

Avaliação Preliminar do Desempenho e Custo Financeiro de Aplicações de HPC em *Clusters* de Instâncias *Burstable* da AWS

Artur Soda¹, Vanderlei Filho¹, Márcio Castro¹

¹Universidade Federal de Santa Catarina (UFSC)
Florianópolis/SC

artursodar@gmail.com, vanderlei.filho@proton.me,
marcio.castro@ufsc.br

Resumo. *Instâncias Burstable da Amazon Web Services (AWS) oferecem desempenho variável com base no consumo de créditos de CPU, o que pode representar desafios na avaliação de custo-benefício para cargas de trabalho intensivas. Este trabalho analisa o desempenho e o custo financeiro de aplicações HPC em clusters com instâncias non-Burstable e Burstable. Os resultados mostram que a escolha do tipo de instância ideal depende do perfil da carga de trabalho, com non-Burstable sendo mais eficiente para cargas intensivas e Burstable Standard adequadas para cargas leves.*

1. Introdução

A Computação em Nuvem tem desempenhado um papel essencial na evolução da Computação de Alto Desempenho (HPC), oferecendo acesso a recursos escaláveis e econômicos para diversas aplicações. Entre as ofertas disponíveis, as instâncias *Burstable* da Amazon Web Services (AWS) têm ganhado destaque devido à sua capacidade de fornecer desempenho variável com base no consumo de créditos de CPU. Embora seja economicamente atraente para cargas leves, o desempenho dessas instâncias em cenários de HPC distribuído ainda é pouco explorado.

Este trabalho dá continuidade à pesquisa realizada por Ferrari *et al.* [Ferrari et al. 2024], que analisou instâncias *Burstable* em tarefas de HPC de único nó. Outros estudos, como o de Leitner e Scheuner [Leitner and Scheuner 2015], também avaliaram instâncias da família T2 (*Burstable*) a outros tipos de instâncias de propósito geral, concluindo uma melhor relação custo-desempenho apenas quando a utilização média da CPU é inferior a 40%. Neste estudo, o escopo foi ampliado para cargas mais pesadas em *clusters* multinós. Foram avaliadas instâncias *Burstable* nos modos *Unlimited* e *Standard*, além de uma instância *non-Burstable* similar. Os resultados sugerem que cargas intensivas são mais eficientes em instâncias *non-Burstable*, enquanto *Burstable Standard* favorecem cargas leves. Já as *Burstable Unlimited* apresentaram custos elevados em cargas de trabalho longas e intensivas devido ao consumo de créditos excedentes de CPU.

O artigo está organizado da seguinte forma. Na Seção 2, é apresentada a fundamentação teórica sobre HPC na nuvem e as diferenças entre instâncias *Burstable* e *non-Burstable* da AWS. Na Seção 3, são detalhados o ambiente experimental e as aplicações utilizadas. A Seção 4 discute os resultados experimentais, enquanto a Seção 5 conclui o artigo, resumindo os principais achados e indicando possíveis trabalhos futuros.

2. Fundamentação Teórica

2.1. HPC na nuvem

A área de HPC busca resolver problemas complexos utilizando múltiplos nós de processamento. Com o avanço da Computação em Nuvem, provedores como a AWS passaram

a oferecer recursos que reduzem a dependência de infraestruturas locais de alto custo. Essa abordagem permite maior flexibilidade, escalabilidade e acesso a recursos de alto desempenho sob demanda, tornando a nuvem uma alternativa atraente e acessível para a execução de cargas de trabalho de HPC em diferentes contextos e aplicações.

2.2. Instâncias da AWS: *Burstable* vs. *non-Burstable*

Instâncias *Burstable* da AWS (família T, como T3 e T4) baseiam-se em uma *baseline* de utilização de CPU, que define uma porcentagem fixa de uso contínuo permitido para cada vCPU da instância. Quando o uso de uma vCPU está abaixo dessa *baseline*, a instância acumula créditos de CPU, que podem ser utilizados para atender a picos de demanda, onde o uso da vCPU excede essa *baseline*. Cada crédito de CPU representa um minuto de uso de uma vCPU a 100% de capacidade.

As instâncias *Burstable* possuem dois modos de operação: no modo *Standard*, o desempenho da instância é reduzido quando os créditos acumulados se esgotam. Já no modo *Unlimited*, a instância pode continuar funcionando acima da *baseline*, mas custos adicionais são gerados pelo uso de créditos excedentes, o que pode impactar o custo total da instância, dependendo da utilização. Por outro lado, as instâncias *non-Burstable* (como as famílias C, M e R) oferecem desempenho constante, sem a limitação ou a necessidade de acumular créditos.

