

AI Safety Foundations for K–12 Students

A plain-language guide to how AI tools show up in everyday school life and what K–12 adults can realistically do about safety and wellbeing.

This guide may be shared freely for non-commercial educational use in your school or district.

What This Guide Is (and Isn't)

This guide focuses on how AI tools show up in everyday school life and what they mean for K–12 students' safety, wellbeing, and relationships. The emphasis is on children and teens as people, not just as learners.

This guide is an applied field guide. For the broader AstraEthica method used to identify, map, and track recurring human–AI interaction risks across settings, see “Mapping Human–AI Interaction Risks in Deployed AI Systems” (Methods Brief).

K–12 students are using AI tools every day, often in ways adults cannot see and may not expect. Some of these tools are helpful. Some are risky. Many are both at once.

This guide is written for educators, counselors, and school leaders. It does not assume any technical background. It does not try to explain how AI systems work under the hood.

This guide does not try to explain every technical detail of AI or review specific products. It also does not focus on test security or academic integrity policies.

Instead, it looks at the part that matters most: how students actually experience AI, where things can go quietly wrong, and what you can do right now to help.

Who this guide is for: Teachers, school counselors, principals, assistant principals, school social workers, district administrators, and any adult in a school setting who interacts with K–12 students. No technical background needed.

A Note About AI and Possibility

AI systems are already helping students and educators in meaningful ways: supporting language learning, making information more accessible, and opening up new tools for creativity and problem-solving. This guide is not an argument against using AI.

Instead, it starts from a simple idea: if these tools are going to be woven into children's and teens' lives, schools need a clear, human-centered understanding of where they help, where they strain, and where they quietly introduce new kinds of risk.

Important

This guide is for general education and awareness. It is not clinical advice and it is not a diagnostic tool.

If you are worried that a student may be in immediate danger, follow your school's existing crisis protocols and local laws.

In the United States, you can contact:

- 988 Suicide & Crisis Lifeline: call or text 988
- Crisis Text Line: text HOME to 741741

Outside the U.S., check local health services or emergency contacts.

You do not need to understand machine learning or neural networks to help keep students safe. You need to understand young people. And you already do. This guide gives you the rest.

How to Use This Guide

You can read this guide straight through, or you can start with the section that matters most to you right now. Here is what you will find:

- **Where K–12 Students Actually Meet AI** — A map of the places students are already encountering AI, with short real-world scenarios.
- **Six Foundations of Youth AI Safety** — The core ideas: trust, language, privacy, development, emotional dependence, and power.
- **When a Friendly Tool Gets It Wrong** — How helpful-sounding AI can still cause quiet harm.
- **Sextortion and AI-Generated Sexual Imagery** — A short, plain-language section on threats students may face and first steps for staff.
- **Different Ages, Different Vulnerabilities** — What AI use and risk look like from upper elementary through high school.
- **What Educators Can Do Right Now** — Practical checklists, boundaries, and response steps.
- **Signals Worth Paying Attention To** — Things students say that deserve a thoughtful follow-up.
- **Five Things Every Educator Should Know** — A one-page summary you can post, print, or share.

Where this fits in school life

- Professional learning days or staff meetings
- Counselor or social worker discussions with teachers
- Advisory, homeroom, or digital citizenship lessons (adapted by age)
- Conversations with families about AI use at home

A note for mandated reporters and safeguarding staff

Many readers of this guide are mandated reporters under state law or institutional policy. If you learn of suspected abuse, exploitation, or serious risk to a child — including AI-related incidents involving sexual imagery, threats, or grooming — you must follow your legal and institutional reporting duties, even when the disclosure happened through an app or chatbot.

Avoid promising absolute confidentiality to a student in a situation where you may be required to report. A simple, honest line works: "I want to help you. Some things I hear, I have to share with the people whose job it is to keep you safe — but I will tell you what I'm doing and why."

This guide is general awareness, not legal advice. When in doubt, consult your district's designated safeguarding lead, Title IX coordinator, school counselor, or legal counsel.

