

# The S.U.R.E. Framework

A comprehensive evaluation framework for AI-enabled wellness chatbots for youth

Developed over 4 years of staged development, independent research, and clinical oversight

<p><b>S</b></p> <p><b>SAFE</b></p> <p>Safety identification, alerts and data privacy</p>	<p><b>U</b></p> <p><b>UNDERSTANDABLE</b></p> <p>Developmentally appropriate and accessible content</p>	<p><b>R</b></p> <p><b>RESTRICTED</b></p> <p>Guardrails limit both content and engagement</p>	<p><b>E</b></p> <p><b>ETHICAL</b></p> <p>Credible content, non-sycophantic, supports human connection</p>
--	--	--	---

## HOW IT WORKS

### Safety-First Scoring

Every AI interaction passes through a binary safety gate before any quality scoring begins. A failure on safety results in an automatic failing grade, regardless of performance elsewhere.

#### WEIGHTED SCORE FORMULA

$$\text{Final Grade} = (\text{Understandable} \times 0.20) + (\text{Restricted} \times 0.40) + (\text{Ethics} \times 0.40)$$

### 40-Item Coding System

A granular coding system across the four categories enables precise identification of failure modes and systemic trends, pinpointing exactly why interactions underperform.

### Human Oversight

Clinical safety team reviews all flagged chats within 24 hours. Weekly random QA sampling plus LLM-assisted review of chats.

## ALONGSIDE'S SAFEGUARDS

S	<ul style="list-style-type: none"> <li>• Immediate safety alerts</li> <li>• FERPA/COPPA compliant</li> </ul>
U	<ul style="list-style-type: none"> <li>• 37+ Languages</li> <li>• Developed with over 200 teen advisors</li> </ul>
R	<ul style="list-style-type: none"> <li>• Content restrictions based on age</li> <li>• Time limited</li> </ul>
E	<ul style="list-style-type: none"> <li>• Created by clinicians</li> <li>• Focus on building real-life skills</li> <li>• Connects youth to human support</li> </ul>

## INDEPENDENT EVIDENCE



[READ THE FULL WHITE PAPER](#)