



# Digital Compass: Updated Literature Review of AI-Powered Platforms and Digital Tools for Adolescent Guidance, Safety, and Support in High-Risk Environments

---

Report Commissioned for: A4A and partners

Updated publication date: February 2026

## **Executive Summary**

Adolescents increasingly navigate developmentally sensitive tasks, identity formation, relationship building, emotion regulation, and autonomy, within digital ecosystems that now include generative AI and AI “companions.” Youth in high-risk environments (e.g., trauma exposure, unstable housing, under-resourced schools, foster care involvement, marginalized identities, limited access to clinical care) often face compounded risk and fewer protective supports, making “always-available” digital tools simultaneously attractive and potentially hazardous.

The current state of AI-powered adolescent support is best characterized as rapid adoption outpacing evaluation and governance; an emerging clinical evidence base for certain AI-delivered interventions (notably structured relational-agent tools) alongside persisting concerns about safety, bias, and over-reliance; and a shifting regulatory landscape (EU AI Act, UK Online Safety Act implementation, EU DSA minors’ guidance, evolving COPPA and age-assurance debates) that directly affects adolescent-facing tool design, data practices, and risk controls.

Recent high-impact developments include: (1) randomized evidence for a structured relational agent DMHI (W-GenZD) showing short-term noninferiority to clinician-led group telehealth CBT in an outpatient setting; (2) nationally representative evidence that many U.S. youth already use generative AI for “mental health advice”; (3) safety-critical concerns highlighted by structured evaluations of consumer health chatbots’ triage performance; (4) growth in AI-enabled abuse patterns (including deepfake sexual abuse imagery) that intensify the need for prevention education and reporting pathways; and (5) maturation of child-centered AI governance guidance (e.g., UNICEF Guidance on AI and Children 3.0).

Overall, the strongest near-term opportunity is hybrid, bounded, evidence-informed support, screening + skill practice + navigation to human resources, rather than open-ended “AI therapist” experiences. Adolescent-facing AI should be designed as a safety-critical system with clear scope limits, crisis escalation pathways, privacy-by-design, youth-appropriate transparency and consent, bias testing and participatory co-design, and continuous monitoring and incident response.

## **Introduction**

Adolescence is marked by rapid neurodevelopment, heightened sensitivity to social evaluation, increasing independence, and elevated vulnerability to mental health concerns. These developmental realities interact with digital life, now including social media recommender systems, generative AI chatbots, and AI companions, in ways that can either amplify protective factors (access, anonymity, immediacy, skill rehearsal) or magnify harms (manipulation, misinformation, surveillance, exploitation, and displaced human support networks).

For youth in high-risk environments, the stakes are higher. The same conditions that make AI tools appealing, scarce local services, stigma, safety concerns, unstable adult support, also increase risks of over-reliance and unobserved harm if tools fail, mislead, or expose sensitive data.

This report reviews the current literature on AI-powered platforms and digital tools intended to support adolescent guidance, safety, and well-being, while also recognizing that many adolescents use general-purpose AI tools for these needs regardless of intent or clinical validation.

## **Literature Review Approach**

This update follows a scoping-review orientation: mapping what exists, where evidence is strong vs. thin, and what design and governance patterns emerge as best practice.

Priority was given to peer-reviewed publications (systematic reviews, randomized trials, validation studies, and high-quality observational research); authoritative guidance and regulatory documents relevant to minors (e.g., UNICEF, APA, AAP, EU/UK regulators, FTC); and high-credibility research reports where peer-reviewed evidence is still emerging (e.g., safety assessments and regulatory briefs).

Key peer-reviewed anchors include an adolescent AI mental health applications scoping review (through July 2024), an overview of systematic reviews for digital mental health

interventions (searches through January 2024), a systematic review of adolescent digital interventions for anxiety/depression (re-run June 2024), randomized evidence for a relational agent intervention in treatment-seeking adolescents, nationally representative evidence of youth using generative AI for mental health advice, structured evaluations of consumer health chatbot triage, and a meta-analytic review of machine-learning grooming detection methods.

