



UNITINS
UNIVERSIDADE ESTADUAL DO TOCANTINS

TOCANTINS
GOVERNO DO ESTADO



UNIVERSIDADE ESTADUAL DO TOCANTINS

SISTEMAS DE INFORMAÇÃO

**ALUNO: JOÃO VICTOR PÓVOA FRANÇA e ANA BEATRIZ MARQUES
MOREIRA**

**PESQUISA CIENTÍFICA: SISTEMAS DISTRIBUÍDOS EM INTELIGÊNCIA
ARTIFICIAL (AI) E MACHINE LEARNING (ML)**

PALMAS-TO 2025

RESUMO

Este trabalho de pesquisa explora a interseção crítica entre Sistemas Distribuídos e os campos emergentes da Inteligência Artificial (AI) e Machine Learning (ML). Com o crescimento exponencial de dados e a complexidade dos modelos de AI/ML, a capacidade de processamento e armazenamento de sistemas centralizados tornou-se insuficiente. Sistemas distribuídos oferecem soluções robustas para escalabilidade, tolerância a falhas e eficiência, tornando-se fundamentais para o desenvolvimento e implantação de aplicações de AI/ML em larga escala. Serão abordados os conceitos fundamentais de sistemas distribuídos, suas metas e características, bem como a aplicação de técnicas de AI para otimizar o desempenho e a resiliência desses sistemas. Além disso, serão detalhadas as arquiteturas e frameworks específicos para o ML distribuído, incluindo paralelismo de dados e de modelos, e os desafios inerentes a essa integração. O objetivo é fornecer uma compreensão abrangente de como a sinergia entre sistemas distribuídos e AI/ML impulsiona a inovação tecnológica e supera as limitações computacionais.

Palavras Chave:



UNITINS
UNIVERSIDADE ESTADUAL DO TOCANTINS

TOCANTINS
GOVERNO DO ESTADO



SUMÁRIO

1. Introdução

2. Fundamentação Teórica: Sistemas Distribuídos

- 2.1. Definição e Características
- 2.2. Metas e Vantagens
- 2.3. Tipos de Sistemas Distribuídos
- 2.4. Classificações Arquitetônicas

3. Sistemas Distribuídos em AI e ML

- 3.1. A Necessidade de Sistemas Distribuídos em AI/ML
- 3.2. Aplicações de AI em Sistemas Distribuídos
- 3.3. Técnicas de AI Utilizadas em Sistemas Distribuídos (Machine Learning, Deep Learning, Detecção de Anomalias, Reinforcement Learning, NLP)

4. Arquiteturas e Frameworks para ML Distribuído

- 4.1. Processamento Paralelo para ML Distribuído
- 4.2. Treinamento de Modelos Distribuídos
- 4.3. Algoritmos de ML Distribuído (Parameter Server, AllReduce)
- 4.4. Frameworks e Plataformas para ML Distribuído

5. Desafios e Considerações

- 5.1. Desafios Comuns
- 5.2. Otimização e Melhores Práticas

6. Conclusão

7. Referências Bibliográficas

1. INTRODUÇÃO

A era digital é marcada pelo volume massivo de dados e pela crescente demanda por soluções inteligentes baseadas em Inteligência Artificial (AI) e Machine Learning (ML). A complexidade computacional e a escala dos dados para treinamento e inferência de modelos de AI/ML frequentemente superam as capacidades de sistemas centralizados. Modelos de Deep Learning, por exemplo, podem exigir terabytes de dados e bilhões de parâmetros, tornando inviável sua execução em uma única máquina.

Nesse contexto, os Sistemas Distribuídos emergem como a arquitetura predominante para atender às exigências de alta disponibilidade, escalabilidade e resiliência. Um sistema distribuído é um conjunto de computadores independentes, interligados por rede, que operam de forma coesa, apresentando-se ao usuário como um sistema único. A sinergia entre sistemas distribuídos e AI/ML é crucial: AI/ML demanda a infraestrutura robusta dos sistemas distribuídos, enquanto a AI pode otimizar a gestão e o desempenho desses sistemas, tornando-os mais eficientes e seguros.

