



Re.Data

Rede para a Gestão de
Dados de Investigação

Gestão de Dados na Era da IA Generativa

Re.Data Lunch Webinar

António Luís Lopes (antonio.luis@iscte-iul.pt)

Iscte – Instituto Universitário de Lisboa



Re.Data
Rede para a Gestão de
Dados de Investigação



Universidade do Minho



UNIVERSIDADE DE
COIMBRA

iscte

INSTITUTO
UNIVERSITÁRIO
DE LISBOA



INSTITUTO DE INVESTIGAÇÃO EM
BIOINFORMÁTICA E BIOQUÍMICA



UNIVERSIDADE NOVA
DE LISBOA

Velocidade de Adopção da IA Generativa

ChatGPT reaches 100 million users two months after launch

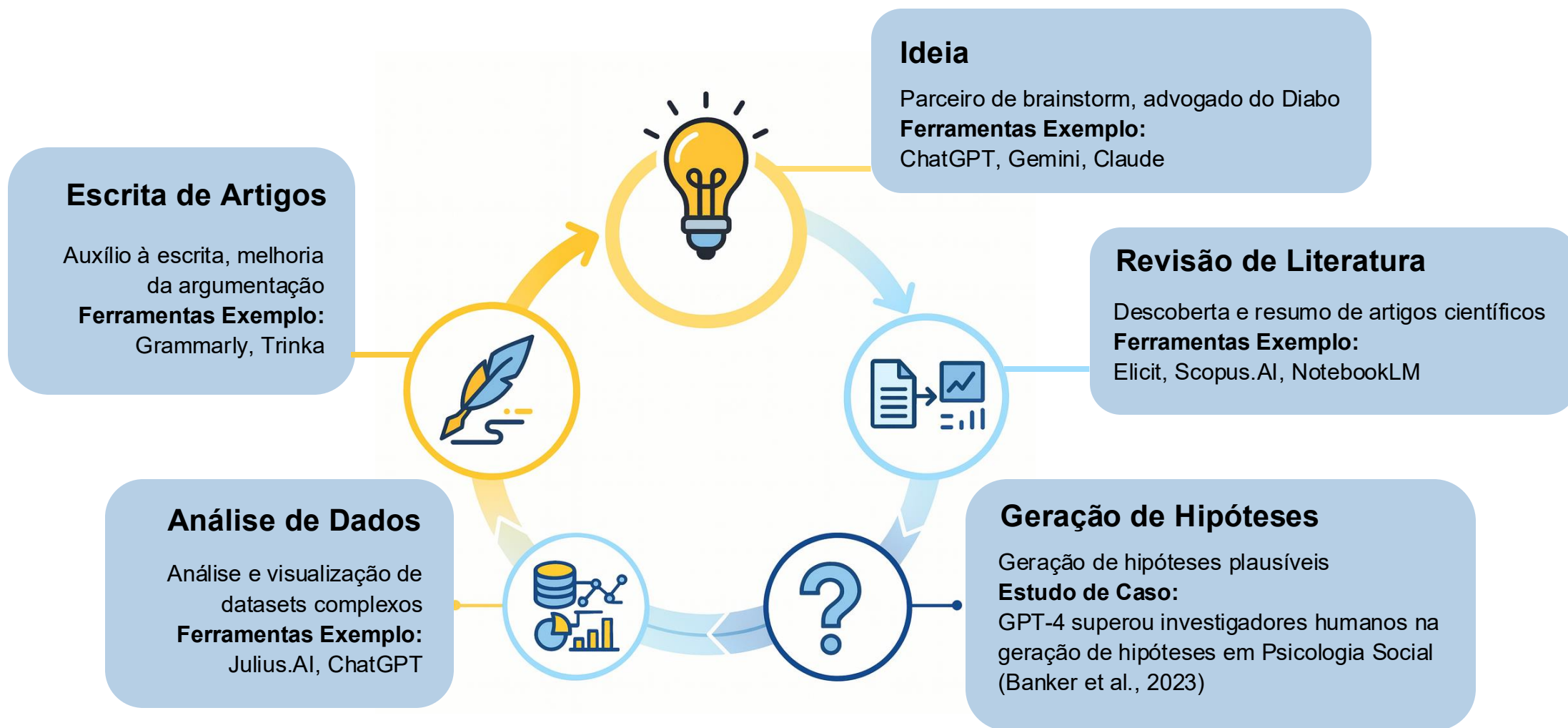
Unprecedented take-up may make AI chatbot the fastest-growing consumer internet app ever, analysts say



[The Guardian, 2023]

Tecnologia	Tempo -> 100M Utilizadores
Telefone	75 anos
Telemóvel	16 anos
Internet	7 anos
Twitter	5 anos
Facebook	4.5 anos
Instagram	2.5 anos
TikTok	9 meses
ChatGPT	2 meses

