

PROTÓTIPO WEB PARA VERIFICAÇÃO DE FAKE NEWS COM API GEMINI**WEB PROTOTYPE FOR FAKE NEWS VERIFICATION WITH GEMINI API**

Rafael Murilo Banin de Sena – rafaelmurilobs@gmail.com
Faculdade de Tecnologia de Taquaritinga (Fatec) – Taquaritinga – São Paulo – Brasil

Prof. Me João de Lucca Filho – joaodelucca@terra.com.br
Fatec Taquaritinga – Taquaritinga – São Paulo – Brasil

DOI: 10.31510/infa.v22i2.2324

Data de submissão: 24/09/2025

Data do aceite: 03/12/2025

Data da publicação: 20/12/2025

RESUMO

A disseminação de notícias falsas na internet representa um risco crescente para a sociedade, comprometendo a confiança das pessoas e dificultando a tomada de decisões informadas. Este artigo apresenta o desenvolvimento e avaliação de um protótipo web simples que utiliza a API Gemini para analisar frases curtas em português, classificando-as como verdadeiras, falsas ou incertas. O sistema também oferece uma justificativa breve e indica fontes confiáveis para consulta. Foi realizado testes com 20 afirmações, tendo o principal objetivo de avaliar a precisão do sistema, sua capacidade de classificar corretamente as frases, sua coerência e o seu nível de confiança. O protótipo alcançou 95% de exatidão, mostrando que o sistema se comporta de forma consistente e confiável, embora apresente mais dificuldade em casos incertos. O estudo sugere melhorias, como ajuste no nível de confiança e integração com buscadores online, reforçando a importância do uso consciente da ferramenta em contextos educacionais e sociais, visando combater de maneira eficaz a propagação da desinformação.

Palavras-chave: Fake News. Verificação de Fatos. Inteligência Artificial. API Gemini.

ABSTRACT

The spread of false information on the internet poses a growing risk to society, undermining trust and hindering informed decision-making. This paper presents the development and evaluation of a simple web prototype that uses the Gemini API to analyze short statements in Portuguese, classifying them as true, false, or uncertain. The system also provides a brief justification for each classification and suggests reliable sources for verification. Tests were conducted with 20 statements to assess the system's accuracy, its ability to classify statements correctly, its consistency, and the level of confidence assigned by the model. The prototype achieved 95% accuracy, showing consistent and reliable performance in identifying both true and false statements, though it had more difficulty with uncertain cases. The study proposes improvements such as adjusting the confidence level and integrating the tool with online search engines for real-time verification. This work emphasizes the importance of responsible use of such tools in educational and social contexts to combat misinformation effectively.

Keywords: Fake News. Fact-checking. Artificial Intelligence. API Gemini.

1 INTRODUÇÃO

O aumento da circulação de informações falsas na internet tem se tornado um desafio urgente. Em poucos minutos, uma notícia enganosa pode se espalhar por redes sociais, aplicativos de mensagens e sites diversos. Durante a pandemia de COVID-19, no Brasil, a propagação de informações falsas teve um impacto profundo nas decisões das pessoas e nas questões políticas, além de afetar a confiança da população nas autoridades. Muitas dessas notícias falsas eram criadas com o objetivo de gerar mais engajamento ou promover interesses políticos específicos. O resultado disso foi prejudicial em áreas como saúde pública, eleições e até mesmo na educação, causando confusão e desinformação generalizada. Diante desse cenário, cresce a demanda por ferramentas que ajudem a verificar a veracidade de frases e notícias de forma rápida e acessível (Lazer, 2018; Vosoughi et al., 2018; Barcelos, 2021).

Para atender a essa necessidade, o artigo apresenta o desenvolvimento de um protótipo web, no qual o usuário insere uma frase curta em português, e a ferramenta, com apoio da API Gemini, informa se a afirmação é verdadeira, falsa ou incerta. Além disso, exibe um resumo explicativo e links para fontes confiáveis, incentivando a checagem por parte do usuário. Nos testes, o sistema alcançou 95% de precisão, mostrando bons resultados em afirmações claramente verdadeiras ou falsas, mas certa dificuldade com frases ambíguas. Isso mostra o potencial da ferramenta, ao mesmo tempo que aponta oportunidades de aprimoramento, como o ajuste da confiança nas respostas e a possível integração com mecanismos de busca online.

