Inteligência Artificial para a Língua Portuguesa e para a Transformação Educativa

Palestra plenária convidada no

Seminário "Educação Inteligente: a Escola na Era da IA", organizado por

Conselho Nacional de Educação, em Lisboa, a 15 de outubro de 2025

António Branco

Diretor Geral

PORTULAN CLARIN – Infraestrutura Nacional de Investigação para a Ciência e Tecnologia da Linguagem

Professor

Universidade de Lisboa, Faculdade de Ciências

META-NET

THE A

Portuguese língua Language in Portuguesa

THE DIGITAL NA ERA

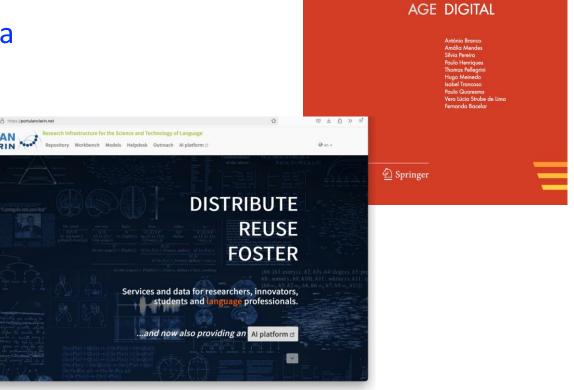
IA & Tecnologia da língua portuguesa

25 anos de I&D

Universidade de Lisboa

Faculdade de Ciências, Dep. de Informática

NLX Grupo de Fala e Linguagem Natural



Logo 2 meses após irrupção do ChatGPT:



OPINIÃO

EDITORIAIS

BARTOON

MAIS V

Público • Quinta-feira, 9 de Fevereiro de 2023 • 9

Língua portuguesa e tecnologia: o futuro é agora



Ana Paula Laborinho e António Branco

s avanços da Inteligência
Artificial têm sido
impressionantes, sobretudo na
sua aplicação à Tecnologia da
Língua. Este progresso é
baseado na aprendizagem
automática com os chamados Grandes
Modelos de Linguagem, como o GPT-3 ou o

um nível inédito de qualidade, como por exemplo tradução, conversação, transcrição de fala e legendagem, geração de texto e fala, análise do conteúdo e extração de informação, etc. Quando integrados em sistemas mais vastos, estão a transformar os diagnósticos e cuidados de saúde, os serviços financeiros e jurídicos, os jogos e o entretenimento, o ensino, a criatividade e a cultura, etc.

Devido ao tamanho dos modelos, estas tarefas de processamento estão disponíveis remotamente como serviços *online*, como é o caso dos motores de busca, e não como os corretores ortográficos de instalação local nos nossos dispositivos. Devido à dimensão dos recursos para a aprendizagem, no imediato esses serviços são disponibilizados pelo

humanas, em particular, afunilada num pequeno oligopólio mundial, gera riscos alarmantes para as soberanias individuais e coletivas.

Impactos indesejáveis de tecnologias emergentes mitigam-se com mais e melhor tecnologia, não com menos. A dispersão do fornecimento destes serviços é crucial para debelar a ameaça que a sua concentração constitui. A resposta encontra-se assim na promoção de um ecossistema de inovação que, em alternativa, permita atempadamente banalizar o acesso aos recursos necessários para a apropriação e exploração da Tecnologia da Linguagem pelo maior número possível de indivíduos e organizações, privadas e públicas, pequenas e grandes, nacionais e internacionais.

e disponibilizar para propiciar tal ecossistema, e perante o mais relevante interesse público em causa, esta é uma incumbência, nova e urgente, para os Estados democráticos, isoladamente e em cooperação.

