WSPPD'09

Em Direção a um Escalonador por Roubo de Tarefas na Criação de Tarefas por Divisão & Conquista com MPI-2

Stéfano **Mór** Nicolas **Maillard**

sdkmor@inf.ufrgs.br
nicolas@inf.ufrgs.br

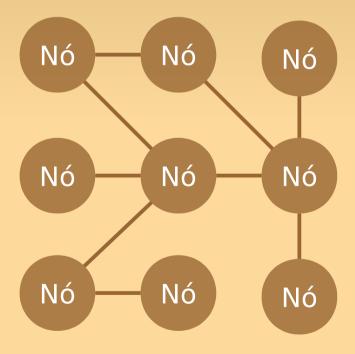


Sumário

- (1) Motivação
- (2) Roubo de Tarefas Aleatorizado
- (3) Adaptando para Memória Distribuída
- (4) Conclusões

cluster

- memória distribuída
- troca de mensagens
- barramento rápido
- MPI
- visão lógica "all-to-all"
- conectividade física varia



desafio: aproveitar ao máximo o hardware paralelo

foco atual: melhorar o escalonamento de processos MPI

escalonamento MPI

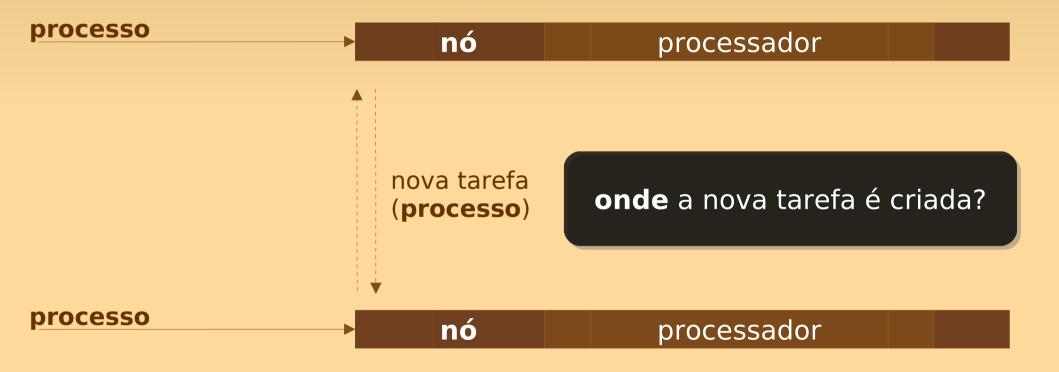
MPI 1 :: estático :: tarefa=processo

processo	nó	processador
processo	nó	processador

processo	nó	processador
processo	nó	processador

escalonamento MPI

MPI 2 :: estático & dinâmico :: processo = tarefa :: "**Task Spawning**" Estado da Arte



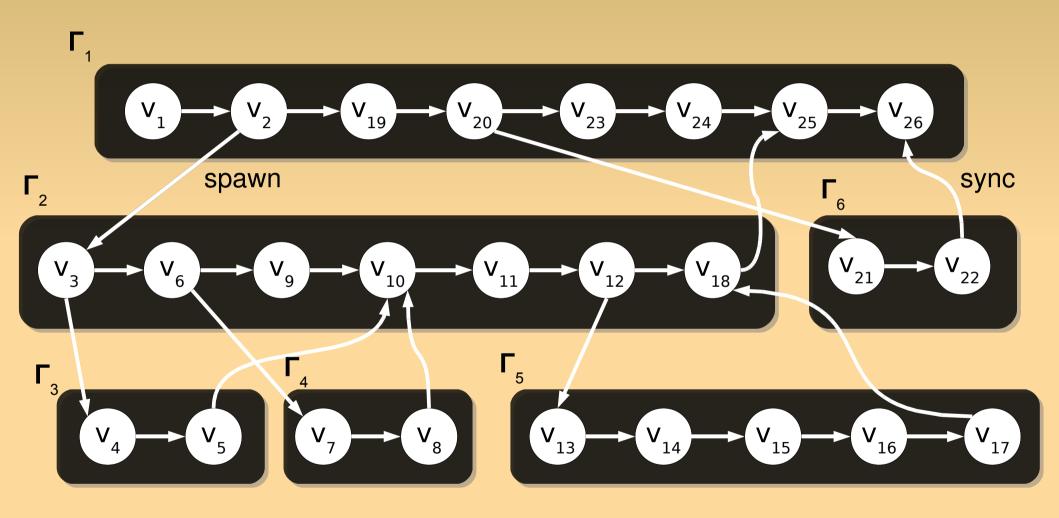
processor-aware: programador decide em qual processador alocar nova tarefa; número de processadores passados como parâmetro.

