

MODELOS, TÉCNICAS E PESQUISA

Curso de Formação de Auditor na ISO | IEC 42.001



Facilitador
Paulo Carvalho
Chief Artificial Intelligence Officer

Inteligência Artificial 2030

A Inteligência Artificial avança em ritmo acelerado, com modelos cada vez mais acessíveis e potentes, exigindo das organizações uma resposta estratégica e adaptativa. Para o Auditor Interno, isso significa a necessidade urgente de **manter-se atualizado sobre as aplicações atuais e emergentes de IA**, compreendendo não só os riscos, mas também o potencial de uso ético e eficaz desses sistemas. A norma ISO/IEC 42001:2024 oferece diretrizes essenciais para estruturar essa governança com foco em segurança, transparência e responsabilidade. Em um cenário de convergência tecnológica, onde IA, biotecnologia e dispositivos conectados se integram, auditar sem esse entendimento é ignorar o contexto real da operação. Assim, o domínio da governança algorítmica torna-se vital para assegurar o uso confiável, sustentável e centrado no ser humano da inteligência artificial.

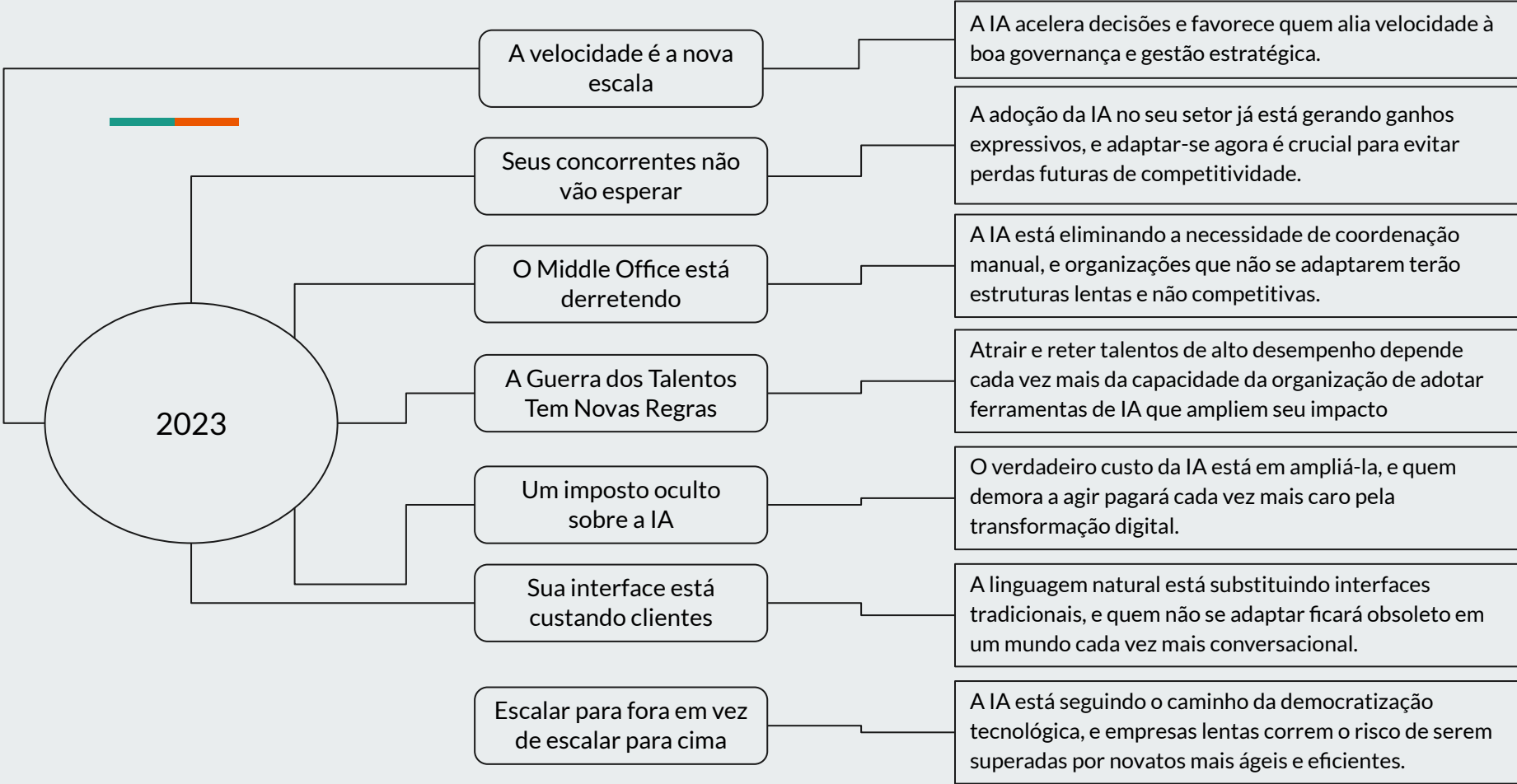
Paulo Carvalho





Quais organizações irão prosperar

IA 2030 - Tendências



2023

A velocidade é a nova
escala

A IA acelera decisões e favorece quem alia velocidade à boa governança e gestão estratégica.

