



**Trusted Tech Africa**

# **Safety by Design Integration for EdTech LLMs**

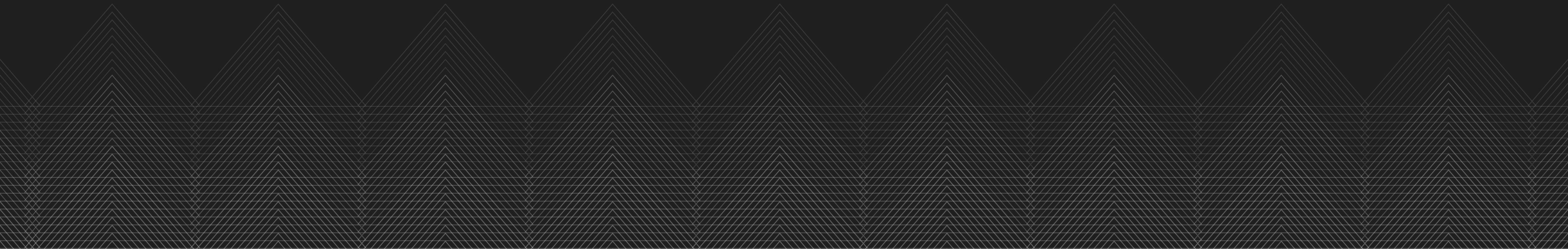




Many education technology platforms globally and increasingly across Africa are integrating artificial intelligence (AI) into their products or building AI-powered educational solutions from the ground up. In Africa alone, over 150+ AI tools are currently being used to support teaching, learning, assessment, and school management, with more emerging each year. This rapid growth presents significant opportunities for improving access to quality education, personalizing learning experiences, and addressing teacher shortages.

However, as AI becomes more embedded in the educational ecosystem, it also introduces new risks particularly for children and young learners such as exposure to harmful content, data privacy violations, algorithmic bias, and overdependence on automated systems. These risks make Safety by Design not just a regulatory concern, but a moral and practical imperative for EdTech startups and developers. Integrating safety from the outset ensures that AI-enhanced education is not only effective, but also trustworthy, inclusive, and developmentally appropriate for all learners.

Below are 9 structured guides on how to apply SbD principles to an LLM-powered EdTech product, whether you're building AI tutors, lesson generators, chatbots, or personalized learning platforms.



## 1. Understand the EdTech Use Case and Risk Profile

Who are the users? (Children, teens, teachers, parents)

What tasks are supported? (Homework help, tutoring, assessment, grading, admin tasks)

What risks may arise?

- Exposure to inappropriate content
- Privacy violations
- Disinformation or hallucinated facts
- Over-reliance on AI
- Biased or discriminatory content

Start with a safety risk matrix to identify high-impact areas.



## 2. Curate Safe and Age-Appropriate Training Data

Use education-specific datasets from trusted sources (curricula, textbooks).

Filter out:

- Hate speech, explicit material
- Biased or misleading content
- Culturally inappropriate examples
- Avoid using webscraped data without strong filters.

Use dataset auditing tools and annotation teams with child safety training.



### 3. Align the Model to Educational and Ethical Goals

Apply Reinforcement Learning from Human Feedback (RLHF) using:

- Teachers
- Child psychologists
- Curriculum experts

Use “Constitutional AI” to bake in norms like:

- Promoting respectful dialogue
- Avoidance of harmful assumptions about students’ abilities or backgrounds
- Providing constructive feedback

Define "Do's and Don'ts" for the model using value-aligned prompts.



## 4. Enforce Privacy and Child Data Protection by Design

Do not store or reuse personal student data without clear consent.

Comply with:

- GDPR (Europe)
- COPPA (US, for under-13s)
- NDPR (Nigeria) or relevant African privacy laws

Ensure anonymization of inputs and outputs.  
Avoid letting LLMs "memorize" or echo back user input.