A Inteligência Artificial (IA), especialmente nas áreas de aprendizado de máquina e Modelos de Linguagem de Grande Escala (LLMs), se tornou um dos maiores avanços tecnológicos dos últimos tempos. A IA permite automatizar tarefas complexas, como identificar padrões e tomar decisões com base em grandes quantidades de dados. Nos últimos anos, com o aumento da disponibilidade de dados e do poder computacional, as tecnologias de IA têm mostrado um desempenho impressionante em tarefas como organizar informações, gerar resumos e validar dados. Suas aplicações têm se destacado em setores como saúde, educação e na verificação de informações. (Norvig e Russell, 2021; Bommasani, 2021).

Este estudo é importante devido à sua relevância social, política e científica no combate à desinformação, utilizando tecnologias que possam tornar a verificação de fatos mais rápida e acessível. Ter uma ferramenta de verificação simples e de fácil acesso é essencial, pois permite que qualquer pessoa, independentemente do seu conhecimento técnico, consiga validar informações de forma prática. A Inteligência Artificial se apresenta como uma parceira fundamental nesse processo, oferecendo soluções escaláveis para identificar e refutar

informações falsas. Isso é especialmente crucial em contextos sociais e educacionais, onde os danos causados pela desinformação podem ser ainda maiores.

A partir disso, este estudo tem como objetivo principal, desenvolver e avaliar um protótipo de verificação de informações, utilizando Inteligência Artificial (IA) com o apoio da API Gemini. O protótipo será capaz de classificar afirmações como verdadeiras, falsas ou incertas, fornecendo ao usuário uma análise detalhada sobre cada caso.

O trabalho está organizado em capítulos que começam com uma discussão sobre os conceitos de Inteligência Artificial (IA), fake news e o papel da IA na verificação de informações. Em seguida, é detalhado o desenvolvimento do protótipo, incluindo as tecnologias utilizadas. Depois, são apresentados os testes realizados com o protótipo, acompanhados de uma análise do seu desempenho. Por fim, a conclusão traz sugestões de melhorias e aponta possíveis direções para o aperfeiçoamento da ferramenta.

2 FUNDAMENTAÇÃO TEÓRICA

2.1 Inteligência Artificial

A Inteligência Artificial (IA) constitui um campo da ciência da computação dedicado ao desenvolvimento de sistemas capazes de desempenhar atividades que, tradicionalmente, dependem de competências humanas, como aprendizado a partir de dados, reconhecimento de padrões e tomada de decisões. O amadurecimento desse domínio resultou de décadas de pesquisa, sendo atualmente impulsionado pela disponibilidade massiva de dados, pelo aumento do poder computacional e pelos avanços em arquiteturas de redes neurais profundas (Norvig e Russell, 2021; Courville et al., 2016; Sichman, 2021).

Dentro da IA, destacam-se diferentes áreas de investigação, incluindo aprendizado de máquina, visão computacional, processamento de linguagem natural e planejamento automatizado. Na prática, o objetivo central é desenvolver tecnologias que reproduzam processos cognitivos de forma autônoma, eficiente e confiável em variados cenários de aplicação (Norvig e Russell, 2021).

O aprendizado de máquina, um dos principais ramos da IA, organiza-se em três categorias predominantes. No aprendizado supervisionado, o modelo é treinado com dados previamente rotulados, permitindo-lhe associar entradas a saídas conhecidas. No aprendizado não supervisionado, o sistema identifica padrões em bases de dados sem rótulos, como ocorre em algoritmos de agrupamento. Já no aprendizado por reforço, a ênfase está na interação de um

agente com um ambiente, de modo que as escolhas realizadas são guiadas por recompensas e penalidades recebidas em função de suas ações (Courville et al., 2016).

Nos últimos anos, os Modelos de Linguagem de Grande Porte (Large Language Models - LLMs) tornaram-se um tema central na área de IA. Esses sistemas são treinados com imensos volumes de dados textuais e possuem a capacidade de compreender e gerar linguagem natural. Diferentemente de abordagens clássicas da IA, geralmente voltadas a funções específicas, os LLMs integram a chamada IA generativa, capaz de produzir conteúdos originais, sintetizar informações e responder a questões abertas (Bommasani, 2021).