3. Método de Análise

3.1. Ambiente Experimental

Os experimentos foram realizados em diferentes *clusters* homogêneos de 1, 2 e 4 nós criados com a ferramenta HPC@Cloud [Munhoz and Castro 2023] utilizando as instâncias da AWS com imagens Linux mostradas na Tabela 1.

Tipo da Instância	Quantidade de vCPUs	Tipo de Instância	Largura de Banda da Rede	Custo por Hora	Custo por Crédito de CPU Excedido
t3.2xlarge	8	<i>Burstable</i>	Até 5 Gbps	0.3328 USD	0.05 USD
m5.2xlarge	8	<i>non-Burstable</i>	Até 10 Gbps	0.384 USD	–

Tabela 1. Características das instâncias utilizadas (us-east-1).

3.2. Aplicações de HPC

Para avaliar o desempenho das instâncias em cargas de trabalho de HPC, foi utilizado o *NAS Parallel Benchmarks* (NPB) [Bailey et al. 1991], um conjunto de *benchmarks* amplamente reconhecido para medir a eficiência de *clusters* e outras infraestruturas de HPC. Dentro do NPB, o foco foi direcionado aos *benchmarks* LU (*Lower-Upper Symmetric Gauss-Seidel*) e EP (*Embarrassingly Parallel*), que representam diferentes tipos de carga de trabalho e abordagens computacionais.

O LU é um *benchmark CPU-bound*, ou seja, sua execução depende principalmente da capacidade de processamento da CPU, envolvendo a resolução de sistemas de equações lineares simétricas. Esse *benchmark* é intensivo em cálculos e exige interação constante entre os nós, sendo adequado para avaliar o desempenho em tarefas com alta dependência de comunicação.

Em contraste, o EP, sendo um *benchmark Embarrassingly Parallel*, possui alto grau de paralelismo, permitindo a execução independente dos processos com pouca ou nenhuma comunicação entre os nós. Esse tipo de carga de trabalho foca no desempenho individual de cada núcleo de CPU.

A combinação dos *benchmarks* LU e EP permite uma avaliação abrangente do desempenho das instâncias da AWS, uma vez que eles englobam tanto cenários com alta dependência de comunicação entre os processos quanto aqueles que exploram o paralelismo de forma mais independente.

4. Resultados Experimentais

Os resultados, apresentados na Figura 1, destacam diferenças significativas no desempenho das configurações para os *benchmarks* EP e LU. No caso do EP, caracterizado pela rápida execução e alta independência entre processos, todas as configurações demonstraram desempenho semelhante em termos de tempo de execução e escalabilidade. Por outro lado, no LU, as instâncias M5 obtiveram o melhor desempenho, possivelmente devido à maior largura de banda disponível. Em contrapartida, as T3 no modo *Standard* apresentaram um desempenho até três vezes inferior, uma limitação atribuída à redução de desempenho imposta quando os créditos de CPU são esgotados.

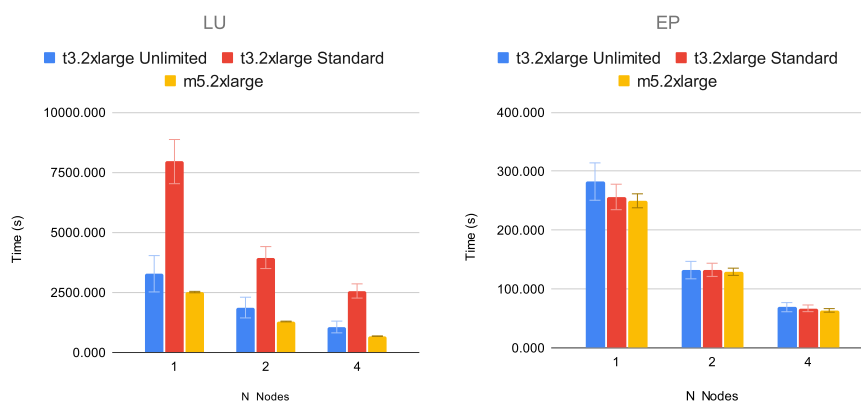


Figura 1. Tempo de execução (segundos) por *benchmark*, tipo de instância e quantidade de nós.

As diferenças nos custos, ilustradas na Figura 2, complementam essa análise. Para o EP, onde os tempos de execução foram semelhantes entre as configurações, as instâncias *Burstable* apresentaram uma ligeira vantagem de custo, graças à tarifa horária reduzida. Já no LU, as M5 mostraram-se superiores em termos de custo-benefício, equilibrando desempenho e custo. No entanto, as T3 no modo *Unlimited* mesmo que não tenham sido as de pior desempenho, apresentaram custos significativamente maiores, devido ao elevado custo associado ao consumo de créditos de CPU excedentes durante cargas de trabalho prolongadas.