Seus concorrentes não
vão esperar

A adoção da IA no seu setor já está gerando ganhos expressivos, e adaptar-se agora é crucial para evitar perdas futuras de competitividade.

O Middle Office está
derretendo

A IA está eliminando a necessidade de coordenação manual, e organizações que não se adaptarem terão estruturas lentas e não competitivas.

A Guerra dos Talentos
Tem Novas Regras

Atrair e reter talentos de alto desempenho depende cada vez mais da capacidade da organização de adotar ferramentas de IA que ampliem seu impacto

Um imposto oculto
sobre a IA

O verdadeiro custo da IA está em ampliá-la, e quem demora a agir pagará cada vez mais caro pela transformação digital.

Sua interface está
custando clientes

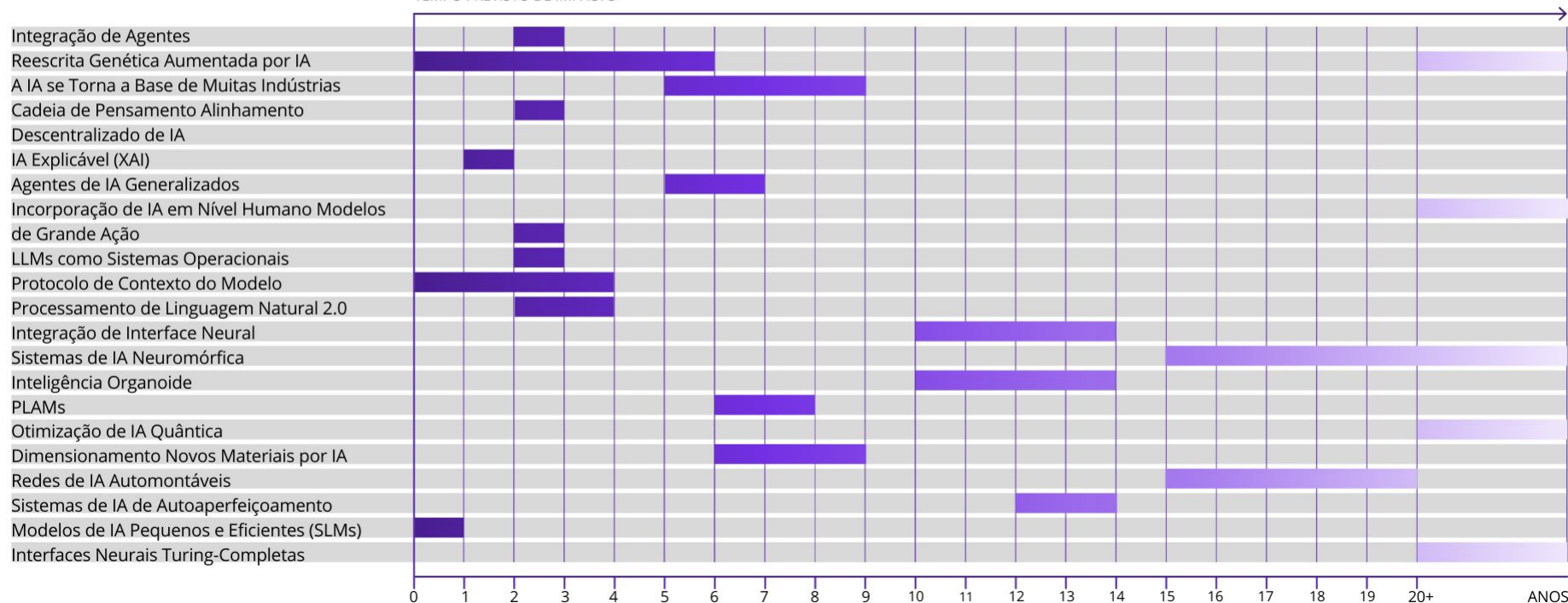
A linguagem natural está substituindo interfaces tradicionais, e quem não se adaptar ficará obsoleto em um mundo cada vez mais conversacional.

Escalar para fora em vez
de escalar para cima

A IA está seguindo o caminho da democratização tecnológica, e empresas lentas correm o risco de serem superadas por novatos mais ágeis e eficientes.

As tecnologias generativas avançam nos próximos anos, enquanto métodos de computação como computação organoide e neuromórfica impulsionam desenvolvimentos a longo prazo.

TEMPO PREVISTO DE IMPACTO





Desenvolvimento de alto nível e curto prazo

IA 2030 - Disrupção

ESCALAÇÃO



A personalização profunda por IA depende do acesso a grandes volumes de dados, mas sistemas fechados e regulações ainda limitam esse potencial.

INVESTIMENTO



O avanço da IA depende do equilíbrio entre investimento em P&D e retorno comercial, moldado pela paciência e visão dos investidores.

RESTRIÇÕES À ADOÇÃO



Mesmo tecnologias maduras enfrentam barreiras culturais e estratégicas que limitam seu impacto, especialmente em setores tradicionais como saúde, seguros e finanças.

