

# PROCESSAMENTO DE LINGUAGEM NATURAL

EPISODE VII : LLM

Alex Marino

14 de novembro de 2024

# Conteúdo

Introdução aos Large Language Models (LLMs)

Arquitetura de Large Language Models (LLMs)

Treinamento de Large Language Models (LLMs)

Aplicações de Large Language Models (LLMs)

Desafios e Limitações dos Large Language Models (LLMs)

Tendências Futuras e Melhorias em Large Language Models (LLMs)

Aplicações de Large Language Models (LLMs)

Tópicos sobre Large Language Models (LLMs) no LifeArchitect.ai

Tamanho e Escala dos Modelos

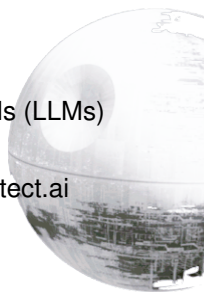
- Capacidades Emergentes e Habilidades dos Modelos

- Comparações de Modelos Populares

- Janelas de Contexto e Computação

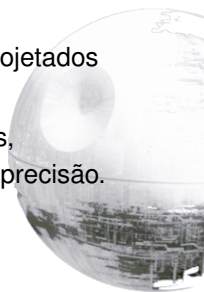
- Datasets e Fontes de Dados

Referências



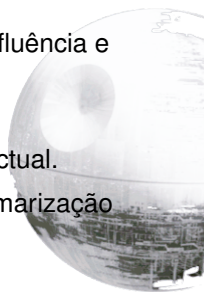
# Introdução aos Large Language Models (LLMs)

- ▶ LLMs são modelos de linguagem de grande escala, projetados para gerar e compreender textos complexos.
- ▶ Caracterizam-se por milhões ou bilhões de parâmetros, permitindo-lhes processar linguagem natural com alta precisão.
- ▶ Exemplos conhecidos: GPT-3, GPT-4, PaLM, LLaMA.



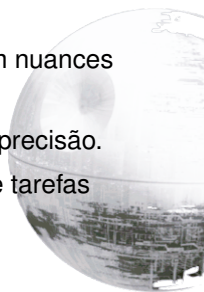
# O que são Large Language Models?

- ▶ **Modelos de Linguagem** capazes de gerar texto com fluência e coerência.
- ▶ Treinados em grandes quantidades de dados textuais, abrangendo desde a gramática até o conhecimento factual.
- ▶ Desempenham tarefas como tradução automática, sumarização e question answering.



# Diferenciais dos LLMs

- ▶ **Escala:** Modelos com bilhões de parâmetros capturam nuances de linguagem e conhecimento contextual.
- ▶ **Precisão:** LLMs realizam tarefas complexas com alta precisão.
- ▶ **Versatilidade:** Aplicáveis em uma ampla variedade de tarefas sem necessidade de re-treinamento.



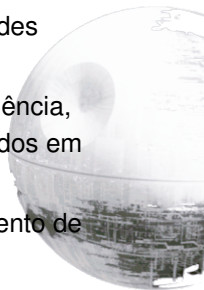
# Arquitetura de Large Language Models

- ▶ Baseada na arquitetura Transformer, introduzida por Vaswani et al. (2017).
- ▶ Utiliza mecanismos de self-attention para analisar relações contextuais em grandes volumes de texto.
- ▶ Escalável para um número muito grande de parâmetros.



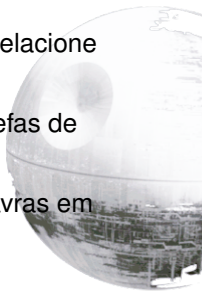
# Transformer como Base para LLMs

- ▶ Estrutura composta por camadas de autoatenção e redes feed-forward.
- ▶ Permite processamento paralelo das palavras na sequência, tornando os LLMs mais eficientes que modelos baseados em RNNs.
- ▶ Alta escalabilidade devido à separação de processamento de tokens por cabeças de atenção.



# Self-Attention e Relações Contextuais

- ▶ **Self-attention** permite que cada token de entrada se relacione com todos os outros da sequência.
- ▶ Captura contexto de longo alcance, essencial para tarefas de linguagem complexas.
- ▶ Essencial para o entendimento do significado das palavras em diferentes contextos.