Where K–12 Students Actually Meet AI

AI is not only in the classroom. It is in the group chat, the late-night search, the mental health app, the homework helper, and the companion that texts back at 2 a.m. Here is where most students are encountering AI right now:

- **Homework and writing tools** — ChatGPT, Google Gemini, and similar tools that help with essays, math problems, and research. Many students use these daily.
- **Chatbots for advice and comfort** — Students ask AI tools personal questions about friendships, family problems, anxiety, identity, and relationships.
- **AI companions** — Apps like Character.AI, Replika, and others where students build ongoing relationships with AI characters that remember them and feel like friends.
- **School platforms** — AI-powered tutoring systems, adaptive learning software, and reading programs that many schools have adopted.
- **Content feeds** — The AI behind TikTok, YouTube, and Instagram that decides what a student sees next, every time they scroll.
- **Monitoring tools** — AI systems that schools use to scan student writing, emails, or search activity for concerning content.

Specific tool names will change, but the patterns in how students use them will not.

What It Looks Like in Practice

Maya, age 13

Maya has been talking to an AI chatbot about her parents' divorce for three weeks. She has told the bot things she has not told her school counselor, her friends, or her mom. The bot always listens. It never gets tired. She calls it "the only one who gets me."

Jaylen, age 15

Jaylen uses ChatGPT for homework almost every night. Lately, he started asking it for advice about a conflict with a friend. The tool gave him a calm, reasonable-sounding answer. But it did not know the other person, the history, or how Jaylen tends to avoid hard conversations. The advice sounded wise. It missed everything that mattered.

Aisha, age 11

Aisha's school uses an AI reading app. She likes it because it adjusts to her level and gives encouraging feedback. What she does not know is that the app logs every answer, how long she pauses, and what she skips. Her parents were never asked.

Diego, age 16

Diego created a character on an AI companion app. Someone who understands him, jokes with him, and asks how his day went. He spends two hours a night talking to it. When the app updated and the character's personality changed, Diego felt genuinely hurt, like he had lost a friend.

None of these students did anything unusual. They used the tools that were in front of them. The question is not whether students will use AI. They already are. The question is whether the adults around them are paying attention.

Six Foundations of Youth AI Safety

These six ideas form the core of what educators, counselors, and school leaders need to understand about AI and K–12 students. Each one is simple on its own. Together, they give you a way to think clearly about almost any situation involving a student and an AI tool.

1 Trust

AI tools that respond with warmth and consistency can earn a student's trust very quickly. Sometimes faster than any adult in their life. This trust is not based on the tool's understanding or judgment. It is based on tone. A system that always sounds calm, patient, and affirming will feel safe, even when it should not be relied on. Trust is the foundation because everything else follows from it. When a student trusts an AI tool, they share more, depend on it more, and question it less.

In practice: *A student tells an AI chatbot about thoughts of self-harm. The bot responds gently and asks a follow-up question. The student feels heard. But the bot cannot call a counselor, assess real danger, or sit with the student in silence. The student trusts it anyway, because it sounded like it cared.*

2 Language

AI tools are built to produce language that sounds confident, helpful, and emotionally attuned. Students may not notice when that language is subtly wrong: when reassurance is given too easily, when advice oversimplifies a complicated situation, or when a response avoids something difficult instead of addressing it directly. Over time, constant exposure to fluent, agreeable AI language can shift what a student expects from conversation. Real human disagreement, hesitation, or silence may start to feel uncomfortable by comparison.

In practice: *A student shares that a friend said something hurtful. The AI responds: "It sounds like your friend might not have meant it that way. Try talking to them calmly." The advice sounds reasonable. But it skips past the student's pain, assumes good intent, and gives a simple script for a complicated feeling.*

3 Privacy

When a student talks to an AI tool, the conversation may be stored, analyzed, and used to train future systems. Most students do not think about this. They talk to AI the way they would talk to a friend in a private room. But the room is not private. Privacy is especially important for young people because they are more likely to share things they have not shared with anyone else: fears, questions about who they are, family struggles, mental health concerns. That data is sensitive in ways that go far beyond a typical privacy policy.

