

Para cada uma das malhas, variamos o número de tarefas entre 2, 4, 8 e 16.

Na versão pura MPI, cada nó SMP inicia sua execução com dois processos MPI e continua com esse número constante até o final da execução.

Na versão híbrida OpenMP/MPI por outro lado, em cada nó SMP apenas um processador inicia sua execução com um thread e à medida que, durante o fluxo do programa, são encontradas regiões paralelas, uma outra thread é criada para no processador disponível executar código em paralelo com a thread principal.

A geometria escolhida para os testes foi um bloco.

Uma das faces do bloco foi fixada para ter seu movimento restringido e foi aplicada uma carga concentrada de 5000 N na direção negativa do eixo z em um nó localizado em um dos vértices de forma a provocar alguma deformação.

Geometria Seleccionada.

Os resultados foram obtidos para quatro malhas tridimensionais formadas por tetraedros lineares.

Malhas utilizadas.

Apresenta a geometria dividida em 933 elementos para a malha 2 e apresenta a deformação obtida no bloco.

Malha dois com 933 nós e 4014 elementos.

Deformação obtida.

As medidas de tempo foram obtidas indiretamente com o uso da instrução assembly RDTSC que conta o número de ciclos de clock ocorridos desde o instante em que a máquina foi ligada ou desde sua última reinicialização.

As grandes vantagens dessa estratégia de instrumentação são a alta resolução nas medidas de tempo obtidas e o baixo overhead introduzido na aplicação para realizá-las.

Apresenta os tempos de execução em segundos para as duas versões da aplicação para as quatro malhas utilizadas.

Apresenta os gráficos com as comparações dos tempos de execução.

Valores de tempo total de execução (em segundos).

Gráfico de comparação do tempo total de execução para o programa híbrido e para o programa puro MPI.

Apresentam o tempo de execução decomposto em tempo de computação e tempo de comunicação, respectivamente para as malhas 1, 2, 3 e 4.

Apresenta o gráfico do tempo decomposto para as quatro malhas.

Para uma melhor visualização, a figura 76 apresenta o gráfico com a comparação apenas do tempo de comunicação para as duas versões da aplicação e para cada uma das malhas utilizadas.

Gráfico de comparação do tempo de comunicação para o programa híbrido e para o programa puro MPI.

Apresenta os valores das medidas dos speedups relativos para as quatro malhas utilizadas na versão híbrida e na versão pura MPI.

Apresenta os gráficos com as comparações dos speedups.

Valores de speedup relativo.

Gráfico de speedup relativo para programa híbrido e puro MPI.

Apresenta os valores das medidas de desempenho para as quatro malhas utilizadas.

Apresenta os gráficos com as comparações de desempenho.

Valores de desempenho.

Os resultados experimentais mostram que na maioria dos casos o desempenho geral da versão pura MPI é melhor que o desempenho da versão híbrida OpenMP/MPI.

Referente ao speedup relativo, pode-se perceber que a versão híbrida alcança um bom desempenho, no entanto na maioria das situações, ele não é melhor do que aquele alcançado pela versão original em MPI da aplicação.

Para o tempo total de execução, ocorre o mesmo e a versão MPI obtém um melhor tempo de execução na maioria dos casos.

Gráfico de desempenho para o programa híbrido e puro MPI.

Entretanto, para as malhas menores e com uma grande quantidade de máquinas, o tempo de execução da versão híbrida foi menor.

Isso se deve ao fato de que para as execuções com essas malhas o tempo consumido é pequeno e a redução no tempo de comunicação é maior que o overhead introduzido.

Essa relação pode ser melhor visualizada observando-se os tempos de comunicação e de computação separadamente.

Pode-se perceber que, em 99% das situações a versão híbrida possui tempo de comunicação menor que a versão pura MPI.

No entanto, o tempo de computação da versão híbrida é maior na maioria dos casos o que indica que o programa híbrido, apesar de diminuir o tempo de comunicação, introduziu um overhead no tempo de execução maior do que o benefício obtido.

A introdução desse overhead, excessivo, pode ser explicada pelas principais fraquezas do modelo OpenMP que são, 1) A multiplicação de regiões paralelas OpenMP implica em um custo de gerência de thread considerável.

As regiões paralelas provocam uma má utilização da hierarquia de memória.

O primeiro fator é especialmente importante nessa aplicação porque foi utilizada uma abordagem de grão fino.

Além disso, pelo fato da aplicação usar um método iterativo para resolver o sistema de equações, threads precisam ser criadas muitas vezes.

A respeito do tamanho das malhas, a maior malha utilizada levou um tempo total seqüencial de aproximadamente 678 segundos.