A base tecnológica dos LLMs é a arquitetura Transformer, que representou um marco ao introduzir mecanismos de atenção, permitindo o processamento paralelo e eficiente de sequências textuais. A inovação viabilizou treinamento de modelos com bilhões de parâmetros, permitiu o surgimento dos sistemas de grande escala utilizados atualmente (Vaswani, 2017).

Na atualidade, os LLMs já são aplicados em diferentes tarefas, como classificação de informações, elaboração de resumos e resposta a perguntas. Essa versatilidade tem favorecido seu uso em contextos práticos, incluindo atividades de apoio à checagem e validação de dados (Bommasani, 2021).

2.2 Agente de IA

Agentes de IA podem ser definidos como sistemas computacionais projetados para perceber o ambiente em que estão inseridos, tomar decisões e agir de maneira autônoma com o propósito de atingir objetivos previamente determinados. Esses sistemas se caracterizam por propriedades como autonomia, capacidade de adaptação, iniciativa e habilidade de comunicação tanto com outros agentes quanto com seres humanos. Tais atributos os distinguem de programas tradicionais que apenas reagem a estímulos de forma pré-programada, possibilitando sua utilização em sistemas multiagentes, nos quais múltiplos agentes interagem cooperando ou competindo para solucionar problemas de elevada complexidade (Nas, 2025).

Segundo Wooldridge (2009), os agentes podem ser classificados em duas categorias fundamentais. Os agentes reativos operam de forma imediata, respondendo diretamente aos estímulos recebidos, sem manter uma representação interna do ambiente. Em contrapartida, os agentes deliberativos constroem modelos do mundo ao seu redor, utilizando-os para planejar e tomar decisões fundamentadas em tais representações.

A proposta apresentada neste trabalho enquadra-se no conceito de agente deliberativo, uma vez que envolve a interpretação das entradas fornecidas pelo usuário e a geração de respostas elaboradas por meio da utilização de modelos de linguagem.

2.3 API e API Gemini

Uma API (Interface de Programação de Aplicações) consiste em um conjunto de protocolos e especificações que permitem a comunicação entre diferentes componentes de um sistema. Com a disseminação da arquitetura REST (Representational State Transfer), proposta por Fielding e Taylor (2002), as APIs RESTful consolidaram-se como padrão para integração entre sistemas. Essa abordagem utiliza o protocolo HTTP, possibilitando que aplicações cliente-servidor compartilhem informações de forma escalável, flexível e independente da linguagem de programação empregada.

No contexto deste trabalho, a API foi empregada para estabelecer a conexão entre o sistema de verificação de fake news e um modelo de linguagem de grande porte, denominado API Gemini. Essa integração viabiliza o envio de informações para análise e a recepção de respostas do modelo de IA de maneira eficiente e ágil (Sommerville, 2016).

A opção pela API Gemini fundamentou-se em sua capacidade de fornecer respostas precisas, apoiadas em modelos de aprendizado de máquina treinados com extensos conjuntos de dados. O modelo de linguagem subjacente a essa API baseia-se na arquitetura Transformer, que utiliza mecanismos de atenção para processar sequências de texto de forma eficiente. Tal estrutura permite uma compreensão aprimorada da linguagem natural, resultando em respostas coerentes e contextualizadas a partir de uma ampla base de conhecimento (Google, 2025).

2.4 Fake News

As fake news consistem em conteúdos falsos apresentados como notícias legítimas. Frequentemente, esses materiais utilizam uma linguagem semelhante à do jornalismo tradicional, buscando transmitir credibilidade e influenciar o público, seja para gerar engajamento, audiência ou obter benefícios políticos e econômicos. A propagação dessas informações ocorre predominantemente por meio de redes sociais, que ampliam seu alcance e dificultam a checagem imediata dos fatos (Lazer, 2018).

No contexto brasileiro, estudos indicam que, apenas em 2020, circularam pelo menos 329 informações falsas relacionadas à COVID-19, com forte presença em plataformas como WhatsApp e Facebook. Muitas dessas notícias envolviam temas políticos ou sanitários, abordando óbitos, métodos de prevenção e tratamentos. (Barcelos, 2021).

Esses episódios geraram picos de buscas na internet por termos relacionados, revelando o alto potencial de disseminação dessas narrativas. Estudos apontam que informações falsas circulam com maior velocidade, alcance e impacto do que notícias verdadeiras, sobretudo quando mobilizam emoções como surpresa ou indignação (Vosoughi et al., 2018).