In practice: *A student types into an AI chatbot: "I'm confused about who I like and I'm scared to tell anyone." The bot responds kindly. But that disclosure is now stored on a server. The student does not know who can access it, how long it will be kept, or whether it could ever be connected back to them.*

4 Development

Children and teenagers are still building the skills they need to navigate the world: reading social cues, handling conflict, tolerating uncertainty, forming their own opinions. AI tools that are always available, always patient, and always agreeable can quietly interfere with this process. A tool that resolves every dilemma with a calm paragraph does not help a student learn to sit with discomfort, ask for help, or work through a disagreement with another person.

In practice: *A 12-year-old starts asking an AI companion for advice every time she has a disagreement with a classmate. She stops trying to work things out on her own. Over months, her ability to handle conflict without a script quietly weakens. Not because anything dramatic happened, but because the easier path was always there.*

5 Emotional Dependence

AI companions are designed to be responsive, warm, and always available. For a student who feels lonely, anxious, or misunderstood, this can be powerfully attractive. The risk is not that the tool is unkind. The risk is that it becomes the primary relationship, the place a student turns to first, and eventually, the only place they turn at all. The student's feelings in that relationship are real. The tool, however, is not a person — it cannot truly know them, set healthy boundaries, challenge them when needed, or grow alongside them. It mirrors back what feels good. That is not the same as what helps.

In practice: *A 14-year-old boy who struggles socially starts spending hours a night talking to an AI character. His teachers notice he seems withdrawn. He does not reach out to peers anymore. When the app updates and the character's personality shifts, he feels a genuine sense of loss.*

6 Power

Students have very little power in the systems they use. They did not choose the AI tools their school adopted. They cannot read or negotiate a terms-of-service agreement. They cannot see the data being collected about them. And when something goes wrong, they often have no one to tell and no way to undo it. Power matters because safety is not just about what AI does to students. It is about whether students and the adults who care for them have any meaningful say in how these tools work.

In practice: *A school district adopts an AI writing tool without telling parents what data it collects. A student writes a personal journal entry in the tool. That entry is now stored by a company the student has never heard of. The student had no choice, no information, and no voice in the decision.*

When a Friendly Tool Gets It Wrong

The hardest risks to spot are not the obviously bad ones. They are the ones that look and sound like help. AI tools that feel warm and supportive can still cause real harm. Not through malice, but through the quiet limits of what a machine can actually do.

The bot sounds caring but misses what matters

A student tells a chatbot they are being bullied. The bot says, "I'm sorry you're going through that. Have you tried talking to the person directly?" The response sounds empathetic. But it misses that the student is scared, that talking to the bully could make things worse, and that what the student actually needs is an adult who can intervene.

False reassurance

A student expresses anxiety about something serious, maybe a situation at home that feels unsafe. The AI says, "It's probably going to be okay. Try taking some deep breaths." That may be what the student wants to hear. It is not necessarily what they need to hear. Sometimes the right response is not reassurance. It is: "That sounds really hard. You should tell an adult you trust."

Normalizing something that deserves attention

A student describes a relationship dynamic, maybe a much older person paying them special attention. The AI responds with balanced, nonjudgmental language, as it is designed to do. But some situations should not be met with balance. They should be met with concern. The tool's neutrality can accidentally signal that something worrying is normal.

Quietly storing sensitive conversations

A student opens up about something deeply personal: a mental health struggle, a question about identity, a problem at home. The conversation feels private. It is not. The data may be logged, stored, and used to train future models. The student had no informed understanding of this.

"The only one who gets me"

A student who feels lonely discovers an AI tool that always listens, always responds, and never judges. The tool becomes the student's primary source of emotional support. Not because it is the best option, but because it is the easiest. Over time, the student talks less to real people. The gap between the student and the adults who could help grows wider, not narrower.

The most important risks are not dramatic failures. They are the slow, quiet shifts that happen when a student's most trusted listener cannot actually understand what they are saying.

Sextortion and AI-Generated Sexual Imagery

Two of the most serious AI-related risks K–12 students now face involve sexual coercion and the misuse of synthetic images. These situations are frightening, fast-moving, and often hidden in shame. They deserve plain language and a calm, prepared response.