Este trabalho visa analisar essa integração, focando nas arquiteturas, técnicas e desafios. Serão abordados os conceitos fundamentais de sistemas distribuídos, a necessidade e as aplicações de AI/ML em ambientes distribuídos, as arquiteturas e frameworks para ML distribuído, e os desafios e melhores práticas para essa integração.

2. FUNDAMENTAÇÃO TEÓRICA: SISTEMAS DISTRIBUÍDOS

2.1. DEFINIÇÃO E CARACTERÍSTICAS

Um sistema distribuído é uma coleção de computadores autônomos interligados por rede, que funcionam como um sistema único e coeso para o usuário. Essa unificação é alcançada pela coordenação de ações e troca de mensagens entre componentes de hardware e software em máquinas distintas. Exemplos incluem serviços de streaming, redes sociais e internet banking. As características incluem múltiplos computadores conectados, dispersão geográfica, ausência de limite de dispositivos e diversidade de configurações. A complexidade é gerenciada por transparência, consistência, colaboração, comunicação e detecção/correção de falhas.

2.2. METAS E VANTAGENS

As metas primárias de um sistema distribuído são otimizar funcionalidade e usabilidade: transparência, compartilhamento de recursos, abertura e escalabilidade.

- **Transparência:** Oculta a complexidade da distribuição dos recursos e processos do usuário e desenvolvedor. As operações parecem locais e simples, mesmo com múltiplos servidores e replicação de dados. Tipos incluem: Acesso, Localização, Migração, Replicação, Concorrência e Falhas.

- **Compartilhamento de Recursos:** Permite o uso conjunto de hardware, dados e serviços, reduzindo custos e facilitando a colaboração. Desvantagens incluem concorrência, consistência de dados, segurança e dependência de rede.
- **Abertura:** Facilita a integração com outros sistemas e a extensibilidade através de padrões e interfaces bem definidas.
- **Escalabilidade:** Capacidade de crescer em usuários, dados e alcance geográfico, mantendo desempenho aceitável. Pode ser de tamanho, administrativa ou geográfica.

2.3. TIPOS DE SISTEMAS DISTRIBUÍDOS

Os sistemas distribuídos classificam-se em três tipos principais:

- **Sistemas de Computação Distribuídos:** Focados em alto desempenho, somam poder de processamento para cargas paralelas. Incluem Cluster Computing, Grid Computing e Cloud Computing.
- **Sistemas de Informação Distribuídos:** Focados em aplicações e dados, garantem consistência em operações críticas. Abrangem Processamento de Transações Distribuídas (com propriedades ACID) e Integração de Aplicações Empresariais (EAI), utilizando chamadas remotas, mensageria e microserviços.
- **Sistemas Pervasivos / IoT:** Integram dispositivos físicos (sensores, atuadores) que coletam dados e interagem com serviços na borda e na nuvem, sendo ubíquos, móveis e utilizando redes de sensores (Internet das Coisas).

2.4. CLASSIFICAÇÕES ARQUITETÔNICAS

A base dos sistemas distribuídos evoluiu de:

- **Sistemas Centralizados:** Processamento e dados em um único computador, simples, mas com ponto único de falha.
- **Sistemas Paralelos Fortemente Acoplados:** Processadores compartilham memória física, comunicação via barramento interno. Baixa latência, mas escalabilidade limitada.
- **Sistemas Paralelos Fracamente Acoplados:** Cada nó tem sua memória, comunicação por rede. Base de clusters e sistemas distribuídos modernos (cliente-servidor, P2P). Oferecem escalabilidade horizontal e flexibilidade, mas com desafios de latência e falhas parciais.

3. SISTEMAS DISTRIBUÍDOS EM IA E ML

3.1. A NECESSIDADE EM SISTEMAS DISTRIBUÍDOS

A demanda computacional de modelos modernos de IA/ML — em especial Deep Learning — estoura facilmente os limites de uma única máquina (CPU/GPU e memória). Desde os primeiros sistemas de treinamento em larga escala, como o DistBelief do Google, a

comunidade já apontava para clusters como caminho obrigatório: “We have developed a software framework called DistBelief that can utilize computing clusters with thousands of machines to train large models”.

