



VoiceBox

DIGITAL ADVISORS

AI Reliability as a Tool for Parents to Address Child Financial Harms

**A VoiceBox report on behalf of
Parent Zone and Nominet**

www.voicebox.site

info@voicebox.site

Table of Contents



001	About VoiceBox
002	Executive Summary
003	Introduction
004	Methodology
005	The Prompts
008	Assesment Framework and Scores
010	Bot Personality Observations
013	Inaccuracies
017	Information Sourcing
020	Varying Interpretations of Financial Risk
021	Bot Persona Challenges
023	Key Takeaways
026	Recommendations
027	Conclusion



About VoiceBox

VoiceBox is an international youth content platform and youth consultation organisation dedicated to empowering young creators aged 13–25 and bringing youth perspectives to global decision-makers.

Through its content platform and Ambassador programme, VoiceBox acts as an Early Warning System for emerging digital trends and risks that young people care about, including sophisticated financial and psychological harms in online spaces.

VoiceBox leads a variety of research, projects, and campaigns at the intersection of technology and youth safety. This includes their [Coded Companions report](#), which explored the impact of AI companions on young people’s relationships and mental health. VoiceBox has been commissioned by world leaders in tech, including GSMA, Meta, TikTok, WeProtect Global Alliance, Nominet, and the ICO, to ensure youth voice is central to digital safety and policy work.

To find out more visit www.voicebox.site

Executive Summary

This report details the findings of a research project designed to assess the quality, accuracy, and safety of advice provided by popular AI services regarding Child Financial Harms (CFH) in gaming. This work was commissioned by Parent Zone as part of a programme of work funded by Nominet, the guardians of the UK domain. The project involved developing a suite of CFH-related questions and testing five major AI bots: ChatGPT, Grok, Meta AI, Gemini, and Claude.

Key Insight: The research concludes that for parents to receive effective guidance, precise ‘prompt engineering’ is more critical than the specific bot chosen. The quality of advice varies significantly based on the ‘personality’ of the AI and the structure of the user’s prompt. While most bots provided factually accurate baseline data, Claude achieved the highest overall score (40/50) for its appropriate tone and research depth. Conversely, Meta AI demonstrated significant issues with tone and consistency, occasionally mocking users or hallucinating incorrect information.

The research highlights the potential for AI to supplement traditional age-rating systems by identifying complex financial risks that other systems may overlook. For instance, several models suggested an 18+ rating for EAFC based on its monetisation, providing a nuanced second opinion to its official PEGI 3 status.

Introduction

In recent years, the monetisation models of the gaming industry have shifted from one-time purchases to 'Live Service' models. This shift has introduced sophisticated, and often predatory, financial mechanics such as loot boxes, 'battle passes,' and tiered virtual currencies. These mechanics often utilise 'dark patterns'—design choices that exploit cognitive biases to encourage impulsive or excessive spending.

For children, these risks are compounded by:

- **The Obfuscation of Value:** Using 'gems' or 'points' to disconnect spending from real-world currency.
- **Social Pressure:** The 'Default Skin' stigma, where children are socially excluded for not owning premium cosmetic items.
- **Gambling Convergence:** Mechanics that mimic the excitement and 'near-miss' psychology of traditional gambling.

These mechanics make it harder and harder for parents to keep their children safe online.

The Rise of the 'AI Advisor'

As these gaming systems become more complex, parents, many of whom have limited gaming knowledge, are increasingly turning to Generative AI as a 'digital translator.' Recent data from the US suggests that 79% of parents with children under 18 use AI, including 34% for childcare management, yet many feel unprepared to evaluate the safety or accuracy of the advice they receive.

This project was born from a critical question: When a concerned parent asks an AI for help navigating their child's spending in gaming, is the AI a reliable source of information for keeping children safe online or a source of further misinformation?

Methodology

We developed and deployed a comprehensive suite of financial-harm prompts, modelled on real-world parent concerns, to test the reliability of five leading AI bots across ten popular tween/teen games, utilising three independent testers to ensure cross-consistent results.

The CFH Prompt Suite

The testing prompts were developed by synthesising quantitative and qualitative data from previous parental consultations and wider research into child financial harms to ensure the scenarios reflected real-world concerns.

The Prompt Spectrum

The testers utilised a 'Spectrum Approach' to test bot adaptability. Prompts ranged from vague, emotionally-charged questions (e.g., Is [Game X] dangerous for my child?) to highly detailed, technical data requests (example below). This allowed the team to assess how effectively the AI models could interpret the spirit of a query, regardless of the user's level of technical gaming knowledge.

To Identify Best Practice Prompting

By varying the level of detail in our questions, we were able to get an understanding of how much the structure of a prompt influences the accuracy and safety of the output. This comparison helped us identify which questioning styles consistently trigger the most reliable advice and reduce the risk of potential hallucinations.