A língua portuguesa, com 250 milhões de falantes em quatro continentes, é uma das grandes línguas internacionais de projeção global. Os indicadores apontam para o seu crescimento até final do século com a maioria dos falantes no continente africano. Contudo, se não acrescentarmos às políticas de língua clássicas uma aposta clara na sua preparação tecnológica, perderá importância e tenderá no limite a ser substituída por outras línguas. Por essa razão, importa congregar esforços para que haja um Plano de Preparação

Rationale 1/2

- Ameaça: intermediação tecnológica pervasiva
 - linguagem natural, ingrediente básico da atividade social e individual
 - utilização da linguagem dependente de pequeno oligopólio de fornecedores de serviços de IA concentrados num país
 - comunicação pode ser interrompida de forma generalizada

- Necessidade: mitigação da dependência tecnológica
 - salvaguarda da liberdade individual
 - salvaguarda da soberania linguística, cultural e coletiva

Rationale 2/2

- Resposta: LLMs abertos
 - LLMs requerem recursos massivos...
 - Vamos desenvolver LLMs abertos e e distribuí-los!
 - democratizar o fornecimento de serviços de IA
 - promover tantos fornecedores independentes quanto for possível

Soberania e liberdade *mesmo*



Associated Press, May 15, 2025

THE HAGUE, Netherlands (AP) — The International Criminal Court's chief prosecutor has lost access to his email, and his bank accounts have been frozen.

. . .

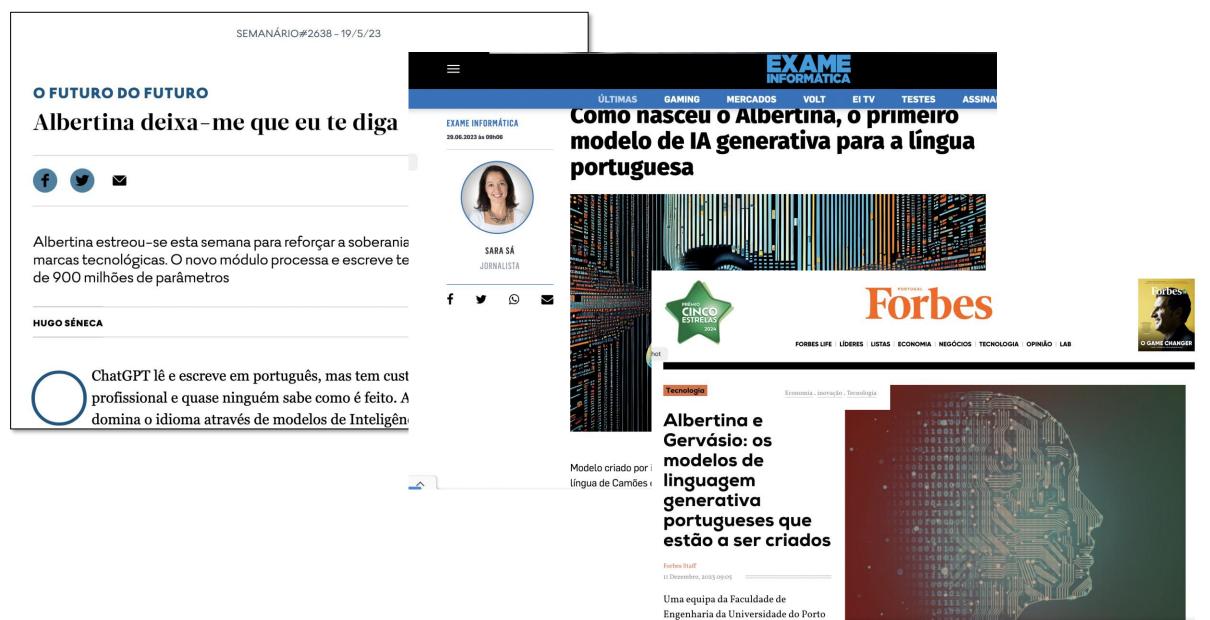
Trump sanctioned the court after a panel of ICC judges in November issued arrest warrants for Israeli Prime Minister Benjamin Netanyahu

. . .

Microsoft, for example, cancelled Khan's email address

LLMs abertos especializados na língua portuguesa precisam ser desenvolvidos e distribuídos: missão!