A crescente complexidade e o volume massivo de dados em aplicações de Inteligência Artificial (AI) e Machine Learning (ML) tornaram os sistemas distribuídos uma necessidade imperativa. Modelos de AI/ML, especialmente os de Deep Learning, frequentemente exigem recursos computacionais que excedem a capacidade de uma única máquina, tanto em termos de poder de processamento quanto de memória para armazenar grandes conjuntos de dados e parâmetros de modelo. A distribuição do trabalho permite que tarefas intensivas sejam divididas e processadas em paralelo, acelerando o treinamento de modelos e a inferência, e possibilitando a manipulação de datasets que seriam inviáveis em um ambiente centralizado.

Na prática, escalabilidade horizontal (adicionar nós) supera a mera escala vertical (mais memória/CPU em um único nó) quando o objetivo é reduzir tempo de treinamento e permitir modelos/datasets que não cabem em um servidor. Em visão computacional, por exemplo, treinamentos clássicos como o ResNet-50 em ImageNet só atingem tempo-para-resultado competitivo ao paralelizar o SGD de maneira síncrona entre muitos GPUs: “Distributed synchronous SGD offers a potential solution” aos tempos de treinamento excessivos.

Além da escalabilidade, os sistemas distribuídos oferecem tolerância a falhas e alta disponibilidade, características cruciais para aplicações de AI/ML em produção. A falha de um único nó em um sistema distribuído não necessariamente interrompe o serviço, pois outros nós podem assumir a carga de trabalho. Isso garante a continuidade das operações e a resiliência das soluções de AI/ML.

Ainda se tem, o desempenho, confiabilidade operacional pesa: sistemas distribuídos fornecem tolerância a falhas (falhas parciais não derrubam o serviço), alta disponibilidade (replicação e failover) e elasticidade (autoescalonamento), o que é vital quando modelos alimentam APIs de predição em tempo real, pipelines de dados 24/7 ou aplicações críticas (ex.: detecção de fraude, monitoramento de infraestrutura). Esses requisitos de produção tornam a distribuição não só desejável, mas imperativa.

3.2. APLICAÇÕES DE IA COM SISTEMAS DISTRIBUÍDOS

A Inteligência Artificial não apenas se beneficia dos sistemas distribuídos, mas também desempenha um papel fundamental na otimização e aprimoramento desses sistemas, tornando-os mais eficientes, confiáveis e seguros. As principais aplicações de AI em sistemas distribuídos incluem:

- **Otimização e Balanceamento de Carga:** Algoritmos de AI são utilizados para distribuir dinamicamente as cargas de trabalho entre os nós do sistema, prevenindo gargalos e garantindo o uso eficiente dos recursos. Diferente dos métodos tradicionais estáticos, a AI pode se adaptar a mudanças em tempo real nos padrões de tráfego e

carga do servidor, tomando decisões instantâneas para otimizar a distribuição de tarefas.

- **Deteção de Falhas e Recuperação:** A AI pode monitorar sistemas distribuídos para identificar rapidamente falhas ou anomalias, como mau funcionamento de hardware ou bugs de software. Modelos de Machine Learning treinados com dados históricos podem prever falhas antes que ocorram, permitindo a implementação de medidas preventivas. Quando falhas são detectadas, sistemas de AI podem automaticamente acionar processos de recuperação, como reiniciar serviços ou redirecionar tráfego.
- **Gerenciamento de Recursos:** A AI otimiza a alocação de recursos como CPU, memória, armazenamento e largura de banda de rede, baseando-se na análise da demanda da carga de trabalho e na previsão de padrões de uso futuros. Modelos de AI consideram fatores como prioridade da aplicação, disponibilidade atual de recursos e tendências históricas para tomar decisões informadas, resultando em economia de custos e melhor desempenho do sistema.
- **Segurança e Deteção de Anomalias:** A AI aprimora a segurança em sistemas distribuídos monitorando continuamente comportamentos anormais que podem indicar violações de segurança. Técnicas como detecção de anomalias, reconhecimento de padrões e análise comportamental ajudam a identificar atividades incomuns, como tentativas de acesso não autorizado ou ataques de negação de serviço distribuído (DDoS). A AI também pode automatizar respostas a essas ameaças, isolando componentes afetados ou bloqueando tráfego malicioso.
- **Manutenção Preditiva:** Com base em dados históricos e monitoramento em tempo real, a AI pode prever quando componentes de um sistema distribuído podem falhar. Ao identificar padrões e correlações, os modelos de AI podem antecipar possíveis falhas e agendar atividades de manutenção antes que os problemas ocorram, minimizando o tempo de inatividade não planejado e estendendo a vida útil dos componentes do sistema.