Indivíduos com acesso limitado a fontes confiáveis, analfabetismo funcional ou inseridos em ambientes com restrição de informações são particularmente suscetíveis aos efeitos da desinformação. Consequentemente, o combate às fake news demanda não apenas ferramentas tecnológicas, mas também ações de educação para o consumo crítico de informações, bem como iniciativas que promovam transparência e verificação rigorosa dos fatos (Wardle e Derakhshan, 2017; Lazer, 2018).

3 PROCEDIMENTOS METODOLÓGICOS

3.1 Abordagem Web e Ferramentas Utilizadas

A metodologia adotada neste estudo enquadra-se na pesquisa aplicada, com foco no desenvolvimento e avaliação de um protótipo funcional. A proposta consiste em uma aplicação web, projetada para proporcionar uma experiência interativa e acessível a usuários interessados em verificar a veracidade de informações. A arquitetura do sistema segue o modelo cliente-servidor: a interface do usuário (*frontend*) encaminha os dados, enquanto o *backend* processa as informações por meio da API Gemini, responsável pela análise e pela classificação das afirmações em três categorias: *verdadeiras*, *falsas* ou *incertas*.

Para a implementação da aplicação, foram empregadas tecnologias atuais e amplamente utilizadas. O Vite foi escolhido como ambiente de desenvolvimento e empacotamento, devido ao seu desempenho elevado e compatibilidade com ferramentas atuais. O projeto foi desenvolvido em TypeScript, que acrescenta tipagem estática ao JavaScript, proporcionando maior robustez e segurança ao código. Já a construção da interface baseou-se no React, uma das bibliotecas mais difundidas para criação de interfaces dinâmicas e responsivas (Vite, 2025; Microsoft, 2025; Meta, 2025).

No desenvolvimento dos componentes visuais, utilizou-se a biblioteca shadcn-ui, integrada ao React e ao Tailwind CSS, o que favorece a criação de interfaces modernas, reutilizáveis e customizáveis. O Tailwind CSS, por sua vez, possibilitou a elaboração de layouts responsivos de forma ágil, garantindo uma experiência visual limpa, intuitiva e consistente para a aplicação (Shadcn, 2025; Tailwind Labs, 2025).

3.2 Criação do Agente de Verificação

Quando o usuário insere uma afirmação na interface e aciona o comando de verificação, a informação é transmitida do *frontend* para o *backend*. Este, por sua vez, realiza uma requisição à API Gemini, encarregada de processar e analisar a frase.

O procedimento envolve o envio de um prompt estruturado ao modelo de IA, que responde classificando a afirmação em uma das três categorias: verdadeira, falsa ou incerta.

No núcleo da aplicação encontra-se um agente de verificação, desenvolvido especificamente para interpretar as afirmações recebidas. Esse agente foi configurado com um prompt direcionado, de modo a assegurar que o modelo execute uma análise criteriosa e apresente, além da classificação, uma justificativa fundamentada para cada resposta gerada.

Ilustração 1- Prompt Agente

```
const prompt = `
Você é um verificador de fatos profissional especializado em análise de informações em português.
Analisar a seguinte afirmação e determine se é verdadeira, falsa ou incerta.

TEXTO PARA VERIFICAR: "${text}"

INSTRUÇÕES DETALHADAS:
1. Procure por informações atuais sobre este tópico na internet
2. Verifique múltiplas fontes confiáveis (sites oficiais, órgãos de imprensa respeitados, instituições)
3. Compare as informações encontradas com a afirmação
4. Seja DECISIVO na sua análise - evite respostas "incertas" quando há evidências claras
5. Para notícias recentes, procure por reportagens de veículos de imprensa conhecidos
6. Para dados científicos, procure por fontes acadêmicas ou órgãos oficiais
7. Para informações sobre pessoas públicas, verifique fontes oficiais

CRITÉRIOS DE CLASSIFICAÇÃO:
- VERDADEIRO (real): Quando há evidências claras e múltiplas fontes confirmam a informação
- FALSO (fake): Quando há evidências que contradizem a afirmação ou não há fontes confiáveis
- INCERTO (uncertain): APENAS quando realmente não há informações suficientes ou fontes conflitantes

Responda EXATAMENTE neste formato JSON:
{
  "status": "verdadeiro|falso|incerto",
  "confidence": [número de 70-95 para verdadeiro/falso, 30-60 para incerto],
  "justification": "Explicação clara e detalhada em português do porquê da classificação, mencionando as fontes verificadas",
  "sources": [
    {
      "title": "Título da fonte",
      "url": "URL da fonte (use URLs reais quando possível)",
      "summary": "Resumo do que a fonte diz sobre o assunto"
    }
  ]
}

IMPORTANTE: Seja confiante na sua análise. Se encontrar evidências claras, classifique como "verdadeiro" ou "falso" com alta confiança.
Use "incerto" apenas quando realmente não há informações suficientes.
`;
```