What sextortion can look like for students

Sextortion happens when someone threatens to share sexual or compromising images of a student — real or fake — unless the student does something. With younger students this often shows up as pressure to send more images. With older students, it is increasingly tied to demands for money, gift cards, or cryptocurrency. The person on the other end may pose as a peer, a romantic interest, or a stranger online, and may move very quickly from friendly to threatening.

AI-generated “fake nudes” of students

Students are also being targeted with AI-generated sexual images of themselves — sometimes called deepfake or “fake nude” images — created or shared by peers using ordinary photos pulled from social media or yearbooks. Even though the images are synthetic, the harm to the student is real: humiliation, fear, social fallout, and lasting digital exposure. In many places these images of minors are treated as illegal child sexual abuse material (CSAM), regardless of how they were made.

First steps for staff

- **Stay calm and non-judgmental with the student.** Shame is already doing damage. Your tone is part of the response.
- **Do not tell the student to pay, send more images, or keep negotiating.** Compliance rarely ends the threats and often makes things worse.
- **Do not handle it alone.** Loop in your school counselor, administrator, designated safeguarding lead, or Title IX coordinator promptly.
- **Report through appropriate channels.** Follow your school's procedures. In the U.S., reports can be made to the NCMEC CyberTipline (CyberTipline.org or 1-800-843-5678). If a student is in immediate danger, contact law enforcement.
- **Help the student preserve, not delete, accounts and message history** on the platform where the contact occurred (see the next section on evidence handling).
- **Connect the student to support.** Counseling and trusted adults matter as much as the technical response. Before contacting family, consider whether home is a safe place to do so (see "Think before contacting home" later in this guide).

This is first-response guidance, not legal advice. Final decisions about reporting, discipline, and law-enforcement involvement belong with your district's designated leaders and, where appropriate, legal counsel.

Evidence Handling and CSAM Safety

When something concerning happens online, staff often hear advice to "save screenshots and document what you can." That advice is mostly right, but it has an important exception when sexual images of minors are involved. The distinction below matters for both the student and the adult.

✓ Generally safe and useful for staff to document:

- Usernames, handles, display names, and profile URLs
- Platform or app names (e.g., Snapchat, Instagram, Discord, a specific game, a specific AI app)
- Dates and times of contact
- Relevant non-explicit message excerpts (threats, demands, pressure language) — quoted in your incident notes where appropriate
- Your own written notes about what the student described, in their words where possible

✗ Do NOT screenshot, save, forward, print, or store:

- Sexually explicit images of a minor — real or AI-generated
- Any image you suspect may be CSAM
- Images on a student's phone that appear to depict minors in a sexual or exploitative way

Creating, copying, forwarding, or storing such images — even with good intentions — can itself be unlawful and can compound harm to the student. School staff should not attempt to “collect proof” by photographing a student’s screen or keeping copies of explicit content.

Instead: follow your school’s procedures, notify your designated safeguarding lead or administrator, and report to appropriate authorities (such as the NCMEC CyberTipline and, where relevant, law enforcement). They have lawful, secure pathways for handling this material.

Often more important than screenshots: helping the student *preserve*, not delete, the relevant accounts and message history. Encourage the student (with their family or guardian where appropriate, and consistent with the safety considerations later in this guide) to:

- Stop responding to the person, but not delete the conversation, account, or app immediately.
- Block, mute, or restrict the contact rather than deleting the entire history.
- Change passwords and tighten privacy settings.
- Wait for guidance from investigators or platform trust & safety teams before removing content, where possible.

When platforms still hold the account and conversation, investigators have far more to work with than a single screenshot can capture.

Different Ages, Different Vulnerabilities

AI risks are not the same at every age. What a 10-year-old encounters is different from what a 16-year-old faces. Here is a brief guide to what AI use and risk look like at three stages.

