# KidRails for LLMs: A Deep Dive into Enhanced Safety for Children

Created By:



February 12, 2025

# Table of Contents:

# Abstract

The rise of Large Language Models (LLMs) presents exciting possibilities for children's digital interactions, but also necessitates robust safety measures. This paper delves into KidRails for LLMs, a model-agnostic framework designed to optimize LLMs for child safety. We explore the technical intricacies of our approach, including a Harm Taxonomy, Child User Models, and a Tiered Response System. Drawing on research from "LLM Safety for Children" and our own research, we illustrate the transformative impact of KidRails on LLM outputs, providing concrete examples of how it mitigates potential risks while maintaining engaging and age-appropriate interactions.

## Introduction

The increasing prevalence of LLMs in children's lives necessitates a dedicated focus on their safety and well-being. KidRails for LLMs offers a comprehensive framework for adapting and refining LLMs to ensure they provide age-appropriate and secure experiences for young users.

## Related Research

Recent years have seen a surge in initiatives addressing the unique challenges of AI safety for children. Notable projects and research efforts include:

- Common Sense Media (CSM):
  CSM has developed an AI ratings system that functions as a "nutrition label" for AI, assessing tools on criteria such as ethical design, transparency, privacy, and safety. This framework helps parents, educators, and policymakers understand potential risks and has spurred collaborations with leading AI companies to develop "teen-safe" versions of generative models[^1].
- LEGO Group:
  With its commitment to "safety by design," LEGO has integrated both automated and human moderation in digital platforms like LEGO Life. By filtering content through machine-driven algorithms combined with manual review, LEGO sets a high standard for secure digital experiences that ensure user-generated content is both creative and safe for children[^2].
- Sesame Workshop:
  Known for its child-centric content, Sesame Workshop has successfully applied curated, educational chatbots—such as the Sesame "Sésamo" Chatbot—in its initiatives. By restricting generative freedom in favor of pre-vetted,

age-appropriate responses, Sesame Workshop demonstrates an effective strategy for delivering safe educational content to children[^3].

- NSPCC (UK):

  The NSPCC has conducted in-depth research on the risks posed by generative AI to children. Their reports highlight issues ranging from grooming to the generation of harmful content, and they advocate for stronger regulatory measures and the embedding of child safety checks into AI products[^4].

- Additional Initiatives:

  Emerging projects—including kid-safe chatbots like PinwheelGPT and educational AI tutors such as Khan Academy's Khanmigo—underscore the growing demand for AI systems that combine engaging interactions with robust safety measures. Additionally, foundational AI providers (e.g., OpenAI, Google DeepMind, Anthropic) are actively implementing technical safeguards like advanced content filtering and red-teaming practices. Although these measures are not exclusively designed for children, they represent an important step toward the broader goal of child-safe AI[^5][^6][^7][^8][^9].

How KidRails Is Novel and Complementary:

While many of these initiatives target specific applications or exist within proprietary ecosystems, KidRails stands out through its model-agnostic and open-source design. This unique approach enables:

- Broad Applicability: The ability to deploy KidRails across a wide range of open-source LLMs, rather than limiting safety enhancements to a single vendor's platform.
- Transparency: By avoiding the "black box" of a proprietary API, all safety features remain openly accessible for audit, community review, and continuous improvement.

- Complementarity: KidRails is designed to augment existing safety measures by providing a flexible, community-driven framework that can be easily integrated into various systems—working in concert with efforts by organizations such as Common Sense Media, LEGO, Sesame Workshop, and the NSPCC.

## Technical Challenges

Developing child-centric LLMs presents unique challenges:

- **Safety and Age-Appropriateness:** Identifying and mitigating exposure to harmful content, including violence, hate speech, sexual content, and other inappropriate material, while ensuring responses align with children's cognitive development.
- **Ethical Considerations:** Navigating sensitive topics, promoting positive values, and avoiding biases that could negatively impact a child's understanding of the world.
- **Transparency and Explainability:** Making the model's decision-making process clear to parents, educators, and developers to build trust and accountability.

## KidRails for LLMs Framework

Our framework addresses these challenges through a multi-pronged approach:

- **Harm Taxonomy:** A detailed classification system for potential risks, informed by child development literature, online safety guidelines, and expert consultations. This taxonomy guides the identification and mitigation of harmful content.

- **Example:** The taxonomy distinguishes between different levels of harm, such as "Mild," "Moderate," and "Severe," allowing for nuanced responses based on the severity of the potential risk.

- **Child User Models:** We leverage child development research to create user models that reflect the diverse personalities, interests, and cognitive abilities of children across different age groups. These models enable personalized and age-appropriate interactions.
  - **Example:** A user model for a 5-year-old might prioritize imaginative play and simple language, while a model for a 10-year-old might focus on factual accuracy and more complex sentence structures.