Fonte: Elaborado pelo autor (2025)

4 RESULTADOS E DISCUSSÃO

Nesta etapa do estudo, são apresentados os resultados obtidos a partir do protótipo desenvolvido para verificação de *fake news*, acompanhados de uma análise visual e reflexiva acerca de seu desempenho. A seção tem como propósito demonstrar o funcionamento prático do sistema, avaliar seu nível de precisão e discutir suas limitações, além de indicar possibilidades de aprimoramento em versões futuras.

4.1 Interface do projeto e exemplo

O protótipo foi desenvolvido como aplicação web, com interface simples e intuitiva. O usuário insere uma afirmação no campo designado e, pode visualizar o resultado da verificação.

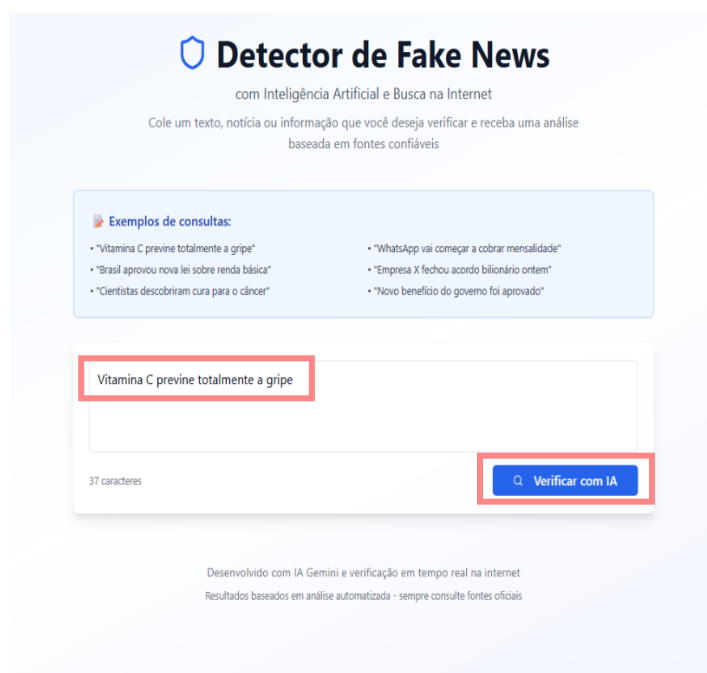
Ilustração 2 - Interface do projeto



Fonte: Elaborado pelo autor (2025)

Para exemplificar o funcionamento do sistema, algumas interações são registradas. Os registros confirmam como protótipo classifica diferentes afirmações inseridas pelos usuários.

Ilustração 3 - Campo para inserção de afirmação e botão de verificação



Fonte: Elaborado pelo autor (2025)

Clicando o botão *verificar com IA*, o sistema exibe, logo abaixo do campo de entrada: a classificação da afirmação (verdadeira, falsa ou incerta), o nível de confiança, um pequeno resumo sobre o conteúdo e uma lista de fontes confiáveis, que ajudam na resposta apresentada.

Ilustração 4 - Afirmação “Vitamina C previne totalmente a gripe”, o sistema classifica como falsa

⊗ **Possível Fake News** 90% confiança

Análise da IA:

A afirmação “Vitamina C previne totalmente a gripe” é FALSA. Embora a vitamina C tenha demonstrado em alguns estudos ter um leve efeito na redução da duração e gravidade dos sintomas da gripe em alguns indivíduos, nenhum estudo científico conclusivo comprova que ela previne “totalmente” a gripe. A prevenção da gripe depende de fatores como higiene, vacinação e fortalecimento do sistema imunológico, sendo a vitamina C apenas um dos componentes que podem contribuir para a saúde geral, mas não uma garantia de imunidade completa contra o vírus influenza. A afirmação de prevenção total é uma generalização excessiva e não condiz com o consenso científico atual. Diversos estudos mostram resultados variados e, em muitos casos, nenhum efeito significativo na prevenção da gripe. A ausência de uma prevenção total é corroborada por órgãos de saúde pública como o CDC (Centers for Disease Control and Prevention) e a OMS (Organização Mundial da Saúde), que recomendam outras medidas de prevenção.