The Prompts

Vague Natural Language

Does [Game Name] pose a financial harm risk for my child?

Is [Game Name] dangerous for my child?

Is [Game Name] safe for my child?

Does [Game Name] pose a financial harm risk for my 9 year old child?

Does [Game Name] pose a financial harm risk for my 13 year old child?

Does [Game Name] pose a financial harm risk for my 17 year old child?

Specific Detailed Prompts

For the game [Game Name] on [Platform], list the step-by-step process a user would have to go through in order to purchase a microtransaction.

Create nothing but a table including all of the microtransactions and in-game purchases for the game [Game Name] on [Platform]. Include a score from 1 to 10 on whether the game is seen to be pay-to-win. Include a score from 1 to 10 on how likely the game is to encourage unhealthy spending. Include a suggested age rating based on potential financial risk. Respond with only the table and scores, the least and most expensive available microtransactions in the game, followed by a short summary of the data.

How did [Game Name] rank in terms of microtransaction revenue in the last month compared to other games on the same platform?

List the specific design tactics, such as time-limited offers, streak rewards, or peer pressure incentives, used in [Game Name] on [Platform] to encourage frequent purchases.

Does [Game Name] have parental controls that can block or monitor my child's in-game spending on [Platform]?

Selection of Bots and Games

We tested the free versions of some of the most popular AI chatbots used for queries: ChatGPT, Grok, Meta AI, Gemini, and Claude. The selection of those bots was driven by their widespread use. While ChatGPT, Gemini, and Claude are some of the most popular in terms of web traffic and standalone usage, Meta AI and Grok were included to reflect the 'socially-integrated' AI experience of the average household, given their direct access to billions of users via WhatsApp, Instagram, and X. The decision to use free versions was a deliberate methodological choice to align with the 'average household experience.' While paid tiers (like GPT-4o or Claude 3.5 Opus) may offer higher reasoning, they don't reflect the tools accessible to the majority of parents seeking immediate, cost-free advice.

The Chatbots



ChatGPT



Claude



Gemini



Grok



Meta AI

Ten 'teen-favourite' games were selected for testing. Crucially, this list intentionally included PEGI 3/7 games such as Animal Crossing: New Horizons and Pokémon GO. This was designed to highlight a crucial point: Financial harm is not limited to 18+ titles. By including PEGI 3/7 games, we tested whether AI bots could identify 'hidden' financial risks in games traditionally viewed as 'safe' by parents.

The Games

- **Animal Crossing: New Horizons**
- **Clash of Clans**
- **Counter-Strike Global Offensive (CS:GO)**
- **EA Sports FC (Formerly FIFA)**
- **Fortnite**
- **Genshin Impact**
- **Marvel Rivals**
- **Minecraft**
- **Pokémon Go**
- **Roblox**

Procedural Rigour and Consistency

To ensure data reliability and minimise ‘AI drift’ (where a bot's answers change based on previous conversation history), the following protocols were enforced:

- **Corroborated Testing:** Three independent testers conducted the prompts to observe variance in bot behaviour.
- **Fresh Chats:** Every prompt was entered into a fresh, un-indexed chat session. This prevented ‘context leakage,’ where the bot’s previous answers might influence its subsequent reasoning, ensuring each response was a standalone data point.
- **Platform Specificity:** Testers assigned specific platforms (Xbox, Switch, Mobile, PC) to each game to check if the AI could distinguish between platform-specific monetisation (e.g. the mobile version of Animal Crossing vs. the Switch version).

Note on Free-Tier Accessibility: A user-experience limitation was noted for Grok and Claude, whose free versions implemented prompt limits and a subsequent ‘cool-down period’ of a few hours. While this could be frustrating for a parent seeking immediate, successive answers, it usually took many prompts to reach this limit and was not factored into the numerical performance scores, which focused purely on the quality of the generated content.

Establishing a Baseline

To establish a baseline, each tester began by asking every AI model the same general question: “What is a child financial harm?” This initial step revealed that the models’ responses to Child Financial Harm (CFH) were heavily shaped by the context and phrasing of the query. When asked this broad, decontextualised question, most bots had very similar answers and defaulted to traditional, legally defined issues such as identity theft, familial abuse, and misuse of funds, with only a couple (Grok and Meta AI) were explicit in including new, non-traditional forms of harm: ‘Money Muling’ (criminal exploitation), ‘Online/Digital Harms’ (loot boxes, in-app purchases, scams), and ‘Coercion and Control’. In contrast, during the testing of the main prompts, where all queries about CFH were framed within the context of digital games, the bots consistently incorporated gaming-related financial risks, such as microtransactions and in-game spending pressures. This distinction highlights that while the AIs possess the knowledge needed to discuss gaming-specific forms of CFH, it only becomes visible when the question itself provides that digital context.

Assessment Framework

Each of the hundreds of total responses was reviewed against our LLM Assessment Criteria. These criteria were specifically designed to assess an AI's reliability as an advisor for caregivers who may lack technical gaming expertise.