REGULAMENTOS



O avanço da IA tem superado a regulação, e a adoção dependerá de como as normas locais se alinham ou entram em conflito entre si.

MENÇÕES NA MÍDIA



O entusiasmo popular e a cobertura da mídia podem acelerar a adoção da IA, mesmo sem avanços técnicos significativos.

PERCEPÇÃO PÚBLICA



A percepção pública sobre a IA moldará sua demanda, especialmente em áreas sensíveis como educação, criatividade, propriedade intelectual e desinformação.

DESENVOLVIMENTOS DE
P&D

Descobertas em IA são imprevisíveis e dependem de fatores como financiamento e equipe, sendo tratadas como variáveis críticas mas incertas no planejamento estratégico.



Desenvolvimento de alto nível e curto prazo

IA 20230 - Termos

AGENTE IA

Agentes de IA são entidades autônomas que percebem o ambiente, tomam decisões e agem para alcançar objetivos, variando de automações simples a sistemas multimodais avançados.

AGI (INTELIGÊNCIA GERAL ARTIFICIAL)

AGI refere-se a sistemas de IA com capacidades cognitivas iguais ou superiores às humanas em todas as tarefas valiosas, ainda teórica, mas com profundas implicações para trabalho, governança e segurança global.

ÉTICA DA IA

A ética da IA é um campo multidisciplinar que analisa os riscos sociais, econômicos e éticos da tecnologia, guiando políticas e regulamentações para alinhar a IA aos valores humanos.

GOVERNANÇA DE IA

A governança de IA engloba sistemas e políticas que regulam seu desenvolvimento e uso, sendo essencial para mitigar riscos, promover concorrência justa e lidar com tensões geopolíticas.

ALGORITMO

Algoritmos são conjuntos estruturados de regras usados em IA para processar dados, gerar insights e tomar decisões de forma automatizada.

ALINHAMENTO

O alinhamento de IA é o processo de garantir que suas decisões e comportamentos reflitam intenções humanas, princípios éticos e normas, prevenindo efeitos indesejados ou prejudiciais.

SUPERINTELIGÊNCIA ARTIFICIAL (ASI)

A Superinteligência Artificial (ASI) é um conceito hipotético de IA que supera a inteligência humana em todos os domínios, levantando sérios desafios de controle, governança e riscos existenciais.

RECONHECIMENTO AUTOMÁTICO DE FALA (ASR)

O Reconhecimento Automático de Fala (ASR) é um sistema de IA que converte fala em texto, viabilizando assistentes virtuais, transcrições e interfaces de voz em diversos aplicativos.

IA AUTÔNOMA

A IA autônoma toma decisões e executa tarefas sem intervenção humana, sendo vital em áreas como robótica e segurança, e exigindo fortes salvaguardas para seu uso responsável.

RACIOCÍNIO EM CADEIA DE PENSAMENTO (COT)

O raciocínio em cadeia é um método de IA que resolve problemas passo a passo, imitando a lógica humana, e aprimora decisões em áreas como matemática, direito e medicina.

VISÃO COMPUTACIONAL

A visão computacional é uma tecnologia de IA que interpreta imagens e vídeos, sendo aplicada em segurança, automação, medicina e veículos autônomos.

IA DE BORDA

A IA de ponta executa modelos diretamente em dispositivos como smartphones e sensores, oferecendo processamento em tempo real, menor latência e maior privacidade dos dados.

MODELO

Modelos de IA são sistemas treinados para analisar dados, fazer previsões e automatizar decisões, variando de regressões simples a redes neurais profundas complexas.

IA GENERATIVA (GERAL)

A IA generativa cria novos conteúdos como texto, imagens e vídeos, revolucionando setores criativos e de serviços, mas também levantando questões sobre direitos autorais e desinformação.

GPU (UNIDADE DE PROCESSAMENTO GRÁFICO)

GPUs são hardwares especializados em computação paralela, fundamentais para acelerar o treinamento de modelos de IA e executar tarefas intensivas em dados em larga escala.

MODELO DE FUNDAÇÃO

Modelos de IA em larga escala, pré-treinados em grandes volumes de dados e adaptáveis a diversas tarefas, servindo de alicerce para aplicações como IA generativa e automação.

IA MULTIMODAL

A IA multimodal integra diferentes tipos de dados, como texto, imagem e áudio, para melhorar a compreensão e decisões, impulsionando assistentes virtuais, chatbots e diagnósticos médicos.

PROCESSAMENTO DE LINGUAGEM NATURAL (PLN)

O Processamento de Linguagem Natural (PNL) permite que máquinas compreendam e gerem linguagem humana, sendo usado em chatbots, traduções, análise de sentimentos e moderação de conteúdo.

**PESQUISA DE ARQUITETURA NEURAL (NAS)**

A Otimização Neural Automatizada (AutoML) é um método de IA que ajusta automaticamente arquiteturas de redes neurais, aumentando o desempenho e minimizando a intervenção humana.

PARÂMETRO

Parâmetros são variáveis ajustadas durante o treinamento de um modelo de IA para melhorar sua precisão, sendo que grandes modelos podem conter bilhões deles, exigindo alto poder computacional.

