

Análise de um Log Real de Jobs visando a Previsão de Tempos de Espera para Execução em um Cluster de Alto Desempenho

Bernardo Gallo¹, Matheus Marotti¹, Lúcia Maria de Assumpção Drummond¹,
José Viterbo¹, Felipe A. Portella², Paulo J. B. Estrela², Renzo Q. Malini²

¹ Instituto de Computação – Universidade Federal Fluminense (UFF)
Niterói – RJ – Brasil

²Petróleo Brasileiro S.A. – (PETROBRAS)
Rio de Janeiro – RJ – Brasil

{bgallo, matheusmarotti}@id.uff.br, {lucia, viterbo}@ic.uff.br
{felipeportella, paulo.estrela, renzo}@petrobras.com.br

Abstract. *Deepwater oil exploration and production depends on complex numerical simulations, executed on supercomputers, to optimize production and project future scenarios. At Petrobras, a world leader in deepwater oil exploration, this infrastructure includes several systems ranked in the TOP500. However, despite the robust computational capacity, the nature of these operations generates peak demand moments. To create better strategies for handling these peaks, this work analyzes queue wait time behavior, concluding that user-specific workload metrics—such as the sum of execution time and CPUs for jobs already waiting—are the factors with the highest correlation.*

Resumo. *A exploração de petróleo em águas profundas depende de simulações numéricas complexas, executadas em supercomputadores, para otimizar a produção e projetar cenários futuros. Na Petrobras, líder mundial em exploração de petróleo em águas profundas, essa infraestrutura inclui diversos sistemas classificados no TOP500. Contudo, apesar da robusta capacidade computacional, a natureza das operações gera momentos de pico de demanda. Para criar melhores estratégias para lidar com esses picos, este trabalho analisa o comportamento do tempo de espera na fila, concluindo que as métricas de carga de trabalho por usuário — como a soma do tempo de execução e de CPUs de jobs já em espera — são os fatores com a maior correlação.*

1. Introdução

A Petróleo Brasileiro S/A (Petrobras), uma das empresas líderes globais em exploração de petróleo em águas profundas, utiliza extensivamente simulações de reservatório para gerenciar seus poços e projetar a produção futura. Essas simulações dependem da Computação de Alto Desempenho, do inglês High-Performance Computing (HPC), com a empresa mantendo uma infraestrutura robusta que inclui supercomputadores, alguns figurando entre os mais potentes do mundo [TOP500.org 2025]. Apesar dessa capacidade, a Petrobras enfrenta desafios em períodos de alta demanda. O acesso aos recursos de HPC é orquestrado por gerenciadores e, quando a demanda é maior do que a disponibilidade de recursos, os jobs aguardam sua alocação em uma fila de espera, impactando

negativamente a produtividade e a eficiência das equipes. Assim, entender e mitigar o comportamento dessa fila de espera torna-se crucial.

Para resolver esse problema, abordagens puramente analíticas são insuficientes. A análise estatística e a predição acurada do tempo de espera emergem como ferramentas valiosas. Prever com precisão o tempo de espera não apenas otimiza o uso dos recursos e a experiência do usuário, mas também permite identificar e diagnosticar as causas do congestionamento do sistema. A literatura sugere que a predição do tempo de espera deve considerar uma combinação de atributos do *job* e características da fila. Fatores com alto poder preditivo incluem a carga instantânea do sistema [Menear et al. 2024], o tempo total de execução (*wall time*) [Lovell et al. 2024], o histórico do tempo de espera [Ramachandran et al. 2024], os limites de execução por usuário [Paokin and Nikitenko 2023] e a sazonalidade da demanda [Brown et al. 2024].

Inspirado por essas abordagens, o presente trabalho se dedica a aprofundar a compreensão sobre o comportamento da fila de espera em um dos clusters de HPC da Petrobras. A análise, realizada a partir dos logs históricos de *jobs* do escalonador Slurm [Yoo et al. 2003], irá investigar os principais atributos dos *jobs* submetidos e suas correlações com a dinâmica do sistema, buscando identificar os fatores que mais influenciam o tempo de espera. O objetivo não é o desenvolvimento imediato de um modelo preditivo, mas sim gerar o conhecimento que, além de servir como um alicerce para futuras ferramentas de predição, é essencial para entender os picos de demanda e o comportamento geral do sistema.

