



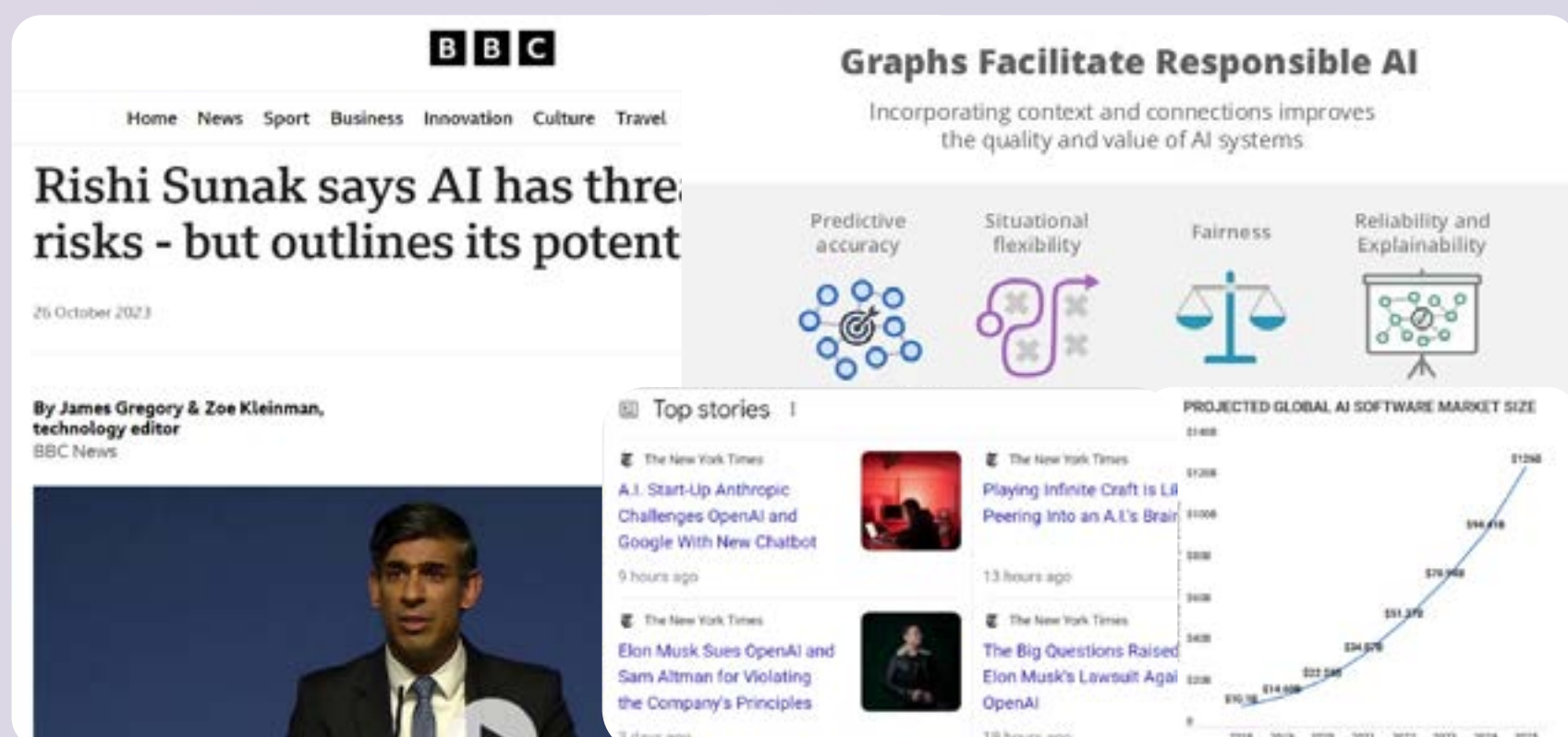
# Evaluation Copilot

Comprehensive evaluation that empowers developers to **DEVELOP WITH TRUST, DEPLOY IN CONFIDENCE.**

## Problem

The rapid integration of Large Language Models (LLMs) in app development introduces challenges in understanding and trusting AI-generated responses.

*A booming market with exponentially increasing news headlines and new developments on a daily basis, ...*



### Key Issues

**AI APP Devs**  
1. Human evaluations are costly and time-consuming;  
2. Struggles to interpret complex evaluation metrics;  
3. Gap between evaluation score and actionable insights.