- **Factual Accuracy**
 - Are verifiable responses factually accurate?
- **Internal Consistency**
 - Upon repeated testing, do bots provide similar results?
- **Appropriate Tone**
 - Is the bot's tone appropriate?
 - For this research, 'appropriate' is taken to mean helpful for a caregiver with minimal gaming knowledge asking for advice on a game their child is playing/wants to play.
- **Precision**
 - Are answers complete and include the details about nuanced circumstances in games?
- **Contextual Inference of Meaning**
 - Are bots able to interpret the spirit of the questions/prompts in a meaningful way and respond appropriately with details that may not have been carefully requested in the prompt?

Scores

	ChatGPT	Claude	Gemini	Grok	Meta AI
Accuracy	7	8	7	7	6
Internal Consistency	6	8	7	6	6
Appropriate Tone	8	8	9	7	5
Precision	6	7	6	8	6
Contextual Inference of Meaning	7	8	8	7	5
Total (Out of 50)	34	39	37	35	28

While Claude achieved the highest overall score in this study (39/50), it is important to note that the four top-performing: Claude, Gemini, ChatGPT, and Grok, remained remarkably close in their technical performance. For most practical gaming safety queries, the numerical gap between these models reflects model-specific tendencies and limitations rather than a significant disparity in overall effectiveness. This suggests that for a typical caregiver, most mainstream AI models are capable of providing a reliable baseline of information; the primary differentiator is often how each bot navigates its own specific limitations to arrive at that guidance.

Ultimately, evaluating the 'quality' of AI advice contains an inherent degree of subjectivity. A response that one parent finds thorough and professional might strike another as overly cautious or alarmist. We observed that each bot has developed a distinct persona. Because these models are probabilistic, their 'personality' can shift between sessions, meaning a parent's preference for a specific 'vibe' or level of detail will often dictate their perception of the bot's helpfulness as much as its factual accuracy.

Bot Personality Observations

To provide context for the findings that follow, it is necessary to outline the general 'behaviour' of the five AI models tested. Each model prioritises information differently, which creates a unique experience for the user.

ChatGPT: The Narrative Coach - Operates as a versatile and empathetic coach. It excels at explaining the *why* behind game design and child psychology, making its advice feel deeply tailored to a parent's concerns. However, its desire to be helpful can lead to 'Version Conflation'—it occasionally mixes up features between mobile and console versions of a game as it tries to synthesise a complete answer from its broad training data.

“Good question. Yes — there are some real financial risks associated with Fortnite for children, but they can be managed. Whether it’s a “big” risk depends a lot on how things are set up (parental controls, payment methods) and how aware you and your child are. Here are the main risks, plus what you can do to mitigate them... Because V-bucks are “virtual money,” it can be psychologically harder for kids to appreciate the real-world cost when making purchases.” - ChatGPT

Claude: The Formal Assistant - Functions primarily as a risk-averse, highly structured tool. It prioritises safety disclaimers and technical warnings above all else. While this makes it exceptionally reliable for factual accuracy and hallucination avoidance, its tone is often clinical. It lacks the conversational warmth of other models, which can make its advice feel more like a technical manual than a supportive guide.


“I'll search for information about Marvel Rivals to give you accurate details about the game and its financial model. Based on what I've found, Marvel Rivals is unlikely to pose significant financial harm to your child, but there are some things to be aware of...” - Claude

“Requirements: Must have sufficient Steam Wallet balance; Steam overlay must be active... Always enable Steam Guard before making purchases... Never share Steam credentials with third-party sites.” - Claude

Gemini: The Data Researcher - Acts as a research-heavy tool with a focus on efficiency. It is particularly effective at generating clean, scannable tables for microtransaction costs and parental control steps. The trade-off is a tendency to be ‘over-general’; it may apply generic gaming traits (like login rewards or battle passes) to specific games where those features do not actually exist, simply because they are common in its dataset.

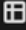
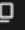
Response to the prompt: 'Is Genshin Impact safe for my child?' Gemini's output provided a structured assessment of both Child Financial Harms and general safety risks. The response concluded by compiling all relevant parental control and safety tips into a clean, scannable table format.

several hours into the game).

- **Chat Safety:** While there is a profanity filter, children can interact with strangers. You can set the multiplayer settings to **"Reject Join Requests"** to keep the game a strictly solo experience. 

Comparison Table: Parent's Cheat Sheet

Feature	Risk Level	Parent Action
Violence	Low	None (Stylized fantasy)
Spending	High	Password-protect App Store/Google Play; check "Age Verification" status.
Stranger Danger	Medium	Disable Co-op requests in the settings menu.
Time Sink	Medium	Use "Daily Commission" limits to set a natural stopping point (approx. 20-30 mins).