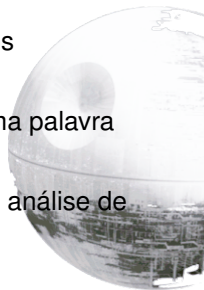
## Escalabilidade e Parâmetros em LLMs

- ▶ LLMs com bilhões de parâmetros são capazes de generalizar conhecimento de grandes volumes de dados.
- ▶ Quanto maior o número de parâmetros, mais capacidade de armazenamento de conhecimento, mas exige maior poder computacional.
- ▶ Exemplo: GPT-3 possui 175 bilhões de parâmetros.



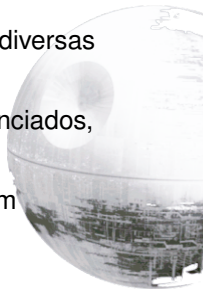
# Treinamento de Large Language Models

- ▶ Treinados em grandes datasets, incluindo livros, artigos científicos e conteúdos online.
- ▶ O treinamento visa ensinar o modelo a prever a próxima palavra com base no contexto.
- ▶ Modelos são ajustados para tarefas específicas, como análise de sentimentos e resposta a perguntas.



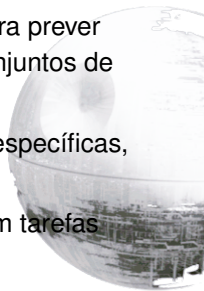
# Dataset em Larga Escala e Questões Éticas

- ▶ Treinamento em grandes volumes de dados de fontes diversas levanta preocupações de privacidade.
- ▶ Sessões de treinamento incluem textos públicos e licenciados, mas ainda há desafios em dados sensíveis.
- ▶ Questões éticas sobre viés e impactos dos modelos em diferentes populações.



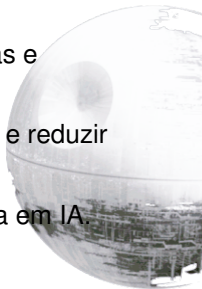
# Pré-treinamento e Fine-tuning

- ▶ **Pré-treinamento:** LLMs são inicialmente treinados para prever palavras de modo não supervisionado em grandes conjuntos de dados.
- ▶ **Fine-tuning:** Os modelos são ajustados para tarefas específicas, como chatbots ou assistentes virtuais.
- ▶ Processo de fine-tuning permite alta especialização em tarefas específicas [2].



# Alinhamento com Intenções Humanas

- ▶ Modelos são ajustados para seguir instruções humanas e responder de maneira ética e segura.
- ▶ O alinhamento visa minimizar respostas inapropriadas e reduzir vieses.
- ▶ Alinhamento ético é um campo emergente na pesquisa em IA.



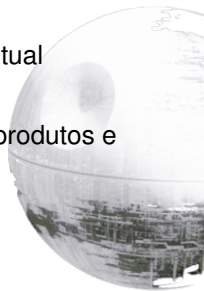
# Aplicações de Large Language Models

- ▶ LLMs são amplamente utilizados em uma variedade de tarefas de processamento de linguagem natural.
- ▶ Capacidades incluem geração de texto, classificação de sentimentos, assistentes virtuais e muito mais.
- ▶ Modelos versáteis que se adaptam a diversas aplicações com alto desempenho.



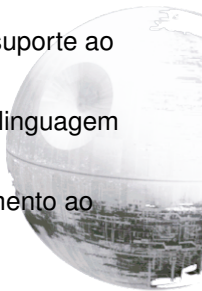
# Geração de Texto

- ▶ LLMs são amplamente usados para criar conteúdo textual coerente e relevante.
- ▶ Aplicações incluem criação de artigos, descrições de produtos e assistentes de escrita.
- ▶ Exemplos: ChatGPT, Copilot, Jasper AI [3].



# Assistentes Virtuais e Chatbots

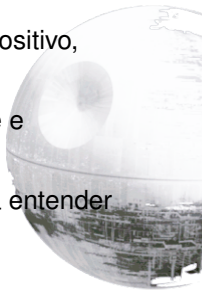
- ▶ Usados para automação de atendimento ao cliente e suporte ao usuário.
- ▶ Capacidade de entender e responder a perguntas em linguagem natural.
- ▶ Exemplo: Aplicações de IA em empresas para atendimento ao cliente (e.g., chatbots bancários).





# Análise de Sentimentos e Classificação de Texto

- ▶ LLMs são usados para identificar o tom de um texto (positivo, negativo, neutro).
- ▶ Análise de sentimentos em redes sociais, e-commerce e feedback de clientes.
- ▶ Exemplo: Aplicações em marketing digital e CRM para entender a percepção do consumidor.