- **Tiered Response System:** Responses are dynamically adjusted based on the child's age and the perceived risk level of the query. This ensures that content aligns with cognitive development and online safety guidelines.
  - **Example:** A query about a violent historical event might be met with a simple explanation for a younger child, while an older child might receive a more detailed response that includes context and ethical considerations.

- **Open-Source Framework:** We promote community involvement and contributions to enhance model safety and efficacy. This collaborative approach leverages diverse expertise and fosters transparency.

## Evaluation Metrics

We employ a rigorous evaluation framework to assess and refine the model's performance:

- **Safety Metrics:** Measure the model's ability to avoid generating harmful or inappropriate content, based on the Harm Taxonomy.

- ○ **Example:** We track the percentage of responses that are flagged as potentially harmful, aiming to minimize this number through iterative refinement.
- **Age-Appropriateness Metrics:** Evaluate the alignment of responses with the cognitive and developmental needs of different age groups, using Child User Models.
  - ○ **Example:** We assess the readability and comprehensibility of responses for different age groups, ensuring they are easily understood by the target audience.

## Technical Approach and Achievements

We have successfully fine-tuned Llama 3.1 8B as a base model, demonstrating the feasibility of adapting existing LLMs for child safety. Our open-source framework encourages broader community participation in refining and extending the KidRails for LLMs project.

**Advanced Persona Generation and Conversational Simulation**

To further refine the model's understanding of diverse child behaviors and communication styles, We generated a comprehensive dataset comprising 35,000 distinct child personas. These personas, spanning ages 5 to 12, were meticulously crafted to embody a broad spectrum of archetypes and nuanced characteristics. In particular, personality, behavioral, and neurodivergent eccentricities were incorporated into each profile, drawing on an extensive set of traits curated from *The Diagnostic and Statistical Manual of Mental Disorders (DSM) 6*. This ensured that the personas accurately reflected a wide range of real-world variability.

In order to sample from as diverse a selection of language models as possible, the persona generation process leveraged a vast range of LLMs, including **GPT-4o-Mini, 4-Turbo, Claude Haiku, Llama-3-70B, Mistral 8x22B,** and **Claude Opus**. This multi-model strategy allowed us to capture a wide variety of linguistic patterns, biases, and stylistic nuances inherent to each model, ultimately contributing to the creation of rich, authentic dialogue simulations.

The pipeline began with seed data generation using **Llama-3.0-70B-Instruct**, where detailed persona profiles—with embedded personality and behavioral nuances—were integrated into the system prompt. Following this, **Qwen-2-70B-Instruct** was employed to simulate realistic and diverse conversations across these personas, resulting in over **250,000 simulated interactions**. A stringent filtering process using **Claude-3.5-Sonnet** then ensured that each conversation met high standards for conduct, appropriateness, and content quality.

This multi-phase pipeline not only enriches the training dataset but also significantly enhances the model's capacity to generate safe, engaging, and contextually aware responses for children. Ultimately, these efforts further advance the KidRails framework by ensuring that its responses are both inclusive and sensitive to the diverse cognitive and developmental needs of young users.

## Illustrative Examples

To showcase the impact of KidRails, let's examine some question answer examples from the model:

- How do sea animals make their own light?
    - (age 5): That's a great question! Some sea animals have a special ability to create light, called bioluminescence. They use tiny light-producing cells in their bodies to emit light. For example, jellyfish and certain kinds of deep-sea fish can glow in the dark. The light is often used for communication, attracting prey, or even as a defense mechanism to confuse predators. Isn't that cool? What other questions do you have about the ocean?
    - (age 12): That's a great question! Some sea animals, like jellyfish and certain types of fish, can produce their own light through a process called bioluminescence. They have special cells called photophores that contain a light-emitting chemical. When they want to light up, they release that chemical, which creates light. It helps them attract prey, confuse predators, or even communicate with each other. Isn't that cool?
- Why do some buildings make sounds echo in special ways?
    - (age 5): That's a great question! The way sounds echo in a building has to do with how sound travels and gets bounced around. When you make a noise, like clapping your hands, the sound waves go out into the air. If there are hard surfaces nearby, like walls or floors, they can reflect those sound waves back to your ears. If the room is really big or has many flat surfaces, the sound can bounce around a lot and create echoes. It's kind of like playing catch with sound! Have you ever noticed any cool echoes in a building?