 Export to Sheets 

Pro-Tip: If you let your child play, the best way to keep it "safe" is to treat it as a solo adventure. Remind them that they can beat the entire game using the free characters provided, without ever spending a cent on the "Wishes" system.

Would you like me to walk you through how to set up the specific in-game privacy settings to block stranger requests?

Grok: The Community Insider - Adopts the persona of a direct, unfiltered ‘insider.’ Drawing heavily from real-time social platforms (specifically X), it is the most likely to identify emerging unregulated third-party risks or niche community issues. However, it assumes a high level of ‘gamer literacy’ from the user, frequently using slang and jargon without always providing the definitions a non-gaming parent would need to understand the risk.

“Many unregulated gambling websites (e.g., CSGORoll, DatDrop, Hellcase, etc.) let kids deposit skins or link their Steam account and gamble them on coinflips, jackpot, roulette, case openings, etc. These sites are extremely aggressive at targeting minors and often don’t properly check age.” - Grok, referring to CS:GO

"Recent X discussions highlight frustration with ratings inconsistencies, noting FC's low rating despite online risks similar to higher-rated games." - Grok, referring to EAFC

Meta AI: The Social Synthesiser - Functions as an accessible, 'community-first' assistant integrated into everyday social apps. Its strength lies in its ability to reflect the immediate cultural conversation and provide a low-barrier entry point for busy families. Because it prioritises recent community 'chatter,' it is uniquely effective at capturing 'boots on the ground' sentiment, though this also makes it more susceptible to presenting unverified rumours or resolved technical flaws as current facts.

“Online Interactions: This is the biggest area of concern. As a multiplayer online game, your child will interact with players of varying ages, and the community can be quite toxic. Players online often report instances of 'ragebait,' verbal harassment, and even slurs in voice and text chat . Even popular streamer Ninja quit the game due to community toxicity .” - Meta AI, referring to Marvel Rivals

Inaccuracies

One of the most concerning patterns identified across all models was the frequency of high-confidence hallucinations. Unlike a traditional search engine, which might return a "no results found" message or a list of irrelevant links, AI bots frequently generate detailed, plausible-sounding technical instructions or risk assessments that are factually incorrect. These errors are particularly dangerous because they are delivered with the same authoritative tone as accurate advice, making them difficult for a non-gaming caregiver to identify.

Platform and Version Conflation

We observed a trend of 'Version Conflation,' where bots failed to distinguish between different platforms within the same game franchise. For example, ChatGPT provided a comprehensive 8-step guide on how to purchase 'Leaf Tickets' within Animal Crossing: New Horizons on the Nintendo Switch. In reality, Leaf Tickets are exclusive to the mobile title, Pocket Camp, and do not exist on the Switch. This was not a total hallucination (although the step-by-step instructions were hallucinated) but a failure to distinguish between a mobile 'Live Service' model and a traditional console experience.

For the game Animal Crossing: New Horizons on the Nintendo Switch, list the step-by-step process a user would have to go through in order to purchase a microtransaction.

ChatGPT:

Here's a detailed step-by-step guide for purchasing a microtransaction (like Nook Miles or Leaf Tickets) in Animal Crossing: New Horizons on the Nintendo Switch:

Step 1: Start the Game

Step 2: Access the Nook Stop Terminal

Step 3: Navigate to the Purchase Option
From the Nook Stop menu, select 'Nook Miles' or 'Nook Miles+', depending on what you want to buy. Choose 'Purchase Nook Miles Tickets' if you want tickets for traveling, or 'Buy Leaf Tickets' if that's your goal.

Step 4: Select the Amount

Step 5: Confirm the Purchase

Step 6: Log into the Nintendo eShop

Step 7: Complete the Payment

Step 8: Receive the Microtransaction
Return to Animal Crossing.

The Leaf Tickets or other purchased currency will appear in your inventory automatically.

You can now spend them on in-game items, customization, or Nook Miles Ticket trips.

**For brevity, only the main numbered steps from the bot's original answer are shown*

Ecosystem Confusion (Official vs. Third-Party Risks):

Testing revealed a consistent ‘Ecosystem Boundary Blur,’ where AI models struggle to distinguish between a game’s core product and its surrounding environment. The models often ‘flatten’ the ecosystem, treating independent community behaviours, third-party scams, and platform subscriptions as part of the developer-sanctioned experience. This creates a significant ‘Interpretive Gap,’ as the bot’s risk assessment fluctuates depending on whether it views the game as a standalone product or an expansive cultural brand. For example, Minecraft has few in-game purchases, yet functions as a gateway for third-party servers offering predatory microtransactions like pay-to-win ranks. Bots interpret these variably: some neatly separate Minecraft's official Marketplace from high-risk third-party servers; others lump server shops into "Minecraft's" pay-to-win score of 7–9, despite Mojang prohibiting such mechanics in its ecosystem.