Fontes de Verificação (3):

Centers for Disease Control and Prevention (CDC) - Flu

O CDC fornece informações abrangentes sobre a gripe, incluindo medidas de prevenção que não incluem a vitamina C como uma forma de prevenção total. Enfatiza a vacinação e medidas de higiene.

Organização Mundial da Saúde (OMS) - Gripe

A OMS também destaca a importância da vacinação e medidas de higiene para a prevenção da gripe, sem mencionar a vitamina C como uma forma de prevenção total.

Exemplo de estudo científico (uma busca por 'vitamin C influenza prevention' em PubMed retornará diversos estudos com resultados variados, nenhum mostrando prevenção total):

Uma busca na base de dados PubMed por estudos sobre a prevenção da gripe pela vitamina C revelará diversos estudos com resultados conflitantes. A maioria não demonstra uma prevenção total, mas sim efeitos possivelmente benéficos em alguns casos, como redução da duração dos sintomas.

💡 **Dica:** Sempre consulte múltiplas fontes confiáveis e veículos de imprensa respeitados para verificar informações importantes.

Fonte: Elaborado pelo autor (2025)

Ilustração 5 - Afirmação for “A água ferve a 100°C ao nível do mar”, o sistema classifica como verdadeira.

✔ **Informação Verdadeira**
90% confiança

Análise da IA:

A afirmação “A água ferve a 100°C ao nível do mar” é verdadeira. Este é um ponto fundamental da física e da química, amplamente ensinado e aceito pela comunidade científica. A temperatura de ebulição da água é afetada pela pressão atmosférica; ao nível do mar, a pressão atmosférica é de aproximadamente 1 atmosfera, e sob essa pressão, a água entra em ebulição a 100°C. Encontrar fontes que contradigam essa informação seria extremamente difícil, pois se trata de um conhecimento básico e amplamente validado. Diversos livros didáticos, sites de ensino de física e química, e até mesmo calculadoras online de ponto de ebulição confirmam essa informação. A ausência de fontes que contradigam essa informação, aliada à sua ampla aceitação científica, justifica a classificação como verdadeira com alta confiança.

Fontes de Verificação (3):

Tabela Periódica - Propriedades da Água

Embora não mencione explicitamente a temperatura de ebulição ao nível do mar, a página descreve as propriedades da água, incluindo seu ponto de ebulição, que é consistentemente citado como 100°C em condições normais de pressão.

Khan Academy - Mudanças de fase da água

A Khan Academy, plataforma educacional respeitada, explica o processo de ebulição da água e o impacto da pressão, corroborando a informação de que a 100°C ao nível do mar a água ferve.

HyperPhysics - Boiling Point of Water

(Em inglês) Esta fonte da Georgia State University fornece uma tabela com pontos de ebulição da água em diferentes pressões, confirmando o valor de 100°C para 1 atm (nível do mar).

💡 **Dica:** Sempre consulte múltiplas fontes confiáveis e veículos de imprensa respeitados para verificar informações importantes.

Fonte: Elaborado pelo autor (2025)

Ilustração 6 - Afirmação for “Empresa X fechou acordo bilionário ontem”, o sistema classifica como incerta

⚠ **Verificação Incerta**
45% confiança

Análise da IA:

A afirmação “Empresa X fechou acordo bilionário ontem” é muito vaga. Para verificar sua veracidade, é necessário saber qual é a Empresa X. Sem identificar a empresa, é impossível procurar por informações em fontes confiáveis como sites de notícias, comunicados oficiais da empresa ou registros financeiros. A falta de especificação da empresa torna a verificação impossível, resultando em uma classificação incerta.

Fontes de Verificação (1):

Análise de Verificação de Fatos

Busca realizada para verificar a veracidade da informação

💡 **Dica:** Sempre consulte múltiplas fontes confiáveis e veículos de imprensa respeitados para verificar informações importantes.