A IA generativa como parceira no processo científico



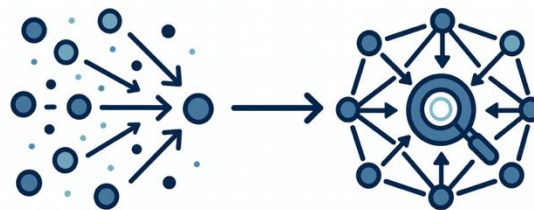
A IA a potenciar a infraestrutura de ciência



Anotação assistida por IA

Os modelos de linguagem podem analisar publicações e sugerir metadados relevantes, aumentando a eficiência e a precisão da anotação de publicações.

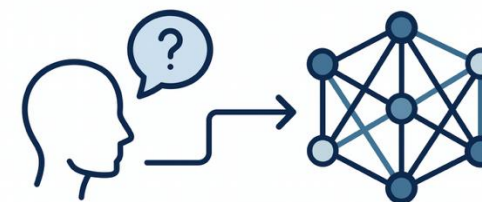
Exemplo: sistema fredato (IMBI)



Enriquecimento de knowledge graphs

Os algoritmos de machine learning podem enriquecer e interligar registos de investigação, extrair entidades, resolver ambiguidades e descobrir relações latentes.

Exemplo: OpenAIRE Graph



Exploração de dados estruturados

Os assistentes de IA e os chatbots podem traduzir linguagem natural em queries complexas, permitindo que se possa interagir com dados sem conhecimentos técnicos.

Exemplo: projeto GRAPHIA

A promessa da IA e os desafios

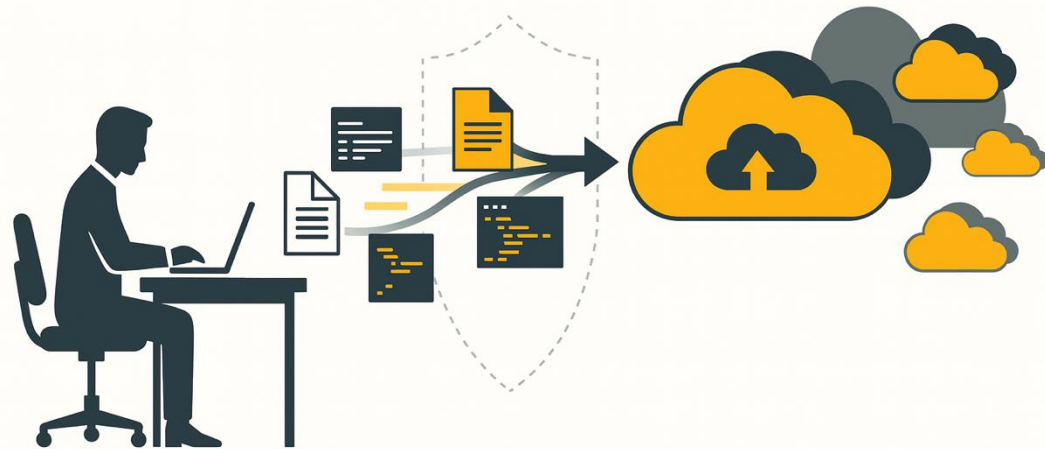
No entanto, esta poderosa capacidade de processar e gerar conhecimento não vem sem custos. A integração da IA nos fluxos de trabalho científico introduz um novo e complexo conjunto de desafios de segurança, legais e éticos que ameaçam os próprios fundamentos da ciência aberta.

Potencial problema de privacidade

Quando se usam estas ferramentas, quaisquer dados que são partilhados, passam a existir na infraestrutura das empresas que as criam.

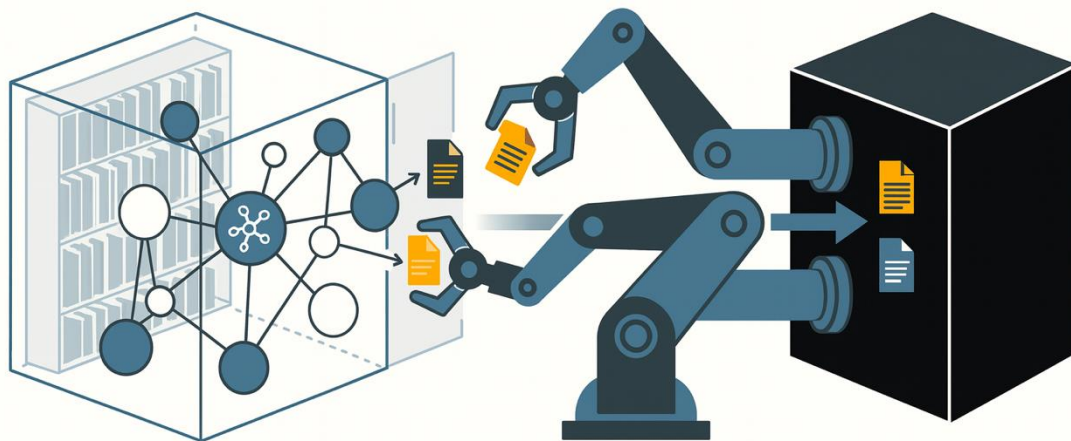
Esses dados podem ser usados para treino de futuros modelos de IA.