**Use edge-processing or encrypted storage for any personal data.**



## 5. Build Age-Appropriate Guardrails

For Example Guardrail with specific age group

Under 13: No open-ended chats, predefined options only

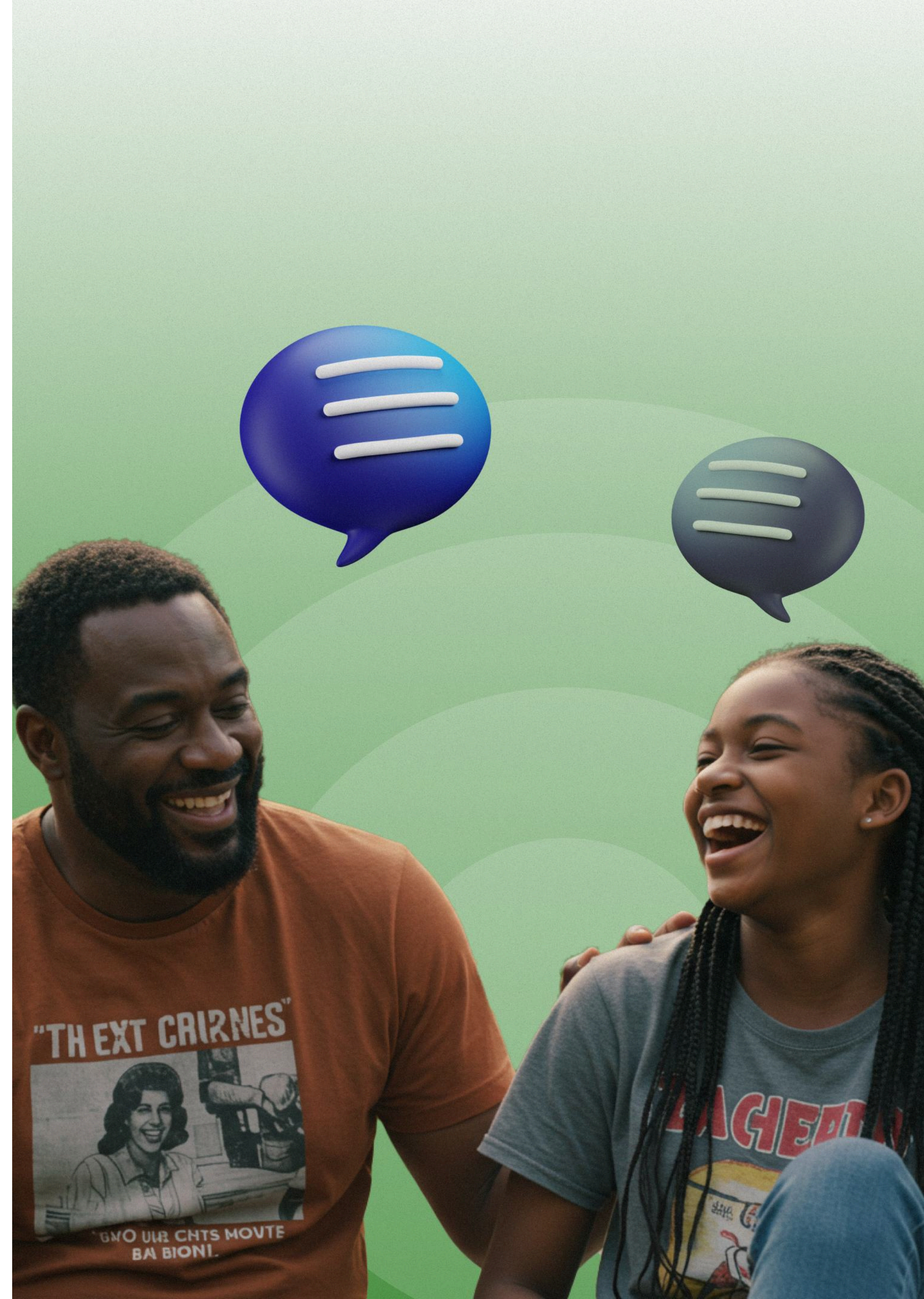
Teens : Curated feedback, fact-checking layer, digital literacy cues

Teachers: Safe content generation + citation support

Implement:

- Content moderation filters (toxicity, bias, hallucination)
- Topic restrictions (e.g., block violent, sexual, or political topics)

Test outputs with students using controlled pilots and safety reviewers.



## 6. Empower Educators and Parents with Controls

Allow visibility into AI-generated content.

Provide dashboards for:

- Reviewing flagged interactions
- Setting learning goals
- Restricting time or topics

Let teachers turn on/off certain features (e.g., open chat vs. fixed prompts).