Fonte: Elaborado pelo autor (2025)

4.3 Relatório de Testes

Para avaliar o desempenho do protótipo, conduziram-se testes com 20 afirmações, categorizadas como *verdadeiro*, *falso* ou *incerto*. O objetivo foi mensurar a precisão do sistema, sua capacidade de classificação, coerência das respostas e o nível de confiança comunicado ao usuário. As afirmações foram inseridas manualmente via interface da aplicação web.

Tabela 1– Relatório de Testes com 20 Afirmações

Afirmação	Rótulo alvo	Predição do app	Confiança	Obs.
Vitamina C previne totalmente a gripe	F	F	90%	OK
O Amazonas é o maior estado brasileiro em área	V	F	90%	Erro
O Brasil adotou o euro como moeda oficial	F	F	90%	OK
A água ferve a 100°C ao nível do mar	V	V	90%	OK
WhatsApp vai começar a cobrar mensalidade amanhã	F	F	90%	OK
Beber água a cada 15 min elimina vírus	F	F	90%	OK
Lei federal perdoa todas as dívidas do FIES	F	F	90%	OK
Cura definitiva do câncer foi anunciada	F	F	90%	OK
A Terra gira em torno do Sol	V	V	95%	OK
Todos os vacinados ficam doentes em um ano	F	F	90%	OK
Fuso de Brasília é UTC-3	V	V	90%	OK
IPTU de SP cairá 90% para todos	F	F	90%	OK
Suplemento X melhora significativamente a memória em adultos	I	I	90%	OK
Governo anunciou benefício de R\$ 2.000 para todos hoje	F	F	90%	OK
Empresa X fechou acordo bilionário ontem	I	I	45%	OK
Instituto Butantã fica em São Paulo	V	V	90%	OK
Plutão é um planeta do Sistema Solar	F	F	90%	OK
População do Brasil ultrapassa 250 milhões	F	F	90%	OK
Vacina da gripe contém microchips	F	F	90%	OK
Em 2023 o presidente dos EUA era Barack Obama	F	F	90%	OK

Fonte: Elaborado pelo autor (2025)

Os testes realizados demonstraram uma taxa de acerto de 95%, correspondendo a 19 classificações corretas em um total de 20 casos analisados. Em grande parte das execuções, o nível de confiança fornecido pelo modelo manteve-se estável em 90%. Foram observadas apenas duas exceções: uma classificação atingiu 95% de confiança, enquanto outra apresentou valor reduzido de 45%. Esses resultados indicam que o sistema apresentou comportamento consistente e desempenho satisfatório na maioria dos cenários avaliados.

5. CONCLUSÃO

Os testes realizados com 20 afirmações demonstraram que o protótipo possui um desempenho bastante satisfatório, alcançando 95% de precisão. O sistema obteve bons resultados tanto na identificação de informações falsas quanto verdadeiras, embora tenha apresentado mais dificuldade em casos considerados incertos. Observou-se também que o nível de confiança retornado pelo modelo permaneceu praticamente constante em 90%, o que sugere a necessidade de ajustes na calibração, principalmente para melhorar a diferenciação em situações mais ambíguas. Entre os poucos erros identificados, apenas um envolveu uma informação real, reforçando a eficiência geral da ferramenta.

Uma limitação significativa é que a API Gemini, embora robusta, não realiza buscas em tempo real. Assim, afirmações recentes podem não ser avaliadas com precisão, pois as respostas se baseiam apenas nos dados utilizados no treinamento do modelo. Para contornar essa

restrição, recomenda-se a integração com fontes externas, como APIs de busca (Google Search, Bing), plataformas especializadas em checagem de fatos (Snopes, PolitiFact, FactCheck.org) e bases acadêmicas atualizadas, de modo a aumentar a confiabilidade das análises.

O objetivo central deste estudo foi desenvolver e avaliar um protótipo de verificação de informações baseado em Inteligência Artificial, com o protótipo cumprindo seu papel como ferramenta de apoio à verificação de afirmações curtas, contribuindo para a circulação de informações mais confiáveis. O sistema apresentou alta taxa de acerto, sobretudo em afirmações diretas e bem estruturadas. Entretanto, a precisão dos resultados depende fortemente da qualidade da entrada: enunciados ambíguos ou carentes de contexto podem reduzir a eficácia da ferramenta, aspecto que deve ser considerado em seu uso.

Como próximos passos, indica-se aprimorar a calibração do nível de confiança, incorporar estratégias para lidar com frases vagas e integrar mecanismos de busca em tempo real, aumentando a precisão e a segurança na verificação de conteúdos atuais.