ENGENHARIA DE PROMPT

A engenharia de prompts consiste em criar entradas estratégicas (PROMPTS) para orientar modelos de IA, sendo essencial para maximizar a eficácia da IA generativa em aplicações comerciais e criativas.

QUANTUM IA

A IA quântica combina computação quântica com inteligência artificial, prometendo acelerar o aprendizado de máquina e transformar áreas como criptografia, ciência dos materiais e otimização.

SISTEMAS DE RECOMENDAÇÃO

Algoritmos de recomendação usam IA para analisar o comportamento do usuário e sugerir conteúdos, produtos ou serviços, sendo amplamente aplicados em e-commerce, streaming e publicidade digital.

APRENDIZAGEM POR REFORÇO A PARTIR DO FEEDBACK HUMANO (RHLF)

O Aprendizado com Feedback Humano (RLHF) é um método de treino em que modelos de IA são ajustados com base em avaliações humanas, aprimorando precisão, alinhamento ético e aplicabilidade prática.

APRENDIZAGEM AUTO-SUPERVISADA

O aprendizado não supervisionado é uma abordagem de machine learning em que modelos identificam padrões em dados não rotulados, reduzindo a dependência de anotações humanas.

APRENDIZAGEM SUPERVISIONADA

O aprendizado supervisionado é um método em que modelos de IA são treinados com dados rotulados, utilizando pares de entrada e saída para aprimorar a precisão em previsões futuras.

IA SIMBÓLICA

A IA simbólica é uma abordagem que representa o conhecimento com símbolos e regras lógicas, permitindo raciocínio explicável e sendo comum em sistemas especialistas.

DADOS SINTÉTICOS

Dados sintéticos são gerados artificialmente para treinar IA quando dados reais são limitados, tendenciosos ou sensíveis, ajudando a melhorar o desempenho e reduzir riscos de privacidade.

DADOS DE TREINAMENTO

Os dados de treinamento são usados para ensinar modelos de IA a identificar padrões e fazer previsões, sendo que sua qualidade e diversidade afetam diretamente a precisão e imparcialidade do modelo.

IA CONFIÁVEL

A IA confiável é projetada com foco em transparência, justiça, responsabilidade e segurança, visando promover a confiança pública e atender às exigências regulatórias.

APRENDIZAGEM NÃO SUPERVISIONADA

O aprendizado não supervisionado permite que modelos de IA identifiquem padrões e estruturas em dados não rotulados, sendo usado em tarefas como agrupamento e detecção de anomalias.

XAI (IA EXPLICÁVEL)

A IA explicável é desenvolvida para oferecer decisões transparentes e compreensíveis, promovendo confiança e responsabilidade em setores críticos como saúde e finanças.

APRENDIZAGEM ZERO-SHOT (ZSL)

O aprendizado zero-shot permite que modelos de IA apliquem conhecimentos prévios para realizar tarefas inéditas sem treinamento específico, sendo usado em tradução automática e reconhecimento de imagens.



Modelos, Técnicas e Pesquisa

IA 2030 -

Os modelos de IA exigem grandes volumes de dados e recursos de computação para desbloquear seu potencial transformador (ou assim pensávamos).

Expansão das modalidades de IA generativa

Modelos

IA multimodal amadureceu e começou a transformar aplicações reais, permitindo que sistemas como o GPT-4o e Claude 3 integrem texto, imagem e som para resolver tarefas complexas com fluidez e precisão — marcando uma nova era de interação conversacional entre humanos e máquinas.

Termos

IA multimodal é um tipo de inteligência artificial que **entende e combina diferentes tipos de informação ao mesmo tempo**, como texto, imagem, áudio e vídeo.

Por exemplo: você pode **mostrar uma foto, falar algo e escrever uma pergunta**, e a IA entende tudo junto para te dar uma resposta mais completa — como se fosse uma pessoa ouvindo, lendo e vendo ao mesmo tempo.

É o que permite, por exemplo:

- Tirar foto de uma receita e pedir para ajustar as porções.
- Mostrar uma imagem médica e pedir uma explicação.
- Falar com um assistente virtual usando voz e imagem.

Ela está tornando a interação com máquinas **mais natural e inteligente**, quase como conversar com alguém que entende vários sentidos ao mesmo tempo.

Afinação

Modelos

O ajuste fino de modelos de linguagem se tornou uma ferramenta estratégica, permitindo personalização precisa, redução de custos e democratização do acesso a LLMs, com aplicações que vão da biotecnologia à conformidade regulatória em setores como saúde e finanças.

Termos

Ajuste fino de IA é como **ensinar uma IA que já sabe muito a se especializar em algo específico**.

Imagine que a IA é como um aluno que já leu milhares de livros sobre tudo (como medicina, direito, culinária, tecnologia...). Com o ajuste fino, você pega esse aluno e o treina só com **os assuntos que interessam pra sua empresa, setor ou necessidade**.

Por exemplo:

- Um hospital pode ajustar a IA para entender termos médicos e ajudar em diagnósticos.
- Uma empresa pode treiná-la com seus produtos e jeitos de falar, pra melhorar o atendimento ao cliente.