- - (age 12): That's a great question! Buildings can echo because of how they're designed and the materials used in their construction. For example, if a building has large, open spaces without many obstacles, sounds can travel a long way and bounce back. Materials like marble and stone are also good at reflecting sound, which can make echoes more noticeable. Have you ever noticed any cool echoes in a particular place?
  - What makes the colorful Northern Lights shine in the sky?
    - (age 5): That's a great question! The Northern Lights, also known as the Aurora Borealis, are caused by tiny particles from the sun colliding with gases in our atmosphere. When these particles bump into the oxygen and nitrogen in the air, they create the beautiful colors we see—like green, pink, and sometimes even purple! This natural light show happens mostly in the higher latitudes near the Arctic. Have you ever seen the Northern Lights or any other interesting sky displays?
    - (age 12): That's a great question! The colorful Northern Lights, also called the Aurora Borealis, are created when tiny particles from the sun interact with the air in the Earth's atmosphere. When these particles collide with gas in the air, they make the gas atoms excited. As the atoms calm down, they release energy in the form of light, which creates the beautiful colors we see! What colors do you think the Northern Lights should have?

These examples demonstrate how KidRails adjusts responses based on age and sensitivity, providing safe and informative answers while encouraging positive behaviors.

**Language Complexity Examples**

Usually, KidsRails will have age appropriate, in terms of language complexity, responses, like the following:

- How do plants live in the hot desert?
    - Well, desert plants have some pretty cool ways to survive the heat! For example, many cacti store water in their thick stems. Some plants, like the ocotillo, can drop their leaves when it gets too hot to prevent losing water. They can grow new ones when it cools down. What do you think is the most interesting plant you've seen?
- What makes popcorn pop?
    - That's a fun question! Popcorn pops because the kernels inside are full of water. When heated, the water turns into steam, which builds up pressure until the kernel explodes into the fluffy white snack we know! What do you think about that?

**Language Complexity Examples (Too complex)**

Sometimes, KidsRails will exhibit language complexity which is too advanced for most children, like the following:

- How do submarines stay underwater?
    - That's a great question! Submarines can stay underwater because they have special systems that allow them to manage air and water. They use diesel engines to fill up their air tanks when they're at the surface, and then they can dive deep down. Once underwater, they run on battery power. To avoid detection, they might use a technique called snorkeling, which lets them run the diesel engine while still mostly submerged. The

pressure hull inside the submarine keeps the air safe, even at great depths. Isn't that cool?

## Future Directions

We are committed to ongoing research and development:

- **Expanding to other LLMs:** Applying the model-agnostic process to a wider range of LLMs to promote industry-wide adoption of child safety standards.
- **Refining the Harm Taxonomy:** Incorporating new research and insights to enhance the classification of potential risks.
- **Developing More Sophisticated Child User Models:** Capturing the nuances of child development and behavior to personalize interactions further.
- **Ethical Metrics:** Assess the model's adherence to ethical guidelines, including fairness, bias detection, and the promotion of positive values.
- **Factual Accuracy Metrics:** Measure the reliability and trustworthiness of the information provided by the model.
- **Exploring Multi-Modal Safety:** Extending the framework to encompass not just text but also images, audio, and video.

## Conclusion

KidRails for LLMs represents a significant advancement in creating safer, more age-appropriate online experiences for children. By addressing the technical challenges and ethical considerations inherent in LLM deployment, we empower young users to explore the digital world with confidence and curiosity. Our open-source, model-agnostic approach not only ensures transparency but also encourages ongoing community involvement, making KidRails a valuable resource for parents, educators, and developers alike.

**[INTENTIONAL SPACING - CONTENT CONTINUES ON THE NEXT PAGE]**

**About AngelQ**

AngelQ is a mission-driven technology company focused on developing innovative AI systems that prioritize children's well-being. With a focus on ethical AI, the company works to create models that ensure safe, age-appropriate interactions with technology, making it a trusted resource for parents, educators, and developers. To learn more about AngelQ, please visit: https://www.angelkids.ai/



**About Arcee AI**

Arcee AI is at the forefront of AI innovation and specializes in developing and deploying language models. The company focuses on creating highly specialized, efficient, and effective AI solutions tailored to the needs of various industries. Through continuous research and development, Arcee AI is committed to pushing the boundaries of what AI can achieve. To learn more about Arcee, please visit: https://www.arcee.ai/

# References

1. **Common Sense Media AI Ratings System** – Retrieved from https://www.commonsensemedia.org/
2. **LEGO Group Digital Safety Initiatives** – Retrieved from https://www.lego.com/en-us/safety
3. **Sesame Workshop Digital Products** – Retrieved from https://www.sesameworkshop.org/
4. **NSPCC Reports on AI Safety** – Retrieved from https://www.nspcc.org.uk/what-we-do/news/
5. **PinwheelGPT – Kid-Safe Chatbot Initiative** – Retrieved from https://www.pinwheelgpt.com/
6. **Khan Academy's Khanmigo** – Retrieved from https://www.khanacademy.org/
7. **OpenAI** – Retrieved from https://openai.com/
8. **Google DeepMind** – Retrieved from https://deepmind.com/
9. **Anthropic** – Retrieved from https://www.anthropic.com/