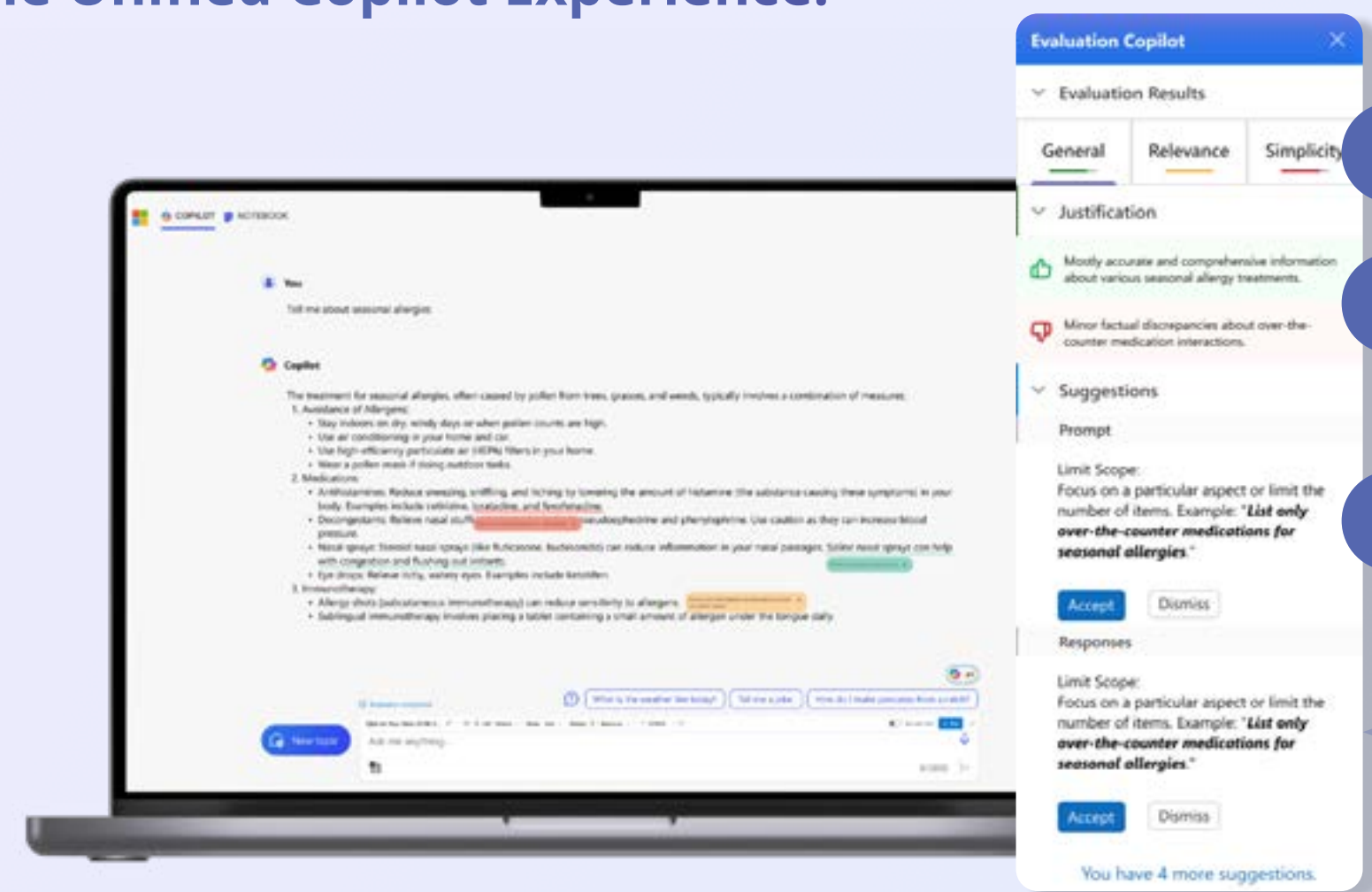
**Business Execs**  
Potential risks in application quality and business outcomes.

## Solution

The "Evaluation Copilot" is a web app that demystifies LLM evaluation metrics for developers, offering an intuitive platform to test, understand, and refine AI-generated text. It provides clear, actionable feedback on how to improve prompts for better LLM responses, ensuring developers can enhance AI reliability and effectiveness in their applications.