LLMs abertos para português





- Nov 2021 ChatGPT
- Abr 2022, 1ª publicação
- Set 2023, EPIA

• 100M, 900M, 1.5B



Advancing Neural Encoding of Portuguese with Transformer Albertina PT-*

João Rodrigues¹(⊠), Luís Gomes¹, João Silva¹, António Branco¹, Rodrigo Santos¹, Henrique Lopes Cardoso², and Tomás Osório²

NLX—Natural Language and Speech Group, University of Lisbon, Faculty of Sciences (FCUL), Dept. Informatics, Campo Grande, 1749-016 Lisboa, Portugal {jarodrigues,luis.gomes,jsilva,antonio.branco,rsdsantos}@fc.ul.pt
Laboratório de Inteligência Artificial e Ciência de Computadores (LIACC)
Faculdade de Engenharia da Universidade do Porto (FEUP), Rua Dr. Roberto Frias, 4200-465 Porto, Portugal

hlc@fe.up.pt tomas.s.osorio@gmail.com

Abstract. To advance the neural encoding of Portuguese (PT), and a fortiori the technological preparation of this language for the digital age, we developed a Transformer-based foundation model that sets a new state of the art in this respect for two of its variants, namely Euro-



- Nov 2021 ChatGPT
- Abr 2022, Albertina
- Jul 2022, 1ª publicação
- Mai 2024, LREC

Advancing Generative AI for Portuguese with Open Decoder Gervásio PT*

Rodrigo Santos, João Silva, Luís Gomes, João Rodrigues, António Branco University of Lisbon

NLX - Natural Language and Speech Group, Department of Informatics Faculdade de Ciências, Campo Grande, 1749-016 Lisboa, Portugal {rsdsantos, jrsilva, luis.gomes, jarodrigues, antonio.branco}@fc.ul.pt

Abstract

To advance the neural decoding of Portuguese, in this paper we present a fully open Transformer-based, instruction-tuned decoder model that sets a new state of the art in this respect. To develop this decoder, which we named Gervásio PT*, a strong LLaMA 2 7B model was used as a starting point, and its further improvement through additional training was done over language resources that include new instruction data sets of Portuguese prepared for

• 1B, 8B, 70B



- Nov 2021 ChatGPT
- Abr 2022, Albertina
- Jul 2022, Gervásio
- Abr 2024, 1º publicação
- Mai 2024, LREC

• 100M, 335M, 900M



Open Sentence Embeddings for Portuguese with the Serafim PT* Encoders Family

Luís Gomes^(⊠), António Branco, João Silva, João Rodrigues, and Rodrigo Santos

NLX—Natural Language and Speech Group, Department of Informatics Faculdade de Ciências, Campo Grande, University of Lisbon, 1749-016 Lisboa, Portugal {luis.gomes,antonio.branco,jrsilva,jarodrigues,rsdsantos}@fc.ul.pt

Abstract. Sentence encoder encode the semantics of their input, enabling key downstream applications such as classification, clustering, or retrieval. In this paper, we present Serafim, a family of open-source

Pioneiros, triplo A, estado da arte

• 3 famílias principais de LLMs

- encoders: Albertinas tarefas de classificação
- decoders: Gervásios tarefas de geração
- embedders: Serafins tarefas de sinonímia

2 variantes

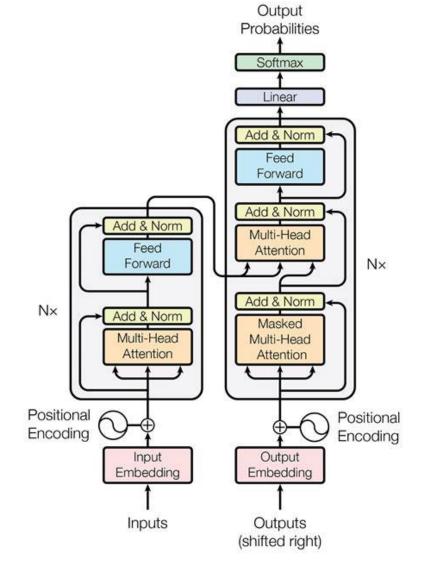
PTPT e PTBR

3x abertos

pesos, licença, distribuição

• Desempenho de topo resp. categorias

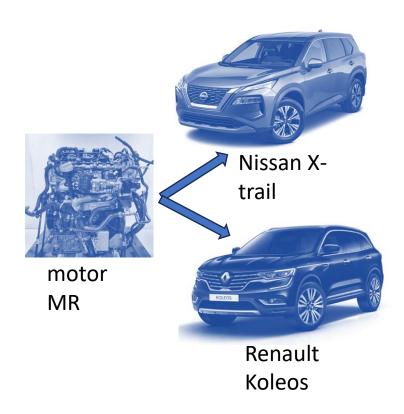
novos/primeiros benchmarks para PT



LLMs abertos especializados na língua portuguesa desenvolvidos (e em desenvolvimento) e distribuídos: e agora ?