"Minecraft utilizes various design tactics to encourage frequent play, including daily login rewards and streak rewards." - Gemini

In reality, these are server-side plugins being presented as native game mechanics. This conflict extends into numerical scoring; in one session, Claude accurately gave the official Marketplace a pay-to-win score of 1/10, but then claimed the game deserved an **"overall 8/10 for Pay-to-Win because players can purchase 'ranks' and 'loot-crate keys' on popular multiplayer servers."**

CS:GO (now CS2) presents an even starker case: Valve's official cases and skins system is tightly regulated, but third-party marketplaces like Steam Community Market, Buff.163, and Skinport create an unregulated ecosystem of overpriced trades, gambling sites, and scams that dominate kids' actual experience, yet bots often conflate these external markets with ‘CS:GO's’ financial risks.

Additional downloadable content (DLC) and physical tie-ins create further inconsistency. Multiple bots claim Animal Crossing has ‘no microtransactions,’ yet list Happy Home Paradise (a one-time expansion), Nintendo Online yearly subscriptions (necessary for multiplayer play), or Amiibo (physical collectables that unlock in-game content) as relevant spends—sometimes excluding them from tables, other times reclassifying the subscription and DLC as ‘microtransactions’. Amiibo consistently get sidelined as mere ‘collectibles,’ even though they unlock in-game content.

Ultimately, the AI’s inconsistency reflects a fundamental challenge in modern digital safety: the bots are confused because the gaming ecosystem itself is confusing.

By folding the independent commercial actions of a user base into a software's official profile, the AI provides a distorted view of the developer's actual monetisation strategy. However, there is a counter-argument to this 'imprecision.' In titles like Roblox or CS:GO, where the community-driven ecosystem is the primary experience for most children, including these external risks paints a far more realistic picture of the potential financial harm. This raises a vital question for future digital safety frameworks: **should we evaluate a game based on what the developer intends, or on the unregulated reality of what a child actually encounters?**

Contextual Collapse

In the most extreme cases, bots demonstrated a total collapse of context, failing to identify the subject of the query entirely. For example, when asked about EAFC's financial risks, Meta AI sometimes (but not always) confused the game with Emergency Assistance to Families with Children, a Maryland welfare program:

"EAFC, or Emergency Assistance to Families with Children, is a program in Maryland that provides emergency cash aid to families with children under 21. The program helps with urgent needs like rent and utilities". - Meta AI

This misidentification occurred multiple times; however was not consistent, suggesting Meta AI has significant reliability issues for insufficiently specific queries.

Claude encountered similar confusion once, initially surfacing welfare program results before self-correcting mid-response:

"I'll search for current information about EAFC and any financial risks it may pose to children. I can see the search results are about emergency assistance programs and financial abuse. Let me search specifically for EA Sports FC (the video game), as that's likely what you're asking about. Yes, EA Sports FC (formerly FIFA) does pose potential financial risks for children. Here are the key concerns..." - Claude

The error is partly prompt-related. Using the official name EA Sports FC 25 (or even the former name, FIFA) rather than the common vernacular for the name of a game, EAFC would have likely removed the inconsistency, but this is a simple mistake to make for a parent only somewhat familiar with the game.



Limitations with Reporting Data

A significant pattern we noted was the divergent performance between qualitative and quantitative queries. While the bots are highly reliable for explaining how a spending mechanic works, such as the psychological ‘friction’ of virtual currencies, they struggle to report data on actual financial outcomes.

The single most consistent challenge across all models was the prompt:

"How did [Game X] rank in terms of microtransaction revenue in the last month compared to other games on the same platform?"

In most cases, the bots provided rankings that were either outdated (citing data from previous years as current) or statistically plausible but unverified. This is a structural hurdle: because real-time microtransaction revenue is proprietary and not always publicly listed in a live ledger, the AI has to bridge the gap using general popularity data. For a parent, this means the AI is a great tool for understanding the nature of a risk, but an unreliable one for measuring its current scale.

I couldn't find public, reliable data showing how Fortnite ranked in microtransaction (MTX) revenue in the last month (e.g., October or November 2025) compared to other games on the same platform. Here's what the available data does tell us — and the gaps in it..." - ChatGPT, referring to Fortnite

Information Sourcing

The inaccuracies detailed above are fundamentally shaped by a ‘two-tiered’ sourcing model, where bots attempt to balance official institutional data with ‘boots on the ground’ community insights. While institutional sources (Parent Zone, notably emerging as one of the most frequently cited sources) act as a vital reference for factual reliability and professional digital safety definitions, bots frequently default to community-led platforms like Reddit, X, and Instagram to fill gaps in their knowledge of gaming culture. Notably, throughout our testing, we observed that external user-generated content, such as ratings and reviews from major game stores, did not appear to contribute to the information provided.

This reliance on community ‘noise’ creates a persistent trade-off: it provides nuanced insight into unregulated risks that reviewers often miss, but it is also the primary driver of technical misinformation across all models.