3.3. TÉCNICAS DE IA UTILIZADAS

Diversas técnicas de Inteligência Artificial são empregadas em sistemas distribuídos para aprimorar seu desempenho, confiabilidade e segurança:

3.3.1 Machine Learning (ML): Inclui Aprendizado Supervisionado (classificação, regressão), Não Supervisionado (clustering, detecção de anomalias) e por Reforço (tomada de decisão, otimização de roteamento).

- **Aprendizado Supervisionado:** Utilizado em tarefas onde dados rotulados estão disponíveis, como classificação (ex: detecção de spam) e regressão (ex: previsão de demanda de recursos).



UNITINS
UNIVERSIDADE ESTADUAL DO TOCANTINS

TOCANTINS
GOVERNO DO ESTADO



- **Aprendizado Não Supervisionado:** Empregado para agrupamento (clustering) e detecção de anomalias em dados não rotulados (ex: identificação de padrões de tráfego incomuns).
- **Aprendizado por Reforço (RL):** Usado para tarefas de tomada de decisão, onde o sistema aprende ações ótimas através de tentativa e erro (ex: otimização de roteamento de rede ou balanceamento de carga dinâmico).

3.3.2 Deep Learning (DL): Utiliza CNNs (análise de imagens), RNNs (dados sequenciais, séries temporais) e Transformers (NLP, análise de logs).

- **Redes Neurais Convolucionais (CNNs):** Eficazes para análise de imagens e vídeos (ex: monitoramento de infraestrutura por visão computacional).
- **Redes Neurais Recorrentes (RNNs):** Adequadas para análise de dados sequenciais, como séries temporais (ex: previsão de desempenho do sistema) e Processamento de Linguagem Natural (NLP).
- **Transformers:** Modelos avançados para tarefas de NLP, permitindo uma melhor compreensão e geração de linguagem humana, aplicáveis em análise de logs e comunicação entre serviços.

3.3.3 Detecção de Anomalias: Emprega Métodos Estatísticos, Clustering e Autoencoders para identificar outliers.

- **Métodos Estatísticos:** Utilizam medidas estatísticas para identificar outliers.
- **Clustering:** Agrupa pontos de dados semelhantes e identifica outliers como anomalias.
- **Autoencoders:** Redes neurais treinadas para reconstruir dados de entrada, onde anomalias são detectadas como dados mal reconstruídos.

3.3.4 Processamento de Linguagem Natural (NLP): Abrange Tokenização, Reconhecimento de Entidades Nomeadas (NER) e Análise de Sentimento (feedback de usuários, logs).

- **Tokenização:** Divide o texto em partes gerenciáveis para análise.
- **Reconhecimento de Entidades Nomeadas (NER):** Identifica e classifica entidades no texto.
- **Análise de Sentimento:** Determina o sentimento expresso no texto, útil para feedback de usuários ou análise de logs.

