

CAP-387(2016) – Tópicos Especiais em Computação Aplicada: Construção de Aplicações Massivamente Paralelas

Aula 40: Tolerância a Falhas

Celso L. Mendes, Stephan Stephany

LAC / INPE

Emails: celso.mendes@inpe.br, stephan.stephany@inpe.br

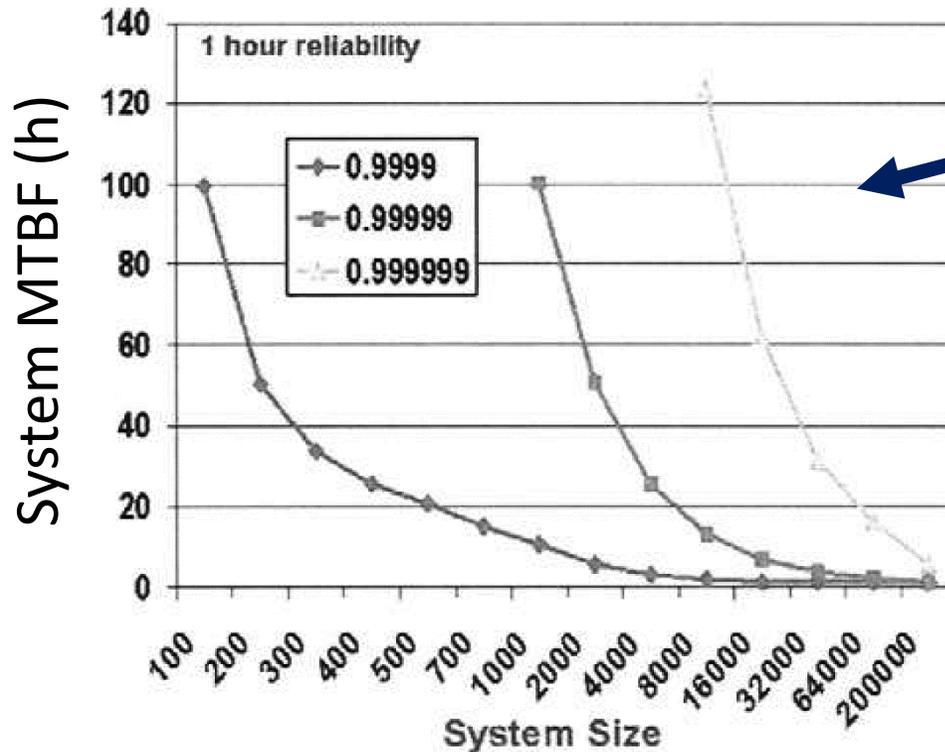


Motivação

- **Sistemas massivamente paralelos atuais**
 - Milhares/milhões de processadores
 - Embora a probabilidade de falhas em *um* processador seja baixa, a probabilidade de alguma falha num job grande pode não ser baixa
 - Em geral, basta que um processador falhe para que todo o job seja interrompido

Motivação

- MTBF: Mean-Time Between Failures



Motivação

- **Origem de falhas em sistemas reais:**
 - Falhas de hardware
 - Falhas de processadores - (figura anterior)
 - Falhas em outros componentes (discos, rede, etc)
 - Falhas de software
 - Erros de programação
 - Erros devido a configurações inapropriadas
 - Erros devido a operação inadequada
 - Falhas ambientais
 - Quedas de energia
 - Esgotamento de recursos (ex: espaço em disco)

Solução Típica

- **Checkpoint/Restart:**
 - Idéia básica: salvar (tipicamente em disco) o estado da computação periodicamente, num *checkpoint*
 - Em caso de falhas:
 - a) Parar a execução
 - b) Restaurar dados do *checkpoint* mais recente
 - c) Reiniciar a execução
 - Overheads implícitos:
 - Tempo gasto para realizar checkpoints periódicos
 - Tempo de execução perdido após uma falha