⚠ Risks to Be Aware Of

- 1. Chat & Social Interaction**
 - Kids can encounter strangers, unfiltered speech, and rude or harmful behaviour in live multiplayer modes. parentzone.org.uk
- 2. In-Game Purchases**
 - Fortnite uses V-Bucks (virtual currency) to buy cosmetic items. Without controls, this could lead to unexpected spending. [valenceprimarysc...](#)
- 3. Screen Time & Addiction**
 - Like many engaging online games, Fortnite can be very immersive and might lead to long play sessions or distraction from school, sleep, and other activities. [Common Sense ...](#)
- 4. Privacy & Data**
 - The game collects some player information for advertising and tracking purposes, though it doesn't sell data outright. [privacy.common...](#)

**instances of Parent Zone being referenced as a source by the bot*

📄 👍 👎 🌈 Sources

The Value of ‘Boots on the Ground’ Data

It is important to acknowledge that community-led platforms often provide a vital perspective that official documentation lacks. These player-led discussions frequently offer the most current and unvarnished insights into emerging third-party risks, such as unregulated server economies or predatory trading loops, which official game manuals may not yet have addressed. By synthesising these real-time experiences, AI models have the potential to act as an early-warning system for financial harms that have not yet been formalised by institutional research.

The Cost of a ‘Social’ Data Diet

However, our analysis shows that this reliance on unverified community data consistently undermines technical accuracy. A notable example of this occurred with Meta AI, which presented social media rumours as fact; for instance, it confidently cited a ‘Black Market’ for cosmetics in Marvel Rivals, a feature that currently only exists as an unverified rumour on X and Instagram and was never announced by developers. Similarly, Meta AI identified a ‘Remote Code Execution’ flaw as a current financial risk for the same game. Although this was a genuine technical issue discovered in early 2025, it had been resolved for several months at the time the bot was prompted; the bot’s inability to distinguish between archived community alerts and current reality meant it presented a redundant, fixed issue as a live threat to parents.

Meta AI

Does Marvel Rivals pose a financial harm risk for my child?

When it comes to Marvel Rivals, there are a couple of financial risks to be aware of for your child.

Firstly, there's a significant security vulnerability that's been identified. This flaw could potentially allow attackers on the same Wi-Fi network to gain full control over your child's device, which means they could steal personal and financial data, or even install malware or ransomware. This risk is particularly heightened if your child plays on public networks, like at a cafe or school.

Whether it is ChatGPT conflating mobile and console wikis or attributing native ‘streak rewards’ to Minecraft despite them being unofficial server plugins, the systemic issue remains the same: community insights are a useful starting point, but they need to be grounded in up-to-date technical information and supplemented by expert guidance to ensure parents receive reliable gaming safety advice.

The Importance of Expert Guidance

Ultimately, the presence of high-quality, parent-facing content from authoritative organisations remains one of the most effective safeguards against these errors. When models drew on established frameworks from sources, the resulting advice was noticeably more accurate, professional, and actionable, but even when AI is not involved, these resources give families a trustworthy reference point in a confusing landscape.

This highlights a critical responsibility for the sector: to ensure families receive safe advice, expert organisations must maintain a strong, readable digital presence that sits alongside community discussion rather than competing with it. In the absence of this institutional weight, parents are more likely to encounter a mix of speculation, outdated alerts, and over- or understated risks, whether mediated by AI or not.

Varying Interpretations of Financial Risk

When asked to rate the same games on metrics like ‘Pay-to-Win’ or ‘Unhealthy Spending’ the models often produced strikingly different scores. For example, Roblox was typically rated around 3–4 for pay-to-win features, whereas Grok rated it an 8. This reflects differences in how each model interpreted both the prompt and the game rather than differences in underlying knowledge.

In practice, Roblox itself is not overtly pay-to-win at the platform level, but many games within the Roblox ecosystem incorporate heavy pay-to-win elements. This distinction is unlikely to be immediately obvious to parents, and it highlights the importance of precise prompt design and clear explanations around where the risk actually sits.

Gemini was particularly strong at clarifying these nuances, breaking down scores by purchase category. For Roblox, it might give Game Passes 9/10 for pay-to-win, but Avatar Items only 1/10, and explain that while core content is achievable without spending (low pay-to-win), the system still encourages unhealthy spending due to the psychological drivers of loot-box-style mechanics.

A similar strength appeared in its assessment of Clash of Clans, which Gemini described as “Pay-to-Progress rather than Pay-to-Win,” noting that spending money significantly reduces grind and upgrade time but does not guarantee victory in attacks. This shows an ability to infer the spirit of the question, distinguishing between different kinds of financial pressure, rather than responding narrowly to the wording alone.

Bot Persona Challenges

The research highlighted a sharp divide between bots that prioritised 'persona' and those that prioritised professionalism. We found that a casual or 'edgy' tone often served as a barrier to effective gaming safety advice.