4. ARQUITETURAS E FRAMEWORKS DE ML EM SISTEMAS DISTRIBUÍDOS

4.1. PROCESSAMENTO PARALELO DE ML DISTRIBUÍDO

Treinar modelos grandes de forma eficiente exige paralelizar computação e I/O. Em DL, isso se materializa em coletivas de comunicação (p.ex., AllReduce para sincronizar gradientes) e estratégias de particionamento que reduzem tempo de época mantendo acurácia. Resultados clássicos mostram que, com ajustes de hiperparâmetros, minibatches gigantes

preservam generalização e levam a speedups expressivos (“no loss of accuracy when training with large minibatch sizes up to 8192 images”).

O processamento paralelo é a espinha dorsal do Machine Learning distribuído, permitindo que modelos complexos sejam treinados em grandes volumes de dados de forma eficiente. A ideia central é dividir as tarefas computacionais e executá-las simultaneamente em múltiplos processadores ou máquinas. Isso é crucial porque o treinamento de modelos de ML e DL envolve operações matemáticas intensivas, como álgebra de matrizes e otimização, em conjuntos de dados massivos.

Historicamente, o aumento da capacidade computacional era alcançado escalando-se verticalmente (aumentando o número de núcleos, memória, etc., em uma única máquina). No entanto, para as demandas atuais, a abordagem mais eficaz é a escalabilidade horizontal, que envolve a adição de mais máquinas (nós) a um sistema distribuído. Essa estratégia permite que tarefas demoradas e improdutivas se tornem mais curtas devido à paralelização, sem a necessidade de investir em hardware de ponta excessivamente caro.

4.2. TREINAMENTOS DE MODELOS DISTRIBUÍDOS (Paralelismo de Dados vs. Paralelismo de Modelos)

No treinamento de modelos distribuídos, o modelo e os dados de treinamento são divididos entre várias máquinas. As duas abordagens básicas para implementar isso são o paralelismo de dados e o paralelismo de modelos:

- **Paralelismo de Dados:** Nesta abordagem, o conjunto de dados é dividido em n partes, onde n é o número de nós de trabalho no cluster. Cada nó de trabalho possui uma cópia completa do modelo e treina essa cópia usando um subconjunto diferente dos dados. Os gradientes calculados por cada nó são então agregados (sincronamente ou assincronamente) para atualizar os parâmetros do modelo principal. Esta técnica é frequentemente utilizada para treinar grandes redes neurais devido à menor necessidade de comunicação entre as máquinas, embora possa apresentar taxas de convergência mais lentas se os gradientes variarem muito.
- **Paralelismo de Modelos:** Quando o modelo neural é muito grande para caber na memória de uma única máquina, seus parâmetros são divididos e distribuídos entre vários computadores. Cada máquina processa uma parte do modelo e uma porção dos dados de entrada, determinando a saída apropriada. A saída final é então combinada.

O paralelismo de modelos reduz os requisitos de memória para cada computador, permitindo o treinamento de modelos maiores. No entanto, exige um design cuidadoso e otimização para minimizar os custos de comunicação, pois as máquinas frequentemente trocam parâmetros do modelo.



UNITINS
UNIVERSIDADE ESTADUAL DO TOCANTINS

TOCANTINS
GOVERNO DO ESTADO



4.3. ALGORITMOS DE ML DISTRIBUÍDO

Para gerenciar a distribuição de tarefas e a sincronização de modelos em ambientes de ML distribuído, algoritmos específicos são empregados:

- **Parameter Server:** Nesta arquitetura, os pesos e bias de um modelo de Machine Learning são distribuídos para vários computadores em um cluster. Cada computador no cluster mantém uma cópia do modelo, e um servidor de parâmetros centralizado gerencia as modificações e atualizações desses parâmetros. Os nós de trabalho calculam os gradientes localmente e os enviam para o servidor de parâmetros, que os agrega e atualiza o modelo, enviando os novos parâmetros de volta aos nós de trabalho.
- **AllReduce:** Este método é utilizado para sincronizar os pesos do modelo entre todos os nós computacionais. Em vez de um servidor central, cada nó calcula seus gradientes e, em seguida, todos os nós se comunicam diretamente para somar os gradientes e atualizar seus próprios pesos do modelo. Isso pode ser mais eficiente para certas topologias de rede e cargas de trabalho, pois evita um gargalo centralizado.