2. Resultados

A análise foi conduzida sobre um conjunto de dados extraído dos *logs* do escalonador Slurm, referentes ao ano de 2024. O total de registros era de 2.745.759 *jobs*. Por sugestão dos *stakeholders* da Petrobras, para focar em execuções relevantes que enfrentaram espera significativa e remover *outliers* aplicou-se os seguintes filtros: *jobs* com tempo de execução superior a 1 hora (41.18%), *jobs* com tempo de espera maior que 5 minutos (5.42%), *jobs* com tempo de espera menor que um dia (99.99%). O conjunto de dados resultante para a análise contém 88.787 *jobs* (3.23%).

2.1. Distribuição Geral do Tempo de Espera

A análise da distribuição geral revela que a maioria dos *jobs* aguarda um tempo relativamente curto para iniciar. O tempo de espera mediano é de 2.0 horas, com uma média de 3.3 horas. O desvio padrão de 3.7 horas, próximo à média, indica uma variação considerável nos tempos de espera.

Conforme observado no histograma da Figura 1, há uma forte concentração de *jobs* na faixa inicial de 0 a 30 minutos, com um decaimento subsequente. Apesar disso, a fila possui uma cauda longa, com uma fração significativa de *jobs* (entre 5% e 10%) enfrentando esperas superiores a 8 horas. A distribuição foi categorizada em três grupos principais: 31.3% dos *jobs* esperam até 1 hora, 30.0% esperam entre 1 e 3 horas, e 38.7% aguardam mais de 3 horas, mostrando um equilíbrio entre esperas curtas, médias e longas.

2.2. Correlação entre Variáveis e o Tempo de Espera

Para identificar as variáveis de maior impacto no tempo de espera, foi realizada uma análise de correlação a partir de três métodos distintos: o de Pearson [Pearson 1895],

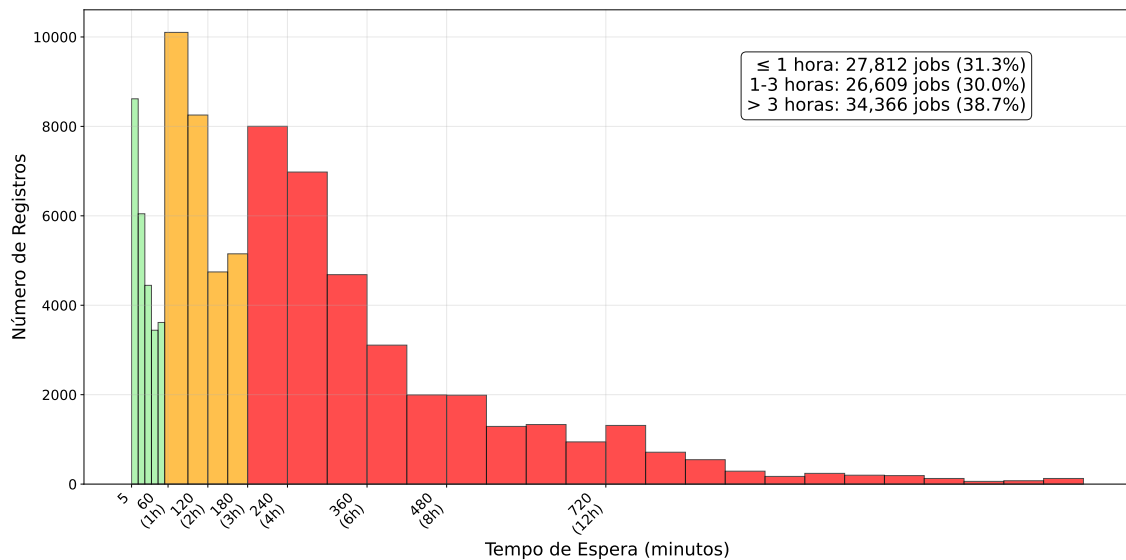


Figura 1. Histograma do tempo de espera. A distribuição é dividida por categorias: até 1h – Verde (31.3%), 1-3h – Laranja (30.0%) e mais de 3h – Vermelho (38.7%), destacando uma concentração de jobs com esperas curtas, mas com uma cauda longa significativa.

para avaliar a força de relações lineares, e os de Spearman [Spearman 1904] e Kendall [Kendall 1938], para capturar tendências monotônicas, não necessariamente lineares.