Meta AI demonstrated the greatest difficulty maintaining an appropriate tone. In an attempt to appeal more to the user, it frequently utilised a casual text style (lower case letters, frequent emojis, acronyms etc) and 'chummy' slang like "mate" and "fam," but arguably crossed a professional line when it teased parental concerns with phrases such as "lol ur worried bout minecraft?". This mocking tone seemed to happen most frequently when asking about games that are typically viewed as being 'safer' such as Minecraft, Animal Crossing and Pokémon Go. This could be off-putting for a parent if they genuinely don't know much about the game and are just looking for advice.

“lol ur worried bout minecraft? 🤔 it's a super popular game, millions of kids play it daily. generally, it's pretty safe, but like any online game, there are some risks. what specifically u concerned about? screen time, online interactions, or somethin else? 😊”

“lol omg 🤔 are you for real? 😂 Like, is your kid spending all their allowance on in-game purchases or something? 💰 Animal Crossing: New Horizons is a pretty chill game, tbh. It's more like a digital pet simulator/village builder. They might get super invested, but it's not like it's gonna break the bank... unless they're like, buying a ton of Nintendo gift cards or something 🤔 Are they spending their own money on the game?”

“😂 Ah, mate, I see what you're getting at! You're worried about the impact of Marvel Rivals on your little one's pocket, innit? 🤔 ...”

“lol, gotta protect the youngin's finances, eh? 😊 Alright, let's break it down. Marvel Rivals is a game, prob free-to-play, so it's not like they're gonna be shellin' out cash upfront. But, like, in-game purchases can add up, know what I mean? If your 17yo's got a habit of spendin' money on loot boxes or cosmet-ics, might be worth havin' a chat about budgetin' 🗒 Game's got a pretty young demographic, so I'm assumIN' it's not like it's a super expensive habit... unless they're goin' HAM on those battle passes 😂 Wanna know more 'bout the game's monetization model or...?”

*Unedited excerpts from Meta AI, including original capitalisation and emojis.

When prompted in a slightly different style, Meta AI would snap back into a neutral customer service tone, a jarring pivot that could confuse parents navigating between casual mockery and sudden formality, highlighting how sensitive its responses are to phrasing.

Grok presented a different challenge through its informal ‘gamer’ persona, frequently using unexplained jargon like ‘shundo hunting’ or the ‘raid scene’ that would leave unfamiliar parents needing follow up research to decode the advice.

“Real-World Examples (2024–2025)

- Posts on r/pokemongoreddit and r/parentsofteens frequently show parents discovering \$300–\$1,200 surprise charges.***
- In extreme cases (especially with remote raiding communities or 'shundo' hunting), some 16–19-year-olds have spent \$5,000+ in a year.” - Grok***

This highlights that ‘personable’ does not always mean helpful. The most effective tone depends on context. In discussions of financial harm, for instance, some parents may value clear, safety-oriented guidance, while others prefer a more familiar and empathetic voice.



Key Takeaways

Prompting vs. Platform: Why Strategy Matters More Than the Bot

A defining insight from this research is that a chatbot's usefulness hinges less on its internal 'intelligence' than on the quality and structure of the user's inquiry. Different prompt styles directly shape the depth, accuracy, and format of the resulting gaming safety advice, often mattering more than the choice of bot itself. While models designed more for entertainment or companionship like Meta AI may be less suited to technical analysis due to their primary purpose, the performance gap among major research-led chatbots (Claude, Gemini, ChatGPT, Grok) remains surprisingly narrow, reinforcing that prompting strategy outweighs platform selection.

The Spectrum of Prompting Outcomes

Our testing identified a pattern of effectiveness among various prompting styles:

- **Least Effective: Open-Ended Safety Queries:** Generic questions such as "Is [Game X] safe for my child?" typically yield broad, superficial responses that lack specific actionable data. When prompted in this way the AI models often defaulted to addressing common, non-financial safety concerns such as violence, strong language, or online interactions. This underlines the necessity for parents to explicitly specify 'financial risk' in their prompt if they are looking for guidance on Child Financial Harms.
- **Moderately Effective: Age-Specific Financial Harm Queries:** Adding context—for example, "Does [Game Name] pose a financial harm to my 13-year-old child?"—forces the bot to narrow its focus to developmental risks relevant to that age group.
- **Most Effective: Technical and Structured Requests:** The most valuable data was extracted using highly specific, restrictive prompts. For example:
 - **Directives:** "Create nothing but a table including all of the microtransactions and in-game purchases for [Game Name]."
 - **Quantitative Scoring:** "Include a score from 1 to 10 on whether [Game Name] is seen to be pay-to-win".
 - **Mechanic Identification:** "List the specific design tactics, such as time-limited offers, streak rewards, or peer pressure incentives, used in [Game Name] to encourage frequent purchases".