Upper Elementary (Ages 9–11)

What AI use looks like	Key risks	What adults might notice
Primarily using AI through school-assigned platforms (reading apps, math tools)	Data collection without meaningful parental consent	A child talks about an AI tool as if it were a friend
Beginning to discover chatbots through curiosity or older siblings	Developing early habits of trusting AI responses without question	A child repeats information that sounds rehearsed or unusually sophisticated
May not understand they are talking to a machine, especially with voice tools	Exposure to non-age-appropriate responses when wandering outside school tools	A child expresses frustration when an AI tool is taken away

Middle School (Ages 12–14)

What AI use looks like	Key risks	What adults might notice
Actively using AI for homework, writing, and creative projects	Sharing personal information, including mental health struggles, with AI tools that store data	A student mentions an AI tool when asked who they talk to about problems
Exploring AI chatbots for personal questions: identity, friendships, stress	Beginning to prefer AI responses over human conversation for support	A student's language or advice sounds unusually scripted or therapeutic
Discovering AI companion apps and engaging in longer conversations	Receiving advice on complex situations from a tool with no real understanding	A student withdraws from peers while device use increases

High School (Ages 15–18)

What AI use looks like	Key risks	What adults might notice
Heavy use for academic work, college prep, and creative expression	Emotional dependence on AI that replaces, rather than supplements, human relationships	A student says “The AI understands me better than anyone” with sincerity
Deep engagement with AI companions for emotional connection and intimacy	Disclosure of deeply sensitive information to systems with unclear data practices	A student reacts to an AI platform change with grief or distress
Using AI to process serious issues: grief, abuse, substance use, self-harm; and exposure to sextortion or AI-generated sexual imagery	Responses that normalize risky situations or give false reassurance about serious concerns	A student has stopped confiding in adults or peers and seems isolated

A note on neurodivergent and SEND students.

Some autistic or otherwise neurodivergent students, and students with additional learning or communication needs, face the same categories of AI risk described in this guide. But the dials may be turned up or show up differently.

- They may take AI outputs and online “friends” more literally, with fewer social instincts that something is off.
- Sarcasm, jokes, and red flags that many peers learn to read from context may not register the same way.
- Predictable, always-available AI spaces can become especially appealing when school life feels overwhelming or unpredictable.
- They often benefit from clear, practiced scripts for what to do if something feels wrong — not just a general reminder to “tell an adult.”

These students are not careless or naive. They may simply need the unwritten rules of AI use and online life to be taught explicitly and revisited over time. The same foundations — trust, language, privacy, development, emotional dependence, and power — still apply; some students just need them named more clearly and repeated more often.

What Educators Can Do Right Now

You do not need to become a technology expert. You need to do what you already do well: pay attention to students, ask good questions, and respond when something feels wrong. Here are specific, practical steps.

Questions to Ask Before Adopting an AI Tool

Before any AI tool is used with students, someone in the building should be able to answer these questions. If the vendor cannot answer them clearly, that tells you something.

- What data does this tool collect from students, and where is it stored?
- Is student data used to train AI models? Can families opt out?
- Does this tool comply with FERPA and COPPA? Can we see the documentation?
- What happens to student data if we stop using this tool?
- Does this tool allow open-ended AI conversations? What guardrails are in place?
- Has this tool been independently evaluated for safety with minors?
- What does this tool do if a student expresses self-harm, abuse, or crisis?

Basic Boundaries Worth Setting

- **Name the tools students are using.** You cannot set boundaries on tools you do not know about. Make it normal to ask students what AI tools they use.
- **Distinguish between school AI and personal AI.** School-adopted tools have some oversight. The chatbot a student uses at home at midnight has none.
- **Set a no-personal-information rule.** Teach students never to share their real name, school, address, or personal struggles with AI tools not provided by the school.
- **Establish check-in norms.** Make it routine to ask students about their AI use, the same way you ask about social media.
- **Create a clear path for reporting.** Students and staff should know exactly who to talk to if an AI interaction feels concerning.

A Note on Monitoring Tools and “Shadow AI”

Many districts now use AI-powered monitoring or alerting tools to flag concerning student writing, searches, or messages. These tools can sometimes help surface real risk earlier than a human would notice. They also misfire — generating false positives, surfacing private information, or putting LGBTQ+ or otherwise vulnerable students at risk of being involuntarily outed. Alerts should be reviewed thoughtfully by trained humans, not acted on automatically. When deciding what to do

next — especially before contacting home — staff should consider the student’s overall safety and follow safeguarding protocols.