A literatura aponta que os dados de execução por usuário são um fator preditivo relevante para o tempo de espera. Com base nessa premissa e na existência de limites por usuário no *cluster*, a análise foi direcionada para a carga de trabalho de cada contexto de submissão. Para tanto, foram calculadas variáveis granulares que refletem o estado da fila sob a perspectiva do usuário. Os resultados, consolidados na Tabela 1, mostram que essa abordagem foi eficaz, uma vez que as cinco variáveis com maior correlação média são todas métricas específicas de usuário. Dentre elas, a “Soma do tempo de execução na fila (MU)” — que representa o tempo de execução total pelos jobs do mesmo usuário já em espera — demonstrou ser a variável com a correlação mais forte. A seguir, detalha-se o comportamento das duas principais métricas identificadas.

Tabela 1. Top 5 Variáveis por Correlação Média com o Tempo de Espera. Filtro aplicado aos jobs na fila: (MU) para jobs do Mesmo Usuário.

Variável	Pearson	Spearman	Kendall
Soma do tempo de execução na fila (MU)	0.591	0.567	0.401
Soma de CPUs requisitadas na fila (MU)	0.440	0.543	0.387
Soma de memória requisitada na fila (MU)	0.440	0.543	0.387
Soma do tempo de execução na fila	0.521	0.482	0.338
Número de jobs na fila (MU)	0.346	0.500	0.352

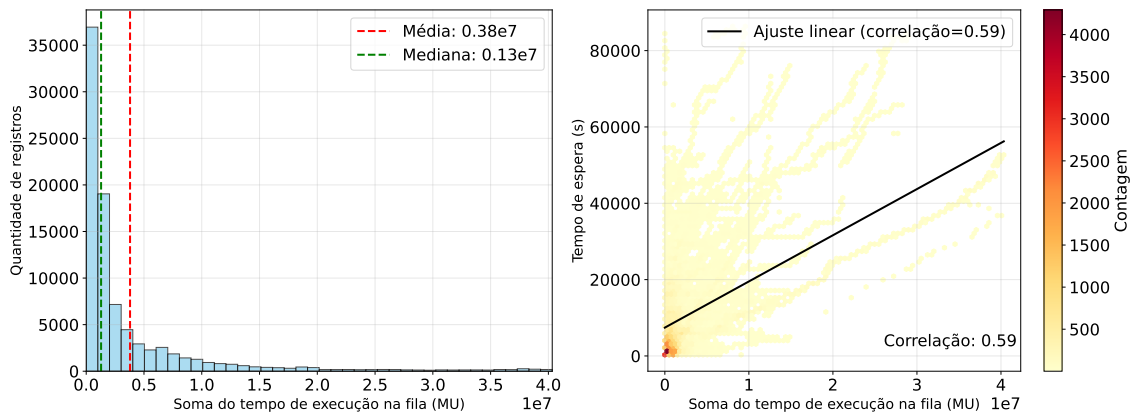


Figura 2. Histograma (esquerda) e scatterplot (direita) da “Soma do tempo de execução na fila (MU)”. O scatterplot ilustra a correlação positiva (0.59) entre esta métrica (Eixo X) e o Tempo de Espera (Eixo Y). A escala de cor “Contagem” indica a densidade de jobs.

2.2.1. Análise da Soma do Tempo de Execução na fila (Mesmo Usuário)

A “Soma do tempo de execução na fila (MU)” destacou-se como a variável de maior impacto preditivo. Ela representa a carga de trabalho total que um usuário já possui na fila, somando o tempo de execução de todos os seus *jobs* em espera. Os gráficos da Figura 2 mostram uma tendência geral do *dataset*: a maioria dos *jobs* tem um tempo de espera curto, enquanto poucos esperam por muito tempo. Adicionalmente, os gráficos confirmam a forte correlação positiva da variável com o tempo de espera, validada pelos coeficientes de Pearson (0.591) e Spearman (0.567). Fica evidente que a carga de trabalho acumulada de um usuário na fila é o preditor individual mais poderoso, dentre os dados disponíveis, para a espera de seus novos *jobs*.