REFERÊNCIAS

- BARCELOS, T. N.; et. al.: **Análise de fake news veiculadas durante a pandemia de COVID-19 no Brasil**. Revista Panamericana de Saúde Pública, v. 45, e65, 2021. Disponível em: <doi.org/10.26633/RPSP.2021.65>. Acesso em: 12 set. 2025.
- BOMMASANI, R. et al. *On the Opportunities and Risks of Foundation Models*. arXiv, 2021. DOI: 10.48550/arXiv.2108.07258. Disponível em: <arxiv.org/abs/2108.07258>. Acesso em: 12 set. 2025.
- COURVILLE, A.; GOODFELLOW, I.; BENGIO, Y.: *Deep Learning*. MIT Press, 2016. Disponível em: <deeplearningbook.org/>. Acesso em: 10 set. 2025.
- FIELDING, R.; TAYLOR, R. *Principled Design of the Modern Web Architecture*. ACM Transactions on Internet Technology, v. 2, n. 2, p. 115-150, 2002. Disponível em: <doi.org/10.1145/514183.514185>. Acesso em: 22 set. 2025.
- GOOGLE. *AI for Developers – Gemini API reference*. [S. l.], 2025. Disponível em: <ai.google.dev/api>. Acesso em: 15 set. 2025.
- LAZER, D. M. J. et al. *The Science of Fake News*. Science, v. 359, n. 6380, p. 1094-1096, 2018. Disponível em: <science.org/doi/10.1126/science.aao2998>. Acesso em: 12 set. 2025.
- META. *React – A JavaScript library for building user interfaces*. [S. l.], 2025. Disponível em: <react.dev/>. Acesso em: 22 set. 2025.
- MICROSOFT. *TypeScript: JavaScript With Syntax for Types*. [S. l.], 2025. Disponível em: <typescriptlang.org/>. Acesso em: 22 set. 2025.
- NAS, E.: **Agentes de IA: riscos e benefícios**. Jornal USP, São Paulo, 23 jan. 2025. Disponível em: <jornal.usp.br/artigos/agentes-de-ia-riscos-e-beneficios/>. Acesso em: 15 set. 2025.
- NORVIG, Peter; RUSSELL, Stuart J. *Artificial Intelligence: A Modern Approach*. 4. ed. London: Pearson, 2021. Disponível em: <aima.cs.berkeley.edu/>. Acesso em: 10 set. 2025.

SHADCN. shadcn/ui – Beautifully designed components. [S. 1.], 2025. Disponível em: <ui.shadcn.com/>. Acesso em: 22 set. 2025.

SICHMAN, J. S.: **Inteligência Artificial e sociedade: avanços e riscos**. 2021. Disponível em: <repositorio.usp.br/item/003108781>. Acesso em: 10 set. 2025.

SOMMERVILLE, I.: *Software Engineering*. 10. ed. Boston: Pearson, 2016. Disponível em: <pearson.com/store/p/software-engineering/P100000254213>. Acesso em: 22 set. 2025.

TAILWIND LABS. *Tailwind CSS – Rapidly build modern websites without ever leaving your HTML*. [S. 1.], 2025. Disponível em: <tailwindcss.com/>. Acesso em: 22 set. 2025.

VASWANI, A. et al. Attention is All You Need. *Advances in Neural Information Processing Systems*, 2017. Disponível em: <arxiv.org/abs/1706.03762>. Acesso em: 22 set. 2025.

VITE. Vite – *Next Generation Frontend Tooling*. [S. 1.], 2025. Disponível em: <vitejs.dev/>. Acesso em: 22 set. 2025.

VOSOUGHI, S.; ROY, D.; ARAL, S.: *The Spread of True and False News Online*. *Science*, v. 359, n. 6380, p. 1146-1151, 2018. Disponível em: <science.org/doi/10.1126/science.aap9559>. Acesso em: 12 set. 2025.

WARDLE, C.; DERAKHSHAN, H.: *Information Disorder: Toward an Interdisciplinary Framework for Research and Policymaking*. Strasbourg: Council of Europe, 2017. Disponível em: <rm.coe.int/information-disorder-report/168076277c>. Acesso em: 10 set. 2025.

WOOLDRIDGE, M. *An Introduction to MultiAgent Systems*. 2. ed. Chichester: Wiley, 2009.