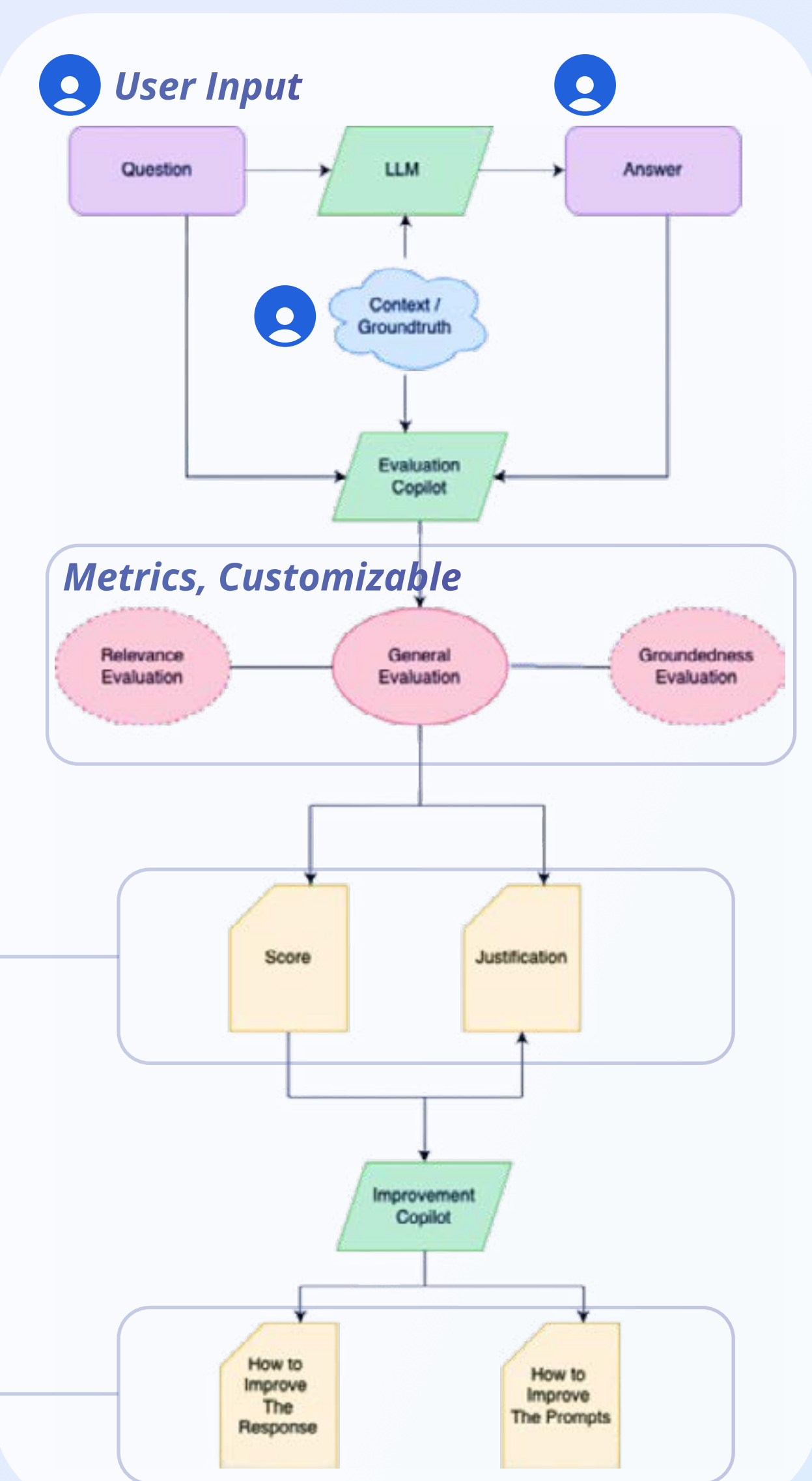
### Prototype

One Unified Copilot Experience.



- 1 Customizable Metrics.
- 2 Comprehensive Justifications.
- 3 Actionable Suggestions on Prompt and Answer, designed to be easily incorporated.
- 4 Contextual Experience.
- 5 Clear CTA, that enables an intuitive and seamless add-on.
- 6 Integrated Workflow, evaluate as the conversation goes.

### Software Architecture



## Process/Approach

Our team combined user-centered design with technological innovation, starting with extensive user research to understand developers' needs and pain points. We iteratively developed a series of prototypes following Agile framework, incorporating feedback from user testing sessions and learnings from the latest research outcomes in the field.