Se esses dados forem sensíveis ou privados, isto pode constituir um potencial problema de privacidade.



A promessa da IA e os desafios

Os repositórios de dados e outputs científicos foram concebidos para acesso humano, difusão de conhecimento e reutilização académica, mas depois acabam por ser vítimas dos bots que extraem quantidades enormes de dados sem qualquer consentimento.



Reutilização não recíproca

As empresas de IA usam estes dados para o treino de modelos. Não só isto cria uma sobrecarga na infraestrutura científica como o conhecimento é reutilizado sem crédito, contexto ou benefício para as comunidades que o criaram. O valor gerado reverte apenas a favor de sistemas proprietários e fechados.

A promessa da IA e os desafios

Os princípios fundamentais da proteção de dados, como a minimização, a limitação da finalidade e o consentimento, colidem frequentemente com o funcionamento dos modelos de linguagem, que podem reter dados dos utilizadores indefinidamente.



Conformidade improvável (ou impossível)

Como pode uma organização cumprir um pedido do “direito ao esquecimento” se os dados de um indivíduo já foram absorvidos por um modelo de linguagem?

A falta de transparência sobre como e onde os dados são armazenados pelos fornecedores de IA torna a demonstração de conformidade quase impossível.

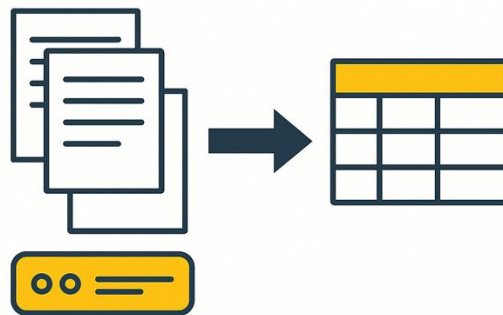
Riscos éticos e epistémicos



Alucinação e Factualidade

Os modelos de linguagem geram textos coerentes mas não baseados na exatidão factual. O texto gerado é o mais provável e não necessariamente o correto.

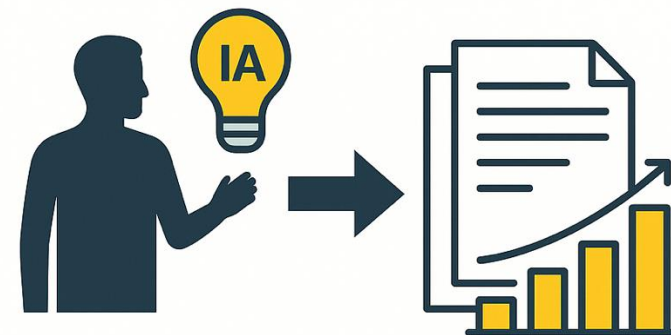
Há um risco de disseminação de informação falsa com uma aparência de autoridade, o que pode minar a confiança na ciência.



Perda de Contexto e Significado

A necessidade de tornar os dados mais legíveis pelas máquinas pode simplificá-los excessivamente, removendo metadados contextuais.

Corremos o risco de se perder contexto e significado na “tradução” de dados qualitativos para formatos padronizados.



Erosão da Integridade Científica

A responsabilidade de cada um (ou a falta dela) na utilização das ferramentas de IA, é um fator determinante para a qualidade do trabalho científico.

A pressão para publicar leva a que investigadores vejam na IA um multiplicador de produtividade, apostando na quantidade em vez da qualidade.

Estratégia de 3 pilares para o uso responsável da IA



Governança

Estabelecer políticas, processos e responsabilidades claras para orientar o uso ético, legal e seguro da tecnologia.



Princípios

Adotar princípios fundamentais para garantir que os dados e modelos são robustos, transparentes e adequados para o propósito.



Ferramentas

Implementar arquiteturas e tecnologias que mitiguem os riscos intrínsecos dos modelos, como a alucinação e a falta de proveniência.

Pilar 1: Governança

1. Desenvolver Plano Gestão Dados Dinâmico

O PGD não é um documento estático. Deve ser um instrumento vivo que descreve o ciclo de vida dos dados, desde a recolha até à preservação ou eliminação, e documenta as decisões relativas à proteção de dados.