motor x modelos

LLMs x Agentes IA







🔀 Evaristo.ai × + ○ A https://evaristo.ai Evaristo.ai BETA plataforma de IA aberta para a língua portuguesa 🔀 Evaristo.ai Modelo ativo 0 Gervásio 8B PTPT chatbot multi-modelo Esta semana e multi-heterónimo Sem 62.6% Cul 17.0% Civ 83.6% Var PTPT Tutebol: Portugal x Espanha Dados de treino até: dez 2023 Website deste modelo ₹ Jacaranda em flor Escreve fado sobre jacarandás Dá-me uma receita com bacalhau Planeia visita a Portugal Aiuda-me a aprender português Ensina-me a teoria heliocêntrica Traduz texto Tempo de qualidade com os Diz-me como pedir aumento Resume texto miúdos Ó Evaristo, tens cá disto? Último jogo Portugal-Espanha? Quem é o atual presidente? Modelos Heterónimos Extensões Pesquisa na internet Definições Modelo: Gervásio 8B PTPT * O resultado gerado pode não ser correto, exato, perfeito e/ou adequado, qualquer que seja a finalidade que o utilizador lhe atribua. · O 6 início contacto imprensa acerca condições privacidade agradecimentos tens cá disto? disponibilização portulan disseminação ousar.ai

junho 2025

- multi-modelo
- multi-heterónimo
- multi-modal
- multilíngue
- agentivo multiextensão
- gratuito, sem registo
- onprem, salvaguarda privacidade
- corre em LLMs abertos triplo A

remplos		
Escreve fado sobre jacarandás	Dá-me uma receita com bacalhau	Planeia visita a Portugal
Ajuda-me a aprender português	Ensina-me a teoria heliocêntrica	Traduz texto
Tempo de qualidade com os miúdos	Diz-me como pedir aumento	Resume texto
Ó Evaristo, tens cá disto?	Quem é o atual presidente?	Último jogo Portugal-Espanha?

Eu sou aluno do ensino secundário. Estou a estudar a teoria heliocêntrica. Explica-ma por favor. Inclui contextualização histórica.



A teoria heliocêntrica é um modelo astronómico que coloca o Sol no centro do sistema solar, com os planetas, incluindo a Terra, orbitando em torno dele. Esta teoria foi proposta por Nicolás Copérnico e melhorada posteriormente por Johannes Kepler ¹. A adopção deste modelo representou um cambio drástico na forma como os cientistas entendiam o movimento dos corpos celestes e a estrutura do universo.

A ideia de que a Terra gira em torno ao Sol foi proposta desde o século III a.C. por Aristarco de Samos, mas não recebeu apoio de outros astrónomos da antiguidade ¹. No entanto, no Renascimento, durante o século XVI, a revolução copernicana começou a tomar forma.

O modelo heliocêntrico é diferente do geocêntrico, que colocava a Terra como centro do universo e todos os corpos celestes girando em torno dela ². A aceitação da teoria heliocêntrica permitiu uma reavaliação das crenças sobre a natureza do universo.

a maneira como o universo era percebido e

3. Além disso, inspirou um maior interesse pela ípios que ainda são utilizados nos modelos atua

ender o movimento dos corpos celestes e a significativo na ciência moderna e continua s atuais ⁴.

ca

l-astronomy

- https://www.space.com/copernicus-helio-centric-theory.html
- 4 https://www.nasa.gov/audience/forstudents/k-4/stories/nasa-knows/what-is-heliocentrism