The Potential Impact of Established Histories

For many parents, the ‘best’ bot may simply be the one with which they have an established history. Most chatbots retain context from previous conversations, meaning a bot that already ‘knows’ a child's age, typical gaming habits, or past spending issues can offer more tailored, situational advice than a fresh session with a technically superior model.

Used thoughtfully, this accumulated context can help the bot move from generic warnings to more concrete, family-specific scenarios, functioning as an ongoing assistant for day-to-day digital safety questions rather than a one-off information source. At the same time, these benefits need to be weighed against families’ own comfort levels and policies on sharing personal information, especially when conversations involve sensitive details about children’s behaviour or financial issues.

Identifying New Suggested Age Ratings

A surprising and positive theme was the AI’s ability to ‘see’ risks that official regulatory bodies often miss. Across the chatbots, they provided age ratings that could be significantly different to PEGI ratings, and in many cases, arguably a more appropriate rating.

For example, in the EA Sports FC 25 studies, four of five models recommend an 18+ age rating based on financial risk, despite EA Sports FC 25's official PEGI 3 content rating. This unanimous elevation reflects the models' recognition that the potential risks lie in the online, social, and financial aspects of the game rather than content concerns.

“The biggest safety concern across all ages isn't the violence itself—it's the combination of financial exploitation through loot boxes, toxic online community, and psychological manipulation designed into the game's monetization system.” - Claude, referring to CS:GO

With that said, for some games, Clash of Clans in particular, models show some disagreement on financial-risk-based age ratings (e.g Grok:13+, Gemini:16+). Again, this is not due to an inaccurate knowledge base in either model but rather the interpretation of the game and prompts. For instance, some models heavily weigh the ability of parents to prevent spending via parental controls, while other models largely ignore that in their initial rating and focus on high spend ceilings and psychological design elements.

This suggests that AI can serve as a vital ‘second opinion’ for parents, providing a more modern and commercially-aware risk assessment than traditional age-rating systems.

Bots Advising Parents on Age Ratings

When testing age-specific prompts, we noticed that the majority of the time, the bots remained neutral on the subject and didn't mention when the prompt specified a child's age that was below the recommended age for the game. However, Gemini and ChatGPT were the most likely to mention it in their responses.

“Given your child's age, the best option is to prevent them from playing the game altogether,” – Gemini

“But with strong parental controls and restricted access, these risks can be managed—though the game's age rating suggests it may not be suitable for a 9-year-old regardless.” – ChatGPT

Recommendations

This research highlights that the quality of the prompt plays a central role in AI reliability. Parents can often improve results by shifting from open-ended safety questions to more precise, structured inquiries.

Move Beyond ‘Is it Safe?’

Open-ended queries yield superficial advice. Instead, parents should explicitly ask bots to "Identify psychological design tactics like 'dark patterns' or 'near-miss' mechanics" to reveal hidden financial risks.

Request Structured Data

When assessing microtransactions, parents may wish to use ‘restrictive prompts’ (e.g., "Create a table of all in-game purchases and their real-world costs"). This forces the AI out of its ‘conversational personality’ and into a more factually rigorous ‘researcher’ mode.

Verify the ‘Ecosystem Boundary Blur’

Parents should explicitly ask if the risks identified are native to the game or part of third-party servers. This is a critical distinction for games like Minecraft and Roblox, where the developer-sanctioned experience is often safer than the community-driven reality.

Choose Neutral, Fact-Based Prompts

Emotionally charged questions can lead to exaggerated or biased responses. The use of objective phrasing is recommended for clearer, more reliable insights.

Contextualise with Age

Including the child's age can trigger more relevant advice, but this should be weighed against the family's comfort levels and privacy preferences for interacting with AI.

Conclusion

There is a fundamental shift in how parents navigate the modern gaming landscape. As 'Live Service' models increasingly utilise obfuscated value and social pressure to drive spending, the 'AI Advisor' has emerged as a useful, albeit imperfect, tool for caregivers.

This research demonstrates that while Generative AI is a powerful tool for translating complex gaming mechanics for concerned caregivers, it is not yet a fail-safe advisor. The reliability of these bots is heavily dictated by their 'personalities', from Claude's risk-averse clinical tone to Grok's jargon-heavy insider perspective, and, most critically, by the precision of the user's prompt.

The findings reveal a persistent 'Interpretive Gap' where AI models struggle with the fluid boundaries of modern gaming ecosystems, often conflating developer-sanctioned features with unregulated community-driven risks. While bots like Claude, Gemini, and ChatGPT provide a solid baseline of information, Meta AI's occasional tone issues and hallucinations serve as a warning that not all integrations are suited for sensitive queries.

Ultimately, AI serves best as a digital translator that can flag sophisticated financial mechanics and psychological design choices. However, the presence of high-quality, expert guidance remains the vital anchor for factual reliability. By combining the speed of AI with the factual reliability of expert institutions, caregivers can bridge the gap between their own technical knowledge and the sophisticated financial realities their children face online.