2. Definir Papéis e Responsabilidades Claras

Identificar o "Responsável pelo Tratamento" (geralmente a instituição), os operadores e os subcontratantes. Em projetos colaborativos, formalizar a "Responsabilidade Conjunta" através de acordos.

3. Implementar Políticas Internas do Uso de IA

Estabelecer diretrizes claras sobre ferramentas autorizadas, usos permitidos/proibidos, e a necessidade de revisão humana. Por exemplo, proibir a inserção de dados pessoais ou confidenciais em ferramentas públicas não aprovadas.

4. Avaliar e Gerir a Cadeia de Fornecedores

Realizar uma avaliação rigorosa dos fornecedores de IA, analisando a origem dos dados de treino, as políticas de retenção de inputs e as cláusulas contratuais sobre titularidade dos outputs.

Pilar 2: Princípios

Para que a IA seja confiável, os dados subjacentes devem ser de alta qualidade. Os princípios FAIR (Findable, Accessible, Interoperable, Reusable) são o ponto de partida, mas a era da IA exige uma nova camada de responsabilidade.

FAIR → FAIR-R

“Responsibly-licensed”

FAIR → FAIR²

“AIR – AI-Readiness”

“RAI – Responsible-AI”

Operacionalizar FAIR para a era da IA

Não basta que os dados sejam reutilizáveis, a sua licença deve ser clara, legível por máquinas e adequada para o treino de modelos de IA.

Mas também devemos explorar novos modelos de partilha mais sofisticados que equilibram a abertura com a proteção:



**Abertura
Condicional**



**Abertura
Recíproca**



**Abertura
Protegida**

Pilar 3: Ferramentas

A governança e os princípios precisam de ser implementados através de arquiteturas tecnológicas que mitiguem os riscos inerentes aos modelos de IA. A solução não é só criar modelos melhores, mas sim controlar melhor os seus inputs e outputs.

Arquitetura RAG: Fundamentar Respostas em Fontes Verificáveis



1. Pergunta

Utilizador faz a pergunta



2. Recuperação (Retrieval)

Pesquisa numa BD vectorial de documentos internos e verificados para encontrar os mais relevantes.



3. Acréscimo (Augmentation)

Os documentos obtidos são inseridos no prompt e junto à pergunta original, servindo de contexto.



4. Geração (Generation)

“Responde a esta pergunta usando APENAS as fontes fornecidas.”

Checklist para a vossa instituição (ideal)

Governança

- Temos uma política interna clara que define ferramentas aprovadas e regras para o manuseamento de dados confidenciais?
- Avaliamos sistematicamente os riscos (jurídicos, de segurança, éticos) de novas ferramentas e projetos de IA?
- Estão definidos os responsáveis pelo tratamento de dados e pela supervisão do uso de IA?

Princípios

- As nossas políticas de gestão de dados incluem orientações sobre licenciamento responsável para reutilização em IA?
- Oferecemos formação aos investigadores sobre IA generativa, os seus riscos (viés, alucinação) e boas práticas?
- Exigimos documentação sobre os modelos de IA que utilizamos, especialmente a origem dos dados de treino?

Ferramentas

- Privilegiamos o uso de ferramentas de IA validadas e seguras, alojadas localmente ou com contratos robustos, em vez de serviços proprietários?
- Estamos a implementar arquiteturas como RAG para fundamentar as respostas da IA em fontes internas e verificáveis?
- Os nossos processos garantem que os outputs da IA, especialmente em decisões críticas, são sempre revistos por humanos?

Deliverable do Projeto Re.Data



Publicado em Novembro 2025

Autoria:

Nuno David, Iscte – Instituto Universitário de Lisboa

Marta Cordeiro, Iscte – Instituto Universitário de Lisboa

Gabriel Cipriano, Iscte – Instituto Universitário de Lisboa

Clara Boavida, Iscte – Instituto Universitário de Lisboa

Jorge Figueiredo, Universidade do Minho

Cláudia Conceição, Instituto de Higiene e Medicina Tropical,
Universidade Nova de Lisboa

Paula Ochôa, Faculdade de Ciências Sociais e Humanas,
Universidade Nova de Lisboa

Kevin Gallagher, Faculdade de Ciências e Tecnologia,
Universidade Nova de Lisboa

Carina Cunha, Iscte – Instituto Universitário de Lisboa

https://redata.pt/questoes_juridicas_protecao_de_-dados_e_licencas/



Re.Data

Rede para a Gestão de
Dados de Investigação

Gestão de Dados na Era da IA Generativa

Re.Data Lunch Webinar

António Luís Lopes (antonio.luis@iscte-iul.pt)

Iscte – Instituto Universitário de Lisboa



Universidade do Minho



UNIVERSIDADE D
COIMBRA