# Resumo e Tradução

- ▶ Resumo automático de grandes volumes de texto, como documentos e artigos científicos.
- ▶ Tradução precisa entre múltiplos idiomas, melhorando a comunicação global.
- ▶ Exemplo: Ferramentas de tradução e serviços de síntese de documentos.



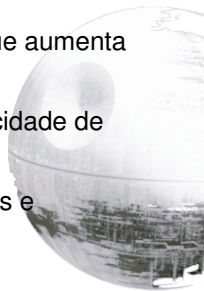
# Desafios e Limitações dos Large Language Models

- ▶ LLMs enfrentam desafios técnicos, éticos e operacionais ao serem implantados em produção.
- ▶ Limitados por infraestrutura e questões de segurança, privacidade e viés.
- ▶ Necessidade de controles e melhorias constantes para garantir uso seguro e ético.



# Escalabilidade e Custos Computacionais

- ▶ LLMs exigem recursos computacionais massivos, o que aumenta o custo de treinamento e inferência.
- ▶ A escalabilidade requer clusters de GPUs e alta capacidade de memória.
- ▶ Exemplos de custo incluem uso de supercomputadores e infraestrutura em nuvem [1].



## Viés nos Dados e Alinhamento Ético

- ▶ LLMs podem herdar vieses dos dados em que foram treinados, levando a respostas enviesadas.
- ▶ Alinhamento ético é essencial para reduzir respostas inadequadas e impactos negativos.
- ▶ Exemplo: Ferramentas que ajustam e monitoram o viés no uso de LLMs.



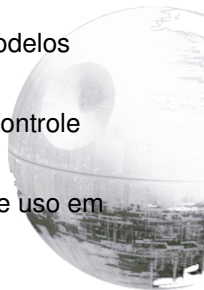
## Interpretação e Controle das Respostas

- ▶ LLMs podem produzir respostas imprevisíveis e difíceis de controlar em tempo real.
- ▶ Ferramentas de controle são usadas para limitar e ajustar as saídas.
- ▶ Exemplo: Intervenção humana em cenários de uso sensível, como atendimento ao cliente.



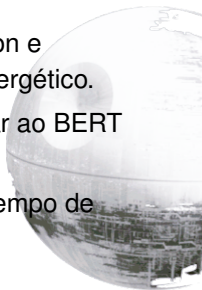
# Tendências Futuras e Melhorias em LLMs

- ▶ A pesquisa em LLMs está avançando em direção a modelos mais eficientes e éticos.
- ▶ Foco em reduzir o custo computacional e melhorar o controle sobre as saídas.
- ▶ Melhoria contínua para atender às novas demandas de uso em diversos setores.



# Modelos mais Eficientes e Técnicas de Redução

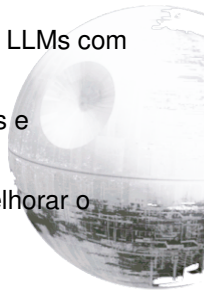
- ▶ Desenvolvimento de modelos menores, como distillation e quantization, para redução de tamanho e consumo energético.
- ▶ Exemplo: DistilBERT, que mantém desempenho similar ao BERT com menor consumo [4].
- ▶ A quantização diminui o uso de memória e acelera o tempo de inferência.





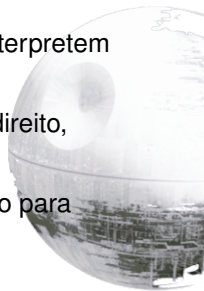
# Aprimoramento do Alinhamento e Controle Ético

- ▶ Pesquisa focada em melhorar o alinhamento ético dos LLMs com as expectativas humanas.
- ▶ Técnicas de treinamento e controle para reduzir vieses e respostas indesejadas.
- ▶ Exemplo: Programas de auditoria para monitorar e melhorar o comportamento dos modelos.



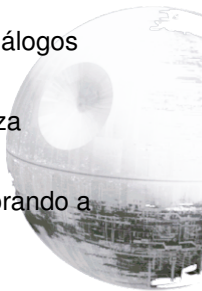
# Aplicações Emergentes e Avanços Tecnológicos

- ▶ Avanços em multimodalidade, permitindo que LLMs interpretem texto, imagens e som.
- ▶ Expansão para novos domínios, como biomedicina e direito, onde o NLP está se tornando crítico.
- ▶ Ferramentas de geração de código e modelos de apoio para desenvolvimento de software.