Sources: w es.wikipedia.org

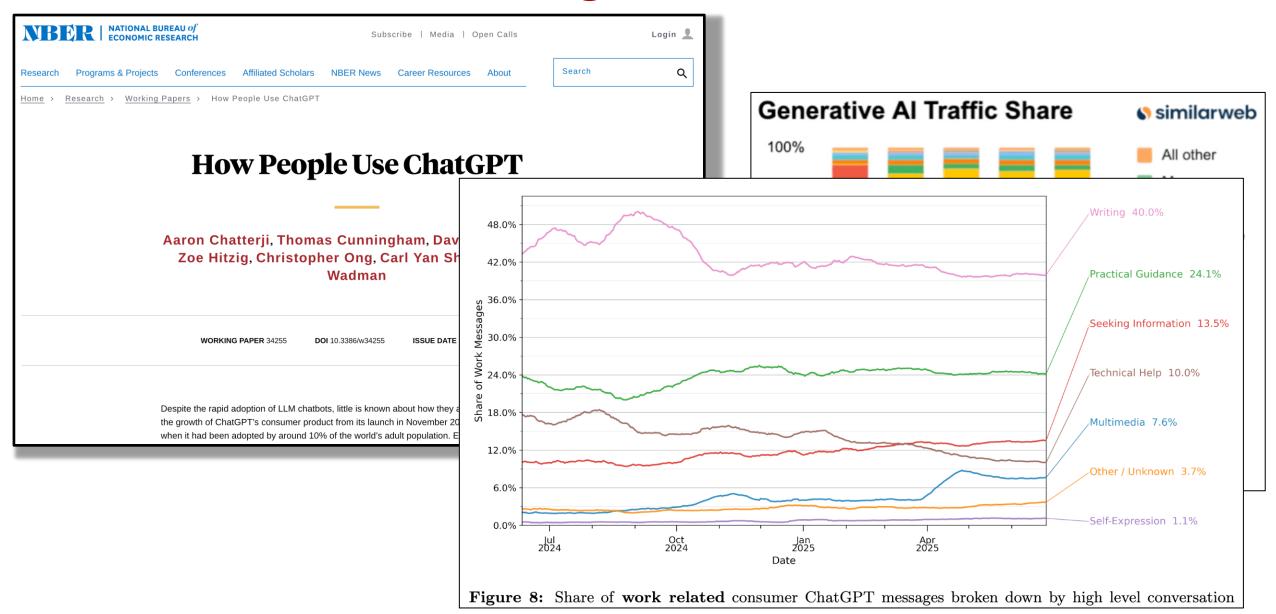




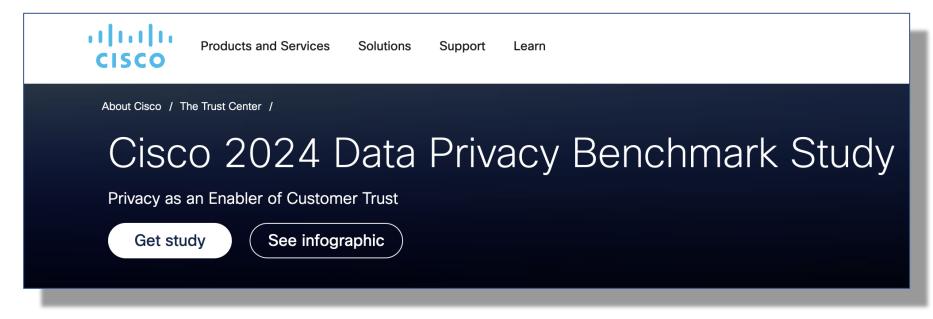
nuevaescuelamexicana.org

Em que medida estes avanços podem ser integrados na e contribuir para a educação?