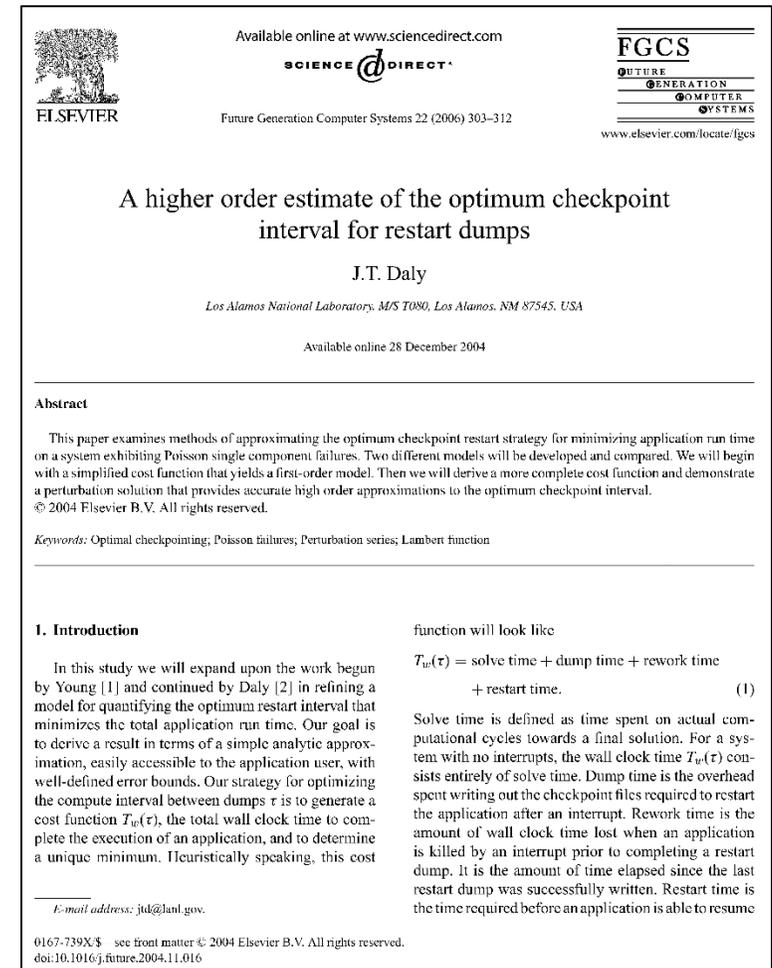


Checkpoint/Restart

- **Detalhes práticos**
 - Normalmente trabalha-se com dois conjuntos de checkpoints, para evitar perda dos dados caso ocorra uma falha durante o processo de escrita do checkpoint
 - Problema: qual a frequência ideal de checkpoint?
 - a) Frequência alta:
 - muito tempo é perdido com o overhead de realizar o checkpoint
 - b) Frequência baixa:
 - muito tempo de execução poderá ser perdido quando houver uma falha
 - Execução pode não avançar caso as falhas sejam frequentes

Frequência de Checkpoint

- **Referência básica:**
 - *Daly: A high order estimate of the optimum checkpoint interval for restart dumps, FGCS, 2004*
 - Calcula o intervalo ideal para se fazer checkpoint e minimizar o tempo de execução do programa
 - Refina resultados de trabalhos anteriores



Modelagem Analítica

- **Tempo de execução com checkpoint/restart:**

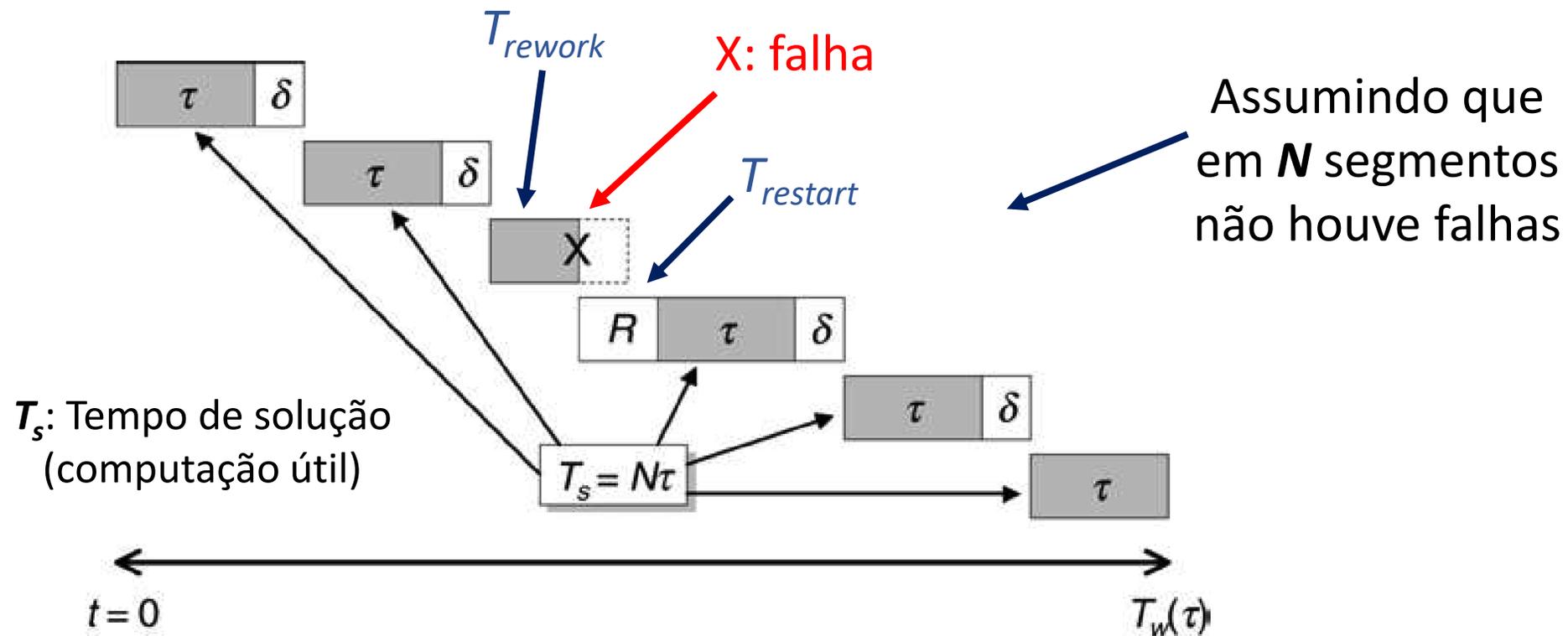
$$T_w(\tau) = T_{solve} + T_{dump} + T_{rework} + T_{restart} \quad \text{onde:}$$

