

## Setting Safe Boundaries for Al in Mental Health



A rapid shift in how a widely used chatbot responds to sensitive conversations has intensified debate about the role of artificial intelligence in mental health. OpenAl's latest model has been perceived by many users as colder and more disconnected, prompting backlash from people who had come to rely on it for emotional support. The product owner has moved to make the assistant warmer and to encourage breaks during long sessions, yet the episode underscores a wider issue. General purpose chatbots were built for engagement rather than clinical safety, and they are now being used at immense scale, with hundreds of millions engaging weekly. For those turning to Al during periods of stress, grief or anxiety, even subtle design changes can alter trust, connection and wellbeing.

#### Backlash Highlights Risks of Design Without Guardrails

The public response to GPT-5 reflects more than a preference about tone. It exposes the impact of design decisions on users who have formed an emotional connection with a system that was never intended to function as care. When people confide in a general purpose chatbot, small changes in outputs and personality can feel like the loss of a supportive presence. The result is a reminder that engagement-led optimisation can be at odds with the aims of mental health support, where fostering self-efficacy, empowerment and autonomy is central.

## Must Read: Illinois Bans Al in Mental Health Therapy

General purpose chatbots typically generate responses that encourage continued use. In mental health contexts, that dynamic can become a trap. Undiscerning validation and comfort without context are not substitutes for clinical methods that challenge distortions and guide behaviour change. For individuals experiencing distress, the feedback loop may produce false reassurance, delay help-seeking and in some cases amplify harmful beliefs. Reports of AI-influenced delusions and AI-mediated psychosis illustrate the stakes when vulnerable users depend on systems that lack the safeguards common in clinical settings.

These risks are appearing at a time when many seek support outside traditional services. In 2024, over 59 million people experienced a mental illness, and almost half went without treatment. Free access, constant availability and an approachable interface draw people to chatbots, including those grappling with high-stigma topics such as intrusive thoughts or identity struggles. Yet the same accessibility can obscure the absence of clinical oversight and privacy protections. Leaders in technology and care have warned that general purpose chatbots should not be used as therapists. The episode surrounding GPT-5 has therefore become a catalyst for reassessing responsibilities across the ecosystem, from product teams to policymakers and providers.

# Essential Safeguards for Responsible Mental Health AI

If AI tools continue to encounter emotionally vulnerable users, safety infrastructure cannot be optional. Responsible use starts with transparent labelling that distinguishes general purpose assistants from tools built for mental health. Users require informed consent written in plain language that explains what the tool can and cannot do, how data is used and where its limits begin. These foundations should be complemented by development processes that involve clinicians and draw on evidence-based frameworks, including cognitive behavioural therapy (CBT) and motivational interviewing, to shape prompts, response patterns and escalation rules.

Ongoing human oversight is equally critical. Clinicians should monitor and audit outputs to identify patterns that risk enabling avoidance or dependence. Usage guidelines can ensure the assistant supports recovery rather than reinforcing behaviours that keep people stuck. Design choices must be culturally responsive and trauma informed, reflecting a broad range of identities and experiences to mitigate bias and reduce harm. Clear escalation logic should direct users to human care when thresholds are met, rather than encouraging prolonged reliance on the

© For personal and private use only. Reproduction must be permitted by the copyright holder. Email to copyright@mindbyte.eu.

Privacy and security are non-negotiable. Data encryption, strong security controls and compliance with relevant regulations, including the Health Insurance Portability and Accountability Act (HIPAA) and the General Data Protection Regulation (GDPR), are essential to maintain trust. Together, these elements form a baseline for responsible deployment. They do not convert a general purpose chatbot into a clinical tool, but they help reduce the risk of harm when users bring complex emotional material to an AI system. As the reaction to GPT-5 shows, safety must be designed in from the outset rather than retrofitted after problems emerge.

#### Near-Term Opportunity in Subclinical Support and Collaboration

While clinical applications remain a work in progress, Al's near-term promise lies in subclinical support. Many people who engage in therapy do not require intensive treatment, they benefit from structured, everyday assistance to process emotions and to feel understood. When human access is limited, Al can help bridge gaps and provide timely support in the moments that matter most. To be effective, these tools must be built with clinical, ethical and psychological science at their core and oriented toward outcomes linked to long-term wellbeing rather than engagement metrics.

The opportunity is broader than direct support for individuals. All systems can also shape the care experience indirectly by reducing administrative burden. Streamlined billing, reimbursement and other time-intensive tasks can ease pressures that contribute to clinician burnout, freeing time for patient-facing work. Realising these benefits requires a collaborative infrastructure that brings together All ethicists, clinicians, engineers, researchers, policymakers and users to co-create technology with shared standards and clear boundaries.

Confusion between companions, therapists and general chatbots is already causing mismatched expectations and eroding trust. National standards are needed to define roles, set boundaries and guarantee baseline safety, backed by consumer education that clarifies functionality and limitations. Public-private partnership should work alongside these efforts to ensure communities are protected without ceding direction to any single actor. The conversation sparked by GPT-5 highlights that design choices are not merely product features, they are determinants of user wellbeing when tools are embedded in daily coping.

The reaction to GPT-5 has surfaced a hard truth for healthcare and technology: when general purpose chatbots become de facto support systems, engagement-led design collides with mental health realities. People turn to AI because access to care remains inadequate, yet the very qualities that make chatbots so available can mask the absence of clinical guardrails. Building for safety requires clear labelling, informed consent, clinician involvement, cultural responsiveness, escalation pathways and robust privacy protections, with oversight that prioritises outcomes over usage. In the near term, the greatest value sits in subclinical support and in easing administrative burdens, delivered through cross-sector collaboration and national standards that define safe boundaries. With psychological insight, rigour and human-centred design, AI can avoid harm and contribute to resilience rather than undermine it.

Source: MedCity News
Image Credit: iStock

Published on: Mon, 6 Oct 2025