**Design interfaces that give adults agency and students autonomy within safe bounds.**



## 7. Enable Real-Time Safety Systems

Use classifiers to detect:

- Bullying
- Self-harm signals
- Unsafe input (e.g., “How to cheat”)
- Discriminatory or hateful speech

Respond with:

- Auto-block or rephrase
- Escalation to a human moderator
- Educational redirects (e.g., “Let’s talk about integrity in schoolwork”)

**Build in a real-time feedback & reporting button.**



## 8. Monitor, Audit, and Evolve

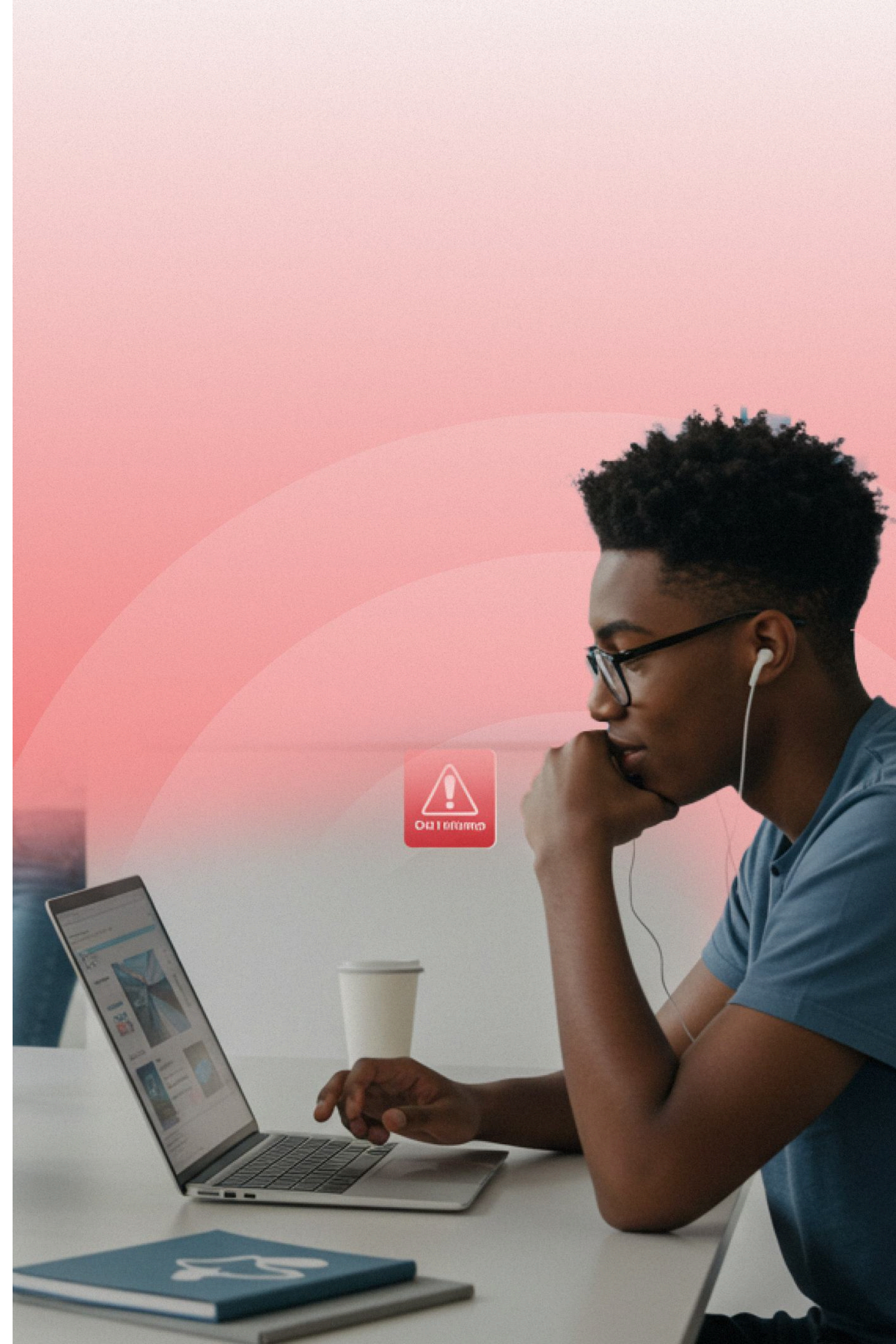
Maintain a Safety Dashboard of :

- Flagged outputs
- Type of misuse
- Feedback from users (students, teachers, parents)

Conduct regular:

- Red-teaming exercises
- Teacher/student feedback loops
- Safety audits with external experts

**Improve model behavior over time, informed by real-world use.**



## 9. Promote Digital Literacy and Well-being

Use the LLM to:

- Teach about AI limitations ("I am an assistant, not a teacher.")
- Encourage critical thinking
- Promote healthy tech habits (e.g., take breaks, question sources)

**Safety isn't just technical, it's educational.**





**Safety by Design is not a checkbox; it is a culture.**

**TrustedTech Africa is your partner to ensure that your EdTech product equips every child to learn, explore, and grow without harm.**

## Contact:

 [info@trustedtechafrica.com](mailto:info@trustedtechafrica.com)

 [trustedtechafrica.com](https://trustedtechafrica.com)