Isso deixa a IA **mais esperta, personalizada e útil para tarefas específicas**, sem precisar criar tudo do zero — economizando tempo e dinheiro. e regulatória em setores como saúde e finanças.

Aprendizagem de reforço automatizada

Modelos

O DeepSeek demonstrou que é possível treinar IA de alto desempenho com **feedback automatizado de computador**, reduzindo custos e tempo em tarefas com respostas objetivas, como matemática e codificação, embora o feedback humano ainda seja essencial em atividades mais subjetivas.

Termos

Aprendizagem de reforço automatizada é uma forma de ensinar a inteligência artificial (IA) **por tentativa e erro, com a ajuda de um computador, sem precisar de humanos o tempo todo.**

Pense assim: a IA faz uma tarefa (como resolver um problema ou jogar um jogo), e o computador dá uma “recompensa” ou uma “punição” dependendo se a resposta foi certa ou errada. Com isso, **a IA vai aprendendo sozinha o melhor caminho**, como se estivesse jogando um jogo em que precisa descobrir a melhor estratégia.

Antes, esse tipo de aprendizado dependia de pessoas avaliando as respostas da IA. Agora, com a **aprendizagem de reforço automatizada**, o próprio sistema avalia e ensina, o que **economiza tempo, dinheiro e acelera o aprendizado**, especialmente em tarefas com respostas claras, como matemática ou programação

Composição Evolutiva

Técnica

A Sakana AI propõe um novo caminho para IA avançada ao combinar modelos existentes com algoritmos evolucionários, evitando a necessidade de treinar sistemas massivos do zero. Essa abordagem torna o desenvolvimento de IA mais acessível, eficiente e sustentável, inspirando-se na complexidade da natureza.

Termos

Composição Evolutiva de IA é como **montar um supertime de diferentes inteligências artificiais já existentes**, combinando as melhores partes de cada uma — e deixando que o próprio sistema descubra, sozinho, **qual é a melhor forma de juntá-las**, usando ideias inspiradas na evolução da natureza.

Em vez de: Treinar um modelo gigante do zero (caro e demorado) Fazemos a IA aprender **como combinar modelos menores que já sabem fazer coisas diferentes** — como entender textos, ver imagens ou resolver contas

Exemplo prático (bem simples):

Imagine que você tem:

- Uma IA que entende **japonês perfeitamente**
- Outra que é ótima em **resolver problemas de matemática**
- E outra que entende **cultura e expressões locais**

A **composição evolutiva** tenta várias combinações dessas IAs, como se fosse testando diferentes receitas, até descobrir uma que junte o melhor de todas — criando um novo modelo que **fala japonês, resolve contas e entende o contexto cultural** melhor do que qualquer uma sozinha.

Pense nisso como montar um time ideal: você não precisa que cada jogador seja perfeito em tudo, mas sim combinar habilidades diferentes para um desempenho excelente em conjunto!

Mistura de especialistas

Técnica

A técnica de Mistura de Especialistas (MoE) divide um modelo em submodelos especializados, com um mecanismo que direciona cada tarefa ao “especialista” mais adequado, tornando a IA mais eficiente. A DeepSeek usou MoE para criar um sistema tipo ChatGPT a um custo bem menor, chamando atenção para essa abordagem. MoE permite escalar IA com desempenho elevado e menos exigência de hardware caro.

Termos

É como montar **uma equipe de especialistas** dentro de um único sistema de IA. Cada parte do sistema é treinada para ser muito boa em uma coisa específica — como matemática, escrita ou análise de imagem — e tem um **“gerente” que escolhe quem faz o quê** dependendo da tarefa recebida.

Imagine que você tem uma escola virtual com:

- Um professor de **matemática**
- Uma professora de **português**
- Um instrutor de **arte**
- E um coordenador que **sabe qual professor deve responder cada pergunta**.

Se um aluno perguntar:

- “Quanto é 12×9 ?” → o coordenador chama o **professor de matemática**.
- “Corrija meu texto de redação.” → ele chama a **professora de português**.
- “Me ajude a desenhar um logotipo.” → ele chama o **instrutor de arte**.

Esse é o jeito que a **IA MoE trabalha**:

Ela **não tenta fazer tudo com um único modelo gigante**, mas **escolhe o melhor “especialista” para cada tarefa** — o que economiza energia, melhora a resposta e reduz custos.

Autonomia dos Especialistas

Pesquisa

Pesquisadores propuseram uma evolução do MoE chamada **Autonomia de Especialistas (AoE)**, onde os próprios especialistas decidem se devem responder a uma entrada. Diferente do MoE tradicional, essa abordagem **dispensa o roteador central** e usa os sinais internos dos especialistas para guiar a escolha. Isso torna o sistema mais ágil e potencialmente mais eficiente.

Termos

É uma nova forma de organizar a inteligência artificial como **uma equipe de especialistas que se autogerencia**, sem precisar de um “chefe” para dizer quem faz o quê.

Como assim?