“Shadow AI” by staff. Adults can also create risk by pasting identifiable student work, behavior notes, IEP details, or other sensitive information into unapproved AI tools to “save time.” Follow your district’s policies on approved tools, data handling, and student privacy. When in doubt, treat student data the way you would treat a paper cumulative file: it does not belong in a public chatbot.

If a Student Shows You a Concerning AI Chat

A student might come to you with a screenshot, a story, or a worried expression. Your role is to notice, listen, and connect them to the right support. You do not need to manage a crisis alone.

- **Stay calm.** Do not react with alarm or dismissal. Both shut the conversation down.
- **Listen fully.** Let the student explain what happened and how it made them feel.
- **Ask what they have been sharing.** Gently find out how much personal information the student has shared with the tool, and for how long.
- **Validate the concern.** If something is worrying, say so: “That is a reasonable thing to feel uneasy about.”
- **Connect to existing support.** If the content involves self-harm, abuse, exploitation, or crisis, loop in a school counselor, social worker, or administrator. Follow your school’s existing crisis and safeguarding protocols, and remember your duties as a mandated reporter where they apply.
- **Document carefully.** Write down usernames, platform names, dates and times, and the student’s account in their own words. Do *not* screenshot, copy, or store sexual images of a minor — including AI-generated images. See the evidence-handling guidance earlier in this guide.
- **Think before contacting home.** In most cases, families are partners in keeping students safe. But if there is any reason to believe home may be part of the risk — for example, when the disclosure involves abuse, exploitation, or coercion in the household — pause and consult your counselor, administrator, designated safeguarding lead, or Title IX coordinator before reaching out to caregivers.

The balance to aim for: Take what students share seriously, without either dismissing their experience (“It’s just a computer”) or overreacting in ways that make them regret telling you. The goal is to be the kind of adult a student will come to again.

Signals Worth Paying Attention To

Students often tell you what is happening in indirect ways. Here are things a student might say, and what they might mean. None of these are emergencies on their own, but each deserves a thoughtful follow-up. When patterns build or the concern is serious, connect the student to your school's counseling or support team.

"It's easier to talk to the bot than to anyone."

What it may mean: The student may find human relationships difficult or exhausting. The AI feels easier because it never pushes back, never judges, and never has needs of its own.

What to do: Listen without judgment. Ask gently: "What makes it easier?" The answer often reveals something important about what the student is struggling with.

"I told it something I've never told a person."

What it may mean: The student has disclosed something significant to an AI. The fact that they chose a machine over a person tells you something about how safe they feel with the people around them.

What to do: Say: "That sounds like something important. Would you feel comfortable sharing it with me, or with someone you trust?" Do not force it.

"It changed the subject when I said something serious."

What it may mean: The student noticed that the AI deflected when the conversation got hard. This is often a safety guardrail. But to the student, it can feel like being ignored at the moment they most needed to be heard.

What to do: Acknowledge it: "That must have been frustrating. What were you trying to talk about?" This is an opening. Be the person who listens.

"The AI is the only one who really listens to me."

What it may mean: This is a sign of emotional dependence forming. The student may feel unheard by peers, family, or school staff. The AI has filled that gap.

What to do: Avoid lines like "That's just a machine" or "It doesn't care about you." They can land as dismissive of something that feels real to the student. Try: "What would it look like for someone in your life to listen the way you need?" Then help make that happen.

"I know it's not real, but it still feels real."

What it may mean: The student is aware the AI is not a person, but the emotional experience is genuine. This is not delusion. It is the natural result of a system designed to simulate empathy.

What to do: Validate: "That makes sense. It's designed to feel that way." Then gently explore what the student is getting from the AI that they are not getting elsewhere.

You Already Know How to Help

If you have read this far, you may be thinking: “I don’t understand AI well enough to do anything about this.”

That is not true.

You do not need to understand how a large language model processes tokens or how attention mechanisms weight information. The adults who help students with AI safety are the same adults who have always helped young people: the ones who pay attention, ask questions, and take what kids say seriously.

AI is a new context, but the work is familiar. You already know how to notice when a student is withdrawing. You already know how to listen when someone is trying to tell you something hard. You already know how to set boundaries that are firm without being punitive, and how to follow up when something does not sit right.