<u>processor-oblivious</u>: tarefa é alocada pelo middleware; não importa ao programador o número de processadores.

Contribuição

idéia central: usar o escalonamento por roubo de tarefas (projetado para memória compartilhada) em um cenário de memória distribuída com MPI

limitação: escalonamento eficiente para computação totalmente estrita (próximo slide) e.g., Divisão & Conquista (nosso foco)

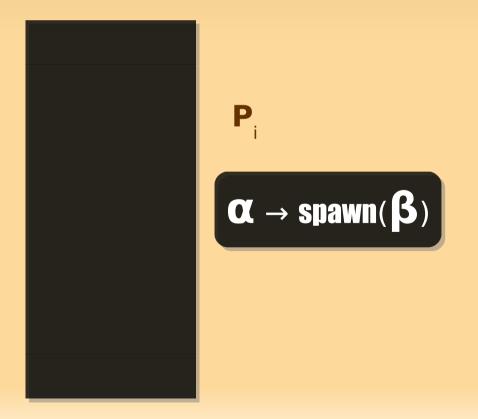


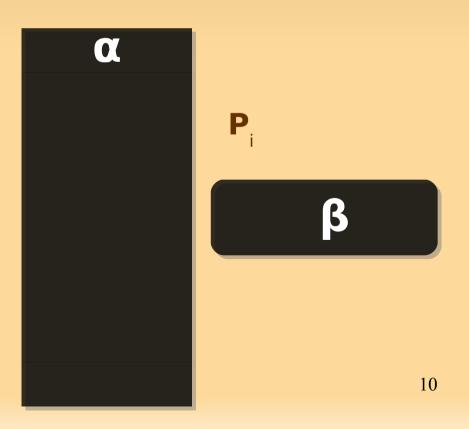
4 situações

uma deque por processador :: em cada processador

thread α faz spawn da thread β

 α é empilhada e β ganha o processador



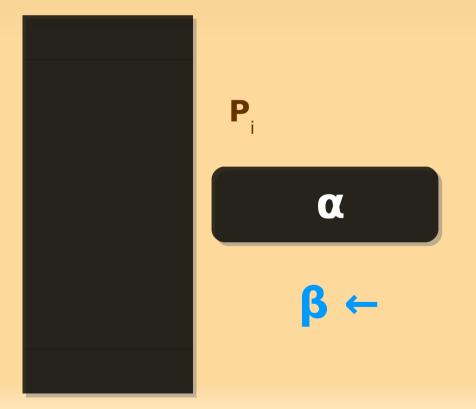


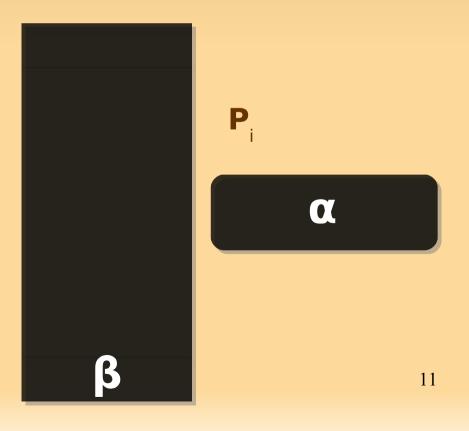
4 situações

uma deque por processador :: em cada processador

thread α está executando e thread β é habilitada

β vai para o fundo da deque



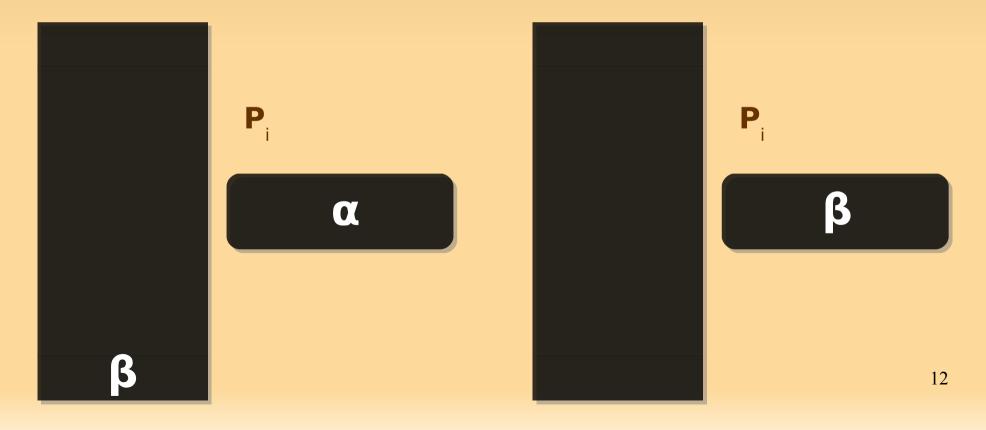


4 situações

uma deque por processador :: em cada processador

thread **a** está executando e morre.

se há thread β no fundo da deque, ela passa a executar.