Chatbot soberano seguro



Privacidade & ghost Al

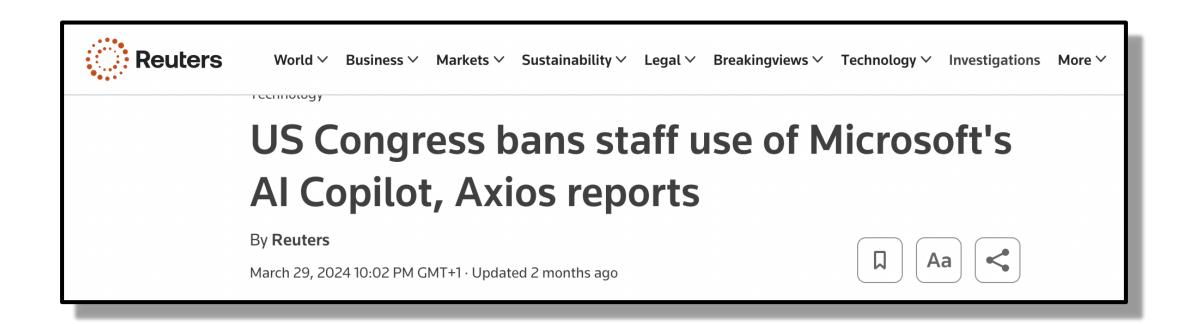


inquérito a 2600+ profissionais de ciber-segurança em 12 países [Jan 25, 2024]:

Austrália, Brasil, China, França, Alemanha, Índia, Itália, Japão, México, Espanha, Reino Unido e Estados Unidos

27% afirmaram que a sua organização proibiu totalmente as aplicações GenAl 63% estabeleceram limitações quanto aos dados que podem ser introduzidos

Segurança e soberania



Agentes IA especializados: internos, externos

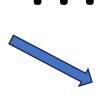






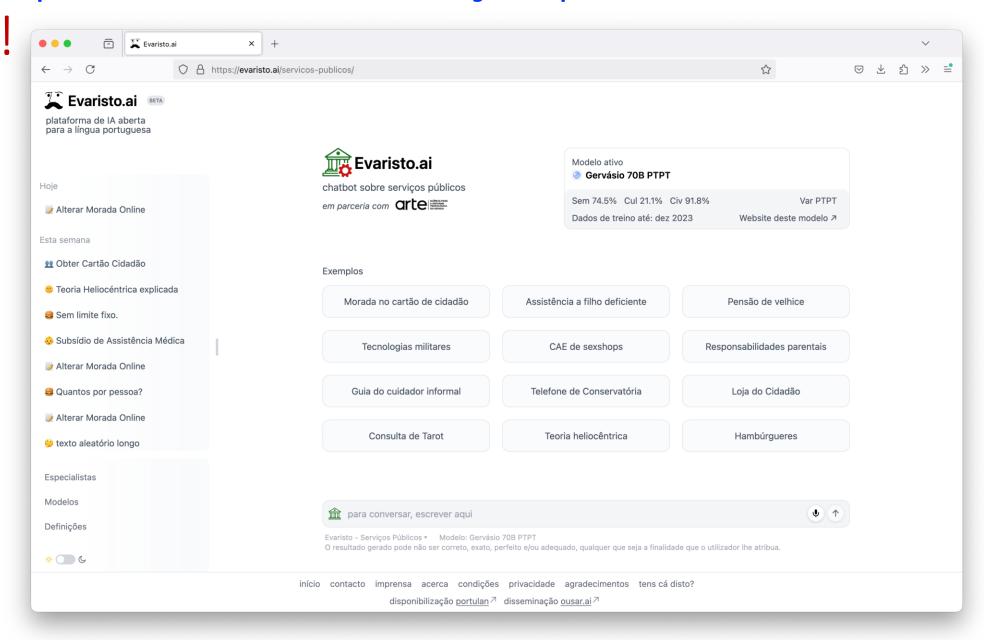






Evaristo – especialista em serviços públicos António Branco

antestreia!!

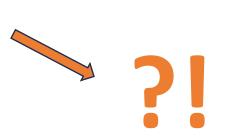


LLMs x Chatbots/Agentes IA









А

Ferramentas para tarefas inteligentes e criativas no processo educativo

B

Soberania/liberdade individual e coletiva

C

Ensino-aprendizagem da IA, ela própria nova área de saber crucial

D

Questões existenciais para a educação e para a humanidade

Obrigado

António Branco

Diretor Geral
PORTULAN CLARIN

Infraestrutura Nacional de Investigação para a Ciência e Tecnologia da Linguagem

Professor

Universidade de Lisboa, Faculdade de Ciências