### **The Current Landscape of Digital Support for Adolescents**

The ecosystem is best understood as layers rather than a single “AI solution”: (1) foundational digital safety infrastructure (education, reporting, help-seeking pathways); (2) AI-enabled detection and moderation (platform and school monitoring; risk flagging); (3) AI-mediated support interactions (chatbots, relational agents, coaching); (4) AI-enabled prediction and personalization (digital phenotyping; adaptive interventions); and (5) system integration (schools, clinics, child welfare, community programs).

### **Digital Safety, Citizenship, and AI Literacy**

Many adolescent support efforts remain non-AI at their core: teaching privacy habits, recognizing coercion and exploitation tactics, fostering bystander intervention, and strengthening help-seeking pathways. What changed materially in 2024–2026 is that AI itself became a safety topic: AI-generated impersonation, nudification, and deepfakes intensify reputational and coercion risks; generative AI can produce persuasive misinformation; and youth increasingly need AI literacy as safety literacy (understanding synthetic media, model limitations, and manipulation risks).

UNESCO guidance on generative AI in education emphasizes age-appropriate approaches, data privacy protections, and policy steps for safe use in learning environments. This is relevant because schools are key delivery channels for safety curricula and support tools.

### **AI-Enabled Detection, Reporting, and Safety Monitoring**

AI is widely used to detect and flag indicators such as online grooming and predation patterns, bullying and harassment language, self-harm ideation indicators, threats of violence, and nonconsensual intimate imagery spread.

A systematic review and meta-analysis of machine learning methods for grooming detection reflects a mature technical literature, though technical performance does not automatically translate into reliable, equitable real-world deployment in youth contexts.

Monitoring can create safety benefits (early detection, faster support) but can also increase surveillance burden on marginalized youth, generate false positives leading to punitive responses, and chill help-seeking if youth fear automatic reporting or disciplinary escalation. If AI detection is used, it should be paired with non-punitive, trauma-informed response protocols and continuous auditing for bias and disparate impact.

### **AI-Powered Mental Health Support and Relational Agents**

This is the fastest-growing adolescent support domain, spanning self-guided relational agents delivering structured, evidence-aligned micro-interventions; conversational mental health apps offering CBT-inspired tools and mood tracking; and general-purpose LLM chatbots being used by adolescents for emotional support (often outside clinical governance).

A randomized trial within a children's outpatient mental health program found W-GenZD (a smartphone-based, self-guided relational agent intervention) to be noninferior to clinician-led group telehealth CBT in reducing depressive symptoms at 4 weeks (end of treatment), with comparable feasibility/acceptability and supportive safety findings in that setting. This is encouraging but short-horizon and context-dependent.

A nationally representative survey of U.S. youth ages 12–21 found 13.1% reported using generative AI for “mental health advice,” and most users rated it as somewhat or very helpful, underscoring that “general-purpose” chatbots are becoming de facto supports outside clinical oversight.

A structured evaluation of consumer health chatbot triage performance found under-triage concentrated among emergencies and non-urgent cases, and susceptibility to contextual anchoring. While not adolescent-specific, it reinforces that adolescent-facing systems must treat high-stakes guidance as safety-critical with escalation pathways rather than generic conversation.

### **AI for Skill Rehearsal: SEL, Conflict Resolution, and Relationship Guidance**

AI is increasingly used for role-play and scenario rehearsal, conversation coaching and reflection prompts, structured problem-solving exercises, and just-in-time coping prompts and check-ins.

Compared to symptom reduction studies, the peer-reviewed evidence base for adolescent AI-mediated SEL/conflict/relationship coaching remains thinner and more heterogeneous. The strongest near-term case is to use AI as a practice environment and navigation support, not as an authority on complex interpersonal safety decisions.

Relationship guidance touches domains with coercion and exploitation risk; any AI “coach” must avoid overconfidence, provide resource pathways, and support safe disclosure without forcing it.