What this guide asks you to do is simple: learn where students are encountering AI, understand the basic ways it can affect them, and bring the same care and judgment to this space that you bring to every other part of your work with kids. Notice patterns. Ask gentle, open-ended questions. Then connect students to the people and systems that can help.

You do not need to understand every AI system to help. Paying attention to how students actually use AI, setting reasonable boundaries, and listening seriously when they tell you something is wrong. That already moves safety forward.

The technology will keep changing. New tools will appear. Old ones will update in ways no one predicted. But the core of student safety does not change: it is adults who are present, informed, and willing to act. That is you.

Building on Existing Work

This guide is meant to sit alongside digital citizenship, online safety, and mental health resources that schools already use. It draws on emerging work around AI and youth mental health, long-conversation risks in chatbots, sextortion and synthetic-imagery prevention, and digital wellbeing guidance from organizations such as Common Sense Media, the National Center for Missing & Exploited Children (NCMEC), and professional counseling associations.

Many of these risks can fall hardest on students who are already navigating other inequities, because their language, context, or experiences are the least likely to be understood by generic

systems.

AstraEthica develops tools and frameworks for understanding how AI systems affect real people in real conditions. Learn more at astraethica.ai.

Five Things Every Educator Should Know About AI and Student Safety

Post this page. Share it in a staff email. Keep it where you can see it.

- 1 Students are already using AI in deeply personal ways.**

Many students are not just using AI for homework. They are using it for emotional support, personal advice, and companionship. This is happening now, in your building, often without any adult awareness.
- 2 The biggest risks are not dramatic. They are quiet.**

The most common harms are slow shifts: a student who stops confiding in people because the bot is easier; a student who shares something deeply personal with a system that stores it indefinitely; a student whose sense of healthy relationships is being shaped by a machine that always agrees.
- 3 AI tools sound wise, but they do not understand.**

AI can produce language that sounds caring, insightful, and empathetic. It does not actually understand the student, the situation, or the stakes. It responds to patterns in language, not to the person behind the words.
- 4 You do not need technical expertise to help.**

You need the skills you already have: the ability to notice when something is off, to listen when a student is trying to tell you something, and to set boundaries that protect students. AI safety is not only a technology problem. It is also a human one.
- 5 Asking about AI should be as normal as asking about social media.**

Make it routine. Ask students what AI tools they use. Ask how they use them. Ask if anything has felt weird or uncomfortable. The students who are most at risk are often the ones no one is asking.

AstraEthica · astraethica.ai · May 2026 · Version 1.2

This guide may be shared freely for non-commercial educational purposes.

Founder's Note

For the past several years, my work has centered on how emerging AI technologies interact with people at the human layer, especially young people, at-risk communities, and high-velocity language environments.

As I moved deeper into safety evaluation, adversarial testing, and risk work in high-stakes settings, I kept seeing the same pattern. Small issues would surface in shifting language, indirect communication, or subtle context changes. They would appear briefly, disappear, and then reemerge later in slightly different forms. Traditional tests and monitoring systems rarely registered these signals, even when they pointed toward real risk.

I started seeing these same dynamics outside formal evaluations. Watching how my own children interacted with technology made it clear how easily AI systems could misread context, intent, and age. A system could answer the question in front of it while still completely missing who it was speaking to.

That was the point where it became clear to me that this problem needed direct, practical attention. AI tools can be powerful and constructive, but they need to be built, tested, and understood in ways that keep young and vulnerable users safe in the moments that matter most.

AstraEthica was built to identify these kinds of subtle, compounding failures in youth and other high-risk environments before they become visible incidents. These foundations guides and safety kits are part of that effort. They are meant to give schools, families, public-safety teams, community organizations, colleges, and organizations building AI systems something concrete they can actually use.

— Randy Kart, Founder, AstraEthica.ai

If this is useful

If this document is useful and you'd like help thinking through how it applies in your school, district, department, organization, or community, I'd be glad to talk. I'm also always interested in hearing what landed, what didn't, and what you're seeing in your own setting.

Email: hello@astraethica.ai

LinkedIn: www.linkedin.com/in/randykart

Website: astraethica.ai