Em vez de ter um coordenador (como no MoE) dizendo:

“Você responde essa pergunta, você faz aquela”,

no **AoE, cada especialista lê a pergunta e decide sozinho se pode responder ou não.**

Se a IA de matemática sentir que a pergunta é sobre números, ela mesma levanta a mão e diz “essa é comigo!”. Se não for, ela se cala — simples assim

Exemplo prático:

Imagine que você tem uma IA com vários especialistas:

Um entende de **comida**, outro de **música**, outro de **viagens**
Você pergunta:

“Qual a melhor época para visitar o Japão?”

Cada especialista “escaneia” a pergunta.

A IA de comida pensa: “Não é comigo.”

A IA de música ignora.

Mas a IA de viagens diz: “Sim! Isso eu sei responder.” Ela responde diretamente, **sem precisar que um coordenador escolha por ela.**

LLMs como sistemas operacionais

Pesquisa

Um novo tipo de sistema operacional, como o AIOS, usa um modelo de linguagem grande (LLM) como seu núcleo, permitindo que os usuários interajam por linguagem natural em vez de comandos técnicos. Esse sistema automatiza tarefas com inteligência, tornando o uso do computador mais intuitivo e eficiente. O LLM também gerencia decisões, recursos e agentes simultâneos, criando uma experiência computacional mais adaptável.

Termos

Tradicionalmente, um **sistema operacional** (como Windows, macOS ou Linux) serve para **controlar tudo o que acontece no seu computador** — arquivos, programas, memória, etc. Mas você precisa clicar, arrastar, abrir janelas ou digitar comandos para fazer o que quer.

Agora imagine que no lugar disso, você **simplesmente conversa com o computador**, como se estivesse falando com um assistente muito inteligente.

Isso é possível com os **LLMs (Modelos de Linguagem de Grande Escala)** sendo o “cérebro” principal do sistema operacional.

Exemplo prático:

Você liga seu computador e diz:

“Organize todos os meus arquivos por tipo e envie os PDFs importantes para o e-mail do João.”

→ O sistema entende, localiza os arquivos certos, organiza as pastas, envia o e-mail e te avisa quando terminar — **sem você clicar em nada**.

Ou ainda:

“Me avise quando tiver uma reunião e prepare uma apresentação com os dados do último relatório.”

→ O computador busca os dados, cria os slides e te avisa na hora da reunião.

LLMs: Maiores e mais caros

Pesquisa

Os LLMs cresceram com base na ideia de que “maior é melhor”, trazendo avanços, mas também custos altíssimos e desafios de eficiência. Estudos recentes mostram que o aumento do tamanho nem sempre gera ganhos proporcionais em desempenho. Agora, o foco está mudando para modelos menores, mais eficientes e especializados, adaptados às necessidades reais das empresas.

Termos

Durante anos, acreditava-se que **quanto maior o modelo de IA, melhor o desempenho**. E de fato, modelos como o **GPT-4** (da OpenAI) e o **Gemini** (do Google) ficaram incríveis — mas também **caríssimos de treinar**. Só o GPT-4 custou mais de **78 milhões de dólares!** E isso exige **muito hardware, energia e dinheiro**.

Mas agora os pesquisadores estão percebendo que:

- 👉 Nem toda tarefa precisa de um modelo gigante.
- 👉 Às vezes, **modelos menores e bem ajustados funcionam melhor e custam menos**.
- 👉 Modelos enormes podem ter “retorno decrescente” — ou seja, crescem, mas não melhoram tanto assim.



Exemplo prático (vida real):

Imagine uma empresa que precisa de uma IA só para **responder dúvidas sobre seus produtos**.

Ela não precisa de um modelo do tamanho do GPT-4, que sabe falar de tudo do mundo.

Ela pode usar um **modelo menor**, treinado apenas com os dados da empresa (ajuste fino), que será:

- Mais barato,
- Mais rápido,
- E até **mais preciso**, porque só foca no que é importante pra ela.

Modelos de cadeia de pensamento

Modelos

Com os limites do crescimento de modelos gigantes, a atenção se volta para abordagens como a Cadeia de Pensamento (CoT), que foca em raciocínio estruturado em tempo real. Modelos como o1 e R1 mostram que gerar etapas lógicas internas melhora a precisão, mesmo que exija mais computação durante o uso. Isso marca uma mudança do foco no pré-treinamento para a eficiência e inteligência na hora da inferência.

Termos

Antes, a ideia era simples:

“Quanto maior o modelo de IA, melhor.”

Ou seja, mais dados, mais memória, mais poder de processamento = melhores respostas.

Mas isso ficou **caro, demorado e com pouco ganho em alguns casos**.

Agora surgiu uma nova abordagem chamada **Cadeia de Pensamento (Chain of Thought, ou CoT)**.

Ela funciona assim: em vez de tentar acertar a resposta direto, a IA **pensa em etapas**, como se estivesse **resolvendo um problema passo a passo em voz alta**, igual a um ser humano.



Exemplo prático (bem simples):

Suponha que você pergunte para a IA:

“Se João tem 12 maçãs e dá 4 para Maria, quantas ele tem?”

- Um modelo antigo tentaria responder direto: **“8”**
- Um modelo com **Cadeia de Pensamento (CoT)** faz assim:



“João começou com 12 maçãs. Ele deu 4 para Maria. 12 menos 4 é 8. Então ele ficou com 8.”