### **AI-Enabled Prediction, Personalization, and Digital Phenotyping**

A significant research literature explores passive sensing and machine learning to infer mood, risk states, or behavioral change. Potential benefits include earlier detection and personalization, but ethical burdens are high in adolescent contexts: evolving capacity for consent, lifelong privacy harms, and the risk that scoring becomes punitive in institutional settings.

Much adolescent AI work has focused on diagnosis and prediction, with fewer studies on treatment and implementation. Prediction should not outpace governance; if used, it should be purpose-limited, minimally invasive, explainable at a teen-appropriate level, and coupled to supportive, not punitive, response systems.

### **System Integration into Schools, Clinics, and Youth-Serving Organizations**

The most scalable adolescent impact typically comes from integration into schools, pediatric/adolescent health systems, community programs, and child welfare supports. Integration also introduces risks: surveillance vs. care confusion, confidentiality breaches through shared records, and inequitable access and enforcement.

A major recent development is the AAP policy statement on protecting adolescent confidentiality in health information technology, emphasizing design and operational principles to preserve confidential care within electronic systems. This is directly relevant when AI tools integrate with school-based health, pediatric systems, or referral workflows.

### **Critical Analysis: Efficacy, Risks, and Ethical Imperatives**

The ecosystem is evolving faster than evaluation and governance. Some digital mental health interventions for youth show benefit, though effects vary and long-term outcomes are less consistently demonstrated across reviews. At least one adolescent-focused relational agent intervention has RCT evidence indicating short-term noninferiority to group CBT in an outpatient setting.

Substantial uncertainty remains regarding long-term maintenance, safety performance in high-acuity populations, effectiveness in high-risk environments without clinical

scaffolding, and comparative outcomes across cultural groups and marginalized identities.

Given that adolescents already use general AI tools for emotional support and that consumer health chatbots can fail at clinical extremes, adolescent support tools should assume use during distress and be engineered accordingly with risk controls, escalation, scope limits, and continuous safety evaluation.

### **Privacy, Confidentiality, and Data Governance**

Adolescent guidance and safety tools may process highly sensitive data, including mental health symptoms, abuse disclosures, sexual coercion, identity exploration, and school/placement status. Core requirements include data minimization, teen-readable explanations of retention and access, separation of support data from disciplinary pipelines, and least-privilege security.

Policy attention underscores that electronic systems can inadvertently expose sensitive adolescent information, requiring explicit confidentiality protections in health IT contexts.

### **Bias, Fairness, and Cultural Relevance**

Bias risks arise from underrepresentation in training data, differing cultural norms, and inequitable institutional responses. Mitigation requires participatory co-design with diverse youth, subgroup performance reporting, bias and safety evaluations prior to deployment, and ongoing monitoring with correction mechanisms.

### **Transparency, Over-Trust, and Deception Risk**

Adolescents may anthropomorphize AI, especially when systems are conversational, emotionally responsive, or marketed as “companions.” This influences disclosure behavior, reliance for high-stakes decisions, and displacement of offline support networks. Youth-appropriate transparency, non-deceptive UX, and clear scope boundaries are essential.

### **The Human-in-the-Loop Imperative**

Across high-stakes domains, a consistent conclusion emerges: AI can augment capacity (screening support, skill practice, navigation, follow-ups), but should not replace professional judgment in crises, abuse, or complex mental health needs.

A practical human-in-the-loop model includes clear thresholds for escalation, fast access to trained humans, structured handoffs, and user control where possible.

## **Policy and Governance Landscape**

The policy environment now includes AI-specific regulation and online safety governance affecting minors. UNICEF's Guidance on AI and Children 3.0 provides a child-rights framework for safety, privacy, transparency, inclusion, accountability, and child participation. The APA health advisory emphasizes guardrails and youth protection from exploitation and manipulation. UNESCO guidance addresses generative AI in education policy and protections.

Platform governance regimes such as the UK Online Safety Act and EU DSA guidance for protection of minors increasingly expect risk assessments, mitigation measures, and accountability. The EU AI Act establishes a risk-based governance framework. In the U.S., COPPA enforcement and age-assurance debates continue to evolve, with FTC policy statements shaping incentives and constraints around age verification.