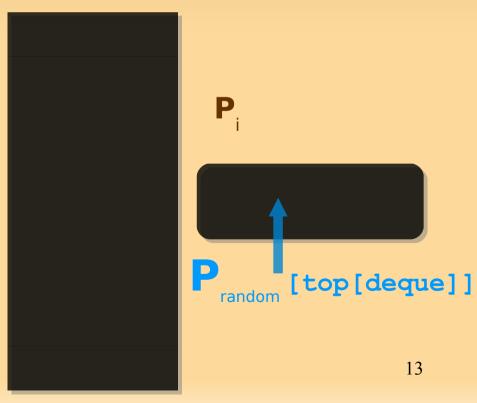
4 situações

uma deque por processador :: em cada processador

thread **a** está executando e morre.

se não há thread no fundo da deque, ela rouba tarefas do topo de outra deque.





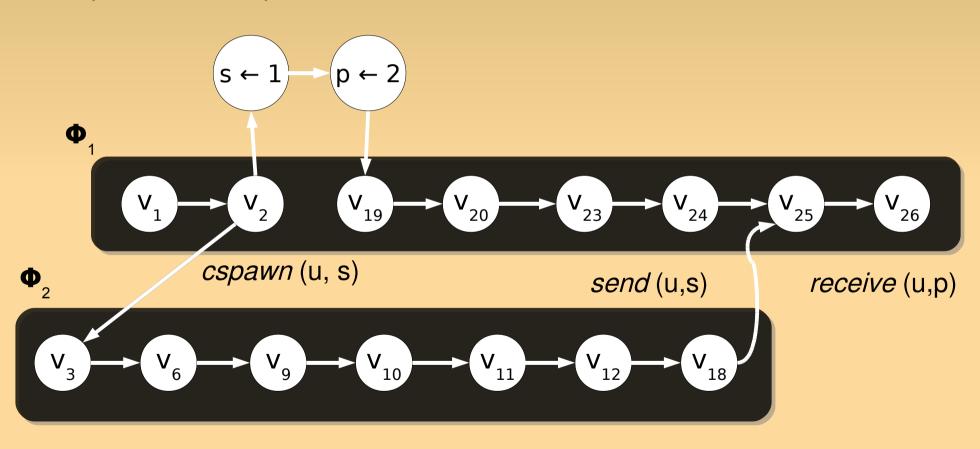
roubo de tarefas

problemas com memória distribuída

- 4 problemas cruciais
 - problema do spawn (+/-)
 - problema da sincronização (✔)
 - problema de acesso concorrente à memória(✔)
 - problema de checkpointing(✔)

roubo de tarefas

problema do spawn



roubo de tarefas

problema da sincronização

- roubo de tarefas deve ser assíncrono em relação à vítima
 - nova thread que atende roubos
 - memória compartilhada com a thread que faz o processamento em si
 - não precisa parar de processar
 - apenas em máquinas SMT

roubo de tarefas

problema da acesso concorrente à memória

- thread de processamento e thread de atendimento devem ter acesso mutuamente exclusivo sobre a deque
 - o mesmo vale para tentativas de roubo sobre a mesma deque
 - implementação do protocolo THE
 - criação de uma nova thread de atendimento apra cada roubo
 - somente funciona com máquina SMT arbitrária

roubo de tarefas

problema do checkpoiting

```
1: V<sub>1</sub>
2: V<sub>2</sub>
3: a ← cspawn t₁
4: b ← cspawn t<sub>2</sub>
5: c ← cspawn t<sub>3</sub>
6: for i ← a to c do
7: receive(u<sub>i</sub>, i)
                                                                    18
8: ProcessResults(u<sub>a</sub>, u<sub>b</sub>, u<sub>c</sub>)
```

Conclusões

2 perguntas questões em aberto

- existe equivalência entre o DAG de memória compartilhada e o DAG de memória distribuída?
 - send/receive bloqueante afeta a eficiência do RSWA?
- é possível implementar o RSWA de maneira eficiente em uma topologia do tipo anel (MPD)?

EOF

Agradecimentos ao CNPq

Obrigado!

Dúvidas?

SEMAC'09

Em Direção a um Escalonador por Roubo de Tarefas na Criação de Tarefas por Divisão & Conquista com MPI-2

Stéfano **Mór** Nicolas **Maillard**

sdkmor@inf.ufrgs.br
nicolas@inf.ufrgs.br