Ela pensa antes de responder, checka o raciocínio, e por isso, erra menos — principalmente em perguntas difíceis ou confusas.

Pequenos modelos de linguagem

Modelos

Pequenos modelos de linguagem (SLMs) estão mostrando que é possível obter alto desempenho com muito menos recursos, superando modelos maiores em tarefas específicas com mais eficiência e menor custo. Exemplos como o GPT-4o mini, Phi-3-mini e OpenELM provam que modelos compactos podem operar em dispositivos locais, garantindo privacidade e rapidez. A tendência aponta para uma IA mais acessível, personalizada e modular, com múltiplos modelos especializados atuando em conjunto.

Por muito tempo, acreditava-se que **só modelos gigantes de IA** (como o ChatGPT-4) eram bons.

Mas agora, os **modelos pequenos de linguagem (SLMs)** estão provando que podem fazer **coisas muito bem — e gastando muito menos!**

Por que isso importa?

1. **Mais barato** – As empresas não precisam investir em supercomputadores.
2. **Mais rápido** – Respostas imediatas, direto no seu dispositivo.
3. **Mais seguro** – Seus dados não saem do aparelho (privacidade!)
4. **Mais personalizado** – A IA pode ser treinada para **tarefas específicas**, como agricultura, manufatura ou saúde, sem precisar ser um modelo gigante que sabe “tudo do mundo



Exemplo prático (vida real):

Imagine que você tem um **smartphone**.

Antes, para usar IA poderosa, seu celular precisava **mandar tudo para a nuvem** (internet), esperar o processamento, e então receber a resposta. Isso gasta tempo, internet e energia.

Agora, com **SLMs**, a IA pode rodar **direto no celular**, sem depender da nuvem!

Você pode:

- Traduzir textos,
- Organizar sua agenda,
- Criar respostas automáticas para e-mails



Tudo com rapidez, privacidade e sem esgotar a bateria.

Aterramento e Aumento de Contexto

Modelos

O NotebookLM do Google usa IA para oferecer respostas personalizadas e confiáveis, conectando-se diretamente aos seus documentos e fontes. Com técnicas como aterramento e aumento contextual, ele reduz erros e melhora a relevância das respostas. Essa abordagem está transformando a IA em uma ferramenta precisa e adaptável para uso pessoal e empresarial.

Termos

O que é o NotebookLM do Google?

É uma ferramenta de IA que funciona como **um assistente de estudos ou de trabalho**, só que **inteligente e personalizado**.

Diferente de chatbots genéricos, o NotebookLM **lê seus documentos do Google Docs**, entende o conteúdo e **responde perguntas com base no que você escreveu**.

Por que isso é útil?

1. **Evita informações erradas (alucinações)** – A IA **só responde com base no que você já escreveu ou carregou**.
2. **Poupa tempo** – Em vez de procurar manualmente nos documentos, a IA **sintetiza tudo pra você**.
3. **Ajuda a conectar ideias** – Ele cruza dados de diferentes fontes para **te dar insights mais inteligentes**.

Exemplo prático (vida real):

Imagine que você está preparando um **TCC, relatório ou planejamento estratégico**.

Você já tem um monte de anotações espalhadas em documentos diferentes.

Com o **NotebookLM**, você pode perguntar:

“*Quais foram os principais pontos da reunião de abril sobre marketing digital?*”

➡ E ele vai **buscar essa informação diretamente nas suas anotações**, e ainda dizer de qual documento tirou a resposta.

Superando a escassez de dados

Modelos

A escassez de dados textuais de alta qualidade ameaça o avanço de grandes modelos de IA, com previsões de esgotamento já em 2026. Para contornar isso, laboratórios estão recorrendo a dados privados, visuais, sintéticos e estratégias como *self-play*, onde modelos treinam contra si mesmos. No entanto, depender demais de dados sintéticos pode prejudicar a diversidade e qualidade dos resultados, exigindo abordagens inovadoras para manter o progresso

Termos

O que está acontecendo com os dados usados para treinar IA?

Para treinar uma **inteligência artificial poderosa**, é preciso dar a ela **muito conteúdo de qualidade** — textos, imagens, vídeos... como se fosse a “alimentação” da IA.

Só que temos um problema:

A internet pública já está ficando sem dados bons para treinar esses modelos — ou seja, os sites com textos bem escritos e confiáveis estão acabando para esse uso.

Estudos dizem que até **2026**, os dados de boa qualidade da internet podem se esgotar.

Os laboratórios estão **buscando novas formas de ensinar IA**. **Comprar dados privados**, como de jornais ou plataformas educacionais. **Usar dados visuais e de áudio**, como vídeos de YouTube ou gravações para ensinar sobre o mundo real.



Exemplo prático:

Imagine que você quer ensinar alguém a cozinhar.

Nos primeiros dias, você dá ótimos livros e vídeos de chefs renomados.

Depois de um tempo, só sobram receitas mal escritas, vídeos de baixa qualidade ou até repetidos.



A pessoa aprende menos, ou até aprende errado.