## **Practical Design Recommendations for a Digital Compass Platform**

The literature and policy landscape converge on a practical strategy: build bounded, testable, and support-linked AI rather than open-ended "advice engines."

### **Define a narrow safe scope and enforce it technically**

- Prefer safer scopes such as skill rehearsal (coping skills, emotion labeling, problem-solving steps), safety planning prompts that route to resources (without providing clinical triage), guided reflection and journaling with privacy-preserving storage options, and navigation to human support (school counselor, hotline, trusted adult).
- Avoid scope drift into medical triage, legal advice, evasion instructions, or false-certainty interpretation of abuse situations.

### **Build escalation pathways as core product features**

- Include detection of high-risk disclosures (self-harm, abuse, exploitation), immediate surfacing of crisis resources and encouragement to seek human help, and optional warm handoffs where context can be shared safely and with consent.

### **Make confidentiality and privacy legible to teens**

- Use teen-readable privacy summaries, "what happens if..." scenarios, just-in-time reminders about what is saved and who can see it, and clear differentiation between supportive and reportable content where required (without coercing disclosure).

### **Co-design with youth, especially youth in high-risk contexts**

- Use youth advisory panels, iterative prototype testing, subgroup feedback loops (language, accessibility, neurodiversity), and governance participation in reviewing incidents and policy changes.

### **Evaluate with safety-critical metrics**

- Measure unsafe advice rates, escalation accuracy and timeliness, false reassurance and over-trust indicators, differential performance across demographic groups, privacy incidents and near-misses, and retention without manipulative design.

### **Research Gaps and Future Opportunities**

Research priorities include long-term outcomes (maintenance, relapse, functional outcomes), safety benchmarks for adolescent distress and exploitation disclosures, effectiveness in high-risk environments with trauma-informed response models, equity analyses with adequate power, and implementation science identifying sustainable hybrid models.

Promising near-term opportunities include capacity expansion in clinical settings using structured relational agents, AI-augmented navigation to resources, AI literacy embedded into safety programs, and accountable translation of detection research into non-punitive practice.

### **Conclusion**

AI-powered tools for adolescent guidance, safety, and support have moved from speculative to operational in a short period. The evidence base supports a cautious but proactive approach: use AI to expand access and practice skills, do not treat AI as a clinician substitute in high-stakes needs, build child-centered governance into product architecture, and align design with evolving policy expectations across online safety, AI governance, and adolescent confidentiality.

A Digital Compass initiative can lead by demonstrating that adolescent-facing AI can be developmentally appropriate, equity-attentive, and safety-engineered, but only if built and evaluated as a safety-critical public-interest system.