Seria interessante a investigação do comportamento da aplicação para malhas maiores e que conseqüentemente levassem mais tempo para serem executadas, porém execuções de malhas maiores que 3717 exigem mais memória do que está disponível atualmente no cluster.

É importante ressaltar que os resultados e observações obtidos a partir dessa versão híbrida da aplicação dizem respeito à abordagem de grão fino utilizada.

Uma abordagem de grão grosso pode levar a resultados diferentes.

Existem na literatura casos que obtiveram bons e maus resultados utilizando cada uma das duas abordagens.

Uma abordagem de grão fino, no entanto possui algumas vantagens sobre a abordagem de grão grosso, o que levou à escolha dessa a técnica para este trabalho.

Dentre essas vantagens, talvez a mais importante é que a abordagem de grão fino possibilita a paralelização incremental do código.

Isso significa que um código MPI já existente pode ser transformado em uma aplicação híbrida sem a necessidade de modificações na estrutura original do programa, apenas com a anotação de trechos do código com diretivas.

Nessa abordagem pode-se inclusive obter as duas versões da aplicação a partir do mesmo código, compilado com ou sem suporte ao padrão OpenMP.

A abordagem de grão grosso, por outro lado exige que a aplicação seja completamente reescrita para adotar um modelo SPMD para as threads OpenMP.

Modelos híbridos de programação para clusters de SMPs foram usados em diversas aplicações.

Apesar de em alguns trabalhos como terem obtido um ganho de desempenho comparado a uma versão pura da mesma aplicação, na maior parte dos casos, observou-se uma perda de desempenho.

Essa perda de desempenho foi reportada em trabalhos como entre outros.

Em, um problema que utiliza o método dos gradientes conjugados, a solução híbrida superou o desempenho da versão pura apenas para um número considerável de nós.

Em, é demonstrado que o ganho de desempenho de uma aplicação híbrida é claramente dependente de aplicação.

Nesse trabalho foram realizadas comparações para 6 diferentes benchmarks e em quatro deles o desempenho da versão original pura MPI foi melhor enquanto em outros dois a versão híbrida apresentou um melhor desempenho.

O desempenho de cada modelo depende da aplicação, do tamanho dos dados e das características dos diferentes componentes da arquitetura (CPU, sistema de memória, sistema de interconexão).

Esse trabalho desenvolveu e avaliou um modelo híbrido de programação paralela para uma aplicação de engenharia baseada no método dos elementos finitos.

Para isso foi utilizada uma abordagem incremental para transformar a aplicação MPI de simulação desenvolvida em em uma aplicação híbrida.

Essa aplicação é capaz de utilizar troca de mensagens entre os nós de um cluster SMP e memória compartilhada dentro do nó SMP.

A hipótese inicial era que o tempo consumido com comunicação na aplicação seria reduzido, pois o uso desse modelo substitui a comunicação por troca de mensagens dentro do nó SMP por uma comunicação de memória compartilhada, mais rápida.

Como consequência poderia se supor que o tempo total de execução também seria reduzido.

Essa hipótese se confirmou em parte já que o tempo de comunicação de fato foi reduzido.

Entretanto, na maioria dos casos, o tempo total de execução aumentou em relação à aplicação pura por troca de mensagens.

Isso acontece porque a transformação da aplicação para usar memória compartilhada introduz um overhead maior do que o ganho obtido com comunicação.

Essa perda foi observada em outros trabalhos na literatura.

Embora outros trabalhos como tenham obtido um ganho de performance, o que se constata é que o ganho de desempenho de uma aplicação híbrida depende de características da aplicação.

Portanto, a escolha entre uma abordagem pura MPI e uma híbrida, para uma aplicação voltada para clusters de SMPs não é trivial.

Apesar disso, a redução no tempo de comunicação obtida, indica que esta pode ser uma estratégia útil para aplicações com requisitos de comunicação mais restritos ou com sérios problemas de escalabilidade.

A aplicação desenvolvida, de fato, obteve um bom desempenho geral, com um speedup satisfatório, apesar de na maioria dos casos ter sido menos eficiente no tempo total de execução do que a versão pura MPI.

O padrão OpenMP se apresentou como uma estratégia de programação paralela extremamente simples e elegante.

Por isso sugere-se que, em trabalhos futuros, esse padrão seja utilizado para resolução de problemas semelhantes de simulação em máquinas SMP como uma maior quantidade de nós.

Outros trabalhos futuros possíveis, especificamente voltados para modelos híbridos de programação, incluem a transformação dessa aplicação em uma versão de grão grosso para avaliar o desempenho dessa estratégia.

E ainda a utilização da abordagem de sobreposição de comunicação apresentada em para buscar um desempenho melhor do que aquele obtido neste trabalho.