Com a IA é igual: **sem conteúdo bom, ela para de evoluir — ou pior, começa a “pirar” nas respostas.**

IA de código aberto

Modelos

O lançamento do modelo R1 pela DeepSeek e o avanço de outros modelos abertos como Llama e Mistral mostram que a IA de código aberto está alcançando e até superando soluções proprietárias, com menor custo. Essa tendência é reforçada pelo crescimento explosivo do desenvolvimento colaborativo em plataformas como o GitHub. Paralelamente, a IA modular vem ganhando destaque ao permitir a criação de sistemas mais eficientes, personalizados e escaláveis, combinando módulos especializados para tarefas específicas.

Termos

O que está acontecendo com a IA de código aberto?

Antes, os modelos de IA mais avançados eram **fechados**, como o ChatGPT ou o Gemini do Google — ninguém podia ver como eram feitos nem usá-los livremente.

Mas agora, empresas como **DeepSeek**, **Meta (Facebook)** e **Mistral** estão lançando modelos **abertos e gratuitos**, que qualquer pessoa pode usar, estudar ou melhorar.

Exemplo:

- O modelo **R1 da DeepSeek** é superpoderoso e custa muito menos que modelos como o GPT-4.
- O **Llama da Meta** e o **Mixtral da Mistral** já estão sendo usados por empresas como o **Wells Fargo** e o navegador **Brave**.

**Exemplo prático:**

Imagine que você é dono de uma empresa e quer uma IA para ajudar seus clientes.

Antigamente, você teria que **pagar caro** para usar uma IA fechada.

Agora, com esses **modelos de código aberto**, você pode **baixar um modelo gratuito**, treinar com os dados da sua empresa e rodar localmente — **sem depender de grandes plataformas**.

Grandes modelos de ação

Pesquisa

Os Modelos de Ação Grande (LAMs) são IAs que não apenas conversam, mas também executam tarefas ativamente em computadores, como clicar, digitar e navegar. Eles representam um passo rumo a sistemas autônomos completos, capazes de agir como usuários humanos em ambientes digitais. Embora promissores, levantam preocupações sobre privacidade, já que precisam de amplo acesso aos dados e sistemas dos usuários.

Termos

O que são LAMs (Modelos de Ação Grande)?

Até agora, a maioria das inteligências artificiais (como o ChatGPT) só **responde com texto** — ou seja, **ela conversa, mas não faz nada no seu computador**.

Os **LAMs** são uma nova geração de IA que **não só falam, mas também agem**.

Eles **clicam, digitam, navegam, arrastam arquivos, preenchem formulários** — como se fossem um humano usando o computador por você.

⚠ **Mas tem um porém...**

Para funcionar assim, a IA precisa **ver sua tela, acessar seus arquivos e comandos**. Isso levanta dúvidas sobre **privacidade e segurança**, pois ela teria acesso a tudo que você faz ou guarda no computador.



Exemplo real:

O **Claude 3.5 Sonnet**, da Anthropic, já consegue:

- Mover o cursor do mouse
- Clicar em botões
- Ler capturas de tela
- Escrever textos
- Navegar em sistemas de computador

Modelos de Ação Grande Pessoal

Pesquisa

Os PLAMs (Modelos de Grande Ação Pessoal) são IAs altamente personalizadas, treinadas com todos os dados da vida digital de uma pessoa, como redes sociais, localização, e-mails e hábitos. Eles seriam capazes de tomar decisões alinhadas ao comportamento e preferências individuais, atuando como uma verdadeira extensão digital do usuário. Embora ainda teóricos, esses modelos prometem transformar a relação entre humanos e IA, criando forte dependência das plataformas que os oferecem.

Termos

O que são PLAMs (Modelos de Grande Ação Pessoal)?

Os **PLAMs** são um tipo de inteligência artificial ainda mais avançada e **muito mais pessoal**.

Diferente de outras IAs que sabem coisas genéricas (como o ChatGPT), um PLAM seria treinado **exclusivamente com os dados da sua vida**.

Ele saberia tudo sobre:

- Suas mensagens, e-mails, calendário 📅
- Seu histórico de compras 🛒
- Onde você está e com quem está 📍
- Seus horários, preferências, rotina e até batimentos cardíacos, se você usar smartwatch ❤️⌚
É como se fosse **um clone digital seu**, que entende você melhor do que qualquer outro sistema.

💡 Exemplo prático:

Imagine que você está saindo do trabalho.

Seu PLAM já sabe:

- Que está chovendo 🌧️
- Que seu carro está com pouco combustível ⛽
- Que você jantou tarde ontem 🍽️
- E que tem uma reunião cedo amanhã 📊

Ele **automaticamente agenda um delivery leve**, ajusta seu despertador e **programa o GPS para abastecer no caminho mais curto para casa** — tudo sem você pedir.



Association de Aquilnestacion end
Lege Compride Orghtation

PRÓXIMA AULA-ISO 42.001 LIDERANÇA



Association for Alier Emication
and Logs, mdutenfor Degistration

PRÓXIMA AULA-IA 2030 SEGURANÇA, ÉTICA E SOCIEDADE