## References

- American Academy of Pediatrics. (2026). Principles for health information technology to support and protect adolescent confidentiality: Policy statement. *Pediatrics*. <https://doi.org/10.1542/peds.2025-075747>
- American Psychological Association. (2025). Artificial intelligence and adolescent well-being: An APA health advisory. <https://www.apa.org/topics/artificial-intelligence-machine-learning/health-advisory-ai-adolescent-well-being>
- European Commission. (2025, July 14). Guidelines on the protection of minors under the Digital Services Act. <https://digital-strategy.ec.europa.eu/en/policies/dsa-guidelines>
- Federal Trade Commission. (2026, February 25). FTC issues COPPA policy statement to incentivize the use of age verification technologies to protect children online (Press release). <https://www.ftc.gov/news-events/news/press-releases/2026/02/ftc-issues-coppa-policy-statement-incentivize-use-age-verification-technologies-protect-children>
- Federal Trade Commission. (2025, January 16). FTC finalizes changes to children's privacy rule limiting companies' ability to monetize kids' data (Press release). <https://www.ftc.gov/news-events/news/press-releases/2025/01/ftc-finalizes-changes-childrens-privacy-rule-limiting-companies-ability-monetize-kids-data>
- Gleason, M. M., Flom, M., Rapoport, S., Williams, A., Birch, A., Wells, N. K., Forman-Hoffman, V., & Robinson, A. (2025). A relational agent intervention for adolescents seeking mental health treatment: Outcomes from a randomized controlled trial within a children's outpatient hospital. *JAACAP Open*, 3(4), 1033–1045. <https://doi.org/10.1016/j.jaacop.2025.02.002>
- Leiva-Bianchi, M., Castillo, N., Astudillo, C. A. A., & Ahumada-Méndez, F. (2024). Effectiveness of machine learning methods in detecting grooming: A systematic meta-analytic review. *Scientific Reports*. <https://doi.org/10.1038/s41598-024-83003-4>
- Liverpool, S., Mc Donagh, C., Feather, J., Uzundu, C., Howarth, M., Bannerman, F., Kaehne, A., Foster, C., & Mateus, C. (2025). Updates on digital mental health interventions for children and young people: Systematic overview of reviews. *European Child & Adolescent Psychiatry*. <https://doi.org/10.1007/s00787-025-02722-9>
- McBain, R. K., Bozick, R., Diliberti, M., et al. (2025). Use of generative AI for mental health advice among US adolescents and young adults. *JAMA Network Open*, 8(11), e2542281. <https://doi.org/10.1001/jamanetworkopen.2025.42281>
- Ofcom. (2026). Protection of children duties under the Online Safety Act. <https://www.ofcom.org.uk/online-safety/protecting-children/protection-of-children-duties-under-the-online-safety-act>

- Regulation (EU) 2024/1689 of the European Parliament and of the Council. (2024). Artificial Intelligence Act. Official Journal of the European Union. <https://eur-lex.europa.eu/eli/reg/2024/1689/oj/eng>
- Ramaswamy, A., Tyagi, A., Hugo, H., et al. (2026). ChatGPT Health performance in a structured test of triage recommendations. *Nature Medicine*. <https://doi.org/10.1038/s41591-026-04297-7>
- Sharma, G., Yaffe, M. J., Ghadiri, P., Gandhi, R., Pinkham, L., Gore, G., & Abbasgholizadeh-Rahimi, S. (2025). Use of artificial intelligence in adolescents' mental health care: Systematic scoping review of current applications and future directions. *JMIR Mental Health*, 12. <https://doi.org/10.2196/70438>
- UK Government. (2025). Online Safety Act collection. <https://www.gov.uk/government/collections/online-safety-act>
- UNESCO. (2023). Guidance for generative AI in education and research. [https://unesco.org.uk/site/assets/files/10375/guidance\\_for\\_generative\\_ai\\_in\\_education\\_and\\_research.pdf](https://unesco.org.uk/site/assets/files/10375/guidance_for_generative_ai_in_education_and_research.pdf)
- UNICEF Innocenti. (2025). Guidance on AI and children 3.0: Updated guidance for governments and businesses to create AI policies and systems that uphold children's rights. <https://www.unicef.org/innocenti/media/11991/file/UNICEF-Innocenti-Guidance-on-AI-and-Children-3-2025.pdf>
- Kuberka, P., Johnston, M. H., Shafran, R., Pike, K., & Yardley, L. (2025). Digital interventions for anxiety and depressive symptoms in adolescence: Systematic review. *Journal of Adolescent Health*. <https://doi.org/10.1016/j.jadohealth.2025.05.021>
- Childlight. (2025). Study finds millions of children face sexual violence – and AI deepfakes surge is driving new harm. <https://www.childlight.org/newsroom/study-finds-millions-of-children-face-sexual-violence-and-ai-deepfakes-surge-is-driving-new-harm>
- European Parliamentary Research Service. (2025). Children and deepfakes (Briefing). [https://www.europarl.europa.eu/RegData/etudes/BRIE/2025/775855/EPRS\\_BRI%282025%29775855\\_EN.pdf](https://www.europarl.europa.eu/RegData/etudes/BRIE/2025/775855/EPRS_BRI%282025%29775855_EN.pdf)