- T_w : Tempo total de execução (*wall time*)
- T_{solve} : tempo para solução do problema em si
- T_{dump} : tempo gasto escrevendo arquivos de checkpoint
- T_{rework} : tempo de execução perdido antes de uma falha
- $T_{restart}$: tempo gasto para reiniciar execução após uma falha
- τ : tempo de computação entre checkpoints sucessivos

Objetivo: encontrar τ tal que $T_w(\tau)$ seja mínimo

Modelagem Analítica

- Modelo de execução com checkpoint/restart:



Modelagem Analítica

- **Resultado inicial, equivalente ao de Young:**
 - $\tau_{opt} = [2\delta(M+R)]^{1/2}$ para $M \gg \tau + \delta$ onde
M é o mean-time between failures
 - Utiliza um modelo simplificado, de primeira ordem, para o tempo de execução
- **Modelo mais sofisticado, de Daly:**
 - Refina aproximações de primeira ordem no modelo
 - Permite ocorrer mais de uma falha por segmento

Modelagem Analítica

- Resultado geral, obtido por Daly:

$$\tilde{\tau}_{opt} = \begin{cases} \sqrt{2\delta M} \left[1 + \frac{1}{3} \left(\frac{\delta}{2M} \right)^{1/2} \right. \\ \left. + \frac{1}{9} \left(\frac{\delta}{2M} \right) \right] - \delta & \text{for } \delta < 2M, \\ M & \text{for } \delta \geq 2M. \end{cases}$$

- τ_{opt} depende apenas de M e δ ; independente de R !
- Notar que τ_{opt} independente também do tempo total

Resultado Prático

- Caso δ seja pequeno em relação a M :
 - Definição de “pequeno”: $\delta < 2M$

$$\tilde{\tau}_{\text{opt}} = \begin{cases} \sqrt{2\delta M} - \delta & \text{for } \delta < \frac{1}{2}M, \\ M & \text{for } \delta \geq \frac{1}{2}M. \end{cases}$$

Exemplo de Uso

- Sistema com 25.000 nós, ~3 falhas de nó por dia:
 - $mtbf(\text{algum_nó}) \approx 8$ horas
- Job-1: $(25.000 \div 5) = 5.000$ nós
 - $mtbf(\text{job}): M = 5 \times 8 = 40$ horas
 - Assumindo $\delta = 6$ minutos = 0,1h (tem-se $\delta < M/2$):
 $\tau_{opt} (2 \times 0,1 \times 40)^{1/2} - 0,1 \approx 2,7h = 162$ minutos
- Job-2: 25.000 nós (máquina inteira)
 - $mtbf(\text{job}): M = 8$ horas
 - Assumindo $\delta = 20$ minutos = 0,33h (ainda tem-se $\delta < M/2$):
 $\tau_{opt} (2 \times 0,33 \times 8)^{1/2} - 0,33 \approx 1,97h \approx 118$ minutos