De forma geral, as M5 se mostraram mais adequadas para cargas de trabalho intensivas e longas, enquanto as T3 *Standard* são mais interessantes para cargas leves, podendo ter picos ocasionais de utilização da CPU. Já as instâncias T3 no modo *Unlimited* são menos atrativas devido aos custos elevados, mas podem ser úteis em cenários específicos onde a aplicação não pode ser limitada durante picos de demanda, garantindo assim um desempenho consistente, mesmo com um custo adicional.

Todas as configurações demonstraram escalabilidade consistente com o aumento do número de nós, reduzindo o tempo de execução, mas elevando proporcionalmente os custos. Esses resultados demonstram a importância de alinhar as características da carga de trabalho às especificidades das instâncias para maximizar a eficiência em termos de desempenho e custo.

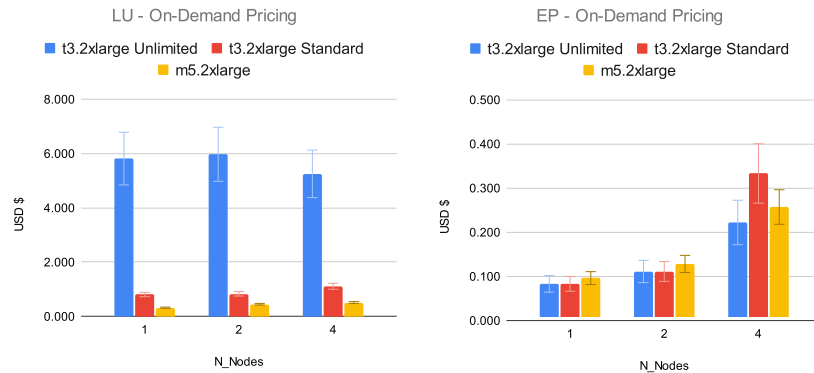


Figura 2. Custo total (USD \$), *On-Demand*, por *benchmark*, tipo de instância e quantidade de nós.

5. Conclusão

Neste trabalho, foi analisado o desempenho e o custo de instâncias *Burstable* e *non-Burstable* da AWS para cargas de trabalho de HPC, utilizando os *benchmarks* LU e EP do NPB. Através dos experimentos, foi possível observar como diferentes tipos de instâncias e configurações impactam o desempenho e os custos em cenários distintos. Os resultados demonstraram que as instâncias M5, devido à sua capacidade de largura de banda e ausência de limitações por créditos, oferecem melhor custo-benefício para cargas de trabalho intensivas e longas. Por outro lado, as instâncias T3 *Standard* se destacaram para cargas de trabalho leves e intermitentes, enquanto as T3 *Unlimited*, apesar de garantirem desempenho contínuo, apresentaram custos significativamente maiores em cargas de trabalho prolongadas.

Como trabalhos futuros, sugere-se ampliar a análise incluindo mais aplicações do NPB, visando uma compreensão mais abrangente das características de diferentes cargas de trabalho. Além disso, planeja-se explorar o SPEChpc, outro conjunto de *benchmarks* amplamente reconhecido. Por fim, pretendemos avaliar instâncias e *clusters* de maior capacidade, possibilitando testes com problemas mais complexos e realistas, expandindo o escopo e a relevância dos resultados apresentados.

Referências

- Bailey, D. H., Barszcz, E., Barton, J. T., Browning, D. S., Carter, R. L., Dagum, L., Fatoohi, R. A., Frederickson, P. O., Lasinski, T. A., Schreiber, R. S., Simon, H. D., Venkatakrisnan, V., and Weeratunga, S. K. (1991). The nas parallel benchmarks. *The International Journal of High Performance Computing Applications*, 5(3):63–73.
- Ferrari, G., Filho, V., and Castro, M. (2024). Comparing burstable and on-demand aws ec2 instances using nas parallel benchmarks. In *Anais da XXIV Escola Regional de Alto Desempenho da Região Sul*, pages 61–64, Porto Alegre, RS, Brasil. SBC.
- Leitner, P. and Scheuner, J. (2015). Bursting with possibilities – an empirical study of credit-based bursting cloud instance types. In *2015 IEEE/ACM 8th International Conference on Utility and Cloud Computing (UCC)*, pages 227–236.
- Munhoz, V. and Castro, M. (2023). Enabling the Execution of HPC Applications on Public Clouds with HPC@Cloud Toolkit. *Concurrency and Computation: Practice and Experience*, pages 1–19.