4.4. FRAMEWORKS E PLATAFORMAS

Diversos frameworks e plataformas foram desenvolvidos para facilitar a implementação e o gerenciamento de ML distribuído, abstraindo grande parte da complexidade subjacente. Alguns exemplos notáveis incluem:

- **TensorFlow Distributed:** Uma extensão do TensorFlow que permite o treinamento de modelos em clusters de máquinas, suportando tanto paralelismo de dados quanto de modelos.
- **PyTorch Distributed:** Oferece ferramentas para computação paralela e distribuída, com foco em flexibilidade e facilidade de uso para pesquisadores.
- **Apache Spark MLlib:** Uma biblioteca de Machine Learning escalável que roda no Apache Spark, ideal para processamento de big data e ML distribuído, utilizando o conceito de RDDs (Resilient Distributed Datasets) para paralelização [2].
- **Horovod:** Um framework de comunicação distribuída para TensorFlow, Keras e PyTorch, que implementa o algoritmo AllReduce para um treinamento de modelo mais eficiente em clusters de GPUs.
- **Ray:** Um framework de computação distribuída de código aberto que fornece uma API simples para construir e executar aplicações distribuídas, incluindo Machine Learning e Deep Learning.



TOCANTINS
GOVERNO DO ESTADO



Esses frameworks e plataformas fornecem as ferramentas necessárias para orquestrar o treinamento e a inferência de modelos em ambientes distribuídos, lidando com aspectos como comunicação entre nós, sincronização de parâmetros e tolerância a falhas.

CONCLUSÃO

A integração entre Sistemas Distribuídos e Inteligência Artificial (AI) e Machine Learning (ML) é um pilar fundamental para o avanço tecnológico. A demanda por processamento de grandes volumes de dados e a complexidade dos modelos de AI/ML tornam os sistemas distribuídos essenciais para escalabilidade, eficiência e resiliência. A AI, por sua vez, otimiza os próprios sistemas distribuídos, aprimorando balanceamento de carga, detecção de falhas, gerenciamento de recursos e segurança.

As arquiteturas e frameworks para ML distribuído, como paralelismo de dados e de modelos, e algoritmos como Parameter Server e AllReduce, são cruciais para lidar com os desafios computacionais. A superação de obstáculos como overhead de comunicação, desbalanceamento de dados e sincronização complexa exige melhores práticas e um robusto sistema de monitoramento. A sinergia entre esses campos é um ciclo virtuoso, impulsionando a inovação e permitindo soluções poderosas e eficientes para os desafios modernos.

REFERÊNCIAS

GeeksforGeeks. Role of AI in Distributed Systems. Disponível em: <https://www.geeksforgeeks.org/artificial-intelligence/role-of-ai-in-distributed-systems/>.

Acesso em: 19 set. 2025.

XenonStack. Distributed Machine Learning Frameworks and its Benefits. Disponível em: <https://www.xenonstack.com/blog/distributed-ml-framework>. Acesso em: 19 set. 2025.

FRANÇA, João Victor Póvoa. Lista-1-Sistemas-Distribuidos.pdf. Material fornecido pelo usuário. Acesso em: 19 set. 2025.

SERGIO. Inteligência Artificial distribuída <https://egcportalantigo.tcm.sp.gov.br/images/conteudo-palestras/X-educontas/26-08-2019/inteligencia-artificial-distribuida.pdf>. Acesso em 19 de set 2025.

Sistemas Distribuídos - Fundamentação. Material fornecido pelo usuário. Acesso em: 19 set. 2025.

Sistemas Distribuídos - Arquiteturas. Material fornecido pelo usuário. Acesso em: 19 set. 2025.

TANYO. Are large-scale distributed systems the future of AI? <https://www.quora.com/Are-large-scale-distributed-systems-the-future-of-AI>. Acesso em 19 de set 2025.

WIKIPEDIA. Inteligência Artificial distribuída https://pt.wikipedia.org/wiki/Intelig%C3%Aancia_artificial_distribu%C3%ADa. Acesso em 19 de set 2025.