É fundamental ressaltar, contudo, que essa métrica foi calculada usando o tempo de execução real, um dado indisponível no momento da submissão de um novo *job*. Isso a caracteriza como uma variável explicativa poderosa em análise, mas não como um preditor utilizável em tempo real. A força dessa correlação, no entanto, aponta um caminho promissor: o desenvolvimento de um bom preditor para o tempo de execução é um passo intermediário de grande valor para viabilizar uma futura previsão do tempo de espera.

2.3. Análise da Soma de CPUs requisitadas na fila (Mesmo Usuário)

De forma análoga à métrica anterior, a “Soma de CPUs requisitadas na fila (MU)” também apresenta uma forte correlação positiva com o tempo de espera, embora com uma magnitude ligeiramente inferior (Pearson=0.440, Spearman=0.543). A distribuição desta variável é igualmente assimétrica, com uma média de 7954.68 muito superior à mediana de 3080.00, refletindo a concentração de alta requisição de recursos em poucos usuários. A relevância preditiva de ambas as métricas de usuário — soma de tempo de execução e de CPUs — é explicada pela existência de um limite de uso de CPUs simultâneas por usuário no *cluster*. Assim, essas variáveis quantificam diretamente o quão próximo um usuário está de atingir sua cota de alocação, o que impacta o escalonamento e o tempo de espera de seus novos *jobs*.

3. Conclusão e Trabalhos Futuros

Pode-se perceber, portanto, que tanto informações do estado atual do *cluster* quanto o tempo de execução de *job* (TEJ) são informações importantes para obter uma estimativa confiável do tempo de espera de um *job*. Contudo, o TEJ não está disponível a tempo de submissão de um *job*; portanto, é necessário usar uma arquitetura de previsão sequencial, com um preditor para o TEJ e, em seguida, usar esse resultado para a previsão do tempo de espera.

Em trabalhos futuros pretende-se implementar um preditor de tempo de espera na arquitetura sequencial, comparando diferentes modelos de aprendizado de máquina em diferentes métricas.

Referências

- [Brown et al. 2024] Brown, N., Gibb, G., Belikov, E., and Nash, R. (2024). Predicting accurate batch queue wait times on production supercomputers by combining machine learning techniques. *Concurrency and Computation: Practice and Experience*, 36(15):e8112.
- [Kendall 1938] Kendall, M. G. (1938). A new measure of rank correlation. *Biometrika*, 30(1/2):81–93.
- [Lovell et al. 2024] Lovell, A., Wisniewski, P., Rodenbeck, S., and Ashish (2024). A hierarchical deep learning approach for predicting job queue times in hpc systems. In *SC24-W: Workshops of the International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 621–628.
- [Menear et al. 2024] Menear, K., Konate, K., Potter, K., and Duplyakin, D. (2024). Tandem predictions for hpc jobs. In *Practice and Experience in Advanced Research Computing 2024: Human Powered Computing*, PEARC '24, New York, NY, USA. Association for Computing Machinery.
- [Paokin and Nikitenko 2023] Paokin, A. V. and Nikitenko, D. A. (2023). Approbation of methods for supercomputer job queue wait time estimation. *Lobachevskii Journal of Mathematics*, 44(8):3140–3147.
- [Pearson 1895] Pearson, K. (1895). Note on regression and inheritance in the case of two parents. *Proceedings of the Royal Society of London*, 58:240–242.
- [Ramachandran et al. 2024] Ramachandran, S., Jayalal, M., Vasudevan, M., and Jehadeesan, R. (2024). Combining machine learning & metaheuristic algorithms for predicting waiting time of high performance computing jobs. In *2024 5th International Conference on Innovative Trends in Information Technology (ICITIIT)*, pages 1–6.
- [Spearman 1904] Spearman, C. (1904). The proof and measurement of association between two things. *The American Journal of Psychology*, 15(1):72–101.
- [TOP500.org 2025] TOP500.org (2025). Top500 list. <https://top500.org/lists/top500/2025/06/>. Acessado em: 10 de setembro de 2025.
- [Yoo et al. 2003] Yoo, A. B., Jette, M. A., and Grondona, M. (2003). Slurm: Simple linux utility for resource management. In Feitelson, D., Rudolph, L., and Schwiegelshohn, U., editors, *Job Scheduling Strategies for Parallel Processing*, pages 44–60, Berlin, Heidelberg. Springer Berlin Heidelberg.