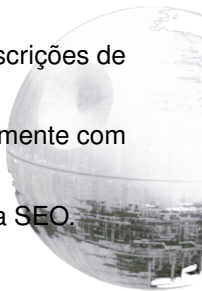
# Assistentes Virtuais e Chatbots Inteligentes

- ▶ **Aplicação:** Utilizados para atendimento ao cliente e diálogos complexos.
- ▶ **Exemplo:** ChatGPT responde a perguntas e automatiza interações com clientes.
- ▶ **Benefícios:** Redução de custos e suporte 24/7, melhorando a experiência do cliente.



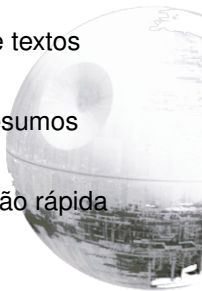
## Geração de Conteúdo e Assistência em Escrita

- ▶ **Aplicação:** Criação de artigos, postagens de blog, descrições de produtos.
- ▶ **Exemplo:** Jasper AI e Copy.ai geram texto automaticamente com instruções específicas.
- ▶ **Benefícios:** Acelera a produção de conteúdo e otimiza SEO.



# Tradução Automática e Sumarização de Textos

- ▶ **Aplicação:** Tradução em tempo real e sumarização de textos extensos.
- ▶ **Exemplo:** GPT-3 e T5 são usados para traduções e resumos com alta precisão.
- ▶ **Benefícios:** Facilita a comunicação multilíngue e revisão rápida de informações.



# Análise de Sentimentos e Classificação de Textos

- ▶ **Aplicação:** Identificação do tom de textos para monitoramento em redes sociais e feedback.
- ▶ **Exemplo:** TextBlob com LLMs monitora a opinião pública e feedback do cliente.
- ▶ **Benefícios:** Permite reação rápida a tendências e insights em tempo real sobre satisfação.



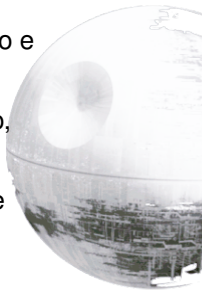
# Pesquisa e Consulta em Documentação Técnica

- ▶ **Aplicação:** Consulta a grandes volumes de documentação técnica e jurídica.
- ▶ **Exemplo:** LLMs são usados em escritórios de advocacia para buscar trechos de leis e regulamentos.
- ▶ **Benefícios:** Acelera a pesquisa, reduzindo o tempo para encontrar informações específicas.



# Desenvolvimento de Software e Assistência na Codificação

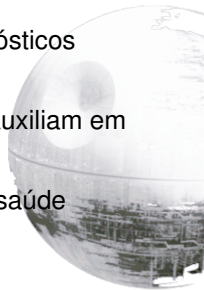
- ▶ **Aplicação:** Ajudam desenvolvedores a escrever código e sugerem exemplos e correções.
- ▶ **Exemplo:** GitHub Copilot oferece sugestões de código, autocompletar e soluções.
- ▶ **Benefícios:** Aumenta a produtividade e reduz erros de codificação.





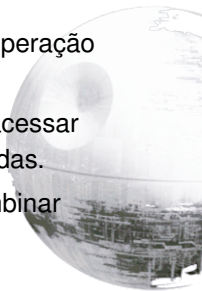
# Medicina e Pesquisa Biomédica

- ▶ **Aplicação:** Análise de publicações científicas e diagnósticos com base no histórico médico.
- ▶ **Exemplo:** LLMs sintetizam informações de artigos e auxiliam em diagnósticos médicos.
- ▶ **Benefícios:** Facilita a atualização de profissionais de saúde sobre avanços científicos.



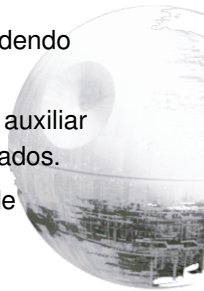
# Sistemas de RAG para Consultoria de Conhecimento

- ▶ **Aplicação:** Integração de LLMs com sistemas de recuperação de informações atualizadas.
- ▶ **Exemplo:** Consultorias financeiras usam RAGs para acessar relatórios financeiros e responder a perguntas detalhadas.
- ▶ **Benefícios:** Melhora a precisão das respostas ao combinar conhecimento atualizado.



## Educação e Tutoria Automatizada

- ▶ **Aplicação:** LLMs atuam como tutores virtuais, respondendo dúvidas e adaptando explicações.
- ▶ **Exemplo:** Plataformas educacionais usam LLMs para auxiliar alunos, oferecendo explicações e materiais personalizados.
- ▶ **Benefícios:** Suporte educacional adaptado ao estilo de aprendizado do usuário.



O site apresenta um panorama detalhado sobre o crescimento do tamanho dos modelos de linguagem, com ênfase na quantidade de parâmetros que aumentou drasticamente ao longo dos anos. Este aumento de escala permite que os modelos capturem mais nuances e informações contextuais.

- ▶ **Aumento no número de parâmetros:** Os LLMs começaram com milhões de parâmetros, mas hoje possuem bilhões, chegando a trilhões em alguns casos.
- ▶ **Impacto da Escala:** Modelos maiores geralmente demonstram melhor desempenho em tarefas complexas de linguagem natural, mas requerem mais recursos computacionais.



Modelo	Ano de Lançamento	Número de Parâmetros
GPT-2	2019	1.5 bilhões
GPT-3	2020	175 bilhões
PaLM	2022	540 bilhões
Gopher	2021	280 bilhões
BLOOM	2022	176 bilhões

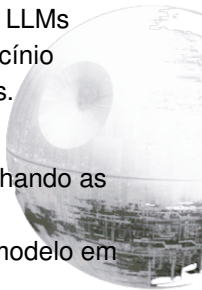
Tabela: Exemplos de Modelos e seus Tamanhos em Parâmetros

O LifeArchitect.ai explora as capacidades emergentes que os LLMs adquiriram, destacando como esses modelos demonstram habilidades complexas, especialmente quando pré-treinados em dados extensos e variados.

- **Aprendizado de Poucas Instruções (Few-Shot Learning):**  
LLMs como GPT-3 conseguem realizar tarefas com um mínimo de exemplos, permitindo adaptação rápida a novas instruções.

- ▶ **Habilidades Compostas:** A capacidade de realizar múltiplas tarefas complexas em linguagens variadas, incluindo geração de texto, tradução e sumarização.
- ▶ **Capacidades de Raciocínio e Interpretação:** Alguns LLMs apresentam habilidade de completar tarefas com raciocínio lógico e contextual, embora ainda enfrentem limitações.

O site fornece comparações entre os principais LLMs, detalhando as diferenças em tamanho, arquitetura e desempenho. Estas comparações ajudam a entender a aplicabilidade de cada modelo em tarefas específicas.



Modelo	Parâmetros	Arquitetura	Tarefas
GPT-3	175 bilhões	Transformer	Geração de texto, tradução, re
PaLM	540 bilhões	Transformer	Aprendizado de instruções, multim
BLOOM	176 bilhões	Multilíngue	Tradução, sumarização
Gopher	280 bilhões	Transformer	Perguntas e respostas, conhecime

Tabela: Comparação entre Modelos Populares de LLMs

O site discute o impacto das janelas de contexto na capacidade dos LLMs em lidar com tarefas que envolvem contexto longo e em escala maior.

- ▶ **Definição de Janela de Contexto:** A janela de contexto representa o número máximo de tokens (palavras ou caracteres) que o modelo pode processar de uma só vez.
- ▶ **Necessidade de Computação e Memória:** Modelos com janelas de contexto maiores exigem mais memória e recursos computacionais, influenciando o custo de execução.

Modelo	Tamanho da Janela de Contexto	Uso Comum
GPT-3	2048 tokens	Geração de texto e diálogo
PaLM	4096 tokens	Tarefas com contexto estendido
Longformer	4096 tokens	Tarefas com documentos longos

Tabela: Janelas de Contexto em Diferentes Modelos

LifeArchitect.ai detalha os principais conjuntos de dados usados para treinar LLMs e aborda as implicações éticas e técnicas do uso de grandes volumes de dados.

- ▶ **Datasets Massivos:** Os LLMs são treinados em datasets como Common Crawl, Wikipedia, e The Pile, que abrangem bilhões de palavras e uma diversidade de tópicos.
- ▶ **Desafios Éticos:** O uso de dados da internet levanta questões sobre privacidade e viés, pois os dados podem conter informações sensíveis ou preconceituosas.



<b>Dataset</b>	<b>Fonte</b>	<b>Descrição</b>
Common Crawl	Internet	Coleta de páginas web, diversidade
Wikipedia	Enciclopédia online	Conteúdo factual e enciclopédico
The Pile	Diversas	Inclui livros, artigos, sites de perguntas e respostas

Tabela: Principais Datasets Utilizados em LLMs



**MAY THE**  
**SOURCE**  
**BE WITH**  
**YOU**



- [1] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [2] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [3] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI Blog*, 1:9